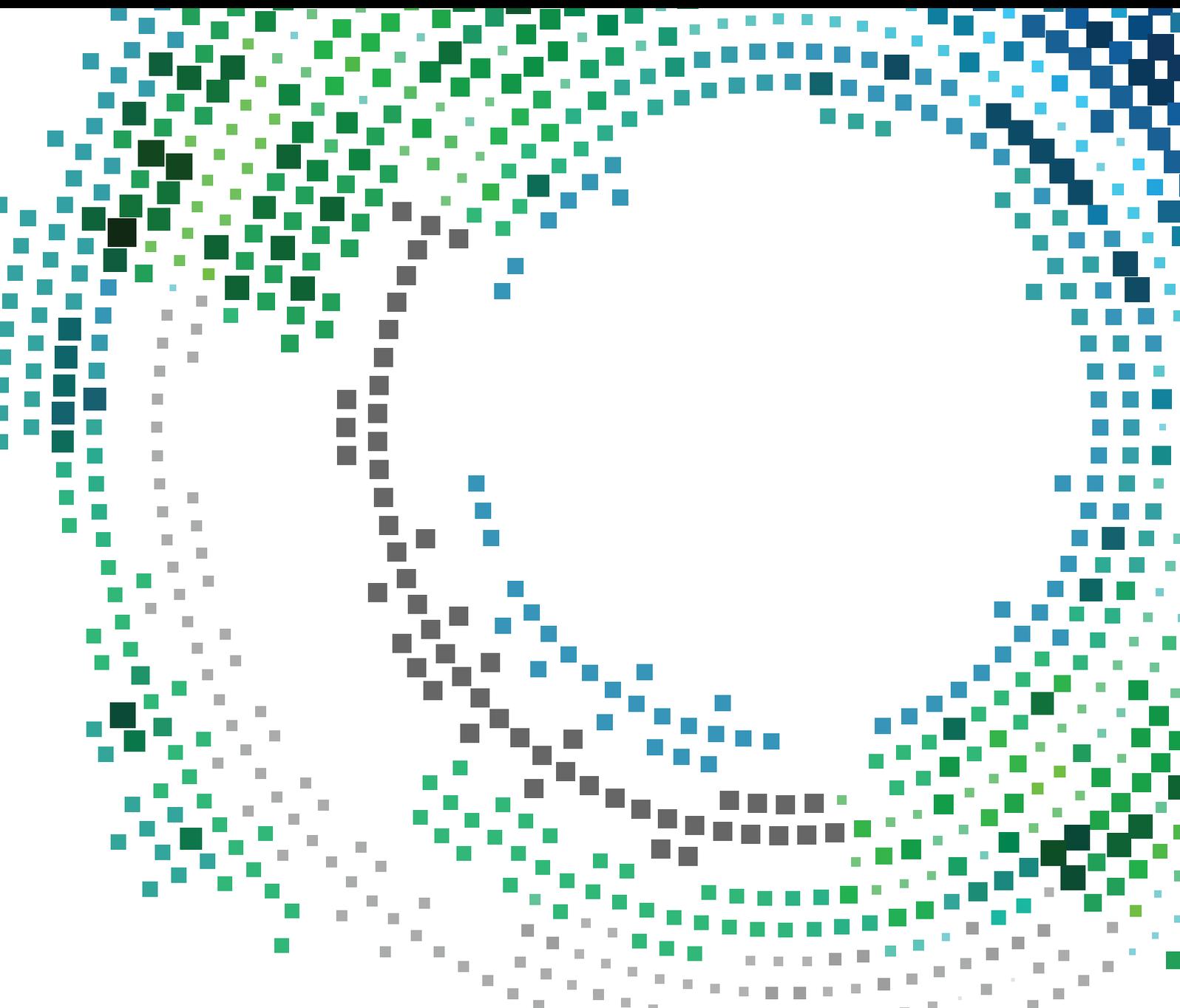


Location-Based Mobile Marketing Innovations 2018

Lead Guest Editor: Jaegeol Yim

Guest Editors: Subramaniam Ganesan and Byeong H. Kang





Location-Based Mobile Marketing Innovations 2018

Mobile Information Systems

Location-Based Mobile Marketing Innovations 2018

Lead Guest Editor: Jaegeol Yim

Guest Editors: Subramaniam Ganesan and Byeong H. Kang



Copyright © 2019 Hindawi. All rights reserved.

This is a special issue published in “Mobile Information Systems.” All articles are open access articles distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Editorial Board

Mari C. Aguayo Torres, Spain
Ramon Agüero, Spain
Markos Anastassopoulos, UK
Marco Anisetti, Italy
Claudio Agostino Ardagna, Italy
Jose M. Barcelo-Ordinas, Spain
Alessandro Bazzi, Italy
Luca Bedogni, Italy
Paolo Bellavista, Italy
Nicola Bicocchi, Italy
Peter Brida, Slovakia
Carlos T. Calafate, Spain
María Calderon, Spain
Juan C. Cano, Spain
Salvatore Carta, Italy
Yuh-Shyan Chen, Taiwan
Wenchi Cheng, China
Massimo Condoluci, Sweden
Antonio de la Oliva, Spain
Almudena Díaz Zayas, Spain

Filippo Gandino, Italy
Jorge Garcia Duque, Spain
L. J. García Villalba, Spain
Michele Garetto, Italy
Romeo Giuliano, Italy
Prosanta Gope, UK
Javier Gozalvez, Spain
Francesco Gringoli, Italy
Carlos A. Gutierrez, Mexico
Ravi Jhavar, Luxembourg
Peter Jung, Germany
Adrian Kliks, Poland
Dik Lun Lee, Hong Kong
Ding Li, USA
Juraj Machaj, Slovakia
Sergio Mascetti, Italy
Elio Masciari, Italy
Maristella Matera, Italy
Franco Mazzenga, Italy
Eduardo Mena, Spain

Massimo Merro, Italy
Aniello Minutolo, Italy
Jose F. Monserrat, Spain
Raul Montoliu, Spain
Mario Muñoz-Organero, Spain
Francesco Palmieri, Italy
José J. Pazos-Arias, Spain
Marco Picone, Italy
Vicent Pla, Spain
Amon Rapp, Italy
Daniele Riboni, Italy
Pedro M. Ruiz, Spain
Michele Ruta, Italy
Stefania Sardellitti, Italy
Filippo Sciarrone, Italy
Florian Scioscia, Italy
Michael Vassilakopoulos, Greece
Laurence T. Yang, Canada
Jinglan Zhang, Australia

Contents

Location-Based Mobile Marketing Innovations 2018

Jaegel Yim , Subramaniam Ganesan , and Byeong Ho Kang 
Editorial (2 pages), Article ID 2164708, Volume 2019 (2019)

A Study on Removing Cloud Drift of Sky-Sea Infrared Image Based on Agent

Jianming Sun, Zhipeng Fan , Zhihui Sun, Qinghua Miao, and Xiaoyan Li
Research Article (9 pages), Article ID 8505219, Volume 2019 (2019)

Consumers Team Detection Model Based on Trust for Multi-Level

Xiaoming Li, Guangquan Xu , Sandhya Armoogum , and Honghao Gao 
Research Article (10 pages), Article ID 4147859, Volume 2019 (2019)

RoC: Robust and Low-Complexity Wireless Indoor Positioning Systems for Multifloor Buildings Using Location Fingerprinting Techniques

Kriangkrai Maneerat  and Kamol Kaemarungsi 
Research Article (22 pages), Article ID 5089626, Volume 2019 (2019)

Research on Precision Marketing Model of Tourism Industry Based on User's Mobile Behavior Trajectory

Jialin Zhang, Tong Wu, and Zhipeng Fan 
Research Article (14 pages), Article ID 6560848, Volume 2019 (2019)

Location-Based Test Case Prioritization for Software Embedded in Mobile Devices Using the Law of Gravitation

Xiaolin Wang , Hongwei Zeng, Honghao Gao , Huaikou Miao , and Weiwei Lin 
Research Article (14 pages), Article ID 9083956, Volume 2019 (2019)

The Identification of Marketing Performance Using Text Mining of Airline Review Data

Jae-Won Hong and Seung-Bae Park 
Research Article (8 pages), Article ID 1790429, Volume 2019 (2019)

From Reputation Perspective: A Hybrid Matrix Factorization for QoS Prediction in Location-Aware Mobile Service Recommendation System

Shun Li , Junhao Wen , and Xibin Wang 
Research Article (12 pages), Article ID 8950508, Volume 2019 (2019)

Dilemma and Solution of Traditional Feature Extraction Methods Based on Inertial Sensors

Zhiqiang Peng  and Yue Zhang 
Research Article (6 pages), Article ID 2659142, Volume 2018 (2019)

A Case Study Analysis of Clothing Shopping Mall for Customer Design Participation Service and Development of Customer Editing User Interface

Ying Yuan and Jun-Ho Huh 
Research Article (19 pages), Article ID 7698648, Volume 2018 (2019)

An Indoor Location-Based Positioning System Using Stereo Vision with the Drone Camera

Young-Hoon Jin , Kwang-Woo Ko, and Won-Hyung Lee 

Research Article (13 pages), Article ID 5160543, Volume 2018 (2019)

Profile-Based Ad Hoc Social Networking Using Wi-Fi Direct on the Top of Android

Nagender Aneja  and Sapna Gambhir 

Research Article (7 pages), Article ID 9469536, Volume 2018 (2019)

CEnsLoc: Infrastructure-Less Indoor Localization Methodology Using GMM Clustering-Based Classification Ensembles

Beenish Ayesha Akram , Ali Hammad Akbar , and Ki-Hyung Kim 

Research Article (11 pages), Article ID 3287810, Volume 2018 (2019)

Malaria Vulnerability Map Mobile System Development Using GIS-Based Decision-Making Technique

Jung-Yoon Kim , Sung-Jong Eun, and Dong Kyun Park 

Research Article (9 pages), Article ID 8436210, Volume 2018 (2019)

Location Privacy Protection Research Based on Querying Anonymous Region Construction for Smart Campus

Ruxia Sun , Jinwen Xi , Chunyong Yin , Jin Wang , and Gwang-jun Kim 

Research Article (11 pages), Article ID 3682382, Volume 2018 (2019)

Editorial

Location-Based Mobile Marketing Innovations 2018

Jaegel Yim ¹, **Subramaniam Ganesan** ², and **Byeong Ho Kang** ³

¹Professor Emeritus, Department of Computer Engineering, Dongguk University, Gyeongju, Republic of Korea

²Professor, Electrical and Computer Engineering, Oakland University, Rochester, USA

³Professor, School of Engineering and ICT, University of Tasmania, Hobart, Australia

Correspondence should be addressed to Jaegel Yim; yim@dongguk.ac.kr

Received 15 April 2019; Accepted 15 April 2019; Published 28 April 2019

Copyright © 2019 Jaegel Yim et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

The increasing complexity of the industry means that marketers must now be experts not only in marketing, but also in people, data, delivery platforms, and mobile location-based marketing.

The objective of location-based marketing via mobile devices is to encourage those activities, as well as to drive foot traffic, share discounts, and build customer loyalty. The mobile devices are used to gather information about nearby businesses including reviews, directions, calling the business, and using the businesses' mobile app. With location-based mobile marketing, the business is easy to find and have skillfully combined location-based marketing with an overall targeted marketing approach that includes social media, push notifications, e-mail newsletters, and even offline marketing.

The objective of this special issue is to bring together research contributions of unpublished research on the recent development and innovations about the location-based service mobile marketing. This aims to facilitate and support research in the current e-commerce innovation that makes business easy, engaging, and at hand of the consumers.

The paper "The Identification of Marketing Performance Using Text Mining of Airline Review Data" aims to firstly extract major keywords using text mining method, secondly to identify prominent keyword from the keywords extracted from text mining analysis, and then to confirm differences in influences of the keywords which affect corporate performance.

In the paper "Dilemma and Solution of Traditional Feature Extraction Methods Based on Inertial Sensors," after analyzing the difference of these indistinguishable movements, the authors propose several new features to improve accuracy of recognition. They compare the traditional features and their custom features. In addition, they examined

whether the time domain features and frequency domain features based on acceleration and angular velocity are different.

In the paper "CEnsLoc: Infrastructure-Less Indoor Localization Methodology Using GMM Clustering-Based Classification Ensembles," the authors propose CEnsLoc, a new easy to train-and-deploy Wi-Fi localization methodology established on GMM clustering and random forest ensembles (RFE). Principal component analysis was applied for dimension reduction of raw data. Conducted experimentation demonstrates that it provides 97% accuracy for room prediction, whereas artificial neural networks, k -nearest neighbors, K^* , FURIA, and DeepLearning4J-based localization solutions provided mean 85%, 91%, 90%, 92%, and 73% accuracy on the collected real-world dataset, respectively. It delivers high room level accuracy with negligible response time, making it viable and befitted for real-time applications.

In "Location Privacy Protection Research Based on Querying Anonymous Region Construction for Smart Campus," the user's query range is introduced to present a novel anonymous region construction scheme. The anonymous server first generates the original anonymous sub-regions according to the user's privacy requirements and then merges them to construct the anonymity regions submitted to LSP based on the size of corresponding querying regions. The security and experiment analyses show that the proposed scheme not only protects the user's privacy effectively but also decreases the area of LSP querying regions and the region-constructing time, improving the quality of service for smart campus.

The paper "Consumers Team Detection Model Based on Trust for Multi-Level" proposes a novel local community detection model E-MLCD. It is jointly based on the multilevel

properties and the strength of similarity of multilevel social interaction among communities. By studying three real-world multilevel social networks and specific QQ zone marketing data, the model defines a new metric of similarity strength based on community structure similarity. Comparison with other state-of-the-art detection methods demonstrates E-MLCD's ability to detect communities more effectively.

In "RoC: Robust and Low-Complexity Wireless Indoor Positioning Systems for Multifloor Buildings Using Location Fingerprinting Techniques," the authors propose a novel integrated framework for wireless indoor positioning systems based on a location fingerprinting technique which is called the robust and low complexity indoor positioning systems framework (RoC framework). The proposed integrated framework consists of two essential indoor positioning processes: the system design process and the localization process. The RoC framework aims to achieve robustness in the system design structure and reliability of the target location during the online estimation phase either under a normal situation or when some reference nodes (RNs) have failed.

The paper "An Indoor Location-Based Positioning System Using Stereo Vision with the Drone Camera" proposes the indoor location-based drone controlling method that does not require the traditional remote controller and can be applied to various services such as a group flight.

In the paper "Research on Precision Marketing Model of Tourism Industry Based on User's Mobile Behavior Trajectory," data mining clustering technology is used to analyze the characteristics of user's mobile behavior trajectories, and the precise recommendation system of tourism is constructed to support for tourism decision making. It can target the tourist group for the precise marketing and make tourist travel smarter.

The paper "A Case Study Analysis of Clothing Shopping Mall for Customer Design Participation Service and Development of Customer Editing User Interface" discusses a service related to the convergence of the traditional clothing industry with IT and a service wherein CT is converged with systems that allow customers to participate in the design work and share the designs they have created. The results show that both production method and production capacity largely affect the user interface of apparel platform services, with customer freedom significantly correlated with their functional roles. Moreover, the lead index is shown to be one of the factors restraining customer freedom.

The paper "Malaria Vulnerability Map Mobile System Development Using GIS-Based Decision-Making Technique" aims to improve the lack of GIS information use and compatibility of multiplatform which represented limits that existing malaria risk analysis tools have. For this, the authors developed mobile web-based malaria vulnerability map system using GIS information. This system consists of system database construction, malaria risk calculation function, visual expression function, and website and mobile application.

In the paper "A Study on Removing Cloud Drift of Sky-Sea Infrared Image Based on Agent," a new shadow extraction

method is proposed. This method tests and removes cloud of infrared images based on cloud characteristics from infrared sky-sea images. Through grey value characteristics of cloud, we can find and use reactive agent layer structure and classify many agents used for local image cloud searching and manage agents used for coordinating many cloud searching agent.

The paper titled "From Reputation Perspective: A Hybrid Matrix Factorization for QoS Prediction in Location-Aware Mobile Service Recommendation System" proposes a hybrid matrix factorization method integrated location and reputation information (LRMF) to predict the unattainable QoS values. The proposed method effectively reduces the impact of unreliable users on QoS prediction and makes credible mobile service recommendation.

The paper "Location-Based Test Case Prioritization for Software Embedded in Mobile Devices Using the Law of Gravitation" uses a smart mall as a scenario to design a novel location-based test case prioritization (TCP) technique for software embedded in mobile devices using the law of gravitation. An empirical evaluation is presented by using one industrial project. The observation, underlying the experimental results, is that the proposed TCP approach performs better than traditional TCP techniques. In addition, besides location information, the level of devices is also an important factor which affects the prioritization efficiency.

The paper "Profile-Based Ad Hoc Social Networking Using Wi-Fi Direct on the Top of Android" presents an architecture and implementation of social networks on commercially available mobile devices that allow broadcasting name and a limited number of keywords representing users' interests without any connection in a nearby region to facilitate matching of interests. The broadcasting region creates a digital aura and is limited by Wi-Fi region that is around 200 meters.

We are very happy to publish this special issue of the mobile information systems. This issue contains 14 articles. Achieving such a high quality of papers would have been impossible without the huge work that was undertaken by the Editorial Board members and External Reviewers. We take this opportunity to thank them for their great support and cooperation.

Conflicts of Interest

The editors declare that they have no conflicts of interest.

*Jaegool Yim
Subramaniam Ganesan
Byeong Ho Kang*

Research Article

A Study on Removing Cloud Drift of Sky-Sea Infrared Image Based on Agent

Jianming Sun,^{1,2} Zhipeng Fan ,^{1,2} Zhihui Sun,¹ Qinghua Miao,¹ and Xiaoyan Li¹

¹Harbin University of Commerce, Harbin, Heilongjiang 150028, China

²Heilongjiang Provincial Key Laboratory of Electronic Commerce and Information Processing, Harbin, Heilongjiang 150028, China

Correspondence should be addressed to Zhipeng Fan; hsdfzp@126.com

Received 25 September 2018; Accepted 26 December 2018; Published 3 April 2019

Guest Editor: Subramaniam Ganesan

Copyright © 2019 Jianming Sun et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

According to the characteristics of cloud-drift area in an infrared sky-sea line image, we put forward a new shadow extraction method by using a layered reactive agent. The layered reactive agent includes a lower cloud search agent and higher management agent. The lower distributed layer search agent can sense local information of the image and retrieve the local cloud area in the image by locating operating point, enlarging scope, moving, marking, perishing, and so on. The higher management agent can sense the information of the whole and restrict and guide search action to the lower distributed search agent. The result of the simulation test shows the efficiency of this method is high, and this method extracts effectively and removes cloud in the infrared sky-sea image.

1. Introduction

To the infrared image of sky-sea background, we can ensure ship attitude by extracting sky-sea line. Cloud appears when the target image is operated by infrared, which will directly influence the success rate of sky-sea line recognition algorithm in that area. Meanwhile, the cloud-drift line often appears on the sea, which will be mistaken for sky-sea line. Therefore, the study on removing cloud has great significance on extracting sky-sea line. To solve this problem, some scholars propose several solutions. They classify the prior knowledge according to the environment condition [1]. They obtain the law of cloud change by statistical characteristics, compare the real image and this change law, and then build a thin cloud layer and cloud-drift removal model [2]. This method needs to use prior knowledge of cloud in some period before operation, so many image data cannot be operated because of the lack of enough prior knowledge. Another kind of method is recognition algorithm of cloud extraction without prior knowledge, like the search method of histogram peak value [3]. This method is simple, so it cannot be used because of its own shortcoming especially in complicated situation. And, this method can only be used in some special images, and hence, it is not fit for other images.

Through observing enough sky-sea infrared images, we can find that brightness of sky outside the cloud decreases significantly, the grey value in the same cloud changes little compared with the noncloud area, and the grey value between clouds distributed around the image has little change.

Agent refers to an active entity which stays in some situation and acts autonomously and cooperates with other Agents in order to achieve the goal of design. Residence, autonomy, and sociality are the basic features [4]. A multi-agent system is composed of Agents which cooperate with. Because of residence and autonomy of the Agent, it can adapt to environmental change effectively, which makes the system more flexible. Sociality of the Agent is an effective way to solve problems. Therefore, MAS is considered as an enabling technology which supports complicated technology and has great potential [5]. This technology offers many engineering ways like high-level abstraction, resolution of problem, and layer and system organization. Some practices show the advantages and potentials of MAS in some complicated system developments, such as production process, intelligence control, distributed network management, distributed control system, and telemedicine system [6].

Agent has been the focus of the research of artificial intelligence and software science because of its features like

perceptions [7], problem-solving ability, and cooperation. People make the program of Agent which can do random search and grey test of similar area. This program is applied to computer tomography, and the results are very good. In light of cloud scope and its distribution, this paper reveals an image processing method based on Agent. This method overcomes the shortcomings of the above methods and can extract efficiently and remove cloud of infrared images in different conditions.

2. Agent Model

2.1. Individual Agent Layer. In individual Agent layer, the research on AOP design abstract and model is to offer description software for autonomous decision action of the Agent. The inside of the Agent has higher concept and abstract, whose core is how to describe decision of the Agent inside. Then, the corresponding individual Agent software model is built. The individual Agent software model supported by AOP has 4 types: knowledge mode, cognitive mode, reactive mode, and hybrid mode. The knowledge mode takes Agent as a knowledge system [8]. The inside structure is composed of abstract and concept based on belief, knowledge, and distributed knowledge. We can use belief modification, knowledge reasoning, situation calculus, and so on to support autonomous decision of Agents by logic tools. AOP language which supports the knowledge mode includes Golog, AGTGolog, and Con Golog. The cognitive mode takes the Agent as a cognitive system. Based on cognitive science and folk psychology like goal, desire, intention, and plan, the inside structure of the Agent can be described by practical reasoning, BDI, and KARO to support Agent's autonomous decision [9]. Most AOP languages can be integrated into an idea, like Agent-0, 3APL, PLACA, Agent Speak (L), AOPLID, GOAL, Dribble, CLAIM, and 2APL [10]. Agent-0 takes the Agent as an active entity which consists of belief, capability, commitment, and action. PLACA expands intention cognitive component to support goal-directed action [11]. Agent Speak (L) and CLAIM are based on the BDI model. The Agent model of 3APL is based on belief, goal, and plan. Furthermore, target concepts of ADP can be divided into procedural one and declarative one. Procedural target corresponds to specific planning, focusing on goal-to-do [12]. Declarative target focuses on goal-to-be, which introduces GOAL. Dribble and 2APL support these two types of targets. The reactive mode takes the Agent as a reactive system [13], which can sense environment and its change and can response to these changes and process them. The reactive Agent model contains events and reactive rules, which can support autonomous decision of individual Agent's action based on events and events processing. For example, SLABSp uses action rule to define the action of Agent and describes the action that the Agent takes when some scene can be met [14]. The hybrid mode integrates the above software types to support structure and realization of the Agent, which uses various abstract and concepts to describe every element that can construct the model.

2.2. Multiagent Layer. In multiagent layer, the research on AOP design abstract and model focuses on how to provide effective concepts and models to support action of the Agent in MAS and organize and adapt to it, which can ensure its operation of cooperation of MAS and acquire the whole action of MAS. In terms of the software development, organization idea provides a feasible decomposition for the design of MAS [15]. And, a diversity of organization structure provides feasible structure for MAS, including layer organization, holonic organization, league, team, gathering, society, federation, market, and matrix organization. Now, organization structure supported by AOP includes team, regulatory organization, structured organization and hybrid organization. Team takes many Agents, which complete a complicated task by cooperation, as a team. This model often uses the traditional BDI model of the Agent to expand joint intention and team planning to build and describe team and guide or control the decision of Agent to achieve the cooperation of multiagent. For example, Simple Team describes multiagent team by describing some concepts like increasing roles to Jack, team ability, team planning, and so on. The Agent can achieve cooperation by executive team. EAMCORE can fulfill coordination between Agents based on team planning and group belief. A BDI structure of the individual Agent expands content and context in order to describe group concept [16]. Regulatory organization takes organization as a group of Agent and a rule set, which can define organization's control and restriction to the action and interaction of the Agent based on laws from sociology. According to different nature, law includes obligation, permission, and prohibition. In light of different contents of restriction, law can be classified action law and status law. And, in light of different enforcement mechanisms, law can be divided into regimentation law and sanction law. Now, typical AOP language supporting regulatory organization has ISLANDER, NOPL, and AOP. ISLANDER defines law based on action, which can prohibit or forbid the action that the Agent executes. All rules of ISLANDER cannot be obeyed. NOPL is a program design language to organization management infrastructure, not to the language for the programmer. Because of the specialty of OMI (violating laws will bring fatal mistakes to platform), NOPL is a simple law program design language which only supports compulsive obligation. That is, the law of NOPL only describes the action which the Agent executes and all the rules cannot be obeyed. Laws defined by status support three kinds of rules based on obligation, permission, and prohibition and provide punishment mechanisms to the Agent violating rules.

2.3. Reactive Agent Model. Reactive Agent can be shown by symbols, which can respond to the changes of external environment. Reactive structure is designed by corresponding action of the assuming Agent. Its complication of action reflects the one of practical environment of the Agent. Structure of the reactive Agent is shown in Figure 1.

Defined environment is the finite set of all discrete and instantaneous status:

$$E = \{e, e^t, \dots\}. \quad (1)$$

Agent has an action set that can complete. These actions can change the status of the environment. The finite set can be shown as follows:

$$Ac = \{\alpha, \alpha^t, \dots\}. \quad (2)$$

R which is an action in some environment of Agent is a sequence between environment status and action alternate replacement:

$$R : e_0 \xrightarrow{\alpha_0} e_1 \xrightarrow{\alpha_1} \dots \quad (3)$$

The action of reactive Agent can be shown as follows:

$$Ag : E \xrightarrow{see} Per \xrightarrow{plan} Per^* \xrightarrow{action} Ac. \quad (4)$$

In equation (4), see, plan, and action stand for environment perceptor, restriction condition, and corresponding action sequence, respectively. Per and Per* stand for the prior sensing range and anticipation set, respectively.

3. Action Design of Reactive Agent in Cloud Testing

3.1. General Frame Structure of Layer Reactive Agent. The structure of layer agent has higher management agent and lower cloud search agent. These two agents fit to reactive agent model, shown in Figure 2.

The function of Search Agent is to sense part of the image environment and then respond to the content of sense. The Search Agent can search and mark cloud scope by its own intelligent action. However, this agent cannot grasp the whole process and information. The management agent can make up the disadvantage of the search agent, which can acquire the whole information of the image and guide and control the search agent.

3.2. Data Definition and Design of the Search Agent. In order to make sure the working condition of the search agent, we should show three measurements of local cloud similarity of the agent. Their influence equals to E of equation (1). Its definition can be shown as (i).

(i) The number of local grey similarities

$$C_{[i,j]_{-region}} = \sum_{\|(i,j)-(k,l)\| \leq R_{(i,j)_{-region}}} \rho(i, j, k, l). \quad (5)$$

Then,

$$\rho(i, j, k, l) = \begin{cases} 1, & \|I(i, j) - I(k, l)\| \leq \delta, \\ 0, & \text{other.} \end{cases} \quad (6)$$

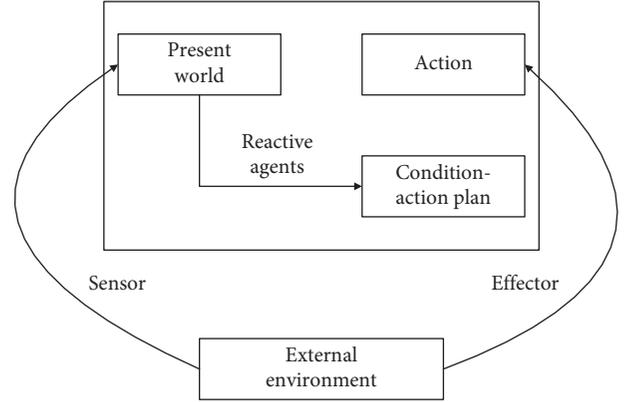


FIGURE 1: Structure of the reactive agent.

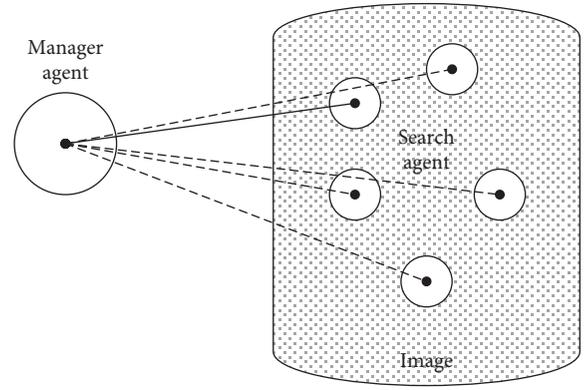


FIGURE 2: Structure of layer agent.

(ii) Expectation of the local area:

$$\text{mean}_{(i,j)_{-region}} = \frac{1}{N} \sum_{\|(i,j)-(k,l)\| \leq R_{(i,j)_{-region}}} I(k, l), \quad (7)$$

$$\begin{aligned} \text{std}_{(i,j)_{-region}} &= \sqrt{\frac{1}{N} \sum_{\|(i,j)-(k,l)\| \leq R_{(i,j)_{-region}}} (I(k, l) - \text{mean}_{(i,j)_{-region}})^2}. \end{aligned} \quad (8)$$

In equation (8), $I(i, j)$ is the grey value of (i, j) . In equation (6), δ is the empirical parameter and $R_{(i,j)_{-region}}$ in Figure 3 is the action radius.

The distributed searching agent directly acts on some pixels of the local image and performs calculation to neighboring environment similarity of local cloud. Cloud-line scope can be searched by locating the working site, expanding scope, moving, marking, and dying out. These actions can be designed to touch off in some condition, which is condition design action of the searching agent. Locate working site: the agent can sense and acquire neighboring information and judge whether this environment adapts to retrieve working. When the environment is fit for retrieving, the agent locates the working site. The rule of locating working site is

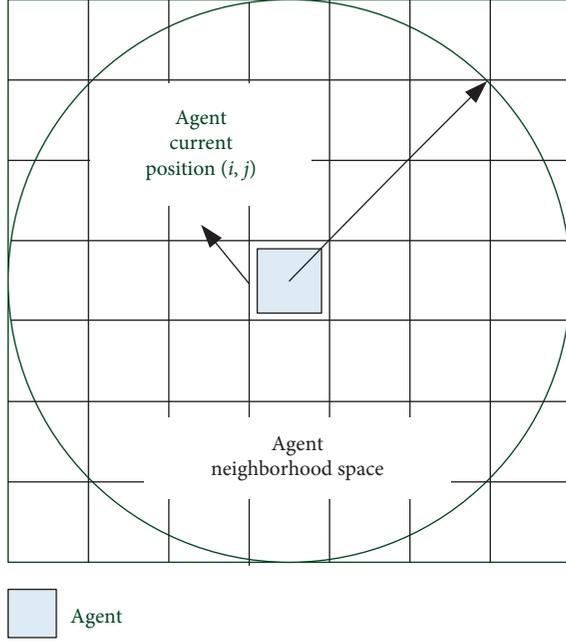


FIGURE 3: Local neighboring region of a shadow-searching agent.

- (i) $R_{(i,j)\text{-region}}$ has characteristics of given images
- (ii) $R_{(i,j)\text{-region}}$ has characteristics of images similar to cloud lines

Rule (i) uses prior knowledge (like the scope of grey value and texture properties) in the environment of practical images in order to guarantee that the agent locates the working site of local images, which can decrease repeat calculation. This paper does not require rule (i) in simulation test algorithm.

Rule (ii) can locate the working site when the agent has pixel of cloud characteristics. Restriction conditions that this algorithm has are as follows:

$$\begin{aligned}
 C_{(i,j)\text{-region}} &\geq C_{\min}, \\
 R_{(i,j)\text{-region}} &\leq M_{\max}, \\
 R_{(i,j)\text{-region}} &\leq S_{\max}.
 \end{aligned} \tag{9}$$

Expanding scope: the searching agent can expand searching space. In this algorithm, the method of expanding searching space is to breed. From the characteristics of cloud, after some searching agent locates working site, we consider its neighboring pixel is the cloud point. The steps by which the searching agent expands scope to breed are as follows. First, we assume the image of A has some agents. A will have 4 offspring after A locates the working site. The prior location of these four offspring is four neighboring points of that agent:

$$A_{(i,j)}^g \longrightarrow \left\{ A_{(i-1,j)}^{(g+1)}, A_{(i+1,j)}^{(g+1)}, A_{(i,j-1)}^{(g+1)}, A_{(i,j+1)}^{(g+1)} \right\}. \tag{10}$$

In this equation, $A_{(i,j)}^{(g)}$ means generation g and working site is the agent of (i, j) .

Moving: the searching agents, which cannot locate working site successfully, need to move in retrieval images to obtain suitable working situation in order to get chances of expanding scope in retrieval work. This is called moving. In the process of moving, the age of them increases gradually. This searching action of the agent can drive the searching agent to find new environment of doing retrieval work.

The specific process of moving in the searching agent can be shown as follows. To the agent that cannot locate the working site and can be defined as $A^{(g+1)}$ to generation $g+1$, we need to search its parent $A^{(g)}$. If its parents do not exist, it shows that agent is the first generation. Then, we can select a direction of moving e and moving distance r randomly. If its parents $A^{(g)}$ exist, we can find $\{B^{(g+1)}\}$ of generation $g+1$ by $A^{(g)}$. We accumulate these points to locate moving direction of agents that have located working site successfully and make up a histogram of counting directions. At last, the ratio between the value of every moving direction and its sum is probability. There is a moving direction e that is made randomly and (e, r) of moving direction r distributed in $[-R, R]$. R is the maximum moving distance, and a new position of $A^{(g+1)}$ in (e, r) appears. This moving mechanism can use family to transfer experience to locate working site of searching site from the aspect of possibility.

Marking: after some searching agent locates working site of a pixel successfully or unsuccessfully, this agent records its result information on the image. If this environment adapts to do retrieval work, then $\text{label}(i, j) = \text{label} A$ or $\text{label}(i, j) = \text{label} B$. $\text{label} A$ and $\text{label} B$ are the given parameters, and label is the marking image.

Dying out: the searching agent will die out after experiencing expanding scope and exceeding maximum of moving. There are two situations in dying-out status:

- (i) After some agent marks and breeds successfully, it can finish algorithm and enter dead state, that is, dying-out status.
- (ii) When some agent cannot still locate the working site of meeting conditions after moving N times, the age of agent is N . If N is greater than maximum age set by the system, then dying-out status appears instantly. Extinction mechanism can eliminate the agent with weak ability and drive population to find new environment suitable to do retrieval work, which avoids optimal solution in local part.

3.3. Management Agent Design. The main task of management agent is to judge whether the searching action of the searching agent can meet requirement. This paper defines the best measure of the whole searching as equation (11) according to the characteristics of cloud:

$$R^{(g)} = \sqrt{\frac{1}{N} \sum_{L(i,j)=\text{label}A} (I(i, j) - M^{(g)})^2}. \tag{11}$$

Then,

$$M^{(g)} = \frac{1}{N} \sum_{L(i,j)=\text{label}A} I(i, j). \quad (12)$$

$I(i, j)$ means the grey value of this image in (i, j) . $\text{label}(i, j)$ is the grey value of the marked image in (i, j) . $\text{label} A$ is the marking value when the searching agent meets the conditions of the working site. $R^{(g)}$ is the best measure of the searching agent of generation g .

The management agent can get $R^{(g)}$ after every iteration according to best measure. If $R^{(g)} > R_M$, it means the current working condition cannot meet the best condition and the working condition needs to be adjusted.

The formulas locating working site include formulas (5)–(7), in which there are three parameters $\{C_{\min}, M_{\max}, S_{\max}\}$. Here, $\{C_{\min}, S_{\max}\}$ often are unchangeable to different sky-sea images. $\{M_{\max}\}$ is related to the brightness of scene and imaging condition. This algorithm is to adjust to $\{M_{\max}\}$, and the process of adjustment can be shown as equation (13).

$K_{\text{threshold}}$ is the threshold of image segmentation when we cluster two kinds of grey K – mean to image $R = \text{label} \& I$.

The existence of the management agent can claim the tolerance of the lower-layer cloud test agent can adjust itself to the environment and make the algorithm suitable to image in some scene.

4. Fulfillment of Algorithm

4.1. Acquisition of Information of Management Agent. When the multiagent handles the similar events, the prior self-testing often executes when we find $|V_m| < (N + 1)/2$ at the end of the task (we call it self-diagnosis). However, results that we get in the recognition decision process described in this paper are instantaneous result. If self-testing can be done after one task period is over, the results that we get in the process no longer exist. In order to know the events happening during the task period, we put forward a method. The corresponding detection code can be inserted in the original task codes dispersedly, and environment information can be extracted real-timely in task operation process. If the original task code sequence is $\{t_1, t_2, \dots, t_n\}$ and detection code sequence is $\{d_1, d_2, \dots, d_m\}$, we get $\{t_1, d_1, t_2, d_2, \dots, d_m, t_n\}$ when $\{d_1, d_2, \dots, d_m\}$ is inserted into $\{t_1, t_2, \dots, t_n\}$ distributedly. If insertion methods are different, sequences are different. Detection code sequence $\{d_1, d_2, \dots, d_m\}$ designed reasonably is the key to voting algorithm. We design a detection sequence and take it as principle verification. This detection sequence is mainly used as an intermediate result that Agent outputs. Detection sequence is as in Algorithm 1.

If recognition result of every Agent in operation process is false (0), Detect Result should be false. If some Agents can recognize right result in some task period q , we should judge the possibility of results according to most voting principles.

When intermediate result of function Agent is greater than one of the common Agents and results are considered real, then the results are real. This voting function can output specific values of results. If three agents consider results real,

the output voting result is 3. The following processing can give further conclusion according to this result. In the system in which reliability requirement of identification recognition is strict, threshold can be improved. Conversely, in the system in which reliability requirement is not strict, this threshold can be decreased appropriately.

4.2. Image Segmentation. In the method of image segmentation based on the agent, the agent points judge whether it meets the consistency criteria by sensing the pixels in its neighborhood. The agent point performs its subsequent behavior through the feedback information of its neighboring regions; it may reproduce offspring and moves to neighboring pixels or from the image disappear.

Usually, we use the following three mathematical criteria to measure whether the consistency region is satisfied: relative contrast, regional average, and regional gray standard variance. Agent point of behavior will be based on the above 3 criteria to determine whether to trigger local stimulation; more detailed consistency criteria are defined as follows:

Definition of contrast formula:

$$g(i, j) = \sum_{\|(i,j)-(k,l)\| \leq R(i,j)} \rho(i, j, k, l), \quad (13)$$

$$\rho(i, j, k, l) = \begin{cases} 1, & \text{if } \|I(i, j) - I(k, l)\| \leq \delta, \\ -1, & \text{others,} \end{cases} \quad (14)$$

where, $g(i, j) \in (\eta_1, \eta_2)$ is a predefined constant; $R(i, j)$ is the neighborhood radius of pixel (i, j) for agent; $I(i, j)$ is the grey scale of pixel (i, j) ; δ is the threshold value of predefinition.

Mean of standard region:

$$\text{mean}_{(i,j)} = \frac{1}{N} \sum_{\|(i,j)-(k,l)\| \leq R(i,j)} I(k, l). \quad (15)$$

Standard region variance:

$$\delta_{\text{stand}(i,j)} = \sqrt{\frac{1}{N} \sum_{\|(i,j)-(k,l)\| \leq R(i,j)} (I(k, l) - \text{mean}_{(i,j)})^2} \quad (16)$$

where $\text{mean}_{(i,j)} \in (\mu_1, \mu_2)$, while μ_1, μ_2 is the constant of predefinition; $\delta_{\text{stand}(i,j)} \in (d_1, d_2)$, d_1, d_2 is the constant of predefinition; N is the number of pixels for region $R(i, j)$.

4.3. Replication and Diffusion. Agent can adopt two different modes of behavior, replication and diffusion, corresponding to different local environments, which become the adaptability of agent. Diffusion is to move up the current pixel to the other ones. Specific process is as follows:

- (1) When an agent searches a pixel satisfying the above three criteria consistency, it will copy the appropriate set of descendant agent in a particular direction on the neighborhood. The copy behavior makes the newly generated offspring locate near the pixels

```

<DetectCode> → <Init>; <DetectPair>*;
<Finish>;
<Init> → DetectI = 0; DetectJ = 0;
<DetectPair> → DetectI = DetectI + 1.0;
DetectJ = DetectJ + 1.0;
<Finish> → DetectResult = DetectI - DetectJ.
(*) means <DetectPair> can be repeated many many times. The typical period task can be described as follows after inserting
detection code:
// states the global variable that detection codes can use out function.
while (1) // period task, repeat operation repeatedly
{
Init // test initialization, test variable zeroes
t1; // the beginning a period, like sensor information
// every task code, insert a detect pair:
Detect Pair
t2;
t3;
Detect Pair
.....
Detect Pair
tn; // when algorithm finishes, voting is prepared
Finish // standardize detect results
r = vote(r, Detect Result); // vote with self-testing
output(r); // output control information
}

```

ALGORITHM 1

which can meet the consistency criteria, for subsequent detection of further regional coherence:

$$\alpha_{(i,j)}^g \Rightarrow \left\{ \alpha_{\nu(w, d_{\text{copy}})} \right\}^{\nu} \quad (17)$$

$$= 1, 2, \dots, m; \omega \in \Omega; d_{\text{copy}} \leq k \}^{g+1}.$$

where g and $g + 1$ is the g and $g + 1$ offspring of the agent, (i, j) is the current position of the agent, α is an agent with active status, ν is the ν offspring of the agent, m is the total amount of the offspring agent, ω is the replication direction of the offspring agent, Ω is the replication direction of a series of possibilities, d_{copy} represents the distance of replicating the offspring agent, and R_c is the replicating radius.

- (2) When an agent finds itself is in a nonuniform region, it will select diffusion pattern, moving along a particular direction to a specific location. Diffusion behavior also plays an important role in the discovery process of consistency area because the diffusion direction is determined by the agents that successfully found consistency area among the parent agent and the agent on behalf of the brothers. The new diffusion agent is still in the neighborhood; its proliferation is not a search for balance but should be seen as looking for a new biased search of consistency pixels:

$$\alpha_{(i,j)}^t \Rightarrow \alpha_{\theta, d_{\text{diffu}}}^{t+1}, \quad \theta \in \Theta, d_{\text{diffu}} \leq R_d. \quad (18)$$

where t and $t + 1$ represent the time of the agent, θ is the spreading direction, Θ is the spreading direction of a series of possibility, d_{diffu} represents the spreading distance, R_d denotes the spreading radius.

- (3) When the agent found a consistency area, which itself is placed in the suppressed state.

Meanwhile, in order to prevent the agent's unlimited searching, its provisions can be a "life cycle," over the life cycle and it is self-suppression.

4.4. Steps of Algorithm

Step one: build an Agent group $\{A_i\}_{i+1}^k$, in which is given the parameter that shows the number of agents. We take the age as 1 and put it to active queue of agent.

Step two: to every agent of active queue, we can make sure whether that agent can locate the working site according to the working condition. If the working site can be located, all corresponding pixel points of that agent can be set as label A , and then the offspring are propagated and added to the active queue of the agent. Meanwhile, it itself will be deleted from the active queue of agent. If the working site cannot be located, corresponding pixel points of that agent can be set as label B first and then we should judge whether its age is greater than the given maximum age. If it is greater, it will be deleted from active queue of agent or it will move. The new coordinate point can be built, and its age will be added to 1.

Step three: Through equation (11), we can get best measure of the current marked image. If measure is greater than the preinstalled parameter, we can correct $\{M_{max}\}$ in the working condition by equation (15) and return to step one. If measure is smaller than the preinstalled parameter, we will begin step four.

Step four: Judge whether active queue of agent is empty or all images are marked. If it is, we can begin step five or return to step two.

Step five: To do subsequent processing of processed images, we can make grey value label A of the original image 1, and other grey values which are not equal to label A are made 0. In this way, we can get an image with two values. Then, we perform dilation to this image and perform dry operation, and we can get target images.

5. Experiment and Result Analysis

In order to test the correctness of this method, we select a group of typical sky-sea pictures as the simulation testing image, including four types of targets. Figure 4 shows a common cloud image including only sea and sky. Figure 4 also reveals a gradual process of removing cloud through this algorithm. The process of operation takes 5.22 seconds.

The parameters that this image uses are as follows: the number of initial agent is 500, the maximum existence age is 5, neighboring search diameter is 4, and the minimum value of similar points of local grey is 3. The maximum initial value in the local area is the grey expectation value of the input image, and the maximum value of local standard deviation is 9.

Simulation platform: simulation software is Matlab 2010a, CPU domain frequency is 3.2 G, and memory is DDR2GB. Figures 5 and 6 are the typical sky-sea image and corresponding cloud test result. Figure 5 is the tilt image for camera angle and the result of testing cloud through that method. Figure 6 is the sky-sea image with occlusion like trees and the result of testing cloud through this method. Furthermore, we select another 9 images to every type of target to do test simulation. In the table of simulation test result, we compare cloud testing result through this method with cloud testing through artificial vision and we perform statistics to ratio in which detection and missed test scopes account for the whole cloud scope, like Table 1.

Through observing simulation results, we find this method can search and remove cloud to different sky-sea images and cannot limit size and shape of cloud in searching. The test result of images with tilt angle and images with occlusions shows that this method can overcome the influence of cloud recognition to complicated scenes. The total operation time of computer is so little, which means time complicity of this algorithm is low. The operation time of a single infrared image (256×256 pixel) can be controlled between 1 and 3 seconds. Therefore, this method can be used widely.



(a)



(b)



(c)

FIGURE 4: Cloud extraction and removal of common images.

6. Conclusion

We can test and remove cloud of sky-sea images, which can improve recognition efficiency of sky-sea line and acquire tilt angle of carriers. This paper puts forward a better algorithm method of testing and removing cloud of infrared images based on cloud characteristics of analyzing infrared sky-sea images. Through grey value characteristics of cloud, we can find and use reactive agent layer structure, classify many agents used for local image cloud searching, and manage agents used for coordinating many cloud searching agents. We can extract and remove cloud of the infrared sky-sea image by intelligent algorithm and cooperation of two agents. Simulation shows this method has good operation effect to various shooting situations. And, time

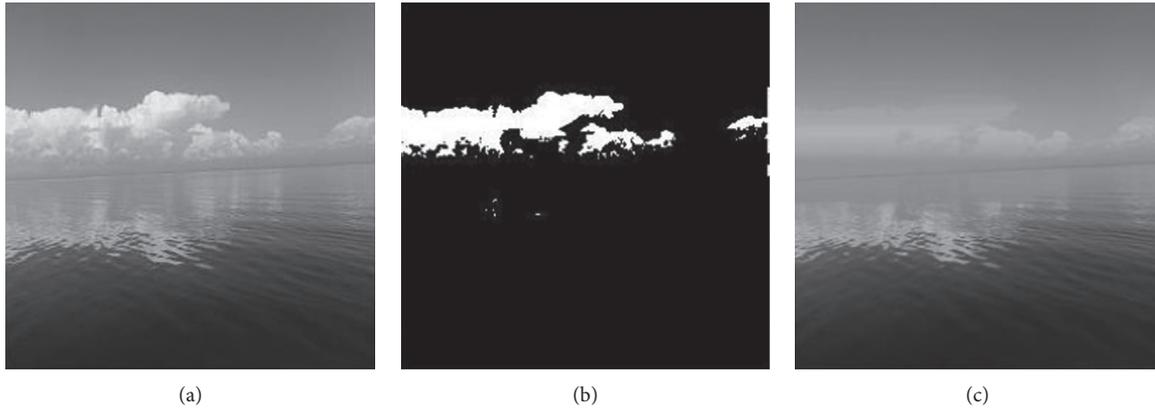


FIGURE 5: Cloud extraction and removal of the sky-sea image with tilt angle.

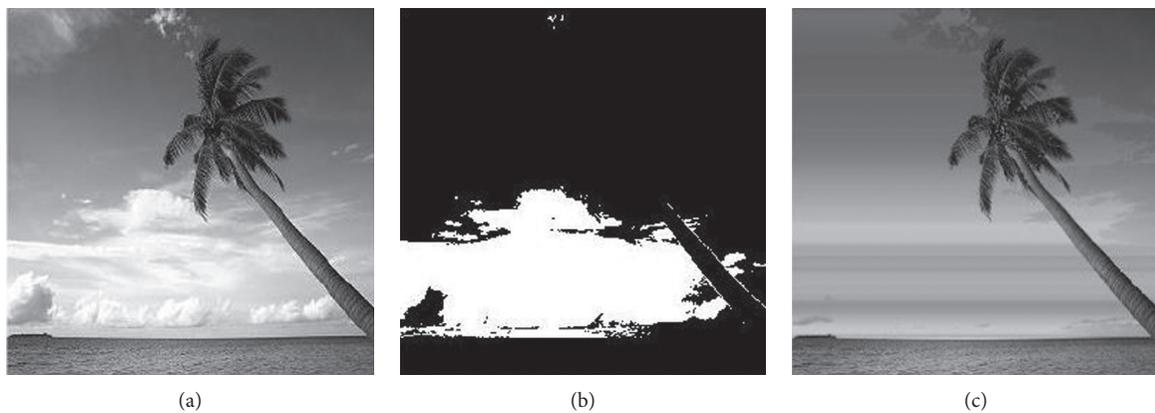


FIGURE 6: Cloud extraction and removal of the sky-sea image with occlusions.

TABLE 1: Experiment result.

Types of image	Total after-detection rate (%)	Total missed detection rate (%)	Average operation time (s)
Common sky-sea image	2.71	3.25	1.56
Image with tilt angle	3.75	4.80	2.12
Image with occlusions	7.53	5.01	2.62

complicity of this algorithm is low, and cloud test and removal task in the infrared sky-sea image can be completed fast and efficiently.

Data Availability

The data used to support the findings of this study are available from the corresponding author upon request.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

This research work was supported by the National Natural Science Foundation of China (51476049) and Harbin

University of Commerce Young Creative Talents Support Project (2016QN050 and 17XN005).

References

- [1] J. J. Yoon, C. Koch, and T. J. Ellis, "ShadowFlash: an approach for shadow removal in an active illumination environment," in *Proceedings of the 13rd British Machine Vision Conference*, pp. 636–645, INSPEC, Cardiff, UK, 2002.
- [2] L. Yan and P. Gong, "Cloud test of images based on DSM cloud simulation and ray tracing in high area," *Journal of Remote Sensing*, vol. 9, no. 4, pp. 357–362, 2005.
- [3] P. Gils, "Remote sensing and cast shadows in mountainous terrain," *Photogrammetric Engineering and Remote Sensing*, vol. 67, no. 7, pp. 833–839, 2001.
- [4] H. Guo, Q. Xu, and B. Zhang, "Building cloud extraction in the multiple constraint," *Wuhan University Journal*, vol. 30, no. 12, pp. 1059–1062, 2005.

- [5] Y. Yang, R. Zhao, and W. Wang, "Test of cloud area in air image," *Signal Processing*, vol. 18, no. 3, pp. 228–232, 2002.
- [6] J. Liu and Y. Y. Tang, "Adaptive image segmentation with distributed behavior-based agents," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 21, no. 6, pp. 544–551, 1999.
- [7] H. He and Y. Q. Chen, "Artificial life for image segmentation," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 15, no. 6, pp. 989–1003, 2001.
- [8] C. Pereira, D. Veiga, J. Mahdjoub et al., "Using a multi-agent system approach for microaneurysm detection in fundus images," *Artificial Intelligence in Medicine*, vol. 60, no. 3, pp. 179–188, 2014.
- [9] Y. Fan, L. Liu, G. Feng, and Y. Wang, "Self-triggered consensus for multi-agent systems with zeno-free triggers," *IEEE Transactions on Automatic Control*, vol. 60, no. 10, pp. 2779–2784, 2015.
- [10] R. H. Baxter, N. M. Robertson, and D. M. Lane, "Human behaviour recognition in data-scarce domains," *Pattern Recognition*, vol. 48, no. 8, pp. 2377–2393, 2015.
- [11] N. C. A. D. Freitas, P. P. R. Filho, C. D. G. D. Moura, and M. P. D. S. Silva, "AgentGeo: multi-agent system of satellite images mining," *IEEE Latin America Transactions*, vol. 14, no. 3, pp. 1343–1351, 2016.
- [12] J.-R. Ruiz-Sarmiento, C. Galindo, and J. Gonzalez-Jimenez, "Scene object recognition for mobile robots through semantic knowledge and probabilistic graphical models," *Expert Systems with Applications*, vol. 42, no. 22, pp. 8805–8816, 2015.
- [13] N. Shiroma, R. Miyauchi, A. Nagafusa, Y. Haga, and F. Matsuno, "Gaze direction based vehicle teleoperation method with omnidirectional image stabilization and automatic body rotation control," *Advanced Robotics*, vol. 29, no. 3, pp. 149–163, 2015.
- [14] D. Wang, L. Liu, X. Wang, and Y. Lu, "A novel feature extraction method on activity recognition using smartphone," *Web-Age Information Management*, vol. 351, pp. 67–76, 2016.
- [15] L. Zhuo, Z. Geng, J. Zhang, and X. G. Li, "ORB feature based web pornographic image recognition," *Neurocomputing*, vol. 173, pp. 511–517, 2015.
- [16] J. M. Beer, C.-A. Smarr, A. D. Fisk, and W. A. Rogers, "Younger and older users' recognition of virtual agent facial expressions," *International Journal of Human-Computer Studies*, vol. 75, pp. 1–20, 2015.

Research Article

Consumers Team Detection Model Based on Trust for Multi-Level

Xiaoming Li,^{1,2} Guangquan Xu ,¹ Sandhya Armoogum ,³ and Honghao Gao ⁴

¹School of Computer Science and Technology, Tianjin University, Tianjin, China

²School of Information Science and Engineering, Zaozhuang University, Shandong, China

³School of Innovative Technologies and Engineering, University of Technology Mauritius, Pointe-Aux-Sables, Mauritius

⁴Computing Center, Shanghai University, Shanghai, China

Correspondence should be addressed to Honghao Gao; gaohonghao@shu.edu.cn

Received 20 July 2018; Revised 11 November 2018; Accepted 5 December 2018; Published 3 February 2019

Guest Editor: Jaegeol Yim

Copyright © 2019 XiaoMing Li et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Due to rapid advances in technology, social networks have become important platforms for daily communication, product marketing, and information dissemination. Targeted delivery of social network advertisement can considerably improve the efficacy of the advertisement and maximize the profits from it. In this context, managing the specific audience of a social network advertisement and achieving targeted advertisement delivery have been the ultimate goals of the social network advertising sector. Identifying user groups with similar properties is critical to increasing targeted sales. When both the scale of mobile social network and the complexity of social network user behaviors grow, similar groups are hidden in user behaviors. In order to analyze community structure with user trust relationship more appropriately in the large-scale multilevel social network environment, a novel local community detection model E-MLCD is proposed in this paper. It is jointly based on the multilevel properties and the strength of similarity of multilevel social interaction among communities. By studying three real-world multilevel social networks and specific QQ Zone marketing data, the model defines a new metric of community trust based on similarity. Comparison between other state-of-the-art detection methods demonstrate E-MLCD's ability to detect communities more effectively.

1. Introduction

Due to the growing diversity of customer requirements and the rapid advances in mobile social network, how to identify exact requirements of customers by distinguishing between customer groups is an important aspect of core competitiveness for enterprises. Identifying specific customer groups from the large-scale mobile social network refers to the issue of community detection.

Since Newman et al. proposed the Girvan–Newman (GN) algorithm in 2002 [1], a lot of attention has been paid across the globe to community detection. In recent years, community detection and community detection algorithms have become a focus of research on complex networks [2–5]. However, the multifaceted correlation between entity and property means that many complex systems are interrelated rather than independent. Consequently, the traditional single-layer network model cannot describe the system very accurately. In this context, a new type called multilevel complex network has not only emerged an extension of

existing network model but also represented a breakthrough in the entire network theory. In the multilayer complex network, the network structure is not completely flat. Instead, interlayer interaction is introduced. As a consequence, most of the traditional single-layer network evaluation methods are no longer suited for multilayer network [6, 7]. Figure 1 shows that the different consumer groups are found on multilayer networks according to geographical location.

The rapid progress in mobile network and electronic communications as well as the growing scale of network, the social networks (e.g., Facebook, Twitter, QQ, Zhihu, and Sina Weibo) have tens of millions of nodes and links [8, 9], making it very difficult to collect data across the entire network. Displaying advertisements and recommending products and friends to users can be more effective if specific groups identified through community detection. However, considering the colossal scale of these social networks, even if the parallel distributed platform is used to analyse the entire networks, the space and time consumption is too high to be acceptable. To make matters even worse, many

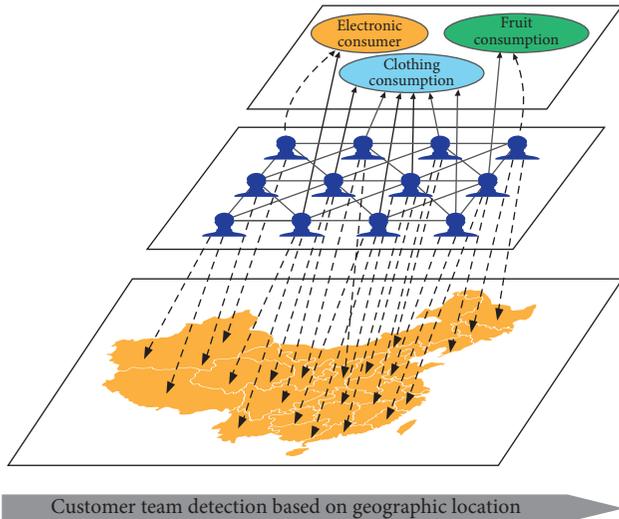


FIGURE 1: Different consumer groups are found on multilayer networks according to the geographical location.

real-world social networks have a stratified structure, which cannot be detected globally without observing network community at different scales. How to detect community very accurately and efficiently using all information tables obtained from the large-scale diverse social networks is a significant challenge. In this context, some works began identifying certain nodes or a local community of several nodes quickly and accurately, given a limited amount of network information. In this way, the prohibitive space and time overhead associated with global calculation can be avoided.

Each user of the social networks has different characteristics, such as the families, circle of friends and classmates, hobbies, and interests. These characteristics are called node trust relationship. The network based on attribute similarity between nodes is called correlation network. The network with node or edge trust relationship is called the attribute network.

In this paper, we jointly consider the trust relationship of all network nodes and the strengths of social similarity. A new metric is defined to describe node attribute and social similarity strength. An algorithm is proposed to detect local community through effective use of the new metric. Experiments are performed to test the proposed algorithm. Unlike the traditional methods which cluster nodes based on topological structure, the proposed algorithm determines node similarity using certain trust relationship of nodes in real life. The relationship between these nodes is then used to obtain the integrated level of similarity. The representative seed node of the community is identified. The community correlated with the attribute is merged by analyzing the similarity between the seed and the community. The strength of social similarity between the attribute and the boundary node is used as a criterion to terminate community expansion. Finally, a scheme based on attribute relationship and social similarity strength is proposed to detect local community within a multilevel network.

The remainder of this paper is organized as follows: in Section 2, the fundamental theory on the detection of

community within multilevel social networks is reviewed. In Section 3, our trust-based method for detection of local community from multilevel network is proposed. In Section 4, the proposed method is evaluated and compared with three other models in real-world multilevel networks and specific QQ Zone marketing data. In Section 5, the conclusion is presented and the future work is recommended.

2. Related Work

The aim of community detection is to segment social network, gain more understanding of social network, identify community members, and find groups with similar ideas, motivations, shopping preferences, and other common characteristics.

2.1. Local Community Detection. The social network is increasingly large, and it contains more information. The method based on local information detects community from a local perspective. By eliminating the need for a global analysis of the network structure, this type of method has become more popular in recent years [10]. In [11], Bagrow and Bollt proposed to start with source node and keep adding continuous shells as the node. In [12], Lancichinetti et al. proposed the F metric for the fitness function to measure the connectivity difference between the inside and outside of community. Their method is easy to implement, but the random selection of the original node usually causes instability in community detection. Moreover, certain parameters of the fitness metric have to be determined in advance. In [13], Chen et al. proposed a method to detect local community using the local-degree central node. In their method, the local community is not identified from the given start node. Instead, it is identified from the local central node correlated with the given start node. In [14], Tabarзад and Hamzeh proposed a heuristic method to detect community by investigating local information. Comparison with other state-of-the-art algorithms demonstrated the ability of their method to detect community and member more effectively and accurately. In [15], Chen et al. proposed a metric called semilocal centrality to find a balance between centrality with a low level of correlation and other time-consuming metrics. In [16, 17], the authors proposed a multiagent algorithm for autonomous community detection from the distributed environment.

The members of a social community usually have multilevel relationships. Most of the traditional community detection algorithms are based on the information of network structure. User behavioral trust relationship is not taken into account in the local community detection methods described above.

2.2. Multilevel Network Community Detection. Detecting community within a multilevel network has drawn a lot of attention in recent years [18–21]. Due to the considerable complexity of the real world, the single-level network is no longer able to describe the community very effectively. In [22], Berlingerio proposed a multilevel network model to

analyze the complex system of the real world and defined the multilevel network relationships. In [23], Gayel et al. proposed a general framework of network quality function, which allows the community in any multilevel network to be studied. In this framework, the network is a combination of coupled links which connects each node of a network slice with each node of other slices. This framework enables us to study the community structure of as many slice networks as we want. In [18], Domenico revealed that any complex system can be represented by a multilevel network. For example, organism genes and the interaction between them can be represented by 7 network layers. In [24, 25], the authors proposed a community detection algorithm based on multilevel network modularity. The concept of modularity was introduced to a broad variety of dynamic and connected networks. In [26, 27], a community detection algorithm based on multilevel network clustering was presented, where the original multilevel networks were first merged into a single-level network under certain strategy before the community detection method for the single-layer network was implemented. In [28–32], a community detection method based on consensus clustering was described. The community detection algorithm for the single-layer network was first implemented on each layer of network. The detection result was then converted into a characteristic matrix of nodes, which was subsequently processed using the traditional clustering algorithm.

The methods described above focus on community detection from a global perspective. It is a new challenge to detect community from the complex and large-scale real-world network systems.

2.3. Similarity Trust. In their paper published in 1998, Buskens [33] derived the assumption of the influence of density, degree centrality, and centralization on the trust degree of client, and according to the test, higher trust can be generated with higher density and higher degree. Sherchan et al. [34] pointed out that the network structure may affect the trust level of social network, and in the network, high density and high interaction between members can generate high-level trust. By utilizing the social network structure and its dynamics, Trifunovic et al. [35] proposed two supplementary methods to build social trust: dominant social trust and implicit social trust. In their article on transactions, Wang et al. [36] obtained the similarity between the same group of neighbours based on the interest of one pair of partners, and they used the two points P_i and P_j to represent two neighbours and define the Jaccard similarity trust relation between two nodes through Jaccard measurement. Jin et al. [37] proposed a trust model based on the evaluation of group similarity; this model can compute the global trust value with similarity between groups as the recommendation reliability, the local trust value, and global trust value can be provided by this model based on the recommended reliability, and they evaluated the model from the perspective of security. In their research, Ziegler and Lausen [38] used trust to supplement or even replace the filtering mechanism, but in order to obtain meaningful result, they

assumed that trust could reflect the user similarity in a certain degree, and they obtained the following conclusion: that is, when the trust network of community is closely integrated with certain application program, there is correlation between trust and user similarity. Boratto et al. [39] proposed the setting to realize trust parameter by computing the average value of top similarity; this parameter represents the minimum value in which the average value of similarity must reach to be recognized as reliable, and the trust parameter obtained from prediction rating must be smaller than 0.85. Among most current works, in [40, 41], the trust degree is generally represented by the user dominance; however, in most online social networks (such as Facebook, Epinions, and Flixster) of real world, there is no specific value to measure the degree of user trust relation, so their paper defined the direct similarity relation between them based on the cosine similarity, and direct trust is represented through similarity relation. Applying existing algorithms to local community detection within a multilevel network has many limitations. For example, global information is not easily available; the node attribute cannot be selected appropriately and utilized effectively; and the model is not as accurate and stable as expected. To address these problems, we propose a model for local community detection within a multilevel network based on trust.

3. Community Detection from Local Network

In this section, we first describe the multilevel network graph model. Next, we propose the E-MLCD model for local community detection from multilevel networks. Finally, the algorithm framework is presented based on the input seed node.

3.1. Problem Formulation. In real life, there is a frequent need to recommend users with appropriate products by identifying their preferences. To achieve this, a core user interested in a certain product is determined and then expanded to find more users who are also interested in the same product. This process is referred to as local community detection, i.e., obtaining a community structure centred on one or more seed users, given limited information on the network. Obviously, the local community detection algorithm is able to effectively identify a community of interest by accessing and manipulating a relatively small part of the network.

After reviewing existing community detection methods, it is learned that no work has been done in the past to jointly consider node similarity strength and node attribute during community detection. The limitations of existing methods can be summarized as follows: (1) detecting community within a multilevel network is very sensitive to the location of start node; (2) information on structural similarity between communities and the topological information are underused; and (3) the expansion direction and breadth within the multilevel network cannot be controlled accurately. In existing methods, the community is expanded outward at a step length of one node. Effective guidance on

expansion is lacking. Therefore, the start node location and community similarity strength should be fully taken into account. The attribute difference between users in different network layers can be used to formulate a new strategy for local community detection within a multilevel network.

3.2. Local Seed Node. In the social network, there is a broad variety of trust, such as age, gender, residence place, shopping preference, and behavioral description. Each of these trust relationship can be seen as an attribute layer. The attribute of the core node in each layer is critical to the node feature description. Therefore, in addition to structure-based social similarity, similarity between nodes in the same layer, which arises from node attribute similarity, should also be taken into account.

Attribute-based node classification can be achieved by classifying nodes with the same attribute into the same category. During local community detection of a network layer, the attribute similarity between node and its neighbouring system in the corresponding structure is defined as the criterion for selecting the seed node. The set of all seeds fulfilling the criterion is determined and then expanded to provide a community.

Let $G_{\mathcal{L}} = (V_{\mathcal{L}}, E_{\mathcal{L}}, V, \mathcal{L})$ denotes the trust-based network model, where $V_{\mathcal{L}}$ denotes the set of nodes in layer \mathcal{L} , V is, $\{V_1, V_2, \dots, V_n\}$ N denotes the number of nodes in the network, and $E_{\mathcal{L}}$ denotes the edges connecting nodes in layer \mathcal{L} .

As the seed node algorithm of the single-layer network cannot be applied to the multilayer network, we describe the core node in the network by defining the local centre degree of the multilayer network. The larger the degree of the node indicates that the more the node is located in the centre of the network, the more important it is in the network and the more it can affect other nodes of the network.

Consider a multitier network with an M layer and N nodes per layer $u = (y, t)$. The node $i \in X (i = 1, \dots, N)$ is expressed as the connection degree or the vector of degree:

$$k_i = (k_i^{[1]}, \dots, k_i^{[M]}), \quad (1)$$

where $k_i^{[1]}$ be the degree of node i in the a layer α . The degree of $k_i^{[\alpha]} = \sum \alpha_j^{[\alpha]}$ is

$$D_i = \sum_{\alpha=1}^M k_i^{[\alpha]}. \quad (2)$$

3.3. Measurement of Intralayer Similarity Trust. During the detection of local community in the multilayer network, it is very important to control the community expansion of network based on the intralayer trust relation of multilayer network and the interlayer trust relation. In this section, we only consider the similarity trust relation when nodes are on the same layer. In order to better express the similarity trust between nodes on the same layer, in this section, the cosine similarity is used to represent the trust relation between nodes on the same layer of network. For an \mathcal{L} -layer network

$G_{\mathcal{L}} = (V_{\mathcal{L}}, E_{\mathcal{L}}, \mathcal{V}, \mathcal{L})$ consisting of the different network layer set \mathcal{L} and node set V , the intralayer trust relation of multilayer network has the following definition.

Definition 1 (intralayer similarity trust relation of multilayer network). For two nodes with cosine similarity on the same layer of multilayer network, according to literatures [42, 43], it is easier for similar nodes to form a community. Obtain the trust measurement of candidate node and target node on the same layer, and the computation method for similarity trust on the same later can be expressed with the following formula:

$$\begin{aligned} T_{\text{int}}^i(C_u, S_v) &= \sum_{u \in C} \sum_{v \in S} \text{sim}_i(C_u, S_v) \\ &= \sum_{u \in C} \sum_{v \in S_v} \frac{|N_i(C_u) \cap N_i(S_v)|}{\sqrt{|N_i(C_u)| |N_i(S_v)|}} \end{aligned} \quad (3)$$

$$\begin{aligned} T_{\text{int}}^i(S_u, S_v) &= \sum_{u \in S} \sum_{v \in S} \text{sim}_i(S_u, S_v) \\ &= \sum_{u \in S} \sum_{v \in S_v} \frac{|N_i(S_u) \cap N_i(S_v)|}{\sqrt{|N_i(S_u)| |N_i(S_v)|}} \end{aligned} \quad (4)$$

In formula (3), $T_{\text{int}}^i(C_u, S_v)$ refers to the trust measurement of node u in local community and node v on the same layer of neighbourhood, in which, C_u is the node in local community C and S_v is the node in neighbourhood set S connected to local community C . $\text{sim}_i(C_u, S_v)$ represents the cosine similarity of two nodes; $N_i(C_u)$ is the neighbour set of node C_u on layer L_i and $N_i(S_v)$ is the neighbour set of node v in set S . In formula (4), $T_{\text{int}}^i(C_u, S_v)$ refers to the trust measurement between nodes S_u, S_v in neighbourhood set S ; $\text{sim}_i(S_u, S_v)$ is the cosine similarity of nodes S_u, S_v in set S .

Because there is no extraattribute, the value of overall trust relation on the same layer is the higher value between $T_{\text{int}}^i(C_u, S_v)$ and $T_{\text{int}}(S_u, S_v)$:

$$T_{\text{int}}(u, v) = \max(T_{\text{int}}^i(C_u, S_v), T_{\text{int}}^i(S_u, S_v)). \quad (5)$$

3.4. Measurement of Interlayer Similarity Trust. In the last section, the similarity trust between nodes on the same layer of multilayer network was defined. However, the nodes on different layers of multilayer network also have different trust relations, and we consider the similarity trust between nodes on different layers as beneficial supplementation to the community detection of the multilayer network.

Definition 2 (interlayer similarity trust relation of multilayer network). For two nodes with cosine similarity on different layers of multilayer network, according to literatures [42, 43], it is easier for similar nodes to form a community. The trust measurement of candidate node and target node on different layers is obtained, and the computation method for interlayer similarity trust can be expressed with the following formula.

According to formula (3), obtain the similarity between nodes u and v on layer L_i and layer L_j . Similar to previous

definition, $\text{sim}_{i,j}(u, v)$ corresponds to the topologic similarity measurement between adjacent sets of u and v , but in this example, they are on different layers:

$$\begin{aligned} T_{\text{out}}^{i,j}(C_u, S_v) &= \sum_{u \in C} \sum_{v \in S} \text{sim}_{i,j}(C_u, S_v) \\ &= \sum_{u \in C} \sum_{v \in S} \frac{|N_i(C_u) \cap N_j(S_v)|}{\sqrt{|N_i(C_u)| |N_j(S_v)|}} \end{aligned} \quad (6)$$

$$\begin{aligned} T_{\text{out}}^{i,j}(S_u, S_v) &= \sum_{u \in S} \sum_{v \in S} \text{sim}_{i,j}(S_u, S_v) \\ &= \sum_{u \in S} \sum_{v \in S} \frac{|N_i(S_u) \cap N_j(S_v)|}{\sqrt{|N_i(S_u)| |N_j(S_v)|}} \end{aligned} \quad (7)$$

In formula (6), $T_{\text{out}}^{i,j}(C_u, S_v)$ refers to the trust measurement of node u in local community and node v in the neighbourhood on different layers, in which C_u is the node in local community C and S_v is the node in neighbourhood set S on different layers connected to local community C . $\text{sim}_i(C_u, S_v)$ represents the cosine similarity of two nodes; $N_i(C_u)$ is the neighbour set of node C_u on layer L_i ; $N_j(S_v)$ is the neighbour set of node v in set S on layer L_j . In formula (7), $T_{\text{out}}^{i,j}(S_u, S_v)$ refers to the trust measurement between nodes S_u, S_v in neighbourhood set S ; $\text{sim}_i(S_u, S_v)$ is the cosine similarity of nodes S_u, S_v in set S . According to formula (6), the value of overall trust relation between the target node and the node to be added on different layers is the higher value between $T_{\text{out}}^{i,j}(C_u, S_v)$ and $T_{\text{out}}^{i,j}(S_u, S_v)$:

$$T_{\text{out}}(u, v) = \max(T_{\text{out}}^{i,j}(C_u, S_v), T_{\text{out}}^{i,j}(S_u, S_v)). \quad (8)$$

Therefore, the seed node of the multilayer network is the node with the largest degree, which is expressed as follows:

$$V_{\downarrow 0} = \max(s_{\downarrow} i^{\uparrow} [1], L, s_{\downarrow} i^{\uparrow} [\mathcal{L}]). \quad (9)$$

3.5. Trust-Based Algorithm for Multilevel Local Community Detection. In this subsection, we propose E-MCLD and a trust-based algorithm for multilevel local community detection. It is valid for any graph with two or more layers. The pseudocode of the general method for trust-based multilevel local community detection is given in Algorithm 1.

The convergence of the proposed algorithm guarantees the dominant relationship between individuals and the upper bound of the two objective functions. According to their definitions, there is an upper bound for the multilevel network to be expanded using the seed node. The community to which the node belongs is determined by the level of similarity of this node with two sets of attribute. Therefore, the proposed algorithm is able to converge. Figure 2 shows the algorithm's trust-based community detection.

4. Experimental Results

The performance of E-MLCD is evaluated through extensive experimentation on single- and multilevel real-world

networks. The experiments are performed on the computer with Windows 7, 3.10 GHz, and 32.00 GB RAM.

4.1. Experimental Datasets. As we know, datasets on social networks with explicit behavioral trust relationship are very rare, and classification trust relationship almost shares the same characteristics as the behavioral trust relationship. Therefore, three real-world networks with classification trust relationship and mobile QQ Zone blog datasets with behavioral trust relationship were used in the experiments to evaluate E-MLCD.

World Championships in Athletics 2013: based on different behavioral trust relationship of this dataset, the network is divided into three layers, i.e., forwarding, mentioning, and replying.

Airline data: different airlines in Europe are defined as attribute. Each layer corresponds to an airline, yielding a network of 37 layers.

Dataset of staff in Aarhus University: five online and offline relationships between university staff are defined as trust relationship, and each attribute is regarded as a relationship. In this way, the network is divided into five layers.

QQ Zone dataset: this dataset was collected from actual QQ Zones of sales personnel. It contains classification of products, as well as gender and age of postissuers, repliers, and praises. The network is divided into three layers, i.e., posting, replying, and praising.

Table 1 shows the main characteristics of the four data we need to use. We use #Nodes to represent the number of multilayer nodes, #Edges to represent the number of edges of multilayer networks, #Layers to represent the number of layers, A_{deg} to represent the average number of nodes considering the average degree of multilateral nodes, and A_{layer} to represent the average number of layers of nodes.

4.2. Evaluation Methods and Metrics

4.2.1. Evaluation Methods. For the purpose of performance evaluation, E-MLCD was compared with three community detection algorithms, inferring multilayer global community structure. In order to compare our E-MLCD method of global community detection, after iterating the algorithm cycle, we finally get a global community structure.

(1) *Louvain (a Modularity-Based Algorithm That Can Detect Community Very Efficiently and Effectively).* In addition, it can detect layered communities. Its optimization objective is to maximize the modularity of the entire graph. On the contrary, Louvain is a nonoverlapping community detection method. By optimizing the value of the modularity function, it allocates each node to the "optimal" cluster, allowing it to process large-scale data efficiently. However, trust relationship and overlapped layers of the network are ignored.

(2) *LCD (a Local Expansion-Based Community Detection Algorithm).* It involves selection of original expansion

Input: Multilayer graph $G_{\mathcal{L}} = (V_{\mathcal{L}}, E_{\mathcal{L}}, V, \mathcal{L})$
 $V_0 \in V_{\mathcal{L}}, E_0 \in E_{\mathcal{L}}, L, N \in G_{\mathcal{L}}$
Output: Local community C for V_0

- 1: $G_{\mathcal{L}} = (V_{\mathcal{L}}, E_{\mathcal{L}}, V, \mathcal{L})$
- 2: $S = V$
- 3: while node $S_i \neq \emptyset$
- 4: v_0 ; // Utilize formula (9) to obtain the seed node
- 5: $S = S - \{v_0\}; C = \{v_0\}$
- 6: for node $v_j \in S$ do
- 7: if $T_{\text{int}}(u, v) \geq T_{\text{out}}(u, v)$ then // Utilize formulas (5) and (8) to determine whether it is intralayer or interlayer expansion
- 8: if $T_{\text{int}}(C_u, S_v) \geq T_{\text{int}}(S_u, S_v)$ then // Utilize formulas (3) and (4) to determine the intralayer nodes to be detected
- 9: $C = C \cup \{v_j\}$ // Combine the nodes to be detected to the local community C
- 10: else $S = S - \{v_j\}$ // Remove the nodes to be detected to set S
- 11: end if
- 12: else // Otherwise, it is interlayer expansion
- 13: if $T_{\text{out}}(C_u, S_v) \geq T_{\text{out}}(S_u, S_v)$ // Utilize formulas (6) and (7) to determine the interlayer nodes to be detected
- 14: $C = C \cup \{v_j\}$ // Combine the nodes to be detected to the local community C
- 15: else $S = S - \{v_j\}$ // Remove the nodes to be detected to set S
- 16: end if
- 17: end if
- 18: $j = j + 1$
- 19: : End for
- 20: $i = i + 1$
- 21: end while
- 22: return C
- 23: end

ALGORITHM 1: Trust-based algorithm for multilevel local community detection.

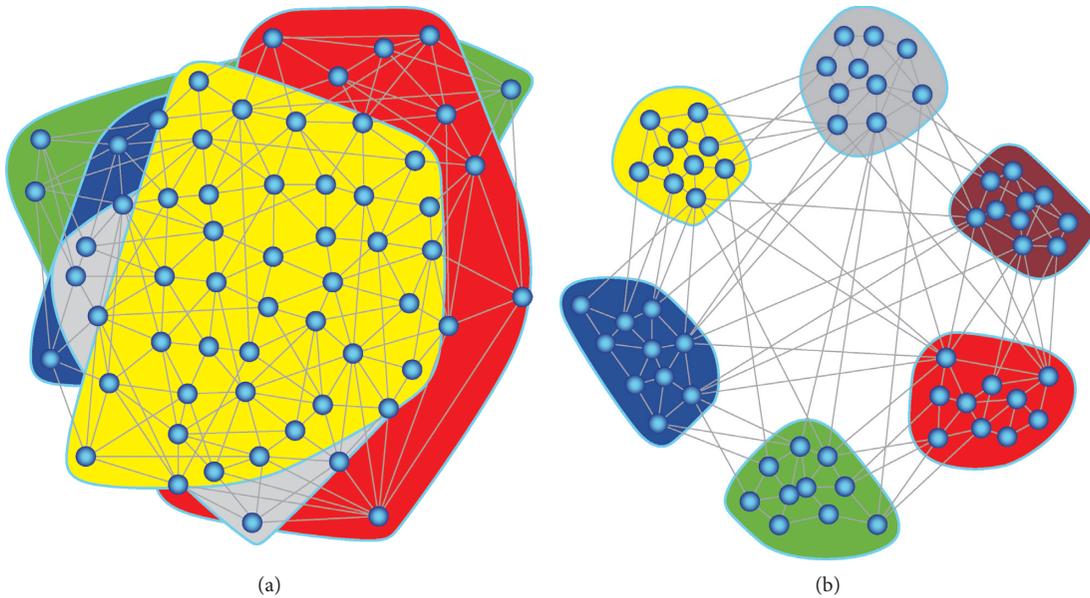


FIGURE 2: Multilevel community detection based on different node trusts.

TABLE 1: Major characteristics of multilevel datasets.

Dataset	#Nodes	#Edges	#Layers	Density	A_{deg}	A_{layer}
MoscowAthletics2013	103319	144591	3	1.68e-3	20.3	1.86
Airlines	417	3588	37	0.023	17.20	4.88
CS-AARHUS	61	620	5	0.122	20.32	3.67
QQ Zone personal selling	820484	16044572	3	2.38e-3	25.56	2.32

subgraphs, expansion strategies, and the conditions for terminating the expansion. In most of the local community detection methods, the expansion process is greedy. In other words, given a fitness function of a local community, the neighbour that can produce the most gains to the fitness is added to the community. The iterative process does not end until no neighbour can improve the local community. Note that LCD is very similar to our method, but LCD is unable to detect community within a multilevel network.

(3) *ML-LCD (a Nonsupervised Method for Community Detection within a Multilevel Network)*. After setting an internal-to-external connectivity ratio, it achieves local community detection by applying ML-LCD-lw, ML-LCD-wlsim, and ML-LCD-clsim to different layers. It supports layer-weight similarity and intralayer and interlayer similarities. This method achieves community detection by comparing connectivity but fails to take into account the level of similarity between communities, resulting in limited performance in the face of communities with many similar trust relationships.

4.2.2. *Evaluation Metrics*. The metrics used to evaluate E-MLCD include the scale of detected community, the scale of multilevel community, and the computational complexity.

(1) *Multilayer Modularity*. Algorithm performance is measured by the community scale. Table 2 shows the detection results of E-MLCD and the other algorithms.

It can be seen that E-MLCD produced the most communities for all datasets except Airlines. Also, the proposed algorithm was more effective in detecting community from dataset with behavioral trust relationship. Consider Moscow Athletics 2013 and QQ Zone datasets, which have obvious behavioral trust relationship such as replying and forwarding. The proposed algorithm achieved performance gain by categorizing these behavioral trust relationships into the same layer. In particular, a large number of attribute information such as gender, geographical location, and age were collected, and more communities were detected in the QQ area. But in the case of the Airlines dataset which has fewer trust relationship, our method was slightly inferior to ML-LCD.

(2) *Average Multilevel Module Club Evaluation*. In the multilayer network, it is impossible to use the network module degree to realize the evaluation, so the multilayer network module degree is used to evaluate the algorithm in this part. In a multitier network, the higher the Q , the better the result of the division of the community.

In the multilayer network, in addition to the largest community, each test result can also reflect the algorithm deviation. Figure 3 shows the size of the module degree obtained by the results of 50 runs by different algorithms under different datasets. Since the LART algorithm cannot display correct results in a large-scale datasets, only GL, PMM, and E-MLCD algorithms are shown in Figure 3.

According to the graphic analysis, in MoscowAthletics2013, the E-MLCD algorithm could obtain great

TABLE 2: Comparison of multilayer modularity for different algorithms.

Dataset	GL	LART	PMM	E-MLCD
MoscowAthletics2013	0.025	NA	0.630	0.680
Airlines	0.037	0.013	0.018	0.023
CS-AARHUS	0.249	0.154	0.222	0.261
QQ Zone personal selling	0.021	NA	0.580	0.620

results, and among 50 operations, most of them achieved great modularity. This is consistent with the results of the maximum modularity in Table 2. According to Table 1, we can see that A_{deg} of the MoscowAthletics2013 dataset was relatively low because the dataset was relatively sparse with large data volume. Therefore, we can see that the algorithm proposed in this section can obtain great modularity in a large-scale large network. In the BioProte dataset, MLCD could obtain similar results as the classic algorithms. In the dataset CS-AARHUS, the GL algorithm had stable performance, and its results were also higher than the results of the E-MLCD algorithm. Therefore, it can be seen that the E-MLCD algorithm does not have advantages in the analysis of small datasets. In the Airlines dataset, the results of our algorithm were lower than the results of the PMM algorithm and slightly lower than the results of the GL algorithm. This is mainly because the Airlines dataset had high A_{deg} and low A_{leyas} , and our algorithm had less trust relationship. The special case was the QQ Zone dataset. The GL and PMM algorithms obtained stable test results, and at most circumstances, the results of the GL algorithm were higher than the results of the E-MLCD algorithm. In the QQ Zone dataset, the algorithm proposed by us obtained high modularity for 6 times.

According to the comparison analysis, our algorithm could obtain good detection results in large-scale community and network with the sparse dataset. However, in closely connected network with small dataset, the detection result of our algorithm was poorer than that of the GL algorithm, and better than that of the PMM algorithm. This is mainly because the PMM algorithm needs to set corresponding prediction parameter in advance, and it cannot accurately detect community in large-scale unknown multilayer network. In the dataset with close network connection, the E-MLCD algorithm had lower or equal performance as the GL algorithm, and this is mainly because our algorithm mainly depends on attribute similarity and social intensity similarity to find community. In the dataset with close connection, its performance is insufficient compared to the GL algorithm.

(3) *Computation Efficiency*. In order to test the time complexity of different algorithms, in this chapter, we choose the dataset in which the E-MLCD algorithm could achieve good results: the QQ Zone dataset. Through comparison, we found that the QQ Zone dataset had a large number of nodes and low A_{deg} , so the time performance of various methods in this dataset was several orders of magnitude higher than that in other datasets, and the average running time was 500 minutes. The running time of other datasets was much

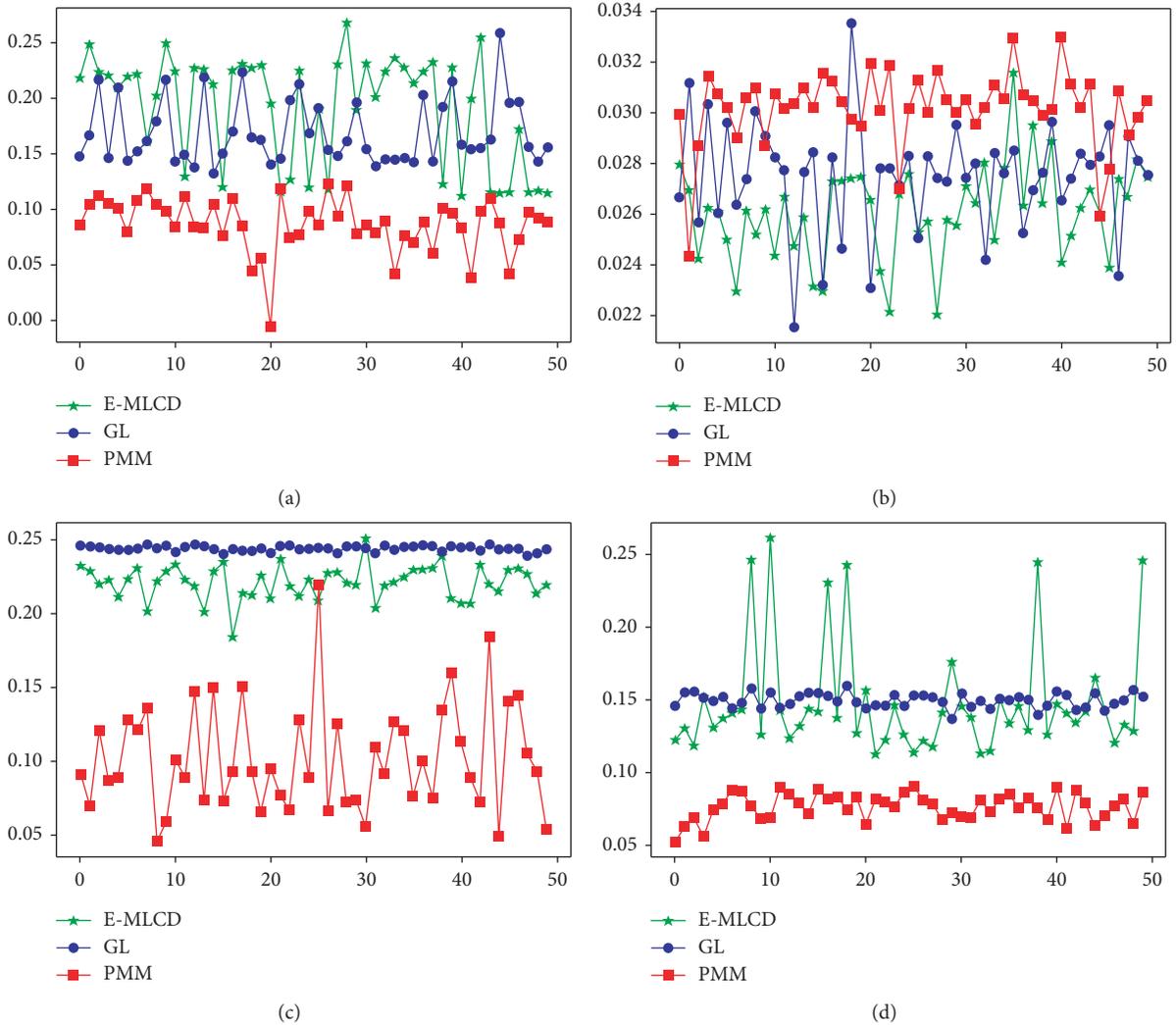


FIGURE 3: Average modularity of different algorithms. (a) MoscowAthletics2013. (b) Airlines. (c) CS-AARHUS. (d) QQ Zone.

shorter, especially the CS-AARHUS dataset, which was measured in second, and it could not well reflect the time performance of algorithm. Therefore, during the comparison of various algorithms, the operating efficiency on sparse dataset QQ Zone was used as the benchmark for comparison. In this paper, different algorithms were used for analysis in the QQ Zone dataset, and Figure 4 shows the final test results. We can see that, for all 3 algorithms, the running time presented direct-proportion distribution with the network scale; however, the running time of E-MLCD was significantly shorter than that of the other 2 algorithms, the running time of PMM algorithm was shorter than that of the GL algorithm, and during several operations, the running time of PMM was also shorter than that of our algorithm.

As shown in Figure 4, the E-MLCD algorithm can be realized in short time in most operations, it has better time performance than the GL and PMM algorithms, and its stability was also better than that of the GL and PMM algorithms. The overall stability of PMM algorithm was better than that of the GL algorithm, and it had similar time performance to our algorithm. The GL algorithm had the

longest operating time, and it has poorer stability than the E-MLCD and PMM algorithms. According to Figure 4, in the QQ Zone dataset, our algorithm could obtain better modularity than the GL and PMM algorithms. Therefore, our algorithm can efficiently process the large-scale sparse network of multiple layers and trust relationship.

5. Conclusions

Detecting user groups of a particular interest from the social network has been the ultimate goal of social network advertisements. This is because targeted advertising can considerably increase its effectiveness and maximize profits. These groups, also known as communities, open up opportunities for a new marketing pattern based on acquaintance circle and six-dimensional space. By clustering users with common interests into customer groups, detecting closely interrelated network users is essential for improving advertisement prompting via social media and extracting potential customers from social network. Based on three real-world social networks and QQ Zone marketing

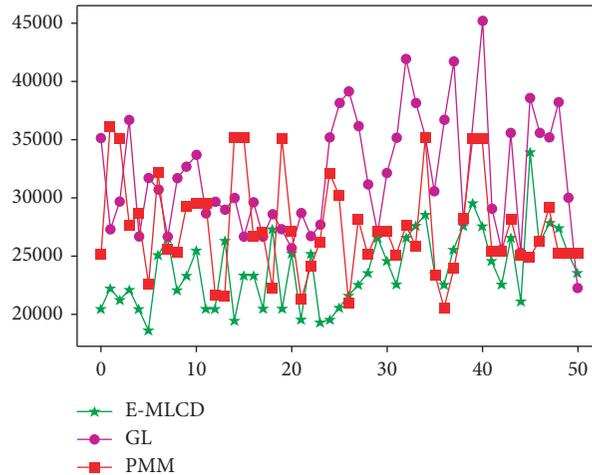


FIGURE 4: Running time of different algorithms in the QQ Zone dataset after running for 50 times.

data, we proposed a new local community detection model E-MLCD, which jointly considers multilevel trust relationship and community structure. To address the problem of expanding local community during the detection process, this chapter proposed the local community detection model based on multilayer trust relationship and community structure (E-MLCD) for the first time. For the local community detection and expansion problem, this model defined new measurement with similar community intensity based on similarity of community structure. The E-MLCD method can fully utilize the structure and attribute information to realize network partition and promote application of multilayer network, such as obtaining better detection in sparse network based on partial attribute and realizing comparison with popular and similar algorithms; the E-MLCD algorithm has advantages in the analysis of a large-scale multilayer network with sparse connection, and it can effectively identify sparse community structure in a large-scale multilayer network and obtain better time performance.

Data Availability

MoscowAthletics2013 datasets [9], Airlines datasets [40], and CS-AARHUS datasets [7] are cited at relevant places within the text as references. The QQ Zone personal selling data used to support the findings of this study are available from the corresponding author upon request.

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

Acknowledgments

This research is partially supported by the Fundamental Research of Xinjiang Corps (2016AC015) and the Applied Basic Research Project of Qinghai Province (No. 2018-ZJ-707).

References

- [1] M. Kivela, A. Arenas, M. Barthelemy, J. P. Gleeson, Y. Moreno, and M. A. Porter, "Multilayer networks," *Journal of Complex Networks*, vol. 2, no. 3, pp. 203–271, 2014.
- [2] C. W. Loe and H. J. Jensen, "Comparison of communities detection algorithms for multiplex," *Physica A: Statistical Mechanics and its Applications*, vol. 431, pp. 29–45, 2015.
- [3] J. Kim and J.-G. Lee, "Community detection in multi-layer graphs: a survey," *ACM SIGMOD Record*, vol. 44, no. 3, pp. 37–48, 2015.
- [4] G. D'Agostino and A. Scala, *Networks of Networks: the Last Frontier of Complexity*, Vol. 340, Springer, Berlin, Germany, 2014.
- [5] D. Cai, Z. Shao, X. He, X. Yan, and J. Han, "Community mining from multi-relational networks," in *European Conference on Principles of Data Mining and Knowledge Discovery*, pp. 445–452, Springer, Porto, Portugal, 2015.
- [6] M. Berlingerio, F. Pinelli, and F. Calabrese, "ABACUS: frequent pAttern mining-based community discovery in multidimensional networks," *Data Mining and Knowledge Discovery*, vol. 27, no. 3, pp. 294–320, 2013.
- [7] M. Magnani, B. Micenkova, and L. Rossi, "Combinatorial analysis of multiple networks," *Computer Science*, 2013, <https://arxiv.org/abs/1303.4986>.
- [8] M. De Domenico, V. Nicosia, A. Arenas, and V. Latora, "Structural reducibility of multilayer networks," *Nature Communications*, vol. 6, no. 1, p. 6864, 2015.
- [9] W. Wang, X. Li, P. Jiao et al., "Exploring intracity taxi mobility during the holidays for location-based marketing," *Mobile Information Systems*, vol. 2017, Article ID 6310827, 10 pages, 2017.
- [10] A. Clauset, "Finding local community structure in networks," *Physical Review E*, vol. 72, no. 2, article 026132, 2005.
- [11] J. P. Bagrow and E. M. Bollt, "Local method for detecting communities," *Physical Review E*, vol. 72, no. 4, article 046108, 2005.
- [12] A. Lancichinetti, S. Fortunato, and J. Kertész, "Detecting the overlapping and hierarchical community structure in complex networks," *New Journal of Physics*, vol. 11, no. 3, article 033015, 2009.
- [13] Q. Chen, T.-T. Wu, and M. Fang, "Detecting local community structures in complex networks based on local degree central

- nodes,” *Physica A: Statistical Mechanics and its Applications*, vol. 392, no. 3, pp. 529–537, 2013.
- [14] M. A. Tabarzad and A. Hamzeh, “A heuristic local community detection method (HLCD),” *Applied Intelligence*, vol. 46, no. 1, pp. 62–78, 2016.
- [15] D. Chen, L. Lü, M.-S. Lü, Y.-C. Zhang, and T. Zhou, “Identifying influential nodes in complex networks,” *Physica A: Statistical Mechanics and its Applications*, vol. 391, no. 4, pp. 1777–1787, 2012.
- [16] P. Wang and J. Liu, “A multi-agent genetic algorithm for local community detection by extending the tightest nodes,” in *Proceedings of 2016 IEEE Congress on Evolutionary Computation*, pp. 3215–3221, Vancouver, BC, Canada, July 2016.
- [17] Z. Bu, Z. Wu, J. Cao, and Y. Jiang, “Local community mining on distributed and dynamic networks from a multiagent perspective,” *IEEE Transactions on Cybernetics*, vol. 46, no. 4, pp. 986–999, 2016.
- [18] X. Li, Q. Tian, M. Tang, X. Chen, and X. Yang, “Local community detection for multi-layer mobile network based on the trust relation,” *Wireless Networks*, 2019, In press.
- [19] X. M. Li, L. Yuan, C. C. Liu et al., “An efficient critical incident propagation model for social networks based on trust factor,” in *Proceedings of International Conference on Collaborative Computing: Networking, Applications and Worksharing*, pp. 416–424, Springer, Cham, Switzerland, December 2017.
- [20] M. De Domenico, A. Lancichinetti, A. Arenas et al., “Identifying modular flows on multilayer networks reveals highly overlapping organization in interconnected systems,” *Physical Review X*, vol. 5, no. 1, p. 011027, 2015.
- [21] R. Lambiotte, J.-C. Delvenne, and M. Barahona, “Random walks, markov processes and the multiscale modular organization of complex networks,” *IEEE Transactions on Network Science and Engineering*, vol. 1, no. 2, pp. 76–90, 2014.
- [22] V. Carchiolo, A. Longheu, M. Malgeri et al., “Communities unfolding in multislice networks,” in *Complex Networks*, pp. 187–195, Springer, Berlin, Heidelberg, Germany, 2011.
- [23] I. Gaye, G. Mendy, S. Ouya, I. Diop, and D. Seck, “Multi-diffusion degree centrality measure to maximize the influence spread in the multilayer social networks,” in *Proceedings of International Conference on e-Infrastructure and e-Services for Developing Countries*, pp. 53–65, Ouagadougou, Burkina Faso, West Africa, December, 2016.
- [24] Q. Han, K. Xu, and E. Airoldi, “Consistent estimation of dynamic and multi-layer block models,” in *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pp. 1511–1520, Lille, France, July 2015.
- [25] M. Porter, “Community structure in time-dependent networks,” in *Proceedings of APS March Meeting*, Portland, Oregon, March 2010.
- [26] L. Tang, X. Wang, and H. Liu, “Community detection via heterogeneous interaction analysis,” *Data Mining and Knowledge Discovery*, vol. 25, no. 1, pp. 1–33, 2011.
- [27] G. Zhu and K. Li, “A unified model for community detection of multiplex networks,” in *Proceedings of International Conference on Web Information Systems Engineering*, pp. 31–46, Thessaloniki, Greece, October 2014.
- [28] P. Bródka, T. Filipowski, and P. Kazienko, “An introduction to community detection in multi-layered social network,” in *Proceedings of World Summit on Knowledge Society*, pp. 185–190, Mykonos, Greece, September 2011.
- [29] X. Li, G. Xu, and M. Tang, “Community detection for multi-layer social network based on local random walk,” *Journal of Visual Communication and Image Representation*, vol. 57, pp. 91–98, 2018.
- [30] Z. Yu, H.-S. Wong, and H. Wang, “Graph-based consensus clustering for class discovery from gene expression data,” *Bioinformatics*, vol. 23, no. 21, pp. 2888–2896, 2007.
- [31] V. Filkov and S. Skiena, “Integrating microarray data by consensus clustering,” *International Journal on Artificial Intelligence Tools*, vol. 13, no. 4, pp. 863–880, 2004.
- [32] F. Saeed, N. Salim, and A. Abdo, “Combining multiple clusterings of chemical structures using cluster-based similarity partitioning algorithm,” *International Journal of Computational Biology and Drug Design*, vol. 7, no. 1, pp. 31–44, 2014.
- [33] V. Buskens, “The social structure of trust,” *Social Networks*, vol. 20, no. 3, pp. 265–289, 1998.
- [34] W. Sherchan, S. Nepal, and C. Paris, “A survey of trust in social networks,” *ACM Computing Surveys (CSUR)*, vol. 45, no. 4, pp. 1–33, 2013.
- [35] S. Trifunovic, F. Legendre, and C. Anastasiades, “Social trust in opportunistic networks,” in *Proceedings of INFOCOM IEEE Conference on Computer Communications Workshops, 2010*, pp. 1–6, San Diego, CA, USA, March 2010.
- [36] G. Wang, F. Musau, S. Guo, and M. B. Abdullahi, “Neighbor similarity trust against sybil attack in P2P e-commerce,” *IEEE transactions on parallel and distributed systems*, vol. 26, no. 3, pp. 824–833, 2015.
- [37] Y. Jin, Y. Zhang, W. Qu et al., “A trust model based on similarity evaluation in P2P networks,” in *Proceedings of IEEE International Symposium on Parallel and Distributed Processing with Applications*, pp. 737–742, IEEE Computer Society, Sydney, NSW, Australia, December 2008.
- [38] C. N. Ziegler and G. Lausen, “Analyzing correlation between trust and user similarity in online communities,” in *Proceedings of International Conference on Trust Management*, pp. 251–265, Springer, Berlin, Heidelberg, June 2004.
- [39] L. Boratto, S. Carta, A. Chessa et al., “Group recommendation with automatic identification of users communities, web intelligence and intelligent agent technologies,” in *Proceedings of WI-IAT’09 2009 IEEE/WIC/ACM International Joint Conference*, vol. 3, pp. 547–550, Milano, Italy, September 2009.
- [40] G. Xu, Z. Feng, H. Wu, and D. Zhao, “Swift trust in a virtual temporary system: a model based on the dempster-shafer theory of belief functions,” *International Journal of Electronic Commerce*, vol. 12, no. 1, pp. 93–126, 2014.
- [41] S. Deng, L. Huang, G. Xu et al., “On deep learning for trust-aware recommendations in social networks,” *IEEE Transactions on Neural Networks & Learning Systems*, vol. 28, no. 5, pp. 1164–1177, 2016.
- [42] X. Chen, C. Xia, and J. Wang, “A novel trust-based community detection algorithm used in social networks,” *Chaos, Solitons & Fractals*, vol. 108, pp. 57–65, 2018.
- [43] V. Nicosia and V. Latora, “Measuring and modeling correlations in multiplex networks,” *Physical Review E Statistical Nonlinear & Soft Matter Physics*, vol. 92, no. 3, article 032805, 2015.

Research Article

RoC: Robust and Low-Complexity Wireless Indoor Positioning Systems for Multifloor Buildings Using Location Fingerprinting Techniques

Kriangkrai Maneerat  and Kamol Kaemarungsi 

National Electronics and Computer Technology Center, NSTDA, Pathumthani, Thailand

Correspondence should be addressed to Kamol Kaemarungsi; kamol.kaemarungsi@nectec.or.th

Received 26 September 2018; Accepted 6 January 2019; Published 3 February 2019

Guest Editor: Jaegeol Yim

Copyright © 2019 Kriangkrai Maneerat and Kamol Kaemarungsi. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Most existing wireless indoor positioning systems have only success performance requirements in normal operating situations whereby all wireless equipment works properly. There remains a lack of system reliability that can support emergency situations when there are some reference node failures, such as in earthquake and fire scenarios. Additionally, most systems do not incorporate environmental information such as temperature and relative humidity level into the process of determining the location of objects inside the building. To address these gaps, we propose a novel integrated framework for wireless indoor positioning systems based on a location fingerprinting technique which is called the Robust and low Complexity indoor positioning systems framework (RoC framework). Our proposed integrated framework consists of two essential indoor positioning processes: the system design process and the localization process. The RoC framework aims to achieve robustness in the system design structure and reliability of the target location during the online estimation phase either under a normal situation or when some reference nodes (RNs) have failed. The availability of low-cost temperature and relative humidity sensors can provide additional information for the location fingerprinting technique and thereby reduce location estimation complexity by including this additional information. Experimental results and comparative performance evaluation revealed that the RoC framework can achieve robustness in terms of the system design structure, whereby it was able to provide the highest positioning performance in either *fault-free* or *RN-failure scenarios*. Moreover, in the online estimation phase, the proposed framework can provide the highest reliability of the target location under the RN-failure scenarios and also yields the lowest computational complexity in online searching compared to other techniques. Specifically, when compared to the traditional weighted k-nearest neighbor techniques (WKNN) under the 30% RN-failure scenario at Building B, the proposed RoC framework shows 74.1% better accuracy performance and yields 55.1% lower computational time than the WKNN.

1. Introduction

Indoor positioning systems (IPSS) refer to wireless network infrastructure systems that provide location information to any requesting user inside an indoor operating area such as airports, shopping centers, and hospitals. These systems are currently experiencing a tremendous growth and becoming a vital part of life in the digital age [1]. The use of unlicensed frequency spectrum ranges and inexpensive wireless communication technologies has facilitated the deployment of IPSS [2]. They can be applied to various domains including for indoor navigation and tracking in the health-care sector, in

industrial areas, and at trade fairs [3, 4]. However, most of these IPSS only have success performance requirements under normal operating situations whereby all of the wireless equipment works properly. There remains a lack of system reliability that can perform in rescue situations, such as localization that supports rescuers in an earthquake scenario or firefighters in a fire situation. Under these emergency situations, some wireless network equipment (i.e., reference nodes) may fail due to some parts of the building being damaged. Therefore, the key performance requirement of the IPS under such node-failure situations is the reliability of the estimated target locations under either a normal situation or

when some reference nodes (RNs) have failed. Moreover, under such unexpected situations, the computational complexity during the online searching process is an important performance requirement for the IPSs.

The general indoor positioning approaches can be divided into three groups: Triangulation, Proximity, and Scene Analysis [5]. First, Triangulation uses the geometric characteristics of triangles to estimate the target location. It has two types, namely, lateration and angulation. The lateration method determines the position of a target by using the relationship of distances between various transmitters and receivers. The accuracy performance of the lateration method depends on the precise time synchronization of all transmitters and receivers. Thus, lateration requires the high cooperation and synchronization of all devices in the system [2, 3]. Examples of IPSs based on the lateration approach are Time of Arrival (TOA) and Time Difference of Arrival (TDOA). Another technique that is based on Triangulation is angulation. This technique estimates the position of an object by exploiting the relationships of angle direction lines between the transmitter and receiver. The angulation method can estimate the position of a target in a three-dimensional (3D) coordinate system, in which no time synchronization between transmitter and receiver is required. However, it requires complex antenna hardware to measure the angle of incidence. Moreover, the precision of the angle measurement may be limited by the multipath effect and the non-line-of-sight (NLOS) propagation of signals [2, 3]. An example of an IPS based on the angulation approach is Angle of Arrival (AOA).

Second, Proximity is a simple indoor positioning approach that provides symbolic relative location information. This approach is relatively simple to implement and requires low complexity hardware. Proximity is often used in the IPSs that deploy Radio Frequency Identification (RFID) technology. However, this method provides the lowest accuracy performance of these three mentioned approaches. It can only identify the approximate position of the target by successfully communicating with one or more transceivers [4, 6]. The proximity is often applied in specific applications such as the exhibitor check-in process and information advertising for visitors in a museum [2, 4].

Third, Scene Analysis uses the characterization of an indoor radio propagation that is associated with the particular location. This technique requires a calibration training phase called the offline calibration phase to compile a radio map (e.g., the location fingerprinting) of the service area, before the indoor positioning algorithm is used to determine the estimated location during the online estimation phase. The advantage of this technique is that it is simple to deploy with no specialized hardware required. It can employ various wireless technologies such as Wireless Local Area Network (WLAN) and Bluetooth Low Energy (BLE). Moreover, this technique can provide high accuracy and precision performance without being affected by multipath fading and shadowing. However, the drawback of this technique is that it is very time-consuming to perform the exhaustive data collection for a wide operation area [1, 2, 7]. Examples of IPSs that are based on Scene Analysis are the

probabilistic, neural networks, and location fingerprinting (i.e., Euclidean distance technique).

From the property aspects of those three indoor positioning approaches, in this article, we focus on the development of IPSs based on Scene Analysis, in which the location fingerprinting technique is considered. This technique is simple to implement and does not require specialized hardware with no time synchronization necessary between transmitter and receiver. Moreover, it may be implemented in the software format which can reduce the complexity and cost significantly compared to the lateration and angulation methods [2, 7]. We employ the IEEE 802.15.4 wireless sensor networks (WSNs) with ZigBee. They are lightweight devices with ultra-low power consumption that can detect several environmental parameters such as temperature and relative humidity [8]. In this paper, we also exploit the association of environmental data, which are the temperature and the relative humidity, as a part of Scene Analysis that can be used to coarsely separate the area of location for an object.

According to the literature reviewed in Section 2, to fulfill the gaps and to extend the ability of the IPSs under unexpected situations (e.g., earthquake and fire situations), reliable estimated target locations and a short time to estimate them are required. Thus, in this article, we propose a novel integrated framework for the IPSs based on the location fingerprinting technique which is called the Robust and low Complexity indoor positioning systems framework (RoC framework). The RoC framework consists of two essential indoor positioning processes: the system design process and the localization process. Our proposed framework can achieve the robustness of the system in terms of the system design structure and the reliability of the target location in the localization process. In particular, the RoC framework can provide low computational complexity during the online searching time without compromising the positioning accuracy. The major contributions of our integrated framework include the following:

- (1) Development of a robust system design that is provisioned to support RN-failure situations
- (2) Development of a robust floor determination algorithm that can efficiently work under RN-failure situations
- (3) Reduction of computational complexity during the location estimation process by using temperature-level classification
- (4) Development of a novel Active Euclidean distance approach that can solve the missing received signal strength (RSS) problem

The remaining organization of this article is as follows. First, Section 2 summarizes and compares related works on design with robustness and low complexity of the IPSs. Second, we describe our proposed RoC framework for the IPSs in Section 3. Third, Section 4 explains about the experimental environment, the hardware specifications used for data collection, and the setup parameters used in this work. Fourth, the experimental results are presented and discussed in Section 5. Finally, Section 6 concludes the

performance achievement of our proposed framework and future research plan.

2. Related Works

Similar to any information technology systems, the development and deployment of wireless indoor positioning system (IPS) can be divided into two processes: the system design process and the localization process. First, in the system design process, a system designer has to gather all necessary requirements such as dimension of the building and determine various wireless network parameters that will affect the performance of the IPS. To ensure a required performance or design goal for the IPS, a well-designed methodology is indeed the essential process. For instance, the required number and placement of reference nodes (RNs) must be determined based on predefined criteria during this process. Second, after the system is successfully installed, the IPS can be used to estimate any target locations in the localization process. Using localization algorithms in this second process, the physical environmental information such as received signal strength information (RSSI) or time of arrival (TOA) measured at the target is analyzed to calculate its position on a coordinate system. Note that the localization process for location fingerprinting approach is further divided into offline calibration phase and online estimation phase. The offline calibration phase is when the location fingerprints are collected to create a location fingerprinting database, while the online estimation phase is when the location determination is actually performed by comparing or matching the location fingerprinting database with the current location pattern such as RSSI pattern of the target.

To ensure the success performance requirements of the IPS based on the location fingerprinting approach, an effective framework that can achieve the consistent relationship between the system design process and the localization process is required. On the other hand, the positioning performance of the system might be dropped if both processes do not have a consistent relationship with each other. For example, the main objective of the node placement design for Triangulation is to place a sufficient number of RNs in optimum locations so that the signal test points (STPs) must be able to receive signals from at least *three* RNs installed in the service area [2]. They are different from the structure design for Scene Analysis (e.g., the location fingerprinting technique), which need as many received signals as possible to make a location estimate. For the Scene Analysis structure design, the STPs must be able to receive signals from at least *four* RNs installed inside the service area to handle the symmetrical RSS problem [9].

For the system design, existing research in the literature has focused on different performance requirements of the system design. The authors of [10–12] presented the system design structure for the IPSs based on Triangulation. They focused on the improvement of the positioning performance under a limited number of RNs and minimizing the number of RNs. Redondi et al. [10] proposed an RN-placement method for the IPSs based on the Cramer–Rao lower bound (CRLB) approach. Their design goal was to minimize

localization errors when operating with a limited number of RNs due to fixed budget constraints. Merkel et al. [11] presented an optimal RN-placement approach for the IPSs based on distributed range-free localization. Achieving the optimal coverage of a certain area while simultaneously minimizing the necessary number of RNs was their focus. Tong et al. [12] proposed an optimum RN placement for the TOA technique by minimizing the Cramer–Rao Bound (CRB) of localization error. The objective of their system design was to achieve the highest localization accuracy when the reference nodes have uniform angular distribution around the service area.

The authors of [13–16] presented their RN-placement design for the IPSs based on Scene Analysis such as the location fingerprinting techniques. They focused on the performance requirements including the reduction of data redundancy and the guarantee of signal coverage. To obtain better location accuracy, Sharma et al. [13] developed an RN-placement method that attempts to minimize the total number of similar fingerprints (SFs). Fang and Lin [14] proposed a framework for relating the positioning performance with the RN placement. To achieve their design objective, their algorithm determined a suitable set of RN locations such as the signal-to-noise ratio (SNR) was maximized. Chen et al. [15] proposed a novel RN-placement algorithm for the IPSs based on the location fingerprinting technique. They focused on maximizing the fingerprint difference (FD) while still guaranteeing the coverage requirement by using the least number of APs. Kondee et al. [16] proposed a novel RN-placement technique using the Binary Integer Linear Programming (BILP) approach to design efficient wireless indoor positioning systems based on the location fingerprinting technique. The main goal of their system design was to install the RNs and improve the location determination performance for a single-floor area and the multifloor building.

In the existing system design literature, most research focused on the IPS performance requirements, including minimizing the number of installed RNs, guaranteeing radio signal coverage requirements in the service area, and minimizing the total number of similar fingerprint locations. A system design for robust IPSs that supports the RN-failure scenario caused by hardware failures or software errors has only been investigated once in our previous work in [17] but not in any other existing literatures. Additionally, physical environment parameters such as temperature and relative humidity inside buildings were not exploited for IPSs. Based on these knowledge gaps, the design of RN placement for robust IPSs and the utilization of temperature and humidity data to improve localization are still open research issues. Therefore, first, in the system design process, we propose a robust system design model for wireless indoor positioning systems based on location fingerprinting techniques. Our proposed model aims to place a suitable number of RNs and to determine their locations whereby their placement is provisioned to support robust system operation both during a normal situation and when some RNs have failed. Table 1 shows a summary of system design approaches for indoor positioning systems.

TABLE 1: Recent system design approaches for indoor positioning systems.

Category	Scheme	Localization algorithm	Optimization method	Design objectives	Design limitations	Service areas	Robustness (support in RN-failure situation)
Triangulation	[10]	RSS-based	Tabu search	To minimize localization error	The design does not consider minimizing the number of RNs	Single floor	No
	[11]	Distance-based	Genetic Algorithm	To maximize signal coverage and minimize number of RNs	The selected fitness evaluation influences the localization results	Single floor	No
	[12]	TOA	—	To minimize localization error	The RN placement has only a uniform angular distribution	Single floor	No
Scene analysis	[13]	KNN	Simulated Annealing	To minimize total number of similar fingerprints	The dynamic nature of the indoor environment is not considered in the system design	Single floor	No
	[14]	WKNN	—	To maximize signal RSS and minimize noise	The design does not consider the signal coverage in the physical surroundings	Single floor	No
	[15]	KNN	Simulated Annealing	To maximize fingerprint difference	Complicated indoor layouts influence the minimum number of RNs	Single floor	No
	[16]	Euclidean distance	Simulated Annealing	To maximize summation of maximum RSS	The RN placement design lacks system reliability	Single-floor and multifloor	No
	Proposed design	Active Fusion	Simulated Annealing	To minimize number of RNs and find their optimum locations to be provisioned to support RN-failure situations	The design considers only the discrete candidate sites for installing RNs	Single floor and multifloor	Yes

Next, in the localization process, we investigate the existing indoor positioning research that has focused on robust localization estimation algorithms. The authors of [18, 19] presented robust floor determination algorithms that considered the movement of objects across the floors. Gupta et al. [18] studied robots moving between floors by using incorporative information from the pressure sensor and the Wi-Fi access points (APs). Their floor determination algorithm is based on the RSS received from the APs that utilize a Maximum Likelihood (ML) to estimate the floor location of the robots. Lee and Park [19] presented a technique to estimate the robot position in each floor by using gyroscopes to recognize the robot motion status on the stairs. The authors of [20–23] presented robust and reliability positioning techniques for the IPSs based on Scene Analysis. They focused on handling the problems of different mobile devices and reducing the impact of environmental dynamics on the accuracy performance. He et al. [20] studied the problem of indoor localization in the case of equipment heterogeneity and an indoor dynamic environment. They proposed a Hierarchical Edit Distance (HED) algorithm based on the location fingerprinting technique. This algorithm can improve the robustness of systems under various environment dynamics and different factors of the target device such as Wi-Fi chipsets. Zayets and Steinbach

[21] proposed a novel location fingerprinting technique which operates by extracting and analyzing individual multipath propagation delays. They presented the localization algorithm based on the Multipath Component Analysis (MCA) to be robust against changes in the environment. Taniuchi and Maekawa [22] proposed a new Wi-Fi indoor positioning algorithm by which the system can be robust over unstable Wi-Fi APs. Their algorithm was based on the location fingerprinting technique that used the ensemble learning approaches to handle the problems of Wi-Fi positioning caused by unstable and uncontrollable infrastructure such as the movement of people. Guan et al. [23] presented a localization algorithm that reduces the influence of symmetry by considering extrarestrictions from redundant APs. They used a novel clustering method to robustly estimate the target location from the fingerprint location candidates.

According to the literature reviewed on the robust localization process, although some existing works have studied the robust and reliable location problem [18–23], they focused on robustness issues, including robust tracking of floor changes when the robots moved in the staircase, improving the robustness of dynamic changes to the indoor environment which lead to the instability of fingerprint information, and reducing the influence of different mobile

devices. In particular, they did not consider the robustness in terms of faulty RNs during the online estimation phase. Thus, in the localization process, we propose a fault tolerance positioning algorithm that can overcome the problem of RN failures. The contributions of our proposed robust positioning algorithm include both the robust floor determination algorithm and the robust location estimation (x, y) .

Other existing IPS research focused on improving the time complexity performance. The authors of [24–27] presented a recent research on the development of the indoor positioning performance in terms of computational complexity for Scene Analysis. Xue et al. [24] proposed an improved neighboring fingerprint location selection method for Wi-Fi indoor localization. They proposed a clustering algorithm based on k-means to classify the nearest neighboring fingerprint location according to its physical distance to the test point. Zhou and Van [25] proposed a location fingerprinting algorithm based on fuzzy c-means (FCM) clustering. Their algorithm employs offline clustering to reduce the online computational complexity of the matching process. Lee et al. [26] proposed a novel Support Vector Machine based Clustering (SVM-C) approach for the large database searching problem. Their proposed algorithm was further able to reduce the mean localization errors and reduce the computational complexity of location estimation. Saha and Sadhukhan [27] presented a novel hierarchical clustering strategy for the location fingerprinting technique that would reduce the searching time in the online estimation phase. Their proposed clustering strategy selects an appropriate cluster in the localization process based on the descending order of APs having the strongest signal strength in the observed RSS vector.

According to the literature reviewed on the development of computational complexity for the IPSs based on location fingerprinting techniques, recent works have aimed at improving overhead searching during the online estimation phase by using a classification approach. They have utilized effective cluster analysis, such as FCM, SVM-C, k-means clustering, and hierarchical clustering. However, their computational complexity approaches could take a long time to search for the closest location solution, in which the worst-case complexity of their clustering algorithms are $O(ndc^2)$, $O(n^3)$, $O(ndc)$, and $O(n^3)$ [28, 29], respectively, where n refers to the number of data points, d represents the number of dimensions, and c represents the number of clusters. Note that we consider the worst-case complexity during the online estimation phase, which is the maximal complexity of the algorithm over all possible inputs. Thus, in this article, we present a low-complexity algorithm based on the classification approach to reduce online searching time without compromising the positioning accuracy. Our proposed clustering algorithm utilizes a temperature-level filter based on a binary search algorithm. The time complexity of the proposed algorithm was lower than those four existing clustering algorithms, in which our algorithm has the worst-case complexity of $O(\log n)$ [29]. Table 2 shows different indoor positioning systems considering robustness and complexity.

In summary, we developed an integrated framework for the IPSs based on a location fingerprinting approach called

the RoC framework. The proposed framework consists of two main processes, the system design process and the localization process. The proposed framework is developed to provide robustness in terms of the system design structure, both during a normal situation and when some RNs fail. Moreover, in the online estimation phase, we have developed the framework to be able to achieve the reliability of the target location under the RN-failure scenarios and provide low-computational complexity using additional temperature and relative humidity data. Furthermore, our proposed framework can be applied to various service area structures ranging from single-floor to multiple-floor environments.

3. Development of Indoor Positioning Systems

In this section, we describe our RoC framework for the IPSs based on location fingerprinting techniques. First, in Section 3.1, an overview of the RoC framework for wireless indoor multifloor positioning systems is described. Then, in Section 3.2, the development of the system design process is explained. Finally, in Section 3.3, the localization process of our proposed framework is described.

3.1. An Overview of RoC Framework for Wireless Indoor Positioning Systems. Figure 1 shows an overview of the RoC framework for the IPSs based on the location fingerprinting techniques, including the system design process and the localization process. Figure 1(a) represents the first process of the framework as the system design process. This process starts with the design of the wireless network infrastructure in which the RNs are the wireless nodes of the network. The aims of the process are to put in place a suitable number of RNs and to determine their optimum locations. Note that the RN's locations may be on a single floor or on multiple floors. In particular, the proposed system designs provisions to support robust operation, both during a normal situation and when there are some RN failures (which will be discussed in Section 3.2). Note that this part of the framework was originally presented in detail in our previous work [17]. Figure 1(b) represents the second process of the RoC framework as the localization process based on location fingerprinting techniques. This process is divided into two phases: the offline calibration phase and the online estimation phase. During the offline calibration phase, each fingerprint location (i.e., represented by the points on a map) records the physical environment of the service area into the fingerprint database. In this article, we use three aspects of physical environment information for creating the database. These aspects consist of the RSS values received from each RN, the temperature value, and the relative humidity value. Note that in [17], only RSS values were considered. Both the temperature value and the relative humidity value are used to classify the similar physical environments of each subarea, which are called the temperature level (TEMP level) (which will be discussed in Section 3.3.1). During the online estimation phase, whenever there is a request for a target location, the mobile target will measure the three environment

TABLE 2: Different indoor positioning systems considering robustness and complexity.

Category	Scheme	Focus on	Location results	Localization algorithm	Robustness	Complexity	Additional sensor
—	[18]	Robustness	Floor	ML	The algorithm provides robustness in terms of the movement of objects across the floors	Low	Pressure sensor
—	[19]	Robustness	Floor	Particle filter	The algorithm provides robust tracking of floor changing when the robots moved in the staircase	High	Gyroscopes
Scene analysis	[20]	Robustness	x, y	HED	The system can provide the robustness of systems under the environment dynamics	Medium	—
	[21]	Robustness	x, y	MCA	The algorithm is can tolerate against changes in the environment	Medium	—
	[22]	Robustness	x, y	Ensemble learning	The system can handle the problems of the Wi-fi positioning caused by the unstable infrastructure	High	—
	[23]	Robustness	x, y	Probabilistic based	The algorithm can reduce the influence of symmetry RSS signal	Low	—
	[24]	Complexity	x, y	k-means	—	Medium	—
	[25]	Complexity	x, y	FCM	The algorithm can classify the data into clusters which are robust to the multipath effect	Medium	—
	[26]	Complexity	x, y	SVM-C	—	High	—
	[27]	Complexity	x, y	Hierarchical clustering	The strategy is more robust in highly fluctuating RSS measurements	High	—
	Proposed algorithm	Robustness and complexity	$x, y, \text{ floor}$	RMoS floor + Active Fusion	The algorithm can overcome the problem of RN failures	Low	Temp. and humid. sensor

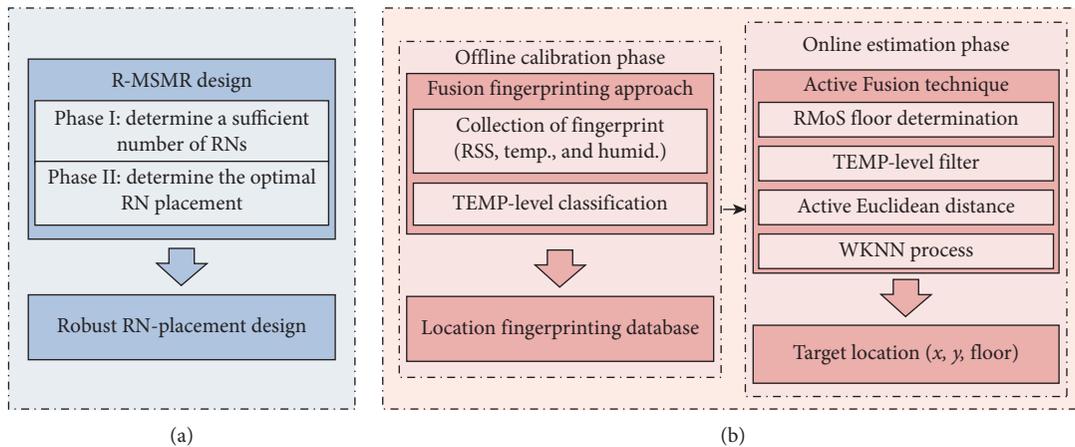


FIGURE 1: An overview of the RoC framework for the IPSs based on location fingerprinting techniques. (a) System design process. (b) Localization process.

parameters and use this online information to calculate the target location by using four processes of online estimation phase (which will be discussed Section 3.3.2).

3.2. System Design Process. In wireless IPS designs, the insufficient placement of RNs can cause an accuracy performance degradation. Furthermore, the system may achieve less than half of its intended performance levels during the online estimation phase in the event of RN failures [30]. The RoC framework recognizes the reliability requirement of the IPSs and implements a system design method that can sustain situations of signal unavailability due to RN failure.

In the first process, the Robust-Maximum Summation of Max RSSI (R-MSMR) developed previously in our mathematical model in [17] is incorporated into the framework. Our system design uses a Binary Integer Linear Programming (BILP) approach, which employs the Simulated Annealing (SA) solution technique, to place a sufficient number of RNs in optimum locations. The design goal of this RN's placement is to achieve a maximizing summation of the maximum RSS at the signal test points (STPs) received from the RNs installed in the service area, as written in equation (1). We will show in Section 5 that our proposed system design can achieve the highest location accuracy during a normal situation and also yields reliable location estimated

results under unexpected situations such as RN failures, where S_{ij} refers to a binary $\{0, 1\}$ variable that equals 1 if the STP i is assigned to RN j , P_{ij} refers to the RSS that STP i receives from RN j , T refers to a set of signal test points, and B represents a set of candidate sites.

$$\text{Maximize } \sum_{i \in T} \max_{j \in B} (S_{ij} P_{ij}). \quad (1)$$

The operation of the R-MSMR design is divided into two phases as shown in Figure 1. Phase I is used to create a good starting solution that defines a sufficient number of RNs to be installed (N_s). In this phase, a set of signal test points (STPs) and a set of candidate sites for installing RNs are assigned as shown in Figure 2(a). An initial number of RNs, which is a minimum starting number of RNs to be installed, is calculated based on the dimension of the service area as shown in Figure 2(b). The initial location of the RNs will be determined by this starting number. Then, the process is repeated until the set of optimization constraints is satisfied. After that, a solution that can provide a sufficient number of RNs is obtained.

In Phase II, the SA approach is used to determine the optimal location of the RNs based on the number of RNs obtained from Phase I. In each iteration, the move operation (in terms of optimization procedure) is used to generate a new possible solution (called a neighbor solution), in which specific attributes of the current solution are modified as shown in Figure 3(a). The cost of each RN-placement solution as the evaluation function for the reliability of the RN-placement structure is calculated in this phase. Then, the process of Phase II continues until a stopping condition is reached. Finally, a robust RN placement solution for the IPSs is attained. Figure 3(b) depicts an example of the RN-placement solution obtained from the proposed R-MSMR design. More details and explanations about the operation of R-MSMR design can be found in [17].

3.3. Localization Process. In the second process, the two main operations of the IPSs based on location fingerprinting techniques are conducted: the offline calibration phase and the online estimation phase. The offline calibration and the online estimation phases are explained in detail as follows.

3.3.1. Offline Calibration Phase. In the offline calibration phase, the fingerprint database is collected by performing a site-survey of the physical environment information inside the service area. In this article, we use the Fusion location fingerprinting approach which has been developed from our previous work [31]. Unlike the traditional location fingerprinting approach, we utilize the combination of three physical environmental parameters to create the information of the fingerprint locations in the database. For each fingerprint location, we record three physical parameters simultaneously, which include the RSS vector, the temperature value, and the relative humidity value. To the best of our knowledge, there is no other research in the literature that utilizes temperature and humidity information in IPSs. After taking those measurements, the TEMP-level classification is

conducted, in which the temperature value (in degrees Celsius) and the relative humidity value (in percentage) are used to classify the similar physical environments into subareas or zones. After that, the TEMP level of each fingerprint location is recorded into the database [31]. Note that the information of each fingerprint location consists of the location coordinates (x, y) , the floor number (*floor*), the TEMP level, and the RSS vector received from the RNs.

Figure 4 illustrates an example of the proposed Fusion location fingerprinting approach in which information of the temperature and the relative humidity are used to classify each fingerprint location in the database. The different colors of the points on the map represent different TEMP levels that are obtained from the TEMP-level classification process. The TEMP-level classification is a classification algorithm based on behavioral observation information of the temperature and the relative humidity recorded from the actual environment. Therefore, the site-survey of the physical environmental information inside the actual service area needs to be calibrated before starting the offline calibration phase. Note that we define the lost RSS observation only during the offline calibration phase with the lowest sensitivity of the wireless transceiver (i.e., -110 dBm).

3.3.2. Online Estimation Phase. In this phase, the indoor positioning technique is used to estimate the target position inside the building. We present the online estimation technique that uses a distance based approach (i.e., Euclidean distance). The proposed positioning technique is called the Active Fusion technique, in which the procedure consists of four steps: the floor determination process, the TEMP-level filter, the Active Euclidean process, and the WKNN process. The descriptions of each online estimation step are as follows.

(1) Floor Determination Process. The first step of the online estimation phase is to determine the floor number on which the target node is located. We use the robust floor determination algorithm called the Robust Mean of Sum-RSS (RMoS) floor algorithm, which was developed in our previous works [32, 33]. The RMoS floor algorithm can accurately determine the floor on which the target nodes are located and can work under either the fault-free scenario or the RN-failure scenarios. The proposed floor algorithm is based on the mean of the summation of the strongest RSSs obtained. Then, the algorithm selects the floor number that has the highest value of the confidence intervals ($\Phi(\Lambda_f)$) as the floor where the target node is situated. Where Λ_f refers to a set of the summations of the strongest RSSs on the f^{th} floor. This can be written as

$$\text{floor} = \arg \max_f (\Phi(\Lambda_f)). \quad (2)$$

Figure 5(a) shows an example of the structure of the IPS in a three-story building under the RN-failure scenario. In this figure, one RN on the 1st floor and another RN on the 2nd floor became unavailable. Under this situation, the percentage of correct floor determination of other existing

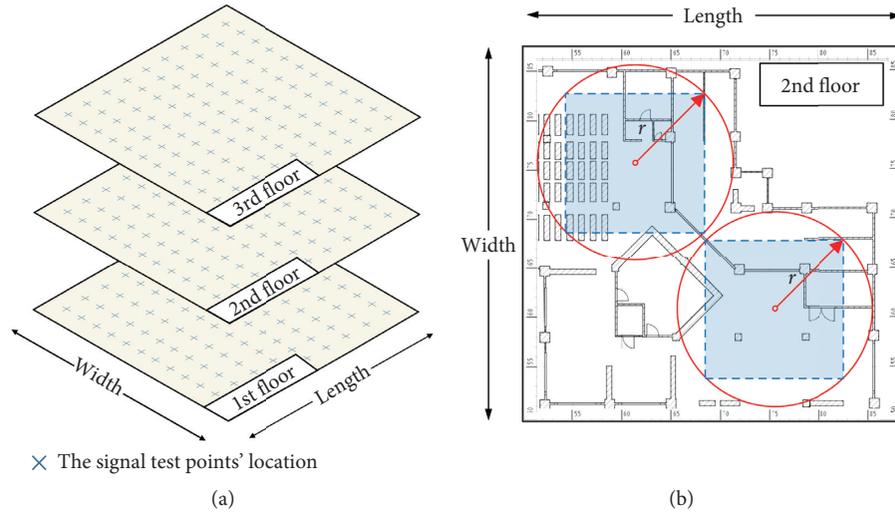


FIGURE 2: An example of R-MSMR design in Phase I. (a) A three-story structure with the location of the assigned STPs. (b) Example of coverage square estimating an RN's coverage area.

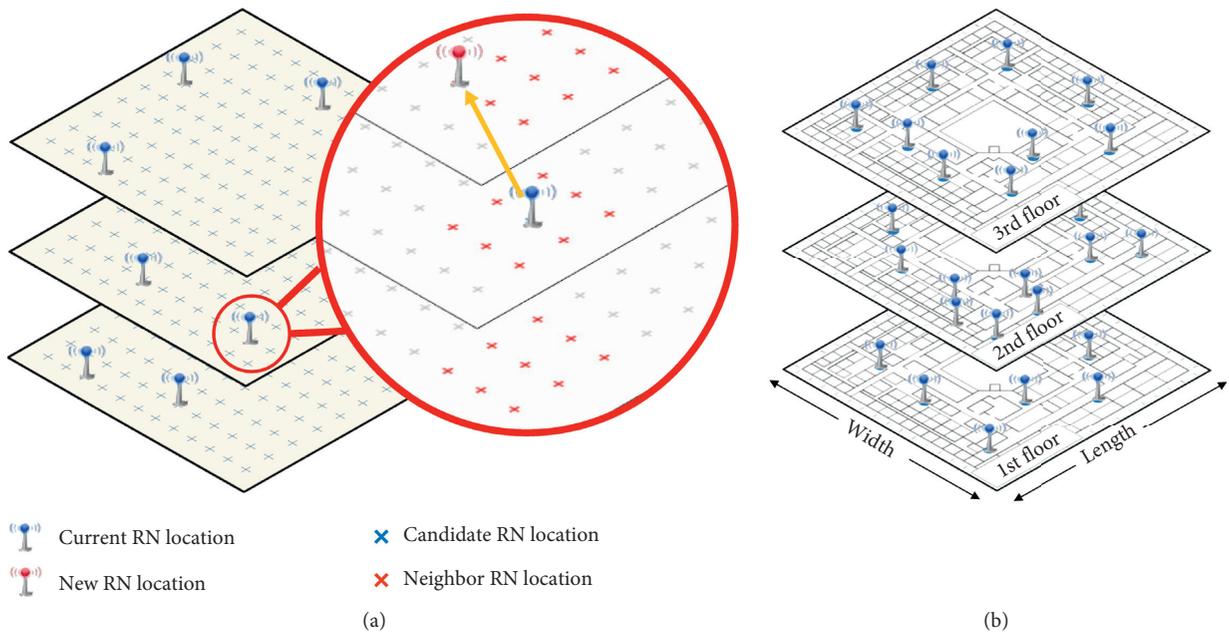


FIGURE 3: An example of R-MSMR design in Phase II. (a) Example of move operation in SA process. (b) An example result of R-MSMR placement design.

floor algorithms could drop from 20% to 50% when two-RNs fail in the system [32]. Unlike other existing floor algorithms, the RMoS floor algorithm can achieve reliable indoor multifloor positioning systems that can provide 100% correct floor determination under such unexpected situations. Figure 5(b) illustrates a 95% confidence interval for the mean of the RSS summations on each floor in a three-story building. It is clear that the confidence interval for the mean of the RSS summations on the actual floor where the target node was located is higher than that of the other floors (i.e., $\Phi(\Lambda_2)$). In this case, the proposed RMoS floor algorithm correctly reports that the target node is on the second

floor of the building. Detailed descriptions of the RMoS floor algorithm can be found in [33].

(2) *TEMP-Level Filter*. The main goal of this second step is to filter the fingerprint locations in the database by using the classification approach. The system will search and select the related fingerprint locations in the database that have the same TEMP level as those on the floor on which the target was located (i.e., the result from previous step). Then, those related fingerprint locations are used to estimate the target location in the next step. This means that the unrelated fingerprint locations are eliminated and are not considered

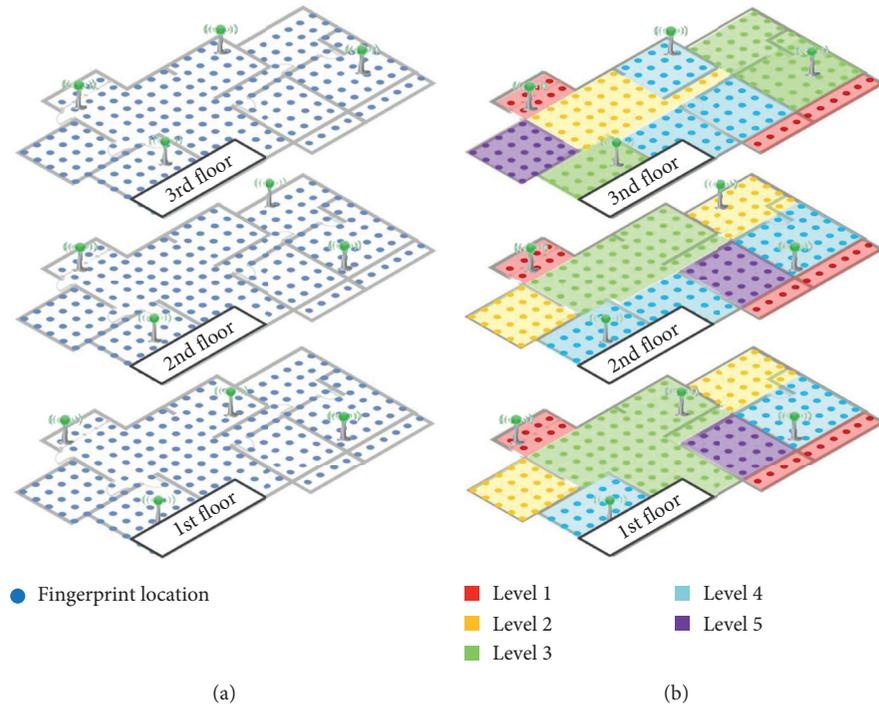


FIGURE 4: An example of Fusion location fingerprinting process. (a) Traditional location fingerprinting approach. (b) The environment parameter classification.

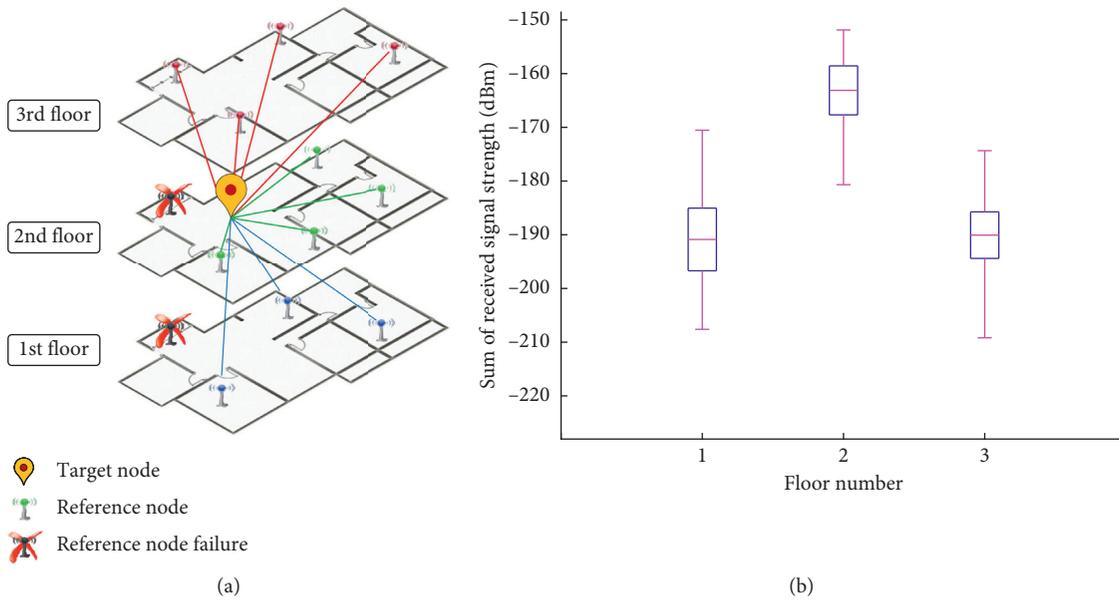


FIGURE 5: An example of the RMoS floor determination process. (a) Floor determination schematic diagram. (b) An example of the result of the RMoS floor determination algorithm.

in the location estimate. Note that our searching process is based on a binary search algorithm that has the time complexity of $O(\log n)$ [29]. Figure 6 illustrates an example of the TEMP-level filter based on the classification approach. We assume that the TEMP level of the target location is level 3 (i.e., green points) and the target’s floor result obtained from the previous process is the second floor. In this case, the

fingerprint locations in the database that are located on the second floor with the 3rd TEMP-level are selected, as shown by the green zone in Figure 6. It is clear that the search space of the location estimate is limited. Consequently, the computational time required during the online estimation phase is reduced. Moreover, this can help the IPSs to reduce the error distance in the location estimate [31].

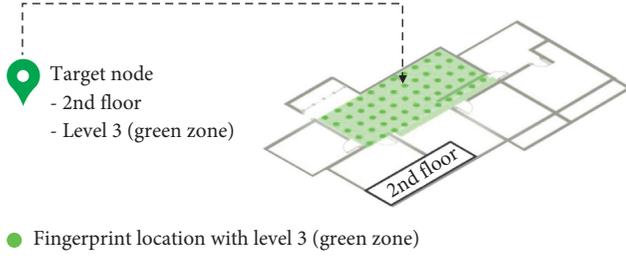


FIGURE 6: Example of the TEMP-level filter process that considers the fingerprint locations from level 3.

(3) *Active Euclidean Process*. Next, the third step of the online estimation phase is the distance metric calculation. The correlation of the RSS vector between fingerprint locations and the current location of the target are computed through the distance metric. The coordinate(s) associated with the fingerprint location that provides the smallest distance (e.g., Euclidean distance) is returned as the estimate of the position. In this article, we propose the novel active process for the Euclidean distance algorithm which is called the Active Euclidean distance. This enhanced Euclidean distance algorithm is developed to handle the problem of the online RSS being missing, either as a result of RN failures during the online estimation phase or of being outside a region of RN coverage. Figure 7 illustrates an example of the IPSs in a three-story building under a missing RSS situation. In this figure, one RN on the 1st floor and another RN on the 2nd floor became unavailable. Either hardware failures or software errors could be the reasons for this unavailability. Moreover, the location of the target is outside the RN coverage area, placed on 3rd floor, which is represented by the red dashed line in the figure. In this situation, the target node cannot observe the signal from those three RNs inside the building. Using the traditional Euclidean distance for matching the RSS pattern, the location accuracy of the system could drop by almost half [30]. The major difference between our proposed Active Euclidean distance technique and the traditional Euclidean distance technique is that we consider only the RSS values that are available (i.e., active RNs) when matching the RSS patterns with the location fingerprints in the database. The Active Euclidean distance for the distance metric calculation is written in the following equation:

$$d_i = \sqrt{\sum_{f=1}^F \left\{ \sum_{n=1}^{N_f} (r_{fn}^i - \ddot{s}_{fn})^2 \right\}}, \quad i = 1, 2, \dots, M, \quad (3)$$

where d_i is the Active Euclidean distance between the fingerprint location i th and the target location, note that $i = 1, 2, \dots, M$, M is the related fingerprint location that has the same TEMP level as the target located on the target's floor, F refers to the total number of floors in the multistory buildings, and N_f is the number of active RNs on the f th floor. The variable r_{fn}^i denotes the average RSS of the fingerprint location i th received from n th active RN on the f th floor. The variable \ddot{s}_{fn} represents the average RSS at the target location received from the n th active RN on the f th

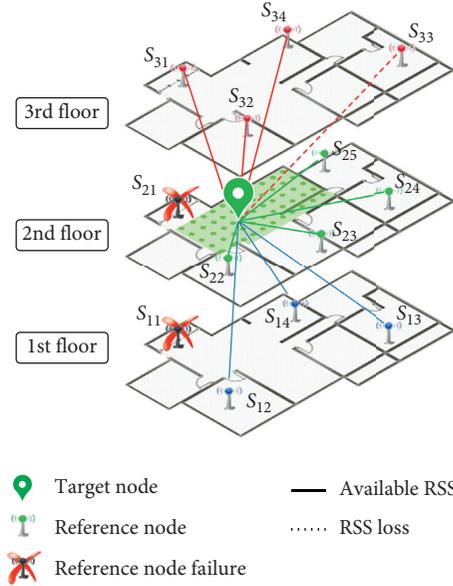


FIGURE 7: Example of the Active Euclidean distance technique.

floor. According to the example of an IPS in a three-story building as shown in Figure 7, the three, four, and three RNs on the 1st, 2nd, and 3rd floors are found during the online scanning, respectively (i.e. $F = 3$, $N_1 = 3$, $N_2 = 4$, $N_3 = 3$). The example of the Active Euclidean distance equation for the 9th fingerprint location (d_9) can be written as

$$\begin{aligned} f = 1, N_1 = 3; & \quad (r_{12}^9 - \ddot{s}_{12})^2 + (r_{13}^9 - \ddot{s}_{13})^2 + (r_{14}^9 - \ddot{s}_{14})^2, \\ f = 2, N_2 = 4; & \quad (r_{22}^9 - \ddot{s}_{22})^2 + (r_{23}^9 - \ddot{s}_{23})^2 + (r_{24}^9 - \ddot{s}_{24})^2 \\ & \quad + (r_{25}^9 - \ddot{s}_{25})^2, \\ f = 3, N_3 = 3; & \quad (r_{31}^9 - \ddot{s}_{31})^2 + (r_{32}^9 - \ddot{s}_{32})^2 + (r_{34}^9 - \ddot{s}_{34})^2. \end{aligned} \quad (4)$$

From this scenario, we do not include the three RNs of \ddot{s}_{11} , \ddot{s}_{21} , and \ddot{s}_{33} in the calculation, which represents a case of the RN failures and being outside a region of RN coverage. Then, the root squared error of RSS between the 9th fingerprinting location and the target location are computed. Finally, the Active Euclidean distance calculation for the 9th fingerprinting location is obtained.

(4) *WKNN Process*. The last step of our online estimation phase is the WKNN calculation. The WKNN algorithm is based on using the distance metric between the measured RSS of the target and the RSS of the related fingerprint locations to calculate the target location. In order to determine the coordinates of the target, the first k order of the related fingerprint locations that have the shortest distance metric (i.e., the shortest Active Euclidean distance) are selected. This closeness fingerprint location is called the nearest neighbor location. Then, those restriction distance metrics of the k -nearest neighbor's locations are used to compute the weighting factor [34]. The calculation of the WKNN method is written in the following equations:

$$x_o = \sum_{j=1}^k x_j \cdot w_j \quad \forall j \in k, \quad (5)$$

$$y_o = \sum_{j=1}^k y_j \cdot w_j \quad \forall j \in k, \quad (6)$$

$$w_j = \frac{(1/d_j^2)}{\sum_{j=1}^k (1/d_j^2)}, \quad (7)$$

where d_j is the smallest Active Euclidean distance j th, for $j = 1, 2, 3, \dots, k$ (i.e., we assign $k = 4$), the coordinate of target (x, y) is estimated by using (5) and (6), while the weighted value for the k -nearest neighbor's coordinates are computed by using (7), where x_o and y_o denote the x - and y -coordinates of location estimation, the variables x_j and y_j refer to the x - and y -coordinates of the nearest neighbor location j th, and w_j refers to the weighted value of the nearest neighbor j th. Finally, we obtain the solution of the target location in three-dimensional space represented by x, y , and $floor$. Figure 8 shows an example of the WKNN process.

4. Experimental Setups

To analyze several aspects of indoor positioning performance, we compared the performance results of the proposed IPSs with the traditional IPSs based on the location fingerprinting approach. Three of the essential performance metrics were used to evaluate indoor positioning performance. The three metrics were accuracy, robustness, and computational complexity. The core of this study can be divided into three objectives:

- (1) To compare the positioning performance of the IPSs based on the location fingerprinting technique in which the system is employed both using and without using the Active process (which will be discussed in Section 5.1)
- (2) To evaluate the performance of the IPSs under a normal situation (accuracy) and the RN-failure situation (robustness) (which will be discussed in Section 5.2)
- (3) To analyze the effect of the indoor positioning technique on computational complexity (which will be discussed in Section 5.3)

4.1. Experimental Settings. In our experimental study, two buildings with different floor structures and with different dimension areas are tested. The first building, labelled Building A, is an office building with dimensions of approximately 75 m (width) \times 75 m (length). Note that this is the same building used in [17]. The second building, labelled Building B, is a laboratory building with dimensions of approximately 30 m (width) \times 44 m (length). Figures 9 and 10 illustrate the floor layouts of Building A and Building B, respectively. The blue dots in the figures represent the RN-placement solution attained from the proposed

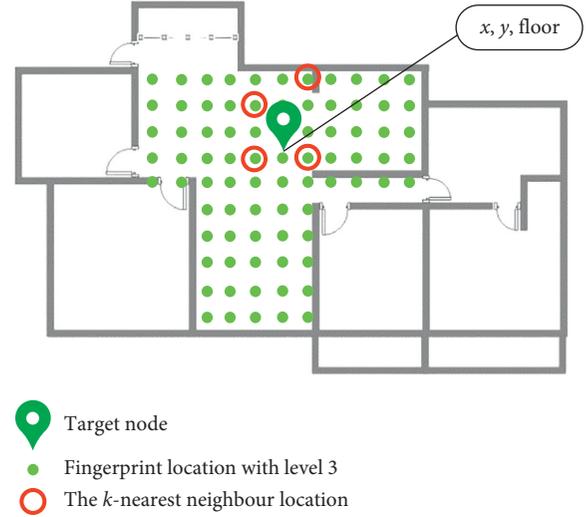


FIGURE 8: Example of the WKNN calculation.

R-MSMR with $R = 2$, in which 27 and 24 RNs are installed at Building A and Building B, respectively. We assigned the five TEMP levels for the proposed Active Fusion technique which is denoted by red, orange, green, blue, and purple for the 1st to 5th TEMP levels, respectively. Note that the white color area is open space in the building and the 1st TEMP level in the TEMP-level classification is the highest temperature and relative humidity value. We assigned the physical environmental classification inside Building A and Building B with three and five TEMP levels, respectively. Table 3 shows a range of temperature level used in this work. We set the grid spacing of the fingerprint locations at four meters for Building A and two meters for Building B. The number of fingerprint locations for Building A and Building B is 984 locations and 1,755 locations, respectively. A total number of 474 and 384 test points (i.e., target locations) were randomly selected for Building A and Building B, respectively. Note that these numbers of test points were obtained by determining the sample size with confidence intervals [35]. Table 4 provides the setting values of the parameters used in our experiments.

4.2. Experimental Equipment. The main hardware components of the IPSs used in this work are shown in Figure 11. They include the RNs, the target node, and the processing unit. The RNs are wireless transceivers that will send out signals upon request from the target node. Their locations are placed by using the system design approaches (discussed in Section 3.2). The target node will collect the RSS values that are transmitted from the RNs in the service area. It is connected to the SHT15 sensor for measuring the temperature and relative humidity value. This measured information (i.e., RSS and environment information) is passed to the processing unit via a UART-USB interface where the indoor positioning algorithm is carried out. The height of the RN is 2.5 meters, while the target node is 0.8 meters. Figure 12 illustrates the RNs and the target node used in this work. The hardware of both RNs and target node is based on the ZigBee evaluation kit of Freescale (now NXP

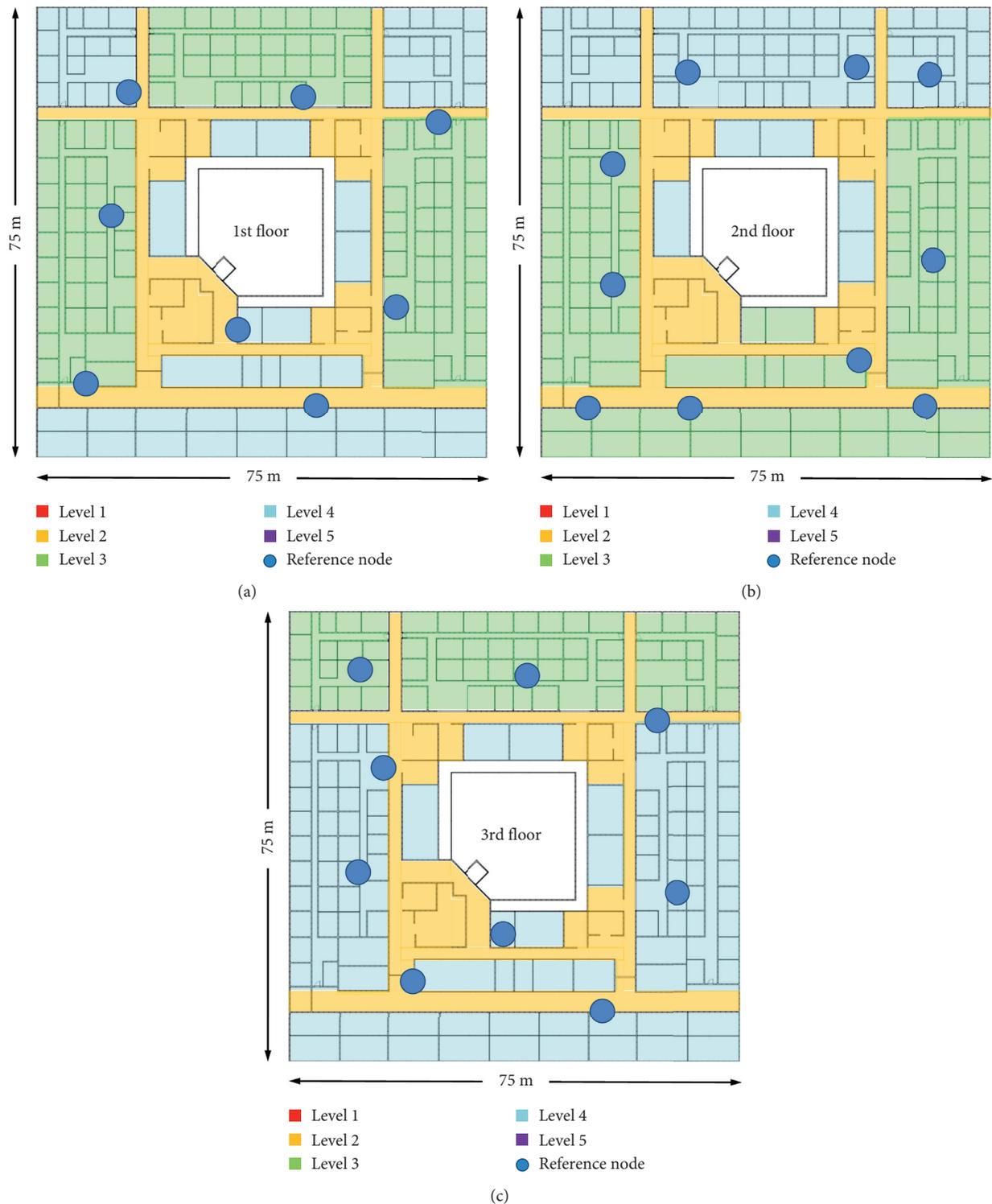


FIGURE 9: The TEMP-level classification for Building A. (a) 1st floor. (b) 2nd floor. (c) 3rd floor.

Semiconductors) called Freescale ZigBee® 1322xEVK. The kit contains ZigBee devices which are equipped with Freescale MC13224V platform-in-package (PiP) for the 2.4 GHz IEEE 802.15.4 standard. Each MC13224V chip has ARM7TDMI-S which is a 32-bit core microcontroller built in with the ZigBee transceiver. We purposely configured our IEEE 802.15.4 wireless transceivers to operate at 2.480 GHz

(i.e., channel 26 according to IEEE standard) [36]. This is to avoid or minimize the interference from Wi-Fi networks inside the buildings. The antennas used in all devices were a mix of the inverted F-shape antennas and external monopole omnidirectional antennas with SMA connectors. The transmit power of all transceivers was set at +3 dBm. Typical receiver's sensitivity of the devices is -95 dBm. With these



FIGURE 10: The TEMP-level classification for Building B. (a) 1st floor. (b) 2nd floor. (c) 3rd floor. (d) 4th floor. (e) 5th floor. (f) 6th floor.

TABLE 3: Range of the temperature level used in this work.

TEMP level	Color	Temperature (°C)	Humidity (%)
1st	Red	≥32	≥80
2nd	Orange	30-31	70-79
3rd	Green	26-29	60-69
4th	Blue	22-25	50-59
5th	Purple	≤21	≤49

setting and specification, the range of these wireless transceivers for indoor with non-line-of-sight is approximately 30 meters. In this work, the wireless transceivers were configured

to collect one RSSI sample at every three seconds (a.k.a. sampling rate). The limitation of our experiment is that we only gather the data for stationary node in this study. Table 5 summarizes the hardware specifications of the device used in this study.

5. Results and Discussion

In this section, major aspects of the indoor positioning performance are discussed. First, Section 5.1 provides a comparative performance evaluation of the IPSs using the

TABLE 4: Values of parameters used in the experiments.

Parameter	Values
Floor dimensions	75 m \times 75 m (for 1st to 3rd floor) for Building A
	30 m \times 44 m (for 1st to 6th floor) for Building B
RN placement	Uniform, MSMR, R-MSMR, $R=2$, and Phi-Uni (discussed in Section 5.1)
Grid spacing of fingerprint locations	4 m \times 4 m for Building A
	2 m \times 2 m for Building B
Number of test points (i.e., target locations)	474 locations for Building A
	384 locations for Building B
Floor determination technique	Group variance floor algorithm [37]
	RMoS floor algorithm [33]
	Enhanced Euclidean distance (E-Euclidean)
Positioning technique	Enhanced WKNN (E-WKNN)
	Active Euclidean distance (Ac-Euclidean)
	Active WKNN (Ac-WKNN)
TEMP level (for Active Fusion technique)	3 levels for Building A
	5 levels for Building B

active process versus those not using the active process. Next, in Section 5.2, the performance of the IPSs under a normal situation and an RN-failure situation is discussed. Finally, in Section 5.3, the computational complexity of different IPSs is analyzed.

5.1. Performance Evaluation on Active Process. In this section, we compare the performance results of the IPSs when using and when not using the active process. The three-story Building A was considered, which employs four different system designs as shown in Figure 13. We compared an average error distance of the proposed RN-placement design (i.e., R-MSMR with $R=2$) with three different system designs, which include the coverage and uniform placement (Uniform) model, the Maximize-Sum of Maximum RSS (MSMR) model [16], and the Phi-Uni [17]. Note that the designs of these four different RN placements for Building A were originally developed by our previous mathematical model in [17]. However, in our previous work [17], the performance comparison was analyzed between the IPSs under two different RN-failure patterns (i.e., similar and across RN-failure patterns). In this work, we compare the positioning performance of two distance based techniques, which consist of the Active Euclidean distance (i.e., the Euclidean distance technique using an active process) and the traditional Euclidean distance (i.e., the Euclidean distance technique not using an active process). Note that the location coordinates (x, y) and the floor number (*floor*) are computed simultaneously for both cases of the traditional Euclidean distance and Active Euclidean distance. The six cases of the RN-failure per floor situations are considered which consist of zero nodes to five nodes. A total number of 474 random locations were tested as the test points (each floor has 158 locations). The average error distance of each RN-failure situation is calculated from four random patterns of RN failure.

First, we consider the performance results of the IPSs with different system designs. Figure 14 reports an average error distance of the IPSs that deploy four different design structures in the three-story Building A. The solid lines represent the IPSs that use the Active Euclidean distance process, while the dash lines represent the IPSs that use traditional Euclidean distance process. From the results, it is clear that the average error distance of all IPSs increases in line with the increase in the number of RN failures. It is also clear from Figure 14 that the proposed R-MSMR design (the blue line with triangles with $R=2_{xx}$ legends) can provide greater location accuracy under both the fault-free scenario and the RN-failure scenario than the other system designs. Notice that these results follow the same trend as shown in our previous experiments in [17]. For example, when considering using the Active Euclidean distance process, the proposed R-MSMR with $R=2$ (the blue solid lines with triangles) has the lowest performance decreasing under the normal situation and the 3RN-failure per floor scenario about 8.4%. On the other hand, the other three different designs have more location accuracy degradation under the same situation. Our results show that they have performance degradation up to 77.3% for the Uniform design, 45.8% for the MSMR design, and 14.9% for the Phi-Uni design. In particular, when we compared the IPSs with R-MSMR with $R=2$ design and the Phi-Uni design that have same number of RNs installed, we found that the proposed R-MSMR with $R=2$ can achieve a better 51.2% fault tolerance than the Phi-Uni design under fifteen of the RNs in a three-story building failed. Note that this was 5RN-failure per floor scenario. The main reason is that the signal coverage availability of the proposed R-MSMR is guaranteed by at least the recommended number of signals of the accuracy index and the reliability index which is also one of the optimization constraints as described in [17]. Thus, the RN placement obtained from the proposed system design can effectively handle the problem of online RSS being missing caused by RN failures during the online estimation phase.

Next, we compare the location accuracy of the IPSs with and without using the active process as shown in Figure 14. The results report that the average error distance of the IPSs not using the active process could have very high error distances (i.e., the dashed lines). For example, under the 4RN-failure per floor scenario, the average error distance of the R-MSMR using the active process (the blue solid line with triangles) is 4.83 meters, whereas the R-MSMR not using the active process was over 20 meters (the blue dashed line with triangles). These results indicate that the location accuracy of the R-MSMR under the 4RN-failure per floor scenario could be dropped by almost 80% if the active process is not employed. The reason why the IPSs employing the Active Euclidean distance outperforms the traditional algorithm is that the proposed Active Euclidean distance is developed to handle the problem of online RSS being missing. Unlike the traditional Euclidean distance which includes all RNs in the localization area in its closest pattern matching with the location fingerprints in the database, the proposed Active Euclidean distance only focuses on the RSS values obtained from available RNs for the RSS matching

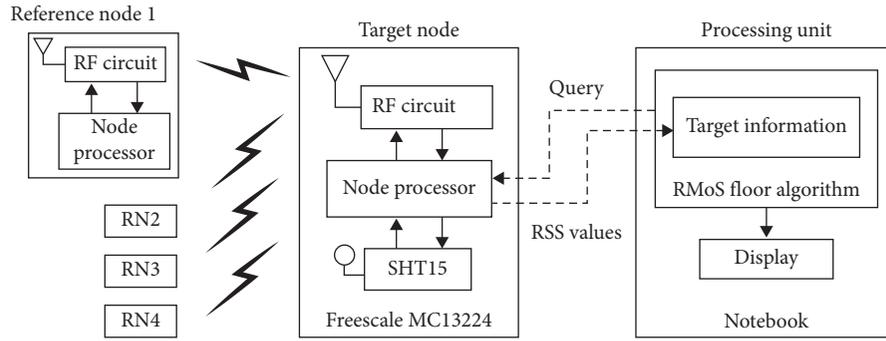


FIGURE 11: Block diagram of the main hardware components in the indoor positioning systems.



FIGURE 12: Experimental equipment used in this work. (a) Reference node (RN). (b) Target node on a cart with laptop PC.

TABLE 5: Hardware specifications of the IEEE 802.15.4 devices.

Specification	Details
Manufacturer	Freescale (now NXP Semiconductors)
Chipset	MC13224V
Frequency range	2.405 GHz to 2.480 GHz
Operating channel	CH 26 (2.480 GHz) according to IEEE 802.15.4
Rx sensitivity	-95 dBm
Transmit power	+3 dBm
Antenna	Inverted-F antenna or external omnidirectional antenna with SMA connector

calculation. Hence, the lost RSS observation during the online estimation phase, which was caused either by RN failures during the online estimation phase or by being outside a region of RN coverage, is handled. Note that under the fault-free scenario, all IPSs using and not using the active

process have an equal average error distance. This means that under a normal situation, the positioning performance of the IPSs using and not using the active process is the same.

Observation from the above result indicates that using the Active Euclidean distance process produces better performance than using the traditional Euclidean distance process, in that the proposed active process can provide greater robustness for the IPSs under the RN-failure situation. Besides the active process, the system design that can satisfy the IPS design requirements for a particular design scenario such as localization that supports rescuers in an earthquake scenario, as the proposed R-MSMR design is also very important for the IPSs that requires system reliability.

5.2. Performance Evaluation under Normal Situation and RN Failure Situation. In this section, we analyze the positioning performance of the IPSs under the fault-free and the RN-

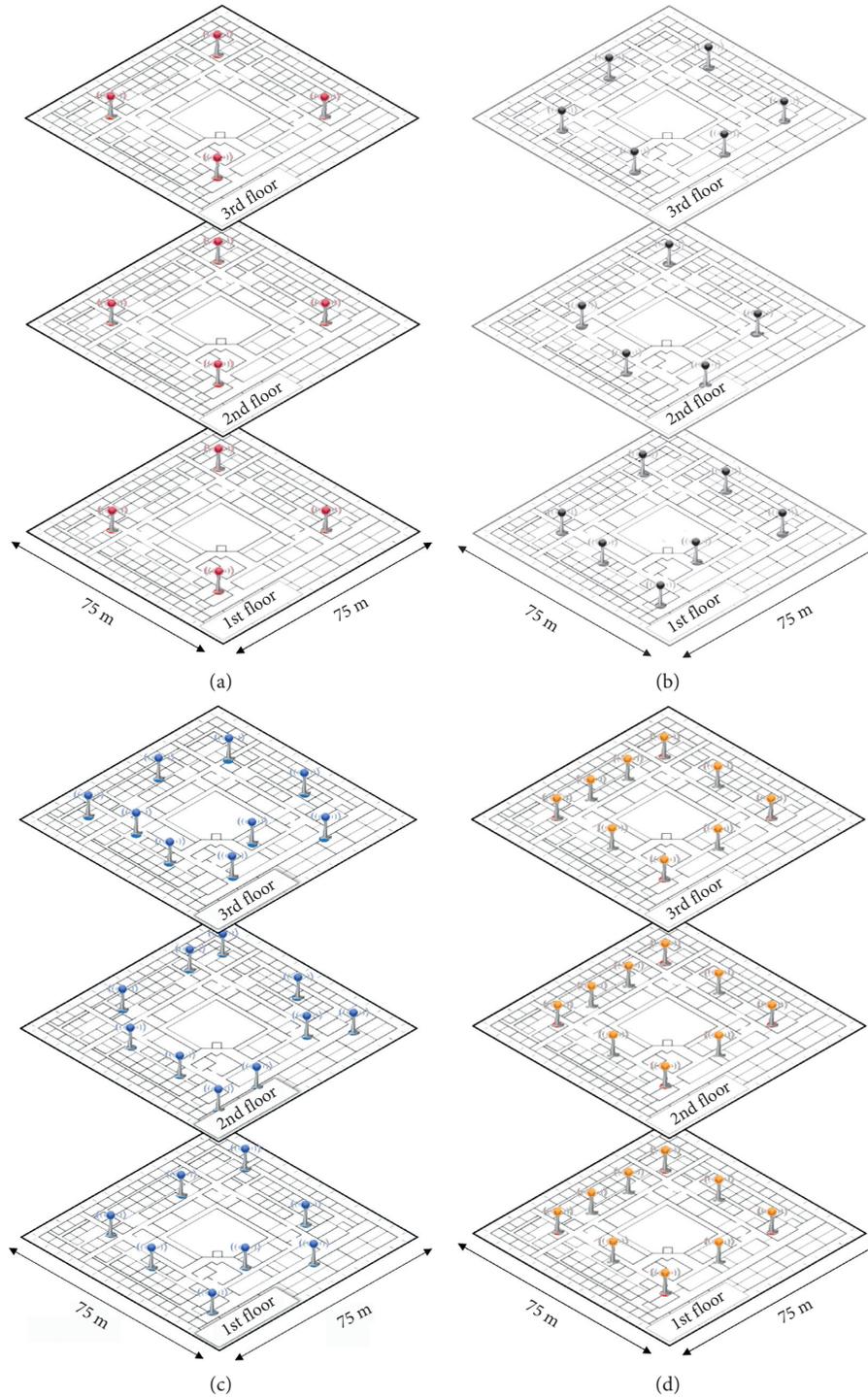


FIGURE 13: The RN placement designed for a three-floor service area from [17]. (a) Uniform (12 nodes). (b) MSMR (18 nodes). (c) R-MSMR, $R=2$ (27 nodes). (d) Phi-Uni (27 nodes).

failure scenarios, in which the key performance studies can be divided into two parts: floor determination and location estimation.

5.2.1. Percentage of Correct Floor Determination. In this section, we evaluate the floor determination performance in the three-story Building A. We compare the performance of two different floor determination algorithms without the use

of the location fingerprinting database: the Group variance floor algorithm [37] and the proposed RMoS floor algorithm [33]. Two system design structures that have the same number of RNs installed are considered, consisting of the proposed R-MSMR with $R=2$ and the Phi-Uni as shown in Figures 13(c) and 13(d), respectively. As mentioned in Section 5.1, the six cases of the RN-failure per floor situations and a total number of 474 test points are tested.

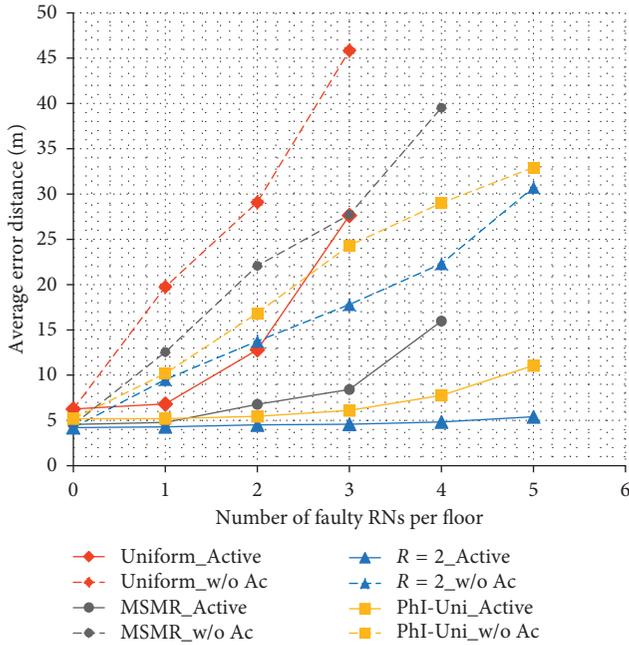


FIGURE 14: Effect of different numbers of RN-failures per floor on average error distance in Building A.

Figure 15 reports the percentage of correct floor determination of the Group variance algorithm and RMoS floor algorithm. We observe that both floor determination algorithms achieve 100% correct floor determination in the fault-free scenario (i.e., all RNs worked properly) for all cases of system designs. However, when some RNs in the system have failed, only the proposed RMoS floor algorithm with the R-MSMR $R = 2$ design (the pink solid line with triangles) can provide fault tolerance that is better than the other IPs. It yields a 100% correct floor determination performance under a 4RN-failure per floor scenario (i.e., the case of 40% of RN failures in the system). Unlike the proposed combination of the RMoS floor algorithm with the R-MSMR, $R = 2$ design, the other IPs such as the Group variance algorithm with the R-MSMR $R = 2$ design (the gray solid line with triangles) could not tolerate the RN-failure scenarios. They fail to achieve 100% correct floor determination performance in all RN-failure scenario cases. The reason is that the missing RSS during the online estimation phase could affect the floor point process of the Group variance algorithm. For example, when considering the floor points which were calculated from one of three online statistical parameters as the availability [37], an error of the availability score occurred when some RNs during the online estimation phase have failed. This is the reason why the floor points are incorrect and cause the Group variance algorithm to determine the wrong floor.

Moreover, we observe that the performance of the proposed RMoS algorithm with R-MSMR, $R = 2$ will drop by almost 0.5% when more than 50% of RNs have failed in the service area (i.e., 5RN-failure per floor scenario). The reason is that the RMoS floor algorithm considers only which 50% of RNs on each floor give the strongest RSS values to be suitable for the RSS summations under RN failure [33].

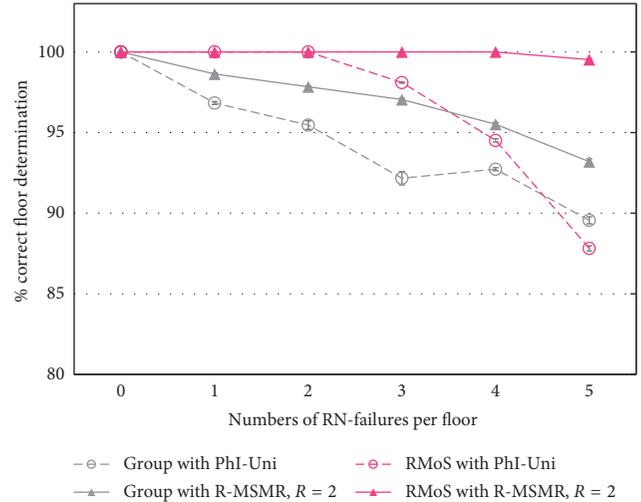


FIGURE 15: Percentage of correct floor determination in different numbers of RN failures per floor at Building A.

Thus, having more than 50% of the RNs in the building fail may directly affect the RSS summations process of the RMoS algorithm and lead to incorrect floor determination. This is a limitation of the proposed RMoS floor algorithm, that is, it cannot support the worst-case scenarios of RN failure, when more than 50% of the RNs installed in the service area fail.

It is clear from the above discussion that the combination of the RMoS floor algorithm with the R-MSMR, $R = 2$ design outperforms the other IPs and provides more robust IPS performance, in which the proposed combination algorithm can achieve a 100% correct floor determination performance when the system encountered 40% RN failure in the system. Once again, the structure of the proposed R-MSMR design can achieve higher fault tolerance than other structures that have same number of RNs installed, such as the PhI-Uni. In the next section, to make a clear comparison of several aspects of indoor positioning performance, we will only use the most robust system design as the R-MSMR with $R = 2$ structure for the performance evaluation.

5.2.2. Estimating Location. The performance of the IPS in terms of the accuracy of the location estimation (x, y) under the fault-free scenario and the RN-failure scenarios is evaluated in this section. Building A and Building B are tested, in which the experimental setups of both buildings were described in Section 4.1. Five indoor positioning techniques are used, which include Enhanced Euclidean distance (E-Euclidean), Enhanced WKNN (E-WKNN), Active Euclidean distance (Ac-Euclidean), Active WKNN (Ac-WKNN), and the proposed Active Fusion technique. Note that only in the case of the proposed Active Fusion technique, the floor number was computed by the RMoS floor algorithm, while in the case of the other four positioning techniques, the floor number was computed simultaneously with location coordinates of x, y . A total number of 474 test points and 384 test points were randomly selected for Building A and Building B, respectively. For the case of the RN-failure scenarios, we created four patterns of

RN failure scenarios by randomly turning off 30% of the RNs in the system. Then, the average error distance of these four scenarios were computed.

Figures 16 and 17 report the average error distance of the IPSs that deploy the R-MSMR, $R=2$ design for Building A and Building B, respectively. The results in the case of a large building, such as Building A, indicate that the five indoor positioning techniques have a similar average accuracy performance under the fault-free scenario of between 3.83 and 4.22 meters. However, under the RN-failure scenarios, in which 30% of the RNs in the system failed, the accuracy performance of the two indoor positioning techniques that did not use the active process dropped by almost 77% (i.e., 76.3% for E-Euclidean and 75.3% for E-WKNN). Unlike the techniques that did not use the active process, the three indoor positioning techniques that did use the active process produced a percentage of difference between their average error distance in the normal situation and the 30% RN-failure scenario of less than 8%. In particular, the proposed RoC framework (i.e., Active Fusion technique) outperformed the other techniques, as the performance of the Active Fusion dropped by less than 4.5%.

Similar results were obtained for the case of the small building (i.e., Building B) as shown in Figure 17. The five indoor positioning techniques had a similar average accuracy performance under the fault-free scenario (i.e., between 3.36 and 3.72 meters). Once again, only the three indoor positioning techniques that used the active process provided fault tolerance under the 30% RN-failure scenario, in which their accuracy performance dropped by less than 5%, while the two techniques that did not use the active process in Building B also reported performance degradation of up to 75% (i.e., 75.1% for E-Euclidean and 74.2% for E-WKNN). It is a major advantage of the proposed active process that it provides a reliable location for the IPSs. It can handle the problem of online RSS being missing caused either by RN failures during the online estimation phase or by being outside a region of RN coverage.

Moreover, we observed that the small building (Building B) produced a higher accuracy performance than the large building (Building A) for all scenarios. The reason is that Building B employs the fingerprinting granularity with 2×2 meters, which has higher resolution of grid spacing than Building A. On the other hand, assigning high fingerprinting granularity for the IPSs based on the location fingerprinting approach can improve the accuracy and precision of the performance [38]. However, the high fingerprinting granularity will be very time-consuming in performing an exhaustive fingerprint collection during the offline calibration phase. This can result in weeks spent on surveying the site and collecting data for a large service area. This is in fact a trade-off between the time consumed and the accuracy of the performance of the IPSs based on the location fingerprinting approach [2].

5.3. Performance Evaluation of Computational Complexity. Based on the experimental results described above, it does not show any tangible difference in those three indoor

positioning techniques that used the active process. Thus, in this section, we investigate another essential performance of the IPSs that explains the location processing time during the online estimation phase as computational complexity. We recorded the computational time of those three indoor positioning techniques at Building A and Building B. The number of fingerprint locations for Building A and Building B is 984 locations and 1,755 locations, respectively. A total number of 474 test points and 384 test points were used for Building A and Building B, respectively. For each test point, the location calculation was run twenty times. Their average computational time was computed and reported as a histogram of computational times, as shown in Figures 18 and 19 for Building A and Building B, respectively.

Consider the results shown for Building A in Figure 18. We can see that the proposed Active Fusion resulted in shorter computational times compared to the other techniques. The average computational time of the Active Fusion was 17.8 milliseconds per location, represented by blue bins, while the Ac-Euclidean and the Ac-WKNN were 24.1 and 25.5 milliseconds per location represented by orange bins and purple bins, respectively. Moreover, when we consider Building B, which has a 40% larger number of fingerprint locations than Building A, we found that the computational time of those two techniques that did not use the Fusion location fingerprinting approach (Ac-Euclidean and Ac-WKNN) grew larger as the number of fingerprint locations increased. Their average computational times increased by almost 40% (i.e., 34.7% for Ac-Euclidean and 37.9% for Ac-WKNN, respectively). Unlike the techniques that did not use the Fusion location fingerprinting approach, the proposed Active Fusion requires less computational time than the other two techniques while still maintaining sufficiently accurate performance. In particular, in the 40% larger search space at Building B, the average computational time of the Active Fusion process increased by 2.7%. From the results considered here, it can be concluded that the proposed RoC framework requires a lower computational time than the Ac-Euclidean and the Ac-WKNN by up to 53.8% and 55.1%, respectively. Note that the computational time results of the Euclidean approaches of the Ac-Euclidean and the E-Euclidean were equal. Similar results were obtained for the case of the Ac-WKNN and the E-WKNN.

The reason that the computational time of the Active Fusion not growing as the size of the search space increased is that our proposed technique deploys the Fusion location fingerprinting approach with a TEMP-level filter based on the classification approach to classify the fingerprint locations before the Euclidean distance process is conducted. Instead of computing all of the huge fingerprint locations in the database as is required with the traditional techniques (e.g., Ac-Euclidean and Ac-WKNN), only the related fingerprint locations that have the same TEMP level as the target are considered. This means that the search space of the system (i.e., the fingerprint locations considered in the Euclidean distance process) is limited. This results in a low computational time during the online estimation phase and also reduces the errors in terms of the location estimation. Thus, an increase of 40% in the search space at Building B

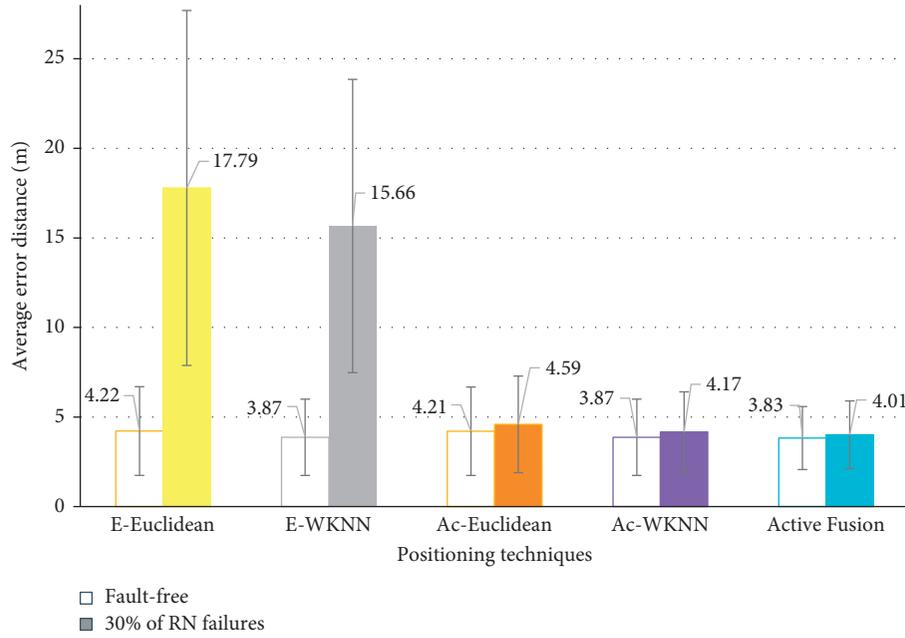


FIGURE 16: An average error distance of wireless indoor positioning systems that deploy R-MSMR, $R=2$ at Building A.

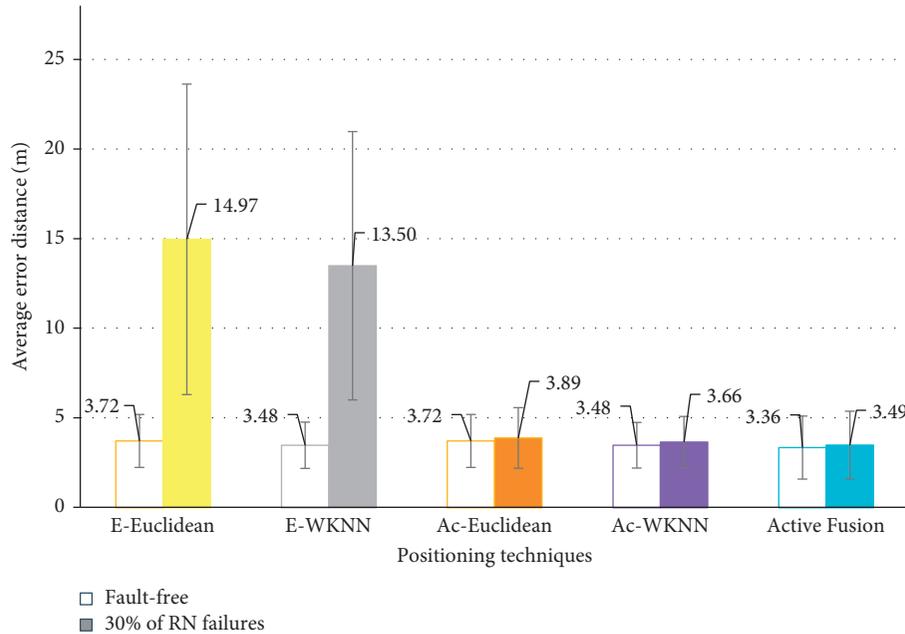


FIGURE 17: An average error distance of wireless indoor positioning systems that deploy R-MSMR, $R=2$ at Building B.

does not significantly affect the computational time of the proposed Active Fusion process, in which less than 4% of all those fingerprint locations need to be calculated in the Euclidean distance process (i.e., they are filtered by floor number and the Temp-level).

Moreover, when we consider the evaluating run-time complexity as a Big O notation which is used to predict the upper bound of the growth rate of the algorithm as the input size grows [39], we found that the worst-case complexity of the proposed TEMP-level filter based on a binary search algorithm was lower than other existing classification approaches. Our algorithm has the worst-case complexity of

$O(\log n)$ [29], while other existing classification approaches such as the k-means clustering and the FCM have the worst-case complexity of $O(ndc)$ and $O(ndc^2)$ [28], respectively, where n refers to the number of data points, d represents the number of dimensions, and c represents the number of clusters. This evaluation could also explain why an increase in the search space does not significantly affect the online computational time of the matching RSS process of the proposed RoC framework. Thus, we can conclude that the IPSs employing the RoC framework are robust and low-complexity wireless IPSs based on the fingerprinting approach in the RN-failure scenarios considered in our study.

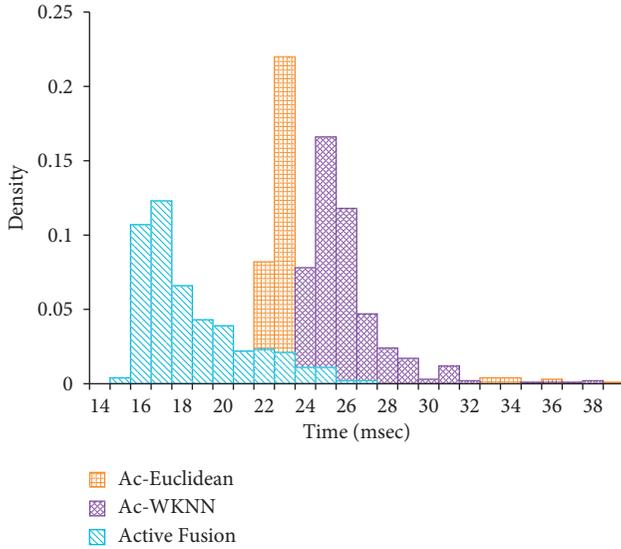


FIGURE 18: Histogram of the computational time at Building A.

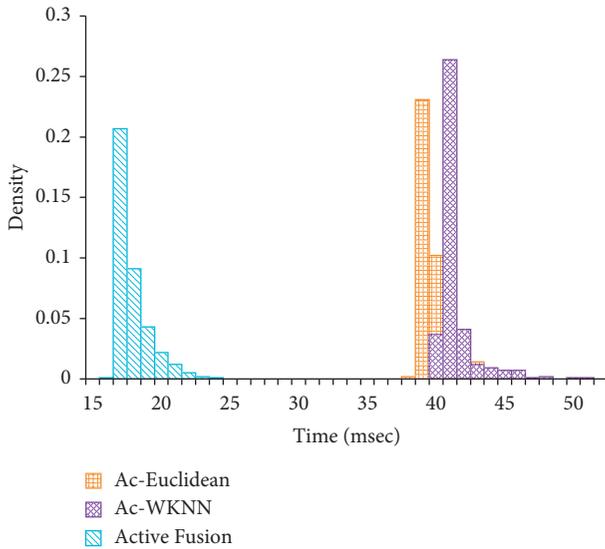


FIGURE 19: Histogram of the computational time at Building B.

6. Conclusion

In this article, we presented an integrated framework for the wireless IPSs based on location fingerprinting techniques that can be applied to variety of indoor scenarios ranging from single-floor to multiple-floor environments. The proposed framework is called the Robust and low Complexity indoor positioning systems framework (RoC framework). This framework consists of two essential indoor positioning processes: the system design process and the localization process. Experimental results reveal that the RoC framework can achieve robustness in terms of the system design structure, in which it was able to achieve the best positioning performance in either *fault-free* or *RN-failure scenarios*. Moreover, in the online estimation phase, the proposed framework can provide target location reliability in the RN-failure scenarios, in which it was able to

attain the highest correct floor determination and the highest location accuracy compared to the other techniques. Furthermore, our proposed RoC framework also yields the lowest computational complexity in online searching time without compromising the positioning accuracy. This can be achieved by exploiting additional environmental information via temperature and relative humidity sensor devices which can help reduce the search space of location fingerprinting. However, this improvement comes at the cost of additional data collection, storage, and filtering process. Specifically, when we compared it to the traditional WKNN under the 30% RN-failure scenario at Building B, the proposed RoC framework demonstrates a better location accuracy than the WKNN by up to 74.1% and yields a lower computational time than the WKNN by about 55.1%.

Our future works will consider system design guidelines which can suggest how and which directions the system designer should take for performance improvement of the Multifloor positioning system. Additionally, the enhanced design framework should still be robust both during the normal situation and when some RNs have failed. These guidelines can also be applied in various other wireless technologies such as Bluetooth Low Energy (BLE) for indoor positioning systems.

Data Availability

The performance comparison results data were used to support the findings of this study are available from the corresponding author upon request.

Disclosure

An earlier version was circulated under the title “Robust and low complexity wireless indoor positioning systems for multi-floor buildings using fingerprinting techniques.” This article was extended from the fourth and fifth chapter of our Ph.D. dissertation.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

The authors would like to express their deepest appreciation to the late Dr. Chutima Prommak, Assistant Professor of the School of Telecommunication Engineering at Suranaree University of Technology, Thailand, whose contribution to this work was of great significance. Without her experience, guidance, and patience, it would have never been possible to complete this work successfully.

References

- [1] S. Yiu, M. Dashti, H. Claussen, and F. Perez-Cruz, “Wi-Fi fingerprint-based indoor positioning: recent advances and comparisons,” *IEEE Communications Surveys & Tutorials*, vol. 18, no. 1, pp. 466–490, 2015.

- [2] K. Pahlavan and P. Krishnamurthy, *Principles of Wireless Access and Localization*, Wiley, Hoboken, NJ, USA, 2013.
- [3] Z. Farid, R. Nordin, and M. Ismail, "Recent advances in wireless indoor localization techniques and system," *Mobile Information Systems*, vol. 2013, Article ID 185138, 12 pages, 2013.
- [4] D. Dardari, P. Closas, and P. M. Djuric, "Indoor tracking: theory, methods, and technologies," *IEEE Transactions on Vehicular Technology*, vol. 64, no. 4, pp. 1263–1278, 2015.
- [5] H. Liu, H. Darabi, P. Banerjee, and J. Liu, "Survey of wireless indoor positioning techniques and systems," *IEEE Transactions on Systems, Man, and Cybernetics, Part C*, vol. 37, no. 6, pp. 1068–1080, 2007.
- [6] Z. Turgut, G. Z. G. Aydin, and A. Sertbas, "Indoor localization techniques for smart building environment," *Procedia Computer Science*, vol. 83, pp. 1176–1181, 2016.
- [7] S. Yiu, M. Dashti, H. Claussen, and F. Perez-Cruz, "Wireless RSSI fingerprinting localization," *Signal Processing*, vol. 131, pp. 253–244, 2017.
- [8] R. F. Brena, J. P. García-Vázquez, C. E. Galván-Tejada, D. Muñoz-Rodríguez, C. Vargas-Rosales, and J. Fangmeyer, "Evolution of indoor positioning technologies: a survey," *Mobile Information Systems*, vol. 2017, Article ID 2630413, 21 pages, 2017.
- [9] O. Baala, Y. Zheng, and A. Caminada, "The impact of AP placement in WLAN-based indoor positioning system," in *Proceedings of ICN'09 Eighth International Conference on Networks*, pp. 12–17, Washington, DC, USA, May 2009.
- [10] A. E. C. Redondi and E. Amaldi, "Optimizing the placement of anchor nodes in RSS-based indoor localization systems," in *Proceedings of 12th Annual Mediterranean Ad Hoc Networking Workshop (MED-HOC-NET)*, pp. 8–13, Ajaccio, France, March 2013.
- [11] S. Merkel, P. Unger, and H. Schmeck, "Evolutionary algorithm for optimal anchor node placement to localize devices in a mobile ad hoc network during building evacuation," in *Proceedings of Genetic and Evolutionary Computation Conference (GECCO'13)*, pp. 1407–1414, Amsterdam, Netherlands, July 2013.
- [12] K. Tong, X. Wang, and A. Khabbazibasmenj, "Optimum reference node deployment for TOA-based localization," in *Proceedings of International Conference on Communications (ICC)*, pp. 3252–3256, Sydney, Australia, September 2015.
- [13] C. Sharma, Y. F. Wong, W.-S. Soh, and W.-C. Wong, "Access point placement for fingerprint-based localization," in *Proceedings of IEEE International Conference on Communication Systems (ICCS)*, pp. 238–243, Singapore, November 2010.
- [14] S.-H. Fang and T.-N. Lin, "A novel access point placement approach for WLAN-based location systems," in *Proceedings of IEEE Wireless Communications and Networking Conference (WCNC)*, pp. 1–4, Sydney, Australia, July 2010.
- [15] Q. Chen, B. Wang, X. Deng, Y. Mo, and L. T. Yang, "Placement of access points for indoor wireless coverage and fingerprint-based localization," in *Proceedings of 10th International Conference on High Performance Computing and Communications*, pp. 2253–2257, Zhangjiajie, China, June 2014.
- [16] K. Kondee, S. Aomumpai, and C. Prommak, "A novel technique for reference node placement in wireless indoor positioning systems based on fingerprint technique," *ECTI Transactions on Computer and Information Technology (ECTI-CIT)*, vol. 9, no. 2, pp. 131–141, 2015.
- [17] K. Maneerat and K. Kaemarungsi, "Robust system design using BILP for wireless indoor positioning systems," *Mobile Information Systems*, vol. 2018, Article ID 4198504, 19 pages, 2018.
- [18] P. Gupta, S. Bharadwaj, S. Ramakrishnan, and J. Balakrishnan, "Robust floor determination for indoor positioning," in *Proceedings of the 20th National Conference on Communications (NCC)*, pp. 1–6, Kanpur, India, March 2014.
- [19] Y. C. Lee and S. H. Park, "Localization method for mobile robots moving on stairs in multi-floor environments," in *Proceedings of the IEEE International Conference on Systems, Man and Cybernetics (SMC)*, pp. 4014–4020, San Diego, CA, USA, December 2014.
- [20] F. He, C. Wu, X. Zhou, and Y. Zhao, "Robust and fast similarity search for fingerprint calibrations-free indoor localization," in *Proceedings of International Conference on Big Data Computing and Communications (BIGCOM)*, pp. 320–327, Chengdu, China, November 2017.
- [21] A. Zayets and E. Steinbach, "Robust Wi-Fi-based indoor localization using multipath component analysis," in *Proceedings of International Conference on Indoor Positioning and Indoor Navigation (IPIN)*, pp. 1–8, Sapporo, Japan, November 2017.
- [22] D. Taniuchi and T. Maekawa, "Robust Wi-Fi based indoor positioning with ensemble learning," in *Proceedings of International Conference on Wireless and Mobile Computing, Networking and Communications (WiMob)*, pp. 592–597, Larnaca, Cyprus, November 2014.
- [23] T. Guan, W. Dong, D. Koutsonikolas, G. Challen, and C. Qiao, "Robust, cost-effective and scalable localization in large indoor areas," in *Proceedings of Global Communications Conference (GLOBECOM)*, pp. 1–6, Washington, DC, USA, February 2016.
- [24] W. Xue, X. Hua, Q. Li, K. Yu, and W. Qiu, "Improved neighboring reference points selection method for Wi-Fi based indoor localization," *IEEE Sensors Letters*, vol. 2, no. 2, pp. 1–4, 2018.
- [25] H. Zhou and N. N. Van, "Indoor fingerprint localization based on fuzzy c-means clustering," in *Proceedings of International Conference on Measuring Technology and Mechatronics Automation*, pp. 337–340, Zhangjiajie, China, April 2014.
- [26] C. W. Lee, T. Lin, S. H. Fang, and Y. C. Chou, "A novel clustering-based approach of indoor location fingerprinting," in *Proceedings of Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, pp. 3191–3196, London, UK, 2013.
- [27] A. Saha and P. Sadhukhan, "A novel clustering strategy for fingerprinting-based localization system to reduce the searching time," in *Proceedings of International Conference on Recent Trends in Information Systems (ReTIS)*, pp. 538–543, Kolkata, India, September 2015.
- [28] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*, Cambridge University Press, Cambridge, UK, 2008.
- [29] D. E. Knuth, *The Art of Computer Programming: Volume 3: Sorting and Searching*, Addison-Wesley Professional, Boston, MA, USA, 1998.
- [30] C. Khauphung, P. Keeratiwintakorn, and K. Kaemarungsi, "On robustness of centralized-based location determination using WSN," in *Proceedings of 14th Asia-Pacific Conference on Communications (APCC)*, pp. 1–5, Tokyo, Japan, February 2009.
- [31] K. Maneerat and C. Prommak, "Low complexity wireless indoor positioning approaches based on fingerprinting

- techniques,” in *Proceedings of International Telecommunication Networks and Applications Conference (ITNAC)*, pp. 244–249, Sydney, Australia, November 2015.
- [32] K. Maneerat and C. Prommak, “Floor determination algorithm with node failure consideration for indoor positioning systems,” in *Proceedings of the 8th International Conference on Signal Processing Systems*, pp. 203–207, Auckland, New Zealand, January 2016.
- [33] K. Maneerat, K. Kaemarungsi, and C. Prommak, “Robust floor determination algorithm for indoor wireless localization systems under reference node failure,” *Mobile Information Systems*, vol. 2016, Article ID 4961565, 12 pages, 2016.
- [34] Q. Wang, Y. Feng, X. Zhang, Y. Sun, and X. Lu, “TWKNN: an effective Bluetooth positioning method based on isomap and WKNN,” *Mobile Information Systems*, vol. 2016, Article ID 8765874, 11 pages, 2016.
- [35] R. Jain, *The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation, and Modeling*, Wiley, Hoboken, NJ, USA, 1st edition, 1991.
- [36] NXP Semiconductors, *MC13224V: 2.4 GHz 802.15.4 RF and 32-bit ARM7™ MCU with 128KB Flash, 96KB RAM*, NXP Semiconductors, Eindhoven, Netherlands, 2010, <http://www.nxp.com/products/wireless-connectivity/2.4-ghz-wireless-solutions/2.4-ghz-802.15.4-rf-and-32-bit-arm7-mcu-with-128kb-flash-96kb-ram:MC13224V>.
- [37] F. Alsehly, T. Arslan, and Z. Sevak, “Indoor positioning with floor determination in multi story buildings,” in *Proceedings of International Conference on Indoor Positioning and Indoor Navigation*, pp. 1–7, Guimarães, Portugal, November 2011.
- [38] K. Kaemarungsi, “Efficient design of indoor positioning systems based on location fingerprinting,” in *Proceedings of Wireless Networks Communications and Mobile Computing*, pp. 181–186, Maui, HI, USA, December 2005.
- [39] R. K. Ahuja, T. L. Magnanti, and J. B. Orlin, *Network Flows Theory Algorithms and Applications*, Pearson, London, UK, 1993.

Research Article

Research on Precision Marketing Model of Tourism Industry Based on User's Mobile Behavior Trajectory

Jialin Zhang,^{1,2} Tong Wu,¹ and Zhipeng Fan ^{1,2}

¹Harbin University of Commerce, Harbin 150028, China

²Heilongjiang Provincial Key Laboratory of Electronic Commerce and Information Processing, Computer and Information Engineering College, China

Correspondence should be addressed to Zhipeng Fan; hsdfzp@126.com

Received 28 September 2018; Revised 4 December 2018; Accepted 13 December 2018; Published 3 February 2019

Guest Editor: Jaegeol Yim

Copyright © 2019 Jialin Zhang et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

With the deep cross-border integration of tourism and big data, the personalized demand of tourist groups is increasingly strong. Precision marketing has become a new marketing mode that the tourism industry needs to pay close attention to and explore. Based on the advantages of big data platform and location-based service, starting from the precise marketing demand of tourism, we design data flow mining technology framework for user's mobile behavior trajectory based on location services in mobile e-commerce environment to get user track data that incorporates location information, consumption information, and social information. Data mining clustering technology is used to analyze the characteristics of users' mobile behavior trajectories, and the precise recommendation system of tourism is constructed to provide support for tourism decision making. It can target the tourist group for precise marketing and make tourists travel smarter.

1. Introduction

1.1. Research Background. Location-based service (LBS) is a kind of value-added service provided by the combination of mobile communication network and satellite positioning system. It obtains the location information of mobile terminal such as latitude and longitude coordinate data through a set of positioning technology and provides it to the communication system, mobile users, and related users to realize various location-related services in military and transportation. As a new mobile computing service in recent years, 80% of the world's information has time and location tags, and location services have developed to the big data stage [1]. Developing location services requires two capabilities: the ability to provide location and the ability to understand location.

The precise marketing information pushed by LBS location service can effectively tap the potential consumer demand and make a scientific and reasonable network marketing strategy based on this, which can further improve the ability of e-commerce enterprises to tap the target customers and potential customers. According to the 2017 China Mobile e-commerce Industry Research Report, the transaction

scale of China's e-commerce market reached 20.2 trillion yuan in 2016, an increase of 23.6% compared with the same period of the year. China's e-commerce market is developing steadily. Among them, the online shopping has a good momentum of development, up from 23.3% in 2015. Huge market potential tempts all walks of life. In 2016, online shopping and B2B e-commerce of small and medium-sized enterprises and enterprises above scale still dominate the Chinese e-commerce market, while online tourism and local life service O2O emerge as bamboo shoots, accounting for 3% and 1.6% of the market, respectively. From 2015 to 2016, the proportion of online tourism market in China's tourism market has greatly increased, the process of product informatization has accelerated, the penetration rate has further improved, the mobile online tourism market has developed rapidly, and consumer understanding and demand and experience of tourism are changing imperceptibly and pursuing higher quality of tourism. With the further development of "Internet+" information technology, the tourism industry has huge room for development, and online travel penetration will also gradually increase.

According to the Statistical Bulletin on National Economic and Social Development of 2016 issued by the State Statistical Bureau, in the whole year of 2016, the number of domestic tourists' trips reached 4.4 billion, an increase of 11.2% over the previous year, and the income of domestic tourism increased by 15.2% to 39390 billion yuan. The number of inbound tourists reached 138.44 million, an increase of 3.5%, and international tourism revenue increased by 5.6% to \$120 billion. The number of domestic residents in China has reached 135 million 130 thousand, an increase of 5.7% [2]. With the continuous promotion of the strategic pace of building a well-off society in an all-round way, tourism has become an important part of the people's daily life in China, marking that China's tourism industry has entered the era of mass tourism.

Based on this background, in the mobile e-commerce environment, based on LBS location service, research and analysis of user's mobile behavior trajectory can extract valuable user's mobile behavior features from a large number of mixed dynamic data and integrate the mobile behavior and consumer behavior of tourism users. Based on LBS location, services will integrate the mobile behavior and consumer behavior of tourism users, then excavate the marketing value of consumers, and timely achieve the marketing objectives of enterprises on the appropriate media, so that mobile e-commerce marketing becomes more accurate and effective. Through the research of this subject, the interests of enterprises, consumers, and media can be maximized at the same time, providing personalized products and services for mobile e-commerce, improving consumer loyalty and core competitiveness of mobile e-commerce and bringing higher profits for e-commerce enterprises [3].

1.2. Presentation of Problems. The data of mobile terminal users' historical consumption behavior and location movement process are recorded and stored according to the time series, forming the user's mobile behavior trajectory data, which can be collected by multiple device terminals. The user's mobile behavior trajectory data contain a lot of useful information. Mobile behavior trajectory can express the behavior activities of mobile users in the real world. These activities imply user's interests, hobbies, experiences, and behavior patterns [4]. For example, a user's activities in a week may start from home to work every day, and a user may go to shopping malls, parks, and other places on weekends. Therefore, how to effectively utilize the user's mobile behavior trajectory and extract useful information from the user's mobile behavior trajectory data is very important for the realization of personalized recommendation service.

From the view of consulting a large number of documents, there are more papers on location services than on mobile marketing. However, most of the previous articles on location services focused on the application of natural science, such as surveying and mapping technology, network development, geographic information, and so on. In recent years, the number of cross-research articles on management science,

medicine, and agriculture combined with location-based services has begun to increase, most of which are the combination of location-based services and related industries to study the application or technology development of specific industries. For example, the combination of location services and logistics technology can track the journey of packages. Combining location services with electronic maps can provide catering, entertainment, discounts, and other information within a certain range according to the location of users. Combining location service with utility technology can quickly find information such as tap water, gas explosion, and so on. The main keywords of literature research include location service technology, location service system, location service terminal, location service strategy, mobile location service, and so on. Based on location, the services industry is currently considered one of the most dynamic industries. With the rapid development of mobile Internet and Internet of Things technology, more debris time has been transferred to mobile phones, tablets, and smart products [5].

The characteristics of mobile marketing, such as precision, interaction, novelty, and effective delivery, are more and more concerned and recognized by various industries. This paper focuses on the main characteristics of the end-users of the tourism industry, such as frequent location movement, strong sense of sharing, and rich demand for services. First of all, the users of online travel must be mobile users who usually do not stay in a location for a long time, and a high probability of frequent location changes will produce a large number of location data. Moreover, in general, traveling users arrive at an unknown location or move in a series of unfamiliar geographical environments, which makes travel users' demand for location-based services take precedence over personal privacy protection and enable them to obtain real-time user location information. These location data provide favorable conditions for our research. Secondly, the behavior of tourist users is quite different from that of ordinary people. In beautiful scenery and not very familiar environment, users will spontaneously produce self-awareness. Most people share location, photos, and moods through social platforms and micromessaging, and travel companies can access these social data to accurately portray users and provide accurate services for them. Third, travel users need high-quality services to obtain high-quality tourism experience. Scenic spots, accommodation, restaurants, transportation, and finance, including tour guides and their fellow travelers, are also important factors in achieving a high-quality tourism experience. These rich demands for services will generate enormous commercial value [6]. Therefore, this article adopts top-down overall analysis to design ideas from bottom to top. By analyzing the user's characteristics through the trajectory of user's mobile behavior, this paper constructs a travel recommendation system in the mobile point-to-point environment and a precise marketing model in the tourism industry based on the trajectory of user's mobile behavior, so as to provide appropriate services for the appropriate users at the right time and place, in order to provide reference for relevant tourism enterprises to achieve precise marketing.

2. User's Mobile Behavior Trajectory

2.1. User's Mobile Behavior Trajectory Definition. User's mobile behavior trajectory is based on the path that users find frequently in the location mobile path generated by daily life. The location information generated by user's daily behavior is acquired by GPS equipment sampling at a certain time interval, and the spatial position of moving object is represented by Euclidean space coordinates, discrete display in electronic map. Through moving sequence pattern mining, we can find the correlation among these discrete location information points and obtain the user's moving behavior trajectory. This will provide effective support for precision marketing in mobile e-commerce [7]. In this paper, we make the following definitions for user's mobile behavior trajectory.

Definition 1. Location information point: the position information points generated by the user's movement can be obtained by receiving devices such as GPS of mobile terminals. Each position information point indicates a position that the user has arrived at. Suppose an independent location information point is represented as two tuple $P = (Z, T)$, among them Z is the position coordinate, and its structure contains longitude $Z.x$ and latitude $Z.y$; T is the time information of arrival position Z .

Definition 2. Mobile behavior trajectory: mobile behavior trajectory can be obtained by GPS log. A mobile behavior trajectory consists of a sequence of position information points arranged in order of time attribute T . Suppose L is user's mobile behavior trajectory, then $L = P_1 \rightarrow P_2 \rightarrow \dots \rightarrow P_n$, where $P_i (0 < i \leq n)$ denotes any sampled position information point. Mobile behavior trajectory L satisfies any $0 < i \leq n, P_i \cdot T < P_{i+1} \cdot T$; n represents the number of location information points and represents it as the length n of mobile behavior trajectory.

Definition 3. Mobile behavior subtrajectory: represents the inclusion or inclusion relationship between two moving behavior trajectories. Suppose that L_1 and L_2 have two trajectories of moving behavior, where $L_1 = a_1 \rightarrow a_2 \rightarrow \dots \rightarrow a_i, L_2 = b_1 \rightarrow b_2 \rightarrow \dots \rightarrow b_n$. If there exists a positive integer m_1, m_2, \dots, m_i , satisfying $1 \leq m_1 < m_2 < \dots < m_i \leq n$, making $a_1 = b_{m_1}, a_2 = b_{m_2}, \dots, a_i = b_{m_i}$, then L_1 is said to be the moving behavior subtrajectory of L_2 , or L_2 is said to be a moving behavior supertrajectory of L_1 . It can be written as $L_1 \subseteq L_2$ or $L_2 \supseteq L_1$. The location information points are adjacent to the mobile behavior subtrajectory and are allowed to be nonadjacent in the original mobile behavior trajectory.

Definition 4. Support degree: the collection of all location information points of moving behavior trajectories constitute a database of mobile behavior trajectories. $DB = \{L_1, L_2, L_3, \dots, L_n\}$, where $L_i (0 < i \leq n)$ is mobile behavior trajectory and $|DB|$ is the number of mobile behavior trajectories in the database. The number of mobile behavior trajectory t contained in DB is t of the support in DB :

$$\text{support}(L_i) = |\{L | L \in D, L_i \subseteq L\}|. \quad (1)$$

Definition 5. Frequent Trajectories: when the support degree of the mobile behavior trajectory is greater than or equal to the minimum support threshold, the mobile behavior trajectory is called the frequent trajectory. $FT = \{l | \text{support}(l) \geq \min, l \subseteq L, L \in D\}$, L represents the mobile behavior trajectory sequence and D represents the mobile behavior trajectory sequence set.

The user's mobile trajectory records the user's activity status in the real world, which can reflect the user's behavior preferences and potential intentions to some extent. For example, if a user moves a lot every day, he may be an outdoor sports enthusiast. Through more fine-grained analysis, we can identify users' occupations, taste habits, and so on from their frequent locations and restaurants. Therefore, mining hot spots and planning roads through multiuser mobile trajectory data sharing is an important research content of this paper.

2.2. Classification of User's Mobile Behavior Trajectory. User's mobile behavior trajectory data refer to the sequence of changes in geographic location information caused by user's own motion behavior in a certain time and space environment. These geographic location information points which change with time series can form a user's mobile behavior trajectory data according to the order of occurrence time [8]. According to the different sampling methods, we can classify these user's mobile behavior trajectory data into three categories.

2.2.1. Location Sampling-Based User's Mobile Behavior Trajectory. A trajectory formed by a change in position during the movement of a user can be sampled sequentially according to the change in position. It focuses on the information of location change when the user moves. The data obtained by this method have abundant semantic information and very detailed location change information. We can record the trajectory data of user's mobile behavior based on position sampling by recording discrete variables. The trajectory of user's mobile behavior can be represented by the sequence of sampling points with the change of moving object's position, and it can be formally expressed as

$$L = \{(x_1, y_1, t_1, \dots), \dots, (x_i, y_i, t_i, \dots), \dots, (x_n, y_n, t_n, \dots)\}. \quad (2)$$

The location $(x_i, y_i), 1 \leq i \leq n$ denotes the geographical location of the mobile user at the time of t_i , and the location (x_i, y_i) of the mobile user at the time of t_i and the location (x_{i+1}, y_{i+1}) of the time of t_{i+1} are not the same.

Trajectory can be divided into three segments according to the information of stopping point, boarding point, and alighting position, and the trajectory can be preserved according to different semantics and application segments. For example, in the prediction of travel time, it is necessary

to delete the stopping point, which may be the vehicle parking or waiting for passengers, in order to measure the trajectory travel time more accurately. For some tasks that analyze the similarity between two users, it is often necessary to use the residence trajectory to reflect the user's region of interest.

2.2.2. Time Sampling-Based User's Mobile Behavior Trajectory. The change of mobile user's behavior is sampled by definite time interval to form the trajectory data of user's mobile behavior, which is called the trajectory of user's mobile behavior sampled according to time. This kind of sampling focused on the change of location information points caused by the change of mobile user's behavior at the same time interval, which has the characteristics of large data volume and wide range. The time-sampled trajectory data of user's mobile behavior is formalized as follows:

$$L = \{(x_1, y_1, t_1, \dots), \dots, (x_i, y_i, t_i, \dots), \dots, (x_n, y_n, t_n, \dots)\},$$

$$t_i = t_1 + (i - 1)\Delta t, \quad (3)$$

where L is a trajectory data of mobile behavior, Δt is equal interval time, (x_i, y_i) , and $1 \leq i \leq n$ denotes the location of the mobile user at any time of t_i . If the time interval between the two sampling points is larger than the threshold value, the trajectory can be divided into two segments through the two sampling points.

2.2.3. User's Mobile Behavior Trajectory Triggered by Events. The trajectory of mobile user's mobile behavior, which is recorded by the system after the sensor event is triggered, is obtained by the event triggering [9]. This sampling method focuses on the event set that triggers the sensor to work. It has the characteristics of short update period and representative sampling objects. Although the behavior of mobile users changes with time, the system does not record the trajectory according to time or position, but only records the trajectory information of mobile users when they produce some specific behavior and trigger sensor events. We can also use discrete variables to record the behavior trajectory of mobile users and formalize it as follows:

$$L = \{(x_1, y_1, t_1, \dots), \dots, (x_i, y_i, t_i, \dots), \dots, (x_n, y_n, t_n, \dots)\}. \quad (4)$$

The location (x_i, y_i) , $1 \leq i \leq n$, denotes the location of the mobile user at the time of t_i , and the location of the mobile user at the time of t_i , (x_i, y_i) and t_{i+1} can be the same (x_{i+1}, y_{i+1}) .

When the trajectory direction changes beyond the threshold value, we can mark the key points according to the direction changes and divide the trajectory into two segments.

2.3. User's Mobile Behavior Pattern Decision. According to the trajectory data of user's movement behavior, the speed of completing the trajectory is calculated by time, and then the user's behavior pattern is determined. Many problems still

need to be considered, such as road congestion, construction, and even traffic accidents, which will affect the speed of user behavior. Vehicles travel much faster than people's walking speed on normal roads, but in congested or abnormal roads, the speed difference between vehicles and people's walking speed is not obvious. Therefore, the identification accuracy of trajectory velocity can only be less than 50% through time calculation [10]. In addition, the user may change several different behavior patterns in the same trip, which makes the same user's moving behavior track contain a variety of different speeds. In the overall calculation, if the average speed is obtained, it is obviously not correct to determine the user's behavior patterns. Therefore, it is necessary to divide the user's moving behavior trajectory into several trajectory segments reasonably. By comparing different trajectory segments, we can analyze whether the user has changed the behavior pattern and further improve the recognition accuracy.

How to realize the reasonable division of user movement behavior segments is the problem we want to study. As shown in Figure 1, the walking user and the driving user travel the same way, but the trajectory data of the user's movement behavior are obviously different. We can analyze the following three aspects:

- (1) Because the trajectory data of user's moving behavior produced by walking often produce direction change or reciprocating motion, we can divide the trajectory segments according to the change of the trajectory data direction of user's moving behavior. In mobile scenes, people get off a bus, walk to another station to continue to take the bus process, and must pass through a section of walking, although the walking section is short, but still can show obvious direction changes.
- (2) The trajectory data of mobile behavior produced by driving users do not change significantly in direction. This kind of characteristic is not affected by traffic conditions. We can train a classification model by the supervised learning method. For example, drivers do not change their direction as freely and frequently as pedestrians do, resulting in a straight line in the trajectory of the user's movement behavior, and the direction of change is not obvious [11].
- (3) We can also judge user behavior patterns by the shape of user behavior trajectory data, especially the trajectory of user behavior generated by different user behavior patterns in a journey, which will have obvious morphological changes of trajectory.

3. Analysis of User's Mobile Behavior Trajectory Data

This paper studies the trajectory of user's mobile behavior generated by online travel users during their mobile process. It contains a lot of information to express the personalized behavior of mobile users. We can use data mining methods such as classification, clustering, frequent itemsets, cycle discovery, and anomaly detection to mine and analyze the trajectory of tourism users' mobile behavior.

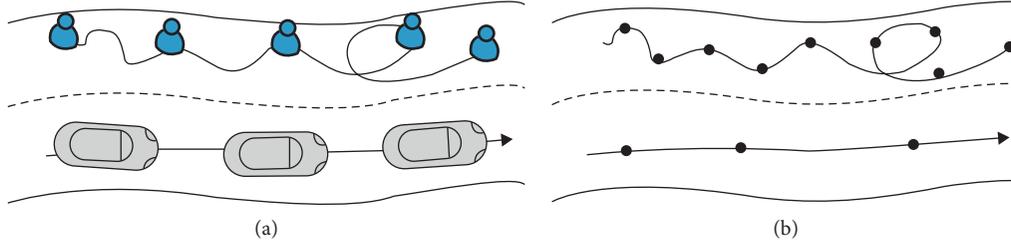


FIGURE 1: Differences in movement behavior between (a) walking users and (b) driving users.

3.1. Dividing Trajectory Segments. Each user movement behavior trajectory can be regarded as an image data. Structural Similarity Index (SSIM) can effectively measure the similarity of two trajectories, and clustering based on the similarity index is more accurate than traditional clustering based on Euclidean distance index [12]. The accuracy of structure similarity matching is closely related to the reasonableness and validity of the segmentation of user motion trajectory. Therefore, this section mainly studies how to detect the large-angle mutation points in the user's moving behavior trajectory, and how to partition and store the user's moving behavior trajectory records at the mutation points, so as to obtain some trajectory fragments which tend to be stable before clustering.

Each user movement behavior trajectory cannot be a straight line. As the precision of position coordinate recording is higher and higher, the direction of each track will change more and more, especially some subtle direction changes, and the angle of rotation can reflect the degree of change of the track direction. The division of track segments is determined according to the size of the track angle. However, if every corner is stored, it is not conducive to reduce the storage of the corner, and it is not conducive to extract it to divide the trajectory segments. Therefore, by storing the large turning point, we can discover and identify the changes of user behavior or abnormal conditions, which is also conducive to retaining the relatively stable local structure features of user trajectory segments.

We define the turning angle of user's moving behavior trajectory as the turning angle caused by the change of direction of adjacent trajectory segments, which can reflect the movement trend of trajectory and the change of user's behavior [13]. As shown in Figure 2, the angle between the direction changes of the user's moving behavior trajectory can be expressed as α , and the angle of rotation can be divided into outer angle and inner angle, expressed as θ_1 and θ_2 , respectively. We set the outer rotation angle θ_1 as a positive value and the inner rotation angle θ_2 as a negative value to facilitate the similarity calculation of the trajectory segments.

As can be seen from Figure 2, the formula for calculating the angle alpha of the direction change is shown in Formula (5), where a , b , and c represent the adjacent and opposite sides of the angle α , respectively.

$$a = \arccos \frac{a^2 + b^2 - c^2}{2ab}. \quad (5)$$

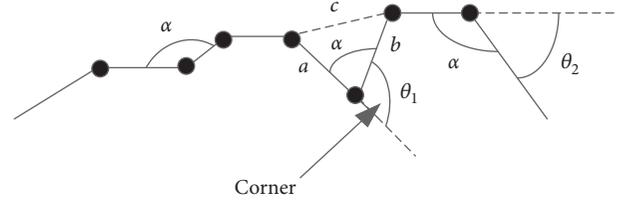


FIGURE 2: Corner of user's mobile behavior trajectory.

According to the above formula, the formula for calculating the angle theta can be obtained (6):

$$\theta = \begin{cases} 180 - \alpha, & \text{if } (a \times b \geq 0), \\ \alpha - 180, & \text{if } (a \times b < 0). \end{cases} \quad (6)$$

This is the first step to partition the trajectory segments of user's mobile behavior. Using formulas (5) and (6), the trajectory segment partitioning algorithm can be implemented (Algorithm 1).

Some trajectory fragments obtained by calculating rotation angle, setting threshold, and partitioning trajectory fragments can be expressed as a set of several feature attribute vectors. These feature attributes can comprehensively express the local features of a trajectory fragment and the global features of user's moving behavior trajectory. In this section, the trajectory fragment is not simply the expression of coordinate information of the position information points, but extracts the speed, shape, position, rotation angle, acceleration, and other characteristic vectors from it. Using these eigenvectors, we can enhance the accuracy of analyzing the trajectory of user movement. We formally represent the trajectory fragment structure as follows: $TS = (D, S, A, L)$. In addition to the above four features, we should also calculate the distance, time, and other features, using vector $W = \{W_D, W_S, W_A, W_L\}$ to represent the weight of the four feature vectors.

Since the weights of feature vectors correspond to the eigenvectors of the trajectory segments, their values should be greater than or equal to zero, and the sum of their weights should be 1; we can generally assume that the weights of all feature vectors are equal probability, and we can take the average value of 0.25 as the weights. Similarly, we can adjust the weights of each feature vector according to the sensitivity of the feature vectors of the trajectory fragments in the actual scene. For example, when analyzing the position-sensitive

Step 1: one by one, scanning the location information point sequence in the user movement behavior track;
 Step 2: formula (5);
 Step 3: formula (6);
 Step 4: Set a threshold ω for corner θ , store the corner satisfying $|\theta| > \omega$ as a mutation point, and then divide the track segment according to the position information point of the corner. n is the number of sampling points, and the time complexity of the algorithm is $O(n)$.

ALGORITHM 1

trajectory fragments, we can focus on the position vectors, and the weights $W_D = W_S = W_A = 0, W_L = 1$ are also feasible.

According to the feature vector and its weight to complete the structural similarity comparison, mainly through the analysis of the differences between the feature vectors of the trajectory segments to complete the comparison [14], according to the definition of the trajectory segment structure, we can define two trajectory segments are $L_i, L_j, 1 \leq i \neq j \leq n$. The comparison function of two trajectory segments is $D(L_i, L_j)$, velocity vector is $S(L_i, L_j)$, angle vector is $A(L_i, L_j)$, and position vector is $L(L_i, L_j)$. The four comparison functions above constitute the calculation of structural similarity of the trajectory segments, as shown in the following Formulas (7) and (8). The function $N(\dots)$ denotes the normalization of the distance. Because the range of each eigenvector in the trajectory segment is different, the normalization of the distance is the normalization of the distance of each eigenvector. The SSIM of structural similarity is represented by 1 minus the normalization of the distance:

$$S(L_i, L_j) = (D \times W_D + S \times W_S + A \times W_A + L \times W_L), \quad (7)$$

$$\text{SSIM}(L_i, L_j) = 1 - N(S(L_i, L_j)). \quad (8)$$

The structural similarity comparison of trajectory fragments can express the structural differences of each trajectory fragment on the feature vectors. Therefore, the smaller the SSIM value of the trajectory fragments, the greater the SSIM value of the trajectory fragments. Moreover, the distance between the structural similarities of the trajectory fragments is symmetrical, that is, $\text{SSIM}(L_i, L_j) = \text{SSIM}(L_j, L_i)$. Therefore, it can be found that the method based on structural similarity can well reflect the structural differences between trajectory segments.

According to structural similarity, the direction information, speed information, angle information, and position information are compared [15].

- (1) The direction vector comparison function $D(L_i, L_j)$ denotes the degree of similarity of two similar trajectory segment L_i, L_j in the direction of motion. As shown in Figure 3(a), ϕ is the angle between the direction of the trajectory, and the formula for calculating direction vector comparison function is as follows:

$$D(L_i, L_j) = \begin{cases} \|L_i\| \times \sin \phi, & \text{if } (0^\circ \leq 90^\circ), \\ \|L_j\|, & \text{if } (90^\circ \leq 180^\circ). \end{cases} \quad (9)$$

If two similar trajectory fragments have the same direction and the angle ϕ is small, the two trajectory fragments tend to be parallel in the same direction, which is called the best state, then the Dir Dist value approaches zero. If two similar trajectory fragments are in opposite directions and the two trajectory fragments with larger angle ϕ tend to be in reverse parallel, the worst condition is that the Dir Dist value is the length of the trajectory fragments involved in the comparison.

- (2) The speed vector comparison function $S(L_i, L_j)$ expresses the trend of user mobility. The velocity vector comparison function is shown in Formula (10), where $S_{\max}(L_i, L_j)$ is $|V_{\max}(L_i) - V_{\max}(L_j)|$, representing the absolute value of the maximum velocity difference between the trajectory segments. Similarly, $S_{\text{avg}}(L_i, L_j)$ and $S_{\min}(L_i, L_j)$ represent the absolute value of the difference between the average velocity and the minimum velocity, respectively. We can judge the difference of velocity vectors from the three aspects of maximum, minimum, and average velocity:

$$S(L_i, L_j) = \frac{1}{3} (S_{\max}(L_i, L_j) + S_{\text{avg}}(L_i, L_j) + S_{\min}(L_i, L_j)). \quad (10)$$

- (3) The angle vector comparison function $A(L_i, L_j)$ expresses the degree of eigenvalue change caused by the change of direction in the trajectory segment. As shown in Formula (11), where the angle of rotation θ is calculated according to Formula (6), the internal rotation angle is positive and the external rotation angle is negative, the angular distance of the trajectory segment is the cumulative value of many internal corners of the trajectory, and the direction of change within the trajectory segment can determine the value of each angle:

$$A(L_i, L_j) = \frac{\sum_{1,1}^{P(L_i), P(L_j)} (|\theta_i - \theta_j|) / (|\theta_i| + |\theta_j|)}{P(L_i) + P(L_j)}. \quad (11)$$

Figure 3(b) shows that if each corner of the two trajectory segments rotates to L_i and L_j matches, the value of the angle vector comparison function is 0, which is the best case. If the two trajectory segments turn to L_i and L_j in opposite directions, that is, the two trajectory segments are

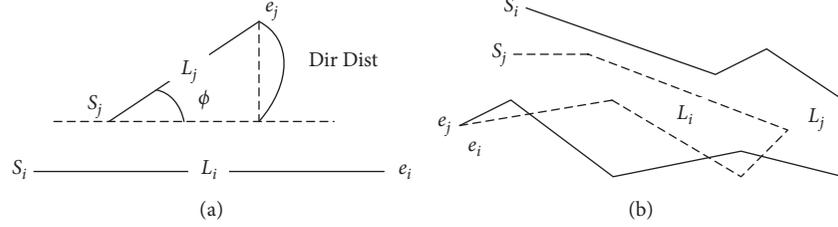


FIGURE 3: Comparison of track direction and rotation angle: (a) direction contrast; (b) corner contrast.

in opposite jagged shape, and the value of the angle vector comparison function is 1, this is the worst case.

- (4) For the position vector comparison function $L(L_i, L_j)$, we can use Hausdorff distance to measure the location distance of the trajectory segment, as shown in the following formula:

$$L(L_i, L_j) = \max(h(L_i, L_j), h(L_j, L_i)), \quad (12)$$

where $h(L_i, L_j) = \max_{a \in L_i} (\min_{b \in L_j} (\text{dist}(a, b)))$ is the direct

Hausdorff distance between L_i and L_j , i.e., the maximum distance from a point in L_i to the nearest L_j and $\text{dist}(a, b)$ represents the Euclidean distance function between points.

3.2. Similarity Computation of User's Mobile Behavior Trajectory. At present, we collect and store the trajectories of tourism users' mobile behavior, cluster the typical similar trajectories from these trajectory data, analyze the behavior patterns of user's mobile behavior trajectories, and predict the personalized needs of tourism users based on structural characteristics. Clustering analysis is to divide user behavior trajectory into several groups with high cohesion and low coupling. It requires high similarity of user behavior trajectory in the same group, and low similarity of user behavior trajectory in different groups. The goal of clustering analysis is to find out the trajectory data with the same or similar behavior patterns from the trajectories of some users' mobile behaviors, analyze the personal preferences, consumer demands, and behavior characteristics of the trajectories of tourism users' mobile behaviors, and accurately determine the similarity between trajectories of users' mobile behaviors. At the same time, the trajectories of users' mobile behaviors with high similarity are gathered into one class [16].

Most of the online travel users are in the same scenic spot, similar routes to carry out activities, and most of the resulting mobile behavior trajectory data have local similarity and global dissimilarity. It is difficult to find the personalized characteristics of tourism users by analyzing the complex and large number of users' mobile behavior trajectories and effectively extract users. The analysis of a part of the mobile behavior trajectory is more conducive to finding the information contained in it [17]. Therefore, the trajectory analysis method based on the whole trajectory in traditional research is easy to cause the inaccuracy of trajectory analysis. In this paper, we use

structural features to calculate the similarity of user movement behavior. This method needs to calculate every corner of the user's mobile behavior trajectory and find the sampling point with larger rotation angle, which is regarded as the sudden change point of the user's mobile behavior, and then divides the trajectory segment by the sudden change point. In this way, the rotation angle of each trajectory segment obtained will not change significantly, and the trajectory structure tends to be stable. Then, a trajectory model of user's mobile behavior is constructed, which is characterized by trajectory direction, trajectory speed, trajectory angle, and trajectory distance. Taking these features as parameters, threshold values are set to express and adjust the weights of each feature according to the actual application scenarios, and a trajectory similarity algorithm is constructed to calculate the user's movement behavior. The object of this paper is to calculate the structural similarity of some trajectory segments which are divided according to the sudden change points of large turning angles by using the trajectory similarity algorithm constructed with structural features as parameters. It is used to judge the similarity degree of each user's moving behavior trajectory and then completes the feature analysis of user's moving behavior trajectory. The simulation results show that the trajectory similarity calculation algorithm is efficient, the weight adjustment of each structural feature is flexible, and the trajectory analysis results are more in line with the needs of practical application scenarios and have higher application value and practical significance.

On the basis of obtaining the feature vector distance of user's moving behavior trajectory segment, the trajectory segment with high similarity is analyzed, and then the clustering algorithm is used to complete the clustering of user's moving behavior trajectory. By comparing the structural similarity between the trajectory segments and other trajectory segments which are not on the same trajectory, a number of ε -nearest neighbor sets of trajectory segments are formed. The number of ε -nearest neighbor sets is used to determine the midpoint of trajectory segment clustering, and then the trajectory segment clustering is realized. A trajectory segment clustering algorithm based on structural similarity is constructed.

The steps of clustering algorithm based on structural similarity are given in Algorithm 2.

From the analysis of the above algorithms, it can be seen that, in the trajectory segment clustering algorithm based on structural similarity, it is very important to determine the

Step 1: first calculate the corner θ of each track segment sampling point P_i ;

Step 2: according to the corner threshold ω , we divide the trajectory of user movement into TS of some track segments.

Step 3: calculate the distance between the trajectory feature vectors based on the weight of the trajectory segment feature vectors.

Step 4: calculate the ε -nearest neighbor set of the track segments with high similarity.

Step 5: the distance clustering segment is centred on the similarity track segment ε -nearest neighbor set.

Step 6: initialize clustering ID and track segment clustering markers.

Step 7: traverse the trajectory fragments, find the core clustering and set the clustering ID, and then add the pointers of these trajectory fragments to a new node in the index tree.

Step 8: determine whether the set center of ε -nearest neighbors meets the specified distance. If it meets the requirement, then add the cluster ID marker to the trajectory fragment, expand the clustering, construct the index tree node, and repeat steps 7 and 8 until all trajectory fragments are traversed.

ALGORITHM 2

threshold value of ω , ε -nearest neighbor, and the threshold value of σ nearest neighbor number, which can directly affect the time complexity and space complexity of the algorithm. It needs to be verified repeatedly and determined according to the actual application fields. Therefore, we mainly analyze the algorithm qualitatively.

Through repeated verification of the algorithm, in the data analysis of trajectory of travel user's movement behavior, the value of ω cannot be set too small, and if set too small, some characteristic details of trajectory segments will be lost. On the contrary, the value of ω cannot be set too large and cannot effectively identify the abrupt change point or sampling abnormality of the trajectory segment, which directly affects the structure of clustering analysis. Similarly, if the threshold value σ of the number of neighbors is set to be large enough, then no trajectory segment can satisfy the requirement of $|N_\varepsilon(L)| \geq \sigma$, and all trajectory segments will be marked as abnormal conditions. On the contrary, if σ is set too small, all the trajectory segments may become the clustering center, so that the trajectory segments will be self-contained and the number of clusters will be too large.

3.3. Discovery of Popular Tourist Attractions. By effectively identifying the location information points in the trajectory data of users' mobile behavior, the feature vectors of the trajectory segments can be extracted, and the semantics of these location information points can be expressed as the route, the scenic spots, and the behavior patterns of an online travel user in the past period of time. By clustering and analyzing the trajectory fragments containing location information points, we can find that the traveling users have a longer time in a certain area, which can be interpreted as the tourist users have a higher degree of interest in a certain scenic spot. Semantic expression is a popular tourist spot with longer stay time for online travel users. In practical scenarios, many traveling users will visit the same or similar scenic spots. From the trajectory of users' mobile behavior and the region of interest, traveling users with similar trajectory and the same region of interest can predict their similar preferences or similar behavior characteristics. These regions of interest frequently stayed by tourist users will appear as overlapping regions in the trajectory of user's mobile behavior. If these

overlapping regions are found, the popular scenic spots concerned by tourist users can be found and the users who like these scenic spots can be clustered. And then, dig out the other characteristics of these users to complete the personalized tourist attractions recommendation of similar tourists. We extract the feature parameters of these overlapping areas, such as overlap time and overlap times, which can reflect the similarity between the traveling users. It can identify the tourist attractions that the tourist users are interested in during the mobile process and recommend the most likely popular tourist inventory for other tourist users who have a higher similarity with their user's mobile behavior trajectory, so as to tap the potential preferences of the tourist users [18]. Assuming that travel user A and travel user B share a higher degree of similarity in the trajectory of users' mobile behavior, it can be found that some scenic spots are visited by users A but not by users B. Through mining, it is known that these scenic spots may be of interest to users B. Then, we can recommend these scenic spots to B users through A users, so that these scenic spots become the potential and most likely scenic spots for B users to visit. We can also use the activity sequence to express the popular scenic spots that tourists often visit, and the trajectory of nearby tourists is more instructive [19].

In the process of analyzing mobile user behavior trajectory data with structured eigenvectors, it is not difficult to find that the moving speed of user behavior trajectory is not the same in different time periods, or it is slower in a certain time period, or it is faster in a certain time period. Figure 4 shows the moving speed of the user in different periods of time, in which the trough is formed during the period when the user moves slowly and the peak is formed during the period when the user moves fast, but both trough and peak can indicate the user's continuous generating activity. And, the slow moving trough time contains more user behavior characteristics, so this paper focuses on the behavior characteristics of mobile users in trough situation.

As shown in Figure 4, by comparing the speed, distance, and time of user's moving behavior trajectory, the structured features of mobile users and the behavior differences of tourist users can be clearly analyzed. We focus on the analysis of two dimensions: the speed and the time of the slower wave trough. As shown in Figure 5, the slower the traveling speed of the tourist user, or the less the change of

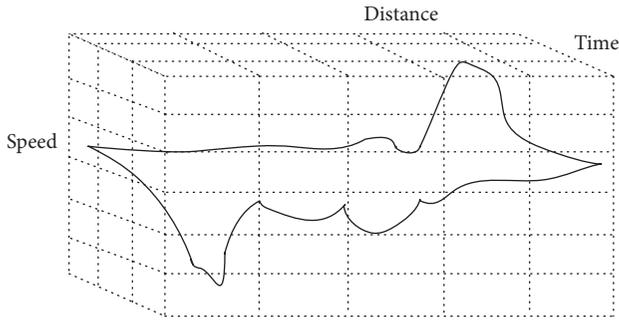


FIGURE 4: Mobile speed of user behavior trajectory.

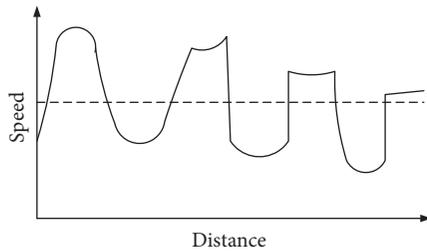


FIGURE 5: Two-dimensional trajectory analysis of user movement behavior.

the active area in a period of time, the most likely the predictable user behavior characteristic is; that is, the traveling user stays at a certain scenic spot for browsing, resting, or taking photos. The longer the trough, the more attractive the scenic spot is. The more tourists are staying at the same scenic spot, the more scenic spots can be designated as popular tourist attractions.

For example, when a tourist visits a scenic spot, he or she forms a trajectory of the user’s movement behavior. Three troughs appear in the trajectory, indicating that the user may have experienced three scenic spots or rest areas, of which the first trough has a shorter experience. It shows that tourists spend less time visiting the first scenic spot, travel faster, continue to move forward at a faster speed after the tour, and spend a little more time watching the tide or taking pictures when they meet the scenic spot of interest. So tourists will slow down, move in a more fixed area, and travel at a slower speed, thus appearing the second trough period, after the tourists continue to move forward; when the formation of the third trajectory speed reached trough state, semantic expression may have two situations. The first is that the tourists reach a certain degree of fatigue or meet a rest area, stop and rest; the second is that the tourists arrive at a well-known scenic spot, gather more tourists, people will stay in a certain position, waiting for sightseeing and photography, moving slowly, and almost stop. The above two semantics can be distinguished by whether the location in the electronic map is a resting area or a scenic spot. However, in the actual tourist attractions, the situation may be more complicated. For example, a tourist is an outdoor sports enthusiast who has good physical strength and likes natural scenery. Because of his fast moving speed, there is little difference between the wave crest and trough of the waveform trajectory formed by the speed and distance.

Although his tour speed is fast and his stay time is short, the location he stays in is still the area of interest. In this way, moving objects with similar frequencies in the velocity-distance waveform can be found not only in the known hot spots of the users, but also in the scenic spots that the potential users may be interested in, even in the preferences, occupations, and personality characteristics of the tourist users. It helps to gather tourists with similar preferences and similar personalities to achieve the confluence module [20].

Popular scenic spots refer to scenic spots with long staying time after arrival [21]. In the user’s mobile behavior trajectory, the hot spots can be marked as $H = \{H_1, H_2, \dots, H_n\}$, $H_j = \{L_j, L_{j+1}, \dots, L_m\}$, H is used to denote a trajectory fragment. When the traveling user passes through a hot spot area with high interest and stays for a long time, the trajectory fragment moves at a speed close to or far below the normal trajectory speed. We can think that the tourist user has conducted a deep browsing in the scenic spot or some behavior activities have taken place in the scenic spot area. We can analyze the information such as the time of arrival, the time of stay, and so on. The region with dense user access points can be expressed as a popular tourist attraction area with high user access frequency [22].

Because GPS receiving equipment receives satellite signals in vast and open areas with high intensity and good positioning effect, satellite signals in indoor areas will be shielded by the wall, resulting in weak positioning signal and reduced positioning accuracy [23]. Therefore, when analyzing tourists’ preference for scenic spots through the status of stay, it is necessary to distinguish between outdoor and indoor scenic spots. The positioning signal of outdoor scenic spots is good and has high precision. It can acquire the location information points at sampling frequency in real time and form the locus of user’s movement with dense location information points. The positioning signal of indoor scenic spots is weak, which affects the positioning accuracy. Even when the signal is lost, the location information points cannot be obtained in time according to the sampling frequency requirement, and the space area of the indoor attractions is small, which makes some location information points overlap. This repetitive activity can also find that the tourists are visiting a certain indoor attractions regularly. The popular scenic spots are divided into two types: one is the outdoor scenic spots, such as natural landscape, gardens, playgrounds, and other broad areas, in a longer period of time, can obtain more dense location information points formed by the user’s mobile behavior trajectory, recorded as HR_{II} ; Another kind is indoor scenic spots, such as restaurants, shopping malls, tourist centers, and other closed areas, in a long period of time, may lose a certain location information point sampling information, but after leaving the area, they can get the location information point again, recorded as HR_I . Firstly, the trajectory data of user’s mobile behavior are obtained by sampling the location information points, and then the trajectory data of user’s mobile behavior is denoised. Finally, according to the characteristics of the location information points, the HR_I and HR_{II} popular scenic spots domain are divided by the density clustering method. The steps for finding popular scenic spots is given in Algorithm 3.

```

Input parameters: user movement behavior trajectory  $Q$ , minimum speed  $S$ , minimum time  $T$ , and maximum disturbance threshold  $MT$ .
Output parameters: collection of popular scenic spots  $HR$ .
Step 1: for  $(i=2, i \leq |Q|, i++)$  /*  $|Q|$  represents the number of location information points */
Step 2:  $D[i-1] \cdot T = \text{cal } T(p_{i-1}, p_i)$ ;
Step 3:  $D[i-1] \cdot S = \text{cal } D(p_{i-1}, p_i) / D[i-1] \cdot T$ ;
Step 4:  $HR = \{\}$ ;  $C = \{\}$ ;  $CO = \text{false}$ ;
Step 5: for  $(j=2, j \leq |T| - 1, j++)$  /* cycle search indoor attractions area  $HRI$  */
Step 6: if  $(D[j-1] \cdot T > T$  and  $D[j-1] \cdot S < n * S)$  then
Step 7:  $C = \{p_{j-1}, p_j\}$ ; /* record location information point stay area */
Step 8: if (not  $CO$ ) then  $CO = \text{true}$ ;
Step 9: else  $C = \text{Update}\{p_{j-1}, p_j\}$ ; /* merge the location information points closer to the collection  $HR$  */
Step 10: else if ( $CO$ ) then
Step 11:  $HR = \{C\}$ ;  $CO = \text{false}$ ;  $C = \{\}$  /* search outdoor attractions area  $HRII$  */
Step 12: if  $(D[j-1] \cdot S \leq S)$  then /* determine activity intensive areas */
Step 13:  $C = \{p_j\}$ ;
Step 14: if (not  $CO$ ) then  $CO = \text{true}$ ;
Step 15: else if ( $CO$ ) then
Step 16: last Index = look Ahead ( $MT, S$ );
Step 17: if (last Index  $\leq j + MT$ ) then
Step 18: for  $k = \text{last Index}$  downto  $j$  do
Step 19:  $C = \{p_k\}$ ;
Step 20:  $j = \text{last Index}$ ;
Step 21: else if (time ( $C$ )  $> T$ ) then
Step 22:  $H\ddot{E} = \{C\}$ ;  $C = \{\}$ ;  $CO = \text{false}$ ;
Step 23: return  $HR$ ;

```

ALGORITHM 3

Because the location information points sampled by GPS are affected by the factors of region, space, and weather, it is easy to have inaccurate positioning or interruption of positioning [24]. When the user enters the indoor scenic spots from the outdoor scenic spots, the short signal interruption will occur, and the positioning data receiving error will easily occur. In order to adapt this error, a maximum disturbance threshold MT is set in the hot spot detection algorithm to enhance the accuracy of location information points.

The outdoor and indoor scenic spots are divided according to the location information of popular scenic spots. The density-based hot spot discovery algorithm is used to retrieve two different types of user residence areas, outdoor and indoor, in the trajectory of user's mobile behavior, and define them as hot spots [25]. The algorithm has four input parameters: trajectory of user's mobile behavior, minimum speed, minimum time, and maximum disturbance threshold. Among them, the threshold of minimum speed is related to the activity speed of tourist users in the scenic spot area. If walking tour, the general speed is 2 to 3 meters per second. If the sudden reaction speed of the location information points slows down significantly, it indicates that the tourist users have arrived at the scenic spot area, and specific user behavior has taken place. On the contrary, if the sudden response speed of the location information points increases over a period of time, it indicates that the travel users have changed their behavior patterns. For example, we can leave the scenic spot and take a sightseeing bus to the next scenic spot, so we can

set it according to the sampling time. There is no absolute fixed value in the setting of the minimum time threshold. Generally, if a tourist user stays in a certain area for more than 30 minutes, it can be considered that the tourist user arrives at a scenic spot or rest area, and changes in user behavior have taken place, resulting in a specific activity. The maximum perturbation threshold is only used to express the number of continuous perturbations when the abnormal location information points are sampled. If the number of abnormal location information points is smaller than the perturbation threshold, the abnormal sampling information points can be merged into the normal location information points of the user's mobile behavior trajectory data set. If the number of abnormal location information points sampled in the former HR is larger than the perturbation threshold, it is necessary to preserve the current user's mobile behavior trajectory and then detect the new hot spots after abnormal location information points sampled to form a new user's mobile behavior trajectory. The settings of the minimum time threshold and the maximum disturbance threshold are related to the sampling frequency of the location information points. In the popular scenic spot area detection algorithm, the trajectory of the user's movement behavior is traversed twice. The algorithm complexity is linear order $O(n)$. Among them, n denotes the number of location information points in the trajectory of the user's mobile behavior, and the algorithm can retrieve the hot spots with frequent activities.

4. Travel Recommendation Model

4.1. Application Scene

- (1) *Tourist Attractions Recommendation.* In the mobile scenario where the user completes the tourist attraction tour, the tourist user will generate multi-dimensional information to realize the application scenario recommended by the tourist attraction, which can hide the rarely used dimensions. Focus on the user's city, mobile behavior trajectory in the corner location information points, hot spots and user behavior patterns and other dimensions, accurate analysis of user preferences, the same characteristics of the user recommend the most likely favorite tourist attractions [26].
- (2) *Hotel Catering-Related Recommendation.* Tourist users need to visit other applications frequently in the process of touring, such as electronic map applications, O2O applications, restaurants, e-commerce, outdoor equipment, etc., a single visit to each application, users need to frequently exit one application login another application, will cause user inconvenience, inefficiency, etc. According to the application relevance assessment model, we analyze the application of tourism users' preferences. Establish the relationship between these applications, so that users can access a travel APP directly related to other applications they are used to daily life, making travel APP a personalized all-round service platform [27].
- (3) *Tourism User Preference Content Recommendation.* Tourist users may often inquire about certain contents, such as outdoor equipment, fitness and health care, restaurants and entertainment, and surrounding scenic spots in the process of using the application. According to the content retrieval evaluation model, users can retrieve keywords, learn user preference content, and discover their hobby characteristics and behavior characteristics. According to these preferences, travel APP can construct a personalized interface for users, giving priority to the information, articles, and news that travel users are interested in [28].

4.2. Travel Recommendation System. With the combination of Bluetooth, WIFI, and other RF communication technologies and mobile terminal devices, mobile point-to-point communication environment has been derived, and many different research topics have also emerged. This research talks about the relationship between mobile attraction recommendation system and social software from the perspective of mobile social software and completes interaction through user comment sharing in mobile point-to-point environment. In this paper, an interactive system of tourism comment information sharing and social networking software is established, which includes three functions: recommendation, reunion, and comment. It is used to explore the interaction between users in mobile point-to-point environment. The preliminary test has been

completed in this paper. The experimental results show that the recommendation, convergence, and comment functions of the system can provide precise services for users and provide a basis for further research on the wide application of user behavior trajectory in precise marketing.

This paper focuses on the problem of information sharing and social interaction of tourism mobile recommendation system in mobile point-to-point environment. The system mainly includes three functions: recommendation, convergence, and information sharing. In the recommendation function section, we assume that users will leave comments and other information after visiting a scenic spot. When other users meet with them, they can exchange comments through RF communication technology. These comments are calculated by system algorithm to recommend scenic spots that meet users' interests. In addition, users can also actively send requests to other people to join the information and find similar interests around users to visit a scenic spot. Of course, users can also actively share location, comments, traffic conditions, tourist density, and other related information.

In order to enable the interaction and sharing of information between remote users, the relay mode under mobile point-to-point can be adopted. Every user in the system plays the role of information transmission, that is to say, each user's mobile terminal is a relay node for information transmission and constantly transfers the information they have mastered to the users at a long distance.

4.2.1. System Architecture. The mobile peer-to-peer environment mainly transmits information through the direct transmission between peer-to-peer users and the relay mode assisted by the third party. Using this feature, the system proposed in this paper mainly provides three services: recommendation, convergence, and review. First of all, the main purpose of recommendation service is to recommend scenic spots similar to user's interests to users through user's evaluation information, so that users can have a reference direction for the next destination in the journey, so that users can travel more smoothly. Secondly, the convergence service allows users to initiate a convening activity, gather other interested users around, visit the scenic spots together, or buy specialty goods together, through group buying to get a better price, or to strive for preferential services. Thirdly, evaluation services are divided into general information and specific information. General information is simply the transmission of personal information and specific information, so that the use of convergence services conveys the convening activities of the department of the offensive, through short messages, and the expression of personal information is incompatible; specific information is only for convening activities issued information. To provide the above services, the system architecture is presented in Figure 6.

(1) Interface Module. This module is responsible for the user and the system function docking; through this module, the system function interface is expressed and the user is guided to carry out the operation of various functions.

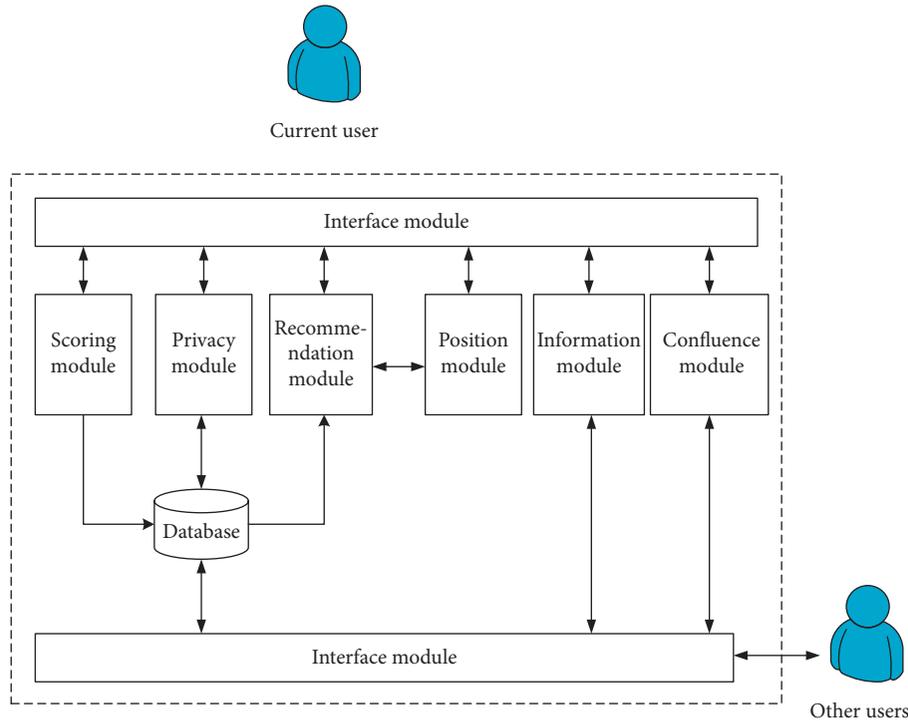


FIGURE 6: Mobile point-to-point tourism recommendation system architecture.

(2) *Scoring Module*. At present, the commonly used recommendation system is based on the scoring mechanism, which collects user's scoring data to calculate and provide recommendation services. The scoring module proposed in this study is mainly responsible for recording the user's evaluation of scenic spots. At the initial stage, the recommendation system often faces the problems of incomplete user scoring information, too many items not scored, and the difficulty of calculation caused by the noise of data, resulting in the decline of recommendation accuracy. Therefore, this study classifies scenic spots, requires users in the initial stage, must be based on the type of scenic spots scoring to ensure that individual users in the initial stage, and needs to score for the type of scenic spots to ensure that individual users' scoring information has been scored by the common column.

(3) *Transport Module*. Because this research system is built in the mobile point-to-point environment, user scoring, information, and other need to be obtained and transmitted through the transmission function; this study uses Bluetooth transmission technology to achieve the transmission of related functions. This module enables the system to automatically exchange scoring data through Bluetooth without interfering with the user when they meet, so as to achieve the purpose of collecting data. In terms of scoring exchange mechanism, this study currently uses unlimited scoring exchange method, when users meet, the exchange of all the scoring data held by both sides. However, information that has not yet been scored is not helpful to the recommendation system. Therefore, in the scoring exchange module, it is assumed that only the user has scored more than five scenic spots before the exchange, while the other scoring data

obtained by others is more than five items before passing on to other users. On the other hand, the transmission module has the search function and can discover other users around the user; when the user wants to pass its ideas to the surrounding users, it can be completed through this module.

(4) *Recommendation Module*. This study analyzes the recommended operation by exchanging accumulated score data. This recommendation module uses collaborative recommendation and Pearson correlation coefficients to perform recommendation operation. The formula is shown in (13). Suppose that $U(i, a)$ is used to predict the possible degree of preference of i to a scenic spots. $F_j(a)$ is the score of user j for a attractions and \bar{F}_i is used to score the average score of holder i . \bar{F}_j is the average score of user j , and $\text{sim}(i, j)$ is the similarity between user i and user j calculated by Pearson correlation coefficient:

$$U(i, a) = \bar{F}_i + \frac{\sum_j \text{sim}(i, j) \times |F_j(a) - \bar{F}_j|}{\sum_j \text{sim}(i, j)}. \quad (13)$$

In the process of recommendation, the recommendation module first calculates Pearson correlation coefficient, calculates the first 20 items of scoring data which are most similar to users, and then runs the subsequent recommendation algorithm. Finally, the user's predicted value of a certain scenic spot is obtained by weighted average of these scoring data and similarity, and five scenic spots with the highest predicted value are recommended to users for reference.

(5) *Position Module*. This module can use Bluetooth GPS receiver to receive satellite signals and select the local latitude

and longitude values to determine the location of the user. Finally, combined with the processing results of the recommendation module, the electronic map shows the location of each recommendation site and the location of the user.

(6) *Information Module*. The concept of user's active comment can be added in the system; through the transmission of information, users can express their personal ideas to other users around. Adding the function of transmitting information in this part, the user can not only transmit the new information but also transmit the received information to other users in the transmission range.

(7) *Convergence Module*. Since the system is designed in a mobile peer-to-peer environment, a mobile convergence function is derived from the concept of mobile social networks. Through this function, travelers can dynamically search for other users with the same goals and preferences. Through the transmission of information, travelers can share the requested information to the surrounding users and thus find travelers willing to act together.

(8) *Privacy Module*. One of the focuses of mobile social software is to explore the interaction between users, but not everyone is willing to interact with others, so this study adds privacy considerations. This module can provide users whether to allow all other users or only allow some friends to search their own location through the system; through the privacy settings, users can not be disturbed by other users to carry out system operations but also to observe whether there is a willingness to interact between users.

J2ME can be chosen as the development platform of the system, and Bluetooth technology is the basic wireless transmission technology commonly available in mobile terminals. Therefore, it is feasible to use Bluetooth technology as a transmission tool. In order to expand the scope of information transmission, WIFI wireless network can also be considered as a transmission medium, which can effectively solve the problem of short transmission distance and unstable signal of Bluetooth.

5. Conclusion

Advanced GPS devices enable people to record their location histories with GPS trajectories. The trajectory of users' mobile behavior means to a certain extent that a person's behavior and interests are related to their outdoor activities, so we can understand the users and their locations and their correlation according to these trajectories. This information enables accurate travel recommendations and helps people to understand a strange city efficiently and with high quality. By measuring the similarity of different user location histories, the similarity between users can be estimated and personalized friend recommendation can be realized. The user stereoscopic user portrait can be portrayed through the integration of user movement behavior trajectory and social information. This paper takes the trajectory data of tourism users' mobile behavior as the research object and constructs the tourism precise marketing model. In the process of

obtaining the trajectory of user movement, the characteristics of mobile user behavior track data are taken into account. The sensitivity of various features in the trajectory analysis process is adjusted by weight. The structured feature vectors and popular scenic spots discovery methods of user's mobile behavior trajectory are fully studied by clustering and collaborative filtering techniques, which lay a foundation for constructing the application model of tourism precision marketing.

Data Availability

The data used to support the findings of this study are available from the corresponding author upon request.

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

Acknowledgments

This research work was supported by the Ministry of Education Humanities and Social Sciences Planning Fund Project (No. 18YJAZH128) and Research Project of Harbin University of Commerce (No. 18XN022).

References

- [1] Y. Yuan and M. Raubal, "Measuring similarity of mobile phone user trajectories—a Spatio-temporal Edit Distance method," *International Journal of Geographical Information Science*, vol. 28, no. 3, pp. 496–520, 2014.
- [2] Z. Sun and X. (Jeff) Ban, "Vehicle classification using GPS data," *Transportation Research Part C: Emerging Technologies*, vol. 37, pp. 102–117, 2013.
- [3] D. Wang, "Approaches for transportation mode detection on mobile devices," in *Proceedings of Seminar on Topics in Signal Processing*, pp. 77–82, 2014.
- [4] S. Hong and A. Vonderohe, "Uncertainty and sensitivity assessments of GPS and GIS integrated applications for transportation," *Sensors*, vol. 14, no. 2, pp. 2683–2702, 2014.
- [5] M. Lin and W.-J. Hsu, "Mining GPS data for mobility patterns: a survey," *Pervasive and Mobile Computing*, vol. 12, pp. 1–6, 2014.
- [6] S. Khajezadeh, H. Oppewal, and D. Tojib, "Mobile coupons: what to offer, to whom, and where?," *European Journal of Marketing*, vol. 49, no. 5-6, pp. 851–873, 2015.
- [7] H. Li, "Review on state-of-the-art technologies and algorithms on recommendation system," in *Proceedings of the International Conference on Mechatronics Engineering and Information Technology (ICMEIT 2016)*, p. 7, Wuhan Zhicheng Times Cultural Development Co., Xi'an, China, 2016.
- [8] Htet Htet Hlaing, "Location-based recommender system for mobile devices on University campus," in *Proceedings of 2015 International Conference on Future Computational Technologies (ICFCT'2015); International Conference on Advances in Chemical, Biological & Environmental Engineering (ACBEE) and International Conference on Urban Planning, Transport and Construction Engineering (ICUPTCE'15)*, p. 7, Universal Researchers in Science and Technology; Universal Researchers in Science and Technology, Singapore, March 2015.

- [9] W. Wörndl and B. Lamche, "User interaction with context-aware recommender systems on Smartphones," *icom*, vol. 14, no. 1, 2015.
- [10] L. O. Colombo-Mendoza, R. Valencia-García, G. Alor-Hernández, and P. Bellavista, "Special issue on context-aware mobile recommender systems," *Pervasive and Mobile Computing*, vol. 38, pp. 444-445, 2017.
- [11] L. O. Colombo-Mendoza, R. Valencia-García, A. Rodríguez-González, G. Alor-Hernández, and J. J. Samper-Zapater, "RecomMetz: a context-aware knowledge-based mobile recommender system for movie showtimes," *Expert Systems with Applications*, vol. 42, no. 3, pp. 1202-1222, 2015.
- [12] W.-S. Yang and S.-Y. H. iTravel, "A recommender system in mobile peer-to-peer environment," *Journal of Systems & Software*, vol. 86, no. 1, pp. 12-20, 2013.
- [13] T. Pessemier, D. Simon, K. Vanhecke, B. Matté, E. Meyns, and L. Martens, *Context and Activity Recognition for Personalized Mobile Recommendations*, Springer, Berlin, Germany, 2014.
- [14] J. Zeng, F. Li, Y. Li, J. Wen, and Y. Wu, "Exploring the influence of contexts for mobile recommendation," *International Journal of Web Services Research*, vol. 14, no. 4, pp. 33-49, 2017.
- [15] A. Majid, L. Chen, G. Chen, H. T. Mirza, I. Hussain, and J. Woodward, "A context-aware personalized travel recommendation system based on geotagged social media data mining," *International Journal of Geographical Information Science*, vol. 27, no. 4, pp. 662-684, 2013.
- [16] D. Gavalas, C. Konstantopoulos, K. Mastakas, and G. Pantziou, "Mobile recommender systems in tourism," *Journal of Network and Computer Applications*, vol. 39, pp. 319-333, 2014.
- [17] S. K. Hui, J. J. Inman, Y. Huang, and J. Suher, "The effect of in-store travel distance on unplanned spending: applications to mobile promotion strategies," *Journal of Marketing*, vol. 77, no. 2, pp. 1-16, 2013.
- [18] L. Liu, J. Xu, S. S. Liao, and H. Chen, "A real-time personalized route recommendation system for self-drive tourists based on vehicle to vehicle communication," *Expert Systems with Applications*, vol. 41, no. 7, pp. 3409-3417, 2014.
- [19] J. P. Lucas, N. Luz, M. N. Moreno, R. Anacleto, A. Almeida Figueiredo, and C. Martins, "A hybrid recommendation approach for a tourism system," *Expert Systems with Applications*, vol. 40, no. 9, pp. 3532-3550, 2013.
- [20] Z. Bahramian, R. Ali Abbaspour, and C. Claramunt, "A CONTEXT-AWARE tourism recommender system based ON a spreading activation method," *ISPRS-International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. XLII-4/W4, pp. 333-339, 2017.
- [21] M. Nilashi, K. Bagherifard, M. Rahmani, and V. Rafe, "A recommender system for tourism industry using cluster ensemble and prediction machine learning techniques," *Computers & Industrial Engineering*, vol. 109, pp. 357-368, 2017.
- [22] I. Cenamor, T. de la Rosa, S. Núñez, and D. Borrajo, "Planning for tourism routes using social networks," *Expert Systems with Applications*, vol. 69, pp. 1-9, 2017.
- [23] K. Meehan, T. Lunney, K. Curran, and A. McCaughey, "Aggregating social media data with temporal and environmental context for recommendation in a mobile tour guide system," *Journal of Hospitality and Tourism Technology*, vol. 7, no. 3, pp. 281-299, 2016.
- [24] Z. Shi and A. B. Whinston, "Network structure and observational learning: evidence from a location-based social network," *Journal of Management Information Systems*, vol. 30, no. 2, pp. 185-212, 2014.
- [25] Q. Lu, "Mobile e-commerce precision marketing model and strategy based on LBS," *E-Business Journal*, no. 4, pp. 20-21, 2014.
- [26] P. J. Danaher, M. S. Smith, K. Ranasinghe, and T. S. Danaher, "Where, when, and how long: factors that influence the redemption of mobile phone coupons," *Journal of Marketing Research*, vol. 52, no. 5, pp. 710-725, 2015.
- [27] N. M. Fong, Z. Fang, and X. Luo, "Geo-conquesting: competitive locational targeting of mobile promotions," *Journal of Marketing Research*, vol. 52, no. 5, pp. 726-735, 2015.
- [28] S. A. Shad and E. Chen, "Precise location acquisition of mobility data using cell-id," *International Journal of Computer Science Issues*, vol. 9, no. 3, 2012.

Research Article

Location-Based Test Case Prioritization for Software Embedded in Mobile Devices Using the Law of Gravitation

Xiaolin Wang ^{1,2}, Hongwei Zeng,¹ Honghao Gao ^{3,4}, Huaikou Miao ¹
and Weiwei Lin ^{1,2}

¹School of Computer Engineering and Science, Shanghai University, Shanghai 200444, China

²Shanghai Key Laboratory of Computer Software Evaluating and Testing, Shanghai 200444, China

³Computing Center, Shanghai University, Shanghai 200444, China

⁴Shanghai Key Laboratory of Intelligent Manufacturing and Robotics, Shanghai 200072, China

Correspondence should be addressed to Honghao Gao; gaohonghao@shu.edu.cn

Received 19 July 2018; Revised 26 October 2018; Accepted 25 November 2018; Published 2 January 2019

Guest Editor: Jaegeol Yim

Copyright © 2019 Xiaolin Wang et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Considering that some intelligent software in mobile devices is related to location of sensors and devices, regression testing for it faces a major challenge. Test case prioritization (TCP), as a kind of regression test optimization technique, is beneficial to improve test efficiency. However, traditional TCP techniques may have limitations on testing intelligent software embedded in mobile devices because they do not take into account characteristics of mobile devices. This paper uses a smart mall as a scenario to design a novel location-based TCP technique for software embedded in mobile devices using the law of gravitation. First, *test gravitation* is proposed by applying the idea of universal gravitation. Second, a specific calculation model of *test gravitation* is designed for a smart mall scenario. Third, how to create a faulted test case set is designed by the pseudocode. Fourth, a location-based TCP using the law of gravitation algorithm is proposed, which utilizes test case information, fault information, and location information to prioritize test cases. Finally, an empirical evaluation is presented by using one industrial project. The observation, underlying the experimental results, is that our proposed TCP approach performs better than traditional TCP techniques. In addition, besides location information, the level of devices is also an important factor which affects the prioritization efficiency.

1. Introduction

Nowadays, the Internet of Things (IoT) develops more and more widely [1]. It is based on wireless sensor networks (WSNs) which combine intelligent software and sensor devices and makes smart home and smart city possible [2, 3]. With the development of hardware (chips) and software (intelligent systems), smart mobile devices (such as smart dust) are gradually emerging, which integrate sensors, processors, intelligent software, and communications. Smart mobile devices not only have an ability to transmit and monitor information but also perform sophisticated intelligent information processing and intelligent prediction by the intelligent software [4] and location-based service (LBS) [5]. The software in each device has its specific functions. For example, some devices are used to monitor

and process temperature information and others are used to process population information. Devices are provided with location-dependent information and interact with other devices in a location-dependent way. It is location information and software complexity of devices that make software testing face a major challenge in IoT.

Regression testing, reusing test suites, is performed on a modified program to install confidence that the system behaves correctly and that modifications have not adversely affected unchanged portions of the program [6]. Test case prioritization (TCP), sorting test cases depending on some criteria, is a way to increase the efficiency of regression testing [7]. It aims at improving the rate of fault detection. Traditional TCP techniques mainly focus on the algorithm design for testing software to improve test prioritization efficiency. However, in IoT, traditional TCP techniques have

limitations because they do not take into account the characteristics of hardware devices, such as location information.

The law of gravitation, according to Newton’s *Philosophiæ Naturalis Principia Mathematica* [8], indicates that there is a force of gravitational attraction existing between any two objects, which is given by the following equation:

$$F = G \frac{m_1 m_2}{r^2}, \quad (1)$$

where G is the universal gravitational constant, m_1 is the mass of one object, m_2 is the mass of the other object, r is the radius of separation between the center of masses of each object, and F is the force of attraction between two objects. The universal gravitation has been applied to the field of data analysis. For example, many research studies make data gravitation (simulating the universal gravitation) applicable to machine learning [9–11]. In IoT, if we can utilize the law of gravitation to prioritize test cases, will it improve the test efficiency?

In this paper, a new location-based TCP using the law of gravitation technique is developed to solve TCP problem of software embedded in mobile devices. This technique is designed for adapting to a smart mall scenario. It is not just a test case prioritization approach but additionally enables to make characteristics of devices utilized, thereby allowing the order of test cases to be beneficial for test efforts. *Test gravitation* is defined in this technique. Under this definition, a specific calculation model of *test gravitation* for a smart mall scenario is designed. First, it calculates the masses of each test case and each faulted test case. The creation of a faulted test case set is related with the occurred faults which detected by those preselected test cases that test different-location-area device representatives. Then, the distance between two specific test cases can be calculated according to location information of devices. For each test case, *test gravitation* is calculated from this test case to each faulted test case. Finally, test cases are prioritized based on *test gravitation*.

The contributions of this work include the following:

- (i) *Test gravitation* is proposed based on the law of gravitation. A specific calculation model of *test gravitation* adapted to a smart mall scenario is given. Specially, the creation of the faulted test case set used in the calculation of *test gravitation* is designed in detail.
- (ii) A location-based TCP using the law of gravitation technique is proposed, and its algorithm is designed by the pseudocode. Its feasibility is illustrated with a small example.
- (iii) An empirical evaluation is presented by using one industrial project. In addition, it discusses whether different evaluation metrics (with or without considering severities of faults) will influence the experimental conclusions. It is also discussed what factors affect the prioritization efficiency.

The rest of this paper is organized as follows: Section 2 describes test case prioritization problem, traditional TCP techniques, special TCP techniques, and TCP problem in

a smart mall scenario. Section 3 presents a location-based TCP using the law of gravitation method and simulates its feasibility with an example. Section 4 describes an empirical evaluation and analyzes the results. Section 5 discusses some related work on test case prioritization and mobile application testing. Finally, the conclusions and future work are given in Section 6.

2. Background

2.1. Test Case Prioritization Methodology. Regression testing, attempting to validate modified version P' of the original program P , checks the results for conformance with requirements [12]. Many techniques have been proposed to improve the cost-effectiveness of regression testing. Test case prioritization is one of these approaches, which rearranges test cases to increase the rate of fault detection during the whole regression testing.

Test case prioritization problem is a research hotspot in the field of software testing. It sorts test cases by using some criteria to detect more faults as fast as possible. A complete definition of TCP problem was first proposed by Rothermel et al. [13]:

Given a test suite already selected (T), the set of all possible prioritizations (orderings) of T (PT), and an objective function from PT to the real numbers (f), which yields an award value for that ordering.

Problem. Find $T' \in PT$ such that $(\forall T'')(T'' \in PT)(T'' \neq T') [f(T') \geq f(T'')]$.

Many test case prioritization techniques have been proposed during the past two decades. Elbaum and Rothermel et al. [13–16] discussed test case prioritization techniques of the fine-grained entity, such as coverage prioritization (statement or branch coverage, etc.) and fault-exposing-potential (FEP) prioritization. Meanwhile, total strategy and additional strategy are proposed [15]. Both are built on a Greedy algorithm which selects a local optimal solution within the search space at each round. Srikanth et al. [17–19] proposed a value-driven approach called PORT which does not require structural coverage information. The PORT algorithm was based on four factors: customer priority, requirements volatility, implementation complexity, and fault proneness of the requirements. Arafeen and Do [20] proposed a test case prioritization technique using requirement-based clustering. It incorporated traditional code analysis information which could improve the effectiveness of test case prioritization techniques.

2.2. Traditional TCP Techniques. Existing TCP methods—prioritizing test cases based on coverage [16] or requirement [19], or even time-aware [21, 22]—are all based on the optimization of software system itself. That is to say, traditional TCP methods focus on improvement of methods themselves. The factors they consider are based on characteristics of software and do not involve characteristics of hardware. Most of the software they test is also cross-platform web application.

There are several classical test case prioritization techniques, introduced as follows:

Random prioritization [16]. Random prioritization orders test cases randomly. It is simple and convenient, but unstable.

Total coverage prioritization [16]. It orders test cases based on the descendent number of units covered by these test cases. When multiple test cases cover the same number of units, the order is determined randomly.

Additional coverage prioritization [16]. It orders test cases to achieve maximized coverage as early as possible. It first picks the test case with the greatest coverage and then successively adds those test cases that cover the most yet uncovered parts.

Prioritization of Requirements for Test (PORT) [17–19]. It orders test cases based on the descending order of weighted priority (WP) values so that the test case with a higher WP value will be ordered in the front.

Optimal prioritization [16]. It prioritizes test cases using the faults, and it can obtain the ordering of test cases that maximizes a test suite's rate of fault detection. It provides an upper bound on the effectiveness of the other heuristics.

2.3. TCP Techniques Utilizing the Execution Information.

When a test case has been executed, it generates execution information, such as its fault detection. As regression testing becomes more complex, scholars have considered the impact of execution information of test history to the current test prioritization.

2.3.1. History-Based TCP. A history-based TCP technique [23] sorts test cases according to the selection probabilities calculated from test history. It defines the selection probability of each test case as follows:

$$\begin{cases} P_0 = h_1, \\ P_k = \alpha h_k + (1 - \alpha)P_{k-1}, \end{cases} \quad (2)$$

where P is selection probability, h is test history, and α is a smoothing constant used to weight individual histories.

Three test histories (based upon each test case's execution history, its fault detection, and the program entities it covers) have been investigated on the effect of test prioritization. Their experimental results show that historical information may be useful in reducing costs and increasing the effectiveness of long-running regression testing processes.

2.3.2. Adaptive TCP. As a main method of dynamic programming [24], the adaptive idea is also used in test case prioritization. Two types of adaptive TCP techniques are introduced here. They all take advantage of the impact of occurred faults to prioritize test cases in current test round.

(1) *Adaptive TCP guided by output inspection.* An adaptive test case prioritization guided by output inspection [25],

which combines the test-case scheduling process and the test-case execution process, prioritizes test cases as the following process:

First, it calculates the initial fault-detection capability of all test cases based on the execution information of the previous output and then selects a test case t with the largest fault-detection capability. Second, t is executed on the modified program, and it records the output of t . Third, it modifies the fault-detection capability of remaining unselected test cases based on t 's output and selects the test case with the largest modified fault-detection capability. Fourth, it repeats the preceding two steps until all the test cases have been prioritized and run.

(2) *TCP based on adaptive sampling strategy.* TCP techniques using cluster filtering [26, 27] select and prioritize test cases as the following process: first, it partitions the test suite based on cluster analysis; then, it selects test cases according to sampling strategy; finally, it prioritizes the selected test cases. In the sampling strategies, the *adaptive sampling* strategy is that it first initially selects one execution at random from each cluster and then all others of its cluster are selected if the first one selected from the cluster is a failure.

2.4. TCP Problem in a Smart Mall Scenario. Figure 1 is a simple distribution diagram of mobile devices for a smart mall scenario. In the figure, a wireless transmitter icon represents a smart mobile device mentioned and studied in this paper. The cloud icon represents central processing. A person with a mobile phone represents a handheld mobile device. Among them, white devices are distributed around specific locations (stores) to monitor and process specific-location information. Black devices are distributed in the middle of the mall to monitor certain types of information and perform distributed information processing. Each mobile device in the mall has its own unique function; that is, its internal intelligent software achieves specific requirements. Integrated testing of the software in devices throughout the mall becomes extremely complicated. For example, in the case like Figure 2, each restaurant has a mobile device that manages information about this restaurant. It can, via Internet, monitor the number of incoming customers/remaining seats, the number of dishes, the temperature, etc., and pushes location preferences, food preferences, etc. to guests (other mobile devices) entering the restaurant. In the hall of the mall, there are restaurant-proxy mobile devices that collect real-time restaurant and people data, via wireless network. It also intelligently pushes the best restaurants (vacant, near, etc.) to the mall customers (other mobile devices) at the current moment via Internet. This can schedule mall customers in real time, which may avoid occurring the case that all customers crowded in front of one restaurant. All data need to be transmitted to the control center for large-scale data processing via wireless network.

When discussing TCP problem in a smart mall scenario, traditional TCP methods can be improved by adding

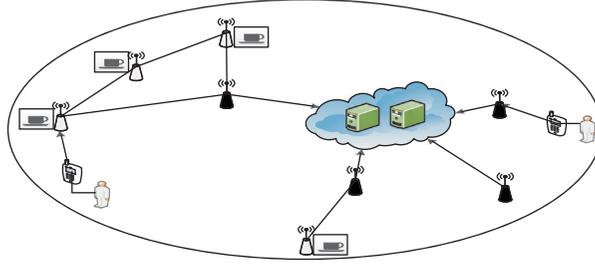


FIGURE 1: A simple mobile device distribution for a smart mall.

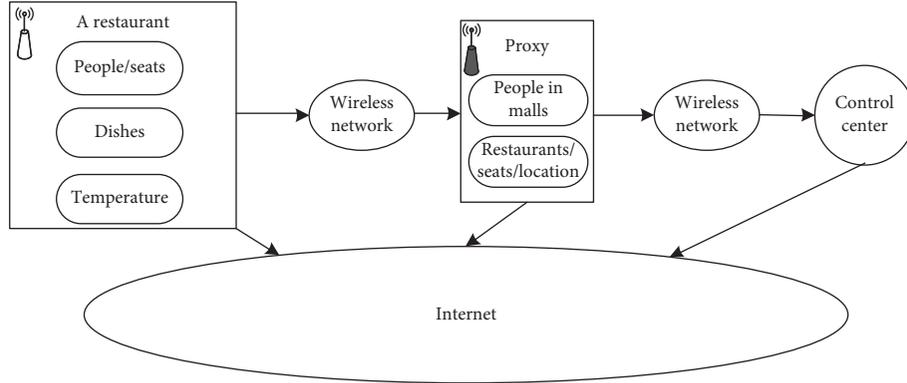


FIGURE 2: Simple software functions and communication structures for restaurant applications.

location information in sorting test cases to adapt the test order for the new scenario. Mobile devices are located in different locations, making them communicate more frequently (functionally interact more closely) with other close-range devices. According to distances between devices, the correlation between functions of intelligent software attached to devices is also strong or weak. As shown in Figure 1, software functions of the black device on the left side should have a greatest relationship with software functions of the other three white devices which communicate with this black device. Test cases test software functions of mobile devices. We set the granularity of a test case as testing all of the functional requirements of a mobile device. In this way, the node graph between traditional test cases becomes a node graph between actual mobile devices (Figure 3). In Figure 3, the left ellipse is a test-case node graph where a circle icon indicates a test case, and the right ellipse is a device node graph where a square icon indicates a mobile device. r represents the distance between test cases or devices. The virtual distance between test cases is mapped to the actual distance between devices.

3. Location-Based TCP Using the Law of Gravitation

This section combines test case information, fault information, and location information to propose a new location-based TCP technique using the law of gravitation.

3.1. Test Gravitation. *Test gravitation (TG)* is introduced to simulate the universal gravitation in our method. *Test*

gravitation F between two test cases t_a and t_b can be defined as follows:

$$F = G \frac{m(t_a)m(t_b)}{r^2}, \quad (3)$$

where G is the test gravitational constant, $m(t_a)$ and $m(t_b)$ are the masses of t_a and t_b , respectively, and r is the distance between t_a and t_b .

G is related to the environment of regression testing. If the criterion of m and r is certain, G should be unique. In this paper, we will not research its influence on the proposed method. So, G is set as 1.

Different attributes of a test case can represent different substances that make up this test case. If this test case detected faults, the attributes of a fault, which is as another type of substances, are also included to make up this test case. The weight of two types of attributes (substances) together makes up the total mass of this test case.

Definition 1. Test case mass m . The mass $m(t)$ of a test case t is defined as follows:

$$m(t) = \mu \sum_{i=1}^n w_i + (1 - \mu) \sum_{j=1}^e \sum_{i=1}^k \bar{w}_i, \quad (4)$$

where n is the total number of attributes of t , w_i is the weight of i th attribute of t , e is the total number of faults which t detects, k is the total number of attributes of a fault, \bar{w}_i is the weight of i th attribute of this fault, and μ is a smoothing constant, which is $0 < \mu \leq 1$.

For instance, in the implementation process, w_1 can be represented as the coverage of a test case and w_2 can be

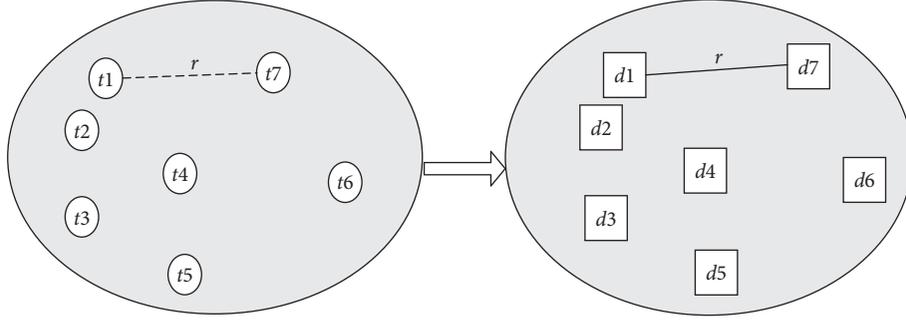


FIGURE 3: A test-case node graph maps to a device node graph.

represented as the importance level of a test case; $\overline{w_1}$ can be represented as the location level of a fault, $\overline{w_2}$ can be represented as the severity of a fault, and so on.

Because faults cannot be known in advance during the actual testing process, there are two ways to obtain faults and their attributes. One way is presetting faults, which can be given based on expert decision or deep learning; the other way is utilizing occurred faults.

Definition 2. Distance r . It indicates the distance between two test cases, denoted as $r(t_a, t_b)$.

For instance, r can be calculated according to the business level (tree relationship) between test cases or according to the spatial distance between devices they are located.

3.2. TG Calculation Model. In a smart mall scenario, according to the above definitions, we design a specific calculation model of TG to make preparations for prioritizing test cases. This model calculates a force F from a test case to a faulted test case.

- (1) Importance level (TI) of a test case t is selected as the only attribute of t . TI is determined by the functional level of the mobile device (DL) which t tests. The mass of t is

$$m(t) = DL(d), \quad (5)$$

where d is the device tested by t . DL is divided into 5 levels. It can use a linear assignment, such as n ($n \in \{1, 2, 3, 4, 5\}$), or a nonlinear assignment, such as a^n ($a \in \mathbb{Z}, n \in \{1, 2, 3, 4, 5\}$).

- (2) Fault severity (FS) is selected as the only attribute of a fault f . The values of FS and $TI(DL)$ compose $m(t_{(f)})$; that is

$$m(t_{(f)}) = \mu DL(d_{(f)}) + (1 - \mu) \sum_{j=1}^e FS(f_j), \quad (6)$$

where $d_{(f)}$ is the device tested by the faulted test case $t_{(f)}$ and f_j is j th fault detected by $t_{(f)}$. FS is divided into 5 levels, like DL .

Occurred faults, detected by preselected test cases in current test round, are used to create a set of faulted

test cases (FTS). The formation of a FTS will be described in detail in Section 3.3.

- (3) Spatial location distance of devices is selected to calculate r . r is the 3-dimensional Euclidean distance between a device which a test case tests and a device which a faulted test case tests, as shown in Figure 4. It is defined as follows:

$$r(t_i, t_{(f)}) = ED(d_i, d_{(f)}), \quad (7)$$

where d_i is the device whose software is tested by t_i , and $d_{(f)}$ is the device tested by $t_{(f)}$.

- (4) From the above, a specific calculation model of TG between a test case t_i and a faulted test case $t_{(f)}$ is as follows:

$$\left\{ \begin{array}{l} F(t_i, t_{(f)}) = G \frac{m(t_i)m(t_{(f)})}{r(t_i, t_{(f)})^2}, \\ G = 1, \\ m(t_i) = DL(d_i), \\ m(t_{(f)}) = \mu DL(d_{(f)}) + (1 - \mu) \sum_{j=1}^e FS(f_j), \\ 0 < \mu < 1, \\ r(t_i, t_{(f)}) = ED(d_i, d_{(f)}), \end{array} \right. \quad (8)$$

where d_i is the device whose software is tested by t_i , $d_{(f)}$ is the device tested by $t_{(f)}$ which detected e faults, and f_j is j th fault detected by $t_{(f)}$.

3.3. Faulted Test Case Set. Occurred faults which detected by some preselected test cases in the current test round are used as the faults mentioned in Section 3.2, so how to collect these occurred faults to create a FTS is an important step. The fault attaches to the device. We use clusters of devices to obtain a FTS . Algorithm 1 describes a clustering process of devices. Euclidean distance is used as the dissimilarity metric.

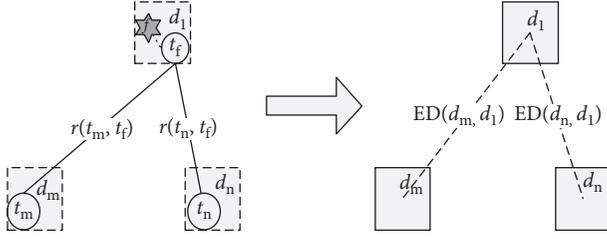


FIGURE 4: The distance between a test case and a faulted test case maps to the distance between devices.

After devices clustered, one device is selected randomly from each cluster as the representative of this cluster. Test cases that test these representatives are put into a test subset ST . ST is executed. If faults occurred, the test cases which detected these faults are combined into a FTS . Algorithm 2 shows the pseudocode of this process.

3.4. Location-Based TCP Using the Law of Gravitation Algorithm

Definition 3. Test case priority P . It indicates the priority of a test case in the execution order. The priority P is defined as follows:

$$P = \sum_{i=1}^h F_i, \quad (9)$$

where h is the number of faulted test cases and F_i is the force F of this test case to the i th faulted test case. The larger the P value is, the earlier this test case will execute.

Algorithm 3 shows the location-based TCP using the law of gravitation approach. Its input is a test suite T . Its output is the prioritized test order T' . First, m of each test case t is calculated according to the level of a mobile device which t tests. Second, a faulted test case set FTS is created according to algorithms 1 and 2. m of each faulted test case $t_{(f)}$ is calculated based on both FS and DL . Third, the distance r between each t and each $t_{(f)}$ can be calculated according to location information of devices. Fourth, for each t , the force F is computed from this t to each $t_{(f)}$. Fifth, the priority P of each t is calculated according to F . Finally, test cases are sorted in descending order of P to obtain a prioritized test execution order T' .

3.5. Example for Simulating Smart Mall. We simulate a smart mall scenario with Figure 5 to explain how to prioritize test cases. There are five mobile devices (d_1-d_5) in the figure, shown by squares. Each device is tested by a test case for its internal intelligent software functions. So, there are five test cases (t_1-t_5), shown by circles. Assume that the devices are clustered into 2 clusters: $c_1(d_1, d_2, d_3)$ and $c_2(d_4, d_5)$. d_2 and d_4 are extracted randomly to be device representatives, and a subset $ST \{t_2, t_4\}$ is formed. After ST run, two faults (f_1 and f_2) are found by t_2 , which are shown by stars in the figure. A dashed line in the figure shows the 3-dimensional Euclidean distance between two devices.

In the above example, let us consider a test case prioritization problem defined over a set of five test cases $T(t_1, t_2, t_3, t_4, t_5)$ with a set of one faulted test case $FTS(t_{(f)1})$ from Table 1. From Figure 5, according to the location of devices, the distances between test cases are obtained, as in Table 1. We suppose that all faults (including their Severity levels) detected in current testing round are shown in Table 2.

We take t_1 as a sample and calculate the force F of t_1 to $t_{(f)1}$, as $F_1 = 0.135$. According to Equation (9), we get the priority value of t_1 which is $P_1 = 0.135$. Similarly, the priority P of t_2, t_3, t_4 , and t_5 are $P_2 = 22.5, P_3 = 0.28125, P_4 = 0.0002$, and $P_5 = 0.000225$. The prioritization order is $t_2-t_3-t_1-t_5-t_4$, and the $APFDc$ [28] value of this order is 78.57%. According to Table 2, the optimal prioritization sorts test cases as the order $t_2-t_3-t_1-t_5-t_4$ (or $t_2-t_3-t_1-t_4-t_5$), whose $APFDc$ value is 78.57%. The random prioritization sorts test cases as one order $t_5-t_1-t_3-t_2-t_4$, whose $APFDc$ value is 41.43%. It can be seen that the effect of our location-based TCP using the law of gravitation has a good effect, which is even consistent with the optimal prioritization.

4. Empirical Evaluation

To investigate the effectiveness of the method, called location-based TCP, using the law of gravitation (L-TCP from now on), an empirical evaluation is performed in terms of the following research questions:

- (i) RQ1: Is L-TCP approach more effective in the rate of fault detection than other traditional prioritization techniques?

This research question aims at understanding whether the L-TCP method can detect faults earlier than other traditional test case prioritization techniques. To answer this question, this paper applies four traditional TCP techniques for comparison.

- (ii) RQ2: When evaluating the efficiency of techniques, is there any difference in the experimental conclusions for whether or not considering faults severities?

Whether or not to consider severity of a fault will undoubtedly make a difference in the judgment of the prioritization effect. This research question mainly discusses the influence of two evaluation metrics on the experimental conclusions.

- (iii) RQ3: In addition to location information, what other factors the prioritization efficiency is also related to?

In the smart mall scenario, test prioritization efficiency of software may be related to the information of mobile devices. This research question combines analyses of the above two questions to discuss factors that influence the efficiency of prioritization.

4.1. Object. The object used in this experimental study is a real industrial project which is for chip testing and has

```

Input:  $D = \{d_1, d_2, \dots, d_n\}$ ,  $k$  //a device set, and the number of clusters
Output:  $C$  //a set of  $k$  clusters
(1)  $C = \emptyset$ ;
(2) put each  $d \in D$  as a cluster  $c$ ;
(3) add all clusters into  $C$ ; //Initialization: get a single-cluster set  $C$ 
(4) Do //Iteration: make clusters merge.
(5)   For each  $c_i, c_j \in C$ 
(6)     If ( $c_i$  and  $c_j$  have the minimum 3-dimensional Euclidean distance)
(7)       merge  $c_i$  and  $c_j$  into a new cluster  $c_{new}$ ;
(8)       delete  $c_i$  and  $c_j$  from  $C$ ;
(9)        $C = C \cup \{c_{new}\}$ ;
(10)    End if
(11)  End for
(12) Until The number of clusters in  $C$  is  $k$  //Break condition

```

ALGORITHM 1: Clustering.

```

Input:
   $C = \{c_1, c_2, \dots, c_m\}$  //a set of device clusters
   $T = \{t_1, t_2, \dots, t_m\}$  //a set of test cases
Output:
   $FTS$  //a set of faulted test cases
(1)  $FTS = \emptyset$ ;
(2) While ( $FTS == \emptyset$ ) do
(3)    $ST = \emptyset$ ;
(4)   For each  $c \in C$ 
(5)     Randomly select one  $d$  (i.e.,  $d_s$ ) from  $c$ ;
(6)     Select the test case  $t \in T$  (i.e.,  $t_s$ ) which tests the software of  $d_s$ ;
(7)      $ST = ST \cup \{t_s\}$ ;
(8)   End for
(9)   For each  $t \in ST$ 
(10)    Execute  $t$ ;
(11)    if ( $t$  detects faults) then
(12)      put  $t$  into  $FTS$ ;
(13)    End if
(14)  End for
(15) End while
(16) Return  $FTS$ 

```

ALGORITHM 2: Creating a *FaS*.

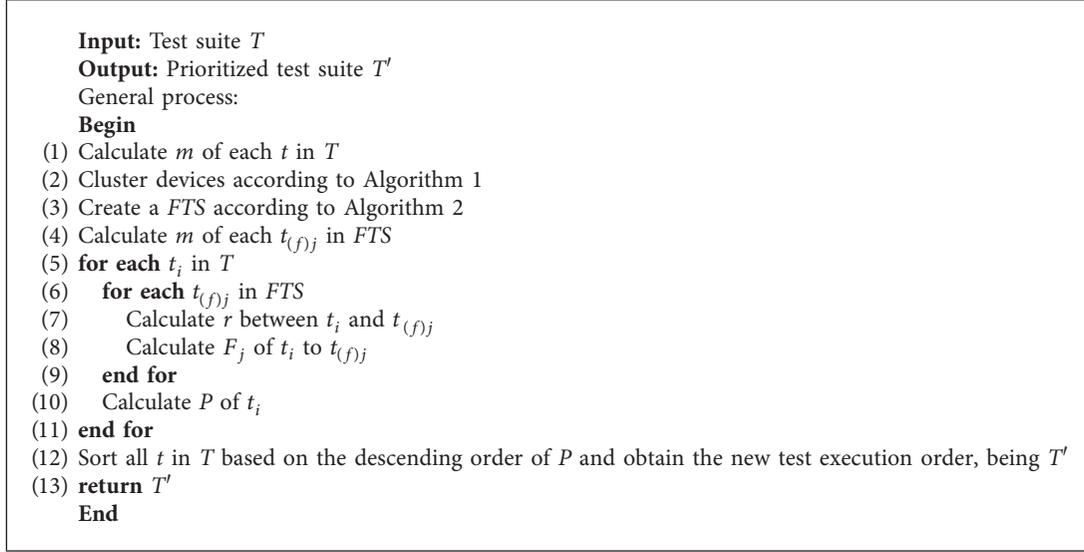
approximately 140,000 lines of codes (LOC), totally. It has many versions, and each version has a few of requirements. The test suite of each version is relatively small. The granularity of test cases is coarse-grained. That is to say, each test case may contain dozens or even hundreds of test scripts, but it tests only one chip function (requirement). These features can be used to simulate test data for a smart mall scenario. First of all, functions of the hardware chip are similar to those of smart devices, so characteristics of the faults may be similar, too. Second, each test case covers only one specific requirement, which can simulate to test one mobile device. Third, there are many rounds (versions) of regression testing, and there are new test cases introduced in each round, which can simulate a step-by-step integration testing environment for the smart mall scenario. The project data include the number of test cases, the functional requirements covered by test cases, the faults detected by test

cases, the fault severities, and so on. We use this project data as a basis and then simulate the distance data between devices. Finally, they are formed into a complete data required for this experiment. There are six versions chosen for this experiment. The basic information is shown in Table 3.

4.2. Variables and Measures

4.2.1. Independent Variables. To address our research questions, one independent variable is manipulated: test case prioritization technique. Besides our proposed L-TCP approach, the following traditional test case prioritization approaches are also implemented for comparison.

- (i) Random (R): this technique uses random prioritization technique to order test cases without using location information of devices in prioritization.



ALGORITHM 3: Location-based TCP using the law of gravitation.

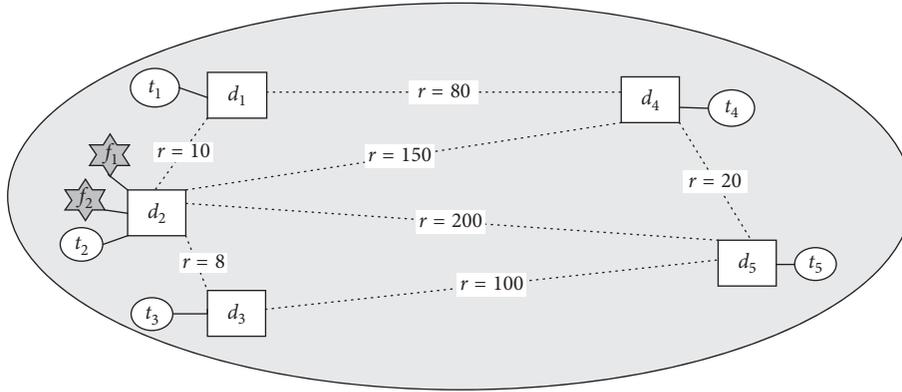


FIGURE 5: A testing example simulating a smart mall.

TABLE 1: Information of T and FTS .

T	m	FTS	m	$r(t, t_{(f)})$	$t_{(f)1}$
t_1	3	$t_{(f)1}$	4.5	t_1	10
t_2	5			t_2	1
t_3	4			t_3	8
t_4	1			t_4	150
t_5	2			t_5	200

TABLE 2: All faults detected by test cases.

	f_1	f_2	f_3	f_4	FS
t_1			√		f_1 3
t_2	√	√			f_2 1
t_3				√	f_3 1
t_4					f_4 2
t_5					

(ii) Total coverage (TC): this technique uses total coverage prioritization technique to order test cases without using location information of devices in prioritization.

TABLE 3: Basic information of ChipTest.

Version	v-9	v-10	v-11	v-12	v-13	v-14
New requirements	5	9	5	8	7	16
New test cases	9	9	5	8	7	16
Total test cases	9	18	23	31	38	54
Faults	7	8	5	8	7	8

- (iii) Additional coverage (AC): this technique uses additional coverage prioritization technique to order test cases without using location information of devices in prioritization.
- (iv) Requirement prioritization (PORT): this technique uses prioritization of requirements for test technique to order test cases without using location information of devices in prioritization.

4.2.2. *Dependent Variable and Metric.* Details on the measures for the dependent variables of these experiments are given here.

APFD. To measure how rapidly a prioritized test suite detects faults, average percentage of fault detection (*APFD*) is used as the dependent variable.

APFD [28], the weighted average of the percentage of faults detected, focuses on the rate of fault detection during the testing life of a test suite. It assumes that the faults severities are equivalent. The equation of *APFD* is as follows:

$$APFD = 1 - \frac{TF_1 + TF_2 + \dots + TF_m}{nm} + \frac{1}{2n}, \quad (10)$$

where n is the number of test cases, m is the number of faults, and TF_i is the index of the first test case that reveals the i th fault in the execution order T . The value of *APFD* varies from 0 to 100%. Since n and m are fixed for any orders, a higher *APFD* value indicates that the faults are detected earlier during the testing process.

APFDc. When considering faults severities, we use *APFDc* [28], the (cost-cognizant) weighted average percentage of faults detected, to reward test case orders proportionally to their rate of units of fault severity detected. We assume that test case costs are identical. The equation of *APFDc* is simplified as follows:

$$APFDc = 1 - \frac{\sum_{i=1}^m (f s_i \times TF_i)}{n \times \sum_{i=1}^m f s_i} + \frac{1}{2n}, \quad (11)$$

where $f s_i$ is the severity of the i th fault and other symbols have the same definition as in the equation of *APFD*.

4.3. Case Study Design. Suppose that there are many smart mobile devices in a mall and each device is responsible for its own unique functions. We now need to test the functionality of their internal intelligence software. We assume that one test case is in charge of testing one device, and it tests all of the software functionality of this device.

First, it collects data. To conduct the comparative experiments, five types of data information are required, including test case, levels of test cases (devices), fault, severities of faults, distance of devices, and coverage information.

The preparation of test case, fault, severities of faults, and coverage information is trivial because it is already available in the original data of the object system. For the preparation of the levels of test cases (devices), we grade test cases according to their name description. The preparation of the distance between mobile devices requires us to give them values by simulating the smart mall scenario.

Second, it performs test case prioritization techniques. This experimental study implements five approaches (TC, AC, R, PORT, and L-TCP) for comparison. Because of the indeterminacy of some prioritization techniques, each technique runs 20 times for each experiment and the average values are presented as results. The smoothing constant is set $\mu = 50\%$.

Third, it calculates *APFD* and *APFDc* for each prioritized test order from each technique. All measure values are compared across different techniques. The results emerging from this comparison are presented in the Section 4.4.

All the experiments are conducted on the same computer which is configured as 64-bit windows 8 operating system, Intel(R) Core(TM) i3-2130 CPU and 4 GB memory.

4.4. Results and Analysis. In this section, we present the results of the experiment(s) and analyze their relevance to our research questions above.

4.4.1. RQ1: Comparison with Traditional TCP Techniques.

Figure 6 shows the box plots of five techniques across all the system versions of ChipTest. The horizontal axis shows versions, and each box in a version presents one TCP technique. The vertical axis presents *APFD* values. Each boxplot shows the median, upper/lower quartile, and max/min *APFD* values achieved by a technique.

From the boxplot, L-TCP, as indicated by *APFD* scores, significantly outperforms the others because its median point reaches up to the highest. Besides L-TCP, PORT performs better with a higher median point. TC, AC, and R have a similar effect, and their median points of *APFD* locate approximately between 40% and 65%.

For instance, let us choose the data of v-9 for analysis. We use $M(\text{Median}, Q1, Q3)$ to denote the median, first, and third quartiles *APFD* values for each technique and $M1-M5$ to denote the five techniques: TC, AC, R, PORT, and L-TCP, respectively. So, results of the five techniques are $M1(38.89, 29.37, 48.02)$, $M2(34.13, 29.37, 43.25)$, $M3(38.1, 31.35, 46.43)$, $M4(45.24, 40.88, 46.83)$, and $M5(78.57, 78.57, 78.57)$, respectively, which clearly indicates that L-TCP overall performs better than the others.

For evaluating the confidence level of the observed results, we test their statistical significance. First, a single sample K-S test is used to check the normal distribution of the data of each technique from 120 executions (20 running \times 6 versions). The significance level is $\alpha = 0.05$. Their results are as follows: in Table 4, the first row is the names of TCP approaches mentioned above. The second row shows the judge of normal distribution for TC, AC, R, PORT, and L-TCP. The third row shows their significance probability values under the null hypothesis.

From Table 4, their results accept the null hypothesis (p values are all greater than 0.05). So, *APFD* values of the five prioritization techniques satisfy a normal distribution.

Next, the paired-samples t -test is employed to obtain sufficient statistical evidence. f_1 and f_2 are defined as the values of *APFD*, which are prioritized by two prioritization approaches, respectively.

The following two hypotheses are considered:

$H_0: f_1 = f_2$, if two techniques have the same effectiveness in the rate of fault detection.

$H_1: f_1 > f_2$, if f_1 is significantly better than f_2 .

If the p value is less than the significance level ($\alpha = 0.05$), we can reject the null hypothesis and accept the alternative hypothesis.

Table 5 reports the results of statistical testing by using the data from 120 executions. Their results show that L-TCP

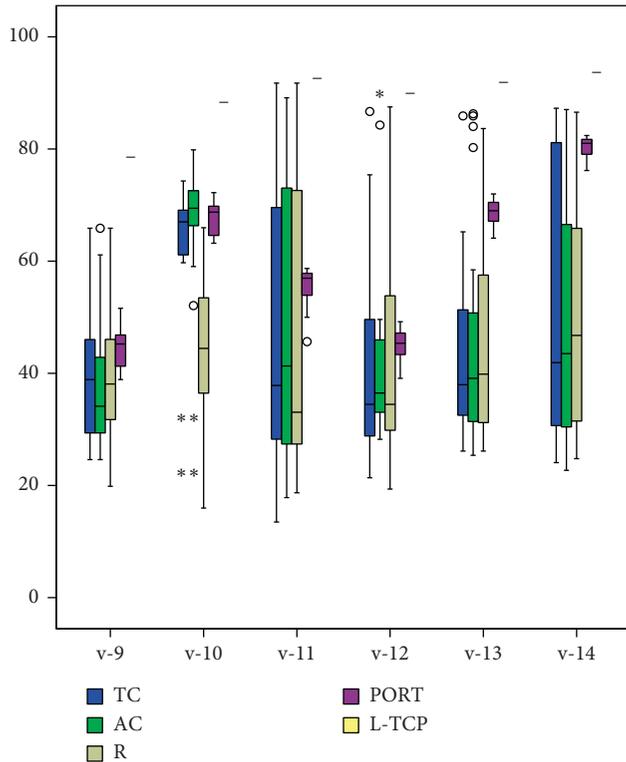


FIGURE 6: *APFD* distributions of different techniques in different versions. *Discrete value; “open circle” indicates an extreme value; –median value.

TABLE 4: Sample K-S test results of *APFD* values of the five TCP techniques.

	TC	AC	R	PORT	L-TCP
Normal distribution?	Y	Y	Y	Y	Y
Sig.	0.976	0.709	1	0.977	0.773

TABLE 5: Statistical test results from comparing *APFD* values of L-TCP to four opponent techniques.

	L-TCP	
	<i>p</i> value	<i>t</i>
TC	0.000	22.947
AC	0.000	21.728
R	0.000	25.486
PORT	0.000	28.295

is statistically significantly better than other TCP techniques because its *t* values are greater than 0 and *p* values are less than 0.05.

For instance, compared with TC, the *p* value of L-TCP equals 0.000 and its *t* value equals 22.947, so we can reject the null hypothesis that L-TCP and TC have the same effectiveness in the rate of fault detection and accept the alternative hypothesis that L-TCP is significantly better than TC.

4.4.2. *RQ2: Effects of Different Evaluation Metrics.* In the evaluation metrics of fault-detection rate, *APFD* is one that does not consider faults severities and *APFDc* is one that considers faults severities.

Figure 7 shows *APFD* and *APFDc* distributions of different techniques in different versions. As can be seen from Figure 7, L-TCP has highest values in both *APFD* and *APFDc* evaluations. That is to say, the prioritization effect of L-TCP is the best among other techniques, regardless of whether or not faults severities are considered in evaluation. In addition, the trend of other techniques is similar in both evaluations (except in version 9 and version 11); that is, PORT is a second best technique besides L-TCP. In version 9, *APFD* evaluation shows that effects of TC, AC, R, and PORT are similar, as shown in Figure 7(a). However, in *APFDc* evaluation, PORT is significantly better than the other three techniques in version 9, as shown in Figure 7(b). In version 11, *APFD* evaluation shows that PORT is significantly better than the others, but in the evaluation of *APFDc*, AC is slightly better than PORT.

Therefore, from the results, whether or not considering faults severities will not affect the conclusion of RQ 1, that is, L-TCP is superior to other techniques. It is just that the degree of excellence varies across different metrics.

4.4.3. *RQ3: Factors Affecting Prioritization Efficiency.* From the analysis of the results of RQ1 and RQ2, it can be seen that L-TCP is the best technique to improve the rate of faults detection. In addition to L-TCP, PORT is the second best performing technique.

In-depth analysis shows that, first of all, according to the characteristics of test data, L-TCP mainly affects the prioritization efficiency by location information of devices. That is, in the smart mall scenario, location information is the main factor affecting the test order efficiency of intelligent software embedded in mobile devices. Second, in the smart mall scenario, PORT sorts test cases according to the priority of software functions of mobile devices in this experiment. The functional priority of a device determines the level of a device, and the device level determines the mass of a test case which tests this device. In retrospect, in the smart mall scenario, the *test gravitation* calculation model considers both *device location information* and *device level*, which are the main factors that influence the test prioritization efficiency. So, this may be the reason why L-TCP can achieve better sorting results in this smart mall scenario.

4.5. *Threats to Validity.* In terms of the internal validity, the choice of the smoothing constant μ can affect the results. In this paper, the selection of this parameter has been based on equalization, that is, $\mu = 50\%$. Further investigations can study the effect of the smoothing constant.

The threats to external validity are from the object, its test data, and its faults used by this experimental study. To reduce this threat in the object, the experimental object we select is the system that tests chips, which is an object that is relatively close to the simulated scenario. Moreover, we select multiple successive versions (6 versions) for

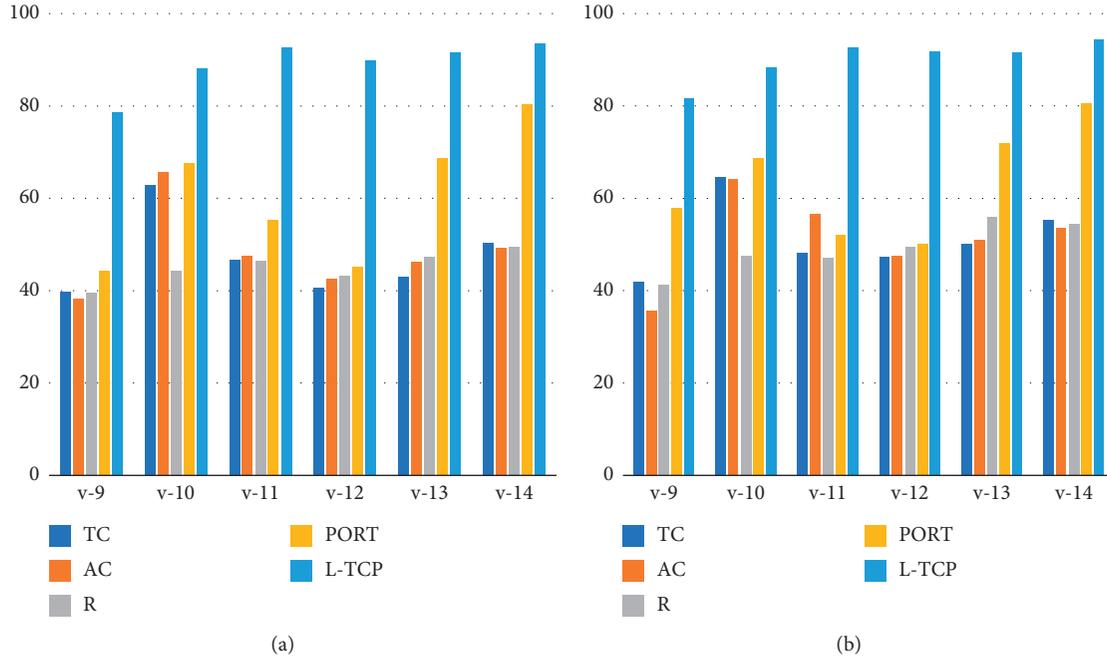


FIGURE 7: *APFD* and *APFDc* distributions of different techniques in different versions: the histograms of (a) *APFD* and (b) *APFDc*.

experiments to simulate step-by-step integration testing of a smart mall scenario. The second external threat lies in the test data in this object. Although the data are relatively real, it is not complete enough for the research in this paper. For incomplete data (such as the lack of distances between devices), we try to simulate the data supplement according to the scenario. The third external threat is the faults. For faults, we use actual real faults in order to be closer to the real scenario.

The threat to construct validity lies in whether the experimental results are measured in a correct way. To reduce this threat, firstly, *APFD* is used to measure the effectiveness of a prioritized test case order since *APFD* can measure the rate of fault detection and has been widely used in the evaluation of the test case prioritization problem. Second, *APFDc* is also used to measure accurately the rate of units of fault severity detected since it considers faults severities.

5. Related Work

Test case prioritization has been an interesting research field for nearly two decades. Rothermel et al. [13] firstly proposed the complete definition of TCP problem which is finding a permutation of T in order to maximize some objective functions. They focus on code-coverage TCP methods at code-level [13–16]. In 2001, Elbaum conducted specific research for TCP metrics, including *APFD* and *APFDc* [28]. *APFD* metric proclaims that all faults have the same severity and all test cases have equal costs. *APFDc*, units of fault severity detected per unit test cost, considers unifying test case costs and fault severities. Their study was primarily focused on white-box testing but not on black-box testing. Zhang et al. [29] considered requirement priorities to TCP and proposed an algorithm called TCP_{RP}_{TC}. The

prioritization technique must predict requirement priorities and test costs before test suite execution, but the prediction was difficult in practice. Chu-Ti et al. [30] presented a history-based TCP method with software version awareness. Yuchi et al. [31] designed and analyzed TCP using weight-based methods for GUI applications. Garg and Datta [32, 33] used test case prioritization in web applications based on modified functionalities or database changes. Saha et al. [34] proposed a fully automated and lightweight test prioritization approach (REPiR) to address the problem of regression test prioritization by reducing it to a standard information retrieval problem so that the differences between two program versions formed the query and the tests constituted the document collection. Some researchers [35] focused on test case prioritization based on mutation analysis. It is an effective method, but the cost is expensive. Another novel refactoring-based approach (RBA) was proposed by Alves et al. [36] which reordered an existing test sequence utilizing a set of refactoring fault models. It promoted early detection of refactoring faults. Wang and Ali et al. [37] proposed a resource-aware multiobjective optimization solution with a fitness function defined based on four cost-effectiveness measures. Prioritizing test cases for the testing of location-aware services was proposed by Zhai et al. [38, 39], and it brings in service selection into a test case prioritization technique for testing the location-based web services.

Mobile application testing is a research direction for testing on mobile devices. However, most of mobile application testing focuses on performance testing or stand-alone testing which sees the software of mobile devices as a stand-alone software. Gao et al. [40] provided a general tutorial on mobile application testing that first examined testing requirements and then looked at current approaches

for both native and Web apps for mobile devices. Muccini, Di Francesco, and Esposito [41] investigated new research directions on mobile applications testing automation, by answering three research questions. Given the first research question (RQ1) *are mobile applications (so) different from traditional ones, so to require different and specialized new testing techniques?*, the natural answer seems to be *yes, they are*. About (RQ2) *what are the new challenges and research directions on testing mobile applications?*, the challenges seem to be many, related to the contextual and mobility nature of mobile applications. As far as concern (RQ3) *which is the role automation may play in testing mobile applications?*, some potentials for automation have been outlined, being aware that a much deeper and mature study shall be conducted. Dantas et al. [42] proposed a set of testing requirements, elicited using the results of an extensive research on how the testing process for mobile applications is done in the literature and in practice. Morla and Davies [43] created a test environment that supports the evaluation of key aspects of location-based applications without the extensive resource investment necessary for a full application implementation and deployment. Zhang and Adipat [44] proposed a generic framework for conducting usability tests for mobile applications through discussing research questions, methodologies, and usability attributes. Vilkomir [45] evaluated the effectiveness of coverage approaches for selecting mobile devices (i.e., smartphones and tablets) to test mobile software applications. Amalfitano et al [46] addressed the problem of testing a mobile app as an event-driven system by taking into accounts both context events and GUI events. Kim, Choi, and Wong [47] proposed a method to support performance testing utilizing a database established through benchmark testing in emulator-based test environment at the unit test level.

6. Conclusion

This paper proposes a location-based TCP using the law of gravitation approach. It introduces *test gravitation*, which combines three factors (test case information, fault information, and location information), to prioritize test cases. Test case information involves the level of mobile device. Fault information includes the severity of fault. In addition, we use occurred faults to create a faulted test case set. It is obtained in three steps: devices clustering, test subset extraction, and running preselected test cases. Location information involves the actual location of devices. It is used to calculate the 3-dimensional Euclidean distance between two devices. Finally, it experimentally verifies the effectiveness of L-TCP technique in comparison with several traditional test case prioritization techniques.

The experimental results show that the median APFD value of L-TCP is 78.57%, which is higher than the values of the baseline methods. When employing the paired-samples *t*-test, L-TCP's *t* values are greater than 0 and *p* values are less than 0.05. Specially, (1) comparing with TC, the *p* value of L-TCP equals 0.000 and its *t* value equals 22.947; (2) comparing with AC, the *p* value of L-TCP equals 0.000 and its *t* value equals 21.728; (3) comparing with R, the *p* value of

L-TCP equals 0.000 and its *t* value equals 25.486; and (4) comparing with PORT, the *p* value of L-TCP equals 0.000 and its *t* value equals 28.295. These results indicate that L-TCP is statistically significantly better than other TCP techniques and it can detect more faults than others at the same time consumption.

When considering the factor of faults severities during the evaluation, the conclusion that L-TCP is superior to other techniques will not be affected. It is just that its degree of excellence varies across different metrics. In the smart mall scenario, location information of devices is the main factor which influences the prioritization performance. Furthermore, the level of devices is also important.

The next step is to expand the scope of empirical evaluation and try to make the conclusion more accurate. Moreover, how to give an appropriate parameter is also a research direction.

Data Availability

The data used to support the findings of this study are available from the corresponding author upon request.

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (grant nos. 61572306 and 61502294), the IIOT Innovation and Development Special Foundation of Shanghai (grant no. 2017-GYHLW-01037), and the CERNET Innovation Project (grant nos. NGII2017051 and NGII20170206).

Supplementary Materials

The ChipTest Data.xlsx file is the data of the project used in the experimental study of my paper. The results.xlsx file is the detailed results of my experimental study. The result analysis.xlsx file includes the original graphs of the results analysis, which are also shown in my paper. (*Supplementary Materials*)

References

- [1] C. Stergiou, K. E. Psannis, B.-G. Kim, and B. Gupta, "Secure integration of IoT and cloud computing," *Future Generation Computer Systems*, vol. 78, pp. 964–975, 2018.
- [2] G. Han, L. Zhou, H. Wang, W. Zhang, and S. Chan, "A source location protection protocol based on dynamic routing in WSNs for the Social Internet of Things," *Future Generation Computer Systems*, vol. 82, pp. 689–697, 2018.
- [3] L. Atzori, A. Iera, and G. Morabito, "The internet of things: a survey," *Computer networks*, vol. 54, no. 15, pp. 2787–2805, 2010.
- [4] H. Liu, Z. Hu, A. Mian, H. Tian, and X. Zhu, "A new user similarity model to improve the accuracy of collaborative

- filtering,” *Knowledge-Based Systems*, vol. 56, pp. 156–166, 2014.
- [5] P. M. Adams, G. W. B. Ashwell, and R. Baxter, “Location-based services—an overview of the standards,” *BT Technology Journal*, vol. 21, no. 1, pp. 34–43, 2003.
- [6] G. Rothermel and M. J. Harrold, “A safe, efficient regression test selection technique,” *ACM Transactions on Software Engineering and Methodology*, vol. 6, no. 2, pp. 173–210, 1997.
- [7] X. Chen, J. H. Chen, X. L. Ju, and Q. Gu, “Survey of test case prioritization techniques for regression testing,” *Ruan Jian Xue Bao/Journal of Software*, vol. 24, no. 8, pp. 1695–1712, 2013.
- [8] I. Newton, *Philosophiae Naturalis Principia Mathematica*, Royal Society, London, UK, 1st edition, 1687.
- [9] L. Peng, B. Yang, Y. Chen, and A. Abraham, “Data gravitation based classification,” *Information Sciences*, vol. 179, no. 6, pp. 809–819, 2009.
- [10] C. Wang and Y. Q. Chen, “Improving nearest neighbor classification with simulated gravitational collapse,” in *Proceedings of International Conference on Natural Computation*, pp. 845–854, Springer, Changsha, China, August 2005.
- [11] M. Indulska and M. E. Orlowska, “Gravity based spatial clustering,” in *Proceedings of 10th ACM international symposium on Advances in geographic information systems*, pp. 125–130, ACM, New York, NY, USA, November 2002.
- [12] G. Rothermel and M. J. Harrold, “Analyzing regression test selection techniques,” *IEEE Transactions on Software Engineering*, vol. 22, no. 8, pp. 529–551, 1996.
- [13] S. Elbaum, A. G. Malishevsky, and G. Rothermel, “Prioritizing test cases for regression testing,” in *Proceedings of International Symposium on Software Testing and Analysis*, pp. 102–112, ACM, Portland, OR, USA, August 2000.
- [14] G. Rothermel, R. H. Untch, C. Chengyun Chu, and M. J. Harrold, “Prioritizing test cases for regression testing,” *IEEE Transactions on Software Engineering*, vol. 27, no. 10, pp. 929–948, 2001.
- [15] G. Rothermel, R. H. Untch, C. Chu et al., “Test case prioritization: an empirical study,” in *Proceedings of IEEE International Conference on Software Maintenance (ICSM’99)*, pp. 179–188, IEEE, Oxford, England, UK, August–September 1999.
- [16] S. Elbaum, A. G. Malishevsky, and G. Rothermel, “Test case prioritization: a family of empirical studies,” *IEEE Transactions on Software Engineering*, vol. 28, no. 2, pp. 159–182, 2002.
- [17] H. Srikanth and L. Williams, “Requirements-based test case prioritization,” *IEEE Transactions on Software Engineering*, vol. 28, 2002.
- [18] H. Srikanth, L. Williams, and J. Osborne, “System test case prioritization of new and regression test cases,” in *Proceedings of 2005 International Symposium on Empirical Software Engineering*, p. 10, Noosa Heads, Queensland, Australia, November 2005.
- [19] H. Srikanth and S. Banerjee, “Improving test efficiency through system test prioritization,” *Journal of Systems and Software*, vol. 85, no. 5, pp. 1176–1187, 2012.
- [20] M. J. Arafeen and H. Do, “Test case prioritization using requirements-based clustering,” in *Proceedings of 2013 IEEE Sixth International Conference on Software Testing, Verification and Validation (ICST)*, pp. 312–321, IEEE, Luxembourg, March 2013.
- [21] K. R. Walcott, M. L. Soffa, G. M. Kapfhammer et al., “Timeaware test suite prioritization,” in *Proceedings of 2006 International Symposium on Software Testing and Analysis*, pp. 1–12, ACM, Shanghai, China, May 2006.
- [22] S. Alspaugh, K. R. Walcott, M. Belanich et al., “Efficient time-aware prioritization with knapsack solvers,” in *Proceedings of the 1st ACM International Workshop on Empirical Assessment of software Engineering Languages and Technologies: held in conjunction with the 22nd IEEE/ACM International Conference on Automated Software Engineering (ASE) 2007*, pp. 13–18, ACM, New York, NY, USA, 2007.
- [23] J. M. Kim and A. Porter, “A history-based test prioritization technique for regression testing in resource constrained environments,” in *Proceedings of 24rd International Conference on Software Engineering ICSE 2002*, pp. 119–129, IEEE, Orlando, FL, USA, May 2002.
- [24] X. Luo, Y. Lv, R. Li, and Y. Chen, “Web service QoS prediction based on adaptive dynamic programming using fuzzy neural networks for cloud services,” *IEEE Access*, vol. 3, pp. 2260–2269, 2015.
- [25] D. Hao, X. Zhao, and L. Zhang, “Adaptive test-case prioritization guided by output inspection,” in *Proceedings of 2013 IEEE 37th Annual Computer Software and Applications Conference (COMPSAC)*, pp. 169–179, IEEE, Kyoto, Japan, July 2013.
- [26] W. Dickinson, D. Leon, and A. Podgurski, “Finding failures by cluster analysis of execution profiles,” in *Proceedings of 23rd International Conference on Software Engineering*, pp. 339–348, IEEE Computer Society, Toronto, ON, Canada, May 2001.
- [27] D. Leon and A. Podgurski, “A comparison of coverage-based and distribution-based techniques for filtering and prioritizing test cases,” in *Proceedings of International Symposium on Software Reliability Engineering*, pp. 442–453, Denver, CO, USA, November 2003.
- [28] S. Elbaum, A. Malishevsky, and G. Rothermel, “Incorporating varying test costs and fault severities into test case prioritization,” in *Proceedings of 23rd International Conference on Software Engineering ICSE 2001*, pp. 329–338, IEEE, Toronto, Canada, May 2001.
- [29] X. Zhang, C. Nie, B. Xu et al., “Test case prioritization based on varying testing requirement priorities and test case costs,” in *Proceedings of Seventh International Conference on Quality Software QSIC’07*, pp. 15–24, IEEE, Portland, Oregon, USA, October 2007.
- [30] C. T. Lin, C. D. Chen, C. S. Tsai et al., “History-based test case prioritization with software version awareness,” in *Proceedings of 2013 18th International Conference on Engineering of Complex Computer Systems (ICECCS)*, pp. 171–172, IEEE, Singapore, July 2013.
- [31] C.-Y. Huang, J.-R. Chang, and Y.-H. Chang, “Design and analysis of GUI test-case prioritization using weight-based methods,” *Journal of Systems and Software*, vol. 83, no. 4, pp. 646–659, 2010.
- [32] D. Garg, A. Datta, and T. French, “A two-level prioritization approach for regression testing of web applications,” in *Proceedings of 19th Asia-Pacific Software Engineering Conference (APSEC)*, vol. 2, pp. 150–153, IEEE, Hong Kong, China, December 2012.
- [33] D. Garg and A. Datta, “Test case prioritization due to database changes in web applications, Software Testing,” in *Proceedings of 2012 IEEE Fifth International Conference on Verification and Validation (ICST)*, pp. 726–730, IEEE, Montreal, QC, Canada, April 2012.
- [34] R. K. Saha, L. Zhang, S. Khurshid et al., “An information retrieval approach for regression test prioritization based on

- program changes,” in *Proceedings of 37rd International Conference on Software Engineering ICSE 2015*, pp. 268–279, IEEE, Firenze, Italy, May 2015.
- [35] R. Just and F. Schweiggert, “Higher accuracy and lower run time: efficient mutation analysis using non-redundant mutation operators,” *Software Testing, Verification and Reliability*, vol. 25, no. 5–7, pp. 490–507, 2014.
- [36] E. L. G. Alves, P. D. L. Machado, T. Massoni, and M. Kim, “Prioritizing test cases for early detection of refactoring faults,” *Software Testing, Verification and Reliability*, vol. 26, no. 5, pp. 402–426, 2016.
- [37] S. Wang, S. Ali, T. Yue et al., “Enhancing test case prioritization in an industrial setting with resource awareness and multi-objective search,” in *Proceedings of 38th International Conference on Software Engineering Companion*, pp. 182–191, ACM, Austin, Texas, USA, May 2016.
- [38] K. Zhai, B. Jiang, and W. K. Chan, “Prioritizing test cases for regression testing of location-based services: metrics, techniques, and case study,” *IEEE Transactions on Services Computing*, vol. 7, no. 1, pp. 54–67, 2014.
- [39] K. Zhai, B. Jiang, W. K. Chan et al., “Taking advantage of service selection: a study on the testing of location-based web services through test case prioritization,” in *Proceedings of 2010 IEEE International Conference on Web Services (ICWS)*, pp. 211–218, IEEE, Miami, FL, USA, July 2010.
- [40] J. Gao, X. Bai, W.-T. Tsai, and T. Uehara, “Mobile application testing: a tutorial,” *Computer*, vol. 47, no. 2, pp. 46–55, 2014.
- [41] H. Muccini, A. Di Francesco, and P. Esposito, “Software testing of mobile applications: challenges and future research directions,” in *Proceedings of 7th International Workshop on Automation of Software Test*, pp. 29–35, IEEE Press, Zurich, Switzerland, June 2012.
- [42] V. L. L. Dantas, F. G. Marinho, A. L. da Costa et al., “Testing requirements for mobile applications,” in *Proceedings of 24th International Symposium on Computer and Information Sciences ISCIS 2009*, pp. 555–560, IEEE, Guzelyurt, Northern Cyprus, Turkey, September 2009.
- [43] R. Morla and N. Davies, “Evaluating a location-based application: a hybrid test and simulation environment,” *IEEE Pervasive computing*, vol. 3, no. 3, pp. 48–56, 2004.
- [44] D. Zhang and B. Adipat, “Challenges, methodologies, and issues in the usability testing of mobile applications,” *International Journal of Human-Computer Interaction*, vol. 18, no. 3, pp. 293–308, 2005.
- [45] S. Vilkomir, “Multi-device coverage testing of mobile applications,” *Software Quality Journal*, vol. 26, no. 2, pp. 197–215, 2017.
- [46] D. Amalfitano, A. R. Fasolino, P. Tramontana et al., “Considering context events in event-based testing of mobile applications,” in *Proceedings of 2013 IEEE Sixth International Conference on Software Testing, Verification and Validation Workshops (ICSTW)*, pp. 126–133, IEEE, Luxembourg, March 2013.
- [47] H. Kim, B. Choi, and W. E. Wong, “Performance testing of mobile applications at the unit test level,” in *Proceedings of Third IEEE International Conference on Secure Software Integration and Reliability Improvement SSIRI 2009*, pp. 171–180, IEEE, Shanghai, China, July 2009.

Research Article

The Identification of Marketing Performance Using Text Mining of Airline Review Data

Jae-Won Hong¹ and Seung-Bae Park² 

¹Department of Global Trade, Gyeongnam National University of Science and Technology, 33 Dongjin-ro, Jinju-si, Gyeongnam 52725, Republic of Korea

²Department of Industrial Management, Seoil University, 28 Youngmasan-ro 90-gil, Jungnang-gu, Seoul, Republic of Korea

Correspondence should be addressed to Seung-Bae Park; sbpark@seoil.ac.kr

Received 1 October 2018; Accepted 13 November 2018; Published 2 January 2019

Guest Editor: Jaegeol Yim

Copyright © 2019 Jae-Won Hong and Seung-Bae Park. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

We aim firstly to extract major keywords using text mining method, secondly to identify prominent keyword from the keywords extracted from text mining analysis, and then to confirm differences in influences of the keywords which affect corporate performance. Results were as following. First, keywords have been found to show distinctive features. Since the keywords posted from the clients showed certain tendency, airlines accordingly need service management by identifying the service property through keyword analysis. Second, prominent keywords have been found out of the keyword extracted from text mining. Some of the keywords have significantly correlated with marketing performance, but others not. This implies that the company could uncover consumers' needs through the prominent keywords and managing the properties related to the prominent keywords would help with improving corporate performance. Third, "recommend" should be treated distinctively with "satisfaction" in terms of service management through the keywords. Results suggest strategic implications to the practical business environment by analyzing keywords around the industry using text mining. We believe this work, which aims to establish common ground for understanding these analyses across multiple disciplinary perspectives, will encourage further research and development of service industry.

1. Introduction

With rapid development of Internet environment and mobile devices, there have been inevitable changes in way of consumers' communication. Consumers form their own community connected with SNS, which refers to social networking service such as Facebook, Twitter, blog, etc. And those network services are treated as reliable message sources over the mass media. As social networks have become increasingly influential, and consumers exploit much WOM information through the social network, the information of online communication from social network plays a critical role in purchase decision making [1]. Moreover, when it is said that traditional WOM communication is implicit or personal, online WOM is more open and collective. Due to the characteristics of the social network, reach of the information is more extensive and the speed of the information

is also even faster. Importance of active responses to consumers' online WOM, therefore, increases in government and social organization as well as corporations.

In this context, online review that consumers produce on SNS can be influential means to grasp a tendency of consumers WOM and recommendation. Consumers, for instance, even read product review to buy a book, which belongs to the section of experiential products. Preceding researches on online review data mainly hire survey method or text mining method [2]. Survey method is used since it gets direct and clear opinions from consumers, but is limited to collect a variety of consumers' opinions and latent opinions. On the other hand, using the text mining method, the limitation of survey method can be overcome by collecting practical opinions from consumers, and through refining data sourced from online, it can also help with deducing consumers' underlying latent opinions as well.

Text mining, on the other hand, is a method of extracting unknown and valuable information from randomly organized text data [3]. Thus, it is described as an automated tool that extracts undisclosed information from text data which are of unstructured format such as mail, reviews, web documents, video clips, or images [4]. Recently, major statistical packages and data mining programs include text mining function to facilitate or simplify the analysis process of those kinds of unstructured data. They provide with function of preprocess to conduct text mining and function of summarizing and categorizing to identify the pattern of the data, but also provides with a variety analysis function of associate analysis, cluster analysis, etc.

It is important to reduce data in handling big data. Large data are not necessarily required to discover important implications. The important thing is the appropriate data. The academic significance of our research is to derive ways to reduce to appropriate data. Current text mining method tends to complement these limitations, accordingly useful to deduce managerial implication on practical decision-making [5]. In this context, by firstly deducing major variables based on the keywords extracted by text mining to reduce large data that is not necessary to discover important implication, and secondly combining them with items in questionnaire, this research will help with suggesting practical implication in the context of aviation industry.

2. Related Works

2.1. Text Mining. Text mining refers to automated methods that extract undiscovered and valuable information from unstructured text by categorizing or structuralizing the text [6]. By extracting information from big data in a variety of field, connectivity within information will be uncovered. This overcomes limitations occurred by simple data analysis and enables to identify underlying meanings from massive text data. In this point, the importance of these methods increases in terms of that the method can be utilized for suggesting practical future strategies.

Through the method of text mining, researchers would take advantage of not only extracting concepts of the text, but also identifying relationship with other concepts and visualizing the relationship among the concepts. Current content analysis relies on items that researchers have arbitrarily selected; accordingly, extensive analysis on gathered data, therefore, is limited, and also external validity is not secured since it relies on coders of the data. Text mining, however, has been considered to surpass the limitation of traditional content analysis and used in a variety of fields using big data analysis, social network analysis, consumer product review analysis, and other useful methods. In other words, text mining extracts the appropriate variables to limit the breakdown of content analysis. Netzer et al. [7] examined the relationship between automobile brands through text mining and analyzed the market structure using the multidimensional scaling method. In addition, Mostafa [8] has classified the lexicon through 3D Map in the research that confirmed the brand sentiments of famous brands such as Nokia, IBM, and DHL through social network text mining.

Text consists of words, and analyzing text can be described as analyzing relationship among the words. In terms of it, text network analysis is also called semantic network analysis. That is to say, depending on the research, it can be called networks of words, network text analysis, semantic nets, networks of concepts, networks of centering words, text network analysis, or semantic networks [9].

Text network analysis, as mentioned above, complements the limitation of traditional content analysis, and extracts underlying meaning that the text delivers. Moreover, the pattern of text can be structurally analyzed to identify the relationship among the meanings and the relationship accordingly can be visualized through the analysis [10]. Text mining process has two phases including the data process phase and data analysis phase. The data process phase is relevant to data gathering and preprocess, while data analysis phase is relevant to text analysis that extracts significant information from the text, and visualizing information and extracting knowledge from the former analysis [6].

Through this process, large-volume data can be made more data-suitable for analysis, enabling continuous research to compare experiments.

2.2. Consumers' Evaluation Criteria on Airlines. Aviation industry can be described as a field where degree of interaction, customization, and labour intensiveness is relatively lower than that of other field [11]. Service, compared to products, has shown distinctive features of intangibility, heterogeneity, inseparability, and perishability [12]. In other words, consumers cannot see or touch purchased "services". In addition, production and consumption take place at the same time, and it cannot be stored, but even extinguishes once it is unused after production. In that properties of service, therefore, perceived service evaluation or recommendation can be crucial indicators in the field of aviation industry. Jia [13] has conducted text mining through a Chinese crowd-sourced online review community, and 49,080 pairs of restaurant ratings and reviews were examined, with high-frequency words, major topics, and sub-topics identified. After text mining, multilinear regression was employed to screen out the most impactful factors that influence taste, environment, and service ratings. Managerially, the idea of triggering the synergistic benefit from customer ratings and reviews is referential for market practitioners both within and beyond the catering industry.

In case of aviation industry, these elements including flight schedule, fare, services, punctuality, comfort of seats, safety, and frequent flyer program can be determinants of service satisfaction on airlines (IATA: International Air Transport Association). And in research of Hong and Park [14], factors that determine consumers' satisfaction on airline services include punctuality, safety, courteous agents, clean equipment, space, desirable schedule, profitability of the airline, reliability on the provided service, and financial costs. In addition, flight fare, boarding process, space and comfort of seats, in-flight meal, baggage delivery, and ticketing process also can affect consumers'

satisfaction. Based on researches, passengers who experience airline services tend to compare and evaluate its overall quality on the basis of tangible and intangible services provided by airlines. Will these factors be derived from the customers' texts in Internet? Internet comments based on anonymity reveal the desire of consumers and reveal more clear facts.

Aviation industry itself can be defined as a service industry where degree of interaction, customization, and labour intensiveness is relatively lower than that of other field [15]. In-flight service which is the interest of this study, however, can be described with high interaction, high customization, and high labour intensiveness. Thus, assessment and recommendation on services which consumers directly have experienced will be important indicator for the industry. In general, in-flight service consists of tangible service and human service that help with passengers' travel experiences; comfort of seats, in-flight employee services, in-flight food and beverage, entertaining services, etc. These elements, which are important in the aviation industry, are identified through text mining and classified through cluster analysis. These studies provide a possibility to become more realistic studies by combining text mining and statistical analysis.

3. Methodologies

3.1. Research Process. In this study, we conducted the following process to find out whether the core keywords can be identified through text mining and whether the core keywords are important for the performance of a company (see Figure 1). First, we gathered online reviews of customers and extracted core keywords from it by text mining method. And, we conducted text clustering analysis to explore the meaning of extracted core keywords. After that, we conducted empirical test for demonstrating the impact of core keywords through analysis of the relevance of the company's marketing performance, such as satisfaction, recommendation.

3.2. Data. This study used the online review data of airline customers, which was provided by global air service evaluation agency Skytrax in United Kingdom. This study set two large air carriers in Korea and Japan. The main contents of the data consist of two types of data. One type is the text data containing the customer's experience after using the air service. The other type is the survey data which include the evaluation of services, satisfaction, and recommendation after using the air service. The questionnaire survey was conducted online for customers who have recently used airline. The respondents are confirmed by presenting their plane ticket. The item of customer satisfaction was measured as 10 point-Likert scale, and recommendation intention was measured as binomial scale.

Data period is 3 years from January 2013 to December 2015. In the data period, 197 reviews were for the Korean carrier and 214 reviews for Japanese carrier. So, this study collected the review data on 411 people in total during data period and used them for analysis. Among the collected data,

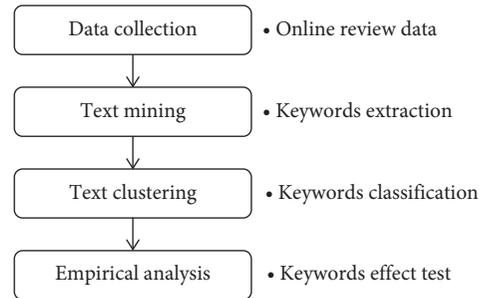


FIGURE 1: Research process.

text data were used for core keyword extraction and, questionnaire data were used for empirical analysis to verify the influence of core keywords.

The number of countries, to which the respondents were affiliated, was 32 in total, and the United States was found to be the most (33.8%), followed by Australia (12.7%), United Kingdom (10.7%), Canada (6.1%), and Singapore (5.1%). Table 1 shows the top 10 countries on the basis of the number of respondents.

3.3. Analytic Process. In order to analyze the text data, this study performed two processes. Text mining process was for extracting words, and text clustering process was for analyzing meaning based on extracted words.

The text mining process was a process of extracting words from a text document and creating word data. And this process included a step of extracting meaningful words based on text data. The detailed procedure is as follows [6]:

3.3.1. Text Mining Process

Step 1 (treat text). This is the process of creating word data based on words contained in a text documents.

$$\vec{t}_d = (\text{tf}(d, t_1)), (\text{tf}(d, t_2)), \dots, (\text{tf}(d, t_m)), \quad (1)$$

where \vec{t}_d is the words vector; D is the documents set $\{d_1, d_2, \dots, d_n\}$; and T is the words set $\{t_1, t_2, \dots, t_m\}$.

Step 2 (extract words). This is the process of extracting meaningful words by assigning weights through the appearance frequency of words between documents set.

$$W(t) = 1 + \frac{1}{\log_2 |D|} \sum_{d \in D} P(d, t) \log_2 P(d, t),$$

$$\text{with } P(d, t) = \frac{\text{tf}(d, t)}{\sum_{l=1}^n \text{tf}(d_l, t)}, \quad (2)$$

where $\text{tf}(d, t)$ is the frequency of word t in a document d .

3.3.2. Text Clustering Process. On the other hand, text clustering process is a process of clustering through the distance between words. The cluster is set up and adjusted

TABLE 1: Characteristics of country.

No.	Country	Korea		Japan		Total	
		N	%	N	%	N	%
1	United States	58	29.4	81	37.9	139	33.8
2	Australia	50	25.4	2	0.9	52	12.7
3	United Kingdom	25	12.7	19	8.9	44	10.7
4	Canada	4	2.0	21	9.8	25	6.1
5	Singapore	3	1.5	18	8.4	21	5.1
6	Japan	3	1.5	15	7.0	18	4.4
7	South Korea	13	6.6	1	0.5	14	3.4
8	Brazil	1	0.5	12	5.6	13	3.2
9	Germany	7	3.6	4	1.9	11	2.7
10	New Zealand	8	4.1	1	0.5	9	2.2

after the cluster was set. In this method, (1) generate k initial seed randomly within data domain, (2) create k clusters by associating every observation with the nearest seed position, and (3) adjust the center position of each of the k clusters using the average of the observations belonging to the cluster. By repeating this process until all observations are associated, it can make clusters which have similar observations corresponding to the number of k .

Step 1 (cluster setting). Set up a cluster with random seeds in the data, and the cluster is grouped by the Euclidean distance between each data and seed based on near distance.

$$\min S = \sum_{i=1}^k \sum_{x \in s_i} \|x - u_i\|^2 \quad (3)$$

where x is the words set $\{x_1, x_2, \dots, x_n\}$, s is the clusters set $\{s_1, s_2, \dots, s_k\}$, and u_i is the average of cluster s_i .

Step 2 (cluster adjust). Reset the value of cluster and adjust the position of cluster by using the average of the data in the cluster.

4. Results

4.1. Keywords Extraction by Text Mining. In this study, several preprocessing steps were performed to improve the accuracy of text mining.

First, all uppercase letters in text data are converted to lowercase letters to unify words. Also, we removed unnecessary special characters such as “@,” “\,” and so on. And the definite article (“the”), the indefinite article (“a,” “an”), prepositions (“of,” “in,” “for,” “through,” etc), and pronouns (“it,” “their,” “his,” etc) were also removed. Through this process, a total of 3,774 keywords were extracted. The keywords were rearranged based on the sparsity which represents the ratio of the whitespace in the matrix, because analyzing all keywords makes it difficult to explain meaningful results.

There were 11 keywords based on the sparsity 0.8, 19 keywords based on the sparsity 0.85, 45 keywords based on the sparsity 0.9, and 132 keywords based on the sparsity 0.95. In this study, a total of 45 keywords were extracted by applying the sparsity 0.9, and it is used for analysis. Too few

keywords are insufficient to understand customer behaviour, but too many keywords can distort results with less important words. In this analysis, 45 keywords were selected considering the usefulness of the interpretation while reflecting the customer behaviour.

The keywords with high frequency were “good,” “seats,” “cabin,” “class,” “excellent,” “comfortable,” “staff,” and so on. Figure 2 shows the frequency distribution of these keywords, and Figure 3 shows the frequency of keyword appearance by the word cloud method.

4.2. Keywords Classification by Text Clustering. In this study, we performed keyword classification analysis for semantic analysis of extracted 45 keywords. The results are as follows. Keyword classification was based on hierarchical clustering, and the distance between the keywords was measured by the Euclidean method.

As a result of clustering analysis, 45 keywords were classified into two clusters. Cluster 1 consisted of service content such as “seats,” “cabin,” “class,” and “staff” who provide the service. And it has service evaluation-related keywords such as “good,” “excellent,” “comfortable,” etc.

On the other hand, Cluster 2 consisted of more detail service content such as “meal,” “movie,” “drinks,” “check,” “served,” and so on. Also, it has service evaluation-related keywords such as “great,” “nice,” “well,” “best,” “clean,” etc (see Figure 4).

Figure 5 shows the result of keyword visualization with two clusters based on the k -means method for better insight. The keywords in Cluster 1 are formed to be spaced apart from each other, and the keywords in Cluster 2 are formed to be relatively wide spacing. And the two components explained 61.83% of the point variability.

4.3. Effect of Core Keywords. In this study, we performed combining analysis with core keywords data extracted by text mining and questionnaire respondent data to overcome the disadvantages of existing text mining research and to provide practical implications to the aviation industry. In other words, since core keywords appear as a result of text mining and cluster analysis can be considered to represent online reviews of airline customers, we attempted to understand the meaning of representative keywords by examining the influence of each evaluation concept on customer satisfaction and customer recommendation based on the evaluation of these customers.

In the analysis, we analyzed the effect of the four keywords, such as “seats,” “staff,” “cabin,” and “class,” which are the main keywords corresponding to the contents of the Cluster 1 on customer satisfaction and recommendation. The reason for this is that adjective and adverbial keywords among the keywords of Cluster 1 are mainly keywords indicating the result of the service. And the keywords appearing in Cluster 2 are keywords that deal with details of the service contents of Cluster 1. So, we define that these keywords are conceptually included in the above four top keywords.

In the analytical model, the questionnaire items were used as variables. Independent variables were the

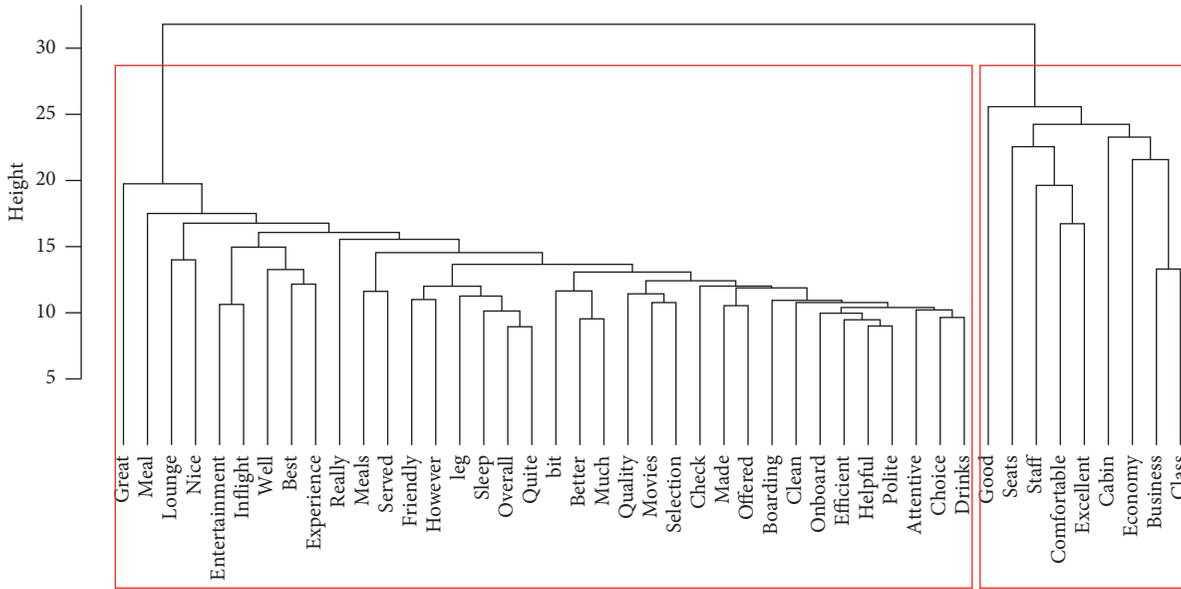


FIGURE 4: Cluster dendrogram.

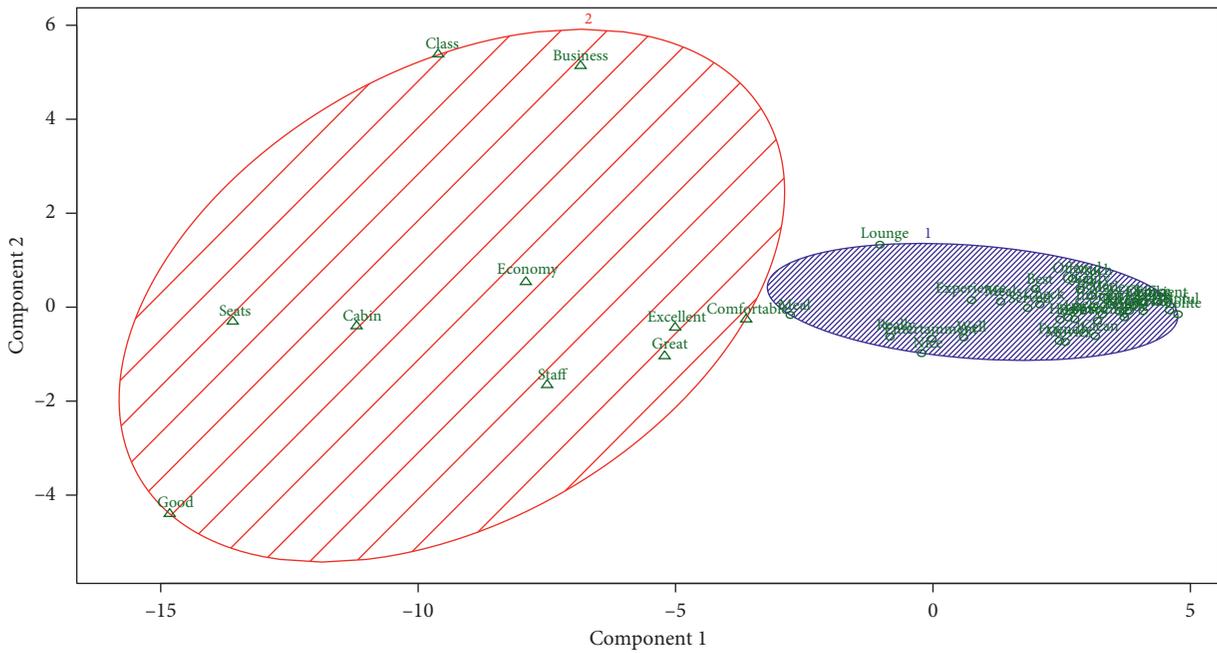


FIGURE 5: Cluster plot. These two components explain 61.83% of the point variability.

TABLE 2: Effects of core keyword on satisfaction.

Variables	<i>B</i>	β	<i>t</i> -value	<i>p</i> -value
Seats	0.987	0.513	18.6***	0.001
Staff	0.874	0.334	11.2***	0.001
Cabin	0.555	0.235	7.4***	0.001
Class	0.056	0.011	0.5	0.645

*** $\rho < 0.001$; *F* value: 342.4; R^2 : 0.771.

TABLE 3: Effects of core keyword on recommendation.

Variables	<i>B</i>	β	<i>t</i> -value	<i>p</i> -value
Seats	0.156	0.536	13.5***	0.001
Staff	0.110	0.275	6.4***	0.001
Cabin	0.020	0.055	1.2	0.224
Class	0.039	0.051	1.5	0.139

*** $\rho < 0.001$; *F* value: 111.0; R^2 : 0.518.

study used text mining to develop a unified understanding of keywords in the aviation industry in a data-driven way. Based on the airline review data, we proposed a two-step process of extracting key keywords by text mining and grouping them into cluster analysis. Specifically, we used a combination of metrics and clustering algorithms to preprocess and analyze text data related to keywords extraction method, including text from the scientific literature and news articles. This study seeks firstly for identifying prominent keywords at consumers' side using text mining method on consumers' online review data and then for confirming influences that the keywords affect corporate marketing performance. This study is not only the research that searches for key keywords [16], but also the research that identifies marketing performance in the aviation industry through text mining in service category level. Conclusion and implication are as follows.

First, keywords are shown to have distinctive cluster characteristics. As a result of identifying characteristics of major keywords through clustering, the keywords have been classified into three sections, including service that airlines provides with, details of each service, and assessment on the services. In other words, since the service provided by airline and the details of each service are differentiated at different levels, it indicates that service management should be centered on core service elements that can be clearly recognized in order to improve service evaluation.

Second, the key keywords extracted from text mining were found to have a different relationship with corporate performance. In other words, the service of the seats and the staff was more important to the company performance than the cabin or class. These results show that comfort is the key customer's needs in long-distance air travel and that the service focused on the comfort of the seat or the comfort of the seat is an important factor in corporate performance.

Third, recommendation and satisfaction should be managed distinctively in the service management of the aviation industry. The results show that keywords that affect consumers' satisfaction and keywords that affect consumers' recommendation are found to differ and the degree of impacts is also shown to be different as well. In other words, seats, staff, and cabin are important factors to improve the satisfaction and recommendation of consumers. Seats are the most important factors in consumer satisfaction and recommendation, but relatively more staff and cabin are important for consumer satisfaction and seats are more important for consumer recommendation. These results show that the human service of the crew, which can be relatively subjective, is more influential in satisfaction. On the other hand, seats are more important for relatively objective recommendations.

This study presents academic implication that the study has extended its application area of text mining. It has currently focused on exploratory study, while this study has extended it to study field of cause and effect. Moreover, the research also presents practical implication for corporations to efficiently manage keywords. In spite of the implications and advantage of text mining [17], this study has limitations. First, the research range is limited to two airlines in Korea

and Japan. Secondly, since the review data are produced spontaneously by clients, collected data can be biased and limited to active clients who are willing to express his/her experience. Therefore, in future study, selecting more representative airlines is needed and also minimizing convenience of the respondents is needed to generalize the results of the study in the future.

Data Availability

The text mining data used to support the findings of this study are available from the authors upon request. The data used in this study are Airline Reviews and Rating data from Skytrax (<http://www.airlinequality.com>).

Conflicts of Interest

The authors declare that they have no conflicts of interest.

References

- [1] P. Nayak and A. K. Chatterjee, "Effects of aluminium exposure on brain glutamate and GABA systems: an experimental study in rats," *Food and Chemical Toxicology*, vol. 39, no. 12, pp. 1285–1289, 2001.
- [2] J. A. Chevalier and D. Mayzlin, "The effect of word of mouth on sales: online book reviews," *Journal of Marketing Research*, vol. 43, no. 3, pp. 345–354, 2018.
- [3] W. Fan, L. Wallace, S. Rich, and Z. Zhang, "Tapping the power of text mining," *Communications of the ACM*, vol. 49, no. 9, pp. 76–82, 2006.
- [4] D. Meyer, K. Hornik, and I. Feinerer, "Text mining infrastructure in R," *Journal of Statistical Software*, vol. 25, no. 5, pp. 1–54, 2008.
- [5] R. V. Kozinets, K. de Valck, A. C. Wojnicki, and S. J. S. Wilner, "Networked narratives: understanding word-of-mouth marketing in online communities," *Journal of Marketing*, vol. 74, no. 2, pp. 71–89, 2010.
- [6] A. Hotho, A. Nürnberger, and G. Paaß, "A brief survey of text mining," *Ldv Forum*, vol. 20, no. 1, pp. 19–62, 2005.
- [7] O. Netzer, R. Feldman, J. Goldenberg, and M. Fresko, "Mine your own business: market-structure surveillance through text mining," *Marketing Science*, vol. 31, no. 3, pp. 521–543, 2012.
- [8] M. M. Mostafa, "More than words: social networks' text mining for consumer brand sentiments," *Expert Systems with Applications*, vol. 40, no. 10, pp. 4241–4251, 2013.
- [9] D. Paranyushkin, "Text network analysis (2010)," in *Conférence du Performing Arts Forum*, 2011, <http://noduslabs.com/research/pathways-meaning-circulation/>.
- [10] R. Feldman and J. Sanger, *The Text Mining Handbook: Advanced Approaches in Analysing Un-Structured Data*, Cambridge University Press, Cambridge, UK, 2007.
- [11] R. W. Schmenner, "How can service businesses survive and prosper?," *Sloan Management Review 1986-1998*, vol. 27, no. 3, p. 21, 1986.
- [12] J. Fitzsimmons and M. Fitzsimmons, *Service management: Operations, Strategy, Information Technology*, McGraw-Hill Higher Education, New York City, PA, USA, 2013.
- [13] S. Jia, "Behind the ratings: text mining of restaurant customers' online reviews," *International Journal of Market Research*, vol. 60, no. 6, pp. 561–572, 2018.
- [14] J.-W. Hong and S.-B. Park, "Study on the extraction of core keywords and its effects through text mining," *International*

- Journal of Web Science and Engineering for Smart Devices*, vol. 3, no. 2, pp. 7–12, 2016.
- [15] D. Jung and S. B. Kim, “The effects of information quality of mobile on consumer behaviour intension; for airline,” *Asia-pacific Journal of Multimedia Services Convergent with Art, Humanities, and Sociology*, vol. 6, no. 9, pp. 19–26, 2016.
- [16] S.-B. Park and J.-W. Hong, “Relationship of core keywords and marketing performance obtained by using text mining: a comparative study,” *New Physics: Sae Mulli*, vol. 67, no. 5, pp. 562–568, 2017.
- [17] X. Wang and S. Liu, “Analysis and research of enterprise technology competent advantage on text mining and correspondence analysis,” *International Journal of Database Theory and Application*, vol. 6, no. 5, pp. 133–140, 2013.

Research Article

From Reputation Perspective: A Hybrid Matrix Factorization for QoS Prediction in Location-Aware Mobile Service Recommendation System

Shun Li ¹, Junhao Wen ¹ and Xibin Wang ²

¹School of Big Data and Software Engineering, Chongqing University, Chongqing 400044, China

²School of Data Science, Guizhou Institute of Technology, Guiyang 550003, China

Correspondence should be addressed to Junhao Wen; jhwen@cqu.edu.cn

Received 27 September 2018; Accepted 19 November 2018; Published 2 January 2019

Guest Editor: Subramaniam Ganesan

Copyright © 2019 Shun Li et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

With the great development of mobile services, the Quality of Services (QoS) becomes an essential factor to meet end users' personalized requirement on the nonfunctional performance of mobile services. However, most of the QoS values in real cases are unattainable because a service user would only invoke some specific mobile services. Therefore, how to predict the missing QoS values and recommend high-quality services to end users becomes a significant challenge in mobile service recommendation research. Previous QoS prediction researches demonstrate that the nonfunctional performance of mobile services is closely related to users' location information. However, most location-aware QoS prediction methods ignore the premise that the obtainable QoS values observed by different users in same location region would probably be untrustworthy, which will lead to inaccurate and unreliable prediction results. To make credible location-aware QoS prediction, we propose a hybrid matrix factorization method integrated location and reputation information (LRMF) to predict the unattainable QoS values. Our approach firstly cluster users into different locational region based on their geographical distribution, and then we compute users' reputation to identify untrustworthy users in every locational region. Finally, the unknown QoS values can be predicted by integrating locational cluster information and users' reputation into a hybrid matrix factorization model. Comprehensive experiments are conducted on a public QoS dataset which contains sufficient real-world service invocation records. The evaluation results indicate that our LRMF method can effectively reduce the impact of unreliable users on QoS prediction and make credible mobile service recommendation.

1. Introduction

Based on the flexibility and expansibility of mobile application development technique, tens of thousands of hybrid mobile services with similar function have been developed and provided in mobile application store. However, this phenomenon automatically leads to information overload problem in mobile service retrieval system. To tackle this challenge, Quality of Service (QoS) is used in service-oriented system to analyse the nonfunctional performance of mobile services [1–4]. QoS has been widely used in service selection, composition, and recommendation research [5–9]. In real-world service invocation scenario, users would only search and select some specific mobile services under

the unpredictable Internet environment. For lots of unknown mobile services, it is impractical to make users invoke each of them and evaluate their nonfunctional performance. Therefore, how to make accurate QoS prediction for unknown mobile services is a critical step to make high-quality service recommendation in mobile service computing paradigm.

Collaborative filtering (CF) is widely utilized in most e-commerce recommender systems to predict miss rating values. Traditional CF models generally fall into two categories: memory-based and model-based. The memory-based CF of QoS prediction process would generally find a subset of similar users for the target user and recommend high-quality mobile services shared by these similar users to the

target user [10]. The model-based methods will train a model by learning users' historical QoS performance and then predict the QoS values for unknown mobile services [11, 12]. Although CF model proved to be effective in QoS prediction on different mobile services, the prediction accuracy is still unsatisfactory because of cold-start problem.

To reduce the impact of cold-start problem, more contextual information is introduced into QoS prediction model for the QoS values are greatly affected by some context factors (e.g., location distribution, invocation time, and so on) in Internet. Based on this realization, the location-aware CF model is proposed to predict unknown QoS values in service recommendation [13]. As we all know, different users in one location region generally share the same set of IT infrastructure and they would suffer from similar Internet usage experience when they invoke mobile services, as it is reported in the work of [14], in which the QoS performance is strongly correlated with the location information of users. In Figure 1, we give an example of service invocation with location information. As mentioned above, the user 1 would have similar QoS records (such as response time) with other users in the US; meanwhile, user 2 may share similar QoS records with other users in India when they invoked services such as YouTube, Twitter, etc.

Previous location-aware CF approaches usually compute the similarity of all different users in one specific location region to find most Top-K similar neighbours for the target user [12, 15]. However, some unreliable users who would submit untrustworthy QoS values will be indiscreetly included in the neighbourhood set. Unreliable users would randomly provide some QoS values or better ones to improve the visibility of their own services and worse values for others' applications [16]. Those untrustworthy QoS values would have a marked negative effect on the prediction accuracy. Therefore, it is essential to introduce credibility of available QoS values in prediction process to enhance the prediction accuracy and persuasiveness of service recommendation mechanism.

Based on above realizations, a hybrid matrix factorization algorithm is proposed by integrating users' reputation and locational information to predict the unattainable QoS values in this paper. Complementary to previous service recommendation method which only adopts available QoS values, our study tends to make credible and accurate QoS prediction for mobile service recommendation by considering the reputation of different users' QoS usage experience. We then exploit personal geographical distance and QoS values to find the locational similarity neighbourhoods and discover the latent connectivity between the target user and his/her neighbours. Meanwhile, we use users' reputation to control the weight of users' latent feature learning. Finally, these constraints are integrated into matrix factorization model to make credibly personalized QoS prediction. The following contributions are achieved in this paper:

- (1) We firstly cluster users into different locational regions based on their geographical distribution and design an iterative method to compute users' reputation score by their provided QoS usage data. Then

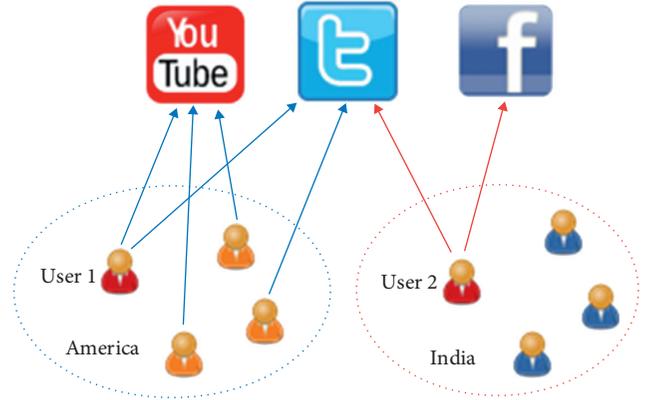


FIGURE 1: An example of mobile service invocation scenario.

a subset of trustworthy users in each locational region can be identified by the rank of reputation score.

- (2) In the next step, a trustworthy neighbourhood can be identified by incorporating both users' location distribution and reputation score. By integrating the latent feature of both available QoS values and those shared by the neighbourhood, a hybrid matrix factorization model is proposed to make high-quality QoS prediction.
- (3) Results of experiments conducted on a public QoS dataset show that by considering data credibility, our method can achieve higher prediction accuracy than other previous studies which involve the untrustworthy impact of available QoS values.

The remainder of this paper is organized as follows: Section 2 introduces related studies; Section 3 shows the basic principles of our method; Section 4 elaborates how to design our proposed method; Section 5 discusses experiments and results analysis; finally, Section 6 concludes the paper and draws future studies.

2. Related Work

2.1. QoS-Based Service Computing. QoS plays a control role in service-oriented architectures, especially in the service discovery and recommendation research. Al-Masri and Mahmoud [17] firstly calculated users' preferences on their historical QoS data, and then proposed a service discovery approach by ranking users' QoS parameters. Kritikos and Plexousakis [18] extracted the contextual information of QoS from the description file of service and designed a service discovery method. Rosario et al. [19] proposed soft probabilistic contracts on QoS parameters to composite web services and validated their method on TOrQuE tool to show its outperformance than other previous studies. Hadad et al. [20] proposed a web service composition framework by exploiting both transactional properties and QoS values. This framework composite plenty of existing web services into a workflow which can satisfy users' preferences on nonfunctional requirements. However, these methods only conducted experiments on synthetic datasets and lack authenticity in real-world service invocation.

2.2. Collaborative Filtering. Collaborative filtering is a common algorithm adopted by many recommendation systems, such as the famous commercial system Amazon [21]. Memory-based collaborative filtering approaches make prediction and recommendation by calculating the similarity of users or items [21, 22]. These methods utilize the whole entire user-item matrix as the input data, which will take a lot of time and memory spaces in the online recommendation system. Model-based collaborative filtering methods will generally train a predefined model on available data and predict the missing values in the test dataset, and then select the appropriate items as candidate list to the target user [23, 24]. Model-based approaches can learn the model quickly with little need for runtime and memory space, which will be often adopted in online recommendation systems.

Matrix factorization model is now widely adopted in many online recommendation systems for its effectiveness and efficiency. Mnih and Salakhutdinov [11] introduced the mathematical theory of matrix factorization in probabilistic analysis and validated the performance of this method on a famous film recommendation system. Zhang et al. [25] designed a personalized recommendation approach by integrating original matrix factorization with a constraint item extracted from their personal information. It identifies users into different clusters by the statistics of user behaviours on different tags and considers this constraint as a regularization term in matrix factorization model to enhance its prediction accuracy. Ma et al. [26] improved the matrix factorization approach with users' social information to enhance the prediction accuracy for social recommendation. This approach uses users' social relationship as an additional constraint which can reflect users' latent judgment of interest on items in the user-item matrix factorization. Recently, the matrix factorization methods have been widely introduced into service recommendation research [10, 15]. Although matrix factorization methods make some improvements in prediction accuracy, none of them realize that the QoS credibility deserve serious consideration.

2.3. Location-Based QoS Prediction. Location information has been widely used for service recommendation in recent years. Ali and Solis [27] presented a novel distributed service architecture that can adapt to the changes of Internet resources and location topology. Wei et al. [15] firstly calculated users' similarity with their real-world distances and then clustered users into geographical sets as a constraint item of matrix factorization to generate the location-aware QoS prediction approach. Tang et al. [28] solved the sparsity issue in QoS aware service recommendation by integrating collaborative filtering with users' geography data. Lee et al. [29] adopt the preference propagation through users in same location region to improve prediction accuracy. This work clusters users and service into different groups by the locational information and then use preference propagation to compute the similarity between different users and services, respectively. Finally, a matrix factorization model is introduced to predict missing QoS values by integrating these

constrains. Gonsalves and Patil [30] exploited users' location information and QoS values to cluster users and web services and then proposed a CF algorithms to make personalized web service recommendation. It firstly uses Pearson correlation coefficient (PCC) to identify different users and service regions and then exploit K-nearest neighbour (KNN) and support vector machine (SVM) in CF algorithm framework to predict missing QoS values. However, above studies do not consider the users' reputation, and neglect the fact that available QoS values may be untrustworthy even though these values are provided by users in same location region.

2.4. Reputation. Based on the achievement of reputation in applications (e.g., YouTube and Twitter) to avoid possible deception risk, some academics introduce the reputation into QoS prediction to enhance the reliability of service-oriented computing. The reputation values evaluated from QoS data can measure whether the available QoS values are trustworthy or not. Qiu et al. [31] designed a QoS prediction method by calculating users' reputation to obtain higher accuracy for service recommendation. In their work, reputation of different users will be computed and ranked to find the subset of unreliable users. Followed by this, the memory-based collaborative filtering model combined with reputation for QoS prediction becomes more remarkable. Based on Qiu's work, Xu et al. [32] presented an improved QoS prediction method with the users' reputation (RMF), which introduced users' reputation weight into of a matrix factorization approaches to make QoS prediction for unknown services and then recommend high-quality services to the target user. Mehdi et al. [33] introduced a stochastic approach to evaluate the reputation of services by leveraging the correlation information among different QoS metrics. Comi et al. [34] proposed a hybrid service composition method by exploiting users' reputation of QoS to help users discover and select high-quality services in multicloud environment. However, these studies do not take location information into consideration on the QoS prediction.

3. Principles and Reputation Analysis

In this section, we would introduce the main principles our LRMF method at the beginning and then present the reputation analysis on different users.

3.1. Principles of LRMF. Previous CF-based QoS prediction approaches (e.g., [7, 35]) only utilize the available QoS values in collaborative filtering model to make personalized service recommendation. However, these methods ignore that users' reputation and location will make great impact on the prediction results. Therefore, we design a novel mobile service recommendation system by considering both of users' reputation and location information when predicting missing QoS values. As presented in Figure 2, historical invocations with QoS and location data will be submitted to server database in the service invocation process. Then, the reputation and location information could be calculated by

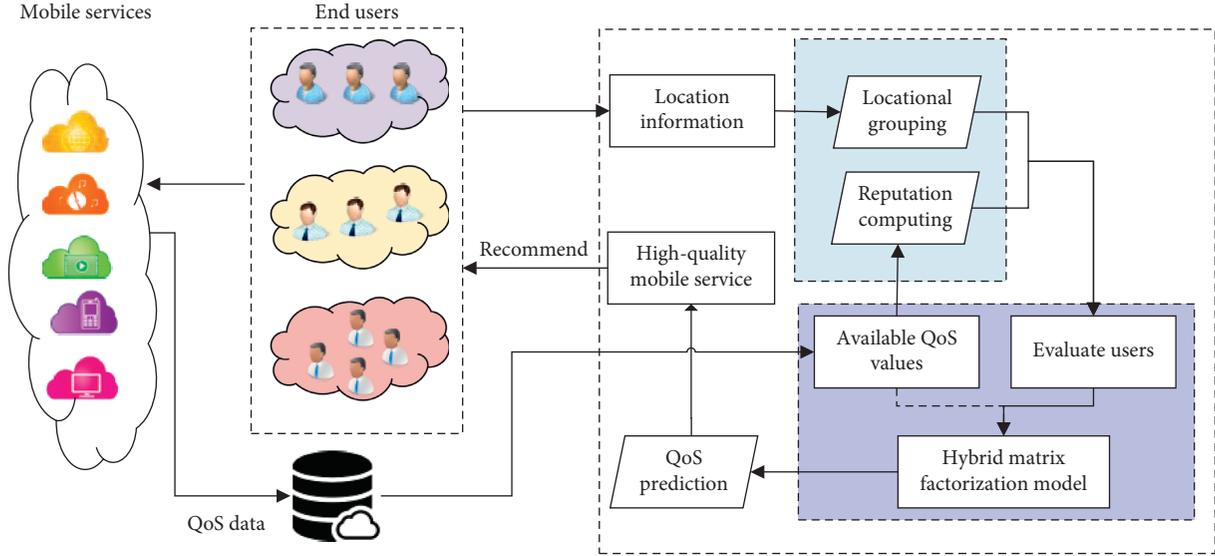


FIGURE 2: Mobile service recommendation framework.

the collected historical QoS dataset, and finally, we can predict missing QoS values with reputation and location information. The main workflow of LRMF is demonstrated as following:

- (1) The multisource invocation records will be collected and submitted into the database when users invoke different mobile services.
- (2) In the data preprocess, the available historical QoS values and users' location info will be extracted as the input data of our method.
- (3) Based on the available QoS values and location data, users' reputation score can be computed by our iterative method and the location region can also be identified by their real-word location distribution. Then, the trustworthiness of users in different location region can be evaluated.
- (4) A hybrid matrix factorization model is proposed by integrating location grouping and users' reputation to predict missing QoS values of unknown mobile services.
- (5) Finally, by combining the predicted results and available QoS data, the high-quality services will be discovered and recommended to the target user.

3.2. User Reputation Analysis. In order to identify untrustworthy users in a given geographical region, we analyse the QoS values from users' historical service invocation records. Figure 3 demonstrates the response time of users in a same region (i.e., the United States in this example) of three randomly selected services from the real-world QoS dataset [36]. As a demonstration, Figure 4 describes the response time of 100 randomly selected services which are invoked by 5 randomly selected users in a same region.

Figure 3 shows that the response time of service invocation varies among users even they are in the same

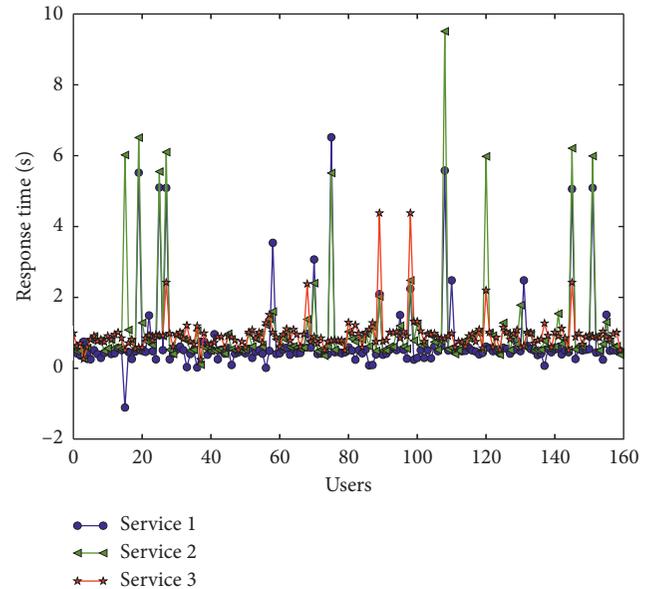


FIGURE 3: The distribution of response time of different users on 3 services in a location region.

location group. Although most of the response time values falls into the normal range, i.e., $[0, 2]$, some users still submit outlier QoS records when they invoke services. It is unlikely that a QoS item would deviate from the normal value too much. For a specific user, if most QoS values significantly deviated from the normal range, he/she is probably an unreliable user. In order to test whether there are some unreliable users or not, we analyse the QoS data submitted by 5 randomly selected users on 100 services.

It is obviously in Figure 4 that the user 4 is an unreliable user because all of his submitted values deviated greatly from the normal range $[0, 2]$. Based on this analysis, we can compute the reputation of users by his past QoS value records and measure whether these users should be regarded

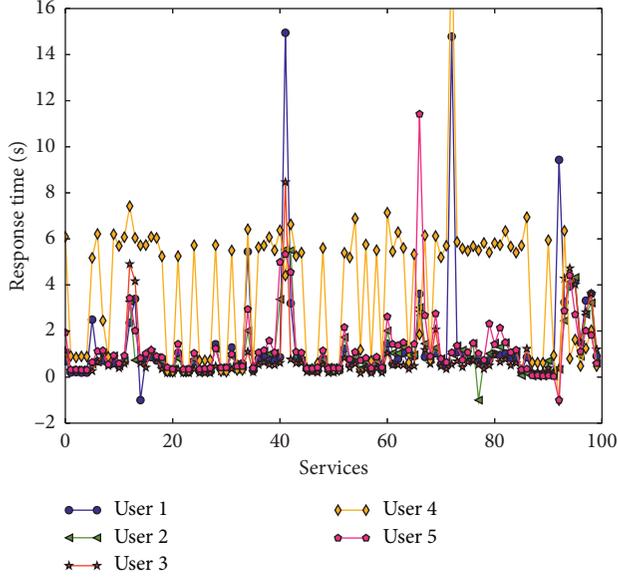


FIGURE 4: The distribution of response time submitted by 5 users on 100 services.

as unreliable users or not. In next steps, the algorithm of users' reputation computation will be introduced in detail.

4. Hybrid Matrix Factorization Based on Location and Reputation

We will firstly present the original matrix factorization approach in this part and then propose our improved matrix factorization model.

4.1. Matrix Factorization Prediction. Matrix factorization utilized low-rank approximations to fit the sparse matrix of user and item. It factorizes original sparse matrix into two low-rank matrices with small number of factors. This factorization is based on the hypothesis that users' latent preference on QoS values would be significantly affected by some latent factors. Then an objective function can be defined as the sum error of original values and the predict values by the conducts of the two low-rank matrices.

We suppose there is a QoS matrix where users in the rows and service in the columns and two low-rank matrices represent user-specific feature and service-specific feature, respectively. The QoS matrix can be regarded as a product of matrix multiplication on low-rank matrix and approximately as follows:

$$R \approx \tilde{R} = U^T S, \quad (1)$$

where $U \in \mathbb{R}^{d \times m}$ and $S \in \mathbb{R}^{d \times n}$. d ($d \ll \min(m, n)$) is the number of latent factors.

Then, the objective function can be defended by minimizing the sum error of available values in original matrix R and the corresponding predicted values in matrix \tilde{R} :

$$\min_{U, S} \psi(U, S) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \|R_{ij} - U_i^T S_j\|_F^2, \quad (2)$$

where R_{ij} represents the available QoS value provided by user i on service j ; U_i denotes the i^{th} row of U ; S_j is j^{th} column of S . However, the available QoS values are limited in real invocations scenario, so an optimal objective function can be defined to solve this issue:

$$\min_{U, S} \psi(U, S) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n I_{ij} (R_{ij} - U_i^T S_j)_F^2, \quad (3)$$

where $I_{ij} = 0$ indicates the QoS value provided by user i on service j is unknown and $I_{ij} = 1$ otherwise. Two regularization terms are introduced in Equation (3) to avoid the overfitting problem as follows:

$$\min_{U, S} \psi(U, S) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n I_{ij} (R_{ij} - U_i^T S_j)_F^2 + \frac{\lambda_u}{2} \|U\|^2 + \frac{\lambda_s}{2} \|S\|^2, \quad (4)$$

where $\|\cdot\|_F$ is the Frobenius norm. Equation (4) is the objective function of matrix factorization approach to minimize the squared error between the predicted values and original values. The gradient decent algorithm is adopted to train our model and get a local minimum value of (4) as follows:

$$U'_i = U_i - \alpha \frac{\partial \psi}{\partial U_i}, \quad (5)$$

$$S'_j = S_j - \alpha \frac{\partial \psi}{\partial S_j}. \quad (6)$$

4.2. Locational Grouping. In this part, we identify different users by their real-world location to acquire a subset of users who have geographical similarity with the target user. Since the location information significantly affects the QoS, it should be considered a significant factor in QoS prediction [13]. We calculate users' physical distance to generate the location region. The Euclidean distance between users i on j can be computed by following definition:

$$\text{dis}(i, j) = \sqrt{(\text{lon}(i) - \text{lon}(j))^2 + (\text{lat}(i) - \text{lat}(j))^2} \times \delta, \quad (7)$$

where $\text{lon}(i)$ and $\text{lat}(i)$ represent the longitude and latitude of user i in the real world, respectively. δ converts the unit of degree into 2D meter with a constant value. In our study, δ takes the value of 111,261.

The geographical region then could be generated by selecting a set of users who are with small distance calculated by Equation (7). On one hand, the size of this set cannot be too small; otherwise, too many similar users would be filtered. On the other hand, it cannot be too large; otherwise the different locations would not be correctly recognized. For a target user i , the region $G(i)$ can be defined as follows:

$$G(i) = \{u_j | \text{dis}(i, j) \leq \varepsilon, i \neq j\}, \quad (8)$$

where ε is a positive variable locational threshold which affects the region size, and $G(i)$ denotes the subset of users who are in the same geographical region with user i . Here, we use the real-world distance (i.e., distance based on longitude and latitude) other than the country-level type (i.e., differentiate users based on their countries).

4.3. Reputation Algorithm. We firstly give a definition r_i as the reputation score of user i . If user i gets a higher reputation score, he/she can be considered as a reliable user. Then we propose an iterative and incorporative method to compute different users' reputation score:

$$\begin{cases} r_i^{k+1} = 1 - \frac{1}{1 + e^{-(d/|\Phi_i|) \sum_{j \in \Phi_i} |R_{ij} - A_j^{k+1}|}}, \\ A_j^{k+1} = \frac{\sum_{i \in \Gamma_j} R_{ij} \times r_i^k}{|\Gamma_j|}, \end{cases} \quad (9)$$

where k is the k^{th} iteration and d represents the damping factor in $[0, 1]$, r_i^k is the k^{th} reputation iteration of user i , Φ_i are user i 's invocated services, A_j^{k+1} represents average value of service j , and Γ_j denotes users who invocated the service j . Each user is assumed to be trustworthy, and the reputation will be assigned an initial value of 1 in first step of the iteration.

From the above discussion, it is obvious that it is an iterative process to calculate users' reputation. The reputation score of each user is computed by the difference between the available QoS values in original matrix R and the average QoS values provided by the target user on all invocated services.

4.4. QoS Prediction Based on Location and Reputation. The method in the previous section is the original matrix factorization method to predict QoS values with historical invocation records. To take full advantage of users' location data and repudiation score, a high-performance hybrid matrix factorization method is proposed as follows:

$$\begin{aligned} \min_{U, S} \psi(U, S) &= \frac{1}{2} \sum_{i=1}^m r_i \sum_{j=1}^n I_{ij} (R_{ij} - U_i^T S_j)_F^2 \\ &+ \frac{\gamma}{2} \sum_{i=1}^m \left\| U_i - \frac{\text{loc_sim}(i, j)}{|G(i)|} \sum_{U_j \in G(i)} U_j \right\| \\ &+ \frac{\lambda_u}{2} |U|^2 + \frac{\lambda_s}{2} |S|^2, \end{aligned} \quad (10)$$

where parameter r_i denotes the reputation score of user i . The previous studies make QoS prediction based on the premise that users in $G(i)$ suffer similar service invocations and observe similar QoS usage experience [15, 28]. However,

there are some unreliable users who would provide untrustworthy QoS values and make bad impact on prediction result. Therefore, we introduce r_i into Equation (10) to regulate the credibility of different users. Meanwhile, users' location information should be introduced into our method to enhance the prediction accuracy as is mentioned in the previous section. Therefore, a locational constraint item is defined as follows:

$$\frac{\gamma}{2} \sum_{i=1}^m \left\| U_i - \frac{\text{loc_sim}(i, j)}{|G(i)|} \sum_{U_j \in G(i)} U_j \right\|, \quad (11)$$

where $\gamma > 0$ denotes the relative proportion of the location grouping, U_i represents the latent factors of user i , $|G(i)|$ represents the subset of users who are near to user i , and $\text{loc_sim}(i, j)$ represents the similarity between user i and user j , defined as

$$\text{loc_sim}(i, j) = 1 - \frac{1}{1 + e^{-\text{dis}(i, j)}}, \quad (12)$$

where $\text{dis}(i, j)$ is the real-world distance between user i and user j , calculated by Equation (7).

The gradient parts of objective function Equation (10) could be calculated by employing the gradient descent method in U_i and S_j :

$$\begin{aligned} \frac{\partial \Psi}{\partial U_i} &= r_i \sum_{j=1}^n I_{ij} (R_{ij} - U_i^T S_j) (-S_j) + \lambda_u U_i \\ &+ \gamma \left(U_i - \frac{\text{loc_sim}(i, j)}{|G(i)|} \sum_{U_j \in G(i)} U_j \right), \end{aligned} \quad (13)$$

$$\frac{\partial \Psi}{\partial S_j} = r_i \sum_{i=1}^m I_{ij} (R_{ij} - U_i^T S_j) (-U_i) + \lambda_s S_j. \quad (14)$$

Based on the update process in Equations (5) and (6), we can update U_i and S_j with the two derivative Equations (13) and (14) until we get the local minimum of objective function Equation (10).

5. Experiments

5.1. Experiment Setup. A series of experiments are conducted on a well-known public QoS dataset, which is provided in the previous related work [7]. The dataset contains 1,974,675 response time records of service invocation, which is collected from 339 users on 5,825 services. Figure 5 shows the distribution of all values in the dataset and nearly 90% of all response time values are in the range $[0, 2]$. To simulate a real-world service invocation scenario, several values of response time records are randomly removed to generate random unreliable users, and different numbers of unreliable users will be introduced into the dataset to study the impact of users' reputation on the prediction accuracy.

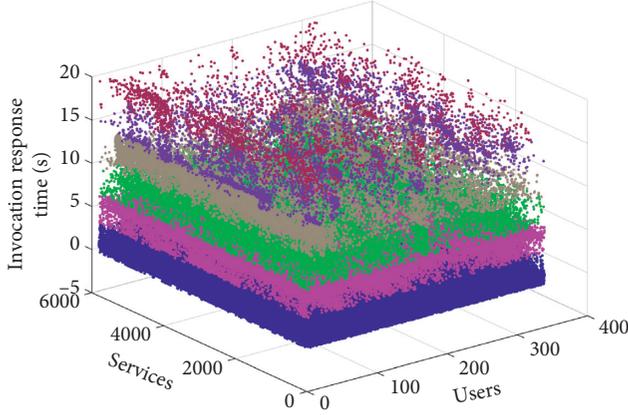


FIGURE 5: The distribution of response time values in real-world service invocation.

5.2. Evaluation Metrics. In the experiment, we utilize two famous statistical metrics, i.e., mean absolute error (MAE) and root mean squared error (RMSE) to evaluate our predicted result. If the prediction result is closer to the actual QoS value, the smaller value of MAE and RMSE will be generated and higher prediction accuracy can be achieved. MAE is given by

$$\text{MAE} = \frac{\sum_{i,j} |R_{ij} - \tilde{R}_{ij}|}{N}, \quad (15)$$

and RMSE is calculated as follows:

$$\text{RMSE} = \frac{\sum_{i,j} \sqrt{|R_{ij} - \tilde{R}_{ij}|^2}}{N}, \quad (16)$$

where R_{ij} represents the available QoS value in original matrix, \tilde{R}_{ij} denotes the predicted QoS value, and N represents the number of unknown QoS values.

5.3. Comparison Study. The proposed approach is compared with following CF approaches:

- (i) UPCC. This approach utilizes PCC to calculate similarity between users [37]. In service computing researches, this approach can be employed to predict QoS values.
- (ii) IPCC. This approach is a common commercial recommendation method. Service recommendation systems usually utilize this method to calculate service similarities and make prediction [38].
- (iii) UIPCC. Both UPCC and IPCC are employed synchronously in the prediction framework [39].
- (iv) PMF. This approach uses probability theory to explain how to use matrix factorization make prediction [11].
- (v) LBR1. This approach predicts QoS values by combining geographical information and matrix factorization approach [15].

- (vi) RMF. This approach utilizes available QoS values to calculate users' reputation and adopt matrix factorization model to predict missing QoS values of unknown services [32].

To simulate service invocation cases, a set number of elements are removed from the original QoS matrix. After this data preprocess, the density of the final matrix is set to 5%, 10%, 15%, and 20%, respectively. The parameter d is set 0.1 to compute users' reputation. We also set $\lambda_u = \lambda_s = 0.001$, dimensionality = 50, and $\gamma = 0.001$. The number of unreliable users is equal to 10 about 2.79 percentage of all users. The overall comparison details are presented in Tables 1 and 2.

The comparison results in Tables 1 and 2 show that the proposed LRMF can achieve higher prediction accuracy than other state-of-the-art approaches, which indicates LRMF has greatly improved QoS prediction accuracy. According to the comparison result, LRMF could achieve the best prediction performance.

5.4. Impact of Unreliable Users. The number of unreliable users determines the untrustworthy QoS values in the training dataset, which will produce great influence on the prediction method. In the experiment, the number of unreliable users is set from 10 to 80 under the condition of dimensionality in 80 and matrix density in 5%, 10%, 15%, and 20%, respectively.

As shown in Figure 6, both MAE and RMSE values of LRMF are significantly smaller when matrix density becomes denser under different conditions. It can also be proven that more available QoS values will make better prediction result. When the number of unreliable users changes in the range of [30, 80], both MAE and RMSE do not increase so much, which demonstrates that our LRMF method could minimize the negative impact produced by unreliable users and improve QoS prediction accuracy.

5.5. Impact of Parameter γ . The parameter γ determines the proportion of the location region factor in our proposed method. For one thing, too large value of γ will create a strong relation between the prediction accuracy and the geographical region; for another, if γ is assigned too small, the location region cannot generate enough contribution to the objective function. The comparative study for γ is introduced with the condition of dimensionality = 50 and matrix density in 5% and 20%. Also, we add 10 unreliable users in experiment settings. The details of the experiment are presented in Figure 7.

As presented in Figures 7(a) and 7(b), the MAE can achieve a minimum when γ reaches a certain value 10^{-3} and when γ is smaller or larger than this specific value, the MAE will fluctuate wildly. It is similar in Figures 7(c) and 7(d), the RMSE will reach a certain threshold when γ increases to 10^{-2} at beginning, but it will increase slightly after the threshold. The analysis shows that users' location data

TABLE 1: Prediction accuracy comparison result of mean squared error (MAE).

Methods	Matrix density			
	5%	10%	15%	20%
UPCC	0.5931	0.5497	0.5195	0.4908
IPCC	0.6239	0.5901	0.5519	0.5227
UIPCC	0.5962	0.5319	0.5029	0.4749
PMF	0.5797	0.5201	0.4872	0.4621
LBR1	0.5647	0.5021	0.4709	0.4498
RMF	0.5639	0.4897	0.4641	0.4529
LRMF	0.5514	0.4719	0.4493	0.4384

TABLE 2: Prediction accuracy comparison result of root mean squared error (RMSE).

Methods	Matrix density			
	5%	10%	15%	20%
UPCC	1.4127	1.3301	1.2697	1.2492
IPCC	1.4371	1.3561	1.2898	1.2139
UIPCC	1.3969	1.3071	1.2598	1.1815
PMF	1.4481	1.2894	1.2319	1.1839
LBR1	1.4431	1.2891	1.2208	1.1641
RMF	1.4372	1.2792	1.2105	1.1621
LRMF	1.4152	1.2584	1.1956	1.1498

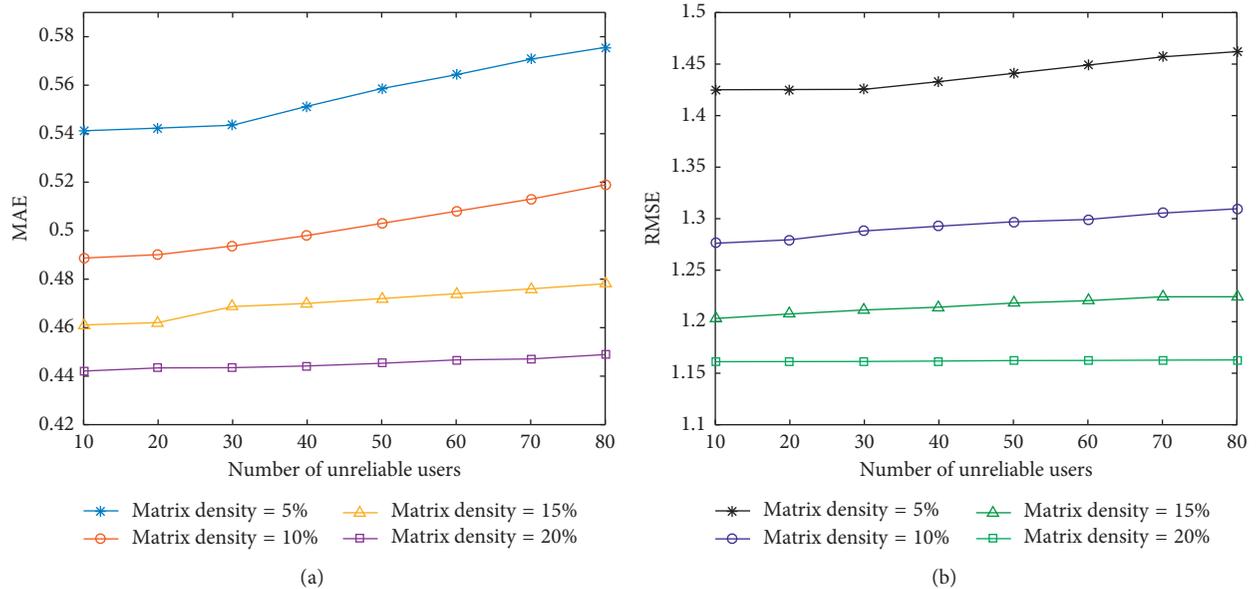


FIGURE 6: Impact of the number of unreliable users.

would make a remarkable impact on the prediction accuracy of LRMF.

5.6. Impact of Parameter ϵ . In LRMF, the location group threshold ϵ determines the geographical region size. If ϵ is assigned too small, the region would be very small and users in the region would have a very short distance. If ϵ is assigned too large, much more users would be identified into the same geographical region, which would lead to more noise data and neglect the local factor.

To study how parameter ϵ impacts on the prediction approach, the value of dimensionality is assigned 50 under the condition of matrix density in 10% and 30%. We also add 10 unreliable users to experiment settings.

Figure 8 illustrates how parameter ϵ affects the prediction accuracy. It is obvious that the evaluation values decrease when parameter ϵ increases firstly, but both two kinds of evaluation values gradually increase when parameter ϵ passes over a threshold. The observed phenomenon could be considered as when parameter ϵ is smaller

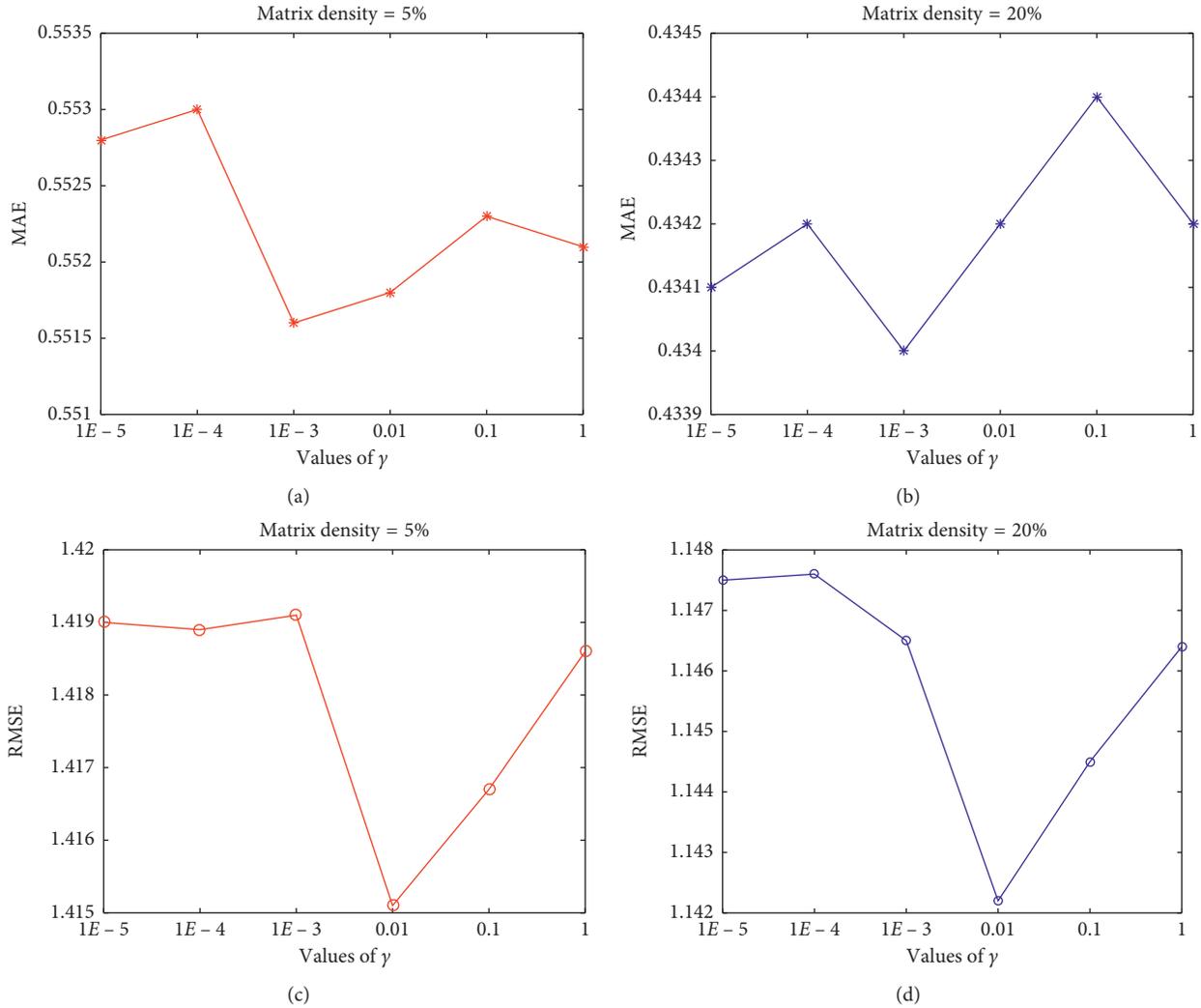


FIGURE 7: Impact of parameter γ .

than a certain value, the location region lacks enough similar users for active user, which prevents the crowds to contribute their collective intelligence. When parameter ϵ is assigned to be a large value, too much noise data will be included in the location region. Both two cases would produce negative impacts on the prediction curacy.

5.7. Impact of Density. To study what impacts does matrix density have on the prediction result, we add 10 unreliable users and set dimensionality to 10 and 50 when assigning matrix density from 5% to 20%.

As shown in Figure 9, both MAE and RMSE decrease firstly when the matrix density increases at first. Then, the curve becomes flat when the matrix density continues to increase. These comparative details show that the sparsity of original data would have great impact on the prediction accuracy. If more additional entries are available, the proposed method could get better prediction result. This observation demonstrates that when the original sparse

matrix becomes denser by collecting more QoS values, the prediction accuracy can be greatly enhanced in our proposed method.

5.8. Impact of Dimensionality. The number of latent feature vectors is regulated by dimensionality in our LRMF method. We vary the dimensionality from 10 to 100 under the condition of matrix density in 5% and 10 and 30 unreliable users, respectively, to conduct the comparative experiments.

Figure 10 shows that a proper value of dimensionality can achieve better prediction result as we can get a smallest prediction error. Both of MAE and RMSE decrease firstly because of more latent factors are added into the factorization process. When the dimensionality overpasses a certain threshold, more noised data may be brought into the training model with overfitting problem. As a result, the threshold of dimensionality in our model is approximately assigned 80.

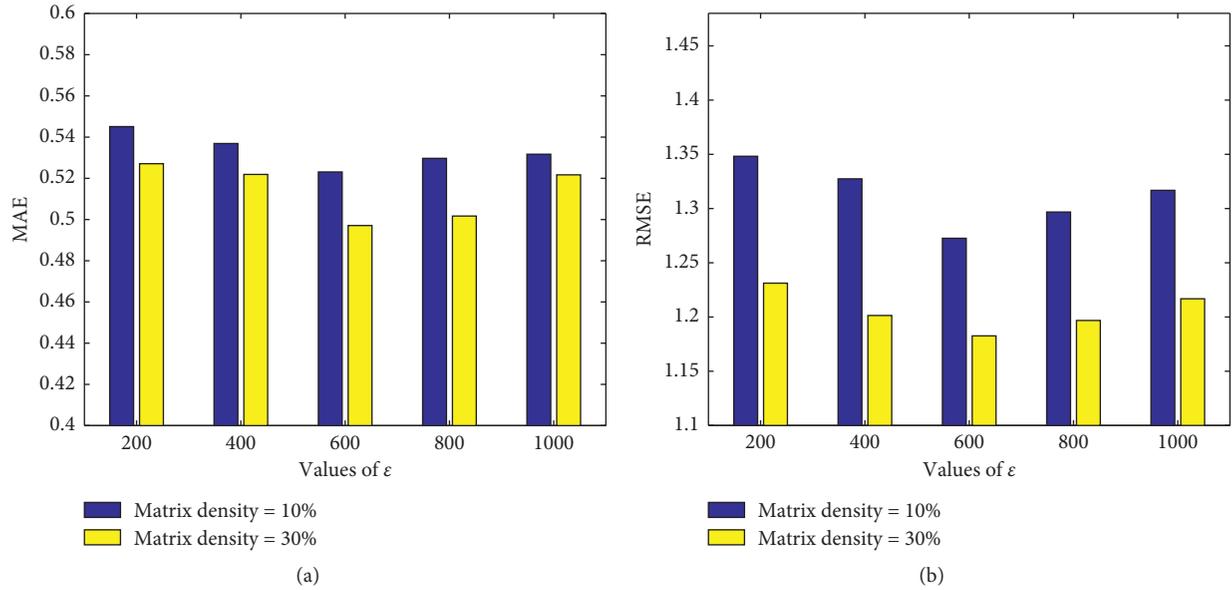
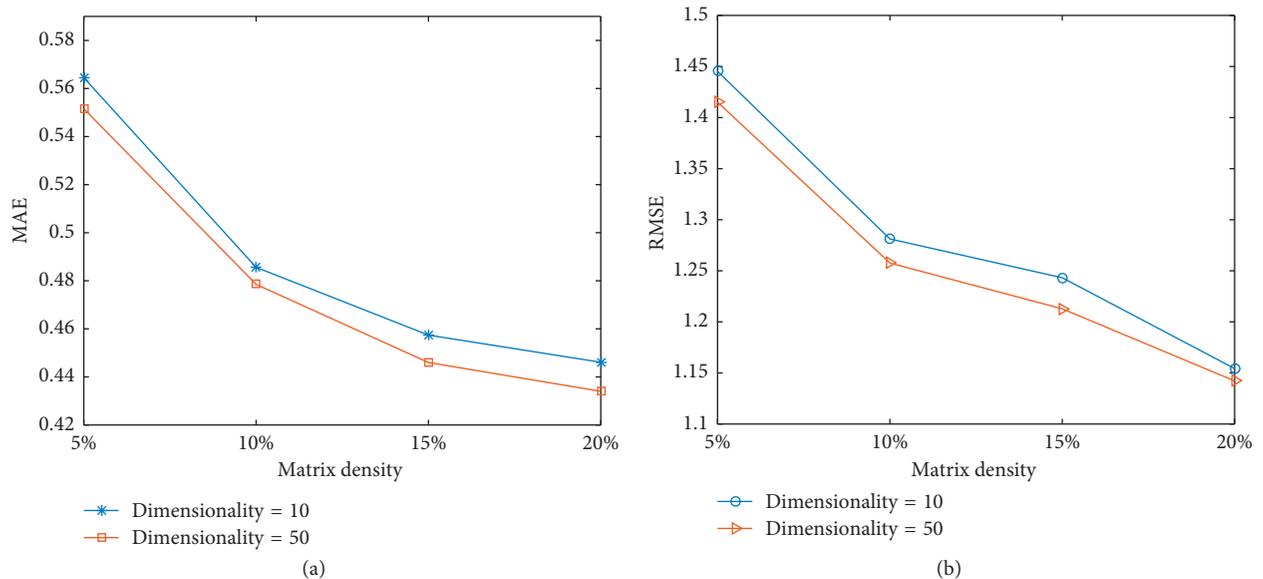
FIGURE 8: Impact of parameter ϵ .

FIGURE 9: Impact of matrix density.

6. Conclusions and Future studies

Based on the premise that the reputation of mobile service users would make great effect on the unknown QoS prediction, this paper propose an efficient prediction method by simultaneously exploring users' reputation and geographic distribution to make personalized service recommendation. We firstly cluster service users into different locational groups by their real-world geographical information and then calculate their reputations by the historical QoS values. At last, a hybrid matrix factorization model is proposed by integrating users' reputation and geographic data to predict

unknown QoS values. Experimental analysis on public QoS dataset demonstrates the high-performance and effectiveness of our LRMF on QoS prediction. The analysis shows that there are some unreliable uses in some location region and they submitted untrustworthy QoS values to gain benefits for their own services. The mobile service recommendation approach proposed in this study could reduce the poor effect of unreliable users and recommend high-quality and credible mobile service to end users.

In this paper, we only consider users' information as a significant issue to predict unknown QoS values for mobile service recommendation. In fact, the geographical distribution

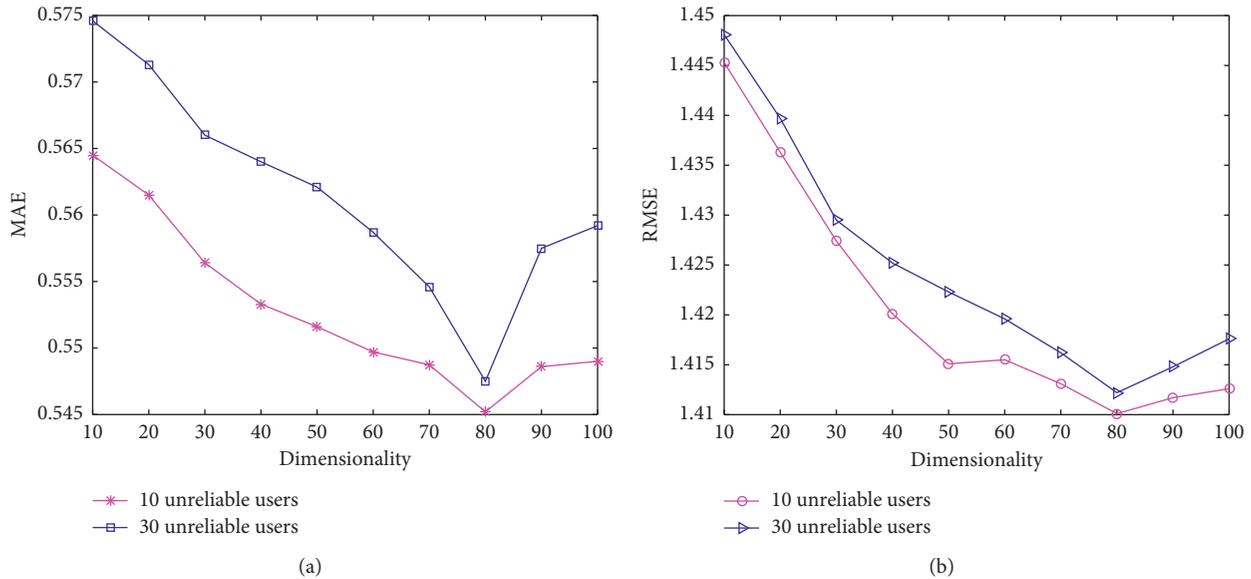


FIGURE 10: Impact of dimensionality.

of services would also provide useful facilities when identifying users' location region based on the distance between user and service. Therefore, it is potential to design more accurate location grouping model by combining both users' and service's location information. Besides, if the algorithms of users' reputation computing have deficiencies, then we would try to introduce intelligence methods, e.g., deep learning, reinforcement learning to design optimal algorithms for QoS prediction. Furthermore, the reliability of user may be affected by their trusted friends, so we will continue to track users' reputation in their social network to make high-quality service recommendation in our future work.

Data Availability

The real number data WS-DREAM used to support the findings of this study are included within the article.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

This research was funded by National Natural Science Foundation of China (nos. 6167060382 and 61602070), New Academic Seedling Cultivation and Exploration Innovation Project (no. [2017]5789-21), and China Scholarship Council (no. 201706050085).

References

- [1] Y. Yin, W. Xu, Y. Xu, H. Li, and L. Yu, "Collaborative QoS prediction for mobile service with data filtering and SlopeOne model," *Mobile Information Systems*, vol. 2017, Article ID 7356213, 14 pages, 2017.
- [2] L. Li, S. Li, and S. Zhao, "Qos-aware scheduling of services-oriented internet of things," *IEEE Transactions Industrial Informatics*, vol. 10, no. 2, pp. 1497–1505, 2014.
- [3] Y. Ma, S. Wang, P. C. K. Hung, C. H. Hsu, Q. Sun, and F. Yang, "A highly accurate prediction algorithm for unknown web service QoS values," *IEEE Transactions on Services Computing*, vol. 9, no. 4, pp. 511–523, 2016.
- [4] L. Xin, M. C. Zhou, L. Shuai, Y. N. Xia, Z. H. You, and Q. S. Zhu, "Incorporation of efficient second-order solvers into latent factor models for accurate prediction of missing qos data," *IEEE Transactions on Cybernetics*, vol. 48, no. 4, pp. 1216–1228, 2017.
- [5] N. Xi, D. Lu, C. Sun, J. Ma, and Y. Shen, "Distributed secure service composition with declassification in mobile clouds," *Mobile Information Systems*, vol. 2017, Article ID 7469342, 13 pages, 2017.
- [6] S. Wang, L. Sun, Q. Sun, X. Li, and F. Yang, "Efficient service selection in mobile information systems," *Mobile Information Systems*, vol. 2015, Article ID 949436, 10 pages, 2015.
- [7] Z. Zheng, H. Ma, M. R. Lyu, and I. King, "Qos-aware web service recommendation by collaborative filtering," *IEEE Transactions on Services Computing*, vol. 4, no. 2, pp. 140–152, 2010.
- [8] M. Silic, G. Delac, and S. Srblic, "Prediction of atomic web services reliability for QoS-aware recommendation," *IEEE Transactions on Services Computing*, vol. 8, no. 3, pp. 425–438, 2015.
- [9] S. Li, J. Wen, F. Luo, M. Gao, J. Zeng, and Z. Y. Dong, "A new qos-aware web service recommendation system based on contextual feature recognition at server-side," *IEEE Transactions on Network and Service Management*, vol. 14, no. 2, pp. 332–342, 2017.
- [10] G. Kang, J. Liu, M. Tang, B. Cao, and Y. Xu, "An effective web service ranking method via exploring user behavior," *IEEE Transactions on Network and Service Management*, vol. 12, no. 4, pp. 554–564, 2015.
- [11] A. Mnih and R. Salakhutdinov, "Probabilistic matrix factorization," in *Proceedings of Advances in Neural Information Processing Systems 2007*, pp. 1257–1264, Kitakyushu, Japan, December 2007.

- [12] M. Tang, Z. Zheng, G. Kang, J. Liu, Y. Yang, and T. Zhang, "Collaborative web service quality prediction via exploiting matrix factorization and network map," *IEEE Transactions on Network and Service Management*, vol. 13, no. 1, pp. 126–137, 2017.
- [13] X. Chen, X. Liu, Z. Huang, and H. Sun, "Regionknn: a scalable hybrid collaborative filtering algorithm for personalized web service recommendation," in *Proceedings of IEEE 17th International Conference on Web Services (ICWS'10)*, pp. 9–16, Miami, FL, USA, July 2010.
- [14] <http://www.google.com/transparencyreport/>.
- [15] L. Wei, J. Yin, S. Deng, Y. Li, and Z. Wu, "Collaborative web service QoS prediction with location-based regularization," in *Proceedings of IEEE 19th International Conference on Web Services (ICWS'12)*, pp. 464–471, Honolulu, HI, USA, June 2012.
- [16] I. Gunes, C. Kaleli, A. Bilge, and H. Polat, "Shilling attacks against recommender systems: a comprehensive survey," *Artificial Intelligence Review*, vol. 42, no. 2, pp. 767–799, 2014.
- [17] E. Al-Masri and Q. H. Mahmoud, "QoS-based discovery and ranking of web services," in *Proceedings of IEEE 16th International Conference on Communications and Networks (ICCCN'07)*, pp. 529–534, Honolulu, HI, USA, August 2007.
- [18] K. Kritikos and D. Plexousakis, "Requirements for QoS-based web service description and discovery," *IEEE Transactions on Services Computing*, vol. 2, no. 4, pp. 320–337, 2009.
- [19] S. Rosario, A. Benveniste, S. Haar, and C. Jard, "Probabilistic QoS and soft contracts for transaction-based web services orchestrations," *IEEE Transactions on Services Computing*, vol. 1, no. 4, pp. 187–200, 2009.
- [20] J. E. Hadad, M. Manouvrier, and M. Rukoz, "TQoS: transactional and QoS-aware selection algorithm for automatic web service composition," *IEEE Transactions on Services Computing*, vol. 3, no. 1, pp. 73–85, 2010.
- [21] G. Linden, B. Smith, and J. York, "Amazon.com recommendations: item-to-item collaborative filtering," *IEEE Internet Computing*, vol. 7, no. 1, pp. 76–80, 2003.
- [22] D. A. Adeniyi, Z. Wei, and Y. Yongquan, "Automated web usage data mining and recommendation system using K-Nearest Neighbor (KNN) classification method," *Applied Computing and Informatics*, vol. 12, no. 1, pp. 90–108, 2016.
- [23] Y. Zuo, J. Zeng, M. Gong et al., "Tag-aware recommender systems based on deep neural networks," *Neurocomputing*, vol. 204, pp. 51–60, 2016.
- [24] B. Engelbert, M. B. Blanken, R. Kruthoff-Bruwer, and K. Morisse, "A user supporting personal video recorder by implementing a generic Bayesian classifier based recommendation system," in *Proceedings of IEEE 19th International Conference on Communications and Networks (ICCCN'11), Workshops*, pp. 567–571, Seattle, WA, USA, March 2011.
- [25] C. X. Zhang, Z. K. Zhang, L. Yu, C. Liu, H. Liu, and X. Y. Yan, "Information filtering via collaborative user clustering modeling," *Physica A Statistical Mechanics and Its Applications*, vol. 396, no. 2, pp. 195–203, 2014.
- [26] H. Ma, H. Yang, M. R. Lyu, and I. King, "SoRec:social recommendation using probabilistic matrix factorization," in *Proceedings of ACM Conference on Information and Knowledge Management*, pp. 931–940, ACM, Napa Valley, CA, USA, November 2008.
- [27] N. Ali and C. Solis, "Self-adaptation to mobile resources in service oriented architecture," in *Proceedings of 2015 IEEE International Conference on Mobile Services*, pp. 407–414, New York, NY, USA, June 2015.
- [28] M. Tang, T. Zhang, J. Liu, and J. Chen, "Cloud service qos prediction via exploiting collaborative filtering and location-based data smoothing," *Concurrency and Computation Practice and Experience*, vol. 27, no. 18, pp. 5826–5839, 2015.
- [29] K. Lee, J. Park, and J. Baik, "Location-based web service QoS prediction via preference propagation for improving cold start problem," in *Proceedings of IEEE 22nd International Conference on Web Services (ICWS '15)*, pp. 400–407, New York, NY, USA, June 2015.
- [30] B. Gonsalves and V. Patil, "Improved web service recommendation via exploiting location and QoS information," in *Proceedings of International Conference on Information Communication and Embedded Systems IEEE*, pp. 1–5, Boca Raton, FL, USA, December 2016.
- [31] W. Qiu, Z. Zheng, X. Wang, X. Yang, and M. R. Lyu, "Reputation-aware QoS value prediction of web services," in *Proceedings of the IEEE 10th International Conference on Service Computing (SCC'13)*, pp. 41–48, Santa Clara, CA, USA, June–July 2013.
- [32] J. Xu, Z. Zheng, and M. R. Lyu, "Web service personalized quality of service prediction via reputation-based matrix factorization," *IEEE Transactions on Reliability*, vol. 65, no. 1, pp. 28–37, 2015.
- [33] M. Mehdi, N. Bouguila, and J. Bentahar, "Trust and reputation of web services through QoS correlation lens," *IEEE Transactions on Services Computing*, vol. 9, no. 6, pp. 968–981, 2016.
- [34] A. Comi, L. Fotia, F. Messina, G. Pappalardo, and D. Rosaci, "A reputation-based approach to improve QoS in cloud service composition," in *Proceedings of IEEE, International Conference on Enabling Technologies: Infrastructure for Collaborative Enterprises*, pp. 108–113, Paris, France, June 2015.
- [35] K. K. Fletcher and X. F. Liu, "A collaborative filtering method for personalized preference-based service recommendation," in *Proceedings of IEEE 22nd International Conference on Web Services (ICWS'15)*, pp. 400–407, New York, NY, USA, June 2015.
- [36] Z. Zheng and M. R. Lyu, *QoS Management of Web Services*, Springer, Berlin, Germany, 2013.
- [37] L. H. Son, "H-FCF: a hybrid user-based fuzzy collaborative filtering method in recommender systems," *Expert Systems with Applications*, vol. 41, no. 15, pp. 6861–6870, 2014.
- [38] K. Y. Chung, D. Lee, and K. J. Kim, "Categorization for grouping associative items mining in item-based collaborative filtering," *Multimedia Tools and Applications*, vol. 71, no. 2, pp. 889–904, 2014.
- [39] M. Papagelis and D. Plexousakis, "Qualitative analysis of user-based and item-based prediction algorithms for recommendation agents," *Engineering Applications of Artificial Intelligence*, vol. 18, no. 7, pp. 781–789, 2005.

Research Article

Dilemma and Solution of Traditional Feature Extraction Methods Based on Inertial Sensors

Zhiqiang Peng  and Yue Zhang 

The Division of Information Science and Technology, Graduate School at Shenzhen, Tsinghua University, Shenzhen, China

Correspondence should be addressed to Zhiqiang Peng; pzqds@126.com

Received 6 July 2018; Accepted 26 September 2018; Published 22 November 2018

Guest Editor: Subramaniam Ganesan

Copyright © 2018 Zhiqiang Peng and Yue Zhang. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Correctly identifying human activities is very significant in modern life. Almost all feature extraction methods are based directly on acceleration and angular velocity. However, we found that some activities have no difference in acceleration and angular velocity. Therefore, we believe that for these activities, any feature extraction method based on acceleration and angular velocity is difficult to achieve good results. After analyzing the difference of these indistinguishable movements, we propose several new features to improve accuracy of recognition. We compare the traditional features and our custom features. In addition, we examined whether the time-domain features and frequency-domain features based on acceleration and angular velocity are different. The results show that (1) our custom features significantly improve the precision of the activities that have no difference in acceleration and angular velocity; and (2) the combination of time-domain features and frequency-domain features does not significantly improve the recognition of different activities.

1. Introduction

The classification of human motion based on inertial sensors has been proven to have many important applications in the medical and health fields. In previous studies, time-domain and frequency-domain features are widely used for feature calculation.

There are many studies that use wavelet transform to extract features to classify human activities. However, the research of Preece et al. [1] shows that the time- and frequency-domain features often exceed wavelet features, indicating that the wavelet feature may be not the most effective method for calculating the human body motion classification features.

Some time-domain features are derived to classify human activities, such as the mean, median, variance, skewness, kurtosis [2], and interquartile range [3]. In order to extract frequency-domain features, the sensor data window is first changed to the frequency domain using discrete Fourier Transform [4]. Then, we can extract some features from the frequency domain to distinguish different activities,

such as power spectral density (PSD) [5], peak frequency [5, 6], entropy [7], DC component [7], median frequency [8], spectral energy [9], and frequency-domain entropy [10]. Of course, there are other methods that process data from accelerometers and gyroscopes. But all in all, to the best of our knowledge, these features are extracted directly from acceleration and angular velocity, which inevitably have some common drawbacks.

We studied 12 kinds of activities and found it easy to confuse elevator up and elevator down. These two kinds of activities do not have obvious differences in acceleration and angular velocity, so the time and frequency-domain features based on acceleration and angular velocity cannot achieve good classification results. Therefore, for those motions that have no significant difference in angular velocity and acceleration, no matter how the features are extracted from the acceleration and angular velocity, it is difficult to achieve good results.

In addition, we begin to wonder if there is any essential difference between time-domain features and frequency-domain features based on acceleration and angular

velocity. In order to solve this discredit, we have separately tested the effects of time-domain features and frequency-domain features. Then, we tested the combination of the two kinds of features and found that the combination of time-domain features and frequency-domain features was slightly higher than only time-domain features. From the experimental results analysis, we believe that, for human activity classification problem, the time-domain features and the frequency-domain features are two aspects of the same rules, and there is no essential difference.

Our contributions in this paper are two-fold: (1) to the best of our knowledge, this is the first time that features of motion classification based on velocity and displacement have been proposed, which solve some problems that cannot be solved by traditional time- and frequency-domain features; (2) As far as we know, for the first time, we have studied the difference between the time-domain features and frequency-domain features.

2. Methods

Indeed, the time-domain and frequency-domain features based on acceleration and angular velocity have achieved some success. However, for activities without obvious difference between acceleration and angular velocity, such as elevator up and elevator down, the traditional method of extracting features based on acceleration and angular velocity is difficult to work. In order to solve this problem, we carefully analyze the two activities of elevator up and elevator down, summarize the differences between them, and propose some new features for distinguishing such activities. After analysis, we sum up the following rules:

- (i) When the elevator goes up, the speed is upward; when the elevator goes down, the speed is downward.
- (ii) When the elevator just starts to move or stops moving, its speed is small and its angular speed is large.
- (iii) When the elevator just starts to rise or fall, the direction of the speed is the same as the direction of the acceleration; when the elevator stops to rise or fall, the direction of the speed is opposite to the direction of the acceleration.

Based on the above evidences, we propose four features to distinguish these activities on each axis of the accelerometer. First, according to the first rule, we introduce three features, namely, the starting speed, the ending speed, and the displacement. Second, according to the second and third rules, we introduce the fourth feature. As some activities have just started and are about to

stop, their speed is small and difficult to distinguish. Therefore, we extracted another feature to enhance the difference between two activities. When the velocity direction is the same as the acceleration direction, we use $v + a$ as the feature; otherwise, we use $v - a$ as feature. In order to describe the movement of the human body in different directions as much as possible, we have introduced the following twelve new features to enhance the difference between different activities. These features are summarized in Table 1.

Suppose the time window we choose is T , the sampling frequency is n , then the total number of samples is Tn . In the experiment, the time window we selected was two seconds. The displacement, time, and acceleration corresponding to the i th sampling interval are $x(i)$, $t(i)$, and $a(i)$. The speed corresponding to the sampling point is $v(i)$. Data section of one time window is shown in Figure 1. The gravity acceleration is g . Due to the high sampling frequency, the time interval between each sample point is short. At the same time, in order to simplify the calculation, we believe that there is uniform linear motion between each sampling point. Now, we derive the speed and displacement formula.

We think these sampling points are equally time-distributed, so we have the following conclusions.

$$t(1) = t(2) = t(3) = \dots = t(Tn) = \frac{T}{Tn} = \frac{1}{n}. \quad (1)$$

First, we derive the formula for the end speed of each axis. According to the kinematics formula, we can easily get the following formula and simplify it by combining it with Equation (1).

$$\begin{aligned} v(Tn) &= v(0) + \sum_{i=1}^{Tn} (a(i) - g(i)) * t(i) \\ &= v(0) + \frac{1}{n} \sum_{i=1}^{Tn} (a(i) - g(i)) \\ &= v(0) + \frac{1}{n} \sum_{i=1}^{Tn} a(i) - \frac{1}{n} \sum_{i=1}^{Tn} g(i). \end{aligned} \quad (2)$$

Next, we introduce the determination of the starting speed. If this window is the first window, we default to a starting speed of zero. Otherwise, the starting speed is the end speed of the previous window marked as $v_{-1}(Tn)$:

$$v(0) = \begin{cases} 0, & \text{the first window,} \\ v_{-1}(Tn), & \text{not the first window.} \end{cases} \quad (3)$$

As for the displacement of each axis, since we think that there is uniform linear motion between each sampling point, the displacement formula can be derived as follows.

TABLE 1: Custom features.

Feature name	Description
End velocity along X	End speed along the x -axis
End velocity along Y	End speed along the y -axis
End velocity along Z	End speed along the z -axis
Starting velocity along X	Starting speed along the x -axis
Starting velocity along Y	Starting speed along the y -axis
Starting velocity along Z	Starting speed along the z -axis
Displacement along X	Displacement along the x -axis
Displacement along Y	Displacement along the y -axis
Displacement along Z	Displacement along the z -axis
Velocity plus acceleration along X	$v + a$ if the product of v along the x -axis and a is positive, $v - a$ otherwise
Velocity plus acceleration along Y	$v + a$ if the product of v along the y -axis and a is positive, $v - a$ otherwise
Velocity plus acceleration along Z	$v + a$ if the product of v along the z -axis and a is positive, $v - a$ otherwise

$t(1)$	$t(2)$...	$t(i)$...	$t(Tn)$	
$a(1)$	$a(2)$		$a(i)$		$a(Tn)$	
$v(0)$	$v(1)$	$v(2)$	$v(i-1)$	$v(i)$	$v(Tn-1)$	$v(Tn)$
$x(0)$	$x(1)$	$x(2)$	$x(i-1)$	$x(i)$	$x(Tn-1)$	$x(Tn)$

FIGURE 1: Data section of one time window.

$$\begin{aligned}
x(Tn) - x(0) &= (x(1) - x(0)) + (x(2) - x(1)) + \dots + (x(Tn) - x(Tn-1)) \\
&= \left(v(0)t(1) + \frac{1}{2}(a(1) - g(1))t^2(1) \right) + \left(v(1)t(2) + \frac{1}{2}(a(2) - g(2))t^2(2) \right) + \dots \\
&\quad + \left(v(Tn-1)t(Tn) + \frac{1}{2}(a(Tn) - g(Tn))t^2(Tn) \right) \\
&= \frac{1}{n}(v(0) + v(1) + \dots + v(Tn-1)) + \frac{1}{2n^2} \sum_{i=1}^{Tn} (a(i) - g(i)) \\
&= \frac{1}{n} \left(v(0) + v(0) + \sum_{i=1}^{Tn-1} (a(i) - g(i))t(i) + v(0) + \sum_{i=1}^{Tn-1} (a(i) - g(i))t(i) + \dots + v(0) + \sum_{i=1}^{Tn-1} (a(i) - g(i))t(i) \right) \\
&\quad + \frac{1}{2n^2} \sum_{i=1}^{Tn} (a(i) - g(i)) \\
&= T * v(0) + \frac{1}{n^2} \sum_{i=1}^{Tn-1} (Tn - i)(a(i) - g(i)) + \frac{1}{2n^2} \sum_{i=1}^{Tn} (a(i) - g(i)). \tag{4}
\end{aligned}$$

In the experiment, in order to calculate the last three features we defined in Table 1, that is, Velocity Plus Acceleration Along X, Velocity Plus Acceleration Along Y, and Velocity Plus Acceleration Along Z, our speed takes the end velocity of the two-second time window, and the acceleration takes the difference between the average acceleration of the two-second time window and the gravitational acceleration in each axis. In this way, we can calculate the last three features, and the expression is shown in (6), in which Velocity Plus Acceleration Along* stands for the last three features.

$$a = \frac{1}{Tn} \sum_{i=1}^{Tn} a(i), \tag{5}$$

$$\text{Velocity plus acceleration along } * = \begin{cases} v + (a - g), & v * (a - g) > 0, \\ v - (a - g), & v * (a - g) < 0. \end{cases} \tag{6}$$

Finally, we introduce the calculation method of gravity acceleration we used in the experiment. Since the tester is just wearing the device for data acquisition, it is generally at

a standstill and the starting speed is zero, which is used in our experiments. Therefore, for the sake of convenience, in the experiment, we believe that the initial acceleration of each axis is the component of gravity acceleration and assume that the component of gravity acceleration in each axis remains unchanged. We take the average of the first 10 sampling points of each axis as the component of gravity acceleration in each axis, recorded as g .

3. Experiments and Results

3.1. Datasets. In order to illustrate the validity of our custom features, we have selected the USC_HAD of University of Southern California as the verification dataset [11]. They use an off-the-shelf sensing platform called MotionNode to capture human activity signals and build their dataset. MotionNode is a 6-DOF inertial measurement unit (IMU) specifically designed for human motion sensing applications, which integrates a 3-axis accelerometer and, 3-axis gyroscope. They selected 14 subjects (7 male; 7 female) to participate in the data collection. The sampling frequency is 100 Hz. Twelve kinds of activities collected are Walking Forward, Walking Left, Walking Right, Walking Upstairs, Walking Downstairs, Running Forward, Jumping Up, Sitting, Standing, Sleeping, Elevator Up, and Elevator Down.

3.2. Results. In order to verify the validity of our custom features, we extracted some common time- and frequency-domain features. In the time domain, we chose the mean, median, variance, skewness, kurtosis, and interquartile range as time-domain features. In the frequency domain, we choose peak frequency, median frequency, power spectral density, DC component, spectral energy, and information entropy as frequency-domain features. Then, we added our custom features to these time- and frequency-domain features and compared their results. In the experiments, we adopted the two commonly used models, SVM and random forest. For SVM, the kernel function we use is a polynomial kernel function. For RF, the number of decision trees we choose is 50.

First, in order to check whether the distinction between the lift of the elevator and the descending of the elevator is achieved, we tested the precision and recall of our custom features on both models. Precision and recall are often used as performance measures for classifiers in classification problems. Table 2 shows the precision and recall of the models' identification of elevators up when adding custom features to time- and frequency-domain features and the combination of time and frequency-domain features without custom features. Table 3 shows the precision and recall of the models' identification of elevators down when adding custom features and no custom features. The left column below each model is the precision rate, and the right column is the recall rate. From Tables 2 and 3, we can see that, after adding the custom features, the model has significantly improved the recognition rate of the elevator up and the elevator down.

Also, we use the ROC (receiver operating characteristics) curve and corresponding AUC (area under the ROC curve)

TABLE 2: Precision and recall of the identification of elevators up when adding custom features and no custom features.

	SVM (%)		RF (%)	
	No custom features	50.60	55.23	56.82
Adding custom features	74.01	74.50	78.05	81.53

TABLE 3: Precision and recall of the identification of elevators down when adding custom features and no custom features.

	SVM (%)		RF (%)	
	No custom features	48.38	53.78	55.24
Adding custom features	71.38	71.61	83.23	85.67

values to check whether the distinction between the lift of the elevator, and the descending of the elevator is achieved. For SVM, we conducted a test. The ROC curves for the two types of classification results for elevator up and elevator down are shown in Figure 2.

From the above figure, we can see that in the SVM, after adding our custom features, the ROC curves of elevator up and elevator down completely cover the ROC curve of the original feature. The AUC value for the no custom features' ROC curve is 0.8572, while the AUC value for adding custom features' ROC curve is 0.9729. The ROC curve shows that our custom features have achieved very good results in distinguishing between the elevator up and the elevator down.

In order to further explain the significance of our custom features we have extracted and verify the difference between time-domain features and frequency-domain features, we have compared our custom features with the time-domain features and frequency-domain features. We performed comparative experiments on four combinations of features over SVM. We conducted five experiments for each experiment. The detailed experimental results are summarized in Table 4. The accuracy in the table is the total classification accuracy of the 12 activities. For convenience, we denote the time-domain feature as 1, the frequency-domain feature as 2, and the custom feature as 3.

From the experimental results, we can see that individual frequency-domain features and time-domain features can achieve good results. However, when frequency-domain features are combined with time-domain features, no significant improvement is obtained, indicating that there is no essential difference between features extracted from the frequency domain and features extracted from the time domain. The superposition of the two did not achieve better results. The time-domain feature is better when combined with our custom features than combined with frequency-domain features. Based on the comprehensive analysis, we believe that our custom feature is a supplement to traditional time-domain features rather than a redundant feature.

The traditional time-domain features and frequency-domain features are all based on the acceleration and angular velocity but there is no essential difference, so the superposition of the two will not bring significant improvement. Our custom feature is the mining of the rules of

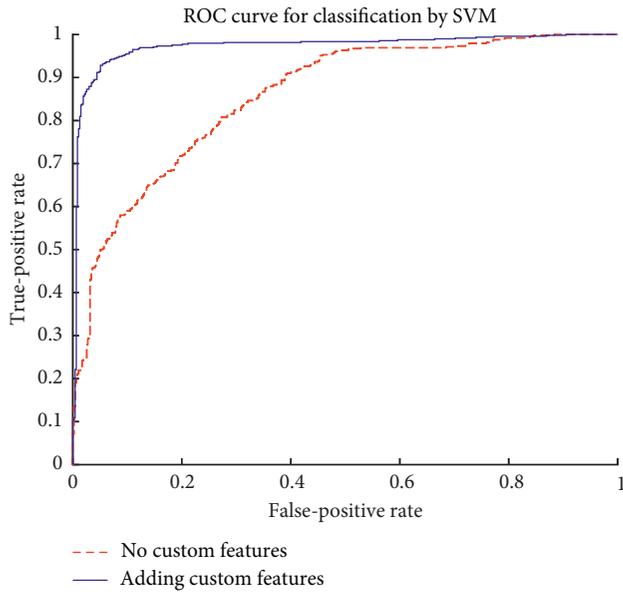


FIGURE 2: ROC curve of elevator up and elevator down over SVM.

TABLE 4: Classification accuracy of 12 activities on four combinations of frequency-domain features, time-domain features, and custom features over SVM.

	1 (%)	2 (%)	1, 2 (%)	1, 3 (%)
First experiment	89.67	83.93	90.46	92.21
Second experiment	90.52	84.11	90.80	92.17
Third experiment	89.78	83.07	91.27	91.94
Fourth experiment	89.34	83.50	90.62	92.61
Fifth experiment	89.29	83.93	90.21	92.25

speed and displacement, which is very different from the traditional mining of acceleration and angular velocity. So, when these features are introduced in the time-domain feature, we can obtain certain promotion. Especially for those motions which have no obvious difference between angular velocity and acceleration but there is a clear difference in speed and displacement, we can achieve good results with these custom features. For example, there is no obvious difference in acceleration and angular velocity in the smooth upward movement of the elevator and the smooth descending of the elevator, but there is a clear difference in speed and displacement. So, when introduce our custom feature, we can obviously increase the recognition rate of the two kinds of motions.

To calculate these features we define, we must know the initial state of motion, especially the initial state of speed. In our experiment, we assumed that the initial state is zero. In the database we use, most of the data are collected after the tester reaches a steady state of various motion postures, which does not satisfy our assumptions. If we can record data when the tester just wears the inertial sensors, so as to meet our assumptions, we believe we can achieve better results. For other types of sports, such as running, walking, and station, there is a difference in speed, which will bring a higher recognition rate.

4. Discussion

In this article, we begin with the elevator up and elevator down, which are indistinguishable based on the existing feature extraction methods and analyze the differences and rules between these two types of movements. Then, we have proposed four features on each axis of the accelerometer that have significantly improved the distinction between the two types of movements, elevator up, and elevator down. In the experiment, we found that the combination of frequency-domain features and time-domain features does not significantly improve the distinction of activities. The two kinds of features are two different aspects of acceleration and angular velocity, and there is no essential difference. From the experimental results, the time-domain features are better than the frequency-domain features and can more fully reflect the differences between different activities. Our custom features are not another response to acceleration, but instead, these features can be used to distinguish movements that differ in the speed of movement. In particular, it is of great significance to distinguish between movements that do not have a significant difference in acceleration and angular velocity but have a significant difference in speed.

In the experiment, for the sake of convenience, we assumed that the component of the gravitational acceleration remains unchanged, which is obviously not in line with the actual situation. Next, related personnel may consider introducing some basic theories of motion analysis in order to accurately calculate the components of the gravitational acceleration and thus more accurately calculate the features we introduce. We believe that when the features of velocity and displacement are introduced, we can make a great breakthrough in the existing human motion classification problem and to some extent get rid of the dilemma that some motions cannot be accurately identified on the features of acceleration and angular velocity.

Data Availability

The USC_HAD data used to support the findings of this study are included in Reference [11].

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

Acknowledgments

This work was supported by the Natural Science Foundation of China (No. 61571268).

References

- [1] S. J. Preece, J. Y. Goulermas, L. P. J. Kenney, and D. Howard, "A comparison of feature extraction methods for the classification of dynamic activities from accelerometer data," *IEEE Transactions on Biomedical Engineering*, vol. 56, no. 3, pp. 871–879, 2009.

- [2] J. Baek, G. Lee, W. Park, and B. J. Yun, "Accelerometer signal processing for user activity detection," in *Lecture Notes in Computer Science*, vol. 3215, pp. 610–617, Springer, Santa Barbara, CA, USA, 2004.
- [3] U. Maurer, A. Rowe, A. Smailagic, and D. Siewiorek, "Location and activity recognition using ewatch: a wearable sensor platform," in *Lecture Notes in Computer Science*, pp. 86–102, Springer, Santa Barbara, CA, USA, 1970.
- [4] F. Attal, S. Mohammed, M. Dedabrishvili, F. Chamroukhi, L. Oukhellou, and Y. Amirat, "Physical human activity recognition using wearable sensors," *Sensors*, vol. 15, no. 12, pp. 31314–31338, 2015.
- [5] B. Nham, K. Siangliulue, and S. Yeung, *Predicting Mode of Transport from Iphone Accelerometer Data*, Machine Learning Final Projects, Stanford University, Stanford, CA, USA, 2008.
- [6] D. Figo, P. C. Diniz, D. R. Ferreira, and J. M. P. Cardoso, "Preprocessing techniques for context recognition from accelerometer data," *Personal and Ubiquitous Computing*, vol. 14, no. 7, pp. 645–662, 2010.
- [7] J. Ho, *Interruptions: Using Activity Transitions to Trigger Proactive Messages*, Massachusetts Institute of Technology, Cambridge, MA, USA, 2004.
- [8] S. J. Preece, J. Y. Goulermas, L. P. Kenney, D. Howard, K. Meijer, and R. Crompton, "Activity identification using body-mounted sensors—a review of classification techniques," *Physiological Measurement*, vol. 30, no. 4, pp. R1–R33, 2009.
- [9] T. M. Huynh and B. Schiele, *Analyzing Features for Activity Recognition*, ACM, New York, NY, USA, 2005.
- [10] B. Ling and S. S. Intille, "Activity recognition from user-annotated acceleration data," in *Lecture Notes in Computer Science*, vol. 3001, pp. 1–17, Springer, Santa Barbara, CA, USA, 2004.
- [11] M. Zhang and A. A. Sawchuk, "USC-HAD: a daily activity dataset for ubiquitous activity recognition using wearable sensors," in *Proceedings of 2012 ACM Conference on Ubiquitous Computing*, pp. 1036–1043, Pittsburgh, PA, USA, September 2012.

Research Article

A Case Study Analysis of Clothing Shopping Mall for Customer Design Participation Service and Development of Customer Editing User Interface

Ying Yuan¹ and Jun-Ho Huh ²

¹Department of Clothing & Textiles, Hanyang University, Seoul, Republic of Korea

²Department of Software, Catholic University of Pusan, Busan, Republic of Korea

Correspondence should be addressed to Jun-Ho Huh; 72networks@pukyong.ac.kr

Received 20 May 2018; Revised 8 August 2018; Accepted 12 September 2018; Published 11 November 2018

Guest Editor: Jaegeol Yim

Copyright © 2018 Ying Yuan and Jun-Ho Huh. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Following the development of networking and mobile devices, the technology of managing the offline information online is being conducted widely. Also, as the social services have become much more active, users are registering and managing their personal information on online websites and sharing it with other users to acquire the information they need. For modern people living in a smart city, the planning of smarter services is required. The convergence of ET and IT or advanced scientific technologies such as AI or Big Data is often mentioned whenever the smart city is discussed. Nevertheless, smart services that could introduce smart solutions to conventional industries or change existing lifestyles should also be considered. Therefore, this paper discusses a service related to the convergence of the traditional clothing industry with IT and a service wherein CT is converged with systems that allow customers to participate in the design work and share the designs they have created. In other words, this study is a case study of CT and IT services in the clothing industry and is inclusive of an apparel shopping mall service that encourages customer participation in design, a customer-oriented editing user interface, and a copyright management system. The results show that both production method and production capacity largely affect the user interface of apparel platform services, with customer freedom significantly correlated with their functional roles. Moreover, the lead index is shown to be one of the factors restraining customer freedom. With this analysis, an apparel shopping mall wherein customers participate in the design work has been developed especially for clothes with more complex designs. The shopping mall emphasizes functionality from the perspective of customer use. At the same time, an online environment for an apparel service appropriate for the smart city has been implemented.

1. Introduction

Following the development of networking and mobile devices, the technology of managing the offline information online is being conducted widely. Also, as the social services have become much more active, users are registering and managing their personal information on online websites and sharing it with other users to acquire the information they need. Recently, the application whose intelligent agent recognizes a user's habits or a lifestyle and provides proper information is increasingly appearing in the online market. For example, a service that provides the information about

the places the user often goes or the available means of transportation to user's destination is provided currently.

This is achieved by searching suitable information from the accumulated data pertaining to the user's history of visiting particular places. People nowadays pay more attention to their appearances especially when it concerns clothes, wearing different clothes depending on with whom they will be meeting or avoiding not to wear the same clothes the next day. Currently, a method which conveniently acquires the data on user's clothing habits and provides useful information to him/her based on the accumulated data is not available. Thus, to discuss a service where the existing

clothing industry, IT-converged services, and CT are integrated to allow customers to participate in the design process and share their designs with others, an application which recommends a suitable design to the designer by accumulating the history of clothing choices of these customers to grasp a particular customer's clothing habits or inclination has been proposed. All of these tasks can be performed and managed with a smart device.

Technical development has accelerated the change of social structure and contemporary lifestyles [1]. Technical advances are transforming us from a postmodern to a participating generation [2]. The normalization of technology and the reduction of computer equipment cost focus on an efficient creation method between developers and users, who tend to focus on the recent small-quantity batch production [3]. Such customization is also being attempted in costumes. Current clothing items for customizing do not create diversity.

Historically, designers used to deliver a one-sided message that has changed in contemporary times, however. This is caused by the democratic desire of customers who have led the way in changing traditional business models [4]. Nonetheless, it is doubtful whether there are participation and creativity of a design for the current clothing design items in customer-participating services. Despite the good customer-participating structure, there is a lack of professional leadership for the service. In 2012, Armstrong and Stojmirovic argued [2] that the point of inducing customer participation is that it must be easy, fun, and less burdensome in terms of price, and it must give excitement to produce creative works under the leadership of an expert. Note, however, that the current t-shirt customizing is insufficient to produce passionate creativity as customers exercise their creativity to the full extent. Such is partly due to the lack of Back End technology. Thus, Design U developed a DTP printing image extracting system that restores the image from the model where the customer designed before. Technically, it prepares a presumption of customer design participation for various clothing items. From the customers' point of view, it is necessary to study the development of a user interface that can be edited in a more complicated content format and a contents format that can communicate with customers. If the designer designs a creative clothing item and suggests lead content, source contents need to be designed for customers to edit. This part can be made by uploading a customer-created image or sharing through the purchase of fee-paid copyright. The design source is sharing content wherein customers can easily and joyfully complete designs. Content sharing with a small amount of copyright is a method created by Richard Stallman who supported the copyright opposition movement and created a concept of free distribution of information. Lessig [5] also supports the flexible copyright method that reuses content information since it directly affects the participating culture.

This study first analyzed the customer editing screen user interface cases for customer design participation in the clothing shopping mall, investigated the pattern copyright cases used for customer editing, and finally developed

a customer-participating editing user interface application for more complicated clothing items.

2. Related Research

Examining the clothing product development process is an essential part in preparing and supplying quality products to a promising market in a timely manner at a reasonable price. The common steps involved in designing clothing products are research, design development, and production [6]. The design development step starts with design sketching and sample production, and most of those who are participating in this process use computer-aided design (CAD) for efficiency and accuracy (Figure 1). The existing CAD systems are often effective in downstream production wherein grading, and marker planning processes take place continuously. Owing to a series of research works carried out by modern engineers and designers, Virtual Sampling, which allows a 2D pattern to be applied to a 3D human model to evaluate its wearability and appearance, has become a standard procedure when designing clothing products [7–10].

Cloth simulations are usually performed to assess the effect of geometrical variation or physical aspects. In most cases, the former draws faster results without considering the physical properties of the cloth being used. This makes it difficult to reproduce the dynamics of the clothes [12]. The latter allows more realistic simulation in understanding the dynamics and provides better accuracy as the cloth material's structural properties will be considered for the simulation. In other words, both the law of dynamics and the law of mechanics are based on discrete dynamics, fluid dynamics, or elasticity theories, all of which determine the cloth behavior and its interaction with external environments [13, 14]. Various methods often categorized as either a continuous physics-based or a discrete physics-based approach have been studied and proposed till now, emphasizing realism or computational efficiency [15]. The former introduces a rigorous, strict representation of a cloth in accordance with the continuum mechanics often adopting either a finite element (FE) or a finite difference model to produce a solution [16–18].

Meanwhile, by using the continuous Lagrange equations to represent the displacements from equilibrium positions, Terozopoulos et al. [16] modeled a surface deformation of a cloth, whereas Eischen et al. [19] employed the nonlinear shell theory and Li and Volkov [20] depicted the image of a cloth immersed in a quasistationary viscous fluid in terms of fluid dynamics. For this, a nonlinear FE method was applied to derive the system equations. This method aimed to produce various types of physical models for computer animation, which are effective in generating the behaviors instead of modeling a certain deformable cloth with high degree of accuracy. This method allowed the qualitative reproduction of similar behaviors without requiring a large number of computations [21]. The instability and high expense are major problems for the continuous physics-based methods, whereas the discrete physics-based methods represent cloth as a grid of particles interacting with each

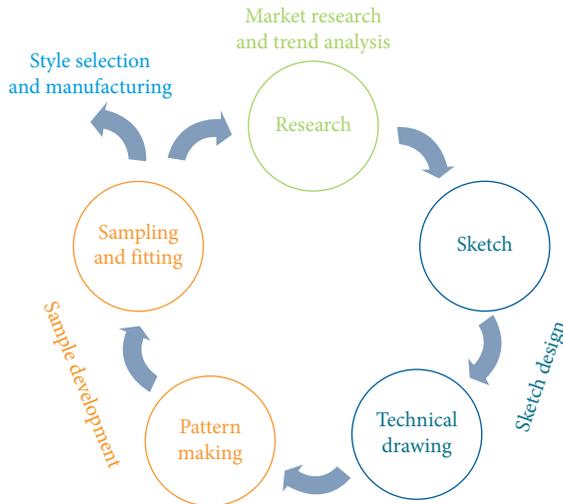


FIGURE 1: Fashion product development cycle [11].

other as well as with the external world following either the force-based Newtonian dynamic laws [22, 23] or the Lagrangian dynamic laws or the energy minimization criteria, all of which are energy-based approaches [21]. Characteristics such as low complexity and simple implementation are the underlying reason for the popularity of the mass-spring system that is also known as the particle system) [15, 24–27]. Nonetheless, an important issue remaining for this system is its accuracy as the physics of a cloth deformation is based on the approximation that is often represented by the mesh topology with a certain discrete physics-based method that influences cloth simulation. A number of meshing and remeshing methods were proposed in the past including Lienhardt [28], Praun and Hoppe [29], and Attene and Falcidieno [30]. Various forms of mesh topologies have been studied by Lu et al. [31], who then validated and proposed an optimal pipeline that is quite adequate for the preparation of the mass-spring model in a scanned garment reconstruction. Meanwhile, some other researchers turned their attention to pattern designing or making. An interesting interactive coevolutionary CAD system for the parametric pattern design was introduced by Chen et al. [32] who produced garment patterns using a neural network along with an immune algorithm.

A fuzzy logic-based optimization of garment pattern design was achieved by Chen et al. [32], whereas Lu et al. [31] proposed an expert knowledge-based approach that is helpful to customized pattern designing. Guo et al. [33] contributed to providing a detailed review of AI applications in the fashion industry. The previous studies suggest that sketch design is still an unexplored field of cloth design and development. Although the use of commercial CAD software such as Adobe Illustrator™ or CorelDraw™ has become a common practice in the sketching process, and their efficiency and effectiveness have been proven, the original idea of a designer starts by creating some sketches, consuming much time and effort. The survey conducted among Hong Kong fashion designers clearly showed that they are continually looking for a user-friendly design support system to

reduce their working hours when designing new clothing [34]. As one of the methods to achieve this, several companies established their own special digital fashion library like SnapFashion™ from which designers can borrow their desired elements to create new designs.

There was a unique development when Mok et al. [35] introduced a customer-oriented design system that allows general customers to create their sketches, adopting some evolutionary computational techniques. Further notable work came from Wan et al. [36] who claimed to have used some shape deformation techniques to create realistic sketches based on standard technical sketches. Fashion illustration and technical sketches are the two main pillars of the design industry (Figure 2). Specifically, fashion illustrations focus on drawing the products that the seller wishes to show and sell to the customers by showing how the products can be arranged and what their uniqueness is. This system translates the technical sketches into fashion sketches without omitting any details. By fitting the same clothing to a different fashion figure that can assume several different positions, the customers could understand the concept and may find the products attractive [11, 37–42].

3. Customized DTP Clothing: Case Analysis

Customized DTP clothing service is drawing much attention, and solutions for copyright issues were analyzed in this study using actual cases. Subjects of the analysis include service methods, user interfaces for the user-edited screens, and pictures used for printing. Among the DTP being serviced, five cases with clear distinctions have been selected for the analysis in order of introduction: My T, Snap T, Design U, Printing Factory, and Adidas. Among these, My T, Snap T, and Printing Factory are mainly selling T-shirts, so they can be regarded as products that originated from IT companies or printing factories instead of being clothing brands. Design U is a service provided by contemporary Korean traditional clothing brand TS, and Adidas is originally one of the clothing brands. In this aspect, the T-shirt business is often run by nonclothing brand firms, and the requirement or level of difficulty of production pertaining to the design of their clothes is quite low. Their interests lie in the files that they need to print. If a DTP service has to be provided for complexly designed clothes like women's clothing, the level of production difficulty will be high, and higher understanding of their design is required.

The DTP mentioned in the case analysis refers to digital textile printing, a method which replaces the conventional dyeing method and saves the time required for cutting patterns so that this method is quite suitable for the modern customized services. The user interface allows convenient communications with the system the user wish to select, aiming to reach the level of communication the customer desires.

Also, from the perspective of recently developing information communication and design technologies, the user interface is an interactive space for the various types of computer-based equipment. Operating the surrounding products in our everyday lives is a normal phenomenon in

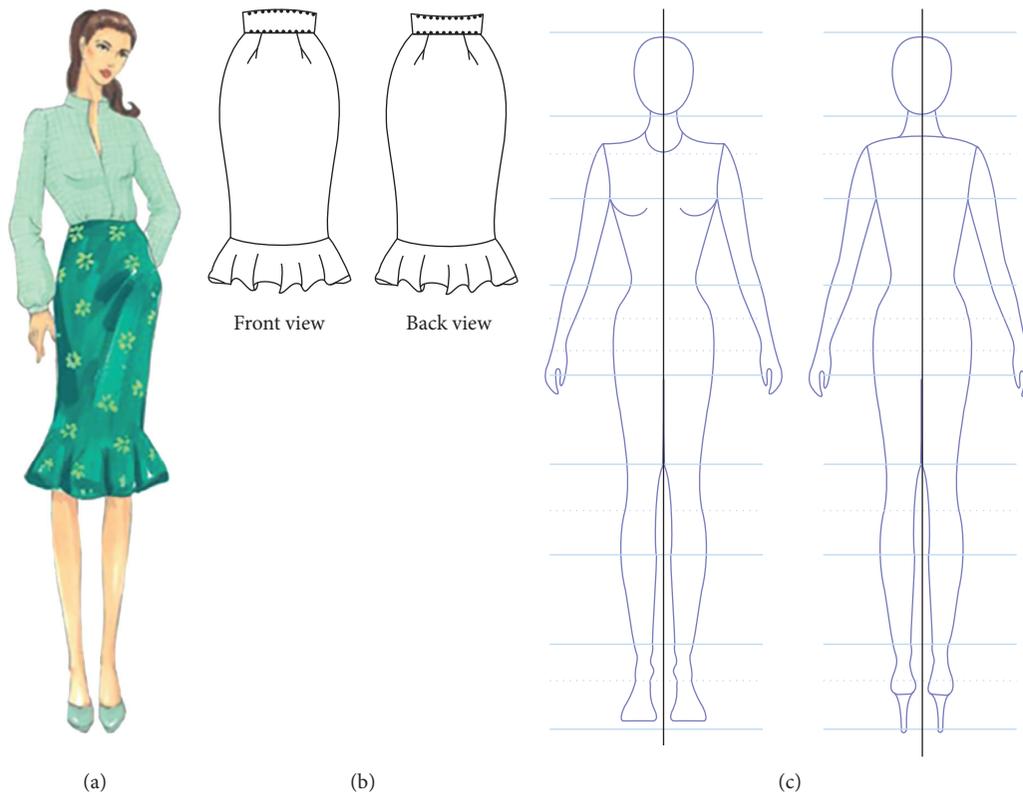


FIGURE 2: Illustration of fashion illustration, technical sketch, and croquis. (a) Fashion illustrations. (b) Technical sketches (flats). (c) Croquis (human figure template).

the environment created between products and user. For instance, the user just needs to use the hardware or the software of the vehicle he/she owns through the embedded user interface even if he/she does not know how it works or what is the principle of it.

A good user interface design makes it easier for the people to operate the products they encounter in their everyday environment. The design includes not only arranging the composition of a computer/similar device's screen or the elements of hardware operation in a convenient way but also includes all the designs of the things the user experiences with the products.

Moreover, since the people who will be providing such service should have a high degree of understanding of the design and form, the work is much different from simply concentrating on the printing files. Although there are many more services utilizing DTP printing, these five services have been selected because they offer customer-participating editing screens (Figures 3–8).

My T (Figure 3) is an app that helps produce personalized T-shirts. It provides an internal service platform that introduces the printable pictures created by artists from various fields and receives a royalty (T-mileage) once they have been sold. The printing house can be arranged, and the optimal printing method will be automatically provided. This service targets consumers who look for T-shirts that are highly individualized but affordable. The customer editing screen is arranged as follows:

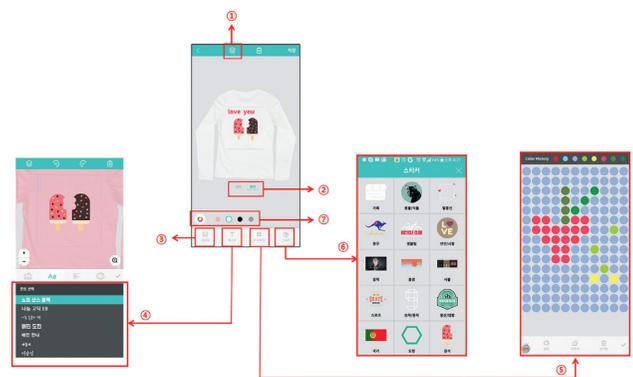


FIGURE 3: My T [11]. ① Layer control. ② Area selection. ③ Import image. ④ Text editing. ⑤ Dart editing. ⑥ Load Sticker. ⑦ Select background color.

- (i) Button ①: *Layer adjustment*. Distinguishes the pictures to be placed in the foreground or background when several pictures overlap.
- (ii) Button ②: *Area selection*. This service can be used for only two areas (front or back), selecting the designated rectangle domains. This shows that the clothing will be produced prior to printing.
- (iii) Button ③: *Image import*. This function usually imports images (jpg, without transparent area) from the customer's mobile phone.

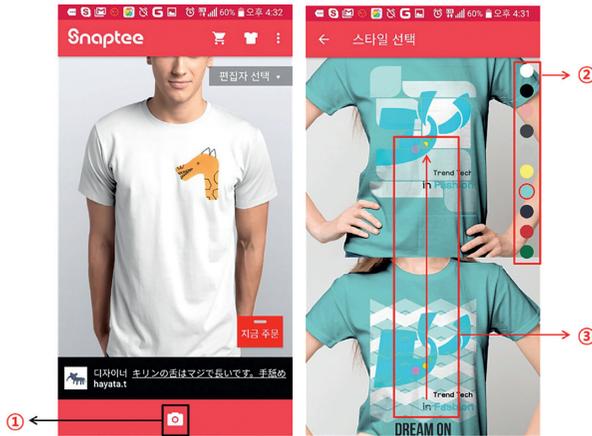


FIGURE 4: Snap T [38]. ① Upload image. ② Select background color. ③ Go down and see various examples.



FIGURE 6: Printing Factory [40]. ① Revert/revive. ② Area selection. ③ Simulation view (2D). ④ Start work. ⑤ Text editing. ⑥ Select background color. ⑦ Image selection. ⑧ Image editing.

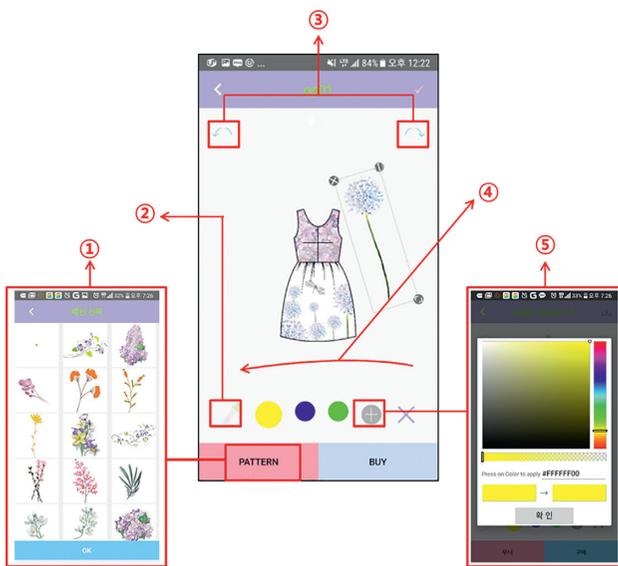


FIGURE 5: Design U [39]. ① Import image. ② Eyedropper. ③ Revert/revive. ④ Swipe: front, back. ⑤ Color selection.

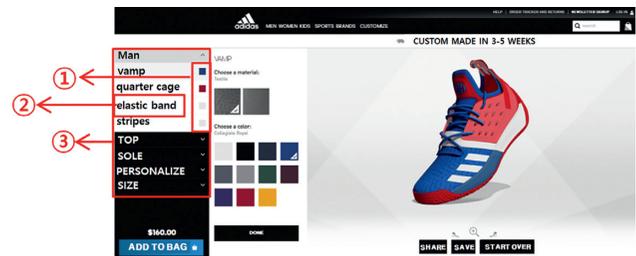


FIGURE 7: My Adidas [41]. ① Color selection. ② Material selection. ③ Area selection.

(iv) Button ④: *Text button*. Font selection options will be displayed once clicked.

(v) Button ⑤: *Dart editing button*. Once the button is clicked, the screen on which the dart can be edited will appear.

The size can be adjusted, and certain figure or letter(s) can be imprinted by coloring each dart. In this area, some game factors can be reflected as well.

(vi) Button ⑥: *The sticker is imported*. Sticker refers to any of the artworks provided by various artists, and it is mainly a PNG file with a transparent domain. Images and stickers have apparently been distinguished in this service based on the source of the image (i.e., user’s mobile phone, the platform itself, etc.) or file format (i.e., transparent or nontransparent).

(vii) Button ⑦: *Background color selection*. Each T-shirt has its own fixed background color. In the case above, however, the color selection option is quite limited as there are only four background colors.

Meanwhile, Snap T (Figure 4) is a service initiated by one of the Hong Kong companies, aiming to reflect collective amusement to Instagram or T-shirt design. Its editing screen has been a little more simplified, but the social function has mainly been strengthened. With this service, the customer can disclose or present his/her design; if another customer purchases the design, a 10% sales royalty will be collected. The editing screen is straightforward as shown in Figure 4: one’s image is uploaded by clicking Button ①; for the uploaded image, various types of draft designs will be displayed when the customer drags the mobile phone screen downward. This is made possible by marking the image with different geometric shapes and inserting varied transparency levels in advance, giving a dream-like design effect. The background can be selected with Button ②; since it will be printed on the preproduced T-shirt, only the existing colors can be selected.

Design U (Figure 5) is an app developed to load clothing contents more intricate than T-shirts, aiming for better communications between the designers preparing a certain designer brand and customers. At the time of development, designers of contemporary Korean traditional clothing had loaded these dress items reflecting modern trends. This service adopts the “print first, produce later” method. Around 100 arbitrary images matching the items have been

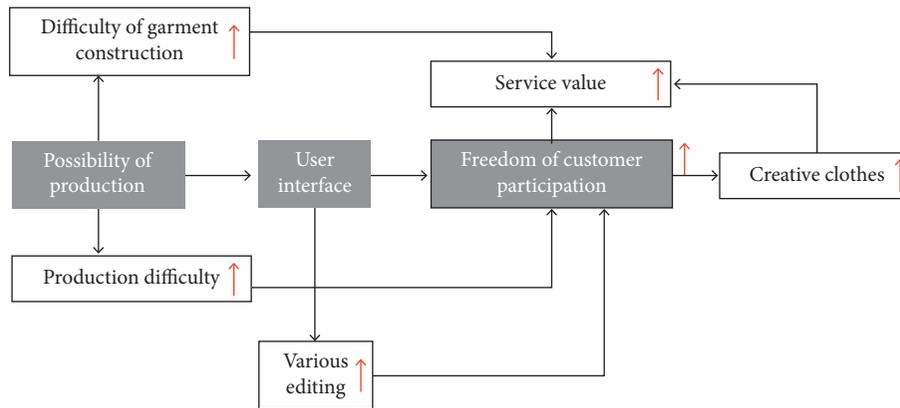


FIGURE 8: Factors affecting the user interface.

uploaded to the app. The customer needs to download the desired design or pattern to import it when editing. The customer editing screen is shown in Figure 5. ① is the image import button with which the customer can apply the downloaded image. ② shows the eyedropper function that elaborately leads the customer when he/she chooses a color. It could be difficult for the customer to find a sensible color in a short period of time among the colors suggested in ⑤. Nevertheless, the customer can have the desired color by importing the image file that includes the same color from the pattern samples and eye-dropping in the desired area. This function can be used to pick out a subtle color shade from a beautiful pattern. There are three circles on the right side of button ⑤ acting as a palette. The customer can start designing by preparing three color combinations in advance with an eyedropper or a color selection function. When the eyedropper in ② has been activated, the palette changes to the color designated by the eyedropper; when this is inactivated, however, it is possible to paint the clothing with the desired color by putting one of the colors in the palette in any area. The function that reverts or restores the working content can be performed with ③. With ④, the front or back of the dress is displayed when the customer moves the screen to the left or right so that the back side can be painted after the front side by shifting the screen to the side.

Although this service provides only two selections (front and back), it is possible to apply the service to all the areas. In other words, the edited arbitrary pattern will be adjusted following the outline of the dress drawing by clicking the check button on the upper right side.

Despite its designing convenience and better looks (prettiness), this method consumes much time since the images have to be put on the clothing sample one by one. In the picture, even though the pattern seems to be pieced together between the different domains, it is little challenging to make the image a perfectly joined and extended one because these domains are actually those of different pieces of fabrics. Therefore, it will be helpful if the work can be performed separately (subdivide) for each domain.

Printing Factory is a Korean T-shirt printing website who undertakes the editing/printing of costumes for fashion shows or K-pop idols by partnering with the relevant

companies. On their website, they offer a customer-editable screen as shown in Figure 6. They are currently planning to extend their service to a DIY business wherein customers can make their clothes. Although their basic business strategy is to cooperate with as many clothing companies as possible continuously, their business is still limited to T-shirts and sportswear. ① is the button for undoing or restoring the performed work. ② selects the area to be designed. For T-shirts, four areas (i.e., front, back, right, and left sleeves) are usually available. In addition, each area is not defined by the rectangular box so that all areas of the cloth can be designed freely. ③ is used to check the simulation.

As the present screen, it is replayed in 2D mode; the difference is that the front (back) side of the front (back) sleeve will be shown just next to the front (back) face after the design has been completed so that the customer will be able to see how the T-shirt will look like after the sleeves have been attached. This is because the designed sleeve area is not recognized as just an ordinary arm part of the clothing but as an unfolded design (i.e., coloring the sleeve pattern). The merit of direct designing on a sleeve pattern is that the customer's design details can be reproduced without any data loss when printing them on a cloth. This offers many conveniences in the production process while reducing customer complaints. The customer may modify the design after seeing the simulation, but this process can be a little inconvenient as well. When the operation start button ④ is pressed, the three major functions such as Image, Text, and Color will appear. Button ⑤ is for text editing that supports font selection, contents input, sorting selection, color coating input, shadow code input, horizontal/vertical position adjustment, and blur effect.

The designing freedom of text itself has been enhanced. For the T-shirts, the text function has been improved as texts are often used more widely compared to women's or men's wear. Button ⑥ is for background color selection. The desired color can be selected from the RGB color scheme. Button ⑦ selects an image. The images are provided by theme, but one's images can be uploaded. Button ⑧ is used for editing the image by enlarging, reducing, and rotating it. The image setting window is provided separately for image cutting, image filter selection, or transparency adjustment.

Figure 7 shows My Adidas, a personalized shoe production service by Adidas. According to Adidas, they have commercialized this online service reflecting the trend among consumers who pursue their unique fashion and personality. This pertains to the highest level of Maslow's Hierarchy of Needs, so this personalized service is a prospective service targeting modern people who can feed and house themselves. When designing the Adidas shoes, the colors can be selected from the fixed set of colors and applied within the fixed areas. The materials can be selected from the given set of materials as well. When selecting an area, the Screen View rotates the shoe 360° degrees to show each area, whereas the applied colors can be viewed through 3D shoe design every time they have been applied. Adidas adopts many fixed functions, especially for the color selection. This is to help the customers match colors more professionally. Although customer freedom may be reduced, limiting customer choice by offering special color combinations would assist them better in terms of designing. The elements affecting user interface in the DTP customer-participating type clothing platform are shown in Figure 8. In the first place, production possibility is the base of the user interface. If possible items are different, the items and functions of the user interface vary. The user interface determined by such production possibility determines the degree of freedom of customer participation. If freedom of customer participation is high, more creative work can be produced. Production possibility, user interface, and customer's degree of freedom form the relation below. Firstly, if production possibility is high, complicated clothes have increased service value. If the production process using DTP is difficult, the customer's degree of freedom increases. In terms of user interface, if the diversity of screen function increases, the degree of freedom increases. If the customer's degree of freedom increases, the creativity of output increases together with the service value. In other words, if the customer's degree of freedom increases, production becomes complicated, and user interface needs a complicated function, but it increases service value. In contrast, the lead index restricts customer participation. It is intended to apply some restriction for professional and refined outputs by the customer with the help of professionals. It is expressed in the form of restriction of the customer selection area and fixing of the location. In the case analysis, ID has a high lead index. Customer participation and leading index apply mutual restriction. While the customer's degree of freedom is intended to restrict creative works, the lead index enables customers lacking creativity to express designs professionally by leading the method of restriction. Moreover, even if there is high freedom of customer participation, it does not have a positive effect on customer participation.

Customer participation may vary by generation based on the experience of Internet device. Generally, younger generations with high level of freedom have a positive reaction to customer participation. Therefore, it is ideal to keep a suitable level of service and price by maintaining the appropriate degree of freedom of customer participation concerning the market. There is a high level of freedom of customer participation and service, so it is important to

create value for various customer groups. If easy clothing items are focused on, it creates severe competition within the item and decreases the value. Therefore, it is necessary to apply various clothing items to the DTP customer Half Design participation service. The complexity of cloth, diversity of customer participation, degree of freedom, and evaluation of user interface are listed in Figure 9. The diversity of user interface is related to customer participation's degree of freedom. Thus, the lead index is not evaluated in Figure 9. Likewise, 5 cases with clear differences including App and Web are analyzed for the evaluation of the index above to figure out the service situation in this industry. There are three analysis cases: App for My T (Korean), Snap T (China), and Design U. As for the web, Printing Factory and Adidas are analyzed.

4. Analysis Result

4.1. Result of Case Analysis for the DTP Clothing Service User Interface. The analysis table below is expected to be used as a reference to evaluate the platform function of DTP customizing clothing service editing screen and production method. In the evaluation items, factors affecting customer participation's degree of freedom are marked in bold. If the computer and mobile environments support customer editing screen in the development environment, the customer's degree of freedom is high. In the design areas, when it moves down, the degree of freedom is high, and the degree of freedom of Whole is highest. In analyzing the user interface editing and customer's degree of freedom, if there is no color restriction, the degree of freedom is high. The degree of freedom for Free Color is highest. In the selectable color, the degree of freedom is high if more than one color can be selected. If there is a spoilt, it has a high lead index. In the image, seven factors including Upload Customer's image, Provide image, zoom/Rotating image, Image repeat, Image game, and Text are considered. If there are many factors, it has a high degree of freedom. Next is a service on production complexity. It is related to service value, which includes the following. If the difficulty of Item goes down, the production complexity of clothing will increase. If Design Area moves down, the complexity of DTP editing will increase. If the number of pattern pieces for design increases, the production complexity in DTP printing editing increases. If there are many pattern pieces overall, production complexity will increase. Regarding the production method, preprinting and postproduction are closer to a traditional customizing production process than preproduction and postprinting. Compared with preproduction and postprinting, the difficulty of production is three times higher.

4.2. Result of Analyzing Creation-Sharing Contents Copyright Cases. Figure 10 shows a chart of the analysis result of the copyright management system with which the creations belonging to each brand can be shared. The shareable creations are divided mainly into art graphics that can be used for designing purposes and finished clothes. In the chart, the

DTP Clothing Half Design Online Platform			1. My T	2. Snap T	3. Design U	4. Printing Factory	5. Adidas	
Environment		App/Web	App	App	App	Web	Web	
		Computer				○	○	
		Mobile	○	○	○		○	
Difficulty of garment Construction	Item	T-Shirt, Eco Bag	○	○		○		
		Sports / Spandex					○	
		Fashion Clothes			○			
	Design Area	Limited Area	○	○			○	
		Free Area						
		Whole			○	○		
	Designable Pattern Quantity	1 sheet		○				
		Within 1/3 Quantity						
		Within 1/2 Quantity	○					
		All			○	○	○	
Whole Pattern Quantity		4	4	7~12	4~6	4		
Editing function	Back ground Color	Color Select	No selection					
			Choose from provided colors	○	○			○
			Free Color			○	○	
		Color Spuit			○			
	Selectable degrees of freedom	Choose only one color	○	○				
		Choice of 1, or more colors			○	○	○	
	Image	Upload Customer's image		○	○	○	○	×
		Provide Image		○	×	○	○	×
		Zoom, Rotating image		○	×	○	○	×
		Image Effect		○	○	×	×	×
		Image repeat		×	×	×	×	×
		Image Game		○	×	×	×	×
		Text		○	×	○	○	×
Production method		Production -> Printing	○	○			○	
		Printing ->Production			○	○		

FIGURE 9: Case analysis of DTP clothing half design platform.

five brands and an additional brand, Real Fabric (Korean fabric printing web service), are included. This service has a robust copyright management system. Even though My T claims to be paying regular returns to the artists when using their designs (pictures), the process is not explicitly shown on their service screen. Meanwhile, Snap T offers a simple editing function, but they have a relatively proper social role

for sharing creations and returns. They publicize the clothes design results, and the profits are shared in case of any purchases made. They do not provide art graphics, so the designers should upload their creations from their own devices. In this aspect, it is an art graphic sharing service. Real Fabric has a selection box showing the art graphics presented by the expert artists and which can be applied to

Copyright management		1. My T	2. Snap T	3. Design U	4. Printing Factory	5. Adidas	6. Real Fabric
Revenue share system	Graphic Design	○	○				○
	Clothing Design		○				
Copyright Flows			○				○
Revenue share automation			○				

FIGURE 10: Analysis of copyright management system by brand.

the fabric, sharing the profit with the artist. The profit-sharing follows their reimbursement rules. Since the copyright holders can be checked in the services provided by Snap T and Real Fabric, their services can be said to offer a better environment wherein the graphic artists can openly create artworks and socialize with their coartists and customers.

5. Development Concept Map

5.1. *Development of DTP Clothing Service User Interface Web (Excluding Knitwear).* Figure 11 shows the designer-participating customer user interface development DB UML for the application of various clothing items. It mainly consists of customer information, image upload information, other information related to image editing, and color information.

5.2. *Creative Work Copyright Application Planning for DTP Half Design Platform.* Figure 12 is a plan for the method of uploading works to use in the platform. If customer A uploads work, customer B buys the work at a minimum copyright fee and saves it to his/her own web file in the platform. While the internal contents of the web file can be applicable to half design participation, the original copy cannot be downloaded. As the method of paid pattern application by customer B in the design participation, when a customer performs design, there is an open image column on the editing screen. Here, free or fee-paid creative works appear in an image for selection. If the desired works are selected, it can be reused for the customer’s work immediately by editing. Such method is quite reasonable as the copyright holder is explicitly specified and payments are made safely when the designs are downloaded to the customer folder.

6. Development Result

6.1. *Result of Development of Customer Editing User Interface Applicable to Various Clothes.* Currently, the case studies show that customized/personalized DTP clothing services are limited to T-shirts and sportswear. Although Design U once attempted to do women’s wear lines, much time was required to reproduce the clothing sample after reflecting the customer design information since the area distinctions were not made clear, but the method used by Printing Factory increases the customer participation level by providing the list of areas for each clothing piece. Note, however, that they use one view for each clothing area. This method does not

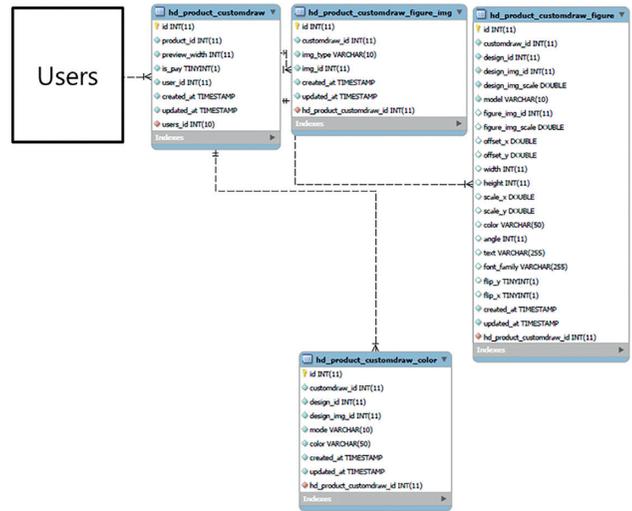


FIGURE 11: User interface web of UML of DTP clothing service.

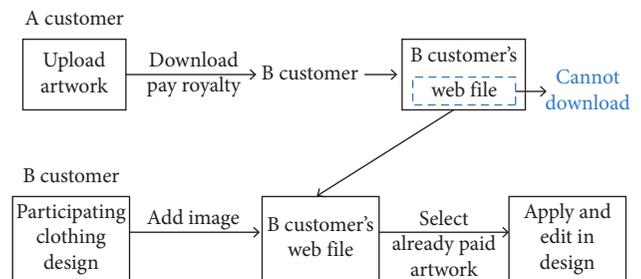


FIGURE 12: Copyrighted artwork in half design platform.



FIGURE 13: Customer design participation screen of design user interface.

allow the customer to see the overall effect, so they added the simulation view function separately. Considering this, for more complicatedly designed clothing, the Design U’s method which connects all the areas in a single view would be more suitable. Therefore, for the customer view, this method has been adopted as in ② of Figure 13. ① in Figure 14 is an Area, Selection List. When an arbitrary area has been selected, the area becomes semitransparent

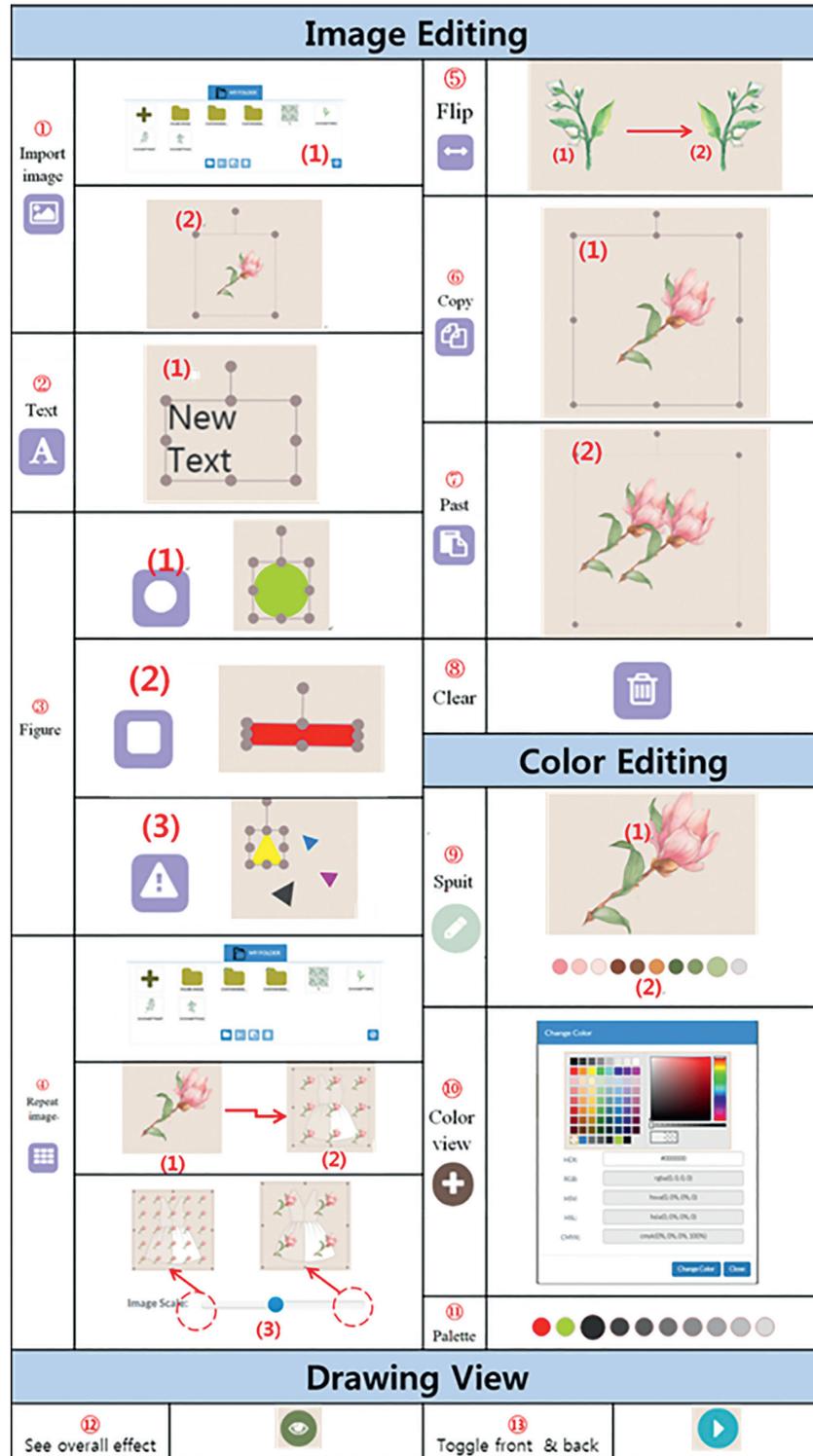


FIGURE 14: Editing function.

(Figure 14) so that the customer can recognize where the selected area is in the picture. When an area has been selected, the entire work output stays in that area. As a result, by pressing the eye-shaped button on the upper left after completing all the areas, the screen will show the result as in Figure 13. In other words, all the images designed in each

area will be masked to show them neatly so that the customer can easily image his/her finished work. This method is being applied to complicated wear with many selection options with which the customer can select each area from the complexly set up areas and represent the designs on a single screen.



FIGURE 15: Customer design participation screen of design user interface.

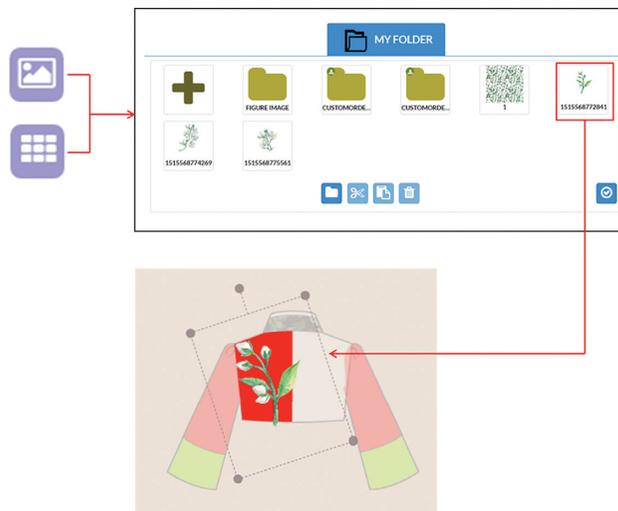


FIGURE 16: Development screen.

The customer design participation screen of the design user interface is divided into five (Figure 13): ① selection area, ② drawing area, ③ image editing, ④ color editing, and ⑤ save. It is divided into selection, coloring, image editing, color editing, and finally saving. Also Figure 15 shows customer design participation screen of design user interface.

As shown in Figure 14, the customer editing screen consists of ① Area Selection List, ② Drawing Area, ③ Image Editing Button, ④ Color Editing Button, and ⑤ Storing and Sharing Buttons. The details of each button are summarized in Figure 14. Button ① is an image import button. Once the button is pressed, the screen shifts to the customer folder where the paid or free pictures have been downloaded. The desired picture can be selected here. (1) of ① allows the user to import the files from his/her mobile phone or PC. (2) of ① shows the picture selected from the folder that has been generated on the screen. (3) is an editing tool generated for each design element, performing five functions such as selection, enlargement, reduction, rotation, and shift. The editing

tool automatically appears when a single item is selected from the elements such as text, geometric pattern, and pattern repetition. Thus, the pattern selected in (1) can be freely edited with editing tool (3) and displayed in (2) by designating the desired location or size.

② is a text button used to input text by clicking the button and positioning it anywhere on the screen. The editing tool is automatically generated when the text has been entered as an element, performing a function similar to image editing. (1) of ② shows the fonts. After the text is selected, one of the available fonts will be selected to decorate the text. ③ shows patterns from which the simplest ones such as (1) circle, (2) rectangle, and (3) triangle were selected. The pattern can be represented freely by drawing it as one big shape or several small ones. By clicking the desired color, the pattern color can be changed. ④ is an image repeating button, and its image import method is the same as ①; after importing an image, however, an arbitrary square should be drawn together with the click-drag function. By doing so, the image selected from (1) repeatedly appears in the square as in

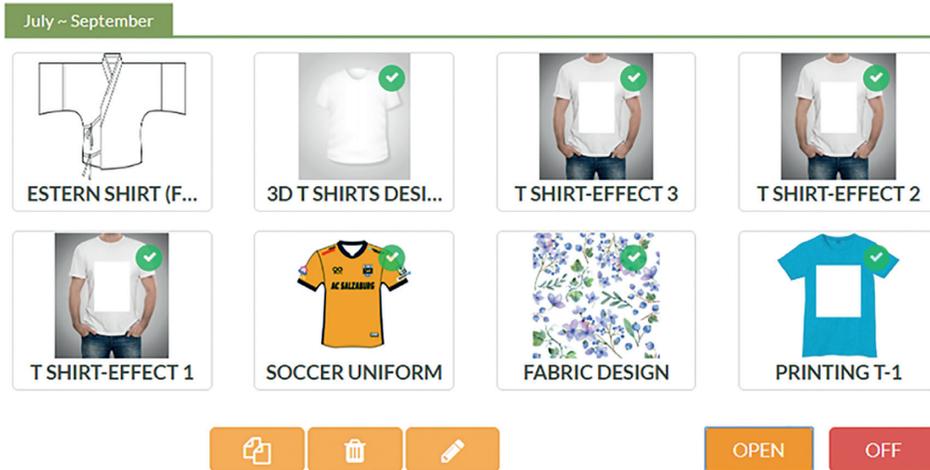


FIGURE 17: The service/product management screens on which the products to be tested are being uploaded.

Example: My T	
Features	<ol style="list-style-type: none"> 1. Service Area: 2 2. Service Area Form: Square 3. Clothing Background Color: Assignable
Simulation Graphic	Picture of item laid flat
My T	
Representation on the platform	<div style="display: flex; justify-content: space-around; margin-top: 5px;"> Before application Adaptation process After application </div>

FIGURE 18: The result obtained from restricting an image within the square.

(2). If one wishes to adjust the size, it will become more substantial when the scale button is moved to the right side. The editing tool is also automatically generated for the finished pattern element so that the direction of pattern repetition can be adjusted by the rotation function of the tool. ⑤ is the flip function allowing only horizontal reversal. In other words, when flip has been performed by clicking one of the elements (e.g., image, repeating image, or text), (1) will disappear and (2) will appear.

⑥ and ⑦ are the copy and paste buttons, respectively. By clicking the ⑥ copy button after clicking the image (1) of ⑥, button ⑦ will be activated. The copied image will be pasted when this button is pressed (image (2) of ⑦).

It is possible to develop the copy function without the attach button; in such case, however, there could be some difficulties when one attempts copying in different areas separately. For this reason, both buttons have been developed as separate ones. ⑧ is a Clear button with which the

Example: Snap T													
Features	1. Service Area: 1 2. Service Area Form: Square x graphic 3. Clothing Background Color: Assignable												
Simulation Graphic	Picture of a person wearing the item												
Snap T	 <p style="text-align: center;">Expressing different effects with the same image</p>												
Result of Application on Platform	<table border="1" style="width: 100%; text-align: center;"> <tr> <td style="width: 15%;">Effect 1</td> <td>  Before Application </td> <td>  Adaptation process </td> <td>  After application </td> </tr> <tr> <td>Effect 2</td> <td>  Before Application </td> <td>  Adaptation process </td> <td>  After application </td> </tr> <tr> <td>Effect 3</td> <td>  Before Application </td> <td>  Adaptation process </td> <td>  After application </td> </tr> </table>	Effect 1	 Before Application	 Adaptation process	 After application	Effect 2	 Before Application	 Adaptation process	 After application	Effect 3	 Before Application	 Adaptation process	 After application
Effect 1	 Before Application	 Adaptation process	 After application										
Effect 2	 Before Application	 Adaptation process	 After application										
Effect 3	 Before Application	 Adaptation process	 After application										

FIGURE 19: The results obtained when the effects have been brought into the restricted area.

design element selected arbitrarily can be deleted. ⑨ is an eyedropper. With this button, the arbitrary elements and all the colors in the image can be put into the palette ⑩. For instance, when one wishes to use a color combination consisting of Pink, Burgundy, and Green for image (1) of ⑨, one needs to activate the eyedropper; after clicking the palette, click the color to be applied on the pattern. Then, as shown by (2), the colors are classified on the palette so that a natural arrangement of colors will be created. In addition, if one wishes to use the color collected with the eyedropper for the background, one needs to click the eyedropper

button again to deactivate the eyedropper. ⑪ is the color view with which one can freely select new colors or input a color code in the HEX box to get the desired color. ⑫ is the palette. There are a total of ten palettes, and the color can be changed in each palette. Customers can have a user-friendly palette by clicking each palette and entering the colors they often use or prefer. ⑬ performs the See Overall Effect. When this button is clicked, the screen in Figure 13 where each area has been selected in advance will be changed into the screen in Figure 15, showing the image applied with all the patterns and colors.

Example: Design U App	
Features	1. Service Area: Entire Area 2. Service Area Form: 2D Form 3. Clothing Background Color: Free Choice
Simulation Graphic	Flat design with a black line (item)
Design U App	
Result of Application on Platform	<div style="display: flex; flex-direction: column; align-items: center;"> <div style="display: flex; justify-content: space-between; width: 100%;"> Before Application  </div> <div style="display: flex; justify-content: space-between; width: 100%; margin-top: 10px;"> Adaptation Process  </div> <div style="display: flex; justify-content: space-between; width: 100%; margin-top: 10px;"> After Application  </div> </div>

FIGURE 20: The result obtained when it was possible to carry out the design in the entire area freely.

The customer editing screen consists of the Area Selection List, Drawing Area, Image Editing Button, Color Editing Button, and Store/Share Buttons.

6.2. *Copyright Image-Applied Screen.* In Figure 16, when clicking the image open and pattern buttons, a customer web folder appears. In the folder, images that customers want are gathered, but they are not available for download. The pattern is free or copyrighted. Moreover, in the copyrighted patten, it can apply the design to the customer web folder.

7. Test Evaluation

To test the efficiency of the platform developed, the methods used for the services have been applied to the

platform to confirm their validity. As examples, the brands such as My T, Snap T, Design U, Real Fabric, and Ninetyplus have been selected as each of them has a distinctive individuality in their respective service function. All of their service approaches were implemented on the platform. Figure 17 shows the service/product management screens on which the products to be tested are being uploaded.

The main feature of My T is that the image is restricted within the squared area (Figure 18). Although such a service does not allow an elevated level of customer freedom, it offers a much more convenient and cost-saving production method, leading to a better marketability because of its direct effect on the customer preference in low-cost products. Through the test, it was verified that the method of defining

Example: My Adidas	
Features	1. Service Area: Entire Area 2. Service Area Form: 3D Panel & Logo 3. Clothing Background Color: Selection from available colors
Simulation Graphic	3D object
My Adidas	
Result of Application on Platform	<div style="text-align: center; background-color: black; color: white; padding: 2px;">Simulation Part</div> 

FIGURE 21: The result obtained when represented in 3D.

service area in a single/multiple squared areas is viable on the platform developed in this study.

Snap T shows various types of effects within the squared area (Figure 19). This is to allow customers to create an effect like a professional artist. Such an effect can be achieved with a png file which determines the transparency, color, and the shape of the area. Through above test, it was verified that the service which allows customers to create an artistic effect just by giving various types of effects on a single photo is viable on the platform developed.

The App version of Design U allows customers to design in every corner of the clothing (Figure 20). This is also a service method in the table presented in Figure 19, which embraces a high degree of customer freedom. Through the test, it was possible to verify that such a service approach is also viable on the platform developed.

The My Adidas product in Figure 21 was represented in 3D where customers cannot insert any images. However, it is possible to set the colors automatically for the logo or the geometric patterns in the system such that the customers

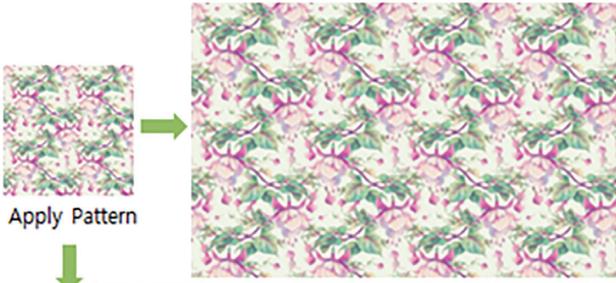
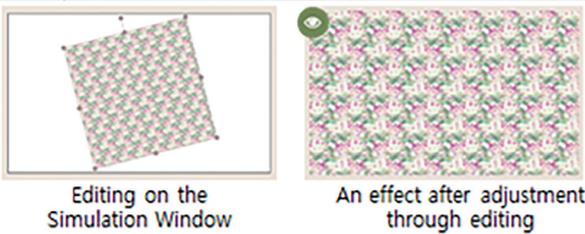
Example: Real Fabric	
Features	1. Service Area: 1 2. Service area Form: Square 3. Background Color: Cannot be assigned
Simulation Graphic	Square
Real Fabric	
Result of Application on Platform	

FIGURE 22: The result obtained from the application on a fabric.

without expert knowledge in designing can also achieve professional-like color matching by making selections from the colors or color combinations prepared by experts. To perform a similar service, a 3D image was used for a T-shirt while making it possible to select the logo color separately on the platform. Although color selection has not been completely limited, the customers will be able to select some professional color combinations based on the color reference given. That is, they can pick their desired colors with a pipette function from the color palette image while using the editing feature. The test result revealed that even though the 3D-based simulation did not produce detailed images, it was verified that such a service can be achieved on the platform.

Figure 22 shows the result obtained from the application on a fabric. The service by Real Fabric can only be applied to fabrics and its simulation screen has a square shape. The information related to the proportion of a pattern based on its length and width can be checked in advance and it is possible to make changes for a single-pattern element as many times as the customer wants as well as the adjust of sizes. The same fabric used for the Real Fabric's simulation was used on the platform for testing. As a result, it was not possible to change the length or the width of the fabric instantaneously such that this problem was solved by selecting a fabric with the desired length and width from the product list in advance. Also, editing was freely performed on the pattern by using the functions (i.e., Enlarge/Reduce, Shift, Rotate functions) having a "Pattern Element Refit"

function. Also, compared to the Real Fabric service where only one pattern can be applied at a time, the proposed platform allows to add several patterns on top of another, making it a superior platform as a customizing service for fabrics.

Figure 23 shows the result obtained from performing a simulation on the existing customized service offering no simulation function. Finally, Ninetyplus (90+) is a company specializing in customized soccer uniforms. This company does not have any simulation systems and adopts an ordering system which requires a customer to download an Excel style order form from the company website and upload it after filling in the specifics. Then, they show the customer their draft design repeatedly for three times. Such a relatively complicated procedure has to be taken as they lack the simulation system. An uploading test was performed on the platform with product of this company by using the simulation technique and the result showed that all the problems can be solved at the same time. One major feature of the color guide used by Ninetyplus is that only a couple or triple of colors can be used. This is to prevent the customers from creating a rainbow-colored design by coloring every corner as they please. By contrast, the proposed platform allows all the predesignated areas will be painted with the same color simultaneously once a single arbitrary area has been painted. This validates that the service similar to one that is being offered by Ninetyplus (i.e., developing a design with just two or small color distribution) can also be achieved on the platform developed in this study.

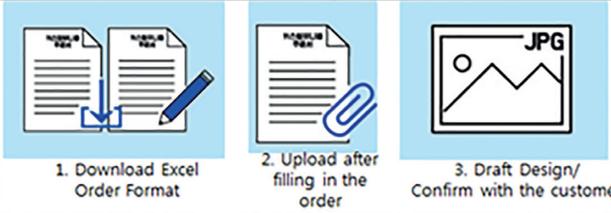
Example: Ninety plus (Soccer uniform customizing service)	
Features	<ol style="list-style-type: none"> 1. Service Area: Entire Area 2. Service Area Form: 2D Panel 3. Background Color: Color-Matching Guide
Simulation Graphic	Nil/Receiving orders with an Excel format
90+	 <ol style="list-style-type: none"> 1. Download Excel Order Format 2. Upload after filling in the order 3. Draft Design/Confirm with the customer
Result of Application on Platform	 <p>Right Arm</p> <p>Adaptation process</p> <p>Right Arm</p> <p>After applying to each area</p>

FIGURE 23: The result obtained from performing a simulation.

8. Conclusion

In this study, research was conducted on the instances of IT-integrated services in the clothing industry to develop an upgraded service version. The case analysis involved the user interface for customer editing and the management system including copyrights. Based on the research, a DTP half design service website with an integrated function to which various methods can be applied has been developed.

Thus, in this study, the DTP clothing half design service was developed. The existing DTP clothing half design service mainly focused on t-shirts. It is related to the production method, and the editing screen user interface by production method is associated with the degree of freedom of customer. In this study, with the premise that there is no DTP design participation service on trendy fashion items other than t-shirts, the user interface to express trendy fashion is developed. Based on Figure 7, trendy fashion includes inner fabric, which has a high level

of difficulty in terms of the composition of cloth. The developed user interface also supports the overall design of pattern, free color, and various editing functions. Moreover, the current user interface enables designing t-shirts and fabric. The user interface covers the preprinting and postproduction method. With this development, digital printing is expected to be applied to women’s clothing with various complexities as well as t-shirts so that customers can express their creativity and personality and it fills the gap with the market.

Data Availability

The data used to support the findings of this study are available from the corresponding author upon request.

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

Acknowledgments

This work was supported by the research fund of Hanyang University (HY-2014).

References

- [1] H. Yoon and J. J. Lee, "A study on ubiquitous fashionable computer design using modules and wear net: with a focus on a detachable modular system," *Journal of Korean Fashion Designers*, vol. 9, no. 1, pp. 1–17, 2009, in Korean.
- [2] H. Armstrong and Z. Stojmirovic, *Participatory Design: The Design Created by Users and Designers*, E. A. Choi, Ed., Beads and Beads, USA, 2012.
- [3] C. Anderson, "In the next industrial revolution, atoms are the new bits," *Wired Magazine*, vol. 1, no. 25, 2010.
- [4] H. Jenkins, *Convergence Culture: Where Old and New Media Collide*, NYU Press, New York, NY, USA, 2006.
- [5] L. Lessig, *The Future of Ideas: The Fate of the Commons in a Connected World*, Vintage, New York, NY, USA, 2002.
- [6] P. Sinha, "The mechanics of fashion," in *Fashion Marketing: Contemporary Issues*, T. Hines and M. Bruce, Eds., pp. 165–189, Butterworth-Heinemann, Oxford, UK, 2001.
- [7] B. K. Hinds and J. McCartney, "Interactive garment design," *The Visual Computer*, vol. 6, no. 2, pp. 53–61, 1990.
- [8] H. Okabe, H. Imaoka, T. Tomiha, and H. Niwaya, "Three dimensional apparel CAD system," *ACM SIGGRAPH Computer Graphics*, vol. 26, no. 2, pp. 105–110, 1992.
- [9] Z. G. Luo and M. M. F. Yuen, "Reactive 2D/3D garment pattern design modification," *Computer-Aided Design*, vol. 37, no. 6, pp. 623–630, 2005.
- [10] F. Durupynar and U. Gudukbay, "A virtual garment design and simulation system," in *Proceeding of 11th International Conference Information Visualization*, pp. 862–870, Zurich, Switzerland, July 2007.
- [11] My T, August 2018, <https://play.google.com/store/apps/details?id=co.snaptee.android>.
- [12] P. Volino, F. Cordier, and N. M. Thalmann, "From early virtual garment simulation to interactive fashion design," *Computer-Aided Design*, vol. 37, no. 6, pp. 593–608, 2004.
- [13] T. J. Kang and S. M. Kim, "Development of three-dimensional apparel CAD system: Part II: prediction of garment drape shape," *International Journal of Clothing Science and Technology*, vol. 12, no. 1, pp. 26–38, 2000.
- [14] N. Metaaphanon and P. Kanongchaiyos, "Real-time cloth simulation for garment CAD," in *Proceedings of the 3rd International Conference on Computer Graphics and Interactive Techniques in Australasia and South East Asia*, pp. 83–89, Dunedin, New Zealand, 2005.
- [15] F. Han and G. K. Stylios, "3D modelling, simulation and visualisation techniques for drape textiles and garments," in *Modelling and Predicting Textile Behaviour*, X. Chen, Ed., vol. 94, pp. 388–421, Woodhead Publishing Ltd., Cambridge, UK, 2009.
- [16] D. Terozopoulos, J. Platt, A. Barr, and K. Fleischer, "Elastically deformable models," *ACM SIGGRAPH Computer Graphics*, vol. 21, no. 4, pp. 205–214, 1987.
- [17] P. Volino and N. Magnenat-Thalmann, "Versatile and efficient techniques for simulating cloth and other deformable objects," in *Proceedings of Annual Conference Series on Computer Graphics, SIGGRAPH*, pp. 137–144, Los Angeles, CA, USA, August 1995.
- [18] K. Y. Sze and X. H. Liu, "Fabric drape simulation by solid-shell finite element method," *Finite Elements in Analysis and Design*, vol. 43, no. 11–12, pp. 819–838, 2007.
- [19] J. W. Eischen, S. Deng, and T. G. Clapp, "Finite element modeling and control of flexible fabric parts," *IEEE Computer Graphics and Applications*, vol. 16, no. 5, pp. 71–80, 1996.
- [20] L. Li and V. Volkov, "Cloth animation with adaptively refined meshes," in *Proceedings of the Twenty-Eighth Australasian Conference on Computer Science*, pp. 107–113, Newcastle, NSW, Australia, January 2005.
- [21] D. H. House and D. E. Breen, *Cloth Modeling and Animation*, A K Peters Ltd., Natick, MA, USA, 2000.
- [22] S. Petrak, D. Rogale, and V. Mandekic-Botteri, "Systematic representation and application of a 3D computer-aided garment construction method," *International Journal of Clothing Science and Technology*, vol. 18, no. 3, pp. 188–199, 2006.
- [23] M. Hauth and O. Eitzmuss, "A high performance solver for the animation of deformable objects using advanced numerical methods," *Computer Graphics Forum*, vol. 20, no. 3, pp. 319–328, 2001.
- [24] D. E. Breen, D. H. House, and M. J. Wozny, "Predicting the drape of woven cloth using interacting particles," in *Proceedings of the 21st Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH'94)*, pp. 365–372, ACM Press/ACM SIGGRAPH, Orlando, FL, USA, July 1994.
- [25] X. Provot, "Deformation constraints in a mass-spring model to describe rigid cloth behaviour," in *Proceedings of Graphics Interface'95*, pp. 147–154, Canadian Information Processing Society, Québec, QC, Canada, May 1995.
- [26] D. H. House, R. W. DeVaul, and D. E. Breen, "Towards simulating cloth dynamics using interacting particles," *International Journal of Clothing Science and Technology*, vol. 8, no. 3, pp. 75–94, 1996.
- [27] M. Fontana, C. Rizzi, and U. Cugini, "3D virtual apparel design for industrial applications," *Computer-Aided Design*, vol. 37, no. 6, pp. 609–622, 2005.
- [28] P. Lienhardt, "N-dimensional generalized combinatorial maps and cellular quasimanifolds," *International Journal on Computational Geometry and Applications*, vol. 4, no. 3, pp. 275–324, 1994.
- [29] E. Praun and H. Hoppe, "Spherical parametrization and remeshing," *ACM Transactions on Graphics*, vol. 22, no. 3, pp. 340–349, 2003.
- [30] M. Attene and B. Falcidieno, "Remesh: an interactive environment to edit and repair triangle meshes," in *Proceedings of the IEEE International Conference on Shape Modelling International (SMI'06)*, pp. 271–276, IEEE Computer Society Press, Silver Spring, MD, USA, 2006.
- [31] J. Lu, M. Wang, C. Chen, and J. Wu, "The development of an intelligent system for customized clothing making," *Expert System with Application*, vol. 37, no. 1, pp. 799–803, 2010.
- [32] Y. Chen, X. Zeng, M. Happiette, P. Bruniaux, R. Ng, and W. Yu, "Optimisation of garment design using fuzzy logic and sensory evaluation techniques," *Engineering Applications of Artificial Intelligence*, vol. 22, no. 2, pp. 272–282, 2009.
- [33] Z. X. Guo, W. K. Wong, S. Y. S. Leung, and M. Li, "Applications of artificial intelligence in the apparel industry: a review," *Textile Research Journal*, vol. 81, no. 18, pp. 1871–1892, 2011.
- [34] B. C. A. M. Chow, "An investigation on the needs and expectations of the fashion industry in relation to design support systems—a qualitative approach," BA thesis, The Hong Kong Polytechnic University, Hung Hom, Hong Kong, 2012.
- [35] P. Y. Mok, J. Xu, X. X. Wang, J. T. Fan, Y. L. Kwok, and H. Xin, "An IGA-based design support system for realistic and practical fashion designs," *Computer-Aided Design*, vol. 45, no. 11, pp. 1442–1458, 2013.

- [36] X. Wan, P. Y. Mok, and X. Jin, "Shape deformation using skeleton correspondences for realistic posed fashion flat creation," *IEEE Transaction on Automation Science and Engineering*, vol. 11, no. 2, pp. 409–420, 2014.
- [37] J. Xu, P. Y. Mok, C. W. M. Yuen, and R. W. Y. Yee, "A web-based design support system for fashion technical sketches," *International Journal of Clothing Science and Technology*, vol. 28, no. 1, pp. 130–160, 2016.
- [38] Snap T, August 2018, <https://play.google.com/store/apps/details?id=com.avennomys.mytee>.
- [39] Design U, August 2018, <https://play.google.com/store/apps/details?id=com.b05studio.designu>.
- [40] Printing Factory, August 2018, <http://printingdiy.co.kr>.
- [41] Adidas, August 2018, <http://www.adidas.com/us/customize>.
- [42] Y. Yuan and J.-H. Huh, "Customized CAD Modeling and design of production process for one-person one-clothing mass production system," *Electronics*, vol. 7, no. 11, p. 270, 2018.

Research Article

An Indoor Location-Based Positioning System Using Stereo Vision with the Drone Camera

Young-Hoon Jin , Kwang-Woo Ko, and Won-Hyung Lee 

Culture Technology Research Institute, Chung-Ang University, 221 Heuksuk-dong, Dongjak-ku, Seoul, Republic of Korea

Correspondence should be addressed to Won-Hyung Lee; whlee@cau.ac.kr

Received 21 May 2018; Revised 7 August 2018; Accepted 29 August 2018; Published 17 October 2018

Academic Editor: Jaegel Yim

Copyright © 2018 Young-Hoon Jin et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Stereo vision is used to reconstruct 3D information of the space by estimating the depth value from the simulation of human eyes. Spatial restoration can be used as a means of location estimation in an indoor area, which is impossible to accomplish using the relative location estimation technology, GPS. By mapping the real world in virtual space, it is feasible to clear the boundary between real space and virtual space. This paper presents a method to control the drone indoors through a positioning system using Structure from Motion algorithm (SfM). SfM calculates the relative relationship between cameras based on images to be acquired from various locations and obtains disparity to enable restoration of 3D space. First, the 3D virtual space is reconstructed using several photographs taken from an indoor environment. Second, the real-time drone position is determined by comparing the 3D virtual space camera with the image displayed on the drone camera. In this case, if the direction of the virtual camera used for 3D virtual space construction is the same as the amount of yaw rotation of the drone, it is possible to quickly find the same position as the image seen in the real drone camera in the virtual space. As a result, if the scale of the actual camera image and the virtual camera image is 1 : 1 matched, then it is possible to know that the drone is in the position of the virtual camera. The proposed indoor location-based drone controlling method can be applied to various drone applications such as group flight in an indoor environment because of its ability to fly the drone without the use of the traditional remote-control and flight trajectory programming.

1. Introduction

In modern society, drones are aviation system with a wide range of applications. The drone developed in the military is efficient equipment that can perform tasks at a low cost and without any risk in industry, agriculture, and disaster prevention. Drones can generally be controlled using equipment such as GPS, cameras, laser scanners, and ultrasonic sensors. For example, GPS can measure position and altitude, cameras can acquire images that are difficult to see with the naked eyes, and laser sensor can measure the distance between objects. The measured data can be used for autonomous flying or object recognition.

The drone can acquire various information while measuring the position as well as altitude, but it is relatively limited in the indoor environment. GPS works only outdoors and a good laser scanner is big and heavy for flying small indoor drones. On the contrary, big drones that could

carry laser scanners are dangerous and difficult to control in the indoor environment. Even though using a light-weight ultrasonic sensor can measure the distance to some extent, the measurement is inaccurate. In order to address this problem, research in the field of computer vision is incorporated in drone technology.

Computer vision can be used for position control of the drone. Distance measurement through a camera generally uses two cameras, which is similar to the principle of the human eyes. Since the two eyes of a human observe an object at different positions, there is a difference between the vantage point angles, which is called disparity, and this disparity can be used to estimate the distance of the observed object from the 'eye' or the camera. Stereo vision is based on this principle.

After measuring the distance, it is possible to determine the relative position of the drone through the camera without using GPS or laser sensors. In order to determine the relative position of the drone, it is necessary to determine the

relative distance from the origin in a 3D virtual environment. The 3D reconstruction of the indoor space can be constructed using Structure from Motion (SfM). SfM reconstructs 3D coordinates from 2D images with disparity, which is the characteristic of stereo vision using multiple images. Multiple images for SfM can be easily acquired with a drone.

As shown in Figure 1, the SfM pipeline finds the feature points in each image and finds the feature points corresponding to the points of interest in one image into the other image, and this process is called correspondence. Various studies have been carried out to find corresponding points; Lowe proposed an algorithm to extract features that are invariant to image size and rotation [1]. The algorithm constructs a scale space and finds the key points through the DoG (difference of Gaussian) operation, then removes any keys that do not meet the criteria, and assigns the directionality to the appropriate key. By configuring the descriptors by assigning a fingerprint to these keys, a unique matching point with scale and rotation invariance is constructed. These matching points have high accuracy but require a large computational complexity and long execution time.

Bay proposed a SURF (Speed-Up Robust Features) algorithm that resolves these drawbacks [2]. To use the integral image, the fast Hessian detector was used. If the determinant is a positive number and the sign of the given value is the same, it is assigned as a key point. Similar to SIFT (scale invariant feature transform), key points are extracted from various scales. Afterwards, the direction is given through the Haar wavelet response. SURF is relatively fast but uses only gray space information.

Once good feature points are found, the path through which the feature points move in each image can be analyzed. Carlo Tomasi proposed a method of tracking strong feature points in every frame through optical flow [3]. It is shown that the optical flow in the region N of the window centered on the pixel is the same and can be traced. Since the above method is a local algorithm, the size of the window is important. It also has a disadvantage because optical flow is affected by the instantaneous changing light.

On the contrary, there is a guided searching method approaching through epipolar geometry. The feature points of an image are mapped to other images by the matrix F and near the epipolar line. In this case, F is called the fundamental matrix, and the rotation and translation $[R|t]$ vectors of the camera can be constructed from F [4–11]. This allows us to identify camera movement.

Wang constructed an indoor space using VisualSfM by obtaining images of multiple angles and the camera trajectory of the acquired images and making “belief maps” through FCNN (fully convolution neural networks) to complete the space [12].

Snavely uses image-based rendering technology to complete 3D model correspondence from various images published on the web and to navigate the virtual space through a smooth transition between images [13].

Ryan et al. conducted an epidemiological investigation at low cost in the western part of the Greenland ice sheet and proved to be efficient in characterizing large jet glaciers [14].

In addition, indoor location-based studies using RSSI (received signal strength indication) or TOF (time of flight) are introduced in this study [15, 16].

Multiview reconstruction is possible, based on the above studies, and spatial reconstruction is utilized efficiently in various areas.

Dense reconstruction is required because the 3D coordinates formed along the camera motion are low in density. In general, dense reconstruction through multiview stereo (MVS) algorithm of Figure 2 has the type of depth map, point cloud, volume scalar field, and mesh [17].

The depth map can be used in various areas such as scene analysis and visualization, but there is a problem in merging the 3D model of the entire area. Likewise, the quality of the 3D model may deteriorate. The point cloud is easy to merge and split, and it can overcome the drawbacks of depth map because it creates a single point cloud from all input images.

The volume scalar field method can be restored from images, depth maps, and point clouds, but integrating into a single mesh is a difficult problem.

This paper aims at 3D reconstruction through an indoor location-based drone control method using a single-view camera mounted on a drone. Various researches in computer vision are examined in the development of this study. In the proposed method, the camera image obtained from the drone is transmitted to the ground station and the spatial coordinates reconstructed by SfM are visualized by the 3D modelling program. If the position of the drone is estimated by comparing the image projected on the virtual camera with the current image of the drone camera, it is possible to control the drone using the virtual spatial coordinates. If the yaw rotation of the drone and the rotation of the virtual camera are equal, the position of the drone can be estimated with high accuracy and it is proven with experimental result (Section 4) in the proposed method.

The rest of this paper describes the proposed method. First, the feature point estimation in each image is explained in Sections 2.1 and 2.2 discusses the epipolar geometry configuration using correspondence. Sections 2.3 and 2.4 describe the 3D reconstruction of correspondence. In Sections 3.1 and 3.2, an environment map is presented and the position estimation of the drone is discussed, respectively.

The experimental verification and discussion are presented in Section 4, and finally, Section 5 draws the conclusion on the proposed method.

2. 3D Reconstruction at Indoor Positioning System

The pipeline for SfM is shown in Figure 1. SIFT extracts the feature from the input image sequence and correspondence is calculated from each image based on the extracted feature. By calculating the fundamental matrix using features with high accuracy, the essential matrix through the camera matrix can be determined. By decomposing the essential matrix into singular value decomposition (SVD), the rotation/translation matrix is calculated. Once calculating $[R|t]$, the movement path of the camera can be known and 3D coordinates can be obtained through triangulation.

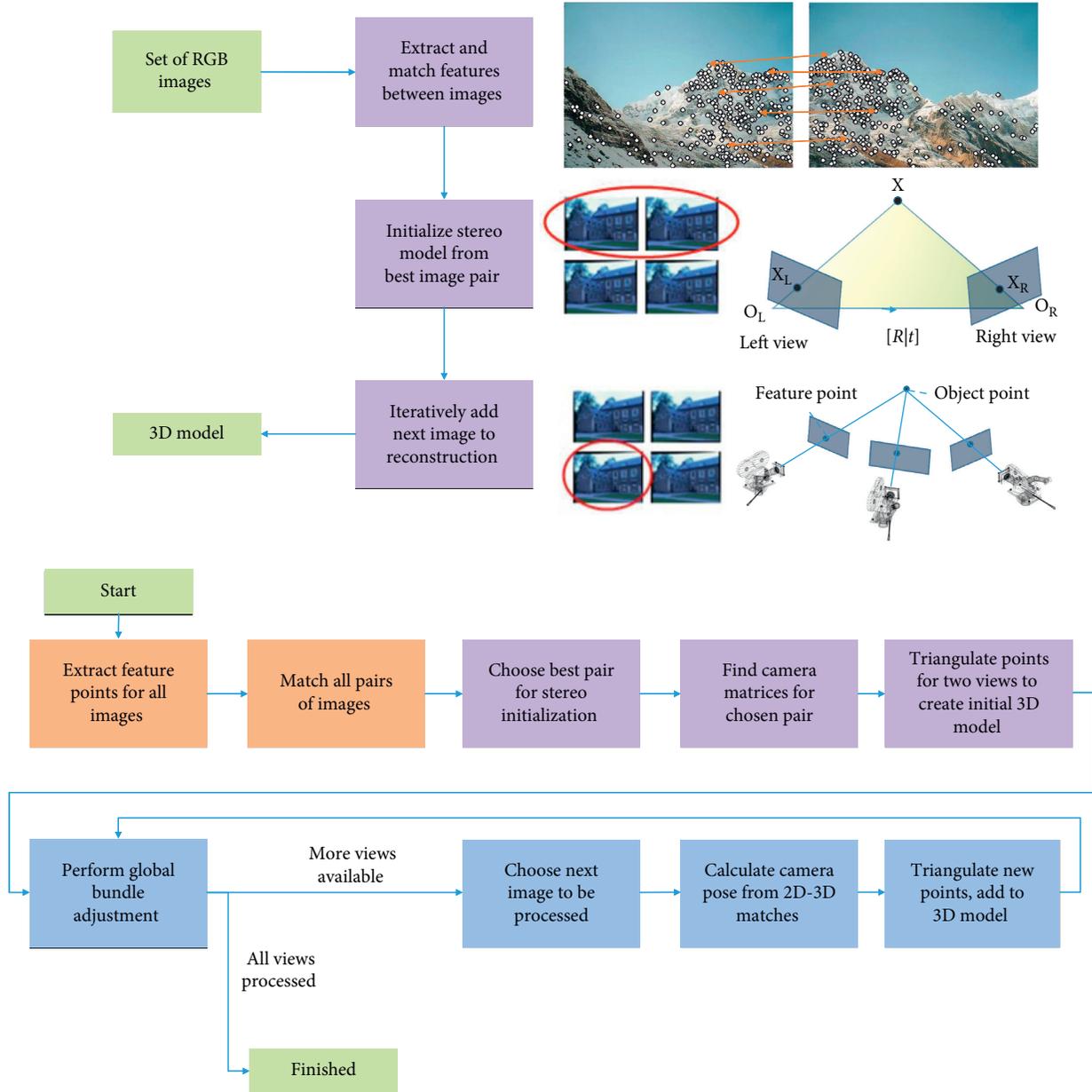


FIGURE 1: Incremental multiview reconstruction [9] and pipeline.

Repeatedly applying this to all image sequences will enable 3D reconstruction.

2.1. Feature Detection and Matching. SIFT is an algorithm that extracts feature points that are robust to scale and rotate. Figure 3 depicts SIFT algorithm procedure.

Scale-space extrema detection generates Gaussian pyramid and calculates DoG to extract strong feature point candidates for scale change. The key point localization step extracts the correct feature points from the candidate group through the Taylor series. Extracted feature points are assigned directionality in the orientation assignment step. The orientation histogram is formed by Gaussian blurring, and then the orientation is estimated.

Finally, the directionality is assigned to a certain area and the descriptor is completed as shown in Figure 4.

The matching method of feature points in the two images finds the same region through comparison of calculated descriptors. In this case, the easiest comparison method is to find the same pair by comparing the distances of two sets as shown in Equation (1) by pairwise matching. Distance matching has advantages of being relatively faster because of being simple to calculate and easy to implement.

Euclidean distance matching:

$$D = \sqrt{\sum_{i=1}^n (a_i - b_i)^2} \quad \begin{matrix} A = (a_1, a_2, a_3 \dots, a_n), \\ B = (b_1, b_2, b_3 \dots, b_n). \end{matrix} \quad (1)$$

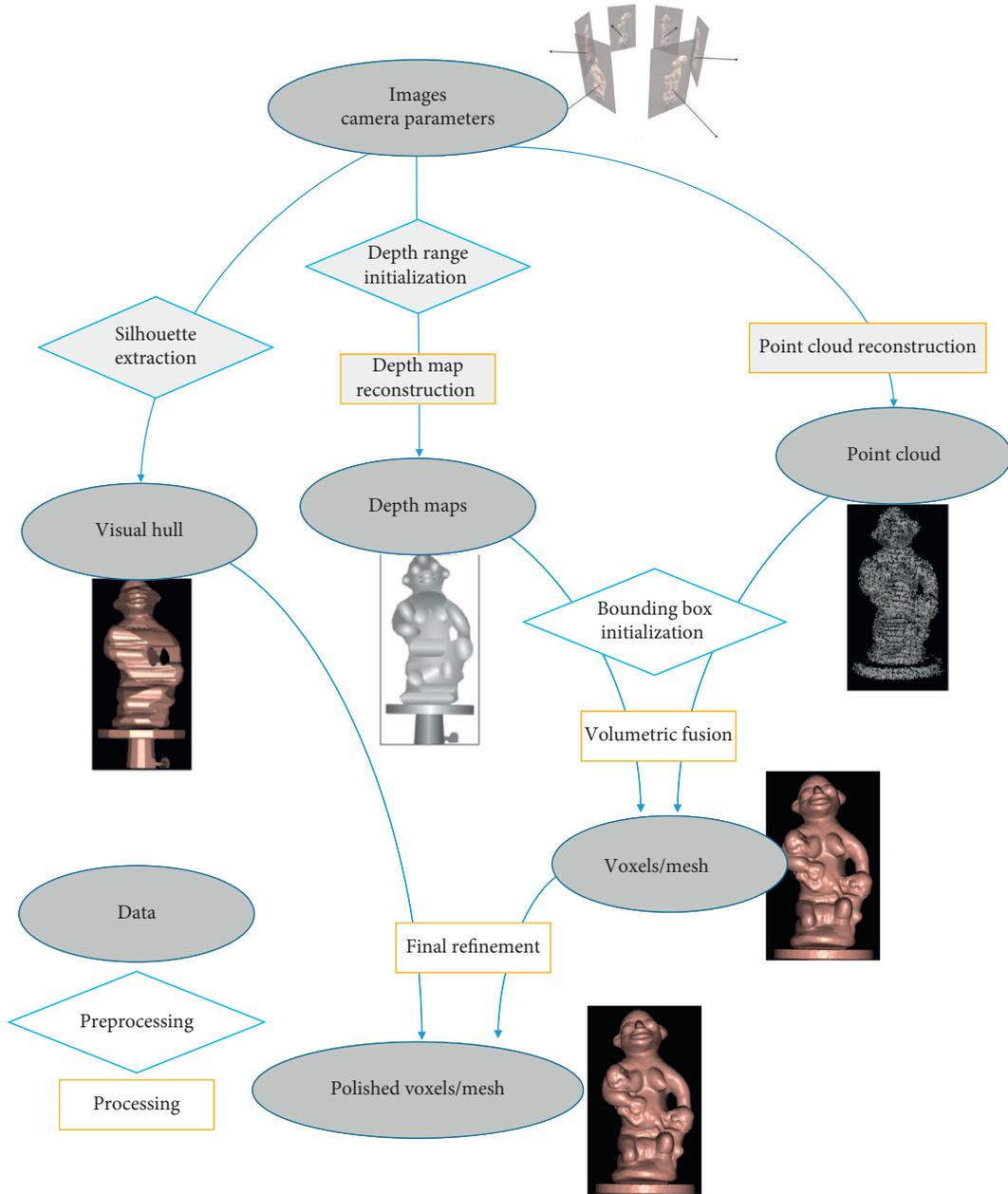


FIGURE 2: MVS processing [17].

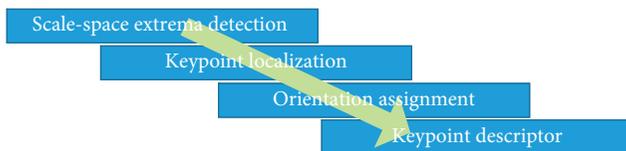


FIGURE 3: SIFT algorithm procedure [1].

After establishing the correspondence of each image, the path of camera movement can be obtained as described in Section 2.2.

2.2. Fundamental and Essential Matrixes. The fundamental matrix is a matrix containing the properties of the camera,

while the essential matrix contains the geometric relationships of the pixel coordinates on the two images. The matrix F in Equation (2) can be calculated when a correspondence is constructed, if there are the eight corresponding points.

Fundamental matrix [4]:

$$p_2^T F p_1 = 0, \quad (2)$$

where F : fundamental matrix. p_1 and p_2 : corresponding points.

The essential matrix implies a geometric relationship in the normalized image plane. This means that the geometric relationship of both cameras is toward a point in the space. In other words, matrix E shows the rotation and movement relationship between the two cameras.

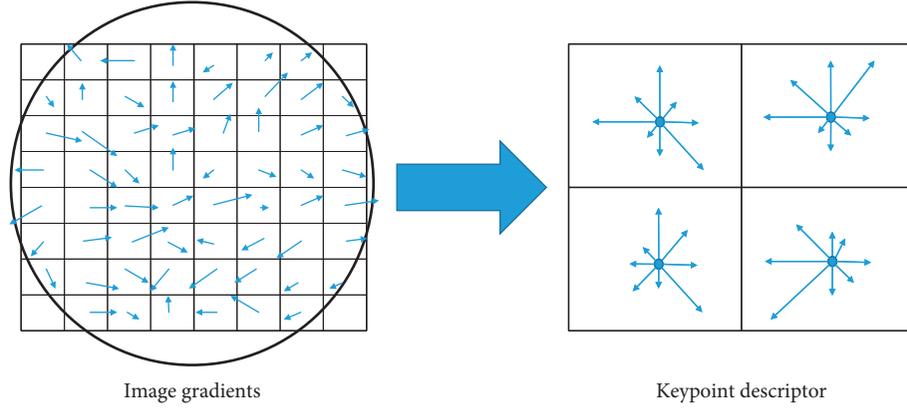


FIGURE 4: Key point descriptor [1].

The matrix E can be computed by $E = K'^T F K$ or the corresponding five points and the camera matrix K , where K is obtained from the EXIF metadata.

Essential matrix [4]:

$$p_2^T K^T E K p_1 = 0, K = \begin{bmatrix} f & 0 & x_{pp} \\ 0 & f & y_{pp} \\ 0 & 0 & 1 \end{bmatrix}, \quad (3)$$

f : focal length and x_{pp} , y_{pp} : the center of the image.

Matrix E can be decomposed into a rotation matrix R and a movement matrix t by SVD. If the camera matrix of the first image is P , then the camera matrix P' of the next image can be obtained as Equation (4).

SVD decomposition of the essential matrix [4]:

$$\begin{aligned} E &= U_{\text{diag}}(1, 1, 0)v^T, P = [I|0], \\ P' &= \begin{bmatrix} UWV^T|+u_3 \\ UWV^T|-u_3 \\ UW^T V^T|+u_3 \\ UW^T V^T|-u_3 \end{bmatrix}, \end{aligned} \quad (4)$$

There are four cases of P' , and the solution in which the feature points exist in front of the two cameras at a single point can be selected.

2.3. Sparse Reconstruction. The motion of the camera in the 3D space can be estimated by obtaining the rotation and movement matrix $[R|t]$. Figure 5 depicts the configuration of the 3D coordinates through camera movement in spatial coordinates. If the epipolar geometry [4, 17] is constructed from the position of the camera estimated from the initial image in the sequence image, a 3D coordinate can be obtained through triangulation [11]. The calculated 3D coordinates are projected on the added camera, and the error is calculated with the feature points formed in the image. This is called the reprojection error [4, 17, 18], and the reconstruction work is performed to correct the error to the minimum. Expression of the reprojection error Equation (5) is as follows. P_i is the projection matrix of the i -th image, X_j

is the j -th 3D point, and x_{ji} is the observation of the 3D point.

Reprojection error:

$$RE(P, X) = \sum_{j=1}^m \sum_{i=1}^n D(x_{ji}, P_i X_j)^2. \quad (5)$$

2.4. Dense Reconstruction. The sparse reconstruction method reconstructs the camera position and direction at the moment of acquiring the image through the corresponding point. The dense reconstruction method, on the contrary, performs dense reconstruction based on known camera motion [17, 19, 20]. Figure 6 depicts the method of dense reconstruction, which constructs the epipolar geometry through the camera's direction and movement path in the 3D space and finds the corresponding points in the epipolar line. It can be quickly found as an epipolar constraint. When the corresponding points are found, triangulation is performed centering on the corresponding points to construct dense 3D space coordinates.

The triangulation can be confirmed by the optical ray intersection of the direction of the feature point mapped to the image from the origin of the camera. However, since the exact intersection cannot be confirmed due to the measurement noise, a pixel with a small error is selected. When a set of selected pixels is referred to as a patch, as shown in Figure 6, the patch that occludes other patches or is hidden from multiple patches is judged to be outliers, thereby increasing the accuracy.

3. Indoor Location Estimation of the Drone Camera

The drone flight usually involves acceleration, gyroscope, and geomagnetic sensors. Acceleration/gyro sensors are involved in the horizontal plane of the drone, and geomagnetic sensors are involved in yaw rotation of the drone. In addition, various sensors can be mounted in the drone. The laser tracking or GPS is often used to locate the drone. But laser tracking is not traceable when the drone is obstructed by obstacles, and the GPS does not work indoors

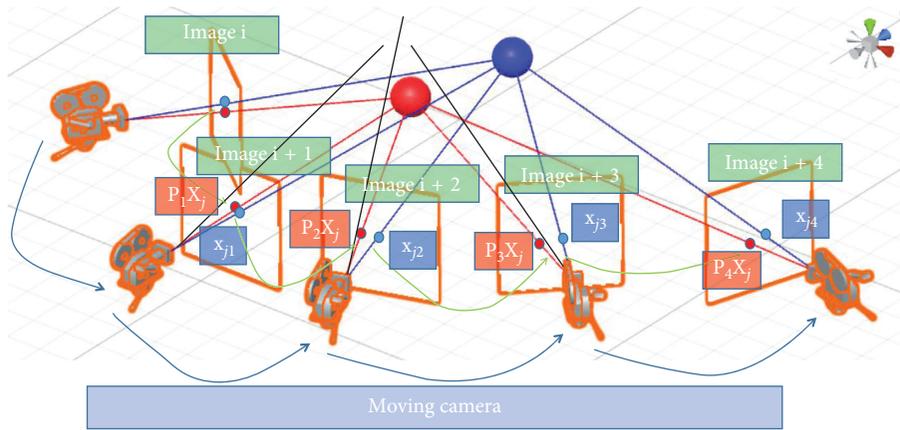


FIGURE 5: Sparse 3D reconstruction with the reprojection error.

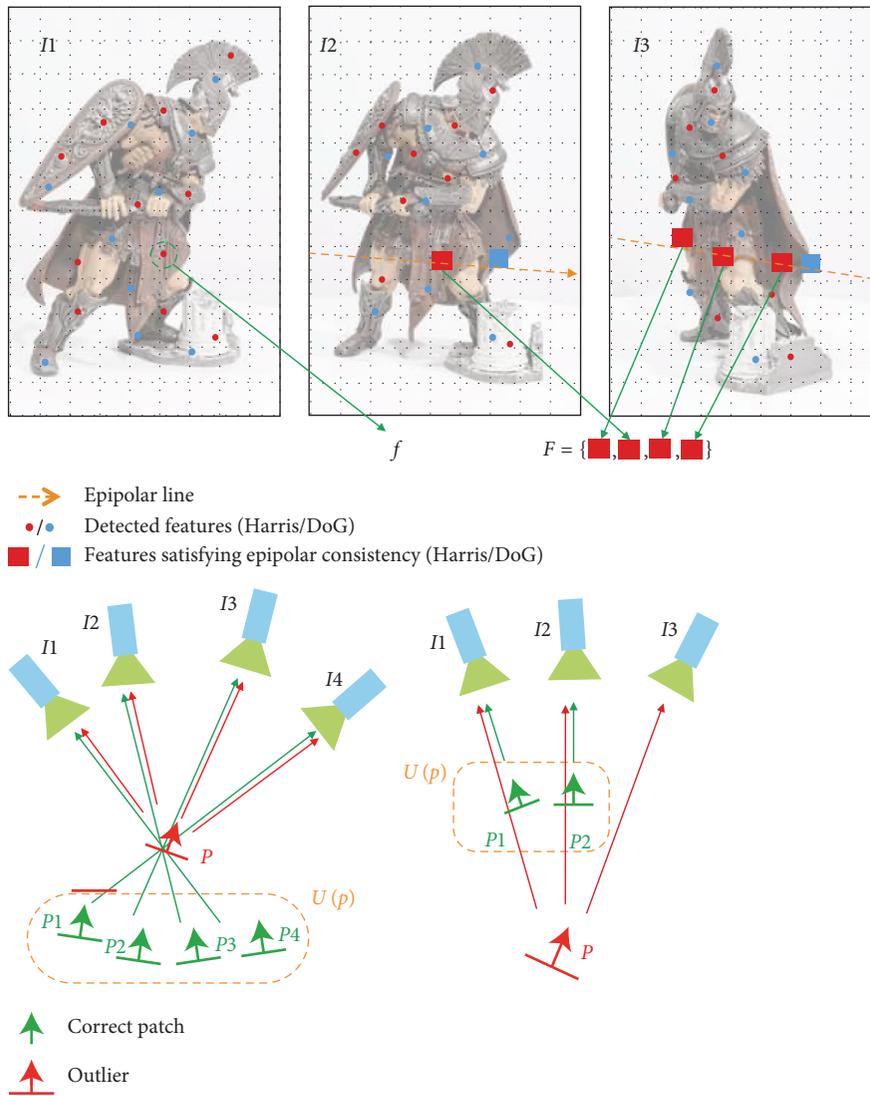


FIGURE 6: Dense reconstruction [19].

because the signal cannot reach. Indoor location estimation can be performed by installing a beacon or an infrared sensor, but accuracy and installation limitations exist. In this paper, the authors propose a method for estimating the indoor position through the computer vision technique.

3.1. Configure Environment Map. In a three-dimensional space, as the position is a relative concept, the distance and direction away from the origin in the space where the origin exists is the position. In other words, in order to locate the drone, an environment map (virtual space) must be constructed first. The origin of the environment map is set as the last position or the average value of the camera positions when acquiring the images for the spatial configuration. The 3D space can be constructed from the sequence images formed by moving along the wall as discussed in Section 2.3. As shown in Figure 7, the spatial configuration is reconstructed through the SfM procedure by obtaining 360° space images through yaw rotation and circular trajectory for a certain time after the takeoff of the drone. The yaw rotation amount of the drone at the time of the last image acquisition is the forward vector.

The partially constructed environment maps are integrated into one coordinate system by registration. Figure 8 illustrates how to make registration with the 12-zone point cloud. The registration can be done through an iterative closest point (ICP) process by constructing key points and looking for correspondence. When two sets of coordinate systems $X = x_1, \dots, x_{N_x}$ and $P = p_1, \dots, p_{N_p}$ exist, the registration as shown in Figure 8 is completed by repeating the calculation of $[R|t]$ where the error of the transformation relation $E(R, t) = (1/N_p) \sum_{i=1}^{N_p} \|x_i - Rp_i - t\|^2$ is minimized [22].

Since the configured environment map is in the form of a point cloud, restoration of the mesh shape is necessary. The mesh is generally constructed by removing the noise properly, smoothly treating the rough surface, and then repeatedly bundling the mesh into the appropriate triangular unit.

Mesh reconstruction uses marching cube, Delaunay triangulation, and Poisson surface reconstruction. The Poisson surface reconstruction, which is strong against matching error and noise, is widely used.

Figure 9 depicts the Poisson surface reconstruction. Since the interior of the model has a negative number and the outer one has a positive number, the boundaries have zero, so a directed sample can be interpreted as a gradient discretization. Thus, we can solve this problem by calculating the divergence of successive vector fields in a point cloud and solving Poisson's equation to find the scalar field with the most appropriate gradient and discretizing the octree. In this case, the screening algorithm is used to reduce the time complexity and improve the accuracy of the solver by interpolating and processing the points in the partial domain rather than the whole domain.

3.2. Location Estimation of the Drone. There are two cases of a location estimation of the drone. First, the position of the

camera can be determined by analyzing the sequence image added to the existing sequence, as shown in Section 2.3. As a result, the position of the camera becomes the position of the drone. In this case, the added sequence should be a position where the corresponding point can be found in the existing sequence. The other is a method of comparing an image captured by a virtual camera in a 3D space with an image captured by the current drone camera as depicted in Figure 10. As shown in the figure, when the drone moves toward the target point C, it approaches as the horizontal-vertical flight. This is to increase the accuracy of the image comparison by reducing the change of the image rather than the straight flight toward C.

Let D be the position of the current drone, C be the position of the virtual camera, and P be the plane parallel to the drone. Then the vector M is an orthogonal projection vector of the vector DC to the plane P. At this time, the angle between the forward vector F of the drone and the vector M can be obtained through $\theta = \cos^{-1}(\vec{F} \cdot \vec{M} / \|\vec{F}\| \|\vec{M}\|)$. If the virtual camera also looks toward the M vector, the image similarity can be estimated.

Image similarity can be achieved by the template matching, which is a method to search and check whether a given small template image exists in a large image. Template matching has the disadvantage of only showing good results with the same scale and direction, but it can produce good results because the direction of the drone is the same as the direction of the virtual camera. The similarity of the image can be calculated by Equation (6).

Template matching similarity:

$$R(x, y) = \sum_{x', y'} (T(x', y') - I(x + x', y + y'))^2. \quad (6)$$

Figure 11 depicts the template matching. Since the drone camera contains more area than the virtual camera image of the target point, it searches the image of the virtual camera in the image of the drone. At this time, when the image of the virtual camera and the image of the drone become the same size, it can be known that the drone has arrived at the target point.

The location estimation through the template matching can reduce the error if the field-of-view of the virtual camera is the same as the drone camera. Also to overcome the error due to the image scale, we construct an image pyramid for the template image. When the drone reaches the target point, a new image sequence is formed and the procedure of Section 2 is repeated to finalize the position. In this case, if the correspondence of the current sequence cannot be found in the existing sequence, it can be estimated as shown in Figure 12. Then, the image of the drone camera is searched from the images projected on the virtual camera in the same direction as the yaw rotation amount of the current drone. At this time, the image of the drone camera becomes a template, the virtual camera moves in the X-Y plane so that the template area is located at the center of the virtual camera image, and the point at which the image scale becomes 1 : 1 by the z-axis is determined as the position of the drone.

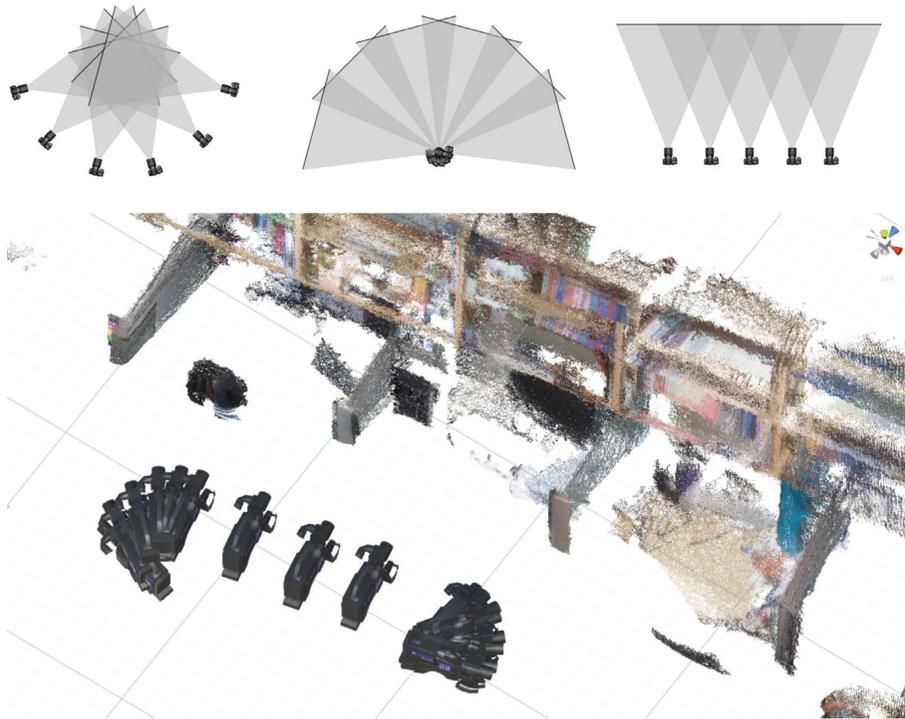


FIGURE 7: Image acquisition method for spatial configuration [21] and result.



FIGURE 8: Registration.

4. Experimental Results and Discussions

Experiments utilize a programmable Arduino drone. Since an Arduino drone is equipped with a low-performance micro-processor, image processing is not possible, so the station is constructed on the ground. The system of the station role consists of Intel 4405U 2.1 GHz CPU, 8 GB ram, uses Unity3D, VisualSFM [25], CloudCompare, and MeshLab for 3D visualization. Figure 13 depicts the entire system. As shown in Figure 13, the drone and the station on the ground exist on the same network and communicate data with each other.

After the drone takes off, it keeps a constant height and scans the space using the image acquisition method of

Figure 7. Precise spatial restoration requires large amounts of photographs with good disparity. The acquired images are sequentially transmitted to the station as shown in Figure 13.

As shown in Figure 14, VisualSFM performs sparse reconstruction from the transmitted image and completes dense reconstruction through clustering views for multiview stereo (CMVS). The point cloud data are meshed by MeshLab, and then the environment map is finally constructed.

Once the environment map is configured, Unity3D visualizes the environment map and takes care of all the control of the drone. The current drone is maintained in the position and orientation in which the last image of the

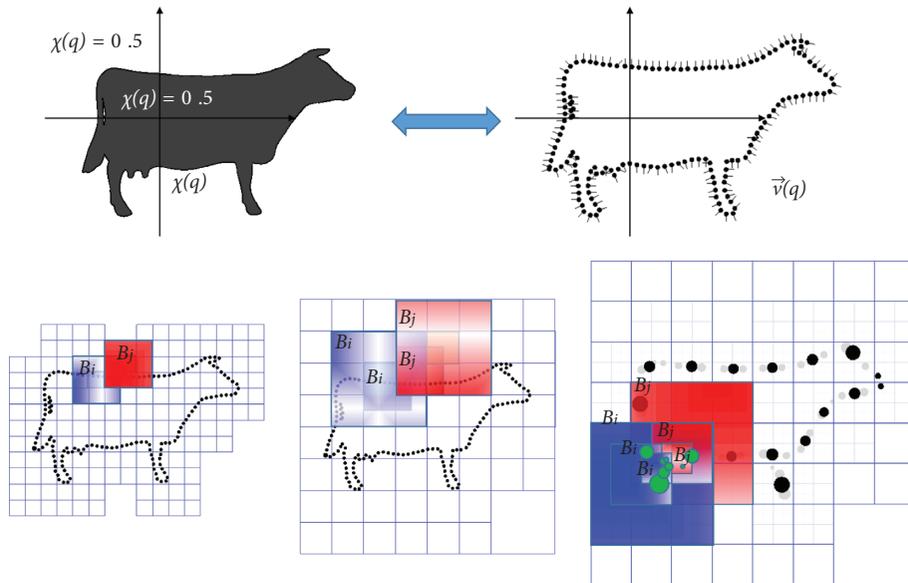


FIGURE 9: Screened Poisson surface reconstruction [23, 24].

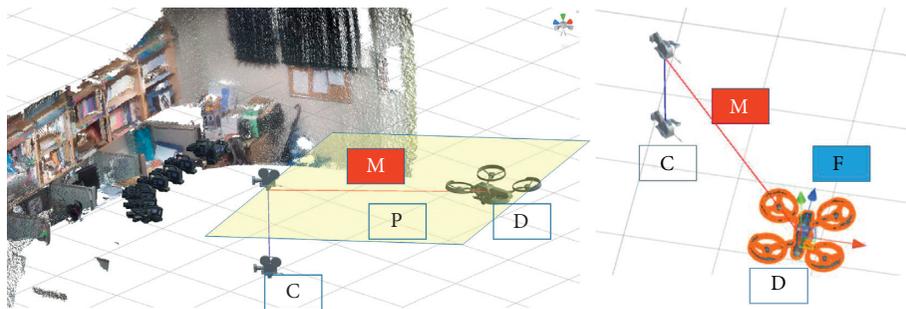


FIGURE 10: Location estimation through virtual camera analysis.

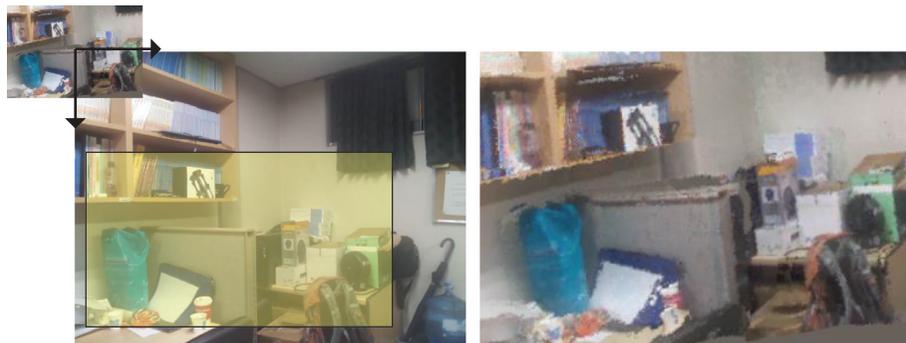


FIGURE 11: Template matching.

environment map configuration was obtained. Based on this, the origin and direction (forward vector) are matched in the world coordinate system of Unity3D. Once the basic setup is complete, the location can be confirmed via the virtual camera.

Figure 12 depicts the current location estimate. The yaw rotation of the drone can be thought of as a camera that reflects a certain range in the space coordinates. As shown in the figure, since the current position is estimated except for

the back surface of the plane P which the current drone cannot be seen, the calculation amount can be reduced to $1/2$. At this time, the normal vector of the plane P is in the same direction as the forward direction of the drone. It also searches at the last position of the drone, allowing faster calculation.

Once the environment map is configured, the current position of the drone is confirmed, and then a controller for controlling the drone is not required. It is possible to set the



FIGURE 12: Search area and projection viewable area by yaw rotation amount.



FIGURE 13: System configuration.

target coordinates of the drone in Unity3D. When the target coordinate is set to P as shown in Figure 15, the horizontal-vertical movement path is calculated as shown in Figure 10.

As shown in Figure 15, if a certain frame image is transferred to the station while moving horizontally, the similarity is calculated as shown in Figure 15 using template matching. In the horizontal movement, the virtual image

(T1) of the horizontal target point is located in the middle region of the actual image, and each image becomes 1:1 scale when the horizontal target point is reached. The initial coordinate of the drone is $x: -0.66, y: 1.149, z: 6.425$ and azimuth is $x: 0, y: -90.8, z: 0$. The target coordinate of the horizontal movement is $x: -0.405, y: 1.026, z: 2.7$, and the azimuth angle is $x: 0, y: -187, z: 0$.

When the horizontal movement is completed, the current image is transmitted to the station and the vertical movement is performed. The vertical movement is sometimes not possible with template matching because it is close to the wall or when there is a large change in the image between the instantaneous arrival point and the final arrival point in a narrow and high space. Therefore, template matching is performed by placing virtual camera images at regular intervals. Since the experimental environment of this paper is not so large, there is no big problem.

As shown in Figure 16, the vertical movement coordinate is $x: -0.405, y: -0.124, z: 2.7$, and the azimuth is equal to the azimuth of horizontal movement.

When the final target is reached, the current image is transmitted to the station and the position of the drone is updated as described in Section 2.3.

As in the experiments, it is possible to obtain the relative position of the drone and control the drone in the virtual space without the controller of the actual drone. But there is some problem, the camera is very sensitive to environmental illumination condition and it is a time-consuming process to construct a detailed virtual space using a sequence of images. Also an error exists. In addition, the small drone is vulnerable to hovering, so the drone camera images may not be clear. In Figure 15, the target point (T2) of the horizontal movement was not suitable for template matching because it was too far from the starting point of the drone. In this case, a good result can be obtained by placing a virtual camera at a coordinate closer (T1) to the starting point of the drone than the coordinate T2 and performing a matching procedure. This is considered to be a phenomenon caused by the qualitative effect of the virtual space or the low resolution of the camera.

An alternative to the proposed method is the RSSI triangulation method. The method of determining the overlapping point based on the three devices that emit the signal is shown in Figure 17(a). It is relatively easy to install, low cost, and easy to implement. However, it is

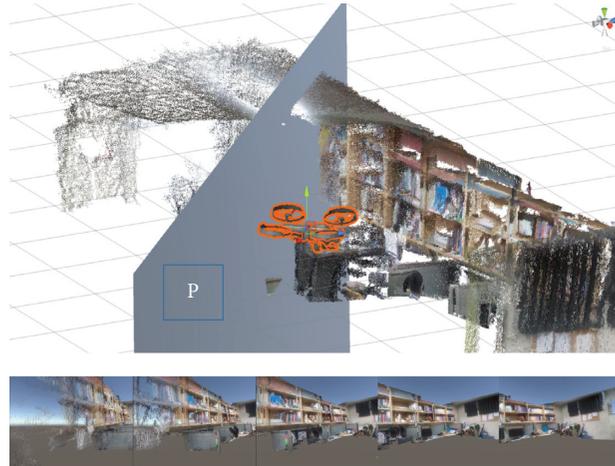


FIGURE 14: Sparse/dense/Poisson reconstruction.



FIGURE 15: After Poisson surface reconstruction, horizontal movement template matching process (virtual/drone camera, drone position). T1, T2: templates.

a two-dimensional estimation and the error rate is high because it generated inadequate noise, and latency is high in real-time tracking. To address this problem, another group estimates the height using ToF camera as shown in Figure 17(b). The ULPS measures the two-dimensional coordinates and the ToF measures the height of the drone

to complete the three-dimensional coordinates. However, there is a limitation in the location estimation of the 2D space and height estimation can be disturbed by the indoor structure. On the contrary, the proposed method takes a lot of time to construct the virtual space, but it has the advantage of being relatively accurate and posing fewer



FIGURE 16: Vertical movement of a virtual camera (drone position, drone/virtual camera) template matching.

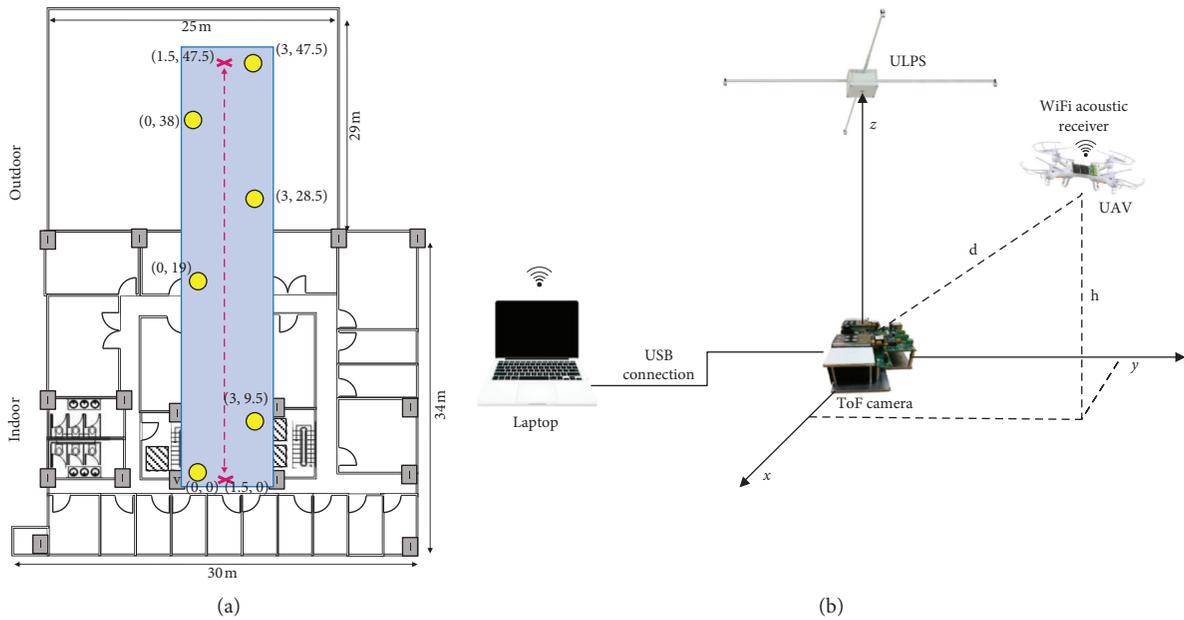


FIGURE 17: (a) Installation for RSSI triangulation [15]. (b) Position estimation with ToF [16].

constraints because the position is obtained based on the similarity between virtual and real space. Of course, direct performance comparison is difficult.

5. Conclusion

The method proposed in this paper shows that it is possible to track the location of the drone using only a single-view camera in the indoor environment. Compared with the position tracking through various sensors method, even

though the 3D restoration process takes a relatively long computational time and cannot be projected in real time, the experimental results guarantee that the accuracy is improved by position correction and image processing. Furthermore, in the proposed method, it is possible to estimate the position of the drone without installing sensors and make 3D reconstruction through additional calculations for shaded areas. Also, by controlling the target point without going through the existing manual controller, the present study may implicate further

research and development on group and/or autonomous flight.

By increasing the power of the station system and using the GPU-based parallel processing, it can be expected to complement the current visual shortcomings of the proposed method and it is expected that better results will be obtained by applying multiview camera or RGB-D sensor in the future.

Data Availability

No data were used to support this study.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

References

- [1] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [2] H. Bay, T. Tuytelaars, and L. V. Gool, "SURF: speeded up robust features," in *Proceedings of European Conference on Computer Vision*, pp. 404–417, Springer, Graz, Austria, May 2006.
- [3] C. Tomasi and T. Kanade, "Detection and tracking of point features," Tech. Rep. CMU-CS-91-132, Carnegie Mellon University, Pittsburgh, PA, USA, 1991.
- [4] R. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, Cambridge University Press, Cambridge, UK, 2nd edition, 2000.
- [5] R. Shah, A. Deshpande, and P. J. Narayanan, "Multistage SFM: revisiting incremental structure from motion," in *Proceedings of 2nd International Conference on 3D Vision (3DV)*, pp. 417–424, IEEE, Qingdao, China, December 2014.
- [6] A. W. Fitzgibbon and A. Zisserman, *Automatic Camera Tracking*, Springer, Boston, MA, USA, 2003.
- [7] M. Han and T. Kanade, "Creating 3D models with uncalibrated cameras," in *Proceedings of Applications of Computer Vision*, pp. 178–185, IEEE, Palm Springs, CA, USA, September 2000.
- [8] M. Pollefeys, L. V. Gool, M. Vergauwen et al., "Visual modeling with a hand-held camera," *International Journal of Computer Vision*, vol. 59, no. 3, pp. 207–232, 2004.
- [9] A. Strupczewski and B. Czupryński, "3D reconstruction software comparison for short sequences," in *Photonics Applications in Astronomy, Communications, Industry, and High-Energy Physics Experiments 2014*, International Society for Optics and Photonics, Bellingham, WA, USA, 2014.
- [10] C. Wu, "Towards linear-time incremental structure from motion," in *Proceedings of 2013 International Conference on 3D Vision*, pp. 127–134, IEEE, Seattle, WA, USA, June 2013.
- [11] R. I. Hartley and P. Sturm, "Triangulation," *Computer Vision and Image Understanding*, vol. 68, no. 2, pp. 146–157, 1997.
- [12] Y. Y. Wang, *Multi-View Indoor Spatial Layout Estimation*, Stanford University, Stanford, CA, USA, 2016.
- [13] N. Snavely, S. M. Seitz, and R. Szeliski, "Photo tourism: exploring photo collections in 3D," in *Proceedings of ACM Transactions on Graphics (TOG)*, pp. 835–846, ACM, New York City, NY, USA, April 2006.
- [14] J. C. Ryan, A. L. Hubbard, J. E. Box et al., "UAV photogrammetry and structure from motion to assess calving dynamics at store glacier, a large outlet draining the Greenland ice sheet," *Cryosphere*, vol. 9, no. 1, pp. 1–10, 2015.
- [15] E. E. L. Lau, B. G. Lee, S. C. Lee, and W. Y. Chung, "Enhanced RSSI-based high accuracy real-time user location tracking system for indoor and outdoor environments," *International Journal on Smart Sensing and Intelligent Systems*, vol. 1, no. 2, pp. 534–548, 2008.
- [16] J. Paredes, F. Álvarez, T. Aguilera, and J. Villadangos, "3D indoor positioning of UAVs with spread spectrum ultrasound and time-of-flight cameras," *Sensors*, vol. 18, no. 2, p. 89, 2017.
- [17] Y. Furukawa and C. Hernández, "Multi-view stereo: a tutorial," *Foundations and Trends® in Computer Graphics and Vision*, vol. 9, no. 1-2, pp. 1–148, 2015.
- [18] C. Engels, H. Stewénius, and D. Nistér, "Bundle adjustment rules," in *Proceedings of Symposium on ISPRS Commission III Photogrammetric Computer Vision PCV'06*, pp. 266–271, Photogrammetric Computer Vision, Bonn, Germany, September 2006.
- [19] Y. Furukawa and J. Ponce, "Accurate, dense, and robust multiview stereopsis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 8, pp. 1362–1376, 2010.
- [20] M. Bujnák, *Dense Reconstruction from Uncalibrated Video*, Comenius University, Bratislava, Europe, 2005.
- [21] J. Dietrich's, "Advanced geographic research blog," 2014, <http://adv-geo-research.blogspot.com/2014/02/camera-geometries-for-structure-from.html>.
- [22] Z. Zhang, *Iterative point matching for registration of free-form curves*, Ph.D. Thesis, Inria, Le Chesnay, France, 1992.
- [23] M. Kazhdan, M. Bolitho, and H. Hoppe, "Poisson surface reconstruction," in *Proceedings of Fourth Eurographics Symposium on Geometry Processing*, pp. 61–70, Eurographics Association, Cagliari, Sardinia, Italy, June 2006.
- [24] M. Kazhdan and H. Hoppe, "Screened poisson surface reconstruction," *ACM Transactions on Graphics (TOG)*, vol. 32, no. 3, p. 29, 2013.
- [25] C. Wu, *VisualSFM: A Visual Structure from Motion System*, 2011, <http://ccwu.me/vsfm/doc.html>.

Research Article

Profile-Based Ad Hoc Social Networking Using Wi-Fi Direct on the Top of Android

Nagender Aneja  and Sapna Gambhir 

YMCA University of Science and Technology, Faridabad, India

Correspondence should be addressed to Nagender Aneja; naneja@gmail.com

Received 1 July 2018; Accepted 16 September 2018; Published 17 October 2018

Guest Editor: Subramaniam Ganesan

Copyright © 2018 Nagender Aneja and Sapna Gambhir. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Ad hoc social networks have become popular to support novel applications related to location-based mobile services that are of great importance to users and businesses. Unlike traditional social services using a centralized server to fetch location, ad hoc social network services support infrastructure-less real-time social networking. It allows users to collaborate and share views anytime anywhere. However, current ad hoc social network applications either are not available without rooting the mobile phones or do not filter the nearby users based on common interests without a centralized server. This paper presents an architecture and implementation of social networks on commercially available mobile devices that allow broadcasting name and a limited number of keywords representing users' interests without any connection in a nearby region to facilitate matching of interests. The broadcasting region creates a digital aura and is limited by the Wi-Fi region that is around 200 meters. The application connects users to form a group based on their profile or interests using the peer-to-peer communication mode without using any centralized networking or profile-matching infrastructure. The peer-to-peer group can be used for private communication when the network is not available.

1. Introduction

Online social networks, for example, Facebook, LinkedIn, or Twitter, are now highly popular among people, and the trend to use the social networking applications on the mobile device is continuously increasing. The pattern of using social networking on the mobile device is being exploited by researchers and service providers to provide location-based social networking [1–3]. Examples of location-based social networking include the Facebook's feature to find exotic locations or friends nearby in a geographical region. However, current social networking applications do not provide location-based services without accessing the present site of a user, and many users consider this a privacy risk. Furthermore, limited data plans for the mobile device and high cost of international roaming constraint users to communicate even in the nearby region. Thus, the current trend is to decentralize the online social network [4].

Ad hoc social network, one-to-one or multipeer connection, can solve this problem of privacy and help facilitate

the communication in a nearby region without using the centralized infrastructure. There are numerous applications; for example, it may also be useful in a business meeting where the distribution of physical business cards is not a convenient method, but the e-cards can be conveniently distributed in a peer-to-peer network to all individuals present nearby. Another application is communication among passengers in an airplane for game playing by children especially when the flight duration is long or for anonymous chatting among interested passengers or with crew members. Furthermore, ad hoc social network can help to communicate in case of natural disasters or government censorship.

Ad hoc wireless peer-to-peer technology connects devices to create a communication group for social interaction. This paper presents a Wi-Fi peer-to-peer-based mechanism called OffAT (OFFline chAT) [5] that helps to find people with similar interests in a nearby region [1, 6, 7] and allows sharing text or handwritten messages without any centralized server. The mechanism can be used to further develop

applications to cater needs of business users to distribute business cards during international conferences, seminars, or meetings for networking when people may have smartphones but do not have access to the Internet.

This paper is organized as follows: Section 2 provides the review of related publications. Section 3 provides research motivation. Section 4 provides architecture and implementation comprising key components used for location-based ad hoc social networking. Section 5 provides results and analysis.

2. Related Work

The ad hoc social network (ASN) has the potential to connect nearby users who should be connected based on similar interests. Eagle and Pentland [8] introduced a service to introduce proximate users by doing profile matching at a central server. The service alerted users of similar interests in a nearby region. Zhang et al. [9] introduced the multihop social network to broadcast promotional offers to nearby users without using the Internet or GPS. The authors used the dynamic source routing ad hoc network protocol to deploy the multihop social network. Zhang et al. [10] considered ASN as the extension of the online social network to further increase global communication. The authors provided an architecture comprising four layers, namely, application, community, network, and device layers. The functions of profile management were implemented in the community layer. Wang et al. [11] implemented the peer-to-peer social network based on Wi-Fi Direct to facilitate communication among nearby users without a centralized infrastructure.

Some researchers have also presented middleware solutions; for example, Aneja and Gambhir [12] published a four-layer architecture with application, transport, ad hoc social, and ad hoc communication layers. Similarly, Bellavista and Giannelli [13] proposed the spontaneous multihop network using a three-layer architecture including the IP layer, spontaneous multihop layer, and semantic dispatching layer. The architecture uses a one-hop network and makes use of the dispatching layer to deliver packets to other nodes without knowing the destination. Bottazzi et al. [14] presented the socially aware and mobile architecture that provides the roaming social network and groups proximate users with similar interests.

Rahman and Hossain [15] developed nine mobile applications using a massive ad hoc social network. The applications are based on Wi-Fi, 4G, Wi-Fi Direct, and other emergency Internet access points. The mobile apps provide nearby services and tools to contact social friends based on the current context and previous history. Ilkhechi et al. [16] provided the decentralized location service scheme that can be used to identify the location in the ad hoc network. The nodes first obtain knowledge of surroundings, and the location of the target is computed based on the knowledge.

Shu et al. [17] presented Talk2Me that is a device-to-device augmented reality social network and allows people

to exchange messages with nearby users. Casetti et al. [18] discussed inter- and intragroup communication technologies using the Wi-Fi Direct protocol. The number of nodes in one group that can participate in the communication depends on the IP class addressing scheme; however, Casetti et al. [18] provided a mechanism to extend the group and the range of the network by having multigroup communication. The multigroup is created by allowing a group owner to become a legacy client or a relay client in another group. However, the groups formed with clients are general groups without doing profile or interests matching. The probability of clients leaving may be high when they do not see interests being matched in the group.

Therefore, prior published works have solved issues related to profile matching at a central server and have been able to provide communication among nearby users without profile matching. As a result, there is a need for a system and a method that provide profile matching of nearby users without using the Internet or central server and allow sharing text or files with the matched users.

3. Research Motivation

Most portable mobile devices today are equipped with Wi-Fi Direct, Bluetooth, and Wi-Fi short-range wireless technologies which can support spontaneous on-the-go social networking by connecting users with similar interests using the ad hoc communication mode. Wi-Fi Direct or P2P technology allows devices to connect to each other to form groups. The devices negotiate roles, and one of them becomes the group owner, and the other devices connect to the owner as a client device. The devices use their MAC address as their device ID for discovery and communication and a temporary MAC address for all frames within a group.

Although direct device-to-device connectivity via the ad hoc mode was available in IEEE 802.11, it is still not widely available in the devices. Furthermore, 802.11z, Tunneled Direct Link Setup enables device-to-device communication, but devices need to be associated with the same access point. Wi-Fi P2P is based on the IEEE 802.11 infrastructure technology, but devices negotiate, and one of the devices becomes the soft access point. In other words, Wi-Fi P2P does not need a centralized fixed physical infrastructure, and any device with Wi-Fi Direct enabled can participate in the negotiation. Wi-Fi Direct or P2P specification is still at an early stage, and the many researchers have started implementing the technology.

The proposed social networking protocol and implementation, OffAT, provides social applications including interests similarity and communication between similar users that has been built using Wi-Fi Direct in the commercially available devices. OffAT allows users to broadcast user's interests in a nearby region without using any centralized infrastructure and performs interests similarity matching and calculates profile similarity to assist users in finding a person of similar interests nearby.

```

// MainActivity.java
public void startOffatAddListener () {
startOffat.setOnClickListener ( new
    View.OnClickListener () {
@Override
public void onClick (View v) {
WifiManager wifi = (WifiManager)
    getApplicationContext ().getSystemService (Context.WIFI_SERVICE);
wifi.reconnect ();
if (!e1.getText ().toString ().equals ("") &&
    !e2.getText ().toString ().equals ("")) {
GlobalData.name = e1.getText ().toString ();
GlobalData.interests =
    e2.getText ().toString ().replaceAll ("\\s+",",").replaceAll (",+",",");
Intent intent = new Intent (v.getContext (),
    WiFiDirectActivity.class);
startActivity (intent);
}
}
});
}
// WiFiDirectActivity.java
public void setNameInterests () {
try {
Method method =
    manager.getClass ().getMethod ("setDeviceName",
    WifiP2pManager.Channel.class, String.class,
    WifiP2pManager.ActionListener.class);
method.invoke ( manager, channel,
    GlobalData.name+"#" +GlobalData.interests,
    new WifiP2pManager.ActionListener () {
public void onSuccess () {}
public void onFailure (int reason) {}
} );
} catch (Exception e) {
Toast.makeText (getApplicationContext (), "Unable
    send User name and
    Interests", Toast.LENGTH_SHORT).show ();
}
}
}

```

ALGORITHM 1

4. Architecture and Implementation

The main components of the proposed mechanism are device discovery, interest-based social network, and intergroup communication.

4.1. Device Discovery. Wi-Fi Direct devices discover each other using device discovery, wherein a P2P device selects a listen channel and alternates between the search state and listen state. The time for each state is allocated randomly between 100 ms and 300 ms but can be configured. Some researchers have used MAC ID as the device ID in the discovery process or some confidential ID that limits to connect with only a specified user [19] or using an Internet connection [20, 21]. For example, MobiClique [22] used Bluetooth for device discovery to locate nearby users and a central server for matching the user profiles. Profile

matching at a central server is only feasible in some cases due to the additional requirement of the server or infrastructure.

The proposed method in this paper provides the mechanism and user interface to configure device ID as user name#user interests. The device ID by default is MAC address and is broadcasted as SSID in the nearby region. This paper proposes an implementation of changing device ID to user name#user interests so that interests can be broadcasted in the nearby region. The total length of name#user interests is restricted to 32 characters as per the IEEE standard for the length of SSID [23]. The broadcasted interests can be representative keywords to be used as the first layer to filter users before actual connection. The representative keywords can be the most frequent keywords that appear in users' prior actions including browsing history.

Using device ID to match user interests helps to reduce network overhead of establishing a social network of users who are not similar to each other and likely to not

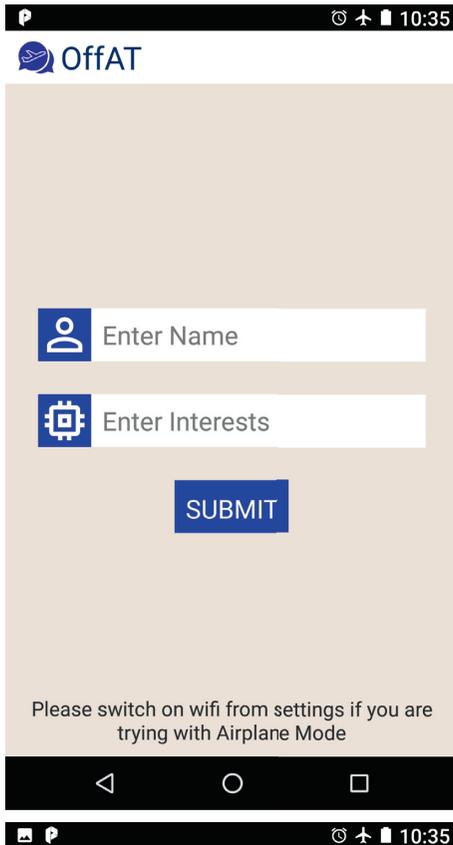


FIGURE 1: User interface to configure device ID.

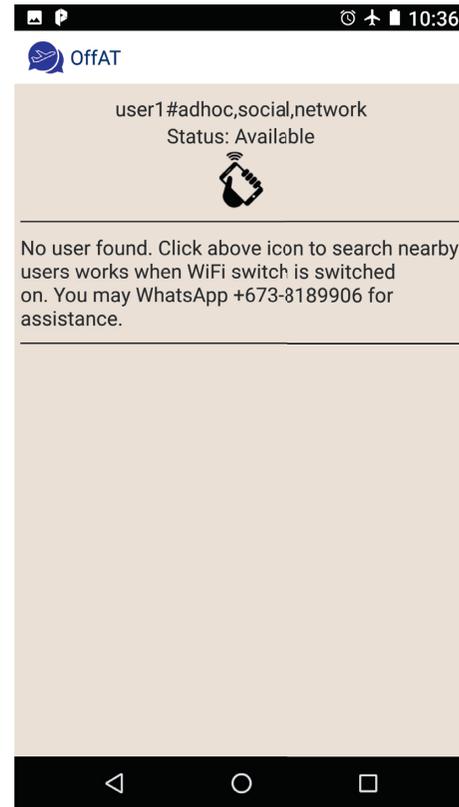


FIGURE 2: Device ID (displaying user name#user interests) in the discovery phase with status as available.

involve in social applications. The below code explains the methodology used to extract user ID and interests in *MainActivity.java* and set the device ID as user name and interests in *WiFiDirectActivity.java* in Java (Algorithm 1).

Figure 1 displays a user interface to configure device ID. Interests are automatically separated in the backend as a list of keywords based on space or comma operator as shown in the code. Modifying device ID to interests solves the following purposes:

- (1) All user devices can access interests of other users without any centralized infrastructure
- (2) Device can compute interests similarity and can decide whether to connect with other users or not
- (3) Computed similarity can also be used to accept or reject any incoming connection request
- (4) Reduces network overhead due to exchanging profiles

A device can have two statuses such as available or connected. A connected device can either be a group owner or an existing group member. Figure 2 shows the device discovery phase and displays the device status as available.

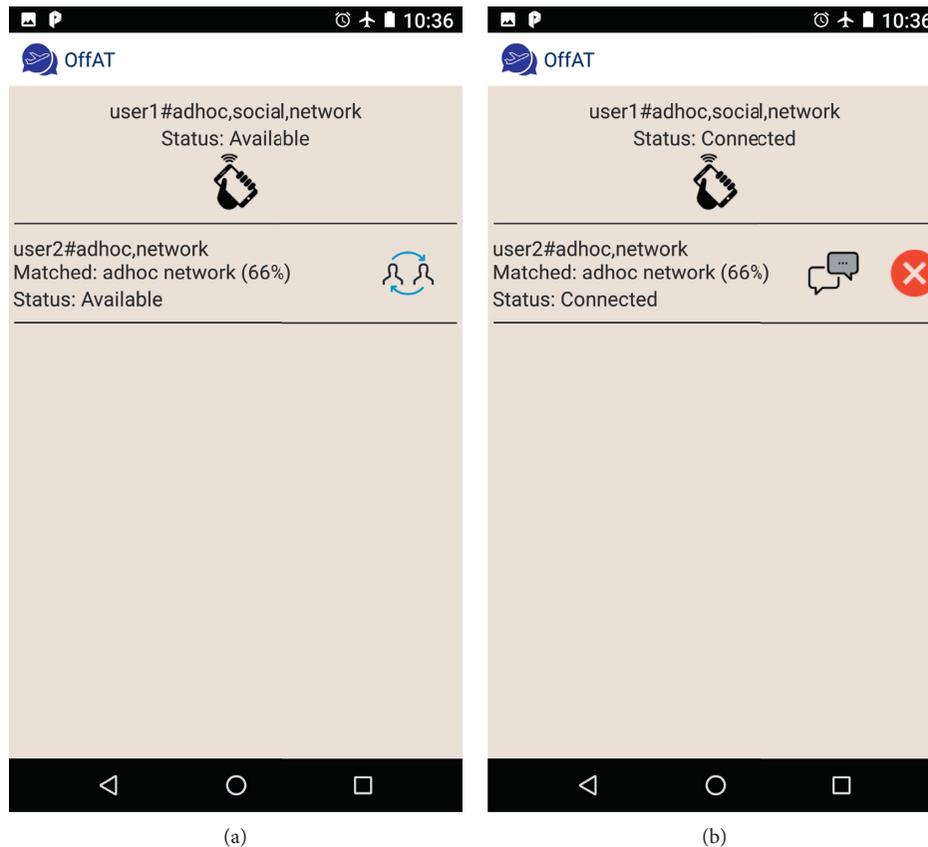


FIGURE 3: Profile/interests similarity with nearby devices as available/connected. (a) Nearby available devices with their interests and profile similarity. (b) Nearby connected devices with an option to chat/disconnect.

4.2. Profile Matching. The profile matching is the core of social networking, especially location-based social networking. Figure 3(a) displays a user device with device ID as user name#interests as available along with nearby available devices, and Figure 3(b) displays a user connected with user2. The discovery process can be manually started by clicking the hand icon available below user name and interests in case no user is found. The OffAT scans interests of nearby users and computes profile similarity based on matched keywords. The option to send a request to a neighboring user, preferably a user with high profile similarity, is displayed along the right side of each user. The list of users can quickly be scrolled to see all users if the number of users is high. Once the other user accepts the connection request, the status of both users will get connected as shown in Figure 3(b). The option to chat or disconnect and start the connection or discovery phase is displayed to the connected user as shown in Figure 3(b).

A number of researchers proposed creating the weighted profile from the user browsing history and/or prior user action [6, 24] and cosine similarity to match the profiles. The cosine similarity of two user profiles varies between 0 and 1. Cosine similarity is 1 when the angle is 0 meaning the two profiles are exactly the same. In a vector space model [24], a user profile is represented by $\{p_1, p_2, \dots, p_n\}$, where the elements are different keywords. A union of all keywords of profiles is considered before computing cosine similarity.

The weights represent the number of times a word appears in a browsing history or prior user action [6, 24]. Gambhir et al. [25] found that cosine similarity is not a correct method especially in case of a weighted profile. For the implementation of the mechanism and simplicity, the profile matching matches keywords to find the profile similarity percentage.

4.3. Social Communication. Once the users are connected based on their profile, they can start chatting or sharing files. This phase may also be used to add more members to the group by sending a connection request to the server designated in negotiation when the request is sent to connect. OffAT also allows exchanging handwritten notes as shown in Figure 4. Multihop communication can be facilitated by configuring a device as a client as well as a group owner where the device alternates roles based on time-sharing.

4.4. Security. Wireless networks are vulnerable to eavesdropping and intrinsically exposed to both active and passive security attacks. This presents a challenge in the incorporation of key-based cryptographic mechanisms for ad hoc networks because there is no trusted authority to provide certification or a centralized key distributor. Since OffAT mobile devices use Wi-Fi Direct, the security protocol implements WPA2 (Wi-Fi Protected Access II)

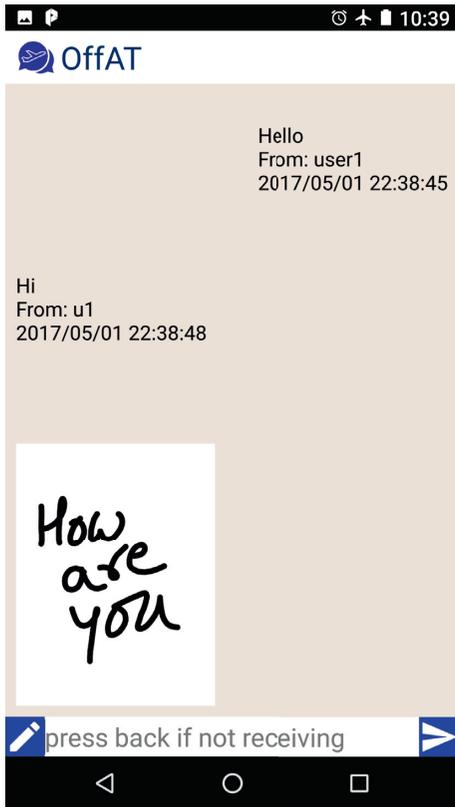


FIGURE 4: Exchange of text and handwritten messages in the airplane mode.

which offers a stronger encryption algorithm known as AES (Advanced Encryption Standard). Mobile devices implement WPA2 (Wi-Fi Protected Access II) as guided by the Wi-Fi Alliance's Wi-Fi Direct certification. With stronger encryption than TKIP (Temporal Key Integrity Protocol), WPA2 also provides AES (Advanced Encryption Standard) support courtesy of the 64-digit hexadecimal keys (preshared key (PSK)). Communication privacy and security are therefore provided through WPA2 and Wi-Fi Direct.

This implementation enables users to maintain privacy because the users can share user ID and limited interests without disclosing their e-mail ID or phone number unless they want to explicitly and are at liberty to decline linking requests (invitations) from other nearby users or disconnect. Furthermore, the connection requests are automatically declined after 30 seconds when the request is not responded. This preventive mechanism further protects data integrity by authenticating sources and sequencing of messages.

5. Results

The app, available at Google Play Store [5], is developed for the Android platform to simulate the proposed ASN. The screenshots shown in this paper are from two devices (Nexus 6P) with one device in the airplane mode and the other using the telecommunication network and also connected to a Wi-Fi router in addition to participating in the ad hoc social

network. During implementation and testing, it was observed that switching on Wi-Fi is required to participate in Wi-Fi Direct communication. Therefore, there is a limitation in the airplane mode to manually switch on the Wi-Fi without even being connected to any centralized infrastructure. Although Android allows activating Wi-Fi programmatically, testing results with users indicated that this is not better since it sometimes delays the process of reconfiguring device ID. Furthermore, asking users to manually switch on Wi-Fi provides the time that an app needs to modify the device ID with user ID and interests.

The users are able to exchange text and handwritten notes among other nearby similar users without any centralized network infrastructure. Implementation results and feedback posted at Google Play Store indicated user preferences towards location-based social networking. There are few suggestions by users to fix user name and interests, and they can be configured by Shared Preferences in the Android. Future work may include to automatically exchange hello packets with nearby or connected users for displaying them as online/offline.

Data Availability

The research paper proposes a mobile app for location-based social networking. The app is available at Google Play Store, and link is available in the references.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

References

- [1] Y. Wang, L. Wei, A. V. Vasilakos, and Q. Jin, "Device-to-device based mobile social networking in proximity (MSNP) on smartphones: framework, challenges and prototype," *Future Generation Computer Systems*, vol. 74, pp. 241–253, 2017.
- [2] H. Li, K. Bok, and J. Yoo, "P2P based social network over mobile ad-hoc networks," *IEICE Transactions on Information and Systems*, vol. 97, no. 3, pp. 597–600, 2014.
- [3] S. Gambhir and N. Aneja, "Ad hoc social network: a comprehensive survey," *International Journal of Scientific and Engineering Research*, vol. 4, no. 8, pp. 156–160, 2013, <http://www.ijser.org/researchpaper%5CAad-hoc-Social-Network-A-Comprehensive-Survey.pdf>.
- [4] B. Guidi, M. Conti, and L. Ricci, "P2P architectures for distributed online social networks," in *Proceedings of International Conference on High Performance Computing and Simulation (HPCS)*, pp. 678–681, IEEE, Helsinki, Finland, July 2013.
- [5] N. Aneja, "OffAT-chat in aeroplane mode," 2016, <http://offat.w3decode.com/>.
- [6] N. Aneja and S. Gambhir, "Geo-social semantic profile matching algorithm for dynamic interests in aAd hoc social network," in *Proceedings of IEEE International Conference on Computational Intelligence Communication Technology*, pp. 354–358, Ghaziabad, India, February 2015.
- [7] Z. Mao, J. Ma, Y. Jiang, and B. Yao, "Performance evaluation of WiFi Direct for data dissemination in mobile social networks," in *Proceedings of Computers and*

- Communications (ISCC)*, pp. 1213–1218, IEEE, New Delhi, India, July 2017.
- [8] N. Eagle and A. Pentland, “Social serendipity: mobilizing social software,” *IEEE Pervasive Computing*, vol. 4, no. 2, pp. 28–34, 2005.
- [9] L. Zhang, X. Ding, Z. Wan, M. Gu, and X.-Y. Li, “WiFace: a secure geosocial networking system using wifi-based multi-hop MANET,” in *Proceedings of 1st ACM Workshop on Mobile Cloud Computing and Services: Social Networks and Beyond-MCS’10*, p. 3, ACM, San Francisco, CA, USA, June 2010.
- [10] D. Zhang, D. Zhang, H. Xiong, C.-H. Hsu, and A. V. Vasilakos, “BASA: building mobile Ad hoc social networks on top of android,” *IEEE Network*, vol. 28, no. 1, pp. 4–9, 2014.
- [11] Y. Wang, A. V. Vasilakos, Q. Jin, and J. Ma, “A wi-fi direct based p2p application prototype for mobile social networking in proximity (MSNP),” in *Proceedings of 12th International Conference on Dependable, Autonomic and Secure Computing (DASC)*, pp. 283–288, IEEE, Dalian, China, June 2014.
- [12] N. Aneja and S. Gambhir, “Middleware architecture for aAd hoc social network,” *Research Journal of Applied Sciences, Engineering and Technology*, vol. 13, no. 9, pp. 690–695, 2016.
- [13] P. Bellavista and C. Giannelli, “Middleware for semantic multicast in spontaneous multi-hop networks,” in *Proceedings of International Conference on Mobile Wireless Middleware, Operating Systems, and Applications*, pp. 45–61, Springer, Berlin, Germany, November 2012.
- [14] D. Bottazzi, R. Montanari, and A. Toninelli, “Context-aware middleware for anytime, anywhere social networks,” *IEEE Intelligent Systems*, vol. 22, no. 5, pp. 23–32, 2007.
- [15] M. A. Rahman and M. S. Hossain, “A location-based mobile crowdsensing framework supporting a massive ad hoc social network environment,” *IEEE Communications Magazine*, vol. 55, no. 3, pp. 76–85, 2017.
- [16] A. R. Ilkhechi, I. Korpeoglu, U. Gdkbay, and . Ulusoy, “PETAL: a fully distributed location service for wireless ad hoc networks,” *Journal of Network and Computer Applications*, vol. 83, pp. 1–11, 2017.
- [17] J. Shu, S. Kosta, R. Zheng, and P. Hui, “Talk2Me: a framework for device-to-device augmented reality social network,” in *Proceedings of International Conference on Pervasive Computing and Communications (PerCom)*, Athens, Greece, December, 2018, <http://www.cse.ust.hk/~panhui/papers/talk2me.percom18.pdf>.
- [18] C. E. Casetti, C. F. Chiasserini, Y. Duan, P. Giaccone, and A. P. Manriquez, “Data connectivity and smart group formation in wi-fi direct multi-group networks,” *IEEE Transactions on Network and Service Management*, vol. 15, no. 1, pp. 245–259, 2018.
- [19] J. Joy, E. Chung, Z. Yuan, J. Li, L. Zou, and M. Gerla, “DiscoverFriends: secure social network communication in mobile ad hoc networks,” *Wireless Communications and Mobile Computing*, vol. 16, no. 11, pp. 1401–1413, 2016.
- [20] V. Smailovic and V. Podobnik, “Bfriend: context-aware ad hoc social networking for mobile users,” in *Proceedings of 35th International Convention MIPRO*, pp. 612–617, IEEE, Piscataway, NJ, USA, May, 2012.
- [21] V. Smailović and V. Podobnik, “BeFriend: a context-aware aAd hoc social networking platform,” *Automatika*, vol. 57, no. 1, pp. 58–65, 2016.
- [22] A.-K. Pietiläinen, E. Oliver, J. LeBrun, G. Varghese, and C. Diot, “MobiClique: middleware for mobile social networking,” in *Proceedings of 2nd ACM workshop on Online Social Networks*, pp. 49–54, ACM, Barcelona, Spain, August, 2009.
- [23] IEEE, *IEEE Standard for Information Technology-Telecommunications and Information Exchange between Systems-Local and Metropolitan Area Networks-Specific Requirements-Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications*, IEEE, Piscataway, NJ, USA, 2007.
- [24] J. Lee and C. S. Hong, “A mechanism for building Ad hoc social network based on user’s interest,” in *Proceedings of 13th Asia-Pacific Network Operations and Management Symposium*, pp. 1–4, Taipei Taiwan, September 2011.
- [25] S. Gambhir, N. Aneja, and L. C. De Silva, “Piecewise maximal similarity for Ad hoc social networks,” *Wireless Personal Communications*, vol. 97, no. 3, pp. 3519–3529, 2017.

Research Article

CEnsLoc: Infrastructure-Less Indoor Localization Methodology Using GMM Clustering-Based Classification Ensembles

Beenish Ayesha Akram ¹, Ali Hammad Akbar ¹ and Ki-Hyung Kim ²

¹Department of Computer Science and Engineering, University of Engineering and Technology, Lahore, Pakistan

²Department of Computer Engineering, Graduate School, Ajou University, Suwon, Republic of Korea

Correspondence should be addressed to Beenish Ayesha Akram; beenish.ayesha.akram@uet.edu.pk

Received 7 June 2018; Accepted 19 August 2018; Published 1 October 2018

Academic Editor: Subramaniam Ganesan

Copyright © 2018 Beenish Ayesha Akram et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Indoor localization has continued to garner interest over the last decade or so, due to the fact that its realization remains a challenge. Fingerprinting-based systems are exciting because these embody signal propagation-related information intrinsically as compared to radio propagation models. Wi-Fi (an RF technology) is best suited for indoor localization because it is so widely deployed that literally, no additional infrastructure is required. Since location-based services depend on the fingerprints acquired through the underlying technology, smart mechanisms such as machine learning are increasingly being incorporated to extract intelligible information. We propose CE_{ns}Loc, a new easy to train-and-deploy Wi-Fi localization methodology established on GMM clustering and Random Forest Ensembles (RFEs). Principal component analysis was applied for dimension reduction of raw data. Conducted experimentation demonstrates that it provides 97% accuracy for room prediction. However, artificial neural networks, k -nearest neighbors, K^* , FURIA, and DeepLearning4J-based localization solutions provided mean 85%, 91%, 90%, 92%, and 73% accuracy on our collected real-world dataset, respectively. It delivers high room-level accuracy with negligible response time, making it viable and befitting for real-time applications.

1. Introduction

Positioning systems aka localization systems both for outside and inside buildings is an ever-exciting area of research and development due to increasing market shares as in smart buildings, assistive and assisted living, safer metropolitans using geographical information systems, and tracking of IoT objects for commercial purposes. User localization is poised to reach 2.6 billion dollars' worth of market share soon [1], specially involving indoor localization solutions. Localization for the outdoors has been successfully commercialized in the form of satellite-based technologies such as GPS, BeiDou, GLONASS, COMPASS, and GALILEO [2]. Indoor positioning cannot be performed using the same technologies because of No-line-of-sight (NLOS) and occlusion. Radio frequency (RF) signals, on the contrary, do not require explicit LOS for operation.

A received signal strength indicator (RSSI) as a measure of RF signal quality for Wi-Fi, RFID, Bluetooth [3], ZigBee

[4] and ultrawideband [5] has been used in several indoor positioning systems. Moreover, several kinds of sensory input such as images [6], video, ambient sound [7], accelerometer [8], magnetometer [9], pedometer, gyroscope readings, and their amalgamation with various aforesaid RF signals [10] have also been explored.

Several approaches and their combinations such as time of arrival (TOA), time difference of arrival (TDOA), pedestrian dead reckoning (PDR), and angle of arrival (AOA) [11] have also been utilized for indoors. Each of these has been shown to have serious limitations. For instance, at successive predictions, PDR suffers from error propagation, and precise clock synchronization between sender and receiver is a requirement for TOA- and TDOA-based systems. In addition, specialized antennas are needed for AOA-based systems.

RSSI-based systems are based on two widely adopted mechanisms; location estimation using formalized propagation models or fingerprints. Location systems that are based on the former suffer from low precision at run time

because of variability in channel behavior including fading and shadowing, and also due to heterogeneity of device types and form factors [12].

Wi-Fi networks are deployed as in Access Points (APs) that are prevalent everywhere. Utilizing these to capture RSSI fingerprint (FP) is conveniently possible with device as simple as smartphones or phablets. Wi-Fi fingerprint-based localization has the following benefits: no requirement of extra hardware at both sender and receiver sides; utilization of already existent infrastructure; easily implementable; and no essential need of propagation model building which may or may not depict real signal propagation at run time [13].

The infrastructure of APs allows us to collect a dataset comprising FPs on selected Reference Points (RPs) that essentially becomes a cue to the physical layout of the building, like a map. It is then utilized to prepare the localization system as in the training phase. Once trained, the system is ready to be used, i.e., for an unseen FP captured anew, a room number or an associated label is returned by the system to estimate the location.

In this paper, a new localization methodology is presented that uses a combination of data reduction technique as in principal component analysis (PCA), soft clustering technique such as Gaussian mixture model (GMM), and bootstrapped aggregated/bagged ensembles of decision trees commonly referred to as Random Forest Ensembles (RFEs). We aspire to provide infrastructure-less indoor localization methodology which is scalable, easy to deploy, and provides real-time response for high room-level accuracy instead of explicit coordinates. First of all, PCA was employed for raw data dimension reduction. Then, clustering was performed to split the data in similar groups to help classifier better learn the data dynamics. Finally, a separate RFE is trained for every single cluster.

The remainder of the paper has been organized as follows. Section 2 presents related work. Preliminary experimentation results are summarized in Section 3. Section 4 provides details of the localization methodology that we have proposed. Section 5 delves into experimental setup, layout and results for validation. Conclusion along with possible future directions is showcased in Section 6.

1.1. Scope of the Study. It is important to declare the scope of the study here to give a prelude to the paper that follows. The paper is aimed at proposing a new methodology that is put to application in a practical and particular environment. The resulting constraints and limitations of our experimental regime are quite natural and fairly generalizable in terms of spatiotemporal aspects, specifically for building size and type, fingerprints' collecting device type, and time of the year. The fingerprints were collected at the Software Engineering Centre of our own university. The building is double-storey and has architectural diversity in terms of rooms, corridors, and an inlaying garden. Nonetheless, an extensive dataset spanned over multiple buildings can be obtained to ensure wider spatial diversity. Our dataset was collected using a single Android phone; however, with a large team using a variety of devices, data collection can be

performed over different times of the year to obtain data both for training and testing of our approach.

2. Related Work

We summarize here the work on IPS based on Wi-Fi that is typically a WLAN standard IEEE802.11 (b, a, g, ac, or any) or a combination of Wi-Fi with another wireless or sensory input. RADAR [14] from Microsoft® labs is the pioneer research work to employ Wi-Fi signals. Wi-Fi signals received at the base stations (Access Points) from a laptop were used for predicting the user's coordinates using k -NN-based method and triangulation, reporting a median error of 2-3 m. Li et al. [12] combined affinity propagation as a message passing-based clustering algorithm with PSA-based artificial neural network (ANN) for the Cartesian coordinates-based prediction. A mean error as low as 1.89 m was reported by them including 2.9 m for 90% of estimates.

Song et al. [13] analyzed FP collection as an AP relevancy problem. Hidden Naïve Bayes (HNB) was used as a mechanism to infer the most relevant APs and suggested that redundant APs may be obviated for each RP through a variant of ReliefF with the Pearson product-moment correlation coefficient (PPMCC). Moreover, clustering was performed on RPs. One HNB was trained per cluster to approximate user coordinates. Cooper et al. [15] employed combination of Wi-Fi and Bluetooth low energy radio signal FPs based on boosting technique targeting the room-level prediction. They trained one classifier per room based on a variant of AdaBoost that conveniently harnessed decision stumps in one-vs-all notion. Using combination of Wi-Fi + BLE, they acquired 96% accuracy with $4.3E-03$ seconds response time. Wang et al. [16], similar to Li et al., presented training of ANN with back propagation that was based on PSA for RSSI measurements of RFID tags. They performed data preprocessing by normalizing dataset to $[0, 1]$ range along with using Gaussian filter. They estimated x and y coordinates reporting a mean error of 0.34 m.

Xu et al. [17] utilized multilayer neural network (MLNN) for Wi-Fi signals along with network boosting. They trained the MLNN in two stages commonly followed in deep learning, namely, pretraining using autoencoders and fine tuning using back propagation algorithm reporting a mean error of 1.09 m. Zhang et al. [18] proposed a coarse localizer composed of four-layered deep neural network using stacked denoising autoencoders, succeeded by the hidden Markov model-based fine localizer, reporting a mean error of 0.39 m.

Calderoni et al. [19] utilized RFID tags' RSSI values targeting room-level accuracy in a hospital environment. They divided the total area into macroregions using k -means variant, followed by a Random Forest trained per macroregion. Multiple random forests for whom the cluster matching score was greater than a particular threshold determined the final prediction with 83% reported accuracy. Jedari et al. [20] investigated room-level prediction using k -NN, rule-based JRip, and Random Forest classifier based on Wi-Fi signals. They concluded that Random Forest produced much better results than k -NN (77.4%) and JRip (72.2%) with 91.3% accuracy.

Mo et al. [21] proposed the usage of kernel PCA (KPCA) algorithm for the coarse-level prediction of manually labelled cluster using Random Forest. They derived trained matrices from extracted KPCA features and prepared sub-radio maps. For prediction, the features extracted from coarse positioning, refined by the trained KPCA matrices were fed to weighted k -NN (WK-NN) for final coordinates estimation. They reported an accuracy of 93% with an error distance of 2 m. Górak et al. [22] employed Random Forest for finding important APs and applied threshold-based elimination. They determined malfunctioning APs during operation based on important APs. They were able to report error rates as low as 4% for detecting the floor, and as for horizontal detection, 2 m error was reported. The results were compared against 30% and 7 m, respectively, when malfunctioning AP were undetected. They [23] divided FP dataset both into subsets which were either overlapping and/or nonoverlapping as per the presence of RSSI from every AP. Furthermore, one Random Forest was trained for each such subset. They compared the results with base Random Forest, signifying around 5% to 9% betterment in average reported error at the floor level. The performance in terms of detection of floors was unchanged.

Aforementioned are some recent efforts in field of indoor localization using RF signals. k -NN works by storing all the data samples along with marked ground truth labels. An unseen sample is compared with complete dataset based on a similarity measure/distance to determine nearest neighbors, where k is the number of neighbors and neighbors' weightage. The final decision of the sample's class is based on the majority of the labels of the k -nearest neighbors. Several IPS such as in Oussalah et al. [11] and Niu et al. [24] based on k -NN and its variants do not scale well when the dataset grows because they require FP matching with whole dataset. Moreover, recently artificial neural networks and deep learning have gained a great deal of focus for indoor localization. Artificial neural networks try to mimic human brain. They form multiple layers of neurons, namely, a single input layer, a single output layer, and one or more hidden layers. At every layer, several neurons are connected to one another according to a specific configuration and triggering function. Every layer affects and triggers neurons in the next layer following rules of the learning function which eventually evolves into the final output. Heavy resource utilization is required during the training phase; however, their response time is negligible due to minimal required computation. IPS employing ANNs and deep learning such as Li et al. [12], Ding et al. [25], León et al. [26], Zhang et al. [18] and Tuncer and Tuncer [27] faces the challenge of finding several tunable parameters such as the optimal architecture, no. of layers, learning function, and no. of neurons at every layer. Moreover, the convergence rate and the final accuracy of various configurations do not follow any specific trend. Sometimes, a 2-layer simpler configuration takes more time to converge than a 4-layer network, and the accuracy achieved by a 6-layer network is lower than the accuracy obtained by a 3-layer ANN. Hence, heuristics are predominantly used for proposing an architecture based on ANNs because Monte Carlo configuration testing is

impossible. An AP location change, or the addition/removal of new APs, will lead to essential retraining of the IPS putting ANN and deep learning at a disadvantage.

Inspired from existing works, we suggest a clustering-based multiclass classifier approach for room-level prediction which is easy to train-and-deploy in terms of computational complexity, provides suitable accuracy, and offers response time appropriate for real-time applications. Clustering follows the divide and conquer approach to help the classifier better learn the group of similar observations instead of the whole dataset. We perform soft classification (clustering) of FPs, not for RPs which has been mostly done in the existing works. PCA-based data dimension reduction helps us to reduce the response time of the system by decreasing the number of predictors.

This approach scales well in terms of response time with an increase in the number of rooms since a single classifier is invoked for each location providing maximum accuracy reported so far to the best of our knowledge.

3. Preliminary Experiments

Preliminary experiments were performed on sample dataset collected from our departmental building to evaluate classifier suitability [28]. The results presented here are on the sample dataset. The detailed experimentation results on the complete dataset of all locations in building are presented in Section 5.

3.1. Dataset Acquisition. A customized app was developed for an Android phone, built to record RSSI as vector data coming from Wi-Fi APs. The Wi-Fi FPs of all observable APs both within 2.4 and 5 GHz bands at a RP were scanned using a commercial off-the-shelf Samsung phone (Version: J5 Galaxy). FPs were obtained at each RP while hovering the smart phone starting at 0° up to right angle (90°) with respect to the floor as shown in Figure 1, so as to make the phone face N, NW, W, SW, S, SE, E, and NE with an effort to keep occlusion as minimum as possible due to the human torso [5, 9]. The user stood at the centre of the RP, held the phone used for FP collection in his hand, and captured multiple FPs in each direction out of the total 8 directions shown in Figure 1. A total of A APs were present in the premises continuously emitting radio signals. These FPs were then stored into DB, each with respective room labels.

3.2. Preprocessing. The resultant dataset was found to be sparsely populated with the identifiers of APs because of the presence/absence of these APs at various RPs labelled in rooms and in the corridors of building. The measured RSSI values varied between -98 dBm and -15 dBm (from being weak to being strong as a result of distance from APs). Being consistent with the well-known practice to keep the missing values slightly weaker than the weakest signal detected in the dataset [12, 14, 18, 19], the missing values were replaced with -100 dBm.

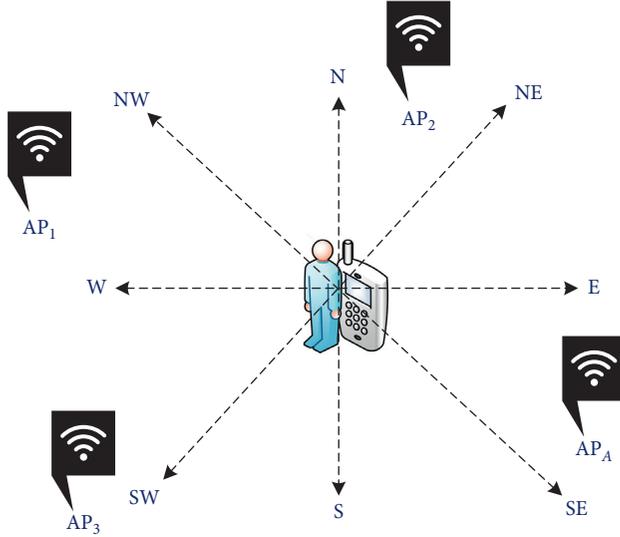


FIGURE 1: User orientation during data collection.

3.3. Classifier Evaluation. We evaluated performance of 60+ classifiers in WEKA for room-level prediction on sample dataset out of which top ten best performing classifier performances are summarized in Table 1.

Taking into account all the performance measures such as accuracy to receiver operating characteristics (ROC) area, the overall performance best attained (descending order) was by K^* , k -nearest neighbors (k -NN), Random Forest Ensemble (RFE), and algorithm for Fuzzy Unordered Rule Induction Algorithm (FURIA), multilayer perceptron, deep learning for JAVA (Dl4jMlpClassifier), support vector machines, Naive Bayes classifier, and finally AdaBoost that uses stumps for decision-making. K^* and k -NN both are instance-based classifiers and produced similar performance results. Followed by RFE, FURIA, and multilayer perceptron (ANN), FURIA and ANN show almost similar performance trend. We selected k -NN and ANN for comparison as many existing works had utilized these which makes comparison with other works easier. Moreover, we chose K^* , FURIA, and DeepLearning4J classifiers for comparison too, as they are state-of-the-art and relatively new machine learning methods. RFE is suited for large datasets to give high accuracy and time efficiency. It is resilient to noise in data and is also capable of dealing with missing values in data. RFE utilizes bootstrapping that reduces variance and keeps the bias in check because creating different subsets of training dataset along with a replacement mechanism ensures that the trees have little or no correlation. Therefore, overfitting is avoided, making it more generalizable. Both training and prediction time of RFE due to parallel computation supported by bagging make it suited for real-time implementation of IPS. Hence, we selected RFE [29] as the suitable classifier module in our proposed methodology.

4. Proposed Localization Methodology

4.1. Problem Formulation. We assume localization/positioning as a combination of clustering and multiclass

classification problem where each room is considered to be a class. A two-dimensional indoor area is partitioned into R square grids of dimensions $C \times D$ m². The centre of each square grid is a reference point (RP). A device equipped with the wireless adapter card can sense wireless signals from a total of A AP-s at a certain RP at a given time, which forms the fingerprint $FP_i = \{RV_i, L_i\}$. $RV_i = \{rssi_{i1}, rssi_{i2}, rssi_{i3}, \dots, rssi_{iA}\}$ where $rssi_{ij}$ symbolizes the RSSI value from j th AP (dBm) in the i th sample of FP collected and L_i is the respective class/room label. Let N such FP constitute the dataset. Localization function, LF, is learnt from the FP dataset to map the observed FP to a certain room label L_x as described by the following equation:

$$L_x = LF(FP_i). \quad (1)$$

There are two phases of the proposed methodology (CEnsLoc); namely, training phase and prediction phase as shown in Figure 2. First of all, a sparse Wi-Fi FP dataset with many missing values is collected. In the training phase, the collected FPs reserved for training the system are pre-processed for missing values replacement, followed by PCA application for dimension reduction. Then GMM-based hard clustering is used for nonoverlapping/disjoint data subsets generation. For each such subset, a separate RFE is trained for location prediction and stored in the database. In the location prediction phase, same steps of missing value replacement and PCA computation are performed on the captured FP. Then, FP is matched with a single cluster using the stored GMM. The final prediction L_x is generated by invoking the respective pretrained RFE for the best matched cluster/subset. These phases are formally elaborated in detail in Sections 4.2 and 4.3 respectively.

4.2. Training Phase. During training, the training dataset is fed to the preprocessing module which replaces empty readings with missing value replacement (MV_r). PCA is then performed on the dataset for dimension reduction. Orthogonal transformation is applied by PCA for redundant information removal to decrease the number of predictors. Principal components (PCs) were obtained by applying PCA, which are a set of linearly uncorrelated variables such that maximum variance by some projection of the data is captured by the first principal component and so on. Choosing the smaller of the number of predictors/APs and number of samples minus one, A PCs are generated $\{PC_1, PC_2, \dots, PC_A\}$. For computation of PCs, first the mean RSSI value of each AP is subtracted from the i th RP using the following equation:

$$X_i = \frac{1}{N} \sum_N FP_i - \frac{1}{R} \frac{1}{N} \sum_R \sum_N FP_i, \quad (2)$$

where N is the total no. of rows of samples/dataset and R is the total no. of RPs. The PCs matrix is computed by the following equation:

$$PC_i = X_i \times E_A, \quad (3)$$

where E_A is the eigenvector matrix of the average RSSI value of each AP in the i th RP. The resulting dataset is then divided

TABLE 1: A comparison of performance measures of various classification algorithms for Wi-Fi-based position estimation.

Algorithm	Time to build model (sec)	Accuracy	Kappa statistic	RMSE	Precision	Recall	F1	MCC	ROC area
K^*	0	99.52	0.98	0.02	0.99	0.99	0.99	0.99	1
k -NN	0	99.06	0.99	0.03	0.99	0.99	0.99	0.98	0.99
RFE	1.11	98.76	0.98	0.04	0.98	0.98	0.98	0.98	1
FURIA	5.92	97.26	0.96	0.05	0.97	0.97	0.97	0.96	0.99
Multilayer perceptron	25.84	97.05	0.96	0.06	0.97	0.97	0.97	0.96	0.99
J48	0.1	95.91	0.95	0.08	0.95	0.95	0.95	0.95	0.98
D4jMlp classifier	26.31	94	0.93	0.08	0.94	0.94	0.94	0.93	0.99
SVM	5.82	90.6	0.89	0.12	0.93	0.90	0.90	0.9	0.94
Naive Bayes	0.02	89.79	0.88	0.12	0.91	0.89	0.89	0.89	0.99
AdaBoost with decision stump	0.1	36.81	0.22	0.25	0.15	0.36	0.21	0.18	0.76

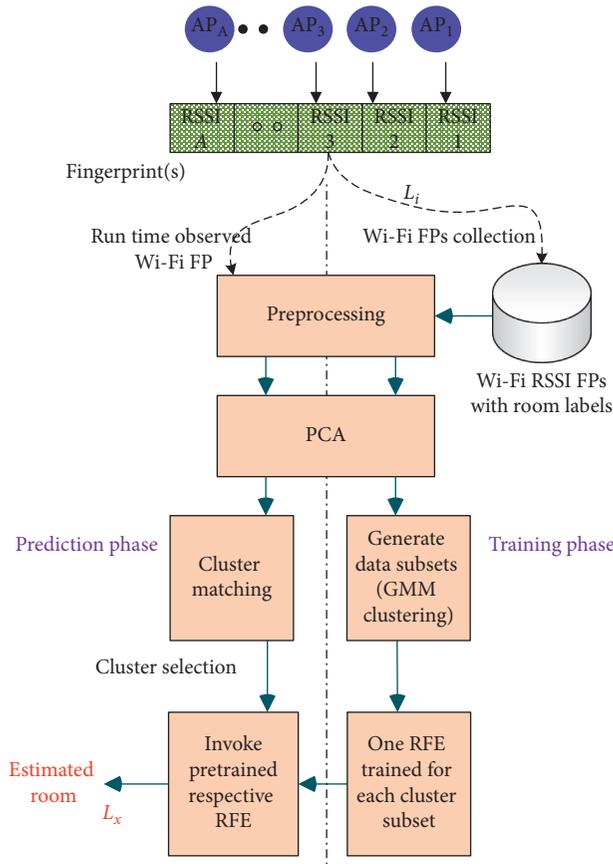


FIGURE 2: Proposed localization methodology (CENSLoc).

into subsets using GMM clustering (K data subsets). Equation (4) is a distribution based on 2D Gaussian where mean is represented by μ and the covariance matrix is Σ . A GMM with N as the no. of overlapping distributions is described by the following equations:

$$N(x | \mu, \Sigma) = \frac{1}{2\pi\sqrt{|\Sigma|}} \exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right\}, \quad (4)$$

$$P(x) = \sum_{k=1}^N \pi_k N(x | \mu_k, \Sigma_k), \quad (5)$$

$$\sum_{k=1}^N \pi_k = 1, \quad (6)$$

where π_k defines the mixing coefficient to express the weight of each mixing element (weighted sum being 1). The resultant shape of 2D Gaussian is the average of distributions individually, in terms of covariance and mixing coefficients. Assuming that a linear mix of weighted coefficients for each of the respective distribution's average and covariance is obtained, and by incorporating sufficient distributions, a final density function may be obtained. The reason behind GMM clustering was the similarity between Gaussian distribution and the radio propagation characteristics of a Wi-Fi AP [30] which makes GMM a highly suitable candidate for clustering Wi-Fi RSSI vectors.

Furthermore, each data subset is used for training a RFE to predict the room label as a multiclass classifier. The trained models of GMM as well as all RFEs are stored for later use in the prediction phase.

The algorithm for training and prediction phases of CENSLoc is given in Algorithm 1.

Let the total number of samples be N_{sample} in the training set, the number of trees is N_{tree} , the number of maximum splits allowed is S_{max} , the number of predictors for the classifier is A , f is the value specifying no. of input predictors that are utilized to split at a tree node, and tc denotes the total no. of classes in the dataset. For finding best split, RFE uses Gini Index as given in the following equation where P_j is the class j 's relative frequency in N_{sample} :

$$\text{Gini}(N_{\text{sample}}) = 1 - \sum_{j=1}^{tc} (P_j)^2. \quad (7)$$

RFE is trained using the method presented in Algorithm 2.

4.3. Prediction Phase. The average of collected RSSI FP is fed to CENSLoc, A PCs are computed by the method similar to the training phase. The saved model of GMM is invoked for cluster matching. Matched cluster's trained RFE is invoked where the final decision is computed by the majority vote described in Algorithm 2.

4.4. Time Complexity of Training and Prediction Phases for CENSLoc. Ceteris paribus, the time complexity of the training phase and the prediction phase for CENSLoc is essentially dependent upon the size of the experimental area,

```

Input: training dataset with total  $A$  predictors
Missing value replacement  $MV_r$ 
Maximum number of clusters  $K_{\max}$ 
Output: predicted location  $L_x$ 
For training:
Replace empty values with  $MV_r$ 
Apply PCA on dataset to generate  $A'$  predictors
For  $k=1 \rightarrow K_{\max}$ 
  Generate clusters
  Generate and save  $k$  data subsets
  For each  $p \in k$  data subsets
    Train  $p$  RFE using Algorithm 2 (training)
    Calculate performance measures
  End for
End for
Choose optimal configuration
Save respective models for GMM, all RFEs
For prediction at a new point  $x$ :
Replace missing values with  $MV_r$ 
Apply PCA on the FP
Match one cluster  $C_{\text{match}}$ 
Invoke RFE of  $C_{\text{match}}$  using Algorithm 2 (prediction)

```

ALGORITHM 1: CEnsLoc training and prediction algorithm.

our acquisition regimen, the resultant dataset of FPs, and how the dataset is manipulated by the tandem of schemes we employ.

4.4.1. Time Complexity of Training. In the training phase for PCA, time complexity is governed by the following equation:

$$O(\min(A^3, N^3)), \quad (8)$$

where A = no. of predictors and N = no. of observations.

For training a decision tree (DT) that has not been pruned, the expression is as follows:

$$O(A \times N \log(N)). \quad (9)$$

As RFE consists of numerous DTs, merely a small no. f is used from total predictors A . Complexity for a single DT in RFE is represented by Equation (10) and the complexity of N_{tree} by Equation (11):

$$O(f \times N \log(N)), \quad (10)$$

$$O(N_{\text{tree}} \times f \times N \log(N)), \quad (11)$$

where N_{tree} = no. of trees in RFE and f = random features that are chosen to get the best split.

While trying to control trees' depth grown using S_{\max} , training complexity of one RFE is as follows:

$$O(N_{\text{tree}} \times f \times N \times S_{\max}). \quad (12)$$

Since K ensembles are grown for predicting room level, time to train complexity is represented by the following equation:

$$O(N_{\text{tree}} \times f \times N \times S_{\max} \times K). \quad (13)$$

For GMM, the complexity is expressed by the following equation:

$$O(N \times K \times D^3), \quad (14)$$

where N = no. of samples, K = no. of components, and D = no. of dimensions.

Incorporating all, the time complexity to train for CEnsLoc is as follows:

$$O(\min(A^3, N^3)) + O(N \times K \times D^3) + O(N_{\text{tree}} \times f \times N \times S_{\max} \times K). \quad (15)$$

4.4.2. Time Complexity of Prediction. For prediction, time complexity for PCA is given by the following equation:

$$O(\min(A^3, N^3)). \quad (16)$$

The complexity for a DT and an ensemble in terms of prediction time are shown by Equations (17) and (18), respectively:

$$O(N \log(N)), \quad (17)$$

$$O(N_{\text{tree}} \times N \log(N)). \quad (18)$$

S_{\max} controls trees depth; therefore, complexity for an ensemble is given by the following equation:

$$O(N_{\text{tree}} \times N \times S_{\max}). \quad (19)$$

Only a single RFE out of K ensembles gets invoked for predicting a room.

Therefore, time complexity for GMM is given by the following equation:

$$O(K \times D^3). \quad (20)$$

Finally, the overall complexity for CEnsLoc in terms of prediction is as follows:

$$O(\min(A^3, N^3)) + O(K \times D^3) + O(N_{\text{tree}} \times N \times S_{\max}). \quad (21)$$

5. Experimental Results and Discussion

This section entails hardware equipment, software used, and particulars of experiments that were used to evaluate the performance of CEnsLoc in light of accuracy, precision, recall, training, and response time.

An Intel machine (64-bit Xeon: X5650) with a master clock at 2.67 GHz with 24 GB RAM, and 64-bit Windows 10 Education was used for experimentation in MATLAB. The real dataset was developed through FP collection at the ground floor of Software Engineering (SE) Centre, University of Engineering and Technology (UET), Lahore Pakistan. The building's dimensions are 39 m × 31 m (1209 m²) containing offices, class rooms, laboratories, and open corridors. Figures 3 and 4 depict the building's floor plan, room labels (L1–L10 closed rooms, L12 open corridor, and L11 a semi-open room), a total of 180 RPs, and the

Input: data subset with total A' predictors for training

No. of tree N_{tree}

Allowed maximum no. of splits S_{max}

Random no. of predictors/features f

Output: estimated location L_x

For system training:

Step 1: for $l = 1$ to N_{tree}

(i) Choose a bootstrap sample set (SS) of size (N_{sample}) with replacement from the training data subset

(ii) Generate a Random Forest Tree (T_l) to SS, via recursively iterating (a-c) for every terminal node of tree, unless the maximum no. of splits (S_{max}) is reached

(a) Randomly select f features/variables from the A' predictors ($f \ll A'$)

(b) Choose the best features/split-point from the f employing Gini Index

(c) Split node forming into two children nodes

Step 2: Produce the resulting ensemble of trees $\{T_l\}_1^{N_{tree}}$

For location prediction at a new point x from RFE L_{rf} :

Let, $L_m(x)$ be the room/class prediction by the m^{th} RFE tree

$L_{rf}^{N_{tree}}(x) = \text{maj. vote } \{L_m(x)\}_1^{N_{tree}}$

ALGORITHM 2: RFE classification algorithm for training and prediction.

number of samples collected per room. RPs that are represented by small coloured dots in rooms in Figures 3 and 4 enlist the total number of samples collected at each location L1–L12. Such notion of rooms was used as walls playing a crucial role in fluctuation of Wi-Fi RSSI values [31, 32]. The area was planned into a grid of cells of $1.5 \times 1.5 \text{ m}^2$. Each cell center was marked as the Reference Point (RP) for FP collection.

The complete dataset consisted of 20087 Wi-Fi RSSI FPs, in which total 40 APs were detected. Figure 4 depicts the number of FPs captured in each room/location marked as L1–L12. It must be noted that all these APs belong to university infrastructure comprising of SE centre and its immediate neighboring buildings. The FPs were pre-processed following the same process described in Sections 3.1 and 3.2. PCA-based dimension reduction was employed with optimal results found with 23 PCs. The resultant dataset was divided with 70:30% stratified ratio for training and testing subsets. The experimental results are discussed using both 10-fold cross validation (10-CV) on training subset and on unseen 30% test subset. GMM clustering then partitioned the training data into further subsets, and the optimal configuration was two clusters, shared covariance kept as true and diagonal covariance. Furthermore, one RFE was trained for each subset with 132 trees, 1024 maximum splits, and 8 random features. CEEnLoc was hosted at a machine, and the mean of observed RSSI FPs was used for run time location estimation after applying the same model of PCA and GMM cluster matching finalized from the training phase. Best matched cluster's respective RFE was invoked for room prediction. Various companion apps can query location of subscribed users using the IPS configured as a server. Response time was computed as an average of the difference between localization query time and prediction generation time. The following formulae were used to compute the performance parameters:

$$\begin{aligned} \text{accuracy} &= \frac{TP + TN}{TP + TN + FP + FN} \\ \text{precision} &= \frac{TP}{TP + FP} \\ \text{recall} &= \frac{TP}{TP + FN} \end{aligned} \quad (22)$$

5.1. Classification Effectiveness and Efficiency. Tables 2 and 3 summarize the 10-CV performance evaluation of CEEnLoc, comparing it with k -NN [33], artificial neural network (ANN) [34], K^* [35], FURIA [36], and DeepLearning4J [36]. The best results are highlighted throughout with boldface. k -NN results were computed averaged over six different configurations based on the number of neighbors, similarity measure employed, and neighbor weightage. Similarly, six different configurations of ANN, namely, 2-, 3-, and 4-layer ANNs each having varying number of neurons per layer, employing two learning algorithms SCG and RBP, were averaged out to obtain the results. For K^* , the entropic blend percentage was varied from 10 to 90 percent. The results obtained for FURIA were obtained by varying the count of folds for growth and pruning as well as through varying minimum instances weight that was for each split. DeepLearning4J results were obtained by varying the number of neurons per layer, number of dense layers in the network, and training algorithm out of many possible variations in the hyper parameters. The results by CEEnLoc were generated by the found optimal configuration of our proposed approach. For all the approaches, the same set of configurations, as used during 10-CV performance evaluation, were used for result generation on 30% unseen stratified test data subset whose results are presented in Table 4. The best performance on both 10-CV results (Table 2) and test dataset (Table 4)

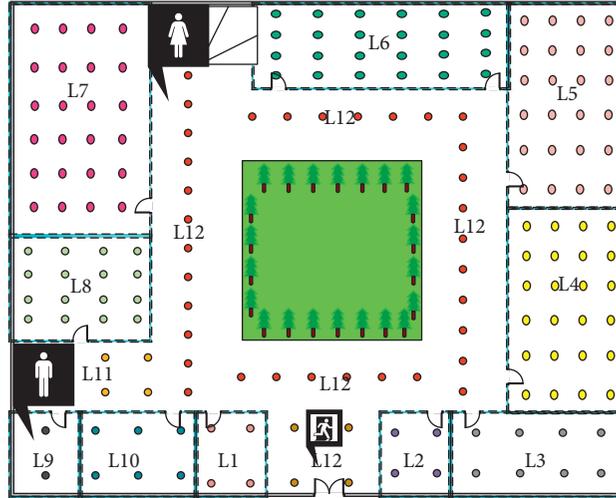


FIGURE 3: Room labels and marked RPs.

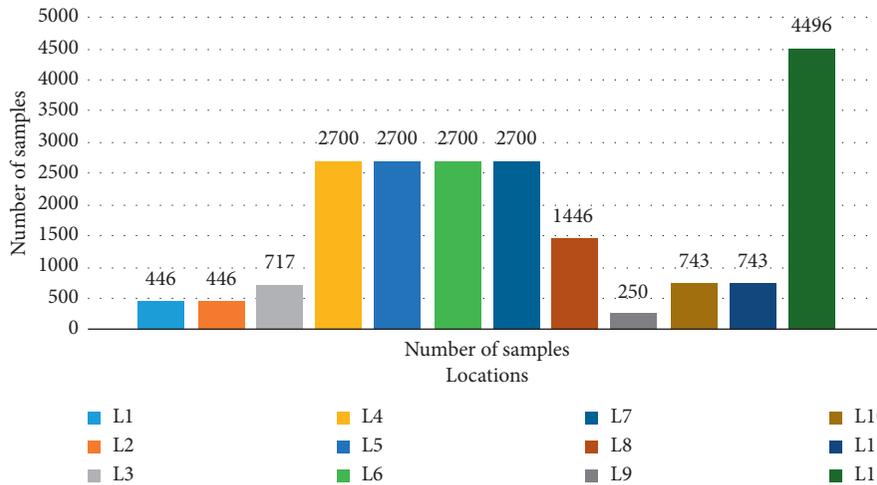


FIGURE 4: Number of samples per room.

were obtained by CEnsLoc with 97% and 95% accuracy followed by FURIA (92%, 90%), k -NN (91%, 89%), K^* (90%, 87%), ANN (85%, 82%), and DeepLearning4J (73%, 71%), respectively. The results are presented on both 10-CV and unseen test data subset to show that, for a small dataset, 10-CV results can provide a good performance estimate, as results on the test data subset were found to be showing similar trends as depicted by the 10-CV results. Although deep leaning and ANN provide good performance results, but finding optimal configuration for a huge number of tunable parameters is a tricky task. Moreover, during experiments, it was found that despite generic guidelines for parameter tuning, the performance measures and convergence rate of ANN as well as deep learning schemes highly fluctuate with even a slight variation in the number of neurons in layers, training algorithm, and number of hidden layers, which makes it harder to find optimal configuration for ANN and deep learning models. Lazy and instance-based approaches such as k -NN and K^* are also good candidates

for indoor localization; however, they have a very limited number of tunable parameters resulting in inability to surpass CEnsLoc performance despite trying different parameter combinations. They do not generalize well as compared to RFE as the end result of prediction is heavily dependent on the majority vote by k closest matched samples in the dataset. They also need template matching with the entire dataset for one location prediction which results in growing response time with increasing number of samples in the dataset which is highly likely in practical real-world scenarios. As in a typical building, there are a quite huge number of visible APs, and a sufficiently large number of samples are also required for classifiers to work properly.

A minimum response time of $2.05E-05$ seconds was obtained by FURIA. DeepLearning4J stood second with $6.82E-05$ seconds response time. ANN, k -NN, and CEnsLoc all had response time on the scale of $E-04$ seconds which is 10 times slower than aforementioned two

TABLE 2: Tenfold cross-validated performance evaluation and comparison of CEnsLoc.

	Accuracy	Precision	Recall
k -NN	0.91	0.88	0.87
ANN	0.85	0.83	0.78
K^*	0.90	0.96	0.62
FURIA	0.92	0.89	0.82
DeepLearning4J	0.73	0.51	0.41
CEnsLoc	0.97	0.98	0.97

TABLE 3: Tenfold cross-validated performance comparison of CEnsLoc in terms of training and response time (seconds).

Time (sec)	Training time (10-fold)	Avg. training time (1-fold)	Response time dataset	Response time (1 sample)
k -NN	—	—	1.97	$1.70E-04$
ANN	267.20	26.72	0.18	$1.41E-04$
K^*	—	—	103.09	$8.17E-02$
FURIA	158.7	15.8	0.03	$2.05E-05$
DeepLearning4J	755.59	75.5	0.09	$6.82E-05$
CEnsLoc	140.07	14.007	2.35	$2.08E-04$

TABLE 4: Performance evaluation and comparison of CEnsLoc on test subset.

	Accuracy	Precision	Recall
k -NN	0.89	0.87	0.85
ANN	0.82	0.81	0.76
K^*	0.87	0.94	0.60
FURIA	0.90	0.86	0.79
DeepLearning4J	0.71	0.48	0.39
CEnsLoc	0.95	0.96	0.94

approaches, but the difference was trifling, which cannot be detected by any human user of the system; the accuracy, precision, and recall provided by CEnsLoc was much greater than all other approaches.

FURIA stands out as the second best performer regarding indoor localization, which is an upgraded version of the RIPPER algorithm, indicating that rule-based algorithms, specially fuzzified versions with good generalization capabilities perform well for location estimation as well. However, it lagged behind CEnsLoc in terms of accuracy, precision, and recall by 5%, 10%, and 15%, respectively.

The details of both 10-CV and test dataset results for all IPS compared and CEnsLoc are depicted in Figure 5 for side by side visual comparison.

5.2. Out-of-Bag (OOB) Error Results. There is a performance measure called OOB error which is peculiar to RFE. During training of RFE, data subsets are generated with replacement, resulting in some repeated and left out observations (OOB observations) for each tree. That particular tree does not train on these left out OOB observations. Prediction capability of RFE can be measured using “OOB Loss” which is the error made on unseen OOB observations

during training. OOB loss measure has been investigated and shown to provide an upper bound on testing error [37], specifically, useful for small-sized datasets. Hence, OOB error can be used just like/instead of unseen test data subset if the available dataset is small or unseen test dataset is unavailable. It provides a very good estimate of the trained classifier’s generalization capability. Table 5 summarizes the OOB loss compared with averaged out 10-CV loss indicating that it indeed bounds it.

6. Conclusion and Future Work

Location prediction/estimation provides derivation of meaningful context for a broad range of services and applications. Indoor localization can open altogether a plethora of new opportunities because humans spend most of their time indoors. CEnsLoc offers shorter response time and an overall improvement in accuracy, precision, and recall. With only a few parameters to be tuned, it is suited for FP-based localization which requires frequent recollection of data and retraining. CEnsLoc was able to attain 97% accuracy in comparison with other IPS averaged over 6 different configurations, namely, FURIA, k -NN, K^* , ANN, and DeepLearning4J with 92%, 91%, 90%, 85%, and 73% accuracies, respectively. It can be utilized for elderly assistance, navigation, smart buildings, and smart transportation to name a few potential applications.

Our future work includes deployment of CEnsLoc across a wide range of civil infrastructures including different floors and/or buildings to understand its performance in more detail as well as its scalability using crowdsourcing. We also aim to build safety, security, and evacuation guide applications with CEnsLoc at their core for users in offices, universities, and retail. Furthermore, integration with GPS, Bluetooth, and PDR to further enhance accuracy, utilizing

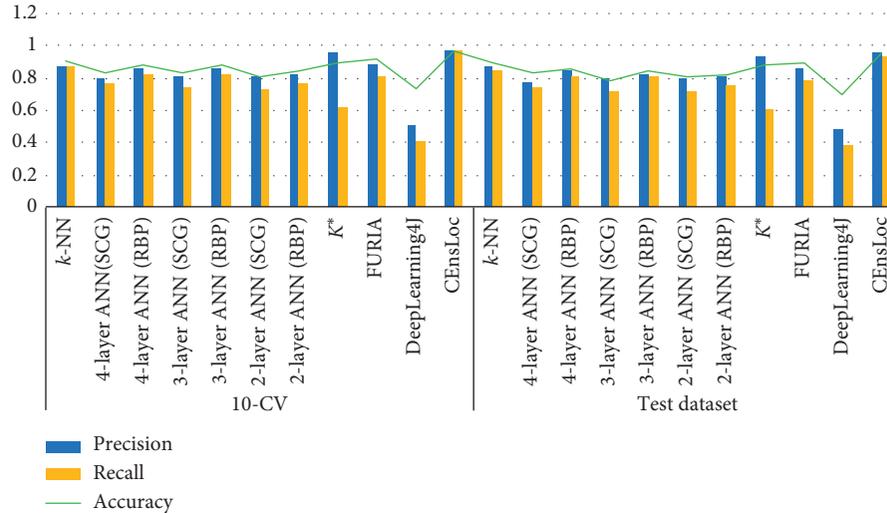


FIGURE 5: Performance measures on both 10-CV and test dataset.

TABLE 5: OOB loss.

Average 10-fold loss	OOB loss
0.0292813	0.0300123

available hybrid technologies at a given time, is also part of the plan.

Data Availability

The data used to support the findings of this study are available from the corresponding author upon request.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

This research was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (NRF-2018R1D1A1B07048697).

References

- [1] J. Torres-Sospedra, R. Montoliu, S. Trilles, Ó. Belmonte, and J. Huerta, "Comprehensive analysis of distance and similarity measures for Wi-Fi fingerprinting indoor positioning systems," *Expert Systems with Applications*, vol. 42, no. 23, pp. 9263–9278, 2015.
- [2] H. Mehmood and N. K. Tripathi, "Cascading artificial neural networks optimized by genetic algorithms and integrated with global navigation satellite system to offer accurate ubiquitous positioning in urban environment," *Computers, Environment and Urban Systems*, vol. 37, no. 1, pp. 35–44, 2013.
- [3] M. M. Soltani, A. Motamedi, and A. Hammad, "Enhancing cluster-based RFID tag localization using artificial neural networks and virtual reference tags," *Automation in Construction*, vol. 54, pp. 93–105, 2015.
- [4] M. L. Rodrigues, L. F. M. Vieira, and M. F. M. Campos, "Fingerprinting-based radio localization in indoor environments using multiple wireless technologies," in *Proceedings of IEEE International Symposium on Personal, Indoor and Mobile Radio Communications*, pp. 1203–1207, Toronto, ON, Canada, September 2011.
- [5] J. Luo and H. Gao, "Deep belief networks for fingerprinting indoor localization using ultrawideband technology," *International Journal of Distributed Sensor Networks*, vol. 12, no. 1, article 5840916, 2016.
- [6] S. Saeedi, L. Paull, M. Trentini, and H. Li, "Neural network-based multiple robot simultaneous localization and mapping," *IEEE Transactions on Neural Networks*, vol. 22, no. 12, pp. 880–885, 2011.
- [7] M. Azizyan, R. R. Choudhury, and I. Constandache, "SurroundSense: mobile phone localization via ambience fingerprinting," in *Proceedings of the 15th annual international conference on Mobile computing and networking-MobiCom'09*, pp. 261–272, Beijing, China, 2009.
- [8] L. Pei, R. Chen, J. Liu et al., "Motion Recognition Assisted Indoor Wireless Navigation on a Mobile Phone," in *Proceedings of the 23rd International Technical Meeting of The Satellite Division of the Institute of Navigation*, pp. 3366–3375, Portland, OR, USA, September 2010.
- [9] P. S. Nagpal and R. Rashidzadeh, "Indoor positioning using magnetic compass and accelerometer of smartphones," in *Proceedings of International Conference on Selected Topics in Mobile and Wireless Networking (MoWNeT)*, pp. 140–145, Montreal, Canada, 2013.
- [10] J. Menke and A. Zakhor, "Multi-modal indoor positioning of mobile devices," in *Proceedings of IEEE International Conference on Indoor Positioning and Indoor Navigation*, pp. 13–16, Alberta, Canada, October 2015.
- [11] M. Oussalah, M. Alakhras, and M. I. Hussein, "Multivariable fuzzy inference system for fingerprinting indoor localization," *Fuzzy Sets and Systems*, vol. 269, pp. 65–89, 2015.
- [12] N. Li, J. Chen, Y. Yuan, X. Tian, Y. Han, and M. Xia, "A wi-fi indoor localization strategy using particle swarm optimization based artificial neural networks," *International Journal of Distributed Sensor Networks*, vol. 12, no. 3, article 4583147, 2016.

- [13] C. Song, J. Wang, and G. Yuan, "Hidden naive bayes indoor fingerprinting localization based on best-discriminating AP selection," *ISPRS International Journal of Geo-Information*, vol. 5, no. 10, p. 189, 2016.
- [14] P. Bahl and V. N. Padmanabhan, "RADAR: an in-building RF based user location and tracking system," in *Proceedings of IEEE INFOCOM 2000. Conference on Computer Communications. Nineteenth Annual Joint Conference of the IEEE Computer and Communications Societies (Cat. No. 00CH37064)*, vol. 2, pp. 775–784, Tel Aviv, Israel, March 2000.
- [15] M. Cooper, J. Biehl, G. Filby, and S. Kratz, "LoCo: boosting for indoor location classification combining Wi-Fi and BLE," *Personal and Ubiquitous Computing*, vol. 20, no. 1, pp. 1–14, 2016.
- [16] C. Wang, F. Wu, Z. Shi, and D. Zhang, "Indoor positioning technique by combining RFID and particle swarm optimization-based back propagation neural network," *Optik (Stuttg.)*, vol. 127, no. 17, pp. 6839–6849, 2016.
- [17] J. Xu, H. Dai, and W. Ying, "Multi-layer neural network for received signal strength-based indoor localisation," *IET Communications*, vol. 10, no. 6, pp. 717–723, 2016.
- [18] W. Zhang, K. Liu, W. Zhang, Y. Zhang, and J. Gu, "Deep neural networks for wireless localization in indoor and outdoor environments," *Neurocomputing*, vol. 194, pp. 279–287, 2016.
- [19] L. Calderoni, M. Ferrara, A. Franco, and D. Maio, "Indoor localization in a hospital environment using Random Forest classifiers," *Expert Systems with Applications*, vol. 42, no. 1, pp. 125–134, 2015.
- [20] E. Jedari, Z. Wu, R. Rashidzadeh, and M. Saif, "Wi-Fi based indoor location positioning employing random forest classifier," in *Proceedings of 2015 International Conference on Indoor Positioning and Indoor Navigation (IPIN)*, pp. 13–16, Banff, Alberta, Canada, October 2015.
- [21] Y. Mo, Z. Zhang, Y. Lu, W. Meng, and G. Agha, "Random forest based coarse locating and KPCA feature extraction for indoor positioning system," *Mathematical Problems in Engineering*, vol. 2014, Article ID 850926, 8 pages, 2014.
- [22] R. Górak and M. Luckner, "Malfunction immune Wi-Fi localisation method," in *Computational Collective Intelligence*, pp. 328–337, Springer, Cham, Switzerland, 2015.
- [23] R. Górak and M. Luckner, "Modified random forest algorithm for Wi-Fi indoor localization system," in *Proceedings of International Conference on Computational Collective Intelligence*, pp. 147–157, Cham, Switzerland, September 2016.
- [24] J. Niu, B. Wang, L. Cheng, and J. J. P. C. Rodrigues, "WicLoc: an indoor localization system based on WiFi fingerprints and crowdsourcing," in *Proceedings of 2015 IEEE International Conference on Communications (ICC)*, pp. 3008–3013, London, UK, June 2015.
- [25] G. Ding, Z. Tan, J. Zhang, and L. Zhang, "Fingerprinting localization based on affinity propagation clustering and artificial neural networks," in *Proceedings of IEEE Wireless Communications and Networking Conference (WCNC)*, pp. 2317–2322, Shanghai, China, April 2013.
- [26] O. León, J. Hernández-Serrano, and M. Soriano, "Securing cognitive radio networks," *International Journal of Communication Systems*, vol. 23, no. 5, pp. 633–652, 2010.
- [27] S. Tuncer and T. Tuncer, "Indoor localization with bluetooth technology using artificial neural networks," in *Proceedings of the IEEE 19th International Conference on Intelligent Engineering Systems*, pp. 213–217, Bratislava, Slovakia, September 2015.
- [28] B. A. Akram, A. H. Akbar, B. Wajid, O. Shafiq, and A. Zafar, "LocSwayamwar: finding a suitable ML algorithm for wi-fi fingerprinting based indoor positioning system," in *Lecture Notes in Electrical Engineering*, A. Boyaci, A. Ekti, M. Aydin, and S. Yarkan, Eds., vol. 504, Springer, Singapore, 2018.
- [29] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [30] K. Kaji and N. Kawaguchi, "Design and implementation of wifi indoor localization based on Gaussian mixture model and particle filter," in *Proceedings of 2012 International Conference on Indoor Positioning and Indoor Navigation (IPIN)*, pp. 1–9, Sydney, Australia, November 2012.
- [31] C. Wu, Z. Yang, Y. Liu, and W. Xi, "WILL: wireless indoor localization without site survey," *IEEE Transactions on Parallel and Distributed Systems*, vol. 24, no. 4, pp. 839–848, 2013.
- [32] S. Yang, P. Dessai, M. Verma, and M. Gerla, "FreeLoc: calibration-free crowdsourced indoor localization," in *Proceedings of IEEE INFOCOM*, pp. 2481–2489, Turin, Italy, April 2013.
- [33] T. Seidl and H.-P. Kriegel, "Optimal multi-step k-nearest neighbor search," *ACM SIGMOD Record*, vol. 27, no. 2, pp. 154–165, 1998.
- [34] W. S. McCulloch and W. Pitts, "A logical calculus nervous activity," *Bulletin of Mathematical Biology*, vol. 52, no. 1-2, pp. 99–115, 1990.
- [35] J. G. Cleary and L. E. Trigg, "K*: an instance-based learner using an entropic distance measure," in *Proceedings of Twelfth International Conference on Machine Learning*, vol. 5, pp. 1–14, Tahoe City, CA, USA, July 1995.
- [36] J. Hühn and E. Hüllermeier, "FURIA: an algorithm for unordered fuzzy rule induction," *Data Mining and Knowledge Discovery*, vol. 19, no. 3, pp. 293–319, 2009.
- [37] S. Ciss, *Generalization Error and Out-of-bag Bounds in Random (Uniform) Forests*, Ph.D. thesis, French University, Paris, France, 2015.

Research Article

Malaria Vulnerability Map Mobile System Development Using GIS-Based Decision-Making Technique

Jung-Yoon Kim ¹, Sung-Jong Eun,² and Dong Kyun Park ³

¹Graduate School of Game, Gachon University, 1342 Seongnam Daero, Sujeong-Gu, Seongnam-Si, Gyeonggi-Do 461-701, Republic of Korea

²Health IT Research Center, Gil Medical Center, Gachon University College of Medicine, Incheon, Republic of Korea

³Department of Gastrointestinal Medicine, Gil Medical Center, Gachon University College of Medicine, Incheon, Republic of Korea

Correspondence should be addressed to Dong Kyun Park; pdk66@gilhospital.com

Received 18 June 2018; Accepted 2 August 2018; Published 18 September 2018

Academic Editor: Jaegool Yim

Copyright © 2018 Jung-Yoon Kim et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

This paper aimed at improving the lack of GIS information use and compatibility of multiplatform which represented limits that existing malaria risk analysis tools had. For this, the author developed mobile web-based malaria vulnerability map system using GIS information. This system consists of system database construction, malaria risk calculation function, visual expression function, and website and mobile application. This system was developed based on Incheon region only. Database includes information on air temperature and amount of precipitation as well. With regard to malaria risk calculation, guideline provided by Korea Centers for Disease Control and Prevention was followed first and then decision-making technique was used. Calculating criteria value for risk index made it possible to estimate more precise risk. With regard to visual expression function, database constructed earlier and risk information were linked to print out graphic map and graphs so that more intuitive and visible expression can be provided based on animation technique. This system allows a user to check information in real time and can be used anywhere anytime. Mobile push function is to enhance user's convenience. Such web map is useful to general users as well as experts.

1. Introduction

Malaria is infectious disease which is mediated and disseminated by *Anopheles sinensis*. Malaria is caused by *Plasmodium* (*Plasmodium* genus) which penetrates in red blood cell. Malaria shows symptom of fever, chill, hepatomegaly, splenomegaly, and so on. Malaria is reported to have occurred in one hundred four countries in the world. It is said that over one billion people around the world are exposed to a danger of malaria and over three hundred million people are infected with malaria and over one million people die of malaria on a yearly basis [1]. *Plasmodium* is *Plasmodium falciparum*, *P. vivax*, *P. malariae*, *P. ovale*, and *P. knowlesi*. It is said that malaria occurs with mosquito as mediator. It is reported that in South Korea, mainly *vivax* malaria caused by *P. vivax* which is mediated

by *Anopheles* genus mosquito has occurred. As from 1979, South Korea was declared as malaria-free area. However, in 1993, *vivax* malaria occurred in soldiers on the front line who served in the army near North Korea. It is reported that from 1993 to 2013, 500–4,000 malaria cases occurred centering on DMZ (demilitarized zone) annually. *Vivax* malaria which reoccurred in 1993 has appeared for about twenty years in northern part of Gyeonggi Province and varied in frequency of occurrence. Until now, *vivax* malaria occurred with three times of peaks (1998~2000, 2007, and 2010, respectively). It is reported that malaria cases have continued to occur until 2017 [1–3]. Figure 1 shows occurrence of malaria in South Korea by years.

Malaria is infectious disease that needs continuous control. Scholars have conducted studies on the cause of malaria and incidence rate by correlation considering meteorological elements and community characters [4–6]. In addition, they

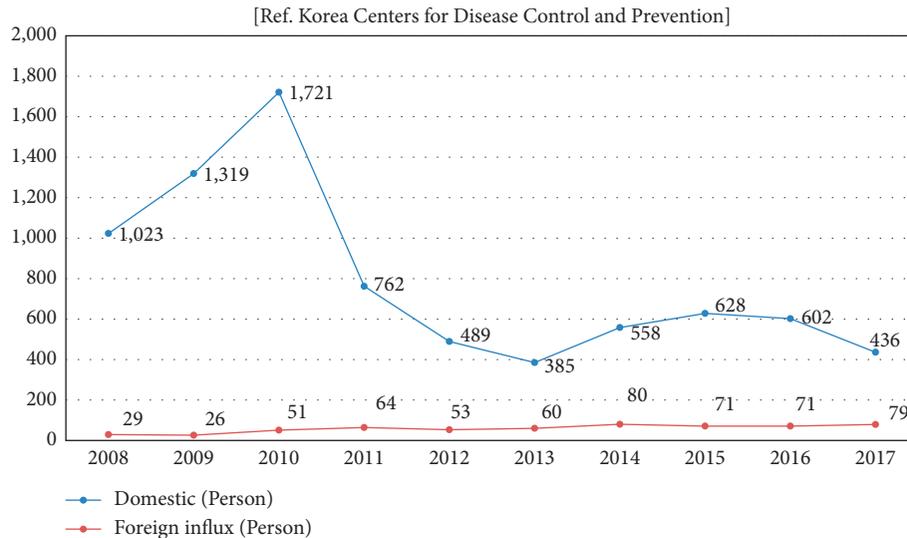


FIGURE 1: Annual malaria infection status.

analyzed route of infection in terms of epidemiology or environmental factors to prevent occurrence of malaria.

Software tool [7] that supports analysis of malaria-related information is being developed. As quantity of data becomes increased, applying technique of processing big data to correlation analysis is being raising with care. Like this, in health science, studies have been conducted in a specialized way making good use of IT-related technologies. Database in the field of health science includes a variety of data most of which is obvious in terms of space or refers to space. Therefore, database in the field of health science has the advantage of using data in GIS. Diversity, size, and complexity of health science data need more application of GIS [8]. Field of controlling infectious disease such as malaria needs GIS-based tools for risk analysis. This is to make good use of GIS-based analysis that can analyze information in an intuitive and visible way. Such tools have been developed, but they are not widely used due to problem of compatibility of independent platform, problem of exclusive institutional policy, and limited information.

This paper built malaria vulnerability map system using GIS information-based decision-making technique in order to limit problems of abovementioned tools and compatibility of multiplatform. Malaria vulnerability map system drew well-founded resulting values by making good use of various data on air temperature, humidity, amount of precipitation, lakes, and population with GIS information. Malaria vulnerability map system was built based on mobile web to enhance users' convenience. In terms of UX/UI, animation function helps enhance users' immersion and accessibility by providing intuitive information. This study consists of present condition of relevant tool, content of development system, conclusion, and investigation.

2. Related Works

Tools that analyze risk of malaria have been developed at home and abroad. In South Korea, abovementioned tools are used in some institutions for internal use only. In foreign

countries, MDAST (Malaria Decision Analysis Support Tool) project is being carried out. MDAST project is developing a tool that analyzes a cause of occurrence of malaria in Kenya, Tanzania, and Uganda where malaria occurs frequently and predicts incidence rate [9].

MDAST aims to provide a platform to analyze influences in terms of health, society, and environment for institution that establishes and evaluates malaria-related policy by using various scientific methods of analyzing information in malaria. This service downloads a tool. Sets it up on local, enters relevant data, sets several factors, applies desired analysis methods, and obtains a result. Analysis methods support various statistical analysis methods. Above tools are widely used but require professional knowledge in using such as factor setting and analysis method selection and impossible to provide online-based real-time support. The major drawback of above tools is that it is difficult for general user to access.

In South Korea, above tools are used through intranet or local installation by health-care-related institutions. They are used by some institutions for their internal use. Since the early 2000, some institutions have developed analysis tools. Environmental factors such as weather and geography and weights are selected on a case-by-case basis to handle system. GIS-based analysis system predicts a risk of malaria and prints out information on a map. For above tools, setting factors of several environmental elements plays an important role. Factors are handled through user definition.

GIS-based analysis system is to cope with unexpected local occurrence of malaria and prevent malaria from spreading by detecting an area where malaria is likely to occur as soon as possible. GIS-based analysis system which is used in South Korea varies in tools among institution, which makes performance deviation for functions supported great. It is difficult to define tools which are used because they are used for internal use.

GIS-based analysis system which is used in South Korea is used by some institutions for their internal use and is based on local. There are few systems which show information on

environmental elements relating to malaria on the website. This study developed map system which delivers GIS-based information on risk of malaria on mobile website in consideration of the foregoing. System was built in Incheon region. DB was built by receiving data from institutions such as Korea Meteorological Administration and K-Weather. This study built web-based malaria vulnerability map system so that general users can be provided with information on malaria risk by region in real time.

3. Developing Malaria Vulnerability Map Mobile Using GIS-Based Decision- Making Technique

The purpose of this paper was to develop a map that can analyze factors which cause malaria by using mobile web-based GIS (Geographic Information System). Decision tree analysis method was applied to model areas where malaria is more likely to occur. GIS and decision-making system to develop Korean-type malaria vulnerability map were finally built.

A map of this system was made based on overlay operation method that can provide new information which is not identified on respective map by combining several maps. Population, humidity, temperature, amount of precipitation, and lakes were considered for decision-making technique. Weight was calculated by using AHP decision-making method. Final malaria vulnerability map was made by using PROMETHEE technique.

Map system was implemented based on mobile web for compatibility with mobile platform. For web-based visual output, frame animation-based reading of areas where malaria appeared, reading of summarized information on areas where malaria appeared, reading of detailed information on areas where malaria appeared, and display of risk according to areas where malaria appeared were expressed visually.

3.1. Details of System Development. System functions consist of system database, malaria risk calculation, visual expression, and website and mobile application. Detail core functions are decision-making technique-based modelling part for malaria vulnerability calculation and mobile web map construction part for visual expression and program extension. The system configuration is shown in Figure 2.

3.1.1. Building Up System Database. Web-based system and MySQL with good compatibility and proven stability were used as database system to organize database.

Database information specific to Incheon region consists of data from 2011 to 2017. Area code was extracted based on longitude and latitude to provide information specific to Incheon so that it can provide detailed information according to region.

Information consisting of DB includes population, humidity, air temperature, amount of precipitation, and lake in Incheon. Information on lakes was used in determining vulnerability because whether there is mosquito can be

a factor on which high weight for malaria risk can be placed. With regard to information on air temperature, relative humidity, and amount of precipitation, data provided by Korea Meteorological Administration and K-Weather were used. Other information was provided by Incheon Institute of Health Environment. For relevant DB, schema was defined in one table without dividing it into several tables according to characteristics of data to cut down on expense caused by reference between tables and extract database information rapidly. Table 1 shows the structure of the system DB, which shows the structure of information on *Anopheles sinensis*, region code, sex, temperature, relative humidity, and precipitation.

For system database organized, unlike existing analysis tools where information on air temperature or amount of precipitation was not provided, various environmental factors were considered.

3.1.2. Developing a Model for Calculating Malaria Risk

(1) Calculating optimum weight using AHP (analytic hierarchy process) technique. Malaria risk is calculated based on various standards. For risk modelling work, various weights are considered. In this paper, population, humidity, air temperature, amount of precipitation, and lake were considered. Weight for each criterion was reflected based on experts' opinion as importance of standards differs in creating malarial risk. For this, AHP (analytic hierarchy process) was used. AHP [10–12] is a calculation model which was devised based on the fact that brain uses stepwise or hierarchical analysis process. AHP reaches a final decision making by dividing whole process of decision making into multiphases and analyzing it on step-by-step basis. The author of this paper carried out AHP based on flowchart as shown in Figure 3.

Order of standards which were considered based on Figure 3 was developed as shown in Figure 4. In phase 2 of order of rank, comparison of standards was made based on judgment of experts in malaria [13]. Table 2 shows comparison of pairs. Lastly, weight for standards was calculated by using eigenvector and eigenvalue in each comparison matrix.

Final weight criteria estimated through AHP decision-making technique were calculated based on temperature of 0.51, humidity of 0.34, population of 0.1, and lake of 0.15.

(2) Calculating decision-making matrix using PROMETHEE technique. Methods that select the best alternative by objectively measuring preference for alternatives considered under each criterion have been studied a lot. Typical techniques are ELECTRE and PROMETHEE. ELECTRE I, II, III, and IV technique proposed by B. Roy can solve difficult problems well and is the most reliable method, but the calculation process is very complex and requires decision maker to decide lots of parameters. On the other hand, PROMETHEE (Preference Ranking Organization Method Enrichment Evaluations) is very simple and easy for decision maker to understand and

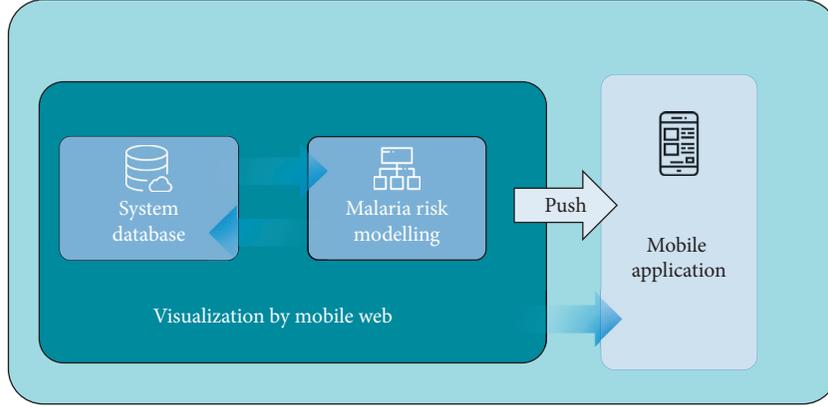


FIGURE 2: Conceptual diagram of malaria vulnerability map mobile system.

TABLE 1: External schema.

Data	Unit	Type
<i>Anopheles sinensis</i>	Count	Integer
Region code	—	Integer
Male	1,000 persons	Integer
Female	1,000 persons	Integer
Temperature	°C	Float
Relative humidity	%	Float
Amount of precipitation	mm	Float

preference intensity for alternatives is expressed with easy concept and requires decision maker to decide up to two parameters, which makes it preferred (Brans, 1985).

PROMETHEE proposed by Brans and Vincke [14] defines preference functions according to each evaluation criterion as Type I (Usual Criterion), Type II (Quasi-criterion), Type III (Criterion with Linear Preference), Type IV (Level Criterion), Type V (Criterion with Linear Preference and Indifference Area), and Type VI (Gaussian Criterion) [15–19]. PROMETHEE evaluates preference relationships among alternatives by using preference index $\pi(a, b)$.

$$\pi(a, b) = \frac{1}{k} \sum_{h=1}^k p_h(a, b), \quad (1)$$

All evaluation criteria must be defined as one of the abovementioned six preference functions. $p_h(a, b)$ in Equation (1) represents preference function value for evaluation criteria h . Decision maker should decide a shape of preference function according to evaluation criteria and assign threshold according to preference functions.

The author of this paper implemented PROMETHEE for each level five times. In each repetition, $n \times k$ decision-making matrix was created. n is the number of points (alternatives) and k is the number of criteria. Decision-making matrix number represents criteria for alternatives (regions). Table 3 represents decision-making matrix drawn when first time was repeated.

Decision matrix elements to be used in decision-making technique are decided based on the foregoing. Applicable components are decision-making matrix. Incheon region was divided into 28 areas. Matrix values to

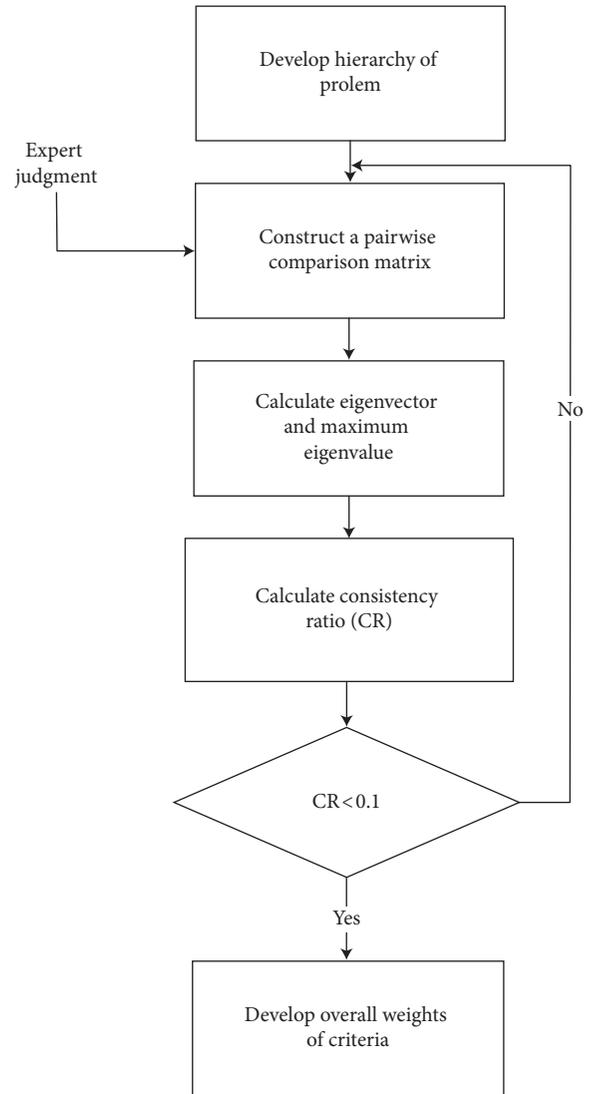


FIGURE 3: AHP for calculating criteria weight.

decide population, humidity, temperature, and lake were as shown above. Table 4 represents final result obtained from PROMETHEE. Five different rankings for 28 areas were obtained by implementing PROMETHEE five times.

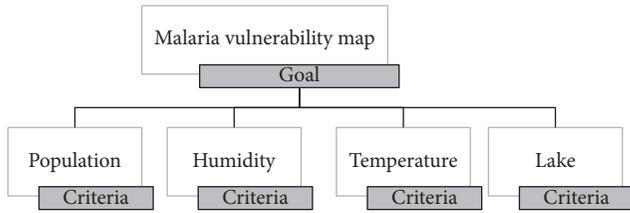


FIGURE 4: Structure of order of rank in criteria.

TABLE 2: External schema.

Criteria	Air temperature	Humidity	Lake	Population
Air temperature	1	3	3	4
Humidity	1/3	1	2	2
Lake	1/3	1/2	1	2
Population	1/4	1/2	1/2	1

TABLE 3: Decision-making matrix.

Zones	Population	Humidity	Temperature	Lake
1	4.085648	5.867935	1	4.064596
2	6.652511	6	1.685264	8.64964
3	6.11661	7.43814	8.964341	5
4	4.140974	7.444163	7.649631	7.165594
5	4	2.487076	1	0
6	3	5	1	0
7	3	5	1	0
8	3.452714	3.34022	1	0
9	5.655222	6	40951.15	6.865466
10	8.097396	6	97781.63	4.169856
11	5.438997	6.478203	1.986046	4.359985
12	7.964342	6.408581	7.865646	3.598465
13	7.019162	7.085386	6.131337	4.006419
14	5.354643	8.078567	8.643161	7.645985
15	8.989335	6	3.641969	4.036452
16	1.969482	9	4.945094	7.264995
17	4.945331	7.766397	5.619694	9.165985
18	7.842272	6.142719	6.940964	4.358643
19	3.928667	2.519822	1	0
20	4.123696	4.496581	1.690846	2.365646
21	4	3.717338	1	0
22	4	2.113713	1	0
23	3	5.402396	1	0
24	3	5	1	0
25	4	3.196849	1	0
26	3.219948	3.921217	1	0
27	3	4.640017	1	0
28	3.361399	3.569278	1	0

3.1.3. Creating Malaria Vulnerability Map

(1) *Raster overlay operation.* Overlay operation is applied to most GIS programs. Overlay operation provides new information by combining several maps. In overlay operation, new special elements are formed based on several maps. Overlay operation is handled based on raster and vector maps. Raster data structure is suitable for such operation because all maps used for analysis have the same georeference. They have the

TABLE 4: PROMETHEE index.

Zones	Population
1	14
2	10
3	2
4	8
5	21
6	0
7	0
8	25
9	11
10	6
11	12
12	5
13	4
14	1
15	9
16	13
17	3
18	7
19	23
20	15
21	17
22	26
23	16
24	18
25	19
26	22
27	20
28	28

same number of pixels and consist of row and column and have the same pixel size and coordinate. Accordingly, when several maps are combined, a program examines each pixel and the same figures can be checked from different maps. In raster overlay, cell numbers are combined in a specific way and figures obtained are assigned to corresponding cells in output layer. Raster overlay is applied to explicit or ordinal number data. Number or characters are stored in each raster cell. Figures in each cell correspond to items of raster variables. In Figure 5, Layer A represents solid data and Layer B shows an example of raster overlay that records land use. Number of potential output items is number of combination that input items can have ($2 * 3 = 6$).

Figure 6 shows number of population, humidity, temperature, lake, and vulnerability map result calculated by raster overlay. Population shape file based on areas of distribution of population has been converted to raster file. As all criteria maps must have common scale, figures of population raster file have been converted to figures of 0 to 10 by using reclassification analysis. Only reclassification analysis was implemented based on humidity data. Figure was obtained by reclassifying humidity raster file. Like humidity map, air temperature map was created by using reclassification analysis. This was handled by reclassifying figures of air temperature raster file. As distance from lake is greater, malaria risk decreased. A number of ring buffer analyzes were implemented to determine area around lake from diverse distances. By using this analysis, various buffers were created from distances of 100, 200, 300, 400, 600, 800, 1000,

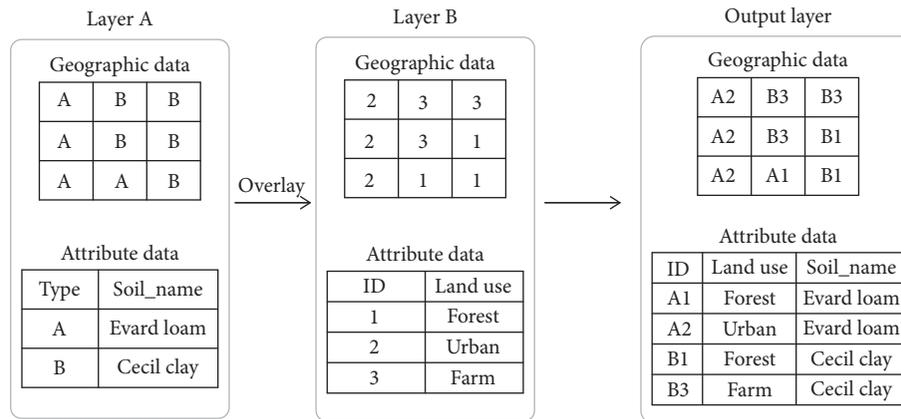


FIGURE 5: Analysis of raster overlay.

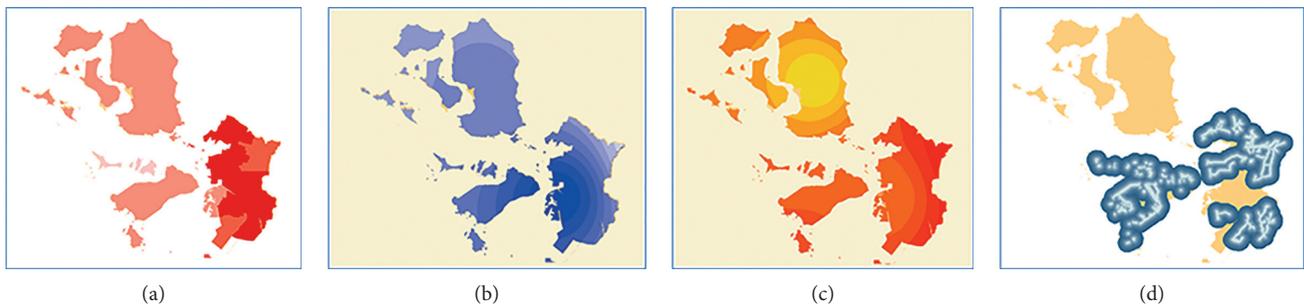


FIGURE 6: A result of raster overlay-based conversion and mapping. (a) Population map. (b) Humidity map. (c) Temperature map. (d) Lake map.

1500, and 2000. Buffer layers were converted to raster file based on distance area. The result that overlaps the map based on the converted raster file is shown in Figure 7.

(2) *Creating malaria vulnerability map through Thiessen polygon.* Incheon region was classified into 28 zones. Procedures for creating polygon were repeated five times, and 28 zones were set in Incheon Metropolitan City at random. Thiessen polygon [20] analysis was conducted by using 28 zones. Figure 7 shows Thiessen polygon created by repetition of five times.

Analysis was conducted by using mean statistics obtained by assigning mean figures in each zone to output cells. This procedure was implemented for Thiessen polygon and shows criteria map according to zones for first repetition. Accordingly, PROMETHEE was applied to criteria figures corresponding to 28 zones. Malaria risk map was created by combining risk maps obtained every time repetition was made. Repetition of applicable procedures helps in creation of more accurate malaria risk map (map closer to reality).

3.2. Developing Mobile Web-Based Malaria Vulnerability Map Application

3.2.1. *Developing Mobile Web-Based Application.* PHP is used as server development language for construction of

website. Client side language was implemented as HTML, CSS3, and JavaScript (Bootstrap, jQuery). Google Maps JavaScript API v3 was used as a map. Marker was put on control area with JavaScript. When marker is selected, information contained in database is presented in graph.

Mobile application was developed as android-based hybrid web app. Push function and GPS sensor were used so that location can be tracked on background. Information on malaria risk according to periods can be provided to a user based on it. Restful API was developed through PHP so that when longitude and latitude are handed over on background, a user can check whether it entered applicable area. For push implementation, GCM (Google Cloud Message) [21, 22] service was used so that push service can be used in a stable manner irrespective of terminal's connection to network.

3.2.2. Developing Animation-Based Visualization Technology.

Visual expression of information was implemented for intuitive and high visibility based on mobile web map. Frame animation [23], one of the animation techniques, was used for immersion. Wind direction, wind speed, spot atmospheric pressure, and sea level pressure in addition to foregoing database information were expressed in graph. Implementation functions consists of reading of areas where malaria appear, reading of summarized information on areas where malaria appear, reading of detailed

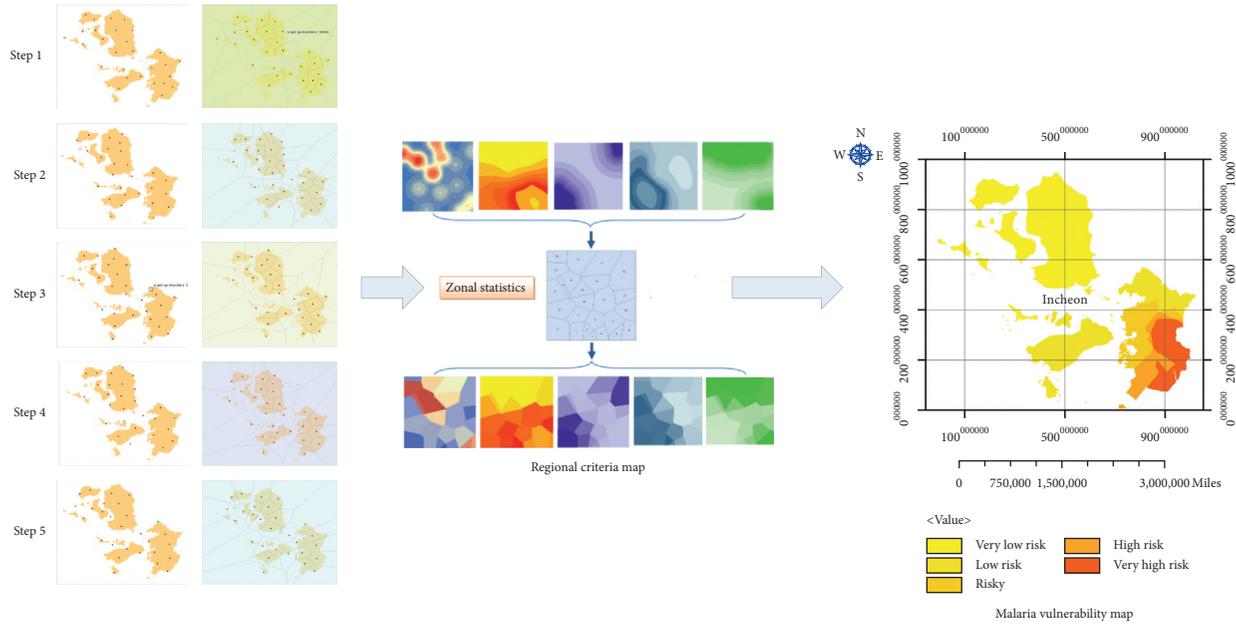


FIGURE 7: Raster overlay analysis.

information on areas where malaria appear, and risk animation function for areas where malaria appear. Visual expression function was implemented by using Google Maps API and JavaScript (Bootstrap, jQuery). Frame animation technique was used for effective view of risk information. Frame animation technique can highlight specific content and provide information to a user by changing images by time and static image to dynamic image stimulating human senses which are sensitive to dynamic motions. The author of this paper attempted to convey risk and environmental information by using the abovementioned functions and make a user have easier access by supporting intuitive UI for higher level of immersion. Applicable implementations are presented in Figure 8. The DB defines the schema as shown in Table 1.

4. Evaluation and Conclusion

4.1. System Model Evaluation. In order to evaluate the system performance, we measured accuracy of the malaria risk model. Comparisons of accuracy criteria were compared with those of the experts' risk criteria. So we use the confusion matrix and calculate the accuracy by true positive (TP), false positive (FP), false negative (FN), and true negative (TN). Table 5 represents confusion matrix result for the evaluation of the system model.

A total of 30 test cases were used for the evaluation, which resulted in a good performance of 91.7% accuracy. Equation (2) for calculating the accuracy is processed through the following confusion matrix:

$$\text{Accuracy} = \frac{(\text{TN} + \text{TP})}{(\text{TN} + \text{TP} + \text{FN} + \text{FP})} \quad (2)$$

For some inaccurate results, there is a problem in weight calculation of the AHP model to create the decision matrix.

Therefore, the issue of weight is in the order of rank in criteria. This section is classified as a special case, and more accurate weighting methods are needed. As a workaround, we expect to give more conditions to the quarter.

4.2. Conclusion. PHP is used as server development language for construction of website. Client side language was implemented as HTML, CSS3, and JavaScript (Bootstrap, jQuery). Google Maps JavaScript API v3 was used as a map.

Scholars have conducted research to develop GIS-based tools that can analyze malaria risk as continuous control of occurrence of malaria is needed. However, analysis tools at home and abroad are not widely used due to the lack of GIS information use, independent platform environment, and exclusive institutional policies.

This paper aimed at improving the lack of GIS information use and compatibility of multiplatform which represented limits that existing malaria risk analysis tools had. For this, the author developed mobile web-based malaria vulnerability map system using GIS information.

This system consists of system database construction, malaria risk calculation function, visual expression function, and website and mobile application. This system was developed based on Incheon region only. Database includes information on air temperature and amount of precipitation as well. With regard to malaria risk calculation, guideline provided by Korea Centers for Disease Control and Prevention was followed first and then decision-making technique was used. Calculating criteria value for risk index made it possible to estimate more precise risk. With regard to visual expression function, database constructed earlier and risk information were linked to print out graphic map and graphs so that more intuitive and visible expression can



FIGURE 8: Final result of system visualization. (a) Marking of malaria generation region. (b) Summation information (tooltip) of malaria generation region. (c) Risk time-series animation of malaria generation region.

TABLE 5: Confusion matrix result for evaluation.

Confusion matrix	Malaria risk model
True positive	27
False positive	3
True negative	28
False negative	2

be provided based on animation technique. This system allows a user to have easier access via website and mobile application which print out above information. This system allows a user to check information in real time and can be used anywhere anytime. Mobile push function is to enhance user's convenience. Such web map is useful to general users as well as experts.

As an evaluation of system model, we can check the good performance of 91.7% accuracy. However, for some inaccurate results, there is a problem in weight calculation of the AHP model to create the decision matrix. This section is classified as a special case, and more accurate weighting methods are needed. To solve the issue, we expect to give more conditions to calculate the weight value.

This system is expected to be used for institutions that operate under exclusive policy and service that prints out information based on web map. This system is expected to be used as basic technology that provides visible and intuitive information in the field of infographics which offer location-based information, GIS, industries as well as health care. The author of this paper will consider environmental elements which were not used in analyzing this system and seek additional mobile functions and a way to enhance user's accessibility. And we will also consider network security technologies for securing user data [24–27].

Data Availability

Previously reported geographic information data were used to support this study and are available at Korea Meteorological Administration, K-Weather, and Incheon Metropolitan City Institute of Health and Environment. These prior studies (and datasets) are cited at relevant places within the text as references [7, 8].

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

Acknowledgments

This research was supported by the Bio & Medical Technology Development Program of the National Research Foundation (NRF) funded by the Ministry of Science, ICT and Future Planning (2017M3A9E2072689).

References

- [1] J. Y. Kim, "Assessments on trend of malaria prevalence in Republic of Korea," *Public Health Weekly Report*, vol. 7, no. 12, pp. 237–242, 2014.
- [2] J.-W. Park, "Status of *vivax* malaria in the Republic of Korea," *Journal of the Korean Medical Association*, vol. 47, no. 6, pp. 521–526, 2004.
- [3] Korea Center For Disease Control and Prevention, *Infectious Diseases Surveillance Yearbook*, Ministry of Health and Welfare, Korea Center For Disease Control and Prevention, Cheongju-si, South Korea, 2014.
- [4] H. S. Shin, "Malaria prevalence rate and weather factors in Korea," *Health and Social Welfare Review*, vol. 31, no. 1, pp. 217–237, 2011.
- [5] S. M. Chae, D. J. Kim, S. J. Yoon, and H. S. Shin, "The impact of temperature rise and regional factors on malaria risk," *Health and Social Welfare Review*, vol. 34, no. 1, pp. 436–455, 2014.
- [6] J. Kwak, J. Lee, H. Han, and H. Kim, "A case study: malaria modeling based on climate variables in Korea," *Health and Social Welfare Review*, vol. 33, no. 4, pp. 547–569, 2013.
- [7] Y. W. Gong, "Design and implementation of outbreak estimation system for infectious disease using GIS (for malaria)," Engineering Master thesis, Kyungwon University, Gyeonggi-do, Republic of Korea, 2001.
- [8] R. Laurini, *Information Systems for Urban Planning: A Hypermedia Cooperative Approach*, CRC Press, Boca Raton, FL, USA, 2014.
- [9] <http://sites.duke.edu/mdast/>.

- [10] T. L. Saaty, "Analytic hierarchy process," in *Encyclopedia of Operations Research and Management Science*, pp. 52–64, Springer, Boston, MA, USA, 2013.
- [11] O. S. Vaidya and S. Kumar, "Analytic hierarchy process: an overview of applications," *European Journal of operational research*, vol. 169, no. 1, pp. 1–29, 2006.
- [12] S. H. Hashemi, A. Karimi, and M. Tavana, "An integrated green supplier selection approach with analytic network process and improved Grey relational analysis," *International Journal of Production Economics*, vol. 159, pp. 178–191, 2015.
- [13] Korea Center for Disease Control and Prevention, *Malaria Administrative Guideline 2013*, Korea Center for Disease Control and Prevention, Cheongju-si, South Korea, 2013.
- [14] J. P. Brans and Ph. Vincke, "Note—A Preference Ranking Organisation Method: (The PROMETHEE Method for Multiple Criteria Decision-Making)," *Management science*, vol. 31, no. 6, pp. 647–656, 1985.
- [15] J. Rezaei, "Best-worst multi-criteria decision-making method," *Omega*, vol. 53, pp. 49–57, 2015.
- [16] J.-J. Wang and D.-L. Yang, "Using a hybrid multi-criteria decision aid method for information systems outsourcing," *Computers & Operations Research*, vol. 34, no. 12, pp. 3691–3700, 2007.
- [17] C. Kahraman, S. C. Onar, and B. Oztaysi, "Fuzzy multicriteria decision-making: a literature review," *International Journal of Computational Intelligence Systems*, vol. 8, no. 4, pp. 637–666, 2015.
- [18] M. Dağdeviren, "Decision making in equipment selection: an integrated approach with AHP and PROMETHEE," *Journal of Intelligent Manufacturing*, vol. 19, no. 4, pp. 397–406, 2008.
- [19] A. Mardani, A. Jusoh, K. M. Nor, Z. Khalifah, N. Zakwan, and A. Valipour, "Multiple criteria decision-making techniques and their applications: a review of the literature from 2000 to 2014," *Economic Research-Ekonomska Istraživanja*, vol. 28, no. 1, pp. 516–571, 2015.
- [20] P. L. Ibisch, M. T. Hoffmann, S. Kreft et al., "A global map of roadless areas and their conservation status," *Science*, vol. 354, no. 6318, pp. 1423–1427, 2016.
- [21] Y. S. Yilmaz, B. I. Aydin, and M. Demirbas, "Google cloud messaging (GCM): an evaluation," in *Proceedings of Global Communications Conference (GLOBECOM)*, pp. 2807–2812, Austin, TX, USA, December 2014.
- [22] N. Saravanan, A. Mahendiran, N. V. Subramanian, and N. Sairam, "An implementation of RSA algorithm in Google cloud using cloud SQL," *Research Journal of Applied Sciences, Engineering and Technology*, vol. 4, no. 19, pp. 3574–3579, 2012.
- [23] M. Achibet, G. Casiez, A. Lécuyer, and M. Marchal, "THING: introducing a tablet-based interaction technique for controlling 3d hand models," in *Proceedings of 33rd Annual ACM Conference on Human Factors in Computing Systems*, New York, NY, USA, May 2015.
- [24] J. Cui, Y. Zhang, Z. Cai, A. Liu, and Y. Li, "Securing display path for security-sensitive applications on mobile devices," *Computers Materials and Continua*, vol. 55, no. 1, pp. 17–35, 2018.
- [25] C. Yin, J. Xi, R. Sun, and J. Wang, "Location privacy protection based on differential privacy strategy for big data in industrial internet-of things," *IEEE Transactions on Industrial Informatics*, vol. 14, no. 8, pp. 3628–3636, 2017.
- [26] J. Wang, C. Ju, H.-J. Kim, R. Simon Sherratt, and S. Lee, "A mobile assist coverage hole patching scheme based on particle swarm optimization for WSNs," *Cluster Computing*, vol. 1, pp. 1–9, 2017.
- [27] Y. Tu, Y. Lin, J. Wang et al., "Semi-supervised learning with generative adversarial networks on digital signal modulation classification," *Computers Materials & Continua*, vol. 55, no. 2, pp. 243–254, 2018.

Research Article

Location Privacy Protection Research Based on Querying Anonymous Region Construction for Smart Campus

Ruxia Sun ¹, Jinwen Xi ¹, Chunyong Yin ¹, Jin Wang ², and Gwang-jun Kim ³

¹School of Computer and Software, Jiangsu Engineering Center of Network Monitoring, Nanjing University of Information Science & Technology, Nanjing, China

²School of Computer & Communication Engineering, Changsha University of Science & Technology, Changsha, China

³Department of Computer Engineering, Chonnam National University, Gwangju, Republic of Korea

Correspondence should be addressed to Gwang-jun Kim; kgj@jnu.ac.kr

Received 15 June 2018; Revised 20 July 2018; Accepted 1 August 2018; Published 16 September 2018

Academic Editor: Jaegel Yim

Copyright © 2018 Ruxia Sun et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Along with the rapid development of smart campus, the deployment of novel learning applications for smart campus requires full consideration of information security issues. Location privacy protection is one of the most important issues, which considers the privacy protection and guarantees the quality of service. The existing schemes did not consider the area of the querying regions for location-based service provider (LSP) during the construction of the anonymous regions, which led to the low quality of service. To deal with this problem, the user's query range was introduced to present a novel anonymous region construction scheme. In the proposal, the anonymous server first generated the original anonymous subregions according to the user's privacy requirements, and then merged these subregions to construct the anonymity regions submitted to LSP based on the size of corresponding querying regions. The security and experiment analysis show that the proposed scheme not only protects the user's privacy effectively but also decreases the area of LSP querying regions and the region-constructing time, improving the quality of service for smart campus.

1. Introduction

With the development and construction of novel learning applications for smart campus, smart devices and services, smart meters, smart terminals, and the like are widely applied to offer real-time learning feedback to students through continuously monitoring and analyzing the status and activities of students with various devices and platforms. As a large number of smart meters and intelligence appliances are accessed, incorporating various technologies and enabling world-changing learning, the network border further extends to the user side. Security risks on the user side for smart campus will become more and more prominent. Data privacy issues, especially location privacy protection and the quality of service, must be considered [1].

Users' location privacy threats refer to the risks that an attacker can obtain unauthorized access to raw location data by locating a transmitting device and identifying the subject

(person) using it. Examples of such risks include spamming users with unwanted advertisements, drawing sensitive inferences from victims' visits to various locations (e.g., students and teachers' offices), and learning sensitive information about them (identity, religious and political affiliations, etc.). Hence, location privacy protection for smart campus is becoming a critical issue [2].

However, location information is consistently sent to service providers without protection when users query LBSs, allowing location providers to collect location information from all users. The collected location information may expose users to customized advertisement or even be sold to third parties. A worse scenario is location information may be leaked to adversaries with criminal intents. Therefore, many researchers focus on creating location protection algorithms to protect the location privacy of users [3].

The European Union's Information Protection Supervision Organization recently said that high-tech equipment

such as smart meters that monitor household energy consumption will pose a huge threat to personal privacy. Smart meters may track personal information, and the vast amounts of information collected can have serious consequences for consumers.

With the popularization of mobile devices and location technology, location-based services (LBSs) are widely used in real life, which refers to the user accesses to its designated location information query and entertainment services through the mobile device [4]. However, the location-based service provider (LSP) may also collect and abuse user's service information while providing the convenient LBS for the user, to illegally obtain the user's confidential information. The location privacy protection in LBS has attracted the extensive attention of researchers [5–8].

In view of the popularization of mobile positioning devices, if these novel learning applications are used in smart campus, the combination of location information and services at different moments in personal privacy may reveal sensitive information such as the user's behavior habits and work nature. For example, if a mobile user is collected near a hospital, the user may be presumed to have any disease or health condition. If the user's starting location and ending location are analyzed since the last few days, the user's home address or work unit, the nature of work, and so on can be speculated. Therefore, the location data of mobile object brings convenience to people and brings about the threat of revealing privacy, which may contain other sensitive information such as home address, personal preference, personality habit, health status, working property, personal income, etc. If this information falls into the hands of normal institutions, it is a tool for information protection, if it falls into the hands of illegal institutions, it will be the weapon of innocent destruction. What we can do is to seek transparency in the use of personal information and to protect user's location privacy not to be exploited by unscrupulous businesses and illegal agencies.

As the most commonly used LBS privacy protection method, the basic idea of k -anonymity [9] is that when a user sends a LBS query, he will send his real location and querying content to a trusted anonymous server, then the anonymous server will remove the user's identification information and generate anonymity regions containing other $k - 1$ users for the real location, finally sending them (all the generated anonymity regions and the real location) along with the querying content to the LSP. Compared with other LBS privacy protection methods (such as pseudo location [10], fuzzification [11], differential privacy [12], and cryptography-based methods [13, 14]), k -anonymity has the following advantages: (1) users can get accurate querying results; (2) the user's cost of computing and communication is small; and (3) this method can confuse the relevance between users and LBS queries. So, k -anonymity is widely used in LBS privacy protection [15, 16].

In smart campus, we can query the nearby points of interest. There are many sensors to be addressed for enabling such novel learning applications and services, which aims to enhance the service quality of novel learning applications. When k -anonymity method is used to protect the privacy of

users in the above query, if the anonymous region generated by the anonymous server is too large, the querying cost of the location-based service provider (LSP) will increase and the service quality deteriorates [17–19]. To solve this problem, the existing methods [20–26] obtain n disjoint anonymous subregions by removing the part that does not contain the users in the regions, to reduce the area of anonymous regions and improve the service quality as shown in Figure 1. However, the quality of service in LBS queries based on k -anonymity is not only related to the size of the anonymous regions but also to the user's query range. If we use the existing methods to construct the anonymous regions, the quality of service cannot be effectively improved. As shown in Figure 2, when using the existing division methods to divide the initial anonymous region, LSP will repeatedly search the points of interest in some regions, reducing the quality of service, and r is the query radius. This paper also proves this point through experiments.

This paper proposes the LBS privacy protection scheme based on querying anonymous region construction. In which, firstly, the anonymous server generates k initial anonymous subregions according to the privacy protection requirements of users and merges the anonymous regions according to the corresponding querying regions so that the anonymous regions finally submitted to the LSP can reduce the querying cost of LSP and improve the service quality without reducing the user's privacy protection level. This is the first k -anonymity privacy protection scheme based on constructing the user querying anonymous regions. The main contributions of this paper are as follows:

- (1) Based on the theoretical analysis, it is concluded that the existing anonymous region demarcation method cannot reduce the LSP querying cost and improve the service quality, and we prove it through experiments.
- (2) We think that the area of querying regions is the judgment criterion to merge the anonymous subregions and propose an anonymous region construction scheme based on the user's query range. Security analysis shows that the proposed scheme can effectively protect the users' location privacy.
- (3) A large number of experiments show that this scheme can effectively reduce the querying cost of LSP and improve the service quality, without causing a large computational cost for anonymous servers.

2. Related Work

At present, the international research on the privacy protection for smart campus is still in the initial stage. The discussion on the privacy protection focuses more on the risk analysis of exposing personal privacy for the wireless applications and devices. Research in the United States is the leader, which publishes relevant documents on this issue.

Privacy laws in the United States do not explicitly address the smart campus and its related data, which is same as the regulations of the existing national Internet of Energy

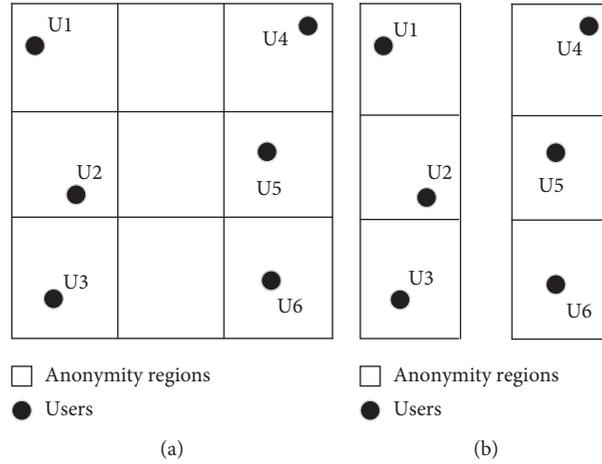


FIGURE 1: The existing anonymous regions division method. (a) Initial anonymity regions. (b) Anonymity regions after division.

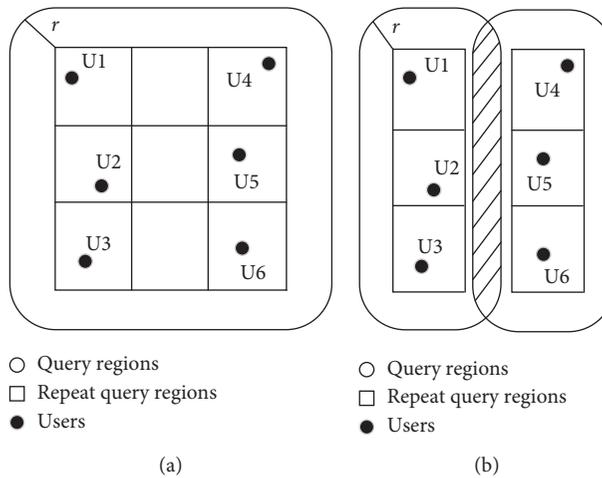


FIGURE 2: The querying regions of nearby points of interest. (a) Initial query regions. (b) Query regions after division.

and power transmission [27]. The existing laws and regulations need to be revised to appeal to the smart campus. At the same time, new data items in the Internet of Energy, as well as new ways of using existing data, require more research and public opinion to adapt to current laws or shape new laws.

US Internet of Energy Information Security is very concerned about privacy issues. NIST (National Institute of Standards and Technology) released “Guidelines for Smart Grid Cyber Security: Vol. 2, Privacy and Smart Grid” [28] in August 2010 and preliminarily analyzed the privacy issues of smart campus. For location privacy protection, many researchers have done a great deal of work.

Grutese applied k -anonymity to the field of location privacy protection for the first time [4]. The anonymous server divides the area with a quadtree structure and stores the users in the corresponding region in the node. When the user requests a service, a quadtree is retrieved upward from a leaf node corresponding to the user’s location. If the current leaf node does not satisfy the k -anonymity, the parent node will be retrieved until no less than other $k - 1$

users to obtain the anonymous region. However, if the number of users in a leaf node does not satisfy the privacy requirements, it needs to retrieve its parent node, which will result in a four-fold increase in the anonymity area, thereby degrading the user’s service quality. To deal with the above problem, Mokbel et al. [29] improved the anonymous region construction method in [4]. In their scheme, if the current leaf node does not satisfy k -anonymity, firstly its sibling nodes will be retrieved, if it still does not satisfy the user’s privacy requirements, then it retrieves the parent node. To further solve the problem of low quality of service due to the large anonymous region based on the quadtree structure, Li et al. [30] proposed to reduce the area of anonymous regions and improve the service quality by suppressing part of users’ requests and deleting the most distant trails.

Subsequently, the researchers proposed Clique Cloak [31] and Hilbert Cloak [32], respectively, to construct the anonymous regions by searching for the nearest users who meet the privacy requirements. In [31], users can customize the privacy protection requirements, but the scheme uses an undirected graph to construct an anonymous region and the

cost of the solution is too large. It may happen that the anonymous region has not been successfully constructed but beyond the anonymous period. In [32], the anonymous server maps all users in two-dimensional space to a one-dimensional array according to the Hilbert curve and divides the users into several sets according to the value of k . When a user makes a service request, the anonymous region is constructed by using the user set to which the user belongs.

Recently, Jin et al. [2] thought the anonymization server could lead to a new class of attacks called location injection attacks which can successfully violate users' indistinguishability (guaranteed by k -anonymity) among a set of users. Therefore, they propose and characterize location injection attacks by presenting a set of attack models and quantify the costs associated with them. Then, they propose and evaluate k -Trustee, which is resilient to location injection attacks and guarantees a low bound on the user's indistinguishability. The experimental results show that the proposed cloaking algorithm guaranteeing k -Trustee is effective against various location injection attacks, but the quality of location services is poor and cannot be guaranteed.

In all the above solutions, if the k users used to construct the anonymous regions are far away from each other, the anonymous regions may still be too large, and the service quality may be low. Tan et al. [33] were the first ones to apply the idea of regionalization to the construction of anonymous regions, which divided the users in the anonymous regions into separate groups through the Hilbert space filling curve. When a user makes a server request, the anonymous server will use the locations of other users in the group to which it belongs to construct the anonymous region. Subsequently, Li and Zhu [34] also used the method of regionalization to research on reducing the area of anonymous regions and improving the quality of service. The anonymous server firstly constructs an anonymous region containing k users and then removes the anonymous regions that do not contain the users according to relationship among the users' locations to form multiple anonymous subregions that do not intersect each other.

However, the existing anonymous region construction schemes [35–37] ignore the impact of the users' query range on LBS querying service quality. When using these methods to construct the anonymous regions, the LSP will repeatedly search the points of interest in the partial regions, reducing the service quality.

3. Improved k -Value Location Privacy Protection Method

3.1. System Structure. This paper uses the centralized system architecture [38], composed of three parts: user, anonymous server, and LSP. When the user requests a service, the anonymous server blurs the real location of the user into an anonymous region and sends it to the LSP, which contains not only the real user but also other at least $k-1$ users. In this case, the correct rate that LSP correlates the service query with the right user will be not more than $1/k$, which achieves k -anonymity. The system structure is shown in Figure 3.

Assume that there is a secure communication channel between the user and the anonymous server. When the user queries the point of interest nearby, the secure channel is used to send the query request $q = \langle \text{ID}, L(x, y), r, \text{POI}, p \rangle$ to the trusted anonymous server. In which, ID represents the identity of the user; $L(x, y)$ represents the location coordinates of the user; r represents the radius of the user query; POI represents the point of interest of the user query; p represents the privacy protection requirement of the user current query; k represents the anonymous region generated by the anonymous server contains at least other $k-1$ users; and A_{\min} represents the minimum area of anonymous regions generated by anonymous server.

After receiving the user's request, the trusted anonymous server will determine the identity through authentication and find other $k-1$ users to generate anonymous regions that the area is not less than A_{\min} according to the user's privacy protection requirements $p = (k, A_{\min})$ and then send the query request $Q = \langle \text{CR}, r, \text{POI} \rangle$ obtained from the anonymization to the semitrusted LSP. CR represents the anonymous region generated by the anonymous server for the current user making the service query.

After receiving the anonymous query request sent by the anonymous server, the LSP will search in the database and return all the query candidate results to the anonymous server. After the anonymous server receives the query result from the LSP, it selects the query result according to the location $L(x, y)$ of the user, and finally returns the accurate query result to the user.

In this system model, we treat the LSP directly as an attacker. The attacking purposes are as follows: (1) the real location of the user could be identified in the anonymous region sent from the anonymous server and (2) the real user will be speculated from the query request.

In addition, in the above model, LBS query quality of service is mainly affected by the following four factors: (1) the time required by the anonymous server to generate anonymous regions; (2) the time that the anonymous server sends the anonymous query request to LSP; (3) the time required by the LSP to retrieve the database for the anonymous query request sent by the anonymous server; and (4) the time it takes for the LSP to send the retrieval result to the anonymous server, and the time required by the anonymous server to finalize the query result. Because the time taken for the LSP to send search results to the anonymous server and the time required for the anonymous server to finalize the query results is affected by the distribution of points of interest, the time required for the anonymous server to send the anonymous query request to the LSP is affected by the transmission bandwidth. Therefore, this paper evaluates the quality of service based on k -anonymous LBS queries only by the time it takes the anonymous server to generate the anonymous regions and the LSP to retrieve the database.

3.2. Querying Region of LSP. After receiving the anonymous query request $Q = \langle \text{CR}, r, \text{POI} \rangle$ sent by the anonymous server, the LSP first calculates the querying region QAR according to the anonymous region CR and the query radius

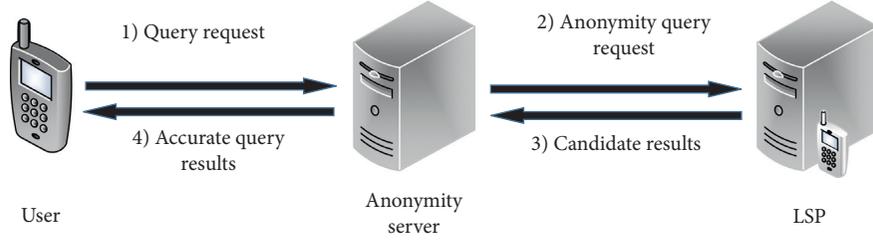


FIGURE 3: System structure.

r and then searches the interest point of the user query in the region QAR. The querying regions corresponding to the anonymous regions are as shown in Figure 4.

When the anonymous region is a circle, the area S_c of the querying region is

$$S_c = \pi(\tilde{r} + r)^2, \quad (1)$$

where \tilde{r} represents the radius of the circular anonymous region and πr^2 is the area of the anonymous region.

When the anonymous region is a rectangle, the area S_r of the querying region is

$$S_r = ab + 2(a + b)r + \pi r^2, \quad (2)$$

where a, b is the side length of the rectangular anonymous region and ab is the area of the anonymous region CR.

It can be learned from the above analysis that after receiving the anonymous query request sent by the anonymous server, the time required by the LSP to retrieve the database is not only related to the size of the anonymous region CR generated by the anonymous server but also to the radius of the user's query, in other words, it is determined by the querying region. However, the existing anonymous region partitioning schemes only take the size of the anonymous regions into account, which makes these solutions not effectively improve the LBS query service quality and even further reduces the service quality.

4. Anonymous Region Construction Scheme Based on Query Range

In this paper, the anonymous server first generates k initial anonymous subregions according to the privacy protection requirements of the users and constructs the corresponding querying regions by calculating and then merges the anonymous subregions to finally obtain the anonymous subregion set. The solution can reduce the LSP query cost and improve the service quality without lowering the user privacy protection level.

4.1. Generation of Initial Anonymous Subregions. After receiving the service request sent by the users, the anonymous server will start to search the other $k - 1$ users according to the user's privacy protection requirement $p = (k, A_{\min})$ and obtain their location information $L_1(x_1, y_1), \dots, L_{k-1}(x_{k-1}, y_{k-1})$. Then, the server will generate k disjoint initial anonymous subregions $AR_0, AR_1, \dots, AR_{k-1}$ to satisfy

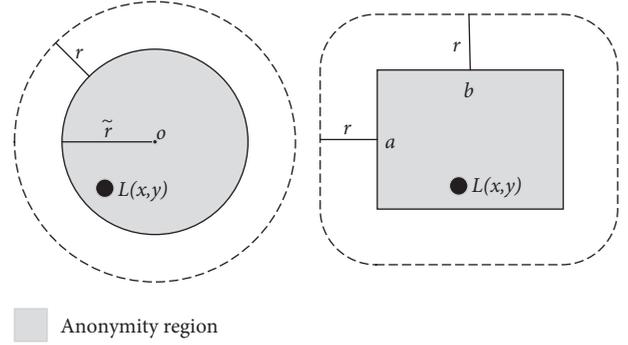


FIGURE 4: Querying region of LSP.

$$\begin{cases} \text{centre}(AR_i) \neq L_i(x_i, y_i), \\ \text{area}(AR_i) \neq A_{\min}, \end{cases} \quad (3)$$

where AR_i represents the initial anonymous subregion containing the i -th location $L_i(x_i, y_i)$, $0 \leq i \leq k - 1$; $L_0(x_0, y_0)$ represents the location of the user who sent the query request; $\text{centre}(AR_i)$ represents the central location of initial anonymous subregion AR_i ; and $\text{area}(AR_i)$ represents the area of initial subanonymous region AR_i .

4.2. Merger of Anonymous Subregions. After generating the k initial anonymous subregions $AR_0, AR_1, \dots, AR_{k-1}$, the anonymous server will calculate the area $\text{area}(QAR_0), \text{area}(QAR_1), \dots, \text{area}(QAR_{k-1})$ of the corresponding querying regions $QAR_0, QAR_1, \dots, QAR_{k-1}$, respectively, according to the query radius of the user and judge whether the anonymous subregions need to be merged on the basis of the area of the querying regions. In order to effectively reduce the cost, LSP retrieves the points of interest, and after the merger of anonymous subregions, it should ensure

$$S_{\min} = \min \sum_{i=0}^l \text{area}(QAR'_i). \quad (4)$$

That is to say, the sum of the area of each querying region QAR'_i corresponding to each anonymous subregion AR'_i in the anonymous regions CS is the smallest. AR'_i represents the i -th anonymous subregion after the merger of anonymous subregions, $0 \leq i \leq l$, $0 \leq l \leq k - 1$.

The merger process of anonymous subregions in the anonymous server is as follows:

- (1) Filtering the anonymous subregions that require region consolidation

To ensure that the querying area of the anonymous region is minimized, the anonymous server chooses the anonymous subregions AR_i and AR_j to merge if the area of querying region is the smallest after the merger of them. The selected anonymous subregions AR_i and AR_j are needed to satisfy

$$\begin{aligned} \forall i, j \in [0, k-1], \text{QAR}_{i,j} &= \operatorname{argmin}\{\operatorname{area}(\text{QAR}_{i,j})\}_{i \neq j}, \\ \operatorname{centre}(\text{QAR}_{i,j}) &\neq L_i(x_i, y_i), \\ \operatorname{centre}(\text{QAR}_{i,j}) &\neq L_j(x_j, y_j). \end{aligned} \quad (5)$$

- (2) Selecting anonymous subregions to merge

For the two anonymous subregions AR_i and AR_j , $AR_{i,j} = \operatorname{gen}(AR_i, AR_j)$ represents the new anonymous subregion formed by merging the anonymous subregion AR_i with AR_j and the corresponding querying region is $\text{QAR}_{i,j}$. If $\operatorname{area}(\text{QAR}_i) + \operatorname{area}(\text{QAR}_j) \leq \operatorname{area}(\text{QAR}_{i,j})$, it does not merge AR_i with AR_j . Else it merges AR_i with AR_j to make the new anonymous subregion $AR_{i,j}$.

- (3) Repeating the process until there is no need to merge anonymous subregions. In this case, the anonymous server will obtain the anonymous set

$$\text{CS} = \{AR'_0, AR'_1, \dots, AR'_j\}. \quad (6)$$

In this scheme, when the anonymous server merges the anonymous subregions, it chooses the anonymous regions AR_i and AR_j to make the area of the querying region the smallest after merging. Therefore, after every merger of the anonymous subregions, the area of the new querying region $AR_{i,j}$ is the smallest, which ensures that the querying area of the final anonymous region set CS is the smallest.

In this paper, the locations of k users are used as the input to generate the corresponding anonymous subregions $AR_0, AR_1, \dots, AR_{k-1}$ and merge them to obtain the final anonymous subregion set CS, which is submitted to the LSP (Algorithm 1).

5. Scheme Analysis and Experimental Results

5.1. Analysis of Security. In our scheme, we construct the disjoint initial anonymous subregions $AR_0, AR_1, \dots, AR_{k-1}$ that the area of each initial anonymous subregion is less than A_{\min} according to k real locations $L_0(x_0, y_0), L_1(x_1, y_1), \dots, L_{k-1}(x_{k-1}, y_{k-1})$ and ensure every real location is not located in the center of the initial anonymous subregion for smart campus, that is, $\operatorname{centre}(AR_i) \neq L_i(x_i, y_i)$. If the anonymous server directly sends the k initial anonymous subregions to the LSP, which cannot correctly identify the real location $L_i(x_i, y_i)$ of the user from the anonymous region AR_i . If the anonymous server uses the proposed scheme to merge the anonymous subregions and sends the final anonymous region set CS to the LSP, the area of every anonymous subregion AR_i satisfies $\operatorname{area}(AR_i) = A_{\min}$, so the

area of every merged anonymous subregion AR'_i will satisfy no less than A_{\min} and for any location $L_i(x_i, y_i) \in AR'_i$, and it satisfies that the center of AR'_i is not $L_i(x_i, y_i)$. In addition, the user's identity has been removed from the anonymous query request sent by the anonymous server. Though the LSP received the anonymous query request $Q = \langle \text{CR}, r, \text{POI} \rangle$, it is unable to know the user who made the service request. Therefore, this scheme proposed in this paper can effectively protect the privacy of users.

5.2. Analysis of Computational Complexity. When the anonymous server receives the service query request sent by the user and uses this solution proposed in this paper to generate an anonymous region for it, it firstly generates k initial anonymous subregions according to the privacy protection requirement of the user. The computational complexity of generating the initial anonymous subregions is $O(k)$. When the anonymous server determines whether any two anonymous subregions AR_i and AR_j need to be merged, it first needs to use the anonymous subregions AR_i and AR_j to generate a new anonymous region $AR_{i,j}$. In this case, the number of new anonymous regions that need to be generated is $C_k^2 = k(k-1)/2$, and the computational complexity is $O(k^2)$. Subsequently, the anonymous server will calculate the area of the querying region corresponding to each newly generated anonymous region $AR_{i,j}$, which requires $k(k-1)/2$ times of computation, and the computational complexity is $O(k^2)$. Finally, comparing the area of the querying region AR_i and AR_j with that of the newly generated anonymous region $AR_{i,j}$ to determine whether the both initial anonymous subregions are needed to merge, which requires $k(k-1)/2$ times of computation, and the computational complexity is $O(k^2)$. So the computational complexity of merging all the initial anonymous subregions is $O(k^2) + O(k^2) + O(k^2) = O(k^2)$. After obtaining the new anonymous region $AR_{i,j}$ merged from the initial anonymous subregions AR_i and AR_j , it also needs to judge whether the new anonymous subregions need to be merged again with other anonymous subregions. During the merging process of the anonymous subregions, the best case is that any one of k initial anonymous subregions does not need to be merged with others. In this case, the computational complexity required to implement this solution is

$$O_{\text{best}} = O(k) + O(k^2) = O(k^2). \quad (7)$$

In contrast, the worst case is that any one of k initial anonymous subregions needs to be merged with others and finally merged to an anonymous region. In this case, it needs to repeat the above anonymous region merger and judge $k-1$ time, and the computational complexity required to implement this solution is

$$O_{\text{worst}} = O(k) + O((k-1)k^2) = O(k^3). \quad (8)$$

5.3. Experimental Results and Analysis. In this paper, a network-based generator of moving objects (NGMO) [39] is used to generate the experimental data. This generator is often

Input: k pieces of locations $L_0(x_0, y_0), L_1(x_1, y_1), \dots, L_{k-1}(x_{k-1}, y_{k-1})$; the query radius r ; the privacy requirement A_{\min} ;
 Output: anonymous region set CS;

```

(1) for  $i = 0$  to  $k - 1$  do
(2)    $AR_i \leftarrow \text{Gen}(L_i(x_i, y_i))$ ;
(3)    $\text{centre}(AR_i) \neq L_i(x_i, y_i), \text{Area}(AR_i) = A_{\min}$ ;
(4)    $CS \leftarrow (AR_i)$ ;
(5)    $QAR_i \leftarrow \text{Gen}(AR_i, r)$ , calculating  $\text{Area}(QAR_i)$ ;
(6) end for
(7) for  $i, j = 0$  to  $ij$  and  $i \neq j$  do
(8)    $AR_{i,j} \leftarrow \text{Gen}(AR_i, AR_j)$ ;
(9)    $CR \leftarrow AR_{i,j}$ ;
(10)   $QAR_{i,j} \leftarrow \text{Gen}(AR_{i,j}, r)$ , calculating  $\text{Area}(QAR_{i,j})$ ;
(11)  if  $\text{Area}(QAR_i) + \text{Area}(QAR_j) \leq \text{Area}(QAR_{i,j})$  then
(12)    $CS \leftarrow CS \setminus \{AR_{i,j}\}$ ;
(13)  end if
(14)  if  $\text{Area}(QAR_i) + \text{Area}(QAR_j) > \text{Area}(QAR_{i,j}), QAR_{i,j} = \text{argmin}\{\text{Area}(QAR_{i,j})\}_{i \neq j}$ ,  $\text{centre}(QAR_{i,j}) \neq L_i(x_i, y_i)$  and  $\text{centre}(QAR_{i,j}) \neq L_j(x_j, y_j)$  then
(15)    $CS \leftarrow CS \setminus \{AR_i, AR_j\}$ ;
(16)  end if
(17) end for
(18) return CS

```

ALGORITHM 1: Anonymous subregion generation algorithm based on querying region division.

used in the existing research of LBS privacy protection, which is based on the Oldenburg map of German city and simulates the location information of users by setting parameters such as the number of moving objects. We set the privacy requirement as k , the generated initial anonymous subregion is rectangular, and its area satisfies $A_{\min} = 160000 \text{ m}^2$. At the same time, to assess the LSP query cost, this paper simulates 500000 points of interest, such as restaurants, hotels, hospitals, parking lots, and so on. In addition, R -tree structure is used to access these points of interest because R -tree is the best-balanced tree for storing high-dimensional data, which can effectively improve the searching efficiency in high-dimensional space [40]. The experimental environment is Intel (R) Core(TM) i7-6700HQ CPU @ 2.60 GHz, 20.0 GB RAM. The algorithms are programmed by Python and the programs run in the Windows 10 Enterprise. The experimental data sets are shown in Table 1.

5.3.1. Problems Existing in Real Anonymity Region Constructions Schemes. To prove that the existing anonymous region construction scheme cannot effectively improve the service quality of LBS query, this paper selects Casper scheme [29] and Fragment scheme [34], respectively, to search nearby points of interest. As the most commonly used anonymous region construction method, the Casper method generates at least an anonymous region that contains all k users. Fragment is to deal with the anonymous region generated by the Casper scheme to reduce the area of the anonymous regions by removing the part that does not contain the user's location according to the location of the user in the anonymous region.

This paper compares the time required to generate the anonymous region, the area of the anonymous region, and the area of the querying region of LSP in the above two

schemes to prove that the existing anonymous region division methods will further reduce the quality of service when there is an overlap querying region between both anonymous regions. In this part of the experiment, this paper sets the user's query radius $r = 500 \text{ m}$. The experimental results are shown in Figures 5 and 6.

As is shown in Figure 5, compared with the Casper scheme, the Fragment scheme uses a region division method to reduce the area of anonymous regions, for example, when $k = 25$, the area of anonymous regions generated by the Casper scheme is $5.71 \times 10^7 \text{ m}^2$, and the time of generating the anonymous regions is 180.275 ms. The area of anonymous regions generated by the Fragment scheme is $3.41 \times 10^7 \text{ m}^2$, and the time of generating the anonymous regions is 175.331 ms. However, in the Casper and Fragment schemes, the time required for the LSP to query the points of interest is, respectively, 9.940 s (the corresponding area of querying regions is $7.234 \times 10^7 \text{ m}^2$) and 10.501 s (the corresponding area of querying regions is $7.329 \times 10^7 \text{ m}^2$). Therefore, when the anonymous server adopts the Casper and Fragment schemes to generate the anonymous regions, the cost time when the user obtains the accurate query results (without considering the transmission delay) is $9.940 + 0.180 = 10.120 \text{ s}$ and $10.501 + 0.175 = 10.676 \text{ s}$, respectively. When the anonymous server adopts the existing region division schemes to construct the anonymous regions, the time for the user to obtain the query result increases instead.

5.3.2. Analysis of Our Scheme's Effectiveness. In this section, we compare our proposed scheme with Casper scheme and prove that the proposed scheme can effectively reduce the query cost of the LSP and improve the query service quality in the smart campus. The experimental results are shown in Figures 5 and 6.

TABLE 1: Experimental data sets.

Status	Id	Reward number	Type id	Time stamp	X-axis	Y-axis	Speed
New point	105	1	0	0	4229.0	16335.0	298.0
New point	106	1	0	0	19065.0	9922.0	132.0
New point	107	1	0	0	3670.0	20230.0	298.0
New point	108	1	1	0	5565.0	18047.0	298.0
New point	109	1	0	0	10567.0	17947.0	298.0
Point	109	2	0	1	10275.4	17638.7	672.0
Point	108	2	1	1	5487.7	17812.9	504.5

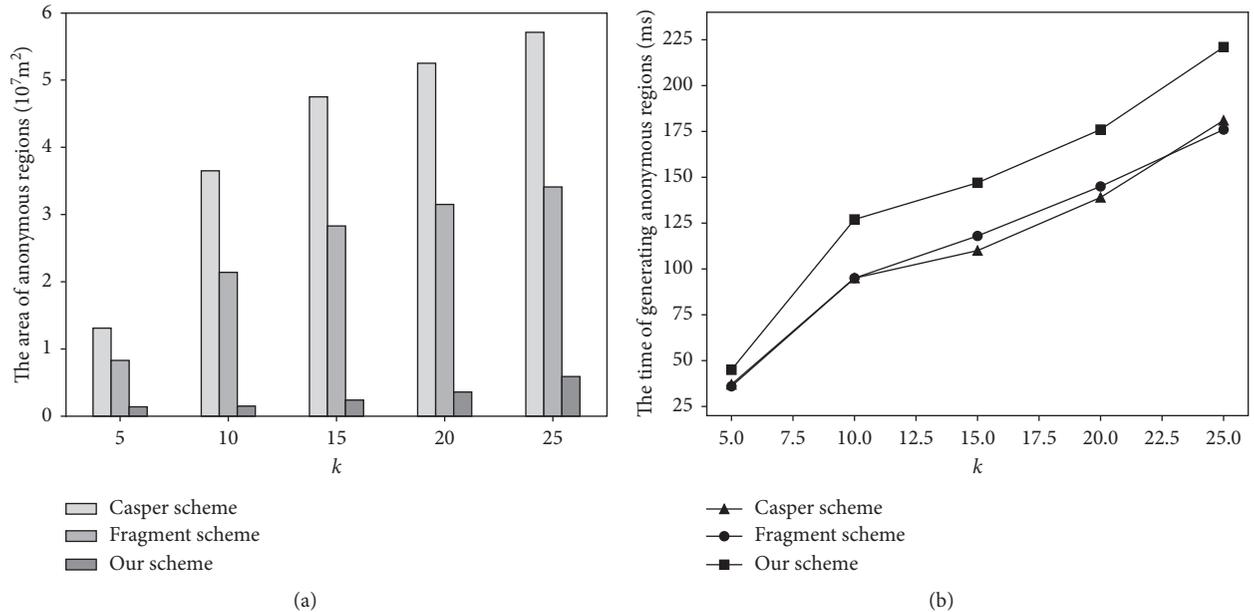


FIGURE 5: Computational cost of anonymity server. (a) The area of anonymous regions. (b) The query cost of LSP.

As is shown in Figures 5 and 6, compared with the existing scheme, the anonymous regions and querying regions generated by our proposed scheme significantly reduced and the same as the time required for LSP to query the points of interest. For example, when $k = 25$, the area of anonymous regions generated by our scheme is $5.89 \times 10^6 \text{ m}^2$, which decreases by $2.80 \times 10^7 \text{ m}^2$ compared to the Casper scheme, and the time to generate the anonymous regions increases from 180.275 ms to 221.034 ms, increasing by 40.759 ms. The area of querying regions generated by our scheme is $2.493 \times 10^7 \text{ m}^2$, which decreases by $4.741 \times 10^7 \text{ m}^2$ compared to the Casper scheme, and the time for query processing increases from 9.940s to 1.961s. Therefore, compared with the Casper scheme, our scheme can effectively improve LBS query service quality.

As can be seen from Figures 5(b) and 6(b), compared with the Casper scheme, as the value of k increases, our scheme needs more extra time to construct the anonymous regions (the time difference between our scheme and the Casper scheme), but less time for query processing. At the same time, the time that the anonymous server generates the anonymous regions is in the milliseconds level, which has little impact on the user delay. The time of query processing is in the second level, which greatly affects the user delay. The

query processing time of the LSP determines the user's service quality to a great extent.

The paper also has made a brief analysis about the influence of the user-specified query radius on the scheme. The experimental results are shown in Figures 7(a) and 7(b). As the user query radius r increases, the number of anonymous subregions that need to be merged also increases, increasing the computational cost of the anonymous server. As a result, the area of querying regions and the time of LSP query processing also increase, which obviously influences the LBS query service quality.

In summary, the proposed scheme not only reduces the computational cost of the anonymous server but also significantly reduces the area of the LSP querying regions and its query time, which effectively improves the LBS query service quality for smart campus.

6. Conclusions

In this paper, we merged novel learning applications technology with mobile technology to research on the location privacy protection technology for smart campus and proposed a new location privacy protection scheme based on querying anonymous region construction which faces the

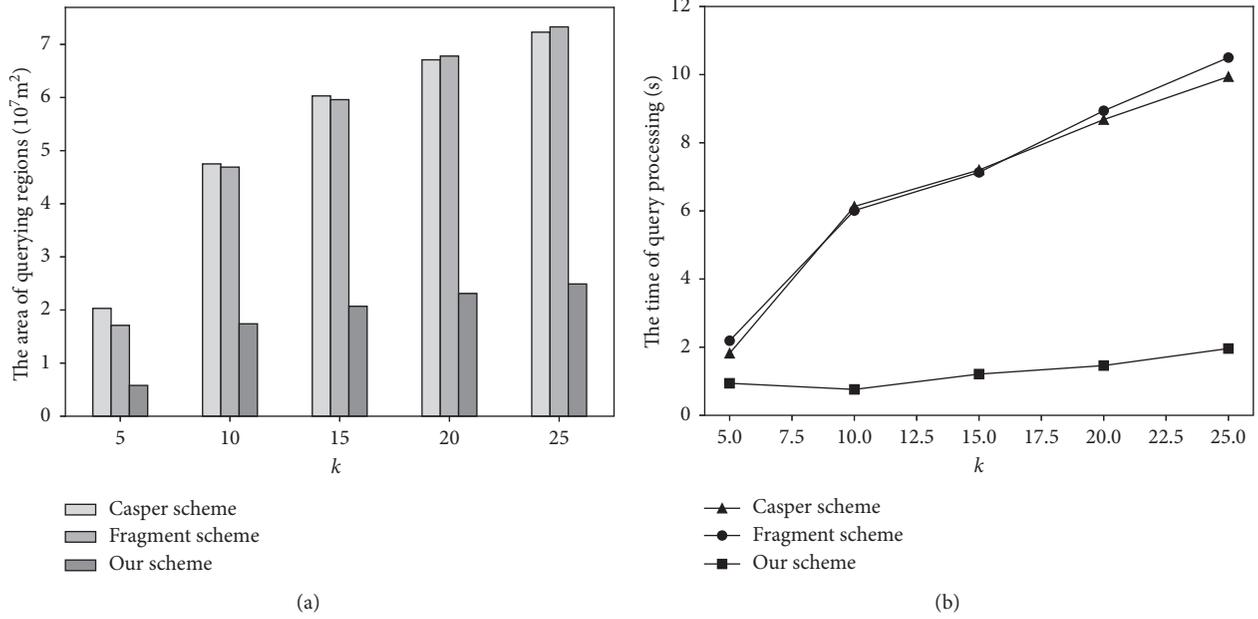


FIGURE 6: The area of querying regions and query cost of LSP. (a) The area of querying regions. (b) The query cost of LSP.

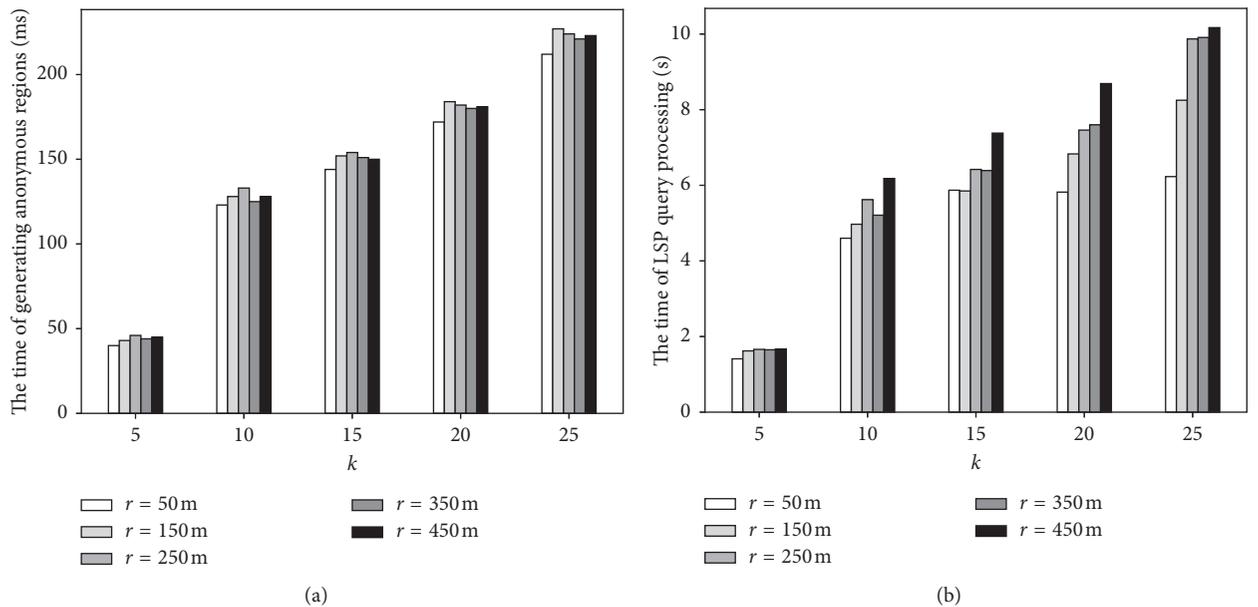


FIGURE 7: The impact of user query radius on the scheme of this paper. (a) The time of generating anonymous regions. (b) The time of LSP query processing.

challenges to achieve the higher quality of smart campus services. Through theoretical analysis and experimental results in the location privacy protection based on k -anonymity, this paper proves that the existing anonymous region construction schemes based on the idea of region division cannot effectively improve the LBS query service quality. The root cause of this problem is that the service quality of sensors in the smart campus is not only related to the anonymous area size constructed by the anonymous server but also to the range of the LSP query. To deal with the

above problem, this paper introduces the user's query range into the construction process of the anonymous region. The anonymous server first generates k initial anonymous subregions based on the user's privacy protection requirement, and then the anonymous subregions will be judged if it needs to be merged based on the size of the querying region. The scheme analysis shows that our proposed scheme can effectively reduce the query cost of LSP and improve the service quality while protecting user's location privacy for smart campus.

Data Availability

The data used to support the findings of this study are available from the corresponding author upon request.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

This work was funded by the National Natural Science Foundation of China (61772282, 61772454, and 61811530332). It was also supported by the Priority Academic Program Development of Jiangsu Higher Education Institutions (PAPD), Postgraduate Research & Practice Innovation Program of Jiangsu Province (KYCX18_1032), Natural Science Foundation of Jiangsu Province (BK20150460), and Jiangsu Collaborative Innovation Center on Atmospheric Environment and Equipment Technology (CICAEET). It was also funded by the open research fund of Key Lab of Broadband Wireless Communication and Sensor Network Technology (Nanjing University of Posts and Telecommunications), Ministry of Education.

References

- [1] J. S. Turner, "New directions in communications," *IEEE Journal on Selected Areas in Communications*, vol. 13, no. 1, pp. 11–23, 1995.
- [2] L. Jin, C. Li, B. Palanisamy, and J. Joshi, "k-trustee: location injection attack-resilient anonymization for location privacy," *Computers & Security*, vol. 78, pp. 212–230, 2018.
- [3] Y. Huang, Z. Cai, and A. G. Bourgeois, "Search locations safely and accurately: a location privacy protection algorithm with accurate service," *Journal of Network and Computer Applications*, vol. 103, pp. 146–156, 2018.
- [4] S. Xiao, "Consideration of technology for constructing Chinese smart grid," *Automation of Electric Power Systems*, vol. 9, no. 33, pp. 1–4, 2009.
- [5] T. Peng, Q. Liu, and G. Wang, "Enhanced location privacy preserving scheme in location-based services," *IEEE Systems Journal*, vol. 11, no. 1, pp. 219–230, 2017.
- [6] X. Chen and Y. Mu, "Preserving user location privacy for location-based service," in *Proceedings of 11th International Conference on Green, Pervasive, and Cloud Computing (GPC 2016)*, pp. 290–300, Xi'an, China, May 2016.
- [7] R. Schlegel, C. Y. Chow, Q. Huang, and D. S. Wong, "User-defined privacy grid system for continuous location-based services," *IEEE Transactions on Mobile Computing*, vol. 14, no. 10, pp. 2158–2172, 2015.
- [8] K. G. Shin, X. Ju, Z. Chen, and X. Hu, "Privacy protection for users of location-based services," *Wireless Communications IEEE*, vol. 19, no. 1, pp. 30–39, 2012.
- [9] M. Gruteser and D. Grunwald, "Anonymous usage of location-based services through spatial and temporal cloaking," in *Proceedings of the 1st International Conference on Mobile Systems, Applications and Services*, pp. 31–42, ACM, San Francisco, CA, USA, May 2003.
- [10] B. Niu, Z. Zhang, X. Li, and H. Li, "Privacy-area aware dummy generation algorithms for location-based services," in *Proceedings of 2014 IEEE International Conference on Communications (ICC)*, pp. 957–962, IEEE, Sydney, NSW, Australia, June 2014.
- [11] X. Shu, D. Yao, and E. Bertino, "Privacy-preserving detection of sensitive data exposure," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 5, pp. 1092–1103, 2015.
- [12] C. Yin, J. Xi, R. Sun, and J. Wang, "Location privacy protection based on differential privacy strategy for big data in industrial Internet-of-Things," *IEEE Transactions on Industrial Informatics*, vol. 14, no. 8, pp. 3628–3636, 2017.
- [13] Z. Fu, F. Huang, K. Ren, J. Weng, and C. Wang, "Privacy-preserving smart semantic search based on conceptual graphs over encrypted outsourced data," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 8, pp. 1874–1884, 2017.
- [14] Z. Fu, K. Ren, J. Shu, X. Sun, and F. Huang, "Enabling personalized search over encrypted outsourced data with efficiency improvement," *IEEE Transactions on Parallel & Distributed Systems*, vol. 27, no. 9, pp. 2546–2559, 2016.
- [15] C. Yin, S. Zhang, J. Xi, and J. Wang, "An improved anonymity model for big data security based on clustering algorithm," *Concurrency and Computation: Practice and Experience*, vol. 29, no. 7, pp. 1–13, 2017.
- [16] K. Vu, R. Zheng, and J. Gao, "Efficient algorithms for K-anonymous location privacy in participatory sensing," in *Proceedings of 31st Annual IEEE International Conference on Computer Communications (INFOCOM, 2012)*, pp. 2399–2407, IEEE, Orlando, FL, USA, March 2012.
- [17] T. Ma, Y. Zhang, J. Cao et al., "KDDEM: a k-degree anonymity with vertex and edge modification algorithm," *Computing*, vol. 97, no. 12, pp. 1165–1184, 2015.
- [18] R. Zhang, Y. Zhang, and C. Zhang, "Secure top-k query processing via untrusted location-based service providers," in *Proceedings of 31st Annual IEEE International Conference on Computer Communications (INFOCOM, 2012)*, pp. 1170–1178, IEEE, Orlando, FL, USA, March 2012.
- [19] R. Zhang, J. Sun, Y. Zhang, and C. Zhang, "Secure spatial top-k query processing via untrusted location-based service providers," *IEEE Transactions on Dependable and Secure Computing*, vol. 12, no. 1, pp. 111–124, 2015.
- [20] C. Y. Chow, M. F. Mokbel, and W. G. Aref, "Casper: query processing for location services without compromising privacy," *ACM Transactions on Database Systems (TODS)*, vol. 34, no. 4, pp. 1–48, 2009.
- [21] C. Yin, J. Xi, and R. Sun, "Location privacy protection based on improved-value method in augmented reality on mobile devices," *Mobile Information Systems*, vol. 2017, Article ID 7251395, 7 pages, 2017.
- [22] C. Yin and J. Xi, "Maximum entropy model for mobile text classification in cloud computing using improved information gain algorithm," *Multimedia Tools and Applications*, vol. 76, no. 16, pp. 16875–16891, 2017.
- [23] H. Rong, T. Ma, M. Tang, and J. Cao, "A novel subgraph K+-isomorphism method in social network based on graph similarity detection," *Soft Computing*, vol. 22, no. 8, pp. 2583–2601, 2017.
- [24] B. Gu, X. Sun, and V. S. Sheng, "Structural minimax probability machine," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 7, pp. 1646–1656, 2017.
- [25] B. Gu, V. S. Sheng, K. Y. Tay, W. Romano, and S. Li, "Incremental support vector learning for ordinal regression," *IEEE Transactions on Neural Networks & Learning Systems*, vol. 26, no. 7, pp. 1403–1416, 2015.
- [26] R. Lu, X. Lin, X. Liang, and X. Shen, "A dynamic privacy-preserving key management scheme for location-based

- services in VANETs,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 13, no. 1, pp. 127–139, 2012.
- [27] A. R. A. Bouguettaya and M. Y. Eltoweissy, “Privacy on the web: facts, challenges, and solutions,” *IEEE Security & Privacy*, vol. 99, no. 6, pp. 40–49, 2003.
- [28] A. Lee, “Guidelines for smart grid cyber security,” NIST Interagency/Internal Report (NISTIR)-7628, National Institute of Standards and Technology, Gaithersburg, MD, USA, 2010.
- [29] M. F. Mokbel, C. Y. Chow, and W. G. Aref, “The new Casper: a privacy-aware location-based database server,” in *Proceedings of IEEE 23rd International Conference on Data Engineering (ICDE 2007)*, pp. 1499–1500, IEEE, Istanbul, Turkey, April 2007.
- [30] X. Li, E. Wang, W. Yang, and J. Ma, “DALP: a demand-aware location privacy protection scheme in continuous location-based services,” *Concurrency and Computation: Practice and Experience*, vol. 28, no. 4, pp. 1219–1236, 2016.
- [31] B. Gedik and L. Liu, “Protecting location privacy with personalized k-anonymity: architecture and algorithms,” *IEEE Transactions on Mobile Computing*, vol. 7, no. 1, pp. 1–18, 2008.
- [32] P. Kalnis, G. Ghinita, K. Mouratidis, and P. Dimitris, “Preventing location-based identity inference in anonymous spatial queries,” *IEEE Transactions on Knowledge & Data Engineering*, vol. 19, no. 12, pp. 1719–1733, 2007.
- [33] K. W. Tan, Y. Lin, and K. Mouratidis, “Spatial cloaking revisited: distinguishing information leakage from anonymity,” in *Proceedings of 11th International Symposium on Spatial and Temporal Databases*, pp. 117–134, Aalborg, Denmark, July 2009.
- [34] T. C. Li and W. T. Zhu, “Protecting user anonymity in location-based services with fragmented cloaking region,” in *Proceedings of 2012 IEEE International Conference on Computer Science and Automation Engineering (CSAE)*, pp. 227–231, IEEE, Zhangjiajie, China, May 2012.
- [35] D. Steiert, D. Lin, Q. Conduff, and W. Jiang, “Poster: a location-privacy approach for continuous queries,” in *Proceedings of the 22nd ACM on Symposium on Access Control Models and Technologies*, pp. 115–117, ACM, Indianapolis, IN, USA, June 2017.
- [36] E. H. Yilmaz, E. Ferhatosmanoglu, and R. C. Aksoy, “Privacy-preserving aggregate queries for optimal location selection,” *IEEE Transactions on Dependable and Secure Computing*, vol. 2017, no. 99, p. 1, 2017.
- [37] I. Memon, “Authentication user’s privacy: an integrating location privacy protection algorithm for secure moving objects in location based services,” *Wireless Personal Communications*, vol. 82, no. 3, pp. 1585–1600, 2015.
- [38] Y. Wang, D. Xu, and F. Li, “Providing location-aware location privacy protection for mobile location-based services,” *Tsinghua Science and Technology*, vol. 21, no. 3, pp. 243–259, 2016.
- [39] T. Brinkhoff, “A framework for generating network-based moving objects,” *GeoInformatica*, vol. 6, no. 2, pp. 153–180, 2002.
- [40] M. Hadjieleftheriou, Y. Manolopoulos, Y. Theodoridis, and V. J. Tsotras, “R-trees—a dynamic index structure for spatial searching,” in *Encyclopedia of GIS*, pp. 993–1002, Springer, Boston, MA, USA, 2008.