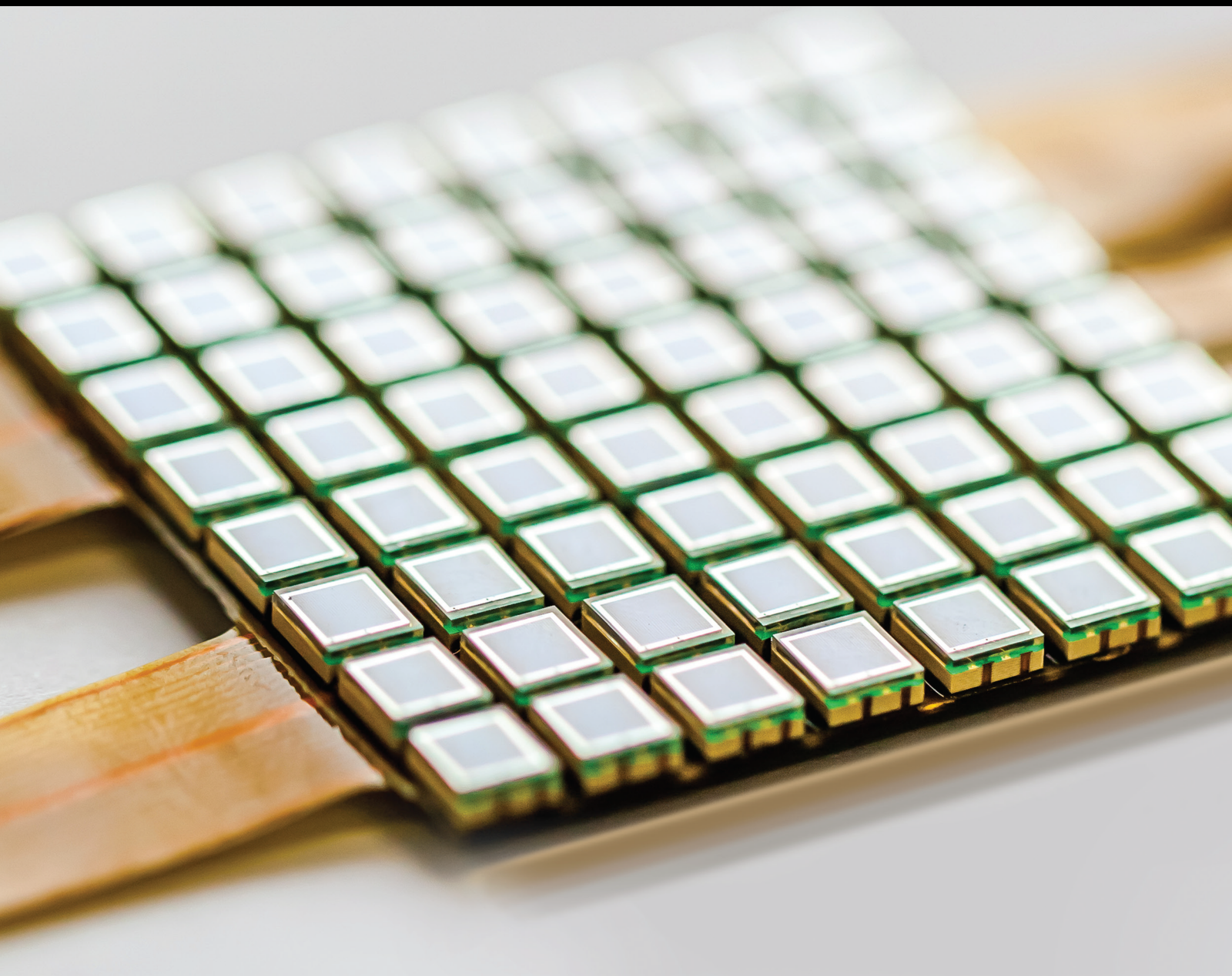# Deep Learning and Artificial Intelligence for Non-Vision Sensors and Imaging

Lead Guest Editor: Yunze He
Guest Editors: Hongjin Wang, Tomasz Chady, and Ruizhen Yang

# Deep Learning and Artificial Intelligence for Non-Vision Sensors and Imaging

# Deep Learning and Artificial Intelligence for Non-Vision Sensors and Imaging

Lead Guest Editor: Yunze He
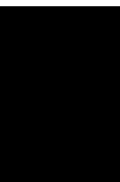Guest Editors: Hongjin Wang, Tomasz Chady, and Ruizhen Yang

Ehsan Namaziandost (iD), Iran
Heinz C. Neitzert (iD), Italy
Sing Kiong Nguang (iD), New Zealand
Calogero M. Oddo (iD), Italy
Tinghui Ouyang, Japan
SANDEEP KUMAR PALANISWAMY (iD),
India
Alberto J. Palma (iD), Spain
Davide Palumbo (iD), Italy
Abinash Panda (iD), India
Roberto Paolesse (iD), Italy
Akhilesh Pathak (iD), Thailand
Giovanni Pau (iD), Italy
Giorgio Pennazza (iD), Italy
Michele Penza (iD), Italy
Sivakumar Poruran, India
Stelios Potirakis (iD), Greece
Biswajeet Pradhan (iD), Malaysia
Giuseppe Quero (iD), Italy
Linesh Raja (iD), India
Maheswar Rajagopal (iD), India
Valerie Renaudin (iD), France
Armando Ricciardi (iD), Italy
Christos Riziotis (iD), Greece
Ruthber Rodriguez Serrezuela (iD), Colombia
Maria Luz Rodriguez-Mendez (iD), Spain
Jerome Rossignol (iD), France
Maheswaran S, India
Ylias Sabri (iD), Australia
Sourabh Sahu (iD), India
José P. Santos (iD), Spain
Sina Sareh, United Kingdom
Isabel Sayago (iD), Spain
Andreas Schütze (iD), Germany
Praveen K. Sekhar (iD), USA
Sandra Sendra, Spain
Sandeep Sharma , India
Sunil Kumar Singh Singh (iD), India
Yadvendra Singh (iD), USA
Afaque Manzoor Soomro (iD), Pakistan
Vincenzo Spagnolo, Italy
Kathiravan Srinivasan (iD), India
Sachin K. Srivastava (iD), India
Stefano Stassi (iD), Italy

Danfeng Sun, China
Ashok Sundramoorthy, India
Salvatore Surdo (iD), Italy
Roshan Thotagamuge (iD), Sri Lanka
Guiyun Tian (iD), United Kingdom
Sri Ramulu Torati (iD), USA
Abdellah Touhafi (iD), Belgium
Hoang Vinh Tran (iD), Vietnam
Aitor Urrutia (iD), Spain
Hana Vaisocherova - Lisalova (iD), Czech
Republic
Everardo Vargas-Rodriguez (iD), Mexico
Xavier Vilanova (iD), Spain
Stanislav Vítek (iD), Czech Republic
Luca Vollero (iD), Italy
Tomasz Wandowski (iD), Poland
Bohui Wang, China
Qihao Weng, USA
Penghai Wu (iD), China
Qiang Wu, United Kingdom
Yuedong Xie (iD), China
Chen Yang (iD), China
Jiachen Yang (iD), China
Nitesh Yelve (iD), India
Aijun Yin, China
Chouki Zerrouki (iD), France

# Contents

*Research Article*

# Multirobot Adaptive Task Allocation of Intelligent Warehouse Based on Evolutionary Strategy

**Yifan Liu** [ID],[1] **Fei Liu** [ID],[1] **Li Tang** [ID],[2] **Chuanzheng Bai** [ID],[1] **and Li Liu** [ID][1]

[1]*School of Information and Electrical Engineering, Ludong University, Yantai 264025, China*
[2]*School of Physics and Optoelectronic Engineering, Ludong University, Yantai 264025, China*

Correspondence should be addressed to Fei Liu; liufeildu@163.com

To solve the dynamic and real-time problem of multirobot task allocation in intelligent warehouse system under parts-to-picker mode, this paper presents a combined solution based on adaptive task pool strategy and Covariance Matrix Adaptation Evolutionary Strategy (CMA-ES) algorithm. In the first stage of the solution, a variable task pool is used to store dynamically added tasks, which can dynamically divide continuous and large-scale task allocation problems into small-scale subproblems to solve them to meet dynamic requirements. And an adaptive control strategy is used to automatically adjust the total number of tasks in the task pool to achieve a trade-off among throughput, energy consumption, and waiting time, which has better adaptability than manually adjusting the size of the task pool. In the second stage of the solution, when the task pool is full, tasks in the task pool will be assigned to robots. For the task allocation problem, this paper regards it as an optimization problem and uses the CMA-ES algorithm to find the optimal task assignment solution for all the robots. By comparing with fixed threshold method under 56 different task pool sizes, the experimental results show that the throughput can be close to reaching the optimal level, and the average distance traveled by robots to handle each unit is lower using adaptive threshold method; so, adaptive task pool solution has better adaptability and can find the optimal task pool size by itself. This method can satisfy the dynamic and real-time requirements and can be effectively applied to the intelligent warehouse system.

## 1. Introduction

In recent years, the orders of various e-commerce platforms have soared, and the scale of distribution centers has become increasingly large, which has brought great challenges to the traditional logistics industry [1]. In the traditional warehouse, 60% to 70% of the workers' time is spent on picking up goods [2], and the efficiency is extremely low. Therefore, more and more automatic machines and equipment have been applied in the field of warehouse [2]. Many companies have started to adopt a new kind of parts-to-picker intelligent warehouse system, such as Kiva system [3]. In the system as shown in Figure 1, robots transport the shelves from storage areas to workstations, and workers need to wait at the stations. When the shelves reach the workstations, they take the needed goods from the shelves or store bundles into the shelves. It has been proved that this kind of the intelligent warehouse system greatly saves labor cost and improves the efficiency of warehouse operation [4].

Cooperative control of multiple mobile robots is the key to realize intelligent warehousing. In a warehouse as shown in Figure 1, there are often numerous tasks such as replenishment and picking, as well as numerous robots to perform these tasks. In addition, the costs of different robots to perform a task are also different. Therefore, the efficiency of the warehouse is determined by selecting suitable robots to perform specific tasks. This is a typical multirobot task allocation (MRTA) problem [5]. With the operation of the warehouse, tasks and the warehouse environment will constantly change. How to find a better task allocation scheme for pick-task and replenishment-task assignment in such a highly dynamic environment [3, 4] is the focus of this paper.

FIGURE 1: Parts-to-picker intelligent warehouse system from ref. [4].

MRTA is one of the most challenging problems in the multirobot system [6]. Market-based methods are the most studied methods at present, such as the single-task auction algorithm proposed in ref. [7]. In order to solve the problem that the single-task auction algorithm is difficult to get the optimal solution, a combined auction algorithm which considers the correlation between tasks was proposed in ref. [8]. When the number of robots and tasks is small, MRTA can be regarded as a zero-one integer linear programming problem and solved by simplex method, branch and bound method, Hungarian algorithm [9], etc. For example, the Hungarian algorithm was adopted in ref. [10] to solve the role assignment problem in robot soccer game. There are also some thresholding based methods such as ALLIANCE [11] and Broadcast of Local Eligibility (BLE) [12], which have good real-time, fault tolerance, and robustness, but usually only local optimal solution can be obtained. For large-scale problems, the heuristic algorithm can effectively reduce solution space and improve search efficiency. For example, in ref. [13], the heuristic algorithm was adopted to solve the task assignment problem in multi-core processor. Evolutionary algorithms are mature global optimization methods with high robustness and wide applicability, which can effectively deal with complex problems that are difficult to be solved by traditional optimization algorithms. Various evolutionary algorithms such as genetic algorithm and simulated annealing algorithm have been widely used in MRTA problem. In ref. [14], the genetic algorithm was used to solve the time-extended multirobot task allocation problem in the case of disaster. A hybrid genetic and ant colony algorithm was proposed in ref. [15] to improve the solving accuracy of the genetic algorithm. In ref. [16], the genetic algorithm was used to solve MRTA problem in the intelligent warehouse. Ref. [17] designed an improved quantum evolutionary algorithm based on the niche coevolution strategy and enhanced particle swarm optimization (IPOQEA) to solve the airport gate allocation problem. In ref. [18], an improved quantum-inspired cooperative coevolution algorithm with

multistrategy is used to solve the knapsack problem and the actual airport gate allocation problem. Refs. [17–20] use the cooperative coevolution framework to divide the complex optimization problem into several subproblems, and these subproblems were solved by independent searching in order to improve the solution efficiency. Similarly, the situation where the number of tasks is variable in an intelligent warehouse can be studied using the idea of divide-and-conquer in Refs. [17–20].

Therefore, we use a task pool to store dynamically added tasks and propose an adaptive control strategy to automatically adjust the task pool size according to the current environment. When the task pool is full, the tasks in the pool will be assigned to the robots. Then, the task allocation problem is regarded as an optimization problem and solved by the CMA-ES algorithm [21].

## 2. Problem Formulation

The intelligent warehouse system consists of many movable shelves and robots as well as some workstations. The robots transport the needed shelves from the storage area to the workstations, and the workers can complete the replenishment and picking without moving. A typical intelligent warehouse layout (a screenshot from the open source software RAWSim-O [22]) is shown in Figure 2. In the figure, the four squares on the left represent the replenishment station, and the replenished bundles are temporarily stored here waiting for shelves. The four squares on the right represent picking stations. After receiving orders, the system will use a special algorithm to assign orders to different stations. There will be an upper limit on the number of orders in the stations [23]. The squares in the middle area are the shelves, in which the goods in the warehouse are stored. Shelves can be lifted and moved by robots. The circles in the figure are robots. A robot can carry a shelf to move. When a robot does not carry a shelf, it can move freely under the shelf.

FIGURE 2: A typical intelligent warehouse layout from ref. [22].

In order to facilitate problem analysis, we make the following assumptions:

(1) Robots are all isomorphic and travel at exactly the same speed. They can only move forward, backward, left, and right.

(2) The time for a robot to lift a shelf and stay at a workstation is very short, which can be ignored.

(3) Every robot carries the required shelf and travels from the position of the shelf to the designated station and then carries the shelf back to its original location.

The shelf selection algorithm will select shelves for each workstation according to requirements. The selected shelves need to be transported from the shelf storage area to the appropriate station for picking up or replenishing goods, and then they are transported back to the original position, which is the task of the robots. If a robot is not assigned a task, it will move to a special resting area for rest. How to reasonably assign tasks to robots is the problem to be studied in this paper.

Referring to ref. [16], suppose that there are $m$ tasks (refers to all tasks from the beginning to the end of the warehouse operation) and $n$ robots in the warehouse, the set of tasks is $T = \{t_1, t_2, t_3, \cdots, t_m\}$, and the set of robots is $R = \{r_1, r_2, r_3 \cdots, r_n\}$. The set of tasks assigned to robot $r_i$ is $T_i$, which is a subset of $T$. $T_1 \cup T_2 \cup T_3 \cup \cdots \cup T_n = T$ and $T_1 \cap T_2 \cap T_3 \cap \cdots \cap T_n = \varnothing$. Let $T_i = \{t_{i1}, t_{i2}, t_{i3}, \cdots, t_{ik}\}$ and $T_i$ is ordered, and then the sequence of tasks to be completed by the robot $r_i$ is $t_{i1} \longrightarrow t_{i2} \longrightarrow t_{i3} \longrightarrow \cdots \longrightarrow t_{ik}$. The cost of robot $r$ to complete its task sequence can be expressed as

$$C(r_i) = I(r_i, t_{i1}) + \sum_{h=1}^{k} S(t_h) + \sum_{h=1}^{k-1} R(t_h, t_{h+1}), \quad (1)$$

where $C(r_i)$ represents the cost of the robot $r_i$ to complete all tasks. Since all robots travel at the same speed, the cost can be expressed as the distance traveled by the robot. The robot can only move forward, backward, left, and right; so, the distance traveled between the two points can be expressed as Manhattan distance.

$I(r_i, t_{i1})$ represents the cost for the robot to get from the initial position to the position of required shelf for the first task $t_{i1}$. Let the initial coordinate of the robot be $(x_r, y_r)$ and the coordinate of the required shelf for the first task be $(x_{t1}, y_{t1})$, and then

$$I(r_i, t_{i1}) = |x_r - x_{t1}| + |y_r - y_{t1}|. \quad (2)$$

$S(t_h)$ represents the cost for the robot to complete task $t_h$, which is only related to task $t_h$ itself. It can be represented by the distance that after the robot carries the required shelf, it travels from the position of the required shelf for the task to the designated station and then returns to the shelf's original position from the station. Let the coordinate of required shelf for task $t_h$ be $(x_p, y_p)$ and the coordinate of target station be $(x_s, y_s)$, and then

$$S(t_h) = \left( |x_p - x_s| + |y_p - y_s| \right) * 2. \quad (3)$$

$R(t_h, t_{h+1})$ represents the cost for the robot to reach the starting position of the next task $t_{h+1}$ after completing task $t_h$. Since the robot needs to transport the shelf back to the original position after completing task $t_h$, it can be directly represented by the Manhattan distance from the position of required shelf for task $t_h$ to the position of required shelf for task $t_{h+1}$. Let the coordinate of required shelf for task $t_h$ be $(x_{p1}, y_{p1})$ and the coordinate of required shelf for task $t_{h+1}$ be $(x_{p2}, y_{p2})$, and then

$$R(t_h, t_{h+1}) = \left| x_{p1} - x_{p2} \right| + \left| y_{p1} - y_{p2} \right|. \tag{4}$$

In order to make the overall allocation scheme as optimal as possible, we consider the following two optimization objectives:

(1) The maximum time taken by all robots to complete all tasks ($C_{\text{time}}$)

(2) The mean distance traveled by all robots ($C_{\text{distance}}$)

where

$$C_{\text{time}} = \max_i C(r_i),$$
$$C_{\text{distance}} = \frac{\sum_{i=1}^n C(r_i)}{n}. \tag{5}$$

$C_{\text{time}}$ describes the efficiency of the robots to complete tasks. The smaller $C_{\text{time}}$ is, the less time the robots take to complete all tasks, and the higher the efficiency is. $C_{\text{distance}}$ describes the power consumption of the multirobot system. The smaller $C_{\text{distance}}$ is, the shorter the total travel distance of all robots is, and the lower the power consumption is. The goal of the method studied in this paper is to reasonably assign all tasks in the system to all robots so that these two values can be as small as possible.

# 3. Method

## 3.1. Architecture.
With the entry of new orders, new tasks are constantly generated and must be completed as soon as possible; so, the warehouse system is a highly dynamic and real-time system. In such a highly dynamic system, it is difficult to find the global optimal solution; so, the problem is divided into many subproblems. Specifically, we created a task pool $P$. When a new task is generated, it is immediately added to $P$. When the number of tasks in the task pool $P$ reaches the threshold value (automatic adjustment of the threshold will be described in Section 3.3), the CMA-ES method in Section 3.2 is used to allocate the tasks in the task pool to robots. The robots insert the new task sequence allocated into the rear of the previous unfinished task sequence, and then the task pool is emptied. The robots execute tasks according to their own task sequence, and the executed tasks are deleted from the sequence. As the new tasks are generated again, the tasks are added to $P$ again. Loop until the warehouse stops running. In Figure 3, the specific steps are as follows:

*Step 1.* Initialize the task pool size and set the task pool $P$ to be empty. For all robots, initialize task sequence $T_i$ of every robot $r_i$.

*Step 2.* The threshold of the task pool size is automatically adjusted using adaptive control strategy in Section 3.3.

*Step 3.* New tasks are constantly added to $P$. Jump to step 4 when the number of tasks in the task pool reaches the threshold.

*Step 4.* The tasks in the task pool are assigned to the robots using the CMA-ES method in Section 3.2, and for all robots, the new task sequence assigned to robot $r_i$ is inserted at the end of the current task sequence $T_i$.

*Step 5.* Clear the task pool $P$ and jump to step 2.

The above solution in Figure 3 is executed by the central controller, and the robot only needs to execute the tasks according to the assigned task sequence. The parallel operation of the two parts enables the robots to be busy all the time, which saves time and meets the requirement of real-time storage system.

## 3.2. CMA-ES Algorithm.
As mentioned in Section 3.1, tasks are assigned to robots when the number of tasks in the task pool reaches the threshold. This problem is regarded as an optimization problem in a static environment. This is a NP-hard problem, and the CMA-ES algorithm is used to find the optimal solution. The successful application in many fields [24–26] proves that the CMA-ES algorithm is a good search algorithm.

### 3.2.1. Representation of Solutions.
Referring to ref. [27], for the task allocation problem with $m$ tasks and $n$ robots, a candidate to represent a task assignment scheme is $X = [x_1, x_2, x_3 \cdots x_m]$. $X$ contains $m$ real numbers, and for each real number $x_i$, it satisfies $1 \leq x_i < n+1, i = 1, 2, 3, \cdots, m$, where $x_i$ means task $i$ is performed by robot $\text{Int}(x_i)$, and $\text{Int}(x_i)$ means the integer of real number $x_i$. If $\text{Int}(x_i) = \text{Int}(x_j), i \neq j$, this means that the task $x_i$ and $x_j$ are both assigned to the same robot, and the task represented by the smaller number between $x_i$ and $x_j$ is executed first. If $x_i = x_j$, the execution order of these two tasks is determined randomly.

For example, there are 8 tasks (represented by numbers 1, 2, 3,..., 8) and 3 robots (represented by numbers 1, 2, 3), and an individual [1.7, 3.8, 2.2, 1.3, 2.8, 1.5, 3.3, 3.7] is generated. Then, the task sequence assigned to robot 1 is $4 \longrightarrow 6 \longrightarrow 1$. The task sequence assigned to robot 2 is $3 \longrightarrow 5$. The task sequence assigned to robot 3 is $7 \longrightarrow 8 \longrightarrow 2$.

### 3.2.2. Fitness Function.
Fitness function is used to evaluate candidates. For the CMA-ES algorithm, individuals with lower fitness value are more excellent. In Section 2, two optimization goals are proposed for the whole system: one is the time $C_{\text{time}}$ for the robots to complete all tasks; the second is the mean driving distance $C_{\text{distance}}$ of all robots. Each planning can be regarded as a subproblem of the whole. For each subproblem, in order to achieve the optimal overall performance, these two goals are still considered; so, fitness function $f$ is calculated through the following equation [16]:

FIGURE 3: The flow chart of the combined solution based on adaptive task pool strategy and CMA-ES.

$$f = \alpha C'_{\text{time}} + (1 - \alpha)C'_{\text{distance}}, 0 \le \alpha \le 1,$$

$$C'_{\text{time}} = \max_i C'(r_i), \quad (6)$$

$$C'_{\text{distance}} = \frac{\sum_{i=1}^{n} C'(r_i)}{n},$$

where $\alpha$ is a constant that can be adjusted according to the actual demand. If more attention is paid to the completion time of a single order, $\alpha$ can be increased. If more attention is paid to the energy consumption of all robots, $\alpha$ can be reduced. $C'(r_i)$ is the cost of robot $r_i$ to execute the tasks in the current task sequence first and then execute the tasks according to the candidate. $C'_{\text{time}}$ is the maximum time taken by the robots. $C'_{\text{distance}}$ is the mean distance traveled by all robots. In the current moment, there may be unfinished tasks in the task sequence. The robot must first complete these tasks before performing the tasks assigned at the current moment. Therefore, for $C'(r_i)$, we divide it into two parts to calculate:

$$C'(r_i) = C'_1(r_i) + C'_2(r_i), \quad (7)$$

where $C'_1(r_i)$ is the cost for the robot to complete the tasks in the current task sequence, and $C'_2(r_i)$ is the cost for the robot

to execute the tasks according to the candidate. $C'_1(r_i)$ and $C'_2(r_i)$ are represented by the distance traveled by the robot and calculated using the method described in Equation (1).

With this fitness function, we try to find the optimal solution at that moment in each optimization and try to approximate the global optimal solution by this method.

3.3. *Automatic Adjustment of Task Pool.* When the number of tasks in the task pool reaches the threshold, the tasks in the task pool will be assigned to the robots. The threshold plays a decisive role in the efficiency of assignment. The larger the threshold is, the more tasks will be involved in the optimization, and then the more the planned scheme will be close to the global optimal solution. If an optimization contains all the tasks in the system, the optimal solution found by the optimization will be the optimal solution of the whole system. But orders in the warehouse are added dynamically over time, so tasks are also generated dynamically. As the threshold increases, the time required for the task pool to be filled will also increase, and this situation will occur: the robot has finished all the tasks assigned to it, but the number of tasks in the task pool has not reached the threshold; so, the next optimization cannot start, and the robot can only wait. This leads to a waste of time and cannot meet the real-time of the warehouse system. Moreover,

**Input:** lastAdjustTime, currentTime, lastTasksCompleted, tasksCompleted, oldThreshold, lastAction
**Output:** newThreshold, lastAction
1: **if** *currentTime − lastAdjustTime > I* **then**
2:    **if** *tasksCompleted* = 0 **then**
3:        *newThreshold* ⟵ *oldThreshold*/2
4:        *lastAction* ⟵ −1
5:    **else if** *tasksCompleted − lastTaskCompleted* ≥ 0 **then**
6:        *newThreshold* ⟵ *newThreshold* + *lastAction*
7:    **else**
8:        *newThreshold* ⟵ *newThreshold* − *lastAction*
9:        *lastAction* ⟵ −*lastAction*
10: **else**
11:    *newThreshold* ⟵ *oldThreshold*
12: **return** newThreshold, lastAction

ALGORITHM 1: Adaptive control strategy.

because each workstation has an order capacity limit, there is also an upper limit on the total number of tasks in the system, and if the task pool size exceeds this upper limit, the number of tasks in the task pool will never reach the threshold, and the system will be stagnant. Therefore, it is very important to set a threshold of appropriate size.

Obviously, for different warehouses, the threshold should be set differently depending on the actual situation. Even for the same warehouse, the number of robots may be adjusted, and the rate of order generation may vary at different times; so, it is not appropriate to set the threshold to a fixed value. Therefore, we design an adaptive control strategy to dynamically adjust the task pool, as shown in Algorithm 1.

First, the setting of the initial threshold is important, which determines the speed of finding the optimal threshold. We believe that the size of the initial threshold should be related to the number of robots and the upper limit number of tasks in the warehouse. The upper limit number of tasks in the warehouse is related to the number of workstations and the capacity of each workstation. So, we propose the following heuristic formula to calculate the initial threshold:

$$\text{initialThreshold} = \frac{(\gamma * \text{stations} + \text{robots})}{2}, \qquad (8)$$

where $\gamma$ is a constant representing the average number of tasks per workstation in unit time, which is set according to the actual situation. stations is the number of stations, and robots is the number of robots. We set a time interval $I$ (It is a constant that can be set according to actual requirements), and every $I$ seconds, the threshold is adjusted (line 1). lastAction is used to record the last adjustment. We counted the total number of tasks completed by the robot from the last adjusted moment to the current moment, and the total number of tasks completed from the penultimate adjusted moment to the last adjusted moment, expressed by tasksCompleted and lastTasksCompleted, respectively. If taskCompleted is 0, indicating that the threshold has been set so high that the number of tasks has not reached the threshold, then simply cut the threshold in half and set lastAction to −1 (line 2, line 3, and line 4). If tasksCompleted is greater than or equal to

lastTasksCompleted, it indicates that the last adjustment has had a positive effect on the system, and the same adjustment will be performed (line 5 and line 6). If tasksCompleted is less than lastTasksCompleted, it indicates that the last adjustment had a negative effect on the system, and the reverse adjustment will be performed (line 7 and line 8). In addition, lastAction will be reversed (line 9).

## 4. Experiments

We used RAWSim-O [22], an open source framework developed by Merschformann et al., as the experimental platform. RAWSim-O is a simulation framework that simulates the operation of an intelligent warehouse system and allows us to test our own methods.

We used the warehouse layout shown in Figure 2. In the warehouse layout, there are 32 robots and 550 shelves. The storage positions of the shelves are at the middle area of the layout. And there are four replenishment stations on the left and four picking stations on the right. To simplify the problem, we set the duration of a robot staying at a workstation to a very small value of 0.1.

For the assessment of performance we take the sum of SKUs (stock keeping unit) in both item bundles stored at the replenishment stations and orders picked at the picking stations as handled units. This represents the throughput of the warehouse, and the higher the better. We also look at the average distance traveled by robots to handle each unit. This can represent the power consumption of the multirobot system.

In order to test the impact of task pool threshold size on the allocation effect, we did 56 experiments, each experiment corresponding to different pool sizes. Each experiment was simulated for 24 hours with 10 repetitions.

Under different task pool sizes, the number of units handled by robots is shown in the blue solid line in Figure 4, and the average distance traveled by robots to handle each unit is shown in the blue solid line in Figure 5. The comparison results among different fixed threshold on handled units and travel distance per unit are shown in Table 1. The maximum number of handled units is 207583 when the fixed threshold is set to 18. The minimum number of travel

FIGURE 4: Comparison between adaptive threshold method and fixed threshold method on handled units. The red dotted line is the adaptive threshold method, and the blue solid line is the fixed threshold method.

distance per unit is 10.73 when the fixed threshold is set to 36, 45, or 47. According to Figures 4 and 5 and Table 1, it is not good to set the threshold too large or too small, which is consistent with our conjecture. If the threshold is set too small, the solution will be too far away from the global optimal solution; therefore, the number of handled units is small, and the travel distance per unit is large. If the threshold is set too large, the solution will be closer to the global optimal solution; so, the travel distance per unit is small, but the robot will have a long waiting time; therefore, the number of handled units will be small.

To sum up, a bad threshold can be very inefficient; so, setting the threshold manually is very risky. Therefore, a method of automatically adjusting threshold is necessary. We used the adaptive control strategy proposed by ourselves to conduct the experiment again, and all conditions were identical except the threshold. According to the workstation capacity, $\gamma$ in Equation (8) was set to 4; so, the initial threshold was calculated as 32. The results are shown in Table 1. We compared the results with the fixed threshold approach,

as shown in Figures 4 and 5. The red dotted line is the adaptive threshold method, and the blue solid line is the fixed threshold method. Compared with fixed threshold 18, the adaptive threshold method gets worse result in handled units but better result in travel distance per unit. Compared with fixed threshold 36, 45, and 47, the adaptive threshold method gets better result in handled units but worse result in travel distance per unit. Taken together, it can be seen from the two figures that the adaptive threshold method can be close to reaching the level when the threshold is set to the optimal in both indexes. The experimental results show that the proposed adaptive control strategy has good application effect.

## 5. Conclusion

In order to solve the dynamic and real-time problem of multirobot task allocation in the intelligent warehouse system, a combined solution based on adaptive task pool strategy and CMA-ES algorithm is proposed in the paper. In the early

FIGURE 5: Comparison between adaptive threshold method and fixed threshold method on travel distance per unit. The red dotted line is the adaptive threshold method, and the blue solid line is the fixed threshold method.

TABLE 1: Comparison between adaptive threshold method and fixed threshold method on handled units and travel distance per unit.

| Method (initial threshold) | Handled units | Travel distance per unit |
|---|---|---|
| Fixed threshold (18) | 207583 | 10.82 |
| Fixed threshold (36) | 204642 | 10.73 |
| Fixed threshold (45) | 201046 | 10.73 |
| Fixed threshold (47) | 200342 | 10.73 |
| Adaptive threshold (32) | 205372 | 10.79 |

stage of the solution, the divide-to-conquer idea is used to design a variable task pool that is used to store dynamically added tasks. The variable task pool is designed to dynamically divide continuous and large-scale task allocation problems into small-scale subproblems to solve them to meet dynamic requirements. And an adaptive control strategy is used to automatically adjust the threshold of the task pool size in real time to achieve a trade-off among throughput, energy consumption, and waiting time, which has better adaptability than manually adjusting the size of the task

pool. In the later stage of the solution, when the task pool is full, tasks in the task pool will be assigned to robots using the CMA-ES algorithm to find the optimal task assignment solution for all the robots according to the fitness function including the maximum time and the mean travel distance required by all robots to complete all the tasks. By comparing with fixed threshold method under 56 different task pool sizes, the experimental results show that the handled units can be close to reaching the optimal level, and the average travel distance per unit is lower using adaptive threshold method; so, adaptive threshold solution indeed has better adaptability. This method can satisfy the dynamic and real-time requirements and can be effectively applied to the intelligent warehouse system.

However, because of the complexity and dynamics of the warehouse environment, it may not be accurate to measure the cost by Manhattan distance. Therefore, how to introduce accurate robot motion model to evaluate the cost will be the next work. Furthermore, the relationships among handled units, travel distance per unit, the maximum time taken by all robots to complete all tasks, and the mean distance traveled by all robots need further study. In addition, the effect of communication quality on allocation is not taken into account and will be deeply studied.

## Data Availability

The data used to support the findings of this study are included within the article.

## Conflicts of Interest

The authors declare that there is no conflict of interest regarding the publication of this paper.

## Acknowledgments

## References

[1] M. Zhou and M. Y. Wang, "Analysis on the development of e-commerce logistics service industry and countermeasures," *Computer and Information Technology*, vol. 20, no. 6, pp. 10–12, 2012.

[2] S. X. Zou, "The present and future of warehouse robot," *Logistics Engineering and Management*, vol. 35, no. 6, pp. 171-172, 2013.

[3] J. J. Enright and P. R. Wurman, "Optimization and coordinated autonomy in mobile fulfillment systems," in *Workshops at the Twenty-Fifth AAAI Conference on Artificial Intelligence*, pp. 33–38, San Francisco, California, 2011.

[4] P. R. Wurman, R. D'Andrea, and M. Mountz, "Coordinating hundreds of cooperative, autonomous vehicles in warehouses," *AI Magazine*, vol. 29, no. 1, p. 9, 2008.

[5] B. P. Gerkey and M. J. Matarić, "A formal analysis and taxonomy of task allocation in multi-robot systems," *International Journal of Robotics Research*, vol. 23, no. 9, pp. 939–954, 2004.

[6] A. Khamis, A. Hussein, and A. Elmogy, "Multi-robot task allocation: a review of the state-of-the-art," *Eds. Cham: Springer International Publishing*, vol. 604, pp. 31–51, 2015.

[7] B. P. Gerkey and M. J. Matarić, "Sold!: auction methods for multirobot coordination," *IEEE Transactions on Robotics and Automation*, vol. 18, no. 5, pp. 758–768, 2002.

[8] M. Berhault, H. Huang, P. Keskinocak et al., "Robot Exploration with Combinatorial Auctions," *In: Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, vol. 2, pp. 1957–1962, 2003.

[9] H. W. Kuhn, "The Hungarian method for the assignment problem," *Naval Research Logistics Quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.

[10] P. MacAlpine, E. Price, and P. Stone, "SCRAM: scalable collision-avoiding role assignment with minimal-makespan for formational positioning," in *Proceedings of the 29th AAAI Conference on Artificial Intelligence*, pp. 2096–2102, Austin, Texas, USA, 2015.

[11] L. E. Parker, "ALLIANCE: an architecture for fault tolerant multirobot cooperation," *IEEE Transactions on Robotics and Automation*, vol. 14, no. 2, pp. 220–240, 1998.

[12] B. B. Werger and M. J. Mataric, "Broadcast of local eligibility: behavior-based control for strongly cooperative robot teams," in *Proceedings of the 4th International Conference on Autonomous Agents*, pp. 21-22, Barcelona, Spain, 2000.

[13] Y. Liu, X. Zhang, H. Li, and D. Qian, "Allocating tasks in multi-core processor based parallel system," in *2007 IFIP International Conference on Network and Parallel Computing Workshops*, pp. 748–753, Liaoning, China, 2007.

[14] E. G. Jones, M. B. Dias, and A. Stentz, "Time-extended multi-robot coordination for domains with intra-path constraints," *Autonomous Robots*, vol. 30, no. 1, pp. 41–56, 2011.

[15] J. Zhang and Y. Q. Cao, "Research on dynamic task allocation for MAS based on hybrid genetic and ant colony algorithm," *Computer Science*, vol. 38, no. S1, pp. 268–270, 2011.

[16] J. J. Dou, C. L. Chen, and P. Yang, "Genetic scheduling and reinforcement learning in multirobot systems for intelligent warehouses," *Mathematical Problems in Engineering*, vol. 2015, 10 pages, 2015.

[17] W. Deng, J. Xu, H. Zhao, and Y. Song, "A novel gate resource allocation method using improved PSO-based QEA," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 3, pp. 1737–1745, 2022.

[18] X. Cai, H. Zhao, S. Shang et al., "An improved quantum-inspired cooperative co-evolution algorithm with muli-strategy and its application," *Expert Systems with Applications*, vol. 171, article 114629, 2021.

[19] W. Deng, J. J. Xu, X. Z. Gao, and H. M. Zhao, "An enhanced MSIQDE algorithm with novel multiple strategies for global optimization problems," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 52, no. 3, pp. 1578–1587, 2022.

[20] W. Deng, S. Shang, X. Cai et al., "Quantum differential evolution with cooperative coevolution framework and hybrid mutation strategy for large scale optimization," *Knowledge-Based Systems*, vol. 224, article 107080, 2021.

[21] N. Hansen, S. D. Müller, and P. Koumoutsakos, "Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (CMA-ES)," *Evolutionary Computation*, vol. 11, no. 1, pp. 1–18, 2003.

[22] M. Merschformann, L. Xie, and H. Li, "RAWSim-O: a simulation framework for robotic mobile fulfillment systems," *Logistics Research*, vol. 11, no. 8, pp. 1–11, 2018.

[23] L. Xie, N. Thieme, R. Krenzler, and H. Y. Li, *Efficient Order Picking Methods in Robotic Mobile Fulfillment Systems*, 2019, https://arxiv.org/abs/1902.03092.

[24] F. Stulp and O. Sigaud, "Path integral policy improvement with covariance matrix adaptation," in *29th International Conference on Machine Learning*, Edinburgh, Scotland, 2012.

[25] T. Geijtenbeek, M. Van De Panne, and A. F. Van Der Stappen, "Flexible muscle-based locomotion for bipedal creatures," *ACM Transactions on Graphics*, vol. 32, no. 6, pp. 1–11, 2013.

[26] P. MacAlpine and P. Stone, "Overlapping layered learning," *Artificial Intelligence*, vol. 254, pp. 21–43, 2018.

[27] H. R. Zhou, W. S. Tang, and H. L. Wang, "Optimization of multiple traveling salesman problem based on differential evolution algorithm," *Systems Engineering Theory & Practice*, vol. 30, no. 8, pp. 1471–1476, 2010.

*Research Article*

# Sensor Fault Diagnosis of Locomotive Electro-Pneumatic Brake Using an Adaptive Unscented Kalman Filter

**Dianzhu Gao,**[1,2] **Jun Peng,**[2] **Yunyou Lu,**[1] **Rui Zhang,**[1] **Yingze Yang** [iD],[2] **and Zhiwu Huang**[1]

[1]*School of Automation, Central South University, Changsha 410075, China*
[2]*School of Computer Science and Engineering, Central South University, Changsha 410075, China*

Correspondence should be addressed to Yingze Yang; yangyingze@csu.edu.cn

Normal operation of the pressure sensor is important for the safe operation of the locomotive electro-pneumatic brake system. Sensor fault diagnosis technology facilitates detection of sensor health. However, the strong nonlinearity and variable process noise of the brake system make the sensor fault diagnosis become challenging. In this paper, an adaptive unscented Kalman filter- (UKF-) based fault diagnosis strategy is proposed, aimed at detecting bias faults and drift faults of the equalizing reservoir pressure sensor in the brake system. Firstly, an adaptive UKF based on the Sage-Husa method is applied to accurately estimate the pressure transients in the equalizing reservoir of the brake system. Then, the residual is generated between the estimated pressure by the UKF and the measured pressure by the sensor. Afterwards, the Sequential Probability Ratio Test is used to evaluate the residual so that the incipient and gradual sensor faults can be diagnosed. An experimental prototype platform for diagnosis of the equalizing reservoir pressure control system is constructed to validate the proposed method.

## 1. Introduction

The electro-pneumatic brake system has shown the extensive applications in passenger trains, metros, and heavy haul trains because of its fast response time and high reliability [1]. Locomotive electro-pneumatic brake is a crucial component which has an important function for the operational safety of the train. Faults in braking systems can lead to a reduction in locomotive braking performance and even induce safety accidents. Therefore, early detection and isolation of faults in the braking system are necessary [2]. Pressure sensors are vital components in the brake system because their reliability and measuring accuracy are crucial to achieving the accurate pressure control and approving braking performance.

The fault diagnosis of the equalizing reservoir pressure sensor is a challenging task. The brake system is composed of the electric, pneumatic, and mechanical subsystem, showing a sophisticated nonlinearity [1]. The energy transmitting medium of the braking force is compressed air, and the compressibility of air makes the system highly nonlinear [3], which makes it difficult to build a precise mathematical model of the brake system. Furthermore, the process noise and measurement noise in the braking process, which are caused by the harsh and noisy working environment, make the fault diagnosis of the brake become more challenging.

Recent years, many studies have developed sensor fault diagnosis methods [4–8]. There are three main categories of sensor fault diagnosis methods: the redundancy method and the knowledge-based method and the model-based approaches. The redundancy method is implemented by the comparison of measurements among several sensors, which has been used in wireless sensor networks [9] or the aerospace system [10], such as satellite attitude control systems [11]. The minimum degree of sensor redundancy necessary to pinpoint the distinction between sensor faults and system faults in the monitoring process is determined in [12]. However, the redundancy methods require additional

hardware sensors, showing less cost-effectiveness, which are not appropriate for the locomotive electro-pneumatic brake system [13].

The development of computer technology has provided a new method for fault diagnosis technology. Knowledge-based method uses an expert system to locate and diagnose sensor faults and does not require a quantitative mathematical model. A fuzzy expert system is established to locate sensor faults [14], and the residual generation and residual evaluation are analyzed in [15], showing its instantaneous handling capability for the fault. The problem of sensor fault recognition is considered pattern recognition in [2]. The sample data is acquired and trained to obtain a classifier, and then, the data is matched according to the classification rules. However, it is general that knowledge-based methods require a large enough amount of data, which means that many types and numbers of sensors need to be added in the brake system.

For model-based methods, the Kalman filter and its enhanced varieties are widely utilized [5, 7, 16] because of their robustness to process and measurement noise and their efficient real-time performance [17]. However, the Kalman filter is not available for the brake system because of its intrinsic nonlinear properties. Therefore, an unscented Kalman filter (UKF) is proposed to address the nonlinear problem. The UKF, which applies the unscented transform to calculate the mean and variance of measurement and process noise, has higher accuracy than the extended Kalman filter [18, 19]. However, the process noise and covariance matrices of measurement for UKF are generally assumed to be stable. And it is difficult to determine the covariance matrices in practical applications. The fault diagnosis method will suffer from performance degradation if the model uncertainty is not well defined by the process noise covariance [20]. To overcome the difficulty, the adaptivity of UKF should be improved. That is, the covariance matrices of measurement and process noise should be adaptively adjusted [8, 21, 22].

This paper proposes an adaptive UKF-based scheme to detect bias faults and drift faults of the equalizing reservoir pressure sensor. For the locomotive electro-pneumatic brake system, different from existing UKF-based fault diagnosis methods, the proposed scheme can detect incipient and gradual sensor faults. The scheme introduced the Sage-Husa mechanism to accurately estimate the pressure transients in the equalizing reservoir by filtering out the measurement noise and the changing process noise of the brake system. Further, the Sequential Probability Ratio Test is utilized to evaluate the residual, the difference between the estimated pressure, and the online sensor measurement. By combining the Sage-Husa mechanism and Sequential Probability Ratio Test, the proposed scheme can detect the incipient and gradual sensor faults of the locomotive electro-pneumatic brake system. The main contributions in this paper include the following:

(i) The mechanism of the electro-pneumatic brake system is analysed adequately, and the accurate analytical pressure model is established

(ii) The adaptive UKF is applied to estimate the system output pressure, improving the robustness of the fault diagnosis approach under the uncertainty and noise

(iii) The Sequential Probability Ratio Test is introduced to evaluate the residual to minimize the occurrence of misinformation or false detection in fault diagnosis

The rest of this paper is organized as follows. Section 2 gives a description of the brake system and builds the mathematical model. Section 3 introduces the theory of adaptive UKF and presents the fault diagnosis scheme of the pressure sensor bias and drift faults. Section 4 shows the experimental results and analysis. Finally, the conclusion is drawn in Section 5.

## 2. System Model and Problem Formulation

*2.1. Principle of the Electro-Pneumatic Brake System.* The electro-pneumatic brake system (see Figure 1) consists of the mechanical, pneumatic, and electric subsystem. The mechanical subsystem is the foundation brake rigging which mainly consists of the brake pads, drum, and shoes. The pneumatic subsystem consists of many components, including the main reservoir, the brake pipe and chamber, an equalizing reservoir, a relay valve, and a compressor. The electrical subsystem mainly contains a brake control unit (BCU), pressure sensors, and solenoid valves (brake valve and release valve).

*2.2. Model of the Equalizing Reservoir Pressure Control System.* The ideal gas law equation is as follows:

$$P = \frac{nRT}{V}, \tag{1}$$

which describes the quantitative relation among pressure, air temperature, and volume of a chamber, where $n$, $R$, $P$, $T$, and $V$ represent the number of moles of the gas, the gas constant, the absolute pressure, the absolute air temperature, and the chamber volume, respectively. Assuming the volume $V$ is invariable, taking the derivative of the equation with respect to time, we can get

$$\dot{P} = \frac{RT}{V} q_m, \tag{2}$$

where $q_m$ is the mass flow in the chamber. According to Bernoulli's equation for adiabatic and isentropic airflow, $q_m$ is calculated as follows [23]:

$$q_m = f(P_u, P_d) = \begin{cases} *20c \dfrac{P_u C_1 A}{\sqrt{RT}} \sqrt{\gamma \left(\dfrac{2}{\gamma+1}\right)^{(\gamma+1)/(\gamma-1)}}, \dfrac{P_d}{P_u} \leq 0.528, \\[4mm] \dfrac{P_u C_2 A}{\sqrt{RT}} \sqrt{\dfrac{2\gamma}{\gamma-1}} \sqrt{\left(\dfrac{P_d}{P_u}\right)^{2/\gamma} - \left(\dfrac{P_d}{P_u}\right)^{(\gamma+1)/\gamma}}, 0.528 < \dfrac{P_d}{P_u} \leq 1, \end{cases}$$

$$\tag{3}$$

Figure 1: The schematic of the locomotive electro-pneumatic brake system.

where $R$ and $T$ have the same meanings as those in the ideal gas law equation. $P_d$ and $P_u$ are the downstream pressure and upstream pressure, respectively. $C_1$ and $C_2$ are the flow rate coefficients. $\gamma$ represents the adiabatic exponent of air, and $A$ represents the orifice passage area. By combining (1) and (2), the equalizing reservoir pressure transients of the brake system in different operating modes can be formulated as (4), (5), and (6).

The equalizing reservoir pressure dynamics in the release process is formulated as

$$
\dot{P} = \begin{cases} \dfrac{P_s C_1 A_1 \sqrt{TR}}{V} \sqrt{\gamma \left(\dfrac{2}{\gamma+1}\right)^{(\gamma+1)/(\gamma-1)}}, \dfrac{P}{P_s} \leq 0.528, \\[4mm] \dfrac{P_s C_2 A_1 \sqrt{TR}}{V} \sqrt{\dfrac{2\gamma}{\gamma-1}} \sqrt{\left(\dfrac{P}{P_s}\right)^{2/\gamma} - \left(\dfrac{P}{P_s}\right)^{(\gamma+1)/\gamma}}, 0.528 < \dfrac{P}{P_s} \leq 1, \end{cases} \tag{4}
$$

where $P$ is the equalizing reservoir pressure, $P_s$ is the main reservoir pressure, $A_1$ represents the orifice passage area of the release valve, and $V$ is the equalizing reservoir volume.

The pressure in the equalizing reservoir remains steady in the hold mode. Then, the pressure dynamics can be described as

$$
\dot{P} = 0. \tag{5}
$$

The equalizing reservoir pressure dynamics in the braking process is formulated as

$$
\dot{P} = \begin{cases} -\dfrac{P C_3 A_2 \sqrt{TR}}{V} \sqrt{\gamma \left(\dfrac{2}{\gamma+1}\right)^{(\gamma+1)/(\gamma-1)}}, \dfrac{P_o}{P} \leq 0.528, \\[4mm] -\dfrac{P C_4 A_2 \sqrt{TR}}{V} \sqrt{\dfrac{2\gamma}{\gamma-1}} \sqrt{\left(\dfrac{P_o}{P}\right)^{2/\gamma} - \left(\dfrac{P_o}{P}\right)^{(\gamma+1)/\gamma}}, 0.528 < \dfrac{P_o}{P} \leq 1, \end{cases} \tag{6}
$$

where $P_o$ is the atmosphere pressure, $A_2$ represents the orifice passage area of the brake valve, and $C_3$ and $C_4$ are the flow rate coefficients.

From (4) and (6), we can know that the brake system is strongly nonlinear; thus, we choose the UKF as the pressure estimator. The mathematical model should be as accurate as possible in order to realize an effective fault diagnosis scheme. Therefore, the parameters in the models should be obtained accurately. The equalizing reservoir volume, the orifice passage areas of the release valve, and the brake valve can be measured directly. However, the flow rate coefficients in (4) and (6) need to be identified. In this paper, the flow rate coefficients $C_1 \sim C_4$ are identified by the least square method [24]. The validity of the model is tested by experiments, which is described in Section 4.

## 3. The Proposed Sensor Fault Diagnosis Method

The theory of adaptive UKF and the proposed sensor fault diagnosis method are introduced in this section. Firstly, the principle of adaptive UKF is developed, and then, the algorithm is applied to fault diagnosis of the equalizing reservoir pressure sensor.

*3.1. The Theory of Adaptive UKF.* Based on the theory of traditional UKF, the prior statistics of the process noise is used to compensate for the changing model uncertainty [25], which is adaptively corrected by the Sage-Husa noise estimator. The UKF is used for the discrete system generally. The general form of a discrete nonlinear system is defined by

$$
\begin{cases} x_{k+1} = f(x_k, y_k) + q_k = F_k x_k + B_k u_k + q_k, \\ y_k = h(x_k) + r_k, \end{cases} \tag{7}
$$

where $u_k$ and $x_k$ are the input vector and $n$-dimensional state vector, respectively. $y_k$ is the $m$-dimensional observation vector. $q_k$ and $r_k$ represent the process noise and measurement noise, respectively, which are the Gaussian white noise with zero mean.

Normally, the statistics of the measurement noise and process noise is unvarying in the UKF. However, the process noise of the brake system is varying and difficult to be determined. The measurement noise lies on the accuracy of the pressure sensor and is relatively constant. The statistics of process noise is described by the covariance matrices $Q$. Similarly, the measurement noise $R$ can be calculated from historical measurements. Then, the Sage-Husa method is applied to tune the covariance matrices $Q$ adaptively. The Sage-Husa suboptimal noise estimator is depicted as follows [26]:

$$\begin{cases} d_k = \dfrac{(1-b)}{\left(1-b^k\right)}, \\ v_k = y_k - h\left(\bar{x}_{k|k-1}\right), \\ Q_k = (1-d_{k-1})Q_{k-1} + d_{k-1}\left[K_v v_k^T K_k^T + P_k - \displaystyle\sum_{i=0}^{2n} \omega_i^c \left(\chi_{i,k|k-1} - \bar{x}_{k|k-1}\right)\left(\rho_{i,k|k-1} - \bar{y}_{k|k-1}\right)^T\right], \end{cases} \tag{8}$$

where $b \in (0.95, 0.99)$ is the forgetting factor and $K_k$ is the Kalman gain.

The more the process and measurement noise change, the higher the value of $b$. The Sage-Husa noise estimator cannot normally work when the prior statistical characteristic of noises is unknown; otherwise, the filter will diverge [27]. The statistics of the measurement noises in the brake system can be obtained according to historical measurements.

*3.2. Residual Evaluation through the Sequential Probability Ratio Test.* In order to minimize the occurrence of misinformation or false detection in fault diagnosis, [28] proposed an improved Sequential Probability Ratio Test (SPRT) method. In this method, statistical hypothesis tests are used, where $H_0$ and $H_1$ are supposed to be the nonfault hypothesis and faulty hypothesis, respectively [29]. The residuals under fault-free condition conform to a normal random variable (variance value $\sigma$ and mean value $\mu_0$), while the residuals in faulty condition have the same variance value $\sigma$, whose mean value is $\mu_1$. The log-likelihood ratio is calculated as follows:

$$L(k) = \ln \frac{p(r_i \mid H_0)}{(r_i \mid H_1)} = \frac{k(\bar{r}_k - \mu_0)^2}{2\sigma^2}, \quad \bar{r}_k = \frac{1}{k}\sum_{i=1}^{k} r_i. \tag{9}$$

The fault detection is then converted into detecting the changes of the residual mean. When there is no fault, $L(k)$ is near to zero. When there is a fault, $\bar{r}_k$ would be away from $\mu_0$ and $L(k)$ would be away from zero.

A fault is detected when $L(k) \geq T(H_1)$, where $T(H_1)$ is the threshold, $P_M$ is the missing report rate, and $P_F$ is the false alarm rate. In this paper, we set $P_M = 0.01$, $P_F = 0.01$, and $T(H_1) = 4.595$. When $k$ is too large, $L(k)$ will exceed $T(H_1)$ even though a small deviation between $\bar{r}_k$ and $\mu_0$ exists. In order to solve the issue, we set an upper bound, 2000 on $k$, and the upper bound is calibrated through experiments.

*3.3. The Proposed Sensor Fault Diagnosis Method.* The schematic of the proposed sensor fault diagnosis approach is described in Figure 2. It is assumed that the process and measurement noise of the equalizing reservoir pressure system is the Gaussian white noise. Thus, an adaptive UKF can be employed to estimate output pressure according to the inputs and outputs of equalizing reservoir pressure system. The inputs of the equalizing reservoir pressure system are generated by the brake control unit. The outputs of the equalizing reservoir pressure system are measured by a pressure sensor. Then, the residual is generated by subtracting the pressure sensor measurement from the adaptive UKF pressure estimation. Afterwards, the residual is passed through the Sequential Probability Ratio Test to increase the sensing sensitivity. By comparing the preset threshold, the fault detection result can be obtained.

The adaptive UKF is applied as the state estimator in this paper. Based on the conventional UKF algorithm, the measurement noise and process noise covariance is adaptively tuned according to (7) by using the Sage-Husa method. The fault diagnosis scheme is based on the equalizing reservoir pressure system models (3), (4), and (5). Since the system is nonlinear, we choose the Runge-Kutta methods to discretize the system models. Because the UKF algorithm can achieve third-order accuracy of the covariance and posterior mean [18], we use the second-order Runge-Kutta method [30], whose local truncation error is $O(h^3)$ and $h$ is the step size. The process of the proposed sensor fault diagnosis scheme is depicted in Algorithm 1, where the equalizing reservoir pressure, $P$, is chosen as the state variable, and the system state equations (3), (4), and (5) are simplified as $\dot{x} = f(x)$. The observation equation is $y = x$, where $y$ is the pressure sensor measurement. $T = 0.02$ (second) is the step size.

## 4. Simulation Results and Discussions

We construct an experimental platform for the equalizing reservoir pressure system and verify the effectiveness and feasibility of the proposed sensor fault diagnosis strategy. The experimental platform (see Figure 3) is part of the real locomotive electro-pneumatic brake.

The detailed parameters of the mathematical models are shown in Table 1, and the flow rate coefficients of the system are identified by the least square method. Firstly, the validity of the mathematical model of the equalizing reservoir pressure system is verified by experiments. Then, bias faults and drift faults are injected into the equalizing reservoir pressure sensor, and fault diagnosis performance of the proposed method is evaluated. Finally, fault diagnosis performance of the proposed method is compared with that of the Luenberger observer.

*4.1. Model Verification.* Figure 4 depicts the equalizing reservoir pressure transients, where the red line represents the equalizing reservoir pressure measured by a normal sensor and the blue line represents the pressure calculated from the mathematical model. The relative error between sensor measurement and model output is described by the green

Figure 2: Schematic of the proposed sensor fault diagnosis method.

1: *Determination of the system state equation and observation equation*: the system state equation $\dot{x} = f(x)$ and the observation equation $y = x$;

2: *Discretization of the system equation and observation equation*: discretized system equation $x_{k+1} = x_k + (T/2) \times (k_1 + k_2), k_1 = f(x_k)$ and $k_2 = f(x_k + T \times k_1), T = 0.02$; discretized observation equation $y_k = x_k$;

3: *Initialization*: for $k = 0$, set: $\bar{x}_0 = E[x_0]$, initial estimation error covariance $P_0 = E[(x_0 - \bar{x}_0)(x_0 - \bar{x}_0)^T]$, initial process noise covariance $Q_0 = 0.1$, measurement noise covariance $R = 0.03$;

4: *Time update and measurement update for $k = 1, 2, \cdots$*;

5: residual generation by making a difference between the pressure sensor measurement and the UKF estimation $r_k = \bar{x}_k - y_k$;

6: residual evaluation through Sequential Probability Ratio Test;

7: **goto***Time update and measurement update*

Algorithm 1: The procedure of sensor fault diagnosis.



Figure 3: The experimental platform for the equalizing reservoir pressure control system, which is part of the real locomotive electro-pneumatic brake, including the following: (a) upper computer, (b) brake control unit, (c) value and installation gas circuit plate, (d) pressure sensor, (e) equalizing cylinder, and (f) cylinder to simulate a train pipe.

Table 1: Parameters of the system model.

| Parameter | Value | Parameter | Value |
| --- | --- | --- | --- |
| $P_s$ | 650 kPa | $P_o$ | 101.33 kPa |
| $A_1$ | 4 mm$^2$ | $A_2$ | 3 mm$^2$ |
| $T$ | 293 K | $R$ | 287 J/(kg·K) |
| $\gamma$ | 1.403 | $V$ | 1.5 L |
| $C_1$ | 0.4593 | $C_2$ | 0.4362 |
| $C_3$ | 0.2505 | $C_4$ | 0.1905 |

line. The equalizing reservoir pressure transients in the braking process are illustrated (see Figure 5), and the relative error between sensor measurement and model output is plotted by the green line. It can be found that the accuracy of the system model is high adequately.

*4.2. Bias Fault Detection.* The residuals resulted from sensor bias faults in the release process are shown (see Figure 6), and the bias faults are injected to the sensor at the third second after the release operation. In this figure, the blue line describes the residuals resulting from the bias fault whose magnitude is 1 kPa, and the residuals resulting from the bias fault of 2 kPa magnitude are depicted by the red line. We see that the amplitude of the residuals changes after the bias fault occurs (see Figure 6). The larger the fault magnitude, the larger the residual magnitude. The fault detection result of sensor bias faults with different magnitudes is shown in the release process (see Figure 7). In this figure, we can see that the bias fault whose magnitude is 2 kPa is detected, while the bias fault of 1 kPa magnitude has not been detected. This is because the model is not accurate enough. To improve the sensitivity of fault detection, the model needs to be sufficiently accurate.

The residuals of sensor bias faults in the braking process are depicted (see Figure 8), and sensor bias faults with different magnitudes occur at the eighth second after the braking operation. The blue line represents the residuals of the bias fault whose magnitude is 1 kPa, and the bias fault of 2 kPa magnitude is represented by the red line. We can know that

FIGURE 4: Pressure transients measured by a sensor and calculated by a mathematical model in the release process.



FIGURE 5: Pressure transients measured by a sensor and calculated by a mathematical model in the braking process.



FIGURE 6: Residuals resulted from sensor bias faults with different magnitudes in the release process.

the residual magnitudes change after the bias fault occurs (see Figure 8). The fault detection result of sensor bias faults with different magnitudes is shown in the braking process

(see Figure 9). From the figure, we can see that the bias fault of 2 kPa magnitude has been detected, while the bias fault whose magnitude is 1 kPa has not been detected.

*4.3. Drift Fault Detection.* In this part, the drift fault detection is implemented in the braking and release processes, which is simulated by injecting a varying error to the measurement process, and the error magnitude increases each sampling period.

The residuals resulting from a sensor drift fault in the release process are described (see Figure 10), and the sensor drift fault occurs at the third second after the release operation. The measurement error increases artificially by 0.02 kPa each sampling period to simulate the drift fault. It can be seen that the residual changes slightly after the drift fault occurs (see Figure 10). The sensor drift fault detection results are shown in the release process (see Figure 11), where the log-likelihood ratio, $L(k)$, changes after the fault occurrence. The drift fault is detected about 2.5 seconds after its occurrence. Then, we can conclude that the Sequential Probability Ratio Test method has excellent performance in detecting the gradual fault. The drift fault detection result

FIGURE 7: Fault detection result of sensor bias faults with different magnitudes in the release process.



FIGURE 8: The influence of different magnitudes on residuals resulted from sensor bias faults.



FIGURE 9: The influence of different magnitudes on sensor bias fault detection result.



FIGURE 10: Residuals resulted from sensor drift fault in the release process.



FIGURE 11: Drift fault detection results in the release process.



FIGURE 12: The comparison of residuals generated by the proposed methods and Luenberger observer.

of the equalizing reservoir pressure sensor in the braking process is similar to that in the release process.

*4.4. Performance Comparison of Sensor Fault Diagnosis Methods.* In order to compare the performance of different fault diagnosis methods, the residual of the proposed fault diagnosis method is compared with that of the Luenberger observer during the brake release procedure. Residual comparison of different methods is carried out when the

locomotive electro-pneumatic brake system is operated under normal conditions. In the experiment, the residuals are generated by subtracting the sensor measurement from the equalization cylinder pressure estimated by the observer and the proposed method based on the adaptive unscented Kalman filter, respectively.

In Figure 12, the blue line describes the residuals of the fault detection method based on the adaptive untraceless Kalman filter. The black line represents the residuals of the Luenberger observer. And the red line represents the fault detection threshold. From Figure 12, it can be seen that there are many false error detections when using the Luenberger observer for fault detection, because the error of the mechanism model is too large. If the threshold is increased to reduce the false positive rate of the Luenberger observer method, the sensitivity of the Luenberger observer method will be reduced. On the contrary, the residual error of the proposed method is much smaller than that of the Luenberger observer method and fluctuates little. By comparing the residuals generated by the two methods and analysing the fault detection results, the conclusion can be drawn that the proposed fault detection method has better accuracy and sensitivity than the Luenberger observer method. This is because the adaptive unscented Kalman filter can filter out the changing process noise and measurement noise and accurately estimate the pressure of the equalizing air cylinder.

## 5. Conclusions

This paper proposes an efficient and novel model-based sensor fault diagnosis algorithm based on UKF for the locomotive electro-pneumatic brake system. For this purpose, the accurate pressure mathematical model is first built. Then, an adaptive UKF is applied to estimate the pressure transients of the equalizing reservoir to improve the algorithm's robustness. The residuals are calculated, and the residual evaluation is implemented by an improved Sequential Probability Ratio Test method. The proposed algorithm can efficiently detect drift faults and bias faults of the equalizing reservoir pressure sensor. Experiments validate the feasibility and effectiveness. The future work that needs to be investigated is to improve the fault detection sensitivity for minor and gradual fault.

### Data Availability

The data used to support the findings of the manuscript are available within the article.

### Conflicts of Interest

The authors declare that they have no conflicts of interest.

### Acknowledgments

## References

[1] C. Siva Chaitanya, S. C. Subramanian, P. Karthikeyan, and N. Jagga Raju, "Modelling an electropneumatic brake system for commercial vehicles," *IET Electrical Systems in Transportation*, vol. 1, no. 1, pp. 41–48, 2011.

[2] G. Zhiwei, C. Cecati, and S. X. Ding, "A survey of fault diagnosis and fault-tolerant techniques part i: fault diagnosis with model-based and signal-based approaches," *IEEE Transactions on Industrial Electronics*, vol. 62, no. 6, pp. 3757–3767, 2015.

[3] A. Ilchmann, O. Sawodny, and S. Trenn, "Pneumatic cylinders: modelling and feedback force-control," *International Journal of Control*, vol. 79, no. 6, pp. 650–661, 2006.

[4] L. Wu and D. Ho, "Fuzzy filter design for Itô stochastic systems with application to sensor fault detection," *IEEE Transactions on Fuzzy Systems*, vol. 17, no. 1, pp. 233–242, 2009.

[5] G. H. B. Foo, X. Zhang, and D. M. Vilathgamuwa, "A sensor fault detection and isolation method in interior permanent-magnet synchronous motor drives based on an extended Kalman filter," *IEEE Transactions on Industrial Electronics*, vol. 60, no. 8, pp. 3485–3495, 2013.

[6] B. Pourbabaee, N. Meskin, and K. Khorasani, "Sensor fault detection, isolation, and identification using multiple-model-based hybrid Kalman filter for gas turbine engines," *IEEE Transactions on Control Systems Technology*, vol. 24, no. 4, pp. 1184–1200, 2015.

[7] Z. Liu and H. He, "Model-based sensor fault diagnosis of a lithium-ion battery in electric vehicles," *Energies*, vol. 8, no. 7, pp. 6509–6527, 2015.

[8] P. Lu, L. Van Eykeren, E. J. Van Kampen, Q. P. Chu, and B. Yu, "Adaptive hybrid unscented Kalman filter for aircraft sensor fault detection, isolation and reconstruction," in *AIAA guidance, navigation, and control conference*, p. 1145, National Harbor, Maryland, 2014.

[9] C. Huang, G. Huang, W. Liu, R. Wang, and M. Xie, "A parallel joint optimized relay selection protocol for wake-up radio enabled WSNs," *Physical Communication*, vol. 47, no. 3, article en8076509, p. 10320, 2021.

[10] R. J. Patton, "Fault detection and diagnosis in aerospace systems using analytical redundancy," in *IEE Colloquium on Condition Monitoring and Fault Tolerance*, p. 1, London, UK, 1990.

[11] R. Wang, Y. Cheng, and M. Xu, "Analytical redundancy based fault diagnosis scheme for satellite attitude control systems," *Journal of the Franklin Institute*, vol. 352, no. 5, pp. 1906–1931, 2015.

[12] M. Taiebat and F. Sassani, "Distinguishing sensor faults from system faults by utilizing minimum sensor redundancy," *Transactions of the Canadian Society for Mechanical Engineering*, vol. 41, no. 3, pp. 469–487, 2017.

[13] D. Li, Y. Wang, J. Wang, C. Wang, and Y. Duan, "Recent advances in sensor fault diagnosis: a review," *Sensors and Actuators A: Physical*, vol. 309, article 111990, 2020.

[14] Y. L. Ou, "Fault diagnosis with fuzzy expert system," *Applied Mechanics and Materials*, vol. 48, 2011.

[15] M. Geetha and J. Jerome, "Fuzzy expert system based sensor and actuator fault diagnosis for continuous stirred tank

reactor," in *2013 International Conference on Fuzzy Theory and Its Applications (iFUZZY)*, pp. 251–257, Taipei, Taiwan, 2013.

[16] K. Xiong, C. Chan, and H. Zhang, "Detection of satellite attitude sensor faults using the UKF," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 43, no. 2, pp. 480–491, 2007.

[17] F. Auger, M. Hilairet, J. M. Guerrero, E. Monmasson, T. Orlowska-Kowalska, and S. Katsura, "Industrial applications of the Kalman filter: a review," *IEEE Transactions on Industrial Electronics*, vol. 60, no. 12, pp. 5458–5471, 2013.

[18] E. A. Wan and R. Van Der Merwe, "The unscented Kalman filter for nonlinear estimation," in *Proceedings of the IEEE 2000 Adaptive Systems for Signal Processing, Communications, and Control Symposium (Cat. No.00EX373)*, pp. 153–158, Lake Louise, AB, Canada, 2000.

[19] R. Van Der Merwe and E. A. Wan, "The square-root unscented Kalman filter for state and parameter-estimation," in *2001 IEEE international conference on acoustics, speech, and signal processing*, pp. 3461–3464, Salt Lake City, UT, USA, 2001.

[20] L. Zhentong and H. Hongwen, "Sensor fault detection and isolation for a lithium-ion battery pack in electric vehicles using adaptive extended Kalman filter," *Applied Energy*, vol. 185, no. 2, pp. 2033–2044, 2017.

[21] A. Mirzaee and K. Salahshoor, "Fault diagnosis and accommodation of nonlinear systems based on multiple-model adaptive unscented Kalman filter and switched MPC and h-infinity loop-shaping controller," *Journal of Process Control*, vol. 22, no. 3, pp. 626–634, 2012.

[22] F. Sun, X. Hu, Y. Zou, and S. Li, "Adaptive unscented Kalman filtering for state of charge estimation of a lithium-ion battery for electric vehicles," *Energy*, vol. 36, no. 5, pp. 3531–3540, 2011.

[23] T. Nguyen, J. Leavitt, F. Jabbari, and J. E. Bobrow, "Accurate sliding-mode control of pneumatic systems using low-cost solenoid valves," *IEEE/ASME Transactions on Mechatronics*, vol. 12, no. 2, pp. 216–219, 2007.

[24] F. Alonge, F. D'Ippolito, and F. M. Raimondi, "Least squares and genetic algorithms for parameter identification of induction motors," *Control Engineering Practice*, vol. 9, no. 6, pp. 647–657, 2001.

[25] K. Myers and B. Tapley, "Adaptive sequential estimation with unknown noise statistics," *IEEE Transactions on Automatic Control*, vol. 21, no. 4, pp. 520–523, 1976.

[26] S. Yong and C. Han, "Adaptive UKF method with applications to target tracking," *Acta Automatica Sinica*, vol. 37, no. 6, pp. 755–759, 2011.

[27] Y. Yang and W. Gao, "An optimal adaptive Kalman filter," *Journal of Geodesy*, vol. 80, no. 4, pp. 177–183, 2006.

[28] L. Kunpeng and Z. Qinghua, "An improved sequential probability ratio test method for residual test," *Electronics Optics & Control*, vol. 16, no. 8, pp. 36–39, 2009.

[29] A. Wald and J. Wolfowitz, "Optimum character of the sequential probability ratio test," *The Annals of Mathematical Statistics*, vol. 19, no. 3, pp. 326–339, 1948.

[30] M. H. Carpenter, D. Gottlieb, S. Abarbanel, and W. S. Don, "The theoretical accuracy of Runge–Kutta time discretizations for the initial boundary value problem: a study of the boundary error," *SIAM Journal on Scientific Computing*, vol. 16, no. 6, pp. 1241–1252, 1995.

*Research Article*

# In-Body Electromagnetic Sensor Combined with AI-Enhanced Electrocardiography for Pacemaker Working Status Telemonitoring

**Wu Lu,[1] Ranran Ding [iD],[1] Bingjie Wu,[1] Wenbin Zhao,[1] Dong Huang [iD],[2] and Xue Zhang[2]**

[1]*The School of Electrical Engineering, Shanghai University of Electric Power, Shanghai 200090, China*
[2]*The Division of Cardiology, The Sixth People's Hospital Affiliated to Shanghai, Shanghai 200233, China*

Correspondence should be addressed to Dong Huang; huangdong1004@126.com

This paper describes the design and implementation of an in-body electromagnetic sensor for patients with implanted pacemakers. The sensor can either be mounted on myocardial tissue and monitor the electrocardiography (ECG) with contact electrodes or implanted under the skin and monitor the ECG with coaxial leads. A 16-bit high-resolution analog front-end (AFE) and an energy-efficient 32-bit CPU are used for instantaneous ECG recording. Wireless data transmission between the sensor and clinician's computer is achieved by an embedded low-power Bluetooth transmitter. In order to automatically recognize the working status of the pacemaker and alarm the episodes of arrhythmias caused by pacemaker malfunctions, pacing mode classification and fault diagnosis on the recorded ECG were achieved based on an AI algorithm, i.e., a resource allocation network (RAN). A prototype of the sensor was implemented on a human torso, and the *in vitro* test results prove that the sensor can work properly for the 1-4-meter transmission range.

## 1. Introduction

According to the WHO's 2019 Global Health Estimates [1], cardiovascular diseases (CVD) have become one of the main sources of human death in the last 20 years, accounting for 16 percent of total death cases. Accurate and advanced screening of heart failure signs is an effective method in reducing the mortality of patients caused by heart failures. Among all the cardiac monitoring technologies, electrocardiography (ECG) signals are most commonly used to assess the state of the heart and indicate irregular heartbeats, due to its high resolution and strong anti-interference abilities.

For continuous heart monitoring, pacing of the heart is generally monitored by a mobile ECG monitoring device, i.e., Holter [2, 3]. However, there are some disadvantages of the Holter. Firstly, it is bulky and inconvenient to carry, which leads to the need for a specific environment for operation. Secondly, patients with CVD need to paste the electrodes on their skin which may lead to allergy. Thirdly, the Holter is not allowed to operate at frequencies above 1 kHz [4], which

makes it difficult to detect abnormal cardiac events in extreme conditions. For example, when a patient with implanted pacemaker is exposed to a transient electromagnetic field, the electromagnetic interference (EMI) with frequency of ~kHz to ~MHz could be created in the pacing loop formed by the leads and the pulse generator. These EMI signals are difficult to detect by a regular Holter, but they need to be properly monitored since these signals could be misunderstood as the normal pacing pulses by the pacemaker and potentially cause pacemaker malfunctions. A sensor which can monitor the normal ECG and the high frequency EMI signals simultaneously is necessary for heart failure protection.

At present, various wearable and implantable sensors have been developed for monitoring heart activity based on surface ECG, but most of them can only detect ECG without classifying the pacing mode [3, 5, 6]. As a result, abnormal ECG signals, such as arrhythmias, can only be recognized by manual visual examination of the ECG by physicians. However, some abnormal signals lack specificity, and the differences between them and normal signals are inconspicuous. Misunderstandings and

FIGURE 1: System overview of the implantable electromagnetic sensing system.

omissions of important information in ECG diagnosis are inevitable. Therefore, an intelligent diagnostic method is necessary to improve the accuracy of ECG diagnosis. Recent research has shown that deep machine learning can establish the mapping relationship of nonlinear functions in ECG and fully explore the information that is difficult to identify manually [6, 7]. For this reason, the introduction of machine learning to assist the identification and diagnosis of ECG signals detected by an implantable sensor can greatly improve the efficiency of diagnosis, reduce the rate of misdiagnosis, and save medical costs [6–8].

In this paper, the design and evaluation of an implantable electromagnetic sensor are introduced, which can monitor and classify ECG with ultralow power consumption, high signal resolution, and automatic pacing mode recognition. This paper is organized as follows: firstly, we describe the architecture and key parameters of the joint ECG telemonitoring system, followed by the embedded forms in various clinical conditions and main circuit blocks of the sensor. Next, a deep learning network based on the RAN model is introduced for automatic fault diagnosis of the measured ECG signal. Eventually, the feasibility of the in-body sensor is tested by *in vitro* tests, where it is indicated that the sensor can work properly for the 1-4-meter transmission range, and two types of the abnormal pacing mode, i.e., pulmonary hypertension (PH) and respiratory sinus arrhythmia (RSA), are successfully validated through the AI-enhanced ECG signals.

## 2. Design of the Implantable Electromagnetic Sensing System

*2.1. System Architecture.* As shown in Figure 1, the implantable electromagnetic sensing system consists of two parts,

i.e., the in-body electromagnetic sensor and the AI diagnostic system. For the in-body sensor, the heart rhythms on the surface of the myocardium are detected by the contact electrodes and amplified by a highly sensitive and low-noise operational amplifier (TLV 9152) integrated with a resistance network and high-pass filter. A time-multiplexed 16-bit SAR ADC digitizes the output signals. The output of the ADC is then sent to the central processing unit (CPU MSP432P4011RGCR), which packs the ADC output data with the proposed frame structure. The recorded data is stored in a 1 GB SPI Flash memory and transmitted wirelessly by an energy-efficient Bluetooth transmitter to the PC. A 3.7 V lithium battery is chosen as the power supply to support up to 1 MHz sampling rate for the 1-4-meter transmission range. The transmitted heart rhythms are in the form of ECG and examined by the AI diagnostic system on the PC. The AI diagnostic system is achieved by offline network training and online pattern recognition. In the offline training section, typical ECG in normal, pulmonary P wave, and arrhythmia cases monitored by Holter were packed into the ECG database. Wavelet packet decomposition (WPD) and principal component analysis (PCA) are used to extract the features of the above three types of ECG. Then, the training was carried out in the model based on the RAN algorithm. For the online model evaluation, the ECG monitored in real time by the in-body sensor are scanned online through WPD and the feature extraction is achieved online by PCA. By using the pretrained weights of the network, the pacing mode can be correctly recognized.

*2.2. Sensor Geometry.* As shown in Figure 2, the implantable electromagnetic sensor mainly consists of a titanium alloy shell and a built-in sampling chip. The titanium alloy shell can be described as a cylindrical shielded container with

(a) Sensor geometry for *in vitro* test

(b) Sensor geometry for *in vivo* test

(c) Built-in sampling chip

Figure 2: The anatomy of implantable electromagnetic sensor.

proper seal cover, which ensures good electromagnetic shielding with frequency up to GHz and is waterproof. It is noted that the contact electrodes for heart rhythm measurement are different for *in vivo* and *in vitro* conditions. For the *in vitro* test, the sensor can stay close to the torso, so that the contact electrodes are designed as a pair of probe needles. For the *in vivo* test, the sensor can only stay in the subcutaneous pocket, so that the contact electrodes are designed as coaxial lead electrodes which can pass through the veins. The built-in sampling chip consists of four parts, i.e., a signal modulation board, a master control board, a wireless communication, and a data storing board. The detailed implementations of each part are described in the following sections.

### 2.3. Circuit Implementation

*2.3.1. Signal Modulation Board.* As shown in Figure 3, the signal modulation board is completely achieved by the implementation of an analog circuit. Its main function is to filter clutter and amplify the heart rhythm signal. The analog circuit is composed of the input circuit, the secondary amplifier circuit, and the trailing circuit. The input circuit is made by two resistances with values of 800 kΩ and 200 kΩ in series, which can realize the function of reducing the input signal by 4/5 through resistor voltage division. The

secondary amplifier is a coamplifier with model TLV9152. The role of the secondary amplifier circuit is to make the signal to the setting threshold value, and ensure that the polarity of the signal can be recorded. The trailing circuit is a resistance network composed of an eight-choice analog switch with model SN74CBTLV3251, which can amplify the signal by 128 times. The trailing circuit can remove clutter which is below the setting threshold. Further, by setting the parameters of the filter, the device can measure signals with frequencies up to 1 MHz.

*2.3.2. Master Control Board.* As shown in Figure 4, the function of the master control system is to transform the amplified and noiseless electrical signals into digital signals. A time-multiplexed 16-bit SAR ADC digitizes the output signals. The output of the ADC is then sent to the central processing unit (CPU MSP432P4011RGCR), which packs the ADC output data with the proposed frame structure. Further, the reference voltage chip (REF4132) provides voltage with the amplitude of 2.5 V to ensure high accuracy and high stability of ADC. This section guarantees that the device has a sampling rate of up to 1 Mbps.

*2.3.3. Wireless Communication and Data Storing Board.* As shown in Figure 5, digital signals are stored in a 1 GB flash

Figure 3: Circuit implementation of signal modulation board.

memory (TC58CVG1S3HRAIJ), which is controlled through the serial interface (SPI) of the CPU. The device can connect to the PC to transmit data through wireless Bluetooth (CC2650). Then, the types of CVD can be recognized through the AI diagnostic system on the PC.

2.4. AI Diagnostic System. The offline network training process provides the pretrained weights of pacing modes for the online classification of the AI diagnostic system. Offline training uses the open data of ECG signals from the PhysioNet database [9]. A typical ECG signal from PhysioNet is shown in Figure 6. Three types of ECG signals are selected as the training target, i.e., normal, pulmonary hypertension (PH), and respiratory sinus arrhythmia (RSA), and the sample size for each pacing mode is defined as 500. For the precision of AI algorithm validation, ECG cases with PH and RAS are selected from people with different degrees of disease. The testing targets for online classification are the ECG signals obtained by the electromagnetic sensor. The detailed evaluations of the AI technology are described in the following sections.

2.4.1. Data Processing. ECG signals processing includes noise elimination, baseline drift, and data enhancement [10, 11]. The empirical decomposition algorithm (EDM) is used to decompose the ECG signal into 10 intrinsic mode functions (IMF) [12]. The wavelet transform algorithm is used to denoise IMF1 and IMF2 with high frequency. IMF9 and IMF10 with low frequency can eliminate baseline drift according to the median filtering algorithm. Then, the processed IMF mode and the remaining unprocessed IMF mode are reconstructed to obtain a smooth and noiseless ECG signal. Finally, data enhancement is performed to prevent the neural network from overfitting and to improve the unbalanced frequency of the ECG [13]. The ECG after data processing is shown in Figure 7.

After the ECG is enhanced by data processing, the pan-Tompkins QRS feature detection algorithm is applied to the processed ECG to locate QRS peak points. According to the position of QRS, the heartbeat signals are cut apart into singles [14]. After the position of QRS is determined, the position of the P wave which is about 200 ms away from QRS can also be located. The peak of the P wave can be obtained in the range of 150 ms-200 ms ahead of the appearance of the peak of QRS [15]. In the process of heartbeat segmentation, 99 sample points were intercepted forward and 100 sample points were intercepted backward at the QRS peak points which had been located. Each sample can contain a P wave, QRS wave, T wave, and other information of a heartbeat cycle. The positioned P wave is shown in Figure 8.

2.4.2. Feature Extraction. In this paper, wavelet packet decomposition (WPD) and principal component analysis (PCA) are used to extract ECG features. After three-layer wavelet packet decomposition of ECG, 8 wavelet packet coefficients of the subfrequency band can be obtained, which help to analyze signals with different frequency bands [16]. The time frequency of ECG cases with normal, PH, and RAS is totally different, so the energy of each ECG is different in various frequency bands. As a result, the energy of each frequency band is calculated as the feature vector to diagnose the ECG type [17]. After the decomposition of the signal wavelet packet, the energy in each subfrequency band is calculated as follows:

$$E_j^* = \sum_{k=1}^{N} \left| d_j^k \right|^2, \tag{1}$$

where $n$ is the number of wavelet packet decomposition layers, $N$ is the number of coefficients in $d_j$, and $d_j^{\ k}$ is the coefficient obtained from the decomposition of the $j$-layer wavelet packet.

The energy in each subfrequency band obtained after $n$-layer wavelet packet decomposition of the signal is calculated as

$$E^* = (E_1^*, E_2^*, \cdots, E_{2^n}^*), \tag{2}$$

FIGURE 4: Circuit implementation of master control board.



FIGURE 5: Circuit implementation of wireless communication and data storing module.



FIGURE 6: Original normal ECG from PhysioNet database.

FIGURE 7: Original normal ECG after data processing.



FIGURE 8: The positioned P wave.

where the signal is decomposed by a 3-layer wavelet packet to get 8 frequency bands, and then, the energy in each sub-frequency band is calculated and normalized to get the feature vector of each ECG.

In order to simplify the calculation process, principal component analysis is used, which can be utilized to remove the redundant feature components and extract the main part. $C_k$ is the corresponding characteristic unit vector, $k = 1, 2, \cdots, K$. $\lambda k$ is the characteristic root. The relationship between the principal component $y_k$ ($k = 1, 2, \cdots, K$), the eigenvector matrix $D$, and the contribution factors of the first $N$ principal components can be expressed as

$$y_k = c'_k D, \tag{3}$$

$$\eta_n = \frac{\sum_{k=1}^{n} \lambda_k}{\sum_{k=1}^{K} \lambda_k}. \tag{4}$$

*2.4.3. Classification Tools.* The resource allocation network (RAN) is a single hidden layer forward network, which can

create a compact network, and has the characteristics of high learning speed [18]. Because of its structure and efficient performance, the RAN network can be used to classify ECG data after feature extraction. The Gaussian activation function is used for each hidden node of the hidden layer in RAN, and the following local mapping is implemented:

$$z_j = \sum_k \left( c_{jk} - I_k \right)^2, \tag{5}$$

$$x_j = \exp \left( -\frac{z_j}{\omega_j^2} \right). \tag{6}$$

$c_{jk}$ is the data center of the hidden node of RAN. $\omega_j$ is the width of the hidden node. $x_j$ is the output of the hidden node. $z_j$ is the connection weight between the hidden node and the output node. In order to accelerate the learning

FIGURE 9: The framework of AI diagnostic system.



FIGURE 10: The profile view of *in vitro* experiment.

speed of the algorithm, the following equation is usually adopted to replace equation (6).

$$
x_j = \begin{cases} 1 - \left(\dfrac{z_j}{q\omega_j^2}\right)^2, & z_j < q\omega_j^2, \\ 0, & \text{other,} \end{cases} \tag{7}
$$

where $q$ is an empirical value of 2.67.

In the offline training process, the ECG cases with normal, PH, and RAS from the PhysioNet database were tagged as the dataset for ECG classification. Then, the ECG features were extracted by WPD and optimized by PCA. Finally, 80% of the ECG data was fed to the RAN network for training, and the remaining 20% of the ECG data was used to test the RAN classifier and establish the ECG diagnostic model.

In the online classification process, the ECG data on the myocardium is collected by the in-body electromagnetic sensor. The features are obtained and selected, and then, the preferred abnormal features are sent to the RAN classifier for online diagnosis. Finally, the ECG type corresponding to the real ECG cases is indicted. The working flowchart of the AI diagnostic system is shown in Figure 9.

## 3. *In Vitro* Experimental Verifications

To further validate our design, an *in vitro* experiment using the prototype sensor is performed, as shown in Figure 10. A human torso phantom with an EMI source is introduced. The human torso phantom is made by a pacemaker connected to a pork heart immersed in saline solution. The EMI source uses the air gap as the on/off switch to charge/discharge a series of capacitors to generate an impulse

(a) The normal ECG



(b) The ECG case with PH



(c) The ECG case with RAS

FIGURE 11: The ECG measured by implantable electromagnetic sensor.

current with duration in microseconds and amplitude of kiloampere, which can create a radiated electromagnetic field in milliteslas surrounding the torso.

During the test, the swine heart is continually excited by pacing pulses at 60 ppm emitted from the pacemaker, with the amplitude of the normal ECG of approximately 0.6-0.8 mV. The implantable electromagnetic sensor was used to monitor ECG signals on the myocardium continually.

The typical ECG signals measured by the in-body sensor are shown in Figure 11. Figure 11(a) shows the normal ECG measured by the sensor. The single ECG cycle is about 800 ms and its amplitude is 0.6 mV, which is the same as the ECG detected by Holter. Moreover, the characteristics of the P wave and the QRS wave can be fully demonstrated by the normal ECG signals measured by the sensor. Figure 11(b) shows the ECG cases with PH measured by the sensor. The key characteristic of this kind of ECG is that the P wave is a peaked wave with an extremely high amplitude, which is about 5 times of the normal P wave. This feature is identical to the clinical ECG cases with PH.

Figure 11(c) shows the ECG cases with RAS recorded by the sensor. The key characteristic of this kind of ECG is that the interval time of the single ECG signal is greater than 1000 ms, which can correspond to clinical ECG cases with RAS.

In order to fully reflect the applicability of the model, we selected three kinds of pacemakers as heartbeat pulse sources in the *in vitro* experiment. Taking 1 kA as the step length, the three pacemakers worked under 1-5 kA, and 5 experiments were conducted at each current level. The identification results are shown in Figure 12 and Table 1. It is shown that the system can automatically distinguish three types of ECG, with an overall classification precision of 83.2%.

## 4. Conclusion

This paper presents a novel sensor system for implantable, wireless communicated, and easy-to-use ECG data acquisition and pacing mode recognition. This system is based on

FIGURE 12: Performance of AI diagnostic system.

TABLE 1: Classification precision of AI system.

| Category | Number of input pictures | Number of recognized pictures | Precision | Mean average precision |
|---|---|---|---|---|
| Normal | 100 | 85 | 85.0% | |
| PH | 78 | 64 | 81.9% | 83.2% |
| RAS | 40 | 33 | 82.6% | |

a cooperative heart-computer-interface technology, i.e., an in-body electromagnetic sensor combined with AI-enhanced ECG. The in-body sensor is formed by a highly sensitive and low-noise analog signal measuring module, a time-multiplexed ADC and high-resolution CPU for data processing, and an energy-efficient Bluetooth transmitter for data transferring to PC. The highly integrated chip-on-chip packaging allows the contact electrode on the sensor in either probe or coaxial lead forms and monitoring normal ECG and EMI signals simultaneously. Further, the deep learning network based on the RAN algorithm is applied on the measured ECG signals from PC, which allows the

precise feature extraction and pattern recognition of both normal and abnormal ECG. The sensor is tested in an *in vitro* experiment, and the results indicate that the system is able to synchronously measure and diagnose ECG signals from pacemakers. Two types of abnormal ECG, i.e., PH and RAS cases, as well as the normal ECG are successfully recognized by the AI diagnosis system with overall classification precision of 83.2%. Overall, the validated and verified design of in-body sensor and AI-enhanced ECG could potentially be used as human-like interpretation of the ECG but also as a powerful tool for long-term and emergency monitoring of cardiac health and diseases.

## 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 is no conflict of interest regarding the publication of this paper.

## Acknowledgments

## References

[1] H. K. Green, O. Lysaght, D. D. Saulnier et al., "Challenges with disaster mortality data and measuring progress towards the implementation of the Sendai framework," *International Journal of Disaster Risk Science*, vol. 10, no. 4, pp. 449–461, 2019.

[2] Z. Ai, Z. Wang, and W. Cui, "Low-power wireless wearable ECG monitoring chestbelt based on ferroelectric microprocessor," *Chinese Control Conference (CCC)*, vol. 2019, pp. 6434–6439, 2019.

[3] C. Xiao, D. Cheng, and K. Wei, "An LCC-C compensated wireless charging system for implantable cardiac pacemakers: theory, experiment, and safety evaluation," *IEEE Transactions on Power Electronics*, vol. 33, no. 6, pp. 4894–4905, 2018.

[4] Z. J. Zhu, X. M. Wu, and Z. X. Fang, "The development of programming and telemetery system for cardiac pacemaker," *Progress in Biomedical Engineering*, vol. 32, 2011.

[5] R. Fensli, J. G. Dale, P. O'Reilly, J. O'Donoghue, D. Sammon, and T. Gundersen, "Towards improved healthcare performance: examining technological possibilities and patient satisfaction with wireless body area networks," *Journal of Medical Systems*, vol. 34, no. 4, 2010.

[6] J. Yoo, L. Yan, S. Lee, Y. Kim, and H.-J. Yoo, "A 5.2 mW self-configured wearable body sensor network controller and a 12 $\mu$W wirelessly powered sensor for a continuous health monitoring system," *IEEE Journal of Solid-State Circuits*, vol. 45, no. 1, pp. 178–188, 2010.

[7] J. C. Jentzer, A. H. Kashou, Z. I. Attia et al., "Left ventricular systolic dysfunction identification using artificial intelligence-augmented electrocardiogram in cardiac intensive care unit patients," *International Journal of Cardiology*, vol. 326, pp. 114–123, 2021.

[8] T. J. W. Dawes, A. de Marvao, W. Shi et al., "Machine learning of three-dimensional right ventricular motion enables outcome prediction in pulmonary hypertension: a cardiac MR imaging study," *Radiology*, vol. 283, no. 2, pp. 381–390, 2017.

[9] "Comparative review of the algorithms for removal of electrocardiographic interference from trunk electromyography," *Journal of Sensors*, vol. 20, no. 17, 2020.

[10] M. H. Islam Chowdhuryy, M. Sultana, R. Ghosh, J. U. Ahamed, and M. Mahmood, "AI assisted portable ECG for fast and patient specific diagnosis," in *2018 International Conference on Computer, Communication, Chemical, Material and Electronic Engineering (IC4ME2)*, pp. 1–4, 2018.

[11] J. Cai, W. Sun, J. Guan, and I. You, "Multi-ECGNet for ECG arrythmia multi-label classification," *IEEE Access*, vol. 8, pp. 110848–110858, 2020.

[12] Y. Wei, J. Zhou, Y. Wang et al., "A review of algorithm & hardware design for AI-based biomedical applications," *IEEE Transactions on Biomedical Circuits and Systems*, vol. 14, no. 2, pp. 145–163, 2020.

[13] A. H. Khandoker, J. Gubbi, and M. Palaniswami, "Automated scoring of obstructive sleep apnea and hypopnea events using short-term electrocardiogram recordings," *IEEE Transactions on Information Technology in Biomedicine*, vol. 13, no. 6, pp. 1057–1067, 2009.

[14] F. Liu, S. Wei, Y. Li et al., "The accuracy on the common pan-Tompkins based QRS detection methods through low-quality electrocardiogram database," *Journal of Medical Imaging and Health Informatics*, vol. 7, no. 5, pp. 1039–1043, 2017.

[15] P. Mohammad and A. Fateme, "Providing an efficient algorithm for finding R peaks in ECG signals and detecting ventricular abnormalities with morphological features," *Journal of medical signals and sensors*, vol. 6, no. 4, pp. 218–223, 2016.

[16] R. F. Atiqah, A. Saidatul, A. A. Azamimi, and N. R. Francis, "The wavelet packet decomposition features applied in EEG based authentication system," *Journal of Physics: Conference Series*, vol. 1997, no. 1, 2021.

[17] R. C. Zhou, Y. P. Huang, K. X. Huang, and J. H. Sun, "Image encryption algorithm based on mixed chaotic system and ECG signal," *Computer Measurement & Control*, vol. 28, no. 12, pp. 191–201, 2020.

[18] T. W. Dawson, M. A. Stuchly, K. Caputa, A. Sastre, R. B. Shepard, and R. Kavet, "Pacemaker interference and low-frequency electric induction in humans by external fields and electrodes," *IEEE Transactions on Biomedical Engineering*, vol. 47, no. 9, pp. 1211–1218, 2000.

*Research Article*

# Improved $Q$-Learning Method for Multirobot Formation and Path Planning with Concave Obstacles

**Zhilin Fan** [ID],[1] **Fei Liu** [ID],[1] **Xinshun Ning,**[1] **Yilin Han,**[1] **Jian Wang,**[2] **Hongyong Yang,**[1] **and Li Liu**[1]

[1]*School of Information and Electrical Engineering, Ludong University, Yantai 264000, China*
[2]*Yantai Municipal People's Procuratorate, Yantai 264000, China*

Correspondence should be addressed to Fei Liu; liufeildu@163.com

Aiming at the formation and path planning of multirobot systems in an unknown environment, a path planning method for multirobot formation based on improved $Q$-learning is proposed. Based on the leader-following approach, the leader robot uses an improved $Q$-learning algorithm to plan the path and the follower robot achieves a tracking strategy of gravitational potential field (GPF) by designing a cost function to select actions. Specifically, to improve the Q-learning, $Q$-value is initialized by environmental guidance of the target's GPF. Then, the virtual obstacle-filling avoidance strategy is presented to fill non-obstacles which is judged to tend to concave obstacles with virtual obstacles. Besides, the simulated annealing (SA) algorithm whose controlling temperature is adjusted in real time according to the learning situation of the $Q$-learning is applied to improve the action selection strategy. The experimental results show that the improved $Q$-learning algorithm reduces the convergence time by 89.9% and the number of convergence rounds by 63.4% compared with the traditional algorithm. With the help of the method, multiple robots have a clear division of labor and quickly plan a globally optimized formation path in a completely unknown environment.

## 1. Introduction

As robots become more and more widely used in various industries, a single robot cannot be competent for complex tasks. Therefore, multirobot formation [1] and path planning have become research hotspots, and they have good applications [2, 3] in collaborative search, exploration, handling, rescue, and group operations. Path planning of multirobot formation requires multiple robots to form a formation and maintain this positional relationship to move to the target. It is necessary not only to avoid obstacles safely but also to find a better path. In addition, compared to the simpler path planning in the known environment, higher requirements on the ability of multiple robots to plan paths are put in the unknown environment. There have been many implementation methods for multirobot formation, including behavior-based method [4], virtual structure method [5], and leader-following method [6]. The behavior-based method is to design

sub-behaviors in advance and choose the behavior to execute according to the changes in the situation, but the accuracy is not enough to integrate various behaviors in a complex environment. The virtual structure method regards the formation as a fixed rigid structure and cannot effectively avoid obstacles. The leader-following method with the advantage of simple and flexible structure mainly realizes collaboration by sharing information of leader. For robot's path planning, A∗ algorithm [7] and reinforcement learning (RL) algorithm [8] are commonly used in global path planning; the former can effectively solve the optimal path, but it needs to know all the environmental information in advance; the latter can learn autonomously in the environment, but it takes more time. The artificial potential field (APF) method [9] is widely used in local path planning, which can cope with the real-time changing environment but lacks the global planning ability.

For the problem of multirobot formation and path planning, Chen et al. [10] proposed a new leader-following

control framework by introducing the RH method, which enables fast convergence of a formation task for a group of mobile robots. Based on the path planning of a single robot, Sruthi et al. [11] designed a nonlinear controller for tracking to achieve the multirobot formation. The above two methods require rigorous modeling of the system and cumbersome theory, which are weak in the application. By mixing formation control with leader-following and priority methods, Sang et al. [12] used the MTAPF algorithm with an improved A∗ algorithm for path planning. Das and Jena [13] implemented collision-free path planning for multiple robots by using an improved particle swarm algorithm and evolutionary operators. Qu et al. [14] used a modified genetic algorithm to plan paths for multiple robots by adding a co-evolution mechanism. Lazarowska [15] used the discrete APF to find crash-free paths for robots in dynamic and static environments. Some of the above methods cannot be carried out in an unknown environment and some cannot plan an optimal path.

At present, $Q$-learning is a widely applied reinforcement learning algorithm. The limitation of $Q$-learning is that it is trial-and-error learning, which requires constant iteration and is time-consuming. Thus, it needs to be improved to quickly plan a globally optimal path. Maoudj and Hentout [16] initialize the $Q$-table to accelerate convergence by presenting a distance-based reward function. Soong et al. [17] integrated the prior knowledge gained from FPA into the traditional $Q$-learning, which provided a good exploration basis for accelerating the learning of mobile robots. Xu and Yuan [18] increased the step length of movement and the direction of the robot to plan a fast and smooth path. Oh et al. [19] specified the initial $Q$-value of the traditional $Q$-learning through the fuzzy rule-based $Q$-learning, which speeded up learning and stabilized convergence. Yan and Xiang [20] initialized the $Q$-table by using inverse Euclidean distance from the current position to the goal position, which improves the efficiency of $Q$-learning. The above methods all initialize the $Q$-value simply by some prior information to improve the algorithm, without considering the avoidance of concave obstacles and the adjustment of the action selection strategy.

In summary, there are still many difficult problems in the formation and path planning of multiple robots in unknown environments. In this paper, we adopt the leader-following approach to study the multirobot dynamic formation problem. The innovation in this paper is as follows: The improved $Q$-learning algorithm is presented to plan paths, in which environmental guidance and virtual obstacle-filling avoidance strategy are added to accelerate convergence and the SA algorithm is applied to improve the action selection strategy; the follower robot can achieve the tracking strategy of GPF by designing the cost function to select actions.

## 2. Related Methods

### 2.1. Q-Learning Algorithm.
The $Q$-learning algorithm [21] is an RL algorithm based on temporal-difference, which combines the Monte Carlo sampling method and the bootstrap-

ping idea of dynamic programming. It is described with the Markov decision process as follows: Firstly, limited state space and action space are given. When the robot needs to accomplish a certain task, it selects and performs the action in the current state, which interacts with the environment. Then, the robot enters the next state and is given an instant reward as feedback by the environment. Finally, the value function is updated according to the update rule by using the reward which is passed to it. One round is continuing the above process until the robot reaches the target, and the rounds are iterated until the cumulative reward is maximum. The update equation of the $Q$ − value function is

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha \left[ r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right], \tag{1}$$

where $s_t$ is the state and $a_t$ is the action at current time $t$, $s_{t+1}$ is the state and $a$ is the action at next time $t + 1$, $r_t$ is the reward obtained by performing action $a_t$ at state $s_{t+1}$, $\alpha \in ( 0, 1)$ is the learning rate, and $\gamma \in (0, 1)$ is the discount factor.

In order to ensure the exploratory nature of the algorithm, the $\varepsilon$-greedy strategy is usually adopted, with the probability $1 − \varepsilon$ of selecting the action that maximizes the value function, and the small probability $\varepsilon$ that is still reserved for random exploration. The mathematical equation of the strategy is

$$\pi(s_t) = \begin{cases} a_{\text{random}}, & \text{if } \delta < \varepsilon, \\ \arg\max_a Q(s, a), & \text{else,} \end{cases} \tag{2}$$

where $\pi(s_t)$ is the selected strategy, $\varepsilon \in (0, 1)$ is the greedy factor, $\delta \in (0, 1)$ is a random number, $a_{\text{random}}$ is a randomly selected action, and $\arg\max_a Q(s, a)$ is an action that maximize $Q$-value at state $s$.

The classical $Q$-learning algorithm is described in Algorithm1.

### 2.2. APF Method.
The APF method is a virtual potential field established artificially, including the GPF and the repulsive potential field. The target generates a gravitational force on the robot to make the robot move towards it. The GPF function is

$$U_{\text{gra}}(q) = \frac{1}{2} \xi d^2 \left( q, q_{\text{goal}} \right), \tag{3}$$

where $\xi$ is the GPF factor, $q$ is the state position of the robot, $U_{\text{gra}}(q)$ is the GPF at $q$, $q_{\text{goal}}$ is the state position of the target to be reached by the robot, and $d(q, q_{\text{goal}})$ is the distance between the robot and the target, which can be measured by specific sensors in practice and is one-dimensional. Gravitation is the negative gradient of the GPF, and the gravitational function is defined as

$$F_{\text{gra}} = -\nabla U_{\text{gra}} = -\xi d \left( q, q_{\text{goal}} \right). \tag{4}$$

```
begin
    Initialization:
        Q(s, a) = {0}, ∀s ∈ S, a ∈ A(s) = {up, down, left, right}  %Initialize Q value with 0, determine the state set and the action set
containing four actions
    for(episode < m) %The episode cannot exceed m which is the maximum number of episodes
        Given initial state s₀;
        while (sₜ ≠ target state)
            (1) Select an action aₜ at state sₜ according to ε-greedy; % ε-greedy is the action selection strategy;
            (2) Execute the action aₜ, then enter state sₜ₊₁ and get a reward rₜ; %Get immediate rewards by performing actions to inter-
act with environment
            (3) Update Q(sₜ, aₜ) using Q(sₜ, aₜ) = Q(sₜ, aₜ) + α[rₜ + γ max Q(sₜ₊₁, a) − Q(sₜ, aₜ)];
                                                                      a
            % Update the value function according to the update equation by using the reward
            (4) sₜ ⟵ sₜ₊₁; %Update state
        end-while
        Episode = episode + 1; % Update episode
    end-for
end
```

ALGORITHM 1: Classical $Q$-learning algorithm.

Obstacles generate a repulsive force on the robot to make the robot move away from it. The repulsive potential field function is

$$U_{\text{rep}}(q) = \begin{cases} \dfrac{1}{2}\eta\left(\dfrac{1}{d(q, q_{\text{obst}})} - \dfrac{1}{d_0}\right), & \text{if } d(q, q_{\text{obst}}) \le d_0, \\ 0, & \text{if } d(q, q_{\text{obst}}) > d_0, \end{cases} \tag{5}$$

where $\eta$ is the repulsive potential field factor, $q$ is the state position of the robot, $U_{\text{rep}}(q)$ is the repulsive potential field at $q$, $q_{\text{obst}}$ is the obstacle state position, $d(q, q_{\text{obst}})$ is the distance between the robot and the obstacle, which can be measured by specific sensors in practice and is one-dimensional. $d_0$ is the influence radius of the obstacle, which is artificially set according to the experimental requirements in practice. Repulsion is the negative gradient of the repulsion potential field, and the repulsion function is defined as

$$F_{\text{rep}} = -\nabla U_{\text{rep}} = \begin{cases} \eta\left(\dfrac{1}{d(q, q_{\text{obst}})} - \dfrac{1}{d_0}\right)\dfrac{1}{d^2(q, q_{\text{obst}})}\nabla d(q, q_{\text{obst}}), & \text{if } d(q, q_{\text{obst}}) \le d_0, \\ 0, & \text{if } d(q, q_{\text{obst}}) > d_0. \end{cases} \tag{6}$$

Therefore, the traditional APF method guides the robot's direction of movement based on the combined force of gravitation and repulsion, but its shortcomings are as follows:

(1) When the robot is far away from the target, the gravitational force is much greater than the repulsive force, and it may hit an obstacle

(2) When the distance between them is relatively close, the obstacles will repel the robot too much to reach the target

(3) When the two reaction forces just cancel out, the phenomenon of local optimum or oscillation may appear

Because of the above shortcomings, the APF method generally cannot be used directly and needs to be improved to use.

### 2.3. SA Algorithm.

The idea of the SA algorithm comes from the solid annealing process, which is an algorithm that jumps out of the local optimum to get the global optimum. The algorithm uses temperature parameters $T$ to control convergence in a finite time. Firstly, the initial temperature and the end temperature are set. The algorithm starts from the initial state and takes it as the current state. Then, it generates a new state in its neighborhood and determines whether to accept the new state based on the Metropolis criterion. The generation process of the new state iterates while the $T$ decays until $T$ is the end temperature. Finally, the algorithm ends with the global approximate optimal solution.

The Metropolis criterion is that when a system enters a state $s_{\text{new}}$ due to a certain change in state $s_{\text{old}}$, the energy of the system correspondingly changes from $E(s_{\text{old}})$ to $E(s_{\text{new}})$ and then the accepted probability equation of the system from $s_{\text{old}}$ to $s_{\text{new}}$ is

$$P(s_{\text{old}} \longrightarrow s_{\text{new}}) = \begin{cases} 1, & E(s_{\text{new}}) \leq E(s_{\text{old}}), \\ e^{E(s_{\text{new}})-E(s_{\text{old}})/T}, & E(s_{\text{new}}) > E(s_{\text{old}}). \end{cases} \tag{7}$$

When $E(s_{\text{new}}) \leq E(s_{\text{old}})$, the new state is accepted as the current state. When $E(s_{\text{new}}) > E(s_{\text{old}})$, if $e^{E(s_{\text{new}})-E(s_{\text{old}})/T} > \delta$, the new state is accepted as the current state; otherwise, the new state is not accepted and the system remains in the current state. $\delta \in (0, 1)$ is the randomly generated number.

# 3. Improved $Q$-Learning Proposed for Path Planning of Leader Robot

*3.1. Environmental Guidance Based on GPF of Target.* The traditional $Q$-learning algorithm has no prior knowledge. In the early learning process, the robot's aimless exploration causes many invalid iterations and slow convergence. So, the idea of the APF method is introduced to guide moving. In this paper, the robot plans a path in an unknown environment where only the start and the target of the task are known. Due to the less environmental information and the shortcomings of the traditional APF method, only the GPF of the target is introduced to initialize the $Q$ value without considering the effect of the repulsive potential field. In order to make the target direction consistent with the increasing direction of the $Q$-value, the GPF function is constructed as

$$U'_{\text{gra}}(s) = \xi * \frac{d_{\text{aim}}(s)}{d_{\text{aim}}(s) + \eta}, \tag{8}$$

where $\xi$ is the GPF factor which is negative and controls the value inversely proportional to the distance, $d_{\text{aim}}(s)$ is the distance from the current position to the target, and $\eta$ is a positive constant that prevents the denominator from being 0.

When the robot moves, the instant reward is detected by sensors and the $Q$-table is initialized at the same time. Therefore, the instant reward of environmental information is added to the $Q$-value initialization. The purpose of RL is to maximize the cumulative reward by maximizing the $Q$-value. The robot always tends to choose the action with the maximum $Q$-value, which will guide the robot to move toward the target while avoiding obstacles. The mathematical equation of $Q$-value initialization with environmental guidance based on GPF of the target is

$$Q(s, a) = k_q * \left( r_q + \gamma * U'_{\text{gra}}(s) \right), \tag{9}$$

where $r_q = \begin{cases} 1, & \text{at target} \\ 0, & \text{else} \\ -1, & \text{at obstacle} \end{cases}$, $k_q$ is the scale coefficient adjusted according to the actual algorithm, $\gamma$ is the discount factor, and $U'_{\text{gra}}(s)$ is the GPF at state $s$.

*3.2. Virtual Obstacle-Filling Avoidance Strategy.* There will be concave obstacles in a more complex environment. The traditional $Q$-learning algorithm can escape from such obstacles through continuous exploration, which greatly extends the learning time. In addition, the robot is more likely to fall into concave obstacles and cannot escape after adding GPF guidance. In the grid map environment, the obstacle grid is the infeasible area and the rest are feasible areas. Setting certain key position grids which is feasible in the path of possibly tending to concave obstacles as infeasible areas can effectively fill and avoid concave obstacles. Therefore, a virtual obstacle-filling avoidance strategy is established for concave obstacles. The strategy is to judge whether the current grid possibly tends to concave obstacles by adding real-time detection information based on the target tendency before the robot takes the next step. Then, it fills non-obstacles on the path of possibly tending to the concave obstacle with virtual obstacles until the concave shape is filled. The filled concave obstacle as a whole becomes an infeasible area, so the robot will not fall into the concave obstacle in subsequent iterations. This strategy makes full use of sensors and the environmental information which have been learned. It not only prevents the robot from falling into concave obstacles but also reduces invalid exploration of some infeasible positions, which improves the efficiency of path planning.

The specific implementation of the virtual obstacle-filling avoidance strategy is as follows.

Firstly, the sensor is used to detect the position status and distance in real time. And a current position-action array is established to store the feasible adjacent positions from the current position. Before the robot moves, the Euclidean distance from the $3 * 3$ grid positions adjacent to the robot's current position to the target position is calculated in turn. Next are the specific judgement steps to judge whether the current grid possibly tends to concave obstacles according to the calculated distances.

If the distance is less than the distance from the current position to the target position, it is further judged whether this adjacent position is an obstacle or not. If the adjacent position is not an obstacle, it is feasible and will be added to the corresponding position of the current position-action array.

If the adjacent position is further away from the target or it is further judged an obstacle, it is an infeasible position and will not be added to the corresponding position.

If the final current position-action array is empty, it indicates that the current position completely tends to infeasible areas which may be in a concave obstacle. The current position will be filled with a virtual obstacle.

Finally, each step of the robot is judged until concave obstacles are filled.

One-step filling of the virtual obstacle-filling avoidance strategy is shown in Figure 1. In the figure, the red grid is the robot, and the yellow grid is the target. As Figure 1(a) shows, the robot enters the grey concave obstacle during the path planning process. According to the distance calculation, the adjacent positions which are down, right, and lower right of the robot's current position are determined

FIGURE 1: One-step filling of the virtual obstacle filling avoidance strategy [22].

```
begin
      Initialization:
            Q(s, a) = ∅, current(s, a) = ∅, n, T, ∀s ∈ S,
            a ∈ A(s) = {up, down, left, right, upleft, upright, downleft, downright}
            %Establish Q-table and current position-action array, define n as times of consecutive iterations, define T as the initial tem-
perature, determine the state set and the action set containing eight actions
      for (episode < m) %The episode cannot exceed m which is the maximum number of episodes
            Given initial state s₀;
            If episode%n == 0 Then use ε = e^(Q(s,a_random)−Q(s,a_max)−q)/T to calculate and update ε; %Adjust ε dynamically using T of the SA
algorithm.
            while (sₜ ≠ target state)
                  (1) If sₜ exists in the Q-table then continue to next step;
                        Else use Q(s, a) = k_q ∗ (r_q + γ ∗ ξ ∗ (d_aim(s′)/(d_aim(s′) + η))) to initialize Q(sₜ, a);
                        %Initialize the Q-table
                  (2) If (sₜ, a) is a feasible area towards the target then add it to the corresponding position of current(sₜ, a);
                        Else the corresponding position of current(sₜ, a) is kept empty;
                        % Add the feasible adjacent positions from the current position to the current position-action array
                  (3) If current(sₜ, a) is empty then (sₜ, a) is completely toward the infeasible area which possibly tends to concave obsta-
cles and fill it with virtual obstacle;
                        Select action aₜ in state sₜ according to ε-greedy which is improved by SA;
                        %Fill concave obstacles using the virtual obstacle-filling avoidance strategy while selecting actions
                  (4) Execute aₜ in sₜ, enter s_{t+1} and get rₜ;
                  (5) Update Q(sₜ, aₜ) using Q(sₜ, aₜ) = Q(sₜ, aₜ) + α[rₜ + γ max_a Q(s_{t+1}, a) − Q(sₜ, aₜ)];
                  (6) sₜ ⟵ s_{t+1};
            end-while
            Episode = episode + 1;
      end-for
end
```

ALGORITHM 2: Improved Q-learning algorithm.

as the positions which are near the target more and are the dark grey grid in Figure 1(b). The three adjacent positions are further judged obstacles which are infeasible positions, indicating that the current position is completely toward the infeasible area which may be in a concave obstacle. Thus, the current position is filled with light grey virtual obstacles, which is seen in Figure 1(c).

3.3. Action Selection Strategy Improved by SA. In the process of path planning with $Q$-learning, the robot expands the

range of movement by exploring the environment and accumulates knowledge of environmental rewards and punishments. Finally, it uses the value function to select the optimal action. In the robot's iterative learning process, more exploration is required in the early stage, but too much or too long exploration will greatly extend the learning time and reduce the learning efficiency. On the contrary, too little exploration will lead to insufficient experience and the action selected finally may be sub-optimal. Thus, it is necessary to balance exploration and utilization. The traditional $\varepsilon$-greedy strategy

FIGURE 2: The multi-robot action, step length, and detection range of sensor[22].

often used in the Q-learning algorithm balances exploration and utilization to a certain extent by setting $\varepsilon$. However, the fixed greedy factor in the learning process makes random actions selected with the same probability each time, which causes slow convergence and fluctuations after convergence. Therefore, the greedy factor needs to be adjusted dynamically with the learning process. One method commonly used in experiments is to set $\varepsilon$ to decrease at a fixed rate, but it is not universal to set a fixed rate of decrease based on experience.

In response to the above problems, the $\varepsilon$-greedy strategy improved by SA which is used to adjust $\varepsilon$ dynamically is proposed. The controlled temperature of SA is adjusted in real time according to the learning situation of the Q-learning algorithm. The algorithm explores as much as possible in the early stage of path planning to increase more prior knowledge and prevent local optimum and cancels unnecessary exploration when approaching convergence later. The steps of the action selection strategy improved by the SA algorithm are as follows:

(1) Define the temperature parameter $T$ and set the initial value $T_0$. Then, use the sample standard deviation of step numbers for $n$ consecutive iterations to control the cooling temperature. The mathematical equation of $T$ is

$$T = i + k * \sqrt{\frac{\left(\text{step}_{m+1} - \text{step}_{\text{avg}}\right)^2 + \cdots + \left(\text{step}_{m+n} - \text{step}_{\text{avg}}\right)^2}{n - 1}}, \tag{10}$$

where $\text{step}_{m+1}, \cdots, \text{step}_{m+n}$, respectively, are the number of steps for $n$ consecutive iterations, $\text{step}_{\text{avg}}$ is the average number of $n$ consecutive iterations, and $k$ is the control factor, which is obtained by repeatedly adjusting according to the experimental effect and controls $T$ in a suitable range. $i$ is a smaller non-zero constant to prevent $T$ from being 0 after convergence

(2) Calculate the accepted probability of randomly selected actions according to the Metropolis criterion. And use it to redefine the greedy factor $\varepsilon$ at $T$. The mathematical equation of $\varepsilon$ is



FIGURE 3: The flow chart of the leader robot's path planning.

$$\varepsilon = e^{(Q(s,a_{\text{random}}) - Q(s,a_{\text{max}}) - q)/T}, \tag{11}$$

where $Q(s, a_{\text{random}})$ is the Q-value of the random action selected at state $s$, $Q(s, a_{\text{max}})$ is the Q-value of the optimal action at state $s$, $q$ is a non-zero constant

FIGURE 4: The flow chart of the follower robot's path planning.

TABLE 1: Implementation details of $Q$-L1 to $Q$-L5 algorithms.

| Algorithm number | Implementation details |
| --- | --- |
| $Q$-L1 | Algorithm 1 |
| $Q$-L2 | $Q$-L1 with the dynamic greedy factor of SA |
| $Q$-L3 | $Q$-L2 with environmental guidance |
| $Q$-L4 | Algorithm 2 |
| $Q$-L5 | $Q$-L4 with a modified reward function |

to prevent the numerator from being 0, and $T$ is the temperature parameter

(3) If $\delta < \varepsilon$, choose the action randomly, otherwise choose $\arg \max_a Q(s, a)$, where $\delta$ is the randomly generated number

### 3.4. Improved Q-Learning Algorithm.
Compared with the original algorithm, there are three innovations in the improved $Q$-learning algorithm proposed in this paper.

Firstly, the $Q$-table of the original $Q$-learning algorithm is initially a zero-value table without any prior knowledge. The improved $Q$-learning algorithm uses the GPF of the known target in the task to initialize the $Q$-table, which adds environmental guidance and reduces invalid exploration.

Secondly, the robot moves immediately after selecting an action in the original algorithm. This algorithm designs a virtual obstacle-filling avoidance strategy for judgment before each step. It fills non-obstacles which is judged to tend to concave obstacles with virtual obstacles.

Finally, the original algorithm uses the traditional $\varepsilon$-greedy strategy to select actions. The strategy improved by the SA algorithm is proposed in the new algorithm. It adjusts $\varepsilon$ dynamically by adjusting the temperature in real time according to the learning situation of $Q$-learning.

The steps of the improved $Q$-learning algorithm are shown in Algorithm 2.

## 4. A Path Planning Method for multirobot Formation

### 4.1. Tracking Strategy Based on GPF for Follower Robot.
The steps of the tracking strategy based on GPF for the follower robot are as follows:

Step 1: if the follower robot obtains the coordinates broadcast by the leader robot, it will determine the next target state according to the formation, i.e., the desired target position at this time. Otherwise, it means that the formation has reached the target position and the path planning ends.

Step 2: the follower robot moves to the target position. Firstly, the robot uses the cost function to calculate the cost

(a) Q-L1

(b) Q-L2

(c) Q-L3

(d) Q-L4

FIGURE 5: Continued.

(e) Q-L5

Figure 5: Path planning maps of 5 $Q$-learning algorithms.

for the eight neighboring states of the current state, which determines the state s with the smallest cost. Then, it selects the corresponding action and executes it. At the same time, it adopts the virtual obstacle-filling avoidance strategy in parallel with the leader robot to share information. Specifically, the cost function is designed by using the idea of GPF. The GPF of the target to the current position is measured by the Euclidean distance from the current position to the target position, which is proportional to the distance. When the checked state is an obstacle, the penalty function $R_{static}$ is given a positive value; otherwise, the value is 0. The equation for measuring the GPF is

$$d_{attr} = \sqrt{\left(x_s - x_{goal}\right)^2 + \left(y_s - y_{goal}\right)^2}. \quad (12)$$

The cost function equation is

$$C(s_t, a_t) = c * d_{attr} + R_{static}(s_t, a_t), \quad (13)$$

where $d_{attr}$ is the GPF which is measured, $x_s$ is the horizontal coordinate of the current state, $y_s$ is the vertical coordinate of the current state, $x_{goal}$ is the horizontal coordinate of the target at this moment, $y_{goal}$ is the vertical coordinate of the target at this moment, $C(s_t, a_t)$ is the cost function at $t$, $s_t$ is the state at $t$, $a_t$ is the action at $t$, $c$ is the adjustment coefficient, and $R_{static}(s_t, a_t)$ is the penalty function.

Step 3: if the state with the minimum cost entered is the target state at this time, return to step 1 and continue. If the state is not the target state at this time, go to step 2 and continue.

4.2. Design Scheme of Path Planning for Leader-following Formation. Adopting the leader-following method, the

design scheme of path planning for leader-following formation proposed in this paper includes three parts:

(1) Initialization: the grid environment is adopted, and the starting position and target position of the multiple robots are determined. A leader-following formation is designed and the robots are divided into two types: leader and follower. Then, one robot is selected as the leader or a virtual robot is supposed to act as the leader, and the rest are follower robots. Multiple robots have eight actions including up, down, left, right, upper left, upper right, lower left, and lower right. Each robot is equipped with a sensor, which can detect the environmental information of the 3 ∗ 3 grids centering on its position. The multirobot action, step length, and detection range of the sensor are shown in Figure 2

(2) Path planning of leader robot: the leader robot is responsible for planning the path. It uses the improved $Q$-learning algorithm to plan a globally optimal path with avoiding simple obstacles and concave obstacles after trial-and-error training. At the same time, it broadcasts the position of each step and some environmental information to the follower robot. The process of the leader robot's path planning is shown in Figure 3

(3) Local following of follower robot: the follower robot is responsible for following the leader robot to maintain the formation according to the requirements. When the follower robot receives the position

(a) Q-L1



(b) Q-L2



(c) Q-L3



(d) Q-L4



(e) Q-L5

FIGURE 6: The cumulative reward change graphs with rounds of 5 Q-learning algorithms.

information broadcast by the leader robot, it determines the desired target depending on the formation. Then, it follows locally using the tracking strategy based on the GPF and avoids obstacles autonomously. The process of the follower robot's path planning is shown in Figure 4

## 5. Experiments Analysis

According to the design scheme of path planning for multiple robots, the method is tested experimentally. The experiment uses Python standard GUI toolkit Tkinter to establish simulation environments.

*5.1. Comparison Experiments of Improved Q-Learning Algorithm.* For the improved Q-learning algorithm, a grid map with three elements: starting point, target point, and obstacles, is first established. The map size is set to $20 \times 20$

grids, and the resolution of each grid is $26 \times 26$ pixels. The starting position of the robot represented by a red grid is set at (0, 0), and the target position represented by a yellow grid is at (19, 19). Obstacles which are black grids are randomly placed on the map, including concave and simple obstacles. To distinguish actual obstacles from virtual obstacles filled during the algorithm operation, virtual obstacles are gray grids.

The experiment is carried out in a comparative way, and five algorithms are implemented: Q-L1 is the traditional Q-learning algorithm, Q-L2 is the Q-learning algorithm with the dynamic greedy factor of SA, Q-L3 adds environmental guidance of GPF on the basis of Q-L2, Q-L4 is the Q-learning algorithm proposed in this paper with the improvements 3.1, 3.2, and 3.3, Q-L5 is the Q-learning algorithm with a modified reward function based on Q-L4. The implementation details of Q-L1 to Q-L5 algorithms are shown in Table 1.

(a) Q-L1

(b) Q-L2

(c) Q-L3

(d) Q-L4

(e) Q-L5

FIGURE 7: The step numbers change graphs with rounds of 5 Q-learning algorithms.

TABLE 2: Comparison table of Q-L1 to Q-L5 algorithm performance.

| Algorithm | Potential field | Filling | r | ε | Convergence time | Convergence round | Steps | Length |
|---|---|---|---|---|---|---|---|---|
| Q-L1 | No | No | r1 | 0.2 | 479.0037 | 2760 | 30 | 37.4558 |
| Q-L2 | No | No | r1 | Dynamic | 379.4773 | 3520 | 30 | 36.0416 |
| Q-L3 | Yes | No | r1 | Dynamic | No | No | No | No |
| Q-L4 | Yes | Yes | r1 | Dynamic | 5.7183 | 120 | 26 | 32.6274 |
| Q-L5 | Yes | Yes | r2 | Dynamic | 48.3259 | 1010 | 22 | 28.6274 |

The same parameter settings of the algorithm are as follows: the maximum number of iteration rounds is 10000, the learning rate $\alpha$ is 0.01, and the discount factor $\gamma$ is 0.9. For using the traditional $\varepsilon$-greedy strategy in the algorithm, the greedy factor $\varepsilon$ is 0.2, and the convergence is determined that the standard deviation of step numbers for 10 consecutive

iterations is less than 5. Parameter settings for using SA in the algorithm are as follows: the initial temperature $T_0$ is set to 10, the number of consecutive iterations $n$ is set to 10, the constant $i$ is set to 0.1, the control factor $k$ is set to 0.03, and the non-zero constant $q$ is set to 1. In the algorithm using the GPF method to improve, the GPF factor $\xi$

(a) In an obstacle-free environment

(b) In a static obstacle environment

FIGURE 8: Map of multi-robot formation's path planning experiments.

is set to -10, the constant $\eta$ is set to 735, and the scale coefficient $k_q$ is set to 0.1. The reward function is set to

$$r1 = \begin{cases} 1, & \text{reach the target,} \\ 0, & \text{else,} \\ -1, & \text{hit an obstacle,} \end{cases}$$

$$r2 = \begin{cases} 1, & \text{reach the target,} \\ -0.05, & \text{else,} \\ -1, & \text{hit an obstacle.} \end{cases}$$

(14)

After setting the parameters, simulation experiments are conducted. The path planning map, the cumulative reward change graph with the round, and the path planning step numbers change graph with the round are obtained. From the figure, the path planning and convergence of each algorithm can be seen. Figures 5(a)–5(e), respectively, show the path planning maps of the robot under algorithms $Q$-L1 to $Q$-L5. Figures 6(a)–6(e), respectively, show the change of cumulative reward with rounds for the robot under algorithms $Q$-L1 to $Q$-L5. Figures 7(a)–7(e), respectively, show the change of step numbers with rounds for the robot under algorithms $Q$-L1 to $Q$-L5.

Figure 5(c) shows that the robot is trapped in a concave obstacle and cannot escape. Figure 6(c) shows that the cumulative reward curve of path planning changes irregularly during the iterative process. Figure 7(c) shows that the step number curve of path planning changes irregularly during the iterative process. The above three results of $Q$-L3 ndicate the algorithm does not converge in the iterative process, and the robot cannot reach the target when only adding the GPF of the target to improve when encountering concave obstacles.

Figures 5(a)–5(e) show that the robot uses the $Q$-L1, $Q$-L2, $Q$-L4, and $Q$-L5 algorithms to effectively avoid black obstacles and plan a red path from the starting to the target, but the path planned by the $Q$-L1 and $Q$-L2 algorithms is more tortuous, the $Q$-L4 algorithm plans a smoother feasible path, and the $Q$-L5 algorithm plans the optimal path. The cumulative reward curve showed by Figures 6(a)–6(e) and the step number curve showed by Figures 7(a), 7(b), 7(d), and 7(e) are both stable after iterating to a certain round, indicating that the algorithms gradually converge as the iterations proceed.

However, curves of Figures 6(a) and 7(a) converges with small fluctuations, curves of Figures 6(b) and 7(b) converges with smoothness, curves of Figures 6(d), 7(d), 6(e), and 7(e) reach smoothness in fewer rounds of iteration. The above shows that by adding the improvements proposed in this paper to $Q$-L4 algorithm and $Q$-L5 algorithm, the experiments achieve better results. The robot moves while initializing the $Q$-value by the environmental guidance based on the GPF of the target, which makes the robot guided by the target direction all the time. It removes invalid movement and speeds up the convergence time. When the robot encounters a concave obstacle, it identifies effectively the infeasible area and fills it with light gray virtual obstacles to prevent the robot from falling into the concave obstacle. The SA method is used to dynamically adjust the greedy factor to accelerate the algorithm convergence and make it stable after convergence. In addition, the $Q$-L5 algorithm adjusts the reward function on the basis of the $Q$-L4 algorithm by giving each step a smaller penalty, and the robot learns the optimal path with the maximum cumulative reward.

Table 2 compares the performance of the above five algorithms after path planning. The data in the table are the average results from conducting several experiments. The analysis is as follows: based on the traditional $Q$-learning

algorithm, $Q$-L2 uses the SA algorithm to improve $\varepsilon$-greedy strategy. Although the convergence round of the algorithm increases, the convergence time is shortened and the overall stability of path planning is improved. Comparing the algorithms $Q$-L3 and $Q$-L4, if the environmental guidance of GPF is added to the algorithm without the virtual obstacle-filling avoidance strategy, the algorithm is difficult to converge when encountering concave obstacles. Comparing the algorithms $Q$-L2 and $Q$-L4, by adding environmental guidance and the virtual obstacle-filling avoidance strategy, the convergence time is reduced by 98.5%, and the convergence rounds is reduced by 96.6%; the total step numbers and the length of the path are stabilized at 26 and 32.6274, respectively. Comparing the algorithms $Q$-L4 and $Q$-L5, adjusting the reward function reasonably on the improved $Q$-learning algorithm proposed in this paper will make the robot plan the optimal path quickly, which is 89.9% shorter than the traditional $Q$-learning algorithm. And the number of convergence rounds is reduced by 63.4%, the step numbers are reduced to 22 and the length of the path is reduced to 28.6274.

*5.2. Experiments of Path Planning for multirobot Formation.*
After experimenting with the improved algorithm of the leader robot's path planning, the follower robot is added to verify the effectiveness of the path planning method for multirobot formation in this paper. The experiment uses a triangular formation and three robots. The leader robot is represented by a red grid, and its initial position is $(2, 2)$. The follower robots are represented by a blue grid and a green grid, respectively, and their initial positions are $(0, 2)$ and $(2, 0)$, respectively. The target position of the leader robot is $(19, 19)$, which also determines the target position of the follower robots. The leader robot uses the improved $Q$-learning algorithm $Q$-L5 to plan the optimal path with the same parameter settings as in 5.1 simulation experiments. The two follower robots, respectively, use the tracking strategy based on the GPF to follow: the penalty function $R_{\text{static}} = 1$ and the adjustment coefficient $c = 0.09$.

The experimental effect is shown in Figure 8(a) that the three robots quickly reach the target with a fixed triangle formation in an obstacle-free environment. In an environment with static obstacles, the leader robot moves first, and the two follower robots immediately move to the corresponding positions in the formation. The green follower robot firstly encounters a black lateral obstacle. It moves along the obstacle to the target direction and smoothly avoids the obstacle. Then, it continues to accelerate to move to the corresponding position of the formation at the current time to maintain the formation. Finally, the leader robot plans a red path, the two follower robots avoid obstacles by themselves during the following process and plan a green path and a blue path, respectively. The three robots reach the target at the same time and complete the formation task. The experimental effect is shown in Figure 8(b).

## 6. Conclusion

In this paper, by combining the improved $Q$-learning algorithm and the idea of the GPF method, a method for multirobot formation and path planning is proposed. The division of labor among multiple robots is clear. The leader robot uses the improved $Q$-learning algorithm to plan the path. It is found that adding environment guidance of the target's GPF and virtual obstacle-filling avoidance strategy effectively accelerates iterative convergence and avoids concave obstacles. It is stable and efficient for the action selection strategy to be improved by the SA method. At the same time, the follower robot uses a tracking strategy based on the improved GPF to follow in real time, which is simple and efficient. This formation method effectively solves the formation and path planning problems of multiple robots in an unknown environment with concave obstacles. In the future, the multirobot formation will be further studied in the context of dynamic environments and privacy protection.

## Data Availability

The data used to support the findings of this study are included within the article.

## Conflicts of Interest

The authors declare that there is no conflict of interest regarding the publication of this paper.

## Acknowledgments

## References

[1] K. K. Oh, M. C. Park, and H. S. Ahn, "A survey of multi-agent formation control," *Automatica*, vol. 53, pp. 424–440, 2015.

[2] A. Muxfeldt, D. Kubus, and F. M. Wahl, "Developing new application fields for industrial robots - four examples for academia-industry collaboration," in *2015 IEEE 20th Conference on Emerging Technologies & Factory Automation (ETFA)*, pp. 1–7, Luxembourg, Luxembourg, 2015.

[3] G. Dissanayake, J. Paxman, J. V. Miro, O. Thane, and H. Thi, "Robotics for urban search and rescue," in *First International Conference on Industrial and Information Systems*, pp. 294–298, Tirtayasa, Indonesia, 2006.

[4] M. Z. Rashid, F. Yakub, S. A. Zaki et al., "Comprehensive review on controller for leader-follower robotic system," *Indian Journal of Geo-Marine Sciences*, vol. 48, no. 7, pp. 985–1007, 2019.

[5] M. A. Lewis and K. H. Tan, "High precision formation control of mobile robots using virtual structures," *Autonomous Robots*, vol. 4, no. 4, pp. 387–403, 1997.

[6] P. Wang and Z. Geng, "Leader-follower formation control of multirobot systems using the dynamic surface approach," in *The 35th China Control Conference (CCC)*, pp. 7757–7762, Chengdu, China, 2016.

[7] W. Yin and X. Yang, "A totally Astar-based multi-path algorithm for the recognition of reasonable route sets in vehicle

navigation systems," *Procedia-Social and Behavioral Sciences*, vol. 96, pp. 1069–1078, 2013.

[8] C. Chen, X. Q. Chen, F. Ma, X. J. Zeng, and J. Wang, "A knowledge-free path planning approach for smart ships based on reinforcement learning," *Ocean Engineering*, vol. 189, article 106299, 2019.

[9] P. Sudhakara, V. Ganapathy, B. Priyadharshini, and K. Sundaran, "Obstacle avoidance and navigation planning of a wheeled mobile robot using amended artificial potential field method," *Procedia Computer Science*, vol. 133, pp. 998–1004, 2018.

[10] Jian Chen, Dong Sun, Jie Yang, and Haoyao Chen, "Leader-follower formation control of multiple non-holonomic mobile robots incorporating a receding-horizon scheme," *International Journal of Robotics Research*, vol. 29, no. 6, pp. 727–747, 2010.

[11] M. Sruthi, K. Rao, and V. Jisha, "Vector field based formation control of multirobot system," *IFAC-PapersOnLine*, vol. 49, no. 1, pp. 189–194, 2016.

[12] H. Sang, Y. You, X. Sun, Y. Zhou, and F. Liu, "The hybrid path planning algorithm based on improved A∗ and artificial potential field for unmanned surface vehicle formations," *Ocean Engineering*, vol. 223, no. 3–4, article 108709, 2021.

[13] P. Das and P. Jena, "Multirobot path planning using improved particle swarm optimization algorithm through novel evolutionary operators," *Applied Soft Computing*, vol. 92, article 106312, 2020.

[14] H. Qu, K. Xing, and T. Alexander, "An improved genetic algorithm with co-evolutionary strategy for global path planning of multiple mobile robots," *Neurocomputing*, vol. 120, no. 23, pp. 509–517, 2013.

[15] A. Lazarowska, "Discrete artificial potential field approach to mobile robot path planning," *IFAC-PapersOnLine*, vol. 52, no. 8, pp. 277–282, 2019.

[16] A. Maoudj and A. Hentout, "Optimal path planning approach based on Q-learning algorithm for mobile robots," *Applied Soft Computing*, vol. 97, no. 2020, article 106796, 2020.

[17] L. E. Soong, O. Pauline, and C. K. Chun, "Solving the optimal path planning of a mobile robot using improved Q-learning," *Robotics and Autonomous Systems*, vol. 115, pp. 143–161, 2019.

[18] X. Xu and J. Yuan, "Path planning for mobile robot based on improved reinforcement learning algorithm," *Journal of Chinese Inertial Technology*, vol. 27, no. 3, pp. 314–320, 2019.

[19] C. H. Oh, T. Nakashima, and H. Ishibuchi, "Initialization of Q-values by fuzzy rules for accelerating Q-learning," in *IEEE World Congress on IEEE International Joint Conference on Neural Networks*, pp. 2051–2056, Anchorage, AK, USA, 1998.

[20] C. Yan and X. Xiang, "A path planning algorithm for UAV based on improved Q-learning," in *2018 2nd International Conference on Robotics and Automation Sciences (ICRAS)*, pp. 1–5, Wuhan, China, 2018.

[21] C. Watkins and P. Dayan, "Q-learning," *Machine Learning*, vol. 8, no. 3-4, pp. 279–292, 1992.

[22] https://link.springer.com/book/10.1007%2F978-981-16-6320-8.

*Research Article*

# Metal Detection of Wood Based on Thermal Signal Reconstruction Algorithm

**Hong Zhang [iD],[1,2] Ruizhen Yang,[3] Wenhui Chen,[1,2] and Ruikun Wu[1,2]**

[1]*Key Laboratory of Non-destructive Testing Technology (Fujian Polytechnic Normal University), Fujian Province University, China*
[2]*School of Electronic and Mechanical Engineering, Fujian Polytechnic Normal University, China*
[3]*College of Civil Engineering, Changsha University, China*

Correspondence should be addressed to Hong Zhang; zhhgw@hotmail.com

In this paper, eddy current thermography is used to detect metal in wood materials, and thermal signal reconstruction (TSR) algorithm has been proposed to solve the problem of low resolution of metal detection. The basic principle of current nondestructive testing technologies for wood materials has been briefly reviewed, and the advantages and disadvantages have been analyzed. TSR algorithm can significantly enhance the contrast ration between metal and surrounding areas, different quantities of metal can be effective identified, and metal positions can be accurately realized. The experimental results show that the proposed eddy current thermography technology can quickly detect metal in wood materials and improve the efficiency and accuracy. The size and quantity of metal can be intuitively observed through thermal images.

## 1. Introduction

With the rapid development of modern industry, China's wood processing industry has developed rapidly. In 2021, the global wood market will reach 41.273 billion US dollars [1]. The wood trading market is an important part of the commodity trading market, which plays an important role in promoting the trading and circulation of wood products and stimulating the local economy [2]. In the process of wood production, there are metal materials in wood, which have a negative impact on the use and commercial value of wood. The existence of metal in wood usually reduces the strength of wood and also affects the appearance and processing process of wood. Therefore, the detection of metal in wood can improve the use safety and maximize the economic benefits.

At present, nondestructive testing of wood materials mainly includes defect detection and mechanical property measurement [3]. The basic methods of nondestructive testing mainly include stress wave method, mechanical stress deformation method, vibration method, microdrilling resistance method, ray method, and radar wave method. According to different wood materials, there are different detection methods, and the metal detection technology for ancient building wood materials is most widely used by stress wave method [4]. The general principle of stress wave for metal detection is that when impact force is applied to wood material, stress waves will be generated inside wood and propagate around. Sensors at both ends are used to receive signals of stress waves. The time difference between two points is calculated in a timer, and then, the propagation speed change of stress waves is obtained to judge the condition of metal inside wood materials. Compared with CT, X-ray, and the like, stress wave in wood materials has the advantages of lower cost, safety and reliability, harmlessness to human body, unaffected by tested materials and sizes, suitability for various environments, and can accurately judge whether there are metals, cavities, and wood knots in wood materials. However, the propagation of stress wave in wood is a complex dynamic process, which is affected by many factors, including

the following aspects: the properties of wood, microstructure of wood, water content of wood, defects of wood, and wood morphology [5]. Additional, stress wave systems can only be used for qualitative testing. Additional NDT methods are required to obtain quantitative results in regard to the size and depth of metal.

This paper proposes eddy current thermography for wood metal detection. Eddy current shielded by vacuum magnetic can effectively improve the sensitivity of ferromagnetic metal. It is a nondestructive and noncontact detection technology based on eddy current effect. It has the characteristics of high linearity, high resolution, fast response, simple structure, and static and dynamic measurement [6]. Eddy current excited thermography technology employed the different thermal radiation physical characteristics of structures or materials to detect various defects and damages on the surface or inside of materials. The obtained thermal images have the disadvantages of fuzzy edge, noise interference, and low resolution. In order to improve the accuracy, efficiency, and resolution of defect detection, different feature extraction algorithms have been used to extract defect information. Xingwang et al. carried out wavelet transform on thermal image sequence, and image fusion algorithm based on pixel level and feature level has been used to process thermal image sequence [7]. The test results of aluminum alloy samples show that the image fusion algorithm can effectively reduce the adverse effects of uneven heating and background noise on defect recognition. To enhance cracks characteristics in the original IR images, Peng et al. applied eddy current pulsed thermography (ECPT) for motor winding defects detection with fast Fourier transform (FFT) and principal component analysis (PCA) by eliminating the nonuniform heating effect [8]. L-shaped ferrite magnetic open sensing structure was proposed for fatigue crack inspection on metallic materials with anomalistic geometry. The modified eddy current pulsed thermography system has better performance in omnidirectional microfatigue crack detection. He et al. discussed the applications of deep learning applied infrared imaging-based machine vision. The principle, cameras, and thermal data of infrared imaging-based machine vision have been reviewed [9].

He et al. used fast Fourier transform (FFT) to process phase-locked thermal imaging data, and its calculation speed is faster than discrete Fourier transform and can observe defect information in frequency domain [10]. The fitting function relationship is used to realize the quantitative recognition of defects in infrared thermal wave detection. Numerical calculation method is used to provide samples for training neural network, which proves the feasibility of the method. Rajic employed principal component analysis (PCA) method to decompose the thermal image sequence into a group of orthogonal statistical patterns by singular value decomposition [11]. PCA is used to reduce redundancy, remove noise, and improve the accuracy of detection. Liang et al. used wavelet transform and PCA to detect the impact defects of composite materials [12]. Sripragash and Sundaresan used thermal signal reconstruction (TSR) to detect the defect depth. Temporal and spatial resolutions of thermal image sequence have been improved [13]. Hyvarinen and Oja employed independent component analysis (ICA) method that is used to extract independent components in thermal image sequence to remove data redundancy and obtain high-order statistical characteristics [14]. Świta and Suszyński used kd-tree algorithm to cluster infrared thermal image sequence to extract depth information and reduce the amount of data [15]. Maldague and Marinetti proposed pulse phase infrared thermography (PPT) algorithm, which transforms the time and space information into the frequency domain through Fourier transform to obtain the phase and amplitude information. The defect information can be extracted through the difference of the phase and amplitude between the defect and nondefect regions [16]. A hybrid multidimensional feature fusion structure of spatial and temporal segmentation model was proposed by Hu et al. for defect detection with thermography. The semantic information can be captured easily. He et al. made a profound study infrared machine vision and infrared thermography with deep learning [17]. Theoretical research and case study method are used in this review paper.

In order to improve the detection accuracy and metal resolution, this paper employed thermal signal reconstruction algorithm to detect metal in wood. The metal materials in wood are measured with eddy current thermography, and the infrared thermal images are analyzed by the proposed TSR algorithm. Compared with the stress wave method, it has advantage of nondestructive testing. Furthermore, the position, size, and number of metal materials are detected.

The rest of paper is organized as follows: Firstly, the proposed method is introduced in Section 2. The experimental set-up is described, and feature extraction and optimization are introduced in Section 3. Then, wood with different metal are characterized. It can prove the accuracy and efficiency of eddy current thermal imaging method in the detection of wood materials. Finally, conclusions are outlined in Section 4.

## 2. Methods and Image Processing

*2.1. Principle of Eddy Current Thermography.* As shown in Figure 1, the measurement device is mainly composed of an excitation system, an excitation coil, IR camera, a cooling system, an excitation system, sample under test, and a PC. Thermal information of eddy current and materials under test is obtained by IR camera. Different types of information can be obtained according to different analysis methods, and corresponding defect information can be obtained by analyzing these information. Eddy current thermography is based on electromagnetic induction, which involves many physical processes such as Joule heating, heat conduction, and infrared radiation. When the excitation coil carrying high-frequency alternating current is close to the conductor to be tested, under the action of the magnetic field of the coil, eddy current will be generated in the place where there are metal bodies or defects in the conductor to be tested, and eddy current will generate heat in the place where there are foreign bodies or defects in the sample under test, causing temperature changes on the surface of the material and from the inside through heat conduction. The information of foreign bodies or defects in materials can be obtained by graphic analysis and processing collected by IR camera.

FIGURE 1: Diagram of eddy current thermography.

The eddy current thermography detection technology can evaluate the metal in the reflection mode and the penetration mode, respectively [18]. With eddy current thermography detection technology in penetration mode can easily detect the surface fracture structure caused by metal. But there are the following disadvantages: (1) Due to the shape of the coil, it will bring uneven heating effect; (2) As time increases, lateral blurring will occur; (3) Periodic wood structure causes thermal abnormalities. Therefore, eddy current thermography in the reflection mode has been employed for the metal evaluation in wood.

*2.1.1. Electromagnetic Induction Heating.* When the excitation coil passes through alternating current with frequency $f$, induced eddy current with the same frequency is generated inside the tested material according to the law of electromagnetic induction. Time-varying equation of eddy current excitation in eddy current pulse thermal imaging is as follows:

$$J_e + \nabla \times \left(\frac{1}{\mu} \nabla \times A\right) - \frac{\sigma}{\sqrt{\mu\varepsilon}} \times (\nabla \times A) = \sigma \frac{V_{\text{loop}}}{2\pi r_d} + J_s. \quad (1)$$

Among them, $\mu$ is the magnetic permeability of the measured material, and $\varepsilon$ is the dielectric constant and eddy current density of the measured material:

$$J_e = \sigma \frac{\partial A}{\partial t}, \quad (2)$$

where $J_e$ is the current density of the excitation coil. $V_{\text{loop}}$ is the loop potential, and $r_d$ is the loop radius, which is the conductivity of the material. $A$ is the magnetic vector potential instead of the magnetic induction intensity $B$ to satisfy:

$$B = \nabla \times A. \quad (3)$$

Due to the resistance inside the material, eddy current is converted from electric energy to heat energy inside the material. According to Joule's law, the generated thermal power $P_w$ is proportional to the eddy current density $J_e$ and the electric field strength $E$:

$$P_w = \frac{1}{\sigma}|J_e|^2 = \frac{1}{\sigma}|\sigma E|^2. \quad (4)$$

*2.1.2. Heat Conduction.* The generated Joule heat $Q$ propagates inside the material, and the propagation process follows the formula.

$$\rho C_p \frac{\partial T}{\partial t} - \nabla(\sigma_T \nabla T) = Q, \quad (5)$$

where $\rho$ is the density of the material, $C_p$ is the specific heat capacity of the material, $T$ is the thermal conductivity of the material, and $t$ is the temperature of the material. In the experiment, the magnetic induction intensity $B$ around the infinite straight wire is defined by the formula:

$$B = \frac{\mu I}{2\pi h}, \tag{6}$$

where $H$ is the distance to the straight wire. It can be obtained that the magnetic field strength of $B$ decays rapidly with the increase of the distance to the coil. The thermal power in induction heating is proportional to the square of eddy current density, which can be obtained from (6).

*2.1.3. Infrared Radiation.* According to Stefan-Boltzmann's law, an object whose temperature is higher than zero degree Kelvin will spontaneously generate infrared radiation outward.

$$J^* = \varepsilon \sigma_{sb} T^4, \tag{7}$$

where $\varepsilon$ is the emissivity of the material, $\sigma_{Sb}$ is Stezmann-Boltzmann constant, and $T$ is the absolute temperature.

*2.2. Thermal Signal Reconstruction (TSR).* The reconstruction of thermal signal sequence is based on one-dimensional heat conduction equation, and the surface temperature response equation of applying instantaneous uniform excitation to thick materials is as follows:

$$\frac{\partial^2 T}{\partial \chi^2} + \frac{1}{k} g(x, t) - \frac{1}{\alpha} \frac{\partial T}{\partial t} = 0, \tag{8}$$

$$g(x, t) = Q\delta(x)\delta(t)\alpha = \frac{k}{\rho c}, \tag{9}$$

where $Q$ is the energy applied to the surface, $K$ is the thermal conductivity, $\rho$ is the density of the material to be detected, and $c$ is the specific heat capacity;

$$T(t) = \frac{Q}{e\sqrt{\pi t}}. \tag{10}$$

Polynomial fitting is performed on it:

$$\ln [\nabla T(t)] = \sum_{n=0}^{N} a_n [\ln (t)]^n. \tag{11}$$

The original data is reconstructed when the coefficient $a_n$ is fitted from Equation (11) as a function of the change of temperature with time at each point

$$\nabla T(t) = \exp \left( \sum_{n=0}^{N} a_n [\ln (t)]^n \right). \tag{12}$$

After reconstruction from Equation (11), differential operation can be performed, so that first-order and second-order differential can be performed. The image and differential obtained after reconstruction of any point of the heat map sequence are obtained by Equation (12). The thermal imaging image processed by TSR increases the spatial and temporal resolution of the thermal image. Between (1) and (5), it is known that the heat generated inside the tested material and its conduction are directly affected by the electrical conductiv-

ity and thermal conductivity of the material, and the temperature of the area where the wood material has metal matter will be significantly different from that of the nondefective area. Radiation energy also has certain influence on thermal conductivity.

The location of metal area in the measured wood material can be observed from infrared thermal imaging to capture the surface temperature thermal image of wood material. At the end, TSR algorithm is used to process data to evaluate metal in wood. The TSR algorithm employs the temporal and spatial variation information of surface temperature to process the temporal information of each pixel in the thermal image sequence and transforms the temperature response curve of each pixel from the time domain to the logarithmic domain. From Equation (12), it can be seen that the temperature change curve of the nonmetal area satisfies the linear relationship, and the temperature change curve of the metal area is nonlinear.

## 3. Experimental Study

The samples under test are two pieces of dry wood materials with a width of 42 mm. In Figure 2, two blocks contain different amounts of metal foreign matter, marked as L1 and L2. The physical diagram of eddy current thermal imaging system is shown as Figure 1. The power source of the excitation induction heating system is MDS-GLY-01, the input voltage is single-phase 220 V/50 Hz, the operation frequency is 150 kHz-250 kHz, and a circular excitation coil is adopted. Specification model of water cooling equipment is MDS-SL-03. For long-wave infrared thermal camera model, FLIR A655SC, its resolution is $640 \times 480$. The speed of full frame 16-bit data is 50 fps. The metal body-containing regions were placed under coil, and the excitation voltage was 58 V, the excitation current was 339 A, and the excitation frequency was 1055 Hz. The excitation time was 1500 ms.

## 4. Results and Discussion

After heating, the frame is selected. The obtained infrared thermal image has been analyzed. During the experiment, the environmental interference is eliminated. As can be seen from Figure 3, the temperature of metal-free wood area is blue area, which means temperature remains constant. The representing metal is at the red dot. The fitting graph of transient temperature is increasing with time. The locations of the metals can be determined from the infrared thermal images due to the effect of thermal diffusion whereas it is difficult to identify the real size of a metal.

The temperature rise in the excitation coil area without metal is almost the same. When there is metal in the specimen, the temperature rise curve with metal is obviously higher than that of without metal, and the temperature rise changes more when the metal is close to the excitation coil. At the end of heating, the temperature rise of the metal-free area and the metal area in L1 is about 14°C, and the temperature rise of L2 is about 24°C. It takes a certain time for metal in wood materials to affect the change of surface temperature. After the excitation time is over, the temperature of the metal drops

(a)                                                    (b)

FIGURE 2: Sample under test L1 (a)/L2 (b).



(a)                                                    (b)

FIGURE 3: Thermal imaging for samples under test ((a) L1; (b) L2).



FIGURE 4: PCA results for sample under test.

slowly, which is due to the poor heat dissipation of the wood material. The heat generated by the metal will be stored in the wood material for a certain period of time. Therefore, when the excitation time is over, the temperature of the metal area drops slowly, and the curve drops slowly. On the other hand, the temperature rise of metal material shows that the

FIGURE 5: ICA results for sample under test.



FIGURE 6: TSR results for sample under test.

metal in wood material has produced eddy current by excitation power supply. As a result, the eddy current density increases in the metal region. With stronger the induction intensity near the excitation coil, the smaller the eddy current density in the metal-free region. Therefore, the temperature characterization of the thermal imaging can be obtained by the eddy current method. It can effectively identify the location of metal in wood.

In this section, in order to evaluate the proposed algorithm, principal component analysis (PCA) and independent component analysis (ICA) algorithms have been selected for comparison. As shown in Figures 4 and 5, the temperature rise of the metal area is not obviously displayed in the image with PCA, and the position and size of the metal are blurred. Despite this, these algorithms can extract features effectively in detecting metals from ECPT system which have more obvious metal characteristics. The metals in the thermal images are relatively easy for human to discern.

After TSR, the results are shown in Figure 6. The temperature rise of metal in wood materials changes more obviously, and the position and size of metal foreign bodies are clearer.

TABLE 1: Evaluation table for proposed method.

| Case | PCA | ICA | TSR |
|---|---|---|---|
| L1 | 84.34% | 86.42% | 86.70% |
| L2 | 84.88% | 86.01% | 86.98% |

The temperature difference between metal and nonmetal decreases from inside to outside, the temperature rise in metal area increases obviously, and the surrounding temperature is decreasing, forming obvious temperature difference. As shown in TSR, that number of metal in L1 is 4 metals, the number of metal in L2 is 5 metals, and the area size represents the size of the metals. Experimental results shows that the location, size, and quantity of metal can be clearly identified.

The proposed method is compared with two state-of-the-art methods by using two samples. The evaluation metrics concern both efficiency (inference time) and effectiveness. The results are the mean of five different infrared thermal datasets. The same platform has been used to run them, and the results are given in Table 1. From Table 1, conclusions can draw that all algorithms can have certain improvements

on metal detection, especially for the size detection. The results are achieving 2.34% and 2.1% gains.

## 5. Conclusions

At the present time, stress wave systems can only be used for qualitative testing. In this paper, eddy current thermography is used to detect metal in wood materials, and the detection principle and thermal signal reconstruction technology (TSR) are analyzed in detail. The conclusions are as follows:

(1) It can accurately detect the presence or absence of metal in wood and other materials and determine the quantity and size of metal

(2) Compared with other nondestructive testing, the effectiveness is reflected in the fact that there is no lift-off effect, the heating is rapid, the rapid detection is convenient, the detection area is large, the sensitivity is high, the use is convenient, and the influence of the shape and structure of the detected object is small

(3) High efficiency is reflected in the ability to accurately determine the position of metal in wood, obtain the size of metal, and greatly improve the production and processing efficiency and detection

However, the main limitation of proposed method is that the overall algorithm is complicated, and the amount of data is large, which requires more calculation and time. The subsequent algorithm and workflow need to be simplified to a certain extent to reduce the amount of calculation. With further research, this problem will be solved in the near future.

## Data Availability

The datasets, codes, and weight files 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

## References

[1] S. Petrovi, "The role of China in sustainable market supply of the EU with wood windows and doors. Sustainability of Forest Based Industries in the Global Economy," *Proceedings of Scientific Papers*, 2020.

[2] A. Dr, L. Kaukale, L. Tupenaite, T. Geys, and A. Knuutila, "Sustainable Public Buildings Designed and Constructed in Wood," in *Sustainable Public Buildings Designed and Constructed in Wood: Handbook*, pp. 10–23, RTU Press, 2021.

[3] L. Zhang, A. Tiemann, T. Zhang et al., "Nondestructive assessment of cross-laminated timber using non-contact transverse vibration and ultrasonic testing," *European Journal of Wood and Wood Products*, vol. 79, no. 2, pp. 335–347, 2021.

[4] X. Yang, Y. Ishimaru, I. Iida, and H. Urakami, "Application of modal analysis by transfer function to nondestructive testing of wood I: determination of localized defects in wood by the shape of the flexural vibration wave," *Journal of Wood Science*, vol. 48, no. 4, pp. 283–288, 2002.

[5] J. D. Stewart and C. S. Mvolo, "Comparison between static modulus of elasticity, non-destructive testing moduli of elasticity and stress-wave speed in white spruce and lodgepole pine wood," *Wood Material Science & Engineering*, pp. 1–11, 2021.

[6] Q. Yi, H. Malekmohammadi, G. Y. Tian, S. Laureti, and M. Ricci, "Quantitative evaluation of crack depths on thin aluminum plate using eddy current pulse-compression thermography," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 6, pp. 3963–3973, 2020.

[7] G. Xingwang, G. Heqing, L. Yingtao, and T. Jia, "Spectrum characteristics and light source selection for infrared thermal imaging testing of semitransparent materials," *Infrared and Laser Engineering*, vol. 46, no. 1, article 104001, 2017.

[8] Y. Peng, S. Huang, Y. He, and X. Guo, "Eddy current pulsed thermography for noncontact nondestructive inspection of motor winding defects," *IEEE Sensors Journal*, vol. 20, no. 5, pp. 2625–2634, 2020.

[9] B. Du, Y. He, and C. Zhang, "Progress and trends in fault diagnosis for renewable and sustainable energy system based on infrared thermography: A review," *Infrared Physics & Technology*, vol. 109, p. 103383, 2020.

[10] Y. Peng, S. Huang, B. Deng et al., "Joint scanning electromagnetic thermography for industrial motor winding defect inspection and quantitative evaluation," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 5, pp. 6832–6841, 2021.

[11] N. Rajic, "Principal component thermography for flaw contrast enhancement and flaw depth characterisation in composite structures," *Composite Structures*, vol. 58, no. 4, pp. 521–528, 2002.

[12] T. Liang, W. Ren, G. Y. Tian, M. Elradi, and Y. Gao, "Low energy impact damage detection in CFRP using eddy current pulsed thermography," *Composite Structures*, vol. 143, pp. 352–361, 2016.

[13] L. Sripragash and M. J. J. N. Sundaresan, "A normalization procedure for pulse thermographic nondestructive evaluation," *NDT & E International*, vol. 83, pp. 14–23, 2016.

[14] A. Hyvarinen and E. J. N. N. Oja, "Independent component analysis: algorithms and applications," *Neural Networks*, vol. 13, no. 4-5, pp. 411–430, 2000.

[15] R. Świta and Z. Suszyński, "Cluster segmentation of thermal image sequences using kd-tree structure," *International Journal of Thermophysics*, vol. 35, no. 12, pp. 2374–2387, 2014.

[16] X. Maldague and S. Marinetti, "Pulse phase infrared thermography," *Journal of Applied Physics*, vol. 79, no. 5, pp. 2694–2698, 1996.

[17] Y. He, B. Deng, H. Wang et al., "Infrared machine vision and infrared thermography with deep learning: a review," *Infrared Physics & Technology*, vol. 116, article 103754, 2021.

[18] R. Yang and Y. He, "Optically and non-optically excited thermography for composites: a review," *Infrared Physics & Technology*, vol. 75, pp. 26–50, 2016.

*Research Article*

# Computer Vision-Based Detection for Delayed Fracture of Bolts in Steel Bridges

**Jing Zhou** [ID] **and Linsheng Huo** [ID]

*State Key Laboratory of Coastal and Offshore Engineering, Dalian University of Technology, Dalian 116024, China*

Correspondence should be addressed to Linsheng Huo; lshuo@dlut.edu.cn

The delayed fracture of high-strength bolts occurs frequently in the bolt connections of long-span steel bridges. This phenomenon can threaten the safety of structures and even lead to serious accidents in certain cases. However, the manual inspection commonly used in engineering to detect the fractured bolts is time-consuming and inconvenient. Therefore, a computer vision-based inspection approach is proposed in this paper to rapidly and automatically detect the fractured bolts. The proposed approach is realized by a convolutional neural network- (CNN-) based deep learning algorithm, the third version of You Only Look Once (YOLOv3). A challenge for the detector training using YOLOv3 is that only limited amounts of images of the fractured bolts are available in practice. To address this challenge, five data augmentation methods are introduced to produce more labeled images, including brightness transformation, Gaussian blur, flipping, perspective transformation, and scaling. Six YOLOv3 neural networks are trained using six different augmented training sets, and then, the performance of each detector is tested on the same testing set to compare the effectiveness of different augmentation methods. The highest average precision (AP) of the trained detectors is 89.14% when the intersection over union (IOU) threshold is set to 0.5. The practicality and robustness of the proposed method are further demonstrated on images that were never used in the training and testing of the detector. The results demonstrate that the proposed method can quickly and automatically detect the delayed fracture of high-strength bolts.

## 1. Introduction

High-strength bolt connections are widely used to assemble the load-bearing members of the steel structure in long-span steel bridges. The popularity of the bolt connections is attributed to their low cost, high reliability, and rapid assembly [1]. However, these bridges are often operated in adverse environments and subject to corrosion [2, 3], vibration and fatigue [4, 5], and thermal cycling, which can contribute to the damage of bolts. The damage types of bolts that occur the most include the looseness and delayed fracture. The delayed fracture of bolts refers to the sudden fracture of bolts under constant tension [6]. Due to the huge energy released by the brittle fracture, the fractured bolts will be missing. The damage of bolts will threaten the safety of the bridges and may even lead to severe accidents. Hence, it is necessary to monitor the bolt condition in the daily operation and maintenance phase.

Over the decades, structural health monitoring methods have attracted lots of attention [7–10], and they have been applied to detect the bolt damage [11, 12]. They mainly rely on the sensor technology to identify the variations of the preload, including piezoelectric active sensing methods [13, 14], the electromechanical impedance methods [15, 16], and the vibroacoustic modulation methods [17, 18]. A "smart washer" with a piezoceramic patch sandwiched between two flat metal rings was developed to monitor the bolted joint through the active sensing method [19]. Further, the fluctuation of the impedance signatures in frequency was utilized to evaluate the bolted joint with the developed "smart washer" [20]. A novel vibroacoustic modulation method was proposed to monitor the early looseness of a bolt in real time [21]. Notably, although the contact sensor-based methods are proposed to detect the decrease of preload induced by initial bolt looseness, they can also be used to detect the delayed fracture of bolts,

which is a kind of brittle damage and can result in the disappearance of the preload [22, 23]. Nonetheless, the contact sensor-based methods face the challenge of the cost scaling up when monitoring multiple bolts, because one sensor can only perform the measurement at one bolt. As a result, most bridges are impossible to equip with enough sensors, and the current monitoring method in practice highly relies on manual inspection. However, the whole inspection process is very dangerous and inefficient. As shown in Figure 1, maintenance workers are inspecting the delayed fracture of high-strength bolts on a long-span steel bridge.

In recent years, computer vision technology has received substantial attention as an interdisciplinary subject, and it has been applied in civil infrastructure inspection and monitoring to improve the accuracy and efficiency of manual vision inspection [24, 25]. It has been applied to detect bolt damage because there are always a huge number of bolts in actual steel structures. Park et al. [26] proposed a vision-based method to evaluate the rotation angle of a bolt nut. Cha et al. [27] utilized the image-processing techniques and the support vector machine to detect the bolt looseness. However, traditional computer vision-based methods have poor robustness and low accuracy. On the other hand, the convolutional neural network (CNN) has gained great success in computer vision [28] with the development of deep learning technology, and CNN-based algorithms have achieved the most advanced performance in various tasks, including image classification [29], object detection [30], and semantic segmentation [31]. This kind of technology has also been applied for bolt damage detection. Huynh et al. [32] proposed a quasiautonomous bolt looseness detection method, where the plausible bolts were detected using a CNN-based object detector and the rotation angle of each bolt was measured by the Hough line transform. Zhao et al. [33] proposed a method for the measurement of the bolt-loosening rotation angle using a CNN-based object detector. Wang et al. [34] designed a computer vision-based method by integrating the perspective transformation to detect the bolt looseness for flange connections. However, most studies have only focused on the detection of bolt looseness, and there is no research on the inspection of the delayed fractures in high-strength bolts, to the authors' best knowledge. The visual characteristics of the delayed fracture of high-strength bolts are totally different from looseness, because the fractured bolts will be missing due to the tremendous amount of energy released by the fracture [22, 23]. Notably, bolt delayed fracture can be more dangerous than bolt looseness in theory, because the former will cause the vanishing of the preload, whereas the latter will only reduce the force. Hence, this paper proposed a computer vision-based inspection method for the delayed fracture of bolts, where the damage was detected and located in an automated fashion using an object detection algorithm.

The task of the object detection is to classify and locate the targets in the image, and various algorithms have been developed with high recognition accuracy. The CNN-based object detection methods can be divided into region-based and region-free classifications according to the differences in the idea of the algorithm. The region-based approaches, such as the region-based convolutional neural network (R-CNN) [35], Fast R-CNN [36], and Faster R-CNN [37], combine



FIGURE 1: Maintenance workers are inspecting the bolt delayed fracture in a long-span steel bridge.

region proposals and CNN to detect objects. The region proposals are produced from the input image, and they are treated as the set of candidate detections. The region-free methods, such as Single Shot MultiBox Detector (SSD) [38], You Only Look Once (YOLO) [39], YOLOv2 [40], YOLOv3 [41], and YOLOv4 [42], frame the object detection task as a regression problem, and these methods directly detect objects from the input image by using CNN. The speed of region-based methods is slower than region-free methods due to the necessity of region proposals. Hence, the region-free methods were selected in our research for the real-time detection. In addition, YOLOv3 boasts improved performance for detecting small objects in wide-scale images [43, 44]. The size of the delayed fracture of bolts is relatively small in an image of a bolt connection. Therefore, YOLOv3 is selected to detect the delayed fracture of bolts.

On the other hand, the performance of the CNN-based object detector heavily relies on extracting information from abundant labeled images, and the performance can be improved with the increase of training data in amount and diversity. However, it is quite difficult to acquire enough labeled images in practice, and then, the performance of the trained detector is always limited to some extent. For the bolt damage detection task in long-span steel bridges, images are difficult to be captured due to the environmental complexity and limitation (such as the positions of fractured bolts in a long-span bridge are inaccessible in most cases), and the manual labeling of the images is laborious due to the concentration of bolts.

Data augmentation is one of the most commonly used methods to alleviate this problem, and it can automatically enlarge the dataset by utilizing the existing images [45, 46]. In recent years, many data augmentation methods have been developed for object detection, and images are augmented by many kinds of image-processing technologies. The widely used technologies include brightness transformation, flipping, noise addition, and perspective transformation [47]. For example, Fast R-CNN and Faster R-CNN use horizontal flipping to augment training data [36, 37]. The perspective transformation was introduced to enlarge the training dataset for transmission-line object detection [48]. Although many augmentation techniques are available, the selection of the techniques is task-specific and primarily depends on the experience. Thus, the augmentation effects are still unclear for each method in the

detector training of fractured bolts, and it is necessary to study the effectiveness of different data augmentation methods.

In this paper, a computer vision-based inspection method is developed to automatically detect and locate the bolt delayed fracture, and a series of data augmentation methods are utilized to improve the performance of the detector without external laboriousness. In addition, the impact of different data augmentation methods on the performance of the detector is analyzed.

## 2. Methodology

### 2.1. Workflow of the Detection Method for the Bolt Delayed Fracture.
As shown in Figure 2, the whole process involves three steps, including dataset preparation, detector training, and damage detection. During dataset preparation, many raw images of high-strength bolt connections are first collected through a camera device. Then, the labeled images can be obtained by artificial labeling and data augmentation, where the damage is labeled by enclosing rectangle bounding boxes. The pairwise images and labels are used to train the YOLOv3 neural network until it can pass the performance checking. Finally, the trained neural network can be used as a damage detector to perform damage detection in the real world.

### 2.2. Overview of YOLOv3 Detector.
YOLOv3 is evolved from its predecessors: YOLO and YOLOv2, which mainly improves the detection accuracy, especially for the detection of small targets. Specifically, a new network, Darknet53, integrating residual networks and Darknet19 (the network used in YOLOv2) was introduced to improve the ability of feature extraction, and the multiscale prediction is used to help simultaneously obtain semantic information and fine-grained information from different feature maps. The architecture of YOLOv3 is shown in Figure 3.

At the beginning of the training process, the image-label pairs are fed into the neural network. Each input image is adjusted to a fixed size, and then, it is divided into $S \times S$ grids using upsampling and feature fusion operations. Each grid is tasked with detecting objects that have their center coordinates enclosed by the grid. Each grid outputs $b$ bounding boxes and $c$ conditional category probability. Each bounding box can be determined by the coordinate information ($x$, $y$, $w$, and $h$) and the confidence score ($S_c$). The coordinates ($x$, $y$) point towards the center of the bounding box. The parameters $w$ and $h$ are, respectively, the width and height of the bounding box. $S_c$ can be obtained according to Equation (1). The loss function value is calculated using the prediction value and label value. The adjustable parameters in the neural network are updated using a backpropagation algorithm. The process is repeated until the loss function value converges at a small value. During the inference process, only the image is fed into the trained neural network, and the prediction of the neural network is regarded as the detection result.

$$S_c = P(\text{Object}) \times \text{IOU}_{\text{pred}}^{\text{truth}}, \quad P(\text{Object}) \in \{0, 1\}, \quad (1)$$

where $P(\text{Object})$ is equal to 0 when no object exists in the grid; otherwise, its value is 1. $\text{IOU}_{\text{pred}}^{\text{truth}}$ is the intersection over union

(IOU) between the predicted bounding box and the ground truth of the object.

The loss function in YOLOv3 consists of three parts: coordinate loss, IOU loss, and classification loss. All of them correspond to the output of the neural network prediction. However, the classification loss is removed in this paper, because the number of classifications is only one. The loss function used in this paper is shown in the following equation:

$$\begin{aligned}
\text{loss} = {}& \lambda_{\text{coord}} \sum_{i=1}^{S \times S} \sum_{j=1}^{B} I_{ij}^{\text{obj}} (2 - w_i \times h_i) \left[ \left( x_i - \widehat{x}_i \right)^2 + \left( y_i + \widehat{y}_i \right)^2 \right] \\
& + \lambda_{\text{coord}} \sum_{i=1}^{S \times S} \sum_{j=1}^{B} I_{ij}^{\text{obj}} (2 - w_i \times h_i) \left[ \left( w_i - \widehat{w}_i \right)^2 + \left( h_i + \widehat{h}_i \right)^2 \right] \\
& - \sum_{i=1}^{S \times S} \sum_{j=1}^{B} I_{ij}^{\text{obj}} \left[ \widehat{S}_{ci} \log (S_{ci}) + \left( 1 - \widehat{S}_{ci} \right) \log (1 - S_{ci}) \right] \\
& - \lambda_{\text{noobj}} \sum_{i=1}^{S \times S} \sum_{j=1}^{B} I_{ij}^{\text{obj}} \left[ \widehat{S}_{ci} \log (S_{ci}) + \left( 1 - \widehat{S}_{ci} \right) \log (1 - S_{ci}) \right],
\end{aligned}$$

$$(2)$$

where $\lambda_{\text{noobj}}$ and $\lambda_{\text{coord}}$ are the efficiencies of the IOU loss and coordinate loss, respectively; $\widehat{x}_i, \widehat{y}_i, \widehat{w}_i, \widehat{h}_i, \widehat{S}_{ci}$ are the ground truth values. The value of $I_{ij}^{\text{obj}}$ is 1, when the target falls into the $j_{\text{th}}$ bounding box of the $i_{\text{th}}$ grid; otherwise, it is equal to 0.

In addition, the average precision (AP) is used as an indicator to estimate the performance of the damage detector. The AP sums up the precision-recall curve by computing the area under the curve [49]. The precision ($P$) and recall ($R$) are defined as follows:

$$\begin{aligned}
P &= \frac{\text{TP}}{\text{TP} + \text{FP}}, \\
R &= \frac{\text{TP}}{\text{TP} + \text{FN}},
\end{aligned} \quad (3)$$

where true positives (TP) indicate the number of fractured bolts correctly detected by the detector. False positives (FP) point to the number of other objects in the background falsely detected as fractured bolts. False negatives (FN) refer to the number of fractured bolts missed by the detector.

### 2.3. Data Augmentation Methods to Improve the Detector Performance.
To improve the performance of the detector, five data augmentation methods were introduced in this paper, including brightness transformation (BT), Gaussian blur (GB), flipping (FL), perspective transformation (PT), and scaling (SC). The data augmentation methods are selected considering the practical change of the captured images in engineering and the label preserving ability after augmentation. The BT can mimic images taken under different light intensity conditions. The GB can simulate vague images taken under some unfavorable situations, such as long-distance, slightly out of focus, and foggy weather. The PT can imitate images taken from different viewpoints, including the positions that the camera device cannot reach. The FL can further produce new images captured from different viewpoints. The SC can simulate the image

FIGURE 2: Flowchart of the proposed bolt delayed fracture-detection method.



FIGURE 3: The architecture of YOLOv3 neural network.

resolution changes. The sample images after data augmentation are shown in Figure 4, where the labels of the image are represented by green rectangular bounding boxes.

The data augmentation methods take an image and its label as input and automatically generate a new augmented image and corresponding label. As shown in Figure 4, the BT and GB do not change the coordinates of the bounding box, whereas the PT, FL, and SC can result in the coordinate change of the bounding box. Hence, the bounding boxes after BT and GB are the same as the original bounding boxes, and the bounding boxes after FL, PT, and SC should be rectified.

The details of the data augmentation methods are described in the following text. The linear BT was used in

FIGURE 4: Sample labeled images with different augmentation methods: (a) original image, (b, c) brightness transformation, (d, e) Gaussian blur, (f, g) flipping, (h, i) perspective transformation, and (j, k) image scaling.

this paper, which directly multiplies the pixel value of the image by a certain coefficient. The GB takes the average value of pixels around a certain point as the pixel value of that point, and the surrounding pixels are assigned different weights according to the distance and the normal distribution. The FL consists of vertical flipping and horizontal flipping, and the rectified bounding box can be obtained easily according to the symmetry. The bilinear interpolation technique was used to change the image resolution, and the rectified image and bounding box can be obtained according to the scaling coefficient. The PT transforms an image from one plane into another plane by a perspective transformation matrix, as shown in Equation (4). The components of the matrix can be obtained according to four pairs of points following Equation (5). As shown in Figure 5, four vertices $(A_0, B_0, C_0, \text{and } D_0)$ of the input image and four random

sampling points $(A_0^*, B_0^*, C_0^*, \text{and } D_0^*)$ of the augmented image were used to calculate the perspective transformation matrix, and the coordinates of each point can be obtained following Equations (6) and (7). After obtaining the perspective transformation matrix, a rectangular bounding box $(A_1B_1C_1D_1)$ in the input image can be transformed to a quadrangle bounding box $(A_1^*B_1^*C_1^*D_1^*)$ in the augmented image following Equation (4). However, the nonrectangular bounding boxes cannot be trained by CNN. To tackle this problem, the nonrectangular bounding boxes are rectified using label alignment to generate the new rectangular bounding boxes automatically. To make sure the generated rectangular bounding box totally contain the bolt damage, we let $x_{1PT} = x_{2PT} = \min \{x_1^*, x_2^*, x_3^*, x_4^*\}$, $y_{1PT} = y_{3PT} = \min \{y_1^*, y_2^*, y_3^*, y_4^*\}$, $x_{3PT} = x_{4PT} = \max \{x_1^*, x_2^*, x_3^*, x_4^*\}$, and $y_{2PT} = y_{4PT} = \max \{y_1^*, y_2^*, y_3^*, y_4^*\}$. The

Figure 5: Flowchart of the perspective transformation data augmentation.

generated rectangular bounding box can be represented by $A_{1PT} = (x_{1PT}, y_{1PT})$, $B_{1PT} = (x_{2PT}, y_{2PT})$, $C_{1PT} = (x_{3PT}, y_{3PT})$, and $D_{1PT} = (x_{4PT}, y_{4PT})$. The augmented image and the original image have the same size, and the blank area in the augmented image is filled with black pixels.

$$u_i = \mathbf{T}v_i = \begin{bmatrix} t_{11} & t_{12} & t_{13} \\ t_{21} & t_{22} & t_{23} \\ t_{31} & t_{32} & 1 \end{bmatrix} v_i, \tag{4}$$

$$\begin{bmatrix} x_1 & y_1 & 1 & 0 & 0 & 0 & -x_1 x_1^* & -y_1 x_1^* \\ x_2 & y_2 & 1 & 0 & 0 & 0 & -x_2 x_2^* & -y_2 x_2^* \\ x_3 & y_3 & 1 & 0 & 0 & 0 & -x_3 x_3^* & -y_3 x_3^* \\ x_4 & y_4 & 1 & 0 & 0 & 0 & -x_4 x_4^* & -y_4 x_4^* \\ 0 & 0 & 0 & x_1 & y_1 & 1 & -x_1 y_1^* & -y_1 y_1^* \\ 0 & 0 & 0 & x_2 & y_2 & 1 & -x_2 y_2^* & -y_2 y_2^* \\ 0 & 0 & 0 & x_3 & y_3 & 1 & -x_3 y_3^* & -y_3 y_3^* \\ 0 & 0 & 0 & x_4 & y_4 & 1 & -x_4 y_4^* & -y_4 y_4^* \end{bmatrix} \begin{bmatrix} t_{11} \\ t_{12} \\ t_{13} \\ t_{21} \\ t_{22} \\ t_{23} \\ t_{31} \\ t_{32} \end{bmatrix} = \begin{bmatrix} x_1^* \\ x_2^* \\ x_3^* \\ x_4^* \\ y_1^* \\ y_2^* \\ y_3^* \\ y_4^* \end{bmatrix}, \tag{5}$$

where $\mathbf{T}$ is the projective transformation matrix, $v_i = (x_i, y_i, 1)$ is a point in the original image, and $u_i = (x_i^*, y_i^*, 1)$ is a point in the PT augmented image.

$$x_{tl}, x_{bl}, x_{br}, x_{tr} \in (0, W \times \lambda), \tag{6}$$

$$y_{tl}, y_{bl}, y_{br}, y_{tr} \in (0, H \times \lambda), \tag{7}$$

where $W$ and $H$ represent the width and height of the image; $x_{tl}$, $x_{bl}$, $x_{br}$, $x_{tr}$, $y_{tl}$, $y_{bl}$, $y_{br}$, and $y_{tr}$ are the distances from points to the corresponding image boundary; and $\lambda$ is the intensity parameter for perspective transformation; and the greater the $\lambda$, the more obvious the perspective phenomenon.

## 3. Experimental Verification

*3.1. Dataset Preparation.* Due to the practical challenges of obtaining large amounts of images depicting the bolt delayed fracture in real bridges, only two images of delayed fractured bolts from an actual suspension bridge were collected in this study. It is impossible to train the YOLOv3 neural network with such a limited number of images; however, these two images can be used to demonstrate the practicability of the proposed method. Thus, training images were generated using two steel plates held together with high-strength bolts. Many images of fractured bolts were collected from the fabricated steel plates to train the neural network.

In this paper, a total of 500 raw images were collected at $3016 \times 3016$-pixel and $3016 \times 4032$-pixel resolutions by a smartphone camera from Xiaomi Mi 6. The distance between the object and the camera is approximately from 0.2 m to 1.5 m. To obtain different lighting intensities of a bolt image in an actual bridge, the images were collected outside during different times of the day (e.g., 9 a.m., 1 p.m., and 5 p.m.). The relationship between the camera's viewing direction and the direction of the sunlight illumination will also influence the brightness of the images. Hence, during the image collection, the conditions of front-lighting, back-lighting, and side-

TABLE 1: The number of images in the datasets.

| Dataset | $DA_{OR}$ | $DA_{BT}$ | $DA_{GB}$ | $DA_{PT}$ | $DA_{FL}$ | $DA_{SC}$ | Validation set | Testing set |
|---|---|---|---|---|---|---|---|---|
| Number | 320 | 960 | 960 | 960 | 960 | 960 | 80 | 100 |



FIGURE 6: The result of $k$-means clustering on the original training set.

TABLE 2: The average precision (AP) of the detectors using different training sets (%).

| Training set | $DA_{OR}$ | $DA_{BT}$ | $DA_{GB}$ | $DA_{PT}$ | $DA_{FL}$ | $DA_{SC}$ |
|---|---|---|---|---|---|---|
| $AP_{0.5}$ | 84.62 | 81.81 | 81.29 | 89.14 | 86.37 | 82.99 |
| Increment | — | -2.81 | -3.33 | 4.52 | 1.75 | -1.63 |
| $AP_{0.6}$ | 70.09 | 75.16 | 67.31 | 83.48 | 75.93 | 69.39 |
| Increment | — | 5.07 | -2.78 | 13.39 | 5.84 | -0.7 |
| $AP_{0.7}$ | 42.11 | 47.55 | 38.73 | 60.56 | 48.87 | 46.69 |
| Increment | — | 5.44 | -3.38 | 18.45 | 6.76 | 4.58 |

$DA_{BT}$, $DA_{GB}$, $DA_{FL}$, $DA_{PT}$, and $DA_{SC}$ are the training sets consisting of augmented images generated by the corresponding augmentation method and manually labeled images. The augmented images were produced before training for convenience.

Two BT coefficients were randomly selected from 0.6 to 1.4 and utilized to adjust the brightness for images in $DA_{OR}$, and as a result, 640 new images were generated. The range of the brightness transformation coefficient was determined based on whether the edge of the target can be identified using naked eyes. The original images in $DA_{OR}$ were also modified using GB to generate 640 additional images. The standard deviation for the Gaussian kernel was randomly selected between 0 and 3.0. The range of standard deviations was set in the same manner mentioned in BT. The images in $DA_{OR}$ were horizontally and vertically flipped, and 640 new flipping images were produced. The scaling coefficient was selected from 0.1 to 1.9, and 640 new augmented images were produced. The PT was applied twice, and 640 new augmented images were generated. The perspective intensity parameter $\lambda$ was selected from 0.1 to 0.3. The number of images in different data sets is shown in Table 1.

lighting are all considered. The direction of the camera viewing was set parallel, antiparallel, and perpendicular to the vector of the sunlight, which can, respectively, simulate front-lighting, back-lighting, and side-lighting. Since the shadow from clouds or the bridge structure will affect the detection accuracy, images were also gathered under scattered tree shade. The apparent shape of the bolt changes based on the viewing angle, and thus, images were taken from multiple viewing angles for the same fractured bolt during image acquisition.

After the image acquisition, the fractured bolts in all 500 images were manually labeled with bounding boxes using the custom code written in Python. And considering the convenience of using the dataset in the future, the dataset is converted to PASCAL VOC format [49]. A ".XML" file including the information of the labeled bounding boxes was generated for each image after successful labeling. The file was then converted into a ".txt" file suitable for the training. The labeled images were then randomly divided into three sets: training set, validation set, and testing set, with 320, 80, and 100 images for each set, respectively. During the labeling process, a total of 439 objects were annotated in the 320 training images. The training set was utilized to train the neural network, and the validation set was used to aid the training and avoid overfitting. After training, the performance of the trained detector was estimated with the testing set. Notably, five extra training sets were generated using five data augmentation methods based on the original training set, and finally, a total of six training sets ($DA_{OR}$, $DA_{BT}$, $DA_{GB}$, $DA_{FL}$, $DA_{PT}$, and $DA_{SC}$) were established in this research and used to train six neural networks. $DA_{OR}$ is the original training set with 320 manually labeled images.

3.2. Implementation Details during Training Process. All experiments were performed on a personal computer: Lenovo R720 (a Core i7-7700HQ CPU @ 2.80 GHz, 8 GB DDR4 memory, and 2 GB memory NVIDIA GeForce GTX 1050 Ti GPU). All the training and testing processes were conducted on the GPU. The YOLOv3 neural network was developed using Python 3.6.5 under TensorFlow 1.8.0 frame.

Before the beginning of experiments, the $k$-means clustering algorithm was applied on the size of the bounding boxes of images in $DA_{OR}$ to obtain the bounding box priors and facilitate network learning and detection results. The clustering results are shown in Figure 6. The number of clusters is set at 9 as follows: $(15 \times 11)$, $(13 \times 14)$, $(20 \times 17)$, $(22 \times 21)$, $(25 \times 26)$, $(35 \times 31)$, $(45 \times 45)$, $(74 \times 74)$, and $(131 \times 160)$.

In order to improve the detection accuracy of the detector and adapt to the required input format of the Darknet53, the size of the input image is set to $416 \times 416$ pixels. Due to

(a)

(b)

(c)

(d)

FIGURE 7: Sample images of detected bolt delayed fracture.

the constraints of the GPU memory, the batch size was set to 8. And the training step was set to 10000 to analyze the training process by the loss curve. To prevent training the neural network from scratch, the internal adjustable parameters was initialized by using a pretrained weight, which can be obtained from the website (https://pjreddie.com/yolo/). The initial learning rate was set to 0.003 through trial and error with the help of the validation set. The $\lambda_{coord}$ and $\lambda_{noobj}$ are set to 5 and 0.5, respectively. To analyze the impact of different data augmentation methods, six neural networks were trained using six training sets. The parameter settings during training process are the same except using the different training sets.

## 4. Result and Discussion

After training, the images in the testing set were used to test the performance of the six detectors, and the IOU threshold is set to 0.5, 0.6, and 0.7. The AP values of six detectors under different IOU thresholds were calculated, as shown in Table 2. The highest AP value is 89.14%, which indicates that the trained detector has a strong generalization and excellent detection performance, and the AP value decreases with the increase of the IOU threshold. The AP of the detector trained using $DA_{OR}$ is used as a benchmark, and the AP increment of other detectors is used to estimate the usefulness of different methods. The PT and FL both improve the AP value on the testing set, and the highest increment of AP is induced by PT, achieving 4.52%, 13.39%, and 18.45% corresponding to three different IOU thresholds. In theory, five data augmentation methods all can improve the richness of the training set, and the performance of the detectors trained by augmented training sets should be better than the detector trained by the original training set. However, the BT, GB, and SC reduce the performance, as shown in Table 2. The reason can be that although the change of lighting intensity, distance, and resolution was considered during image collection, the number of the collected raw images in the testing set is too small to represent the entire image sample. Hence, the promotion of the ability of the detector to detect vague images, brightness changes, and resolution changes cannot be reflected on the existing testing set, whereas the improvement of the ability to detect objects captured from different viewpoints is the most obvious, because all images were captured from different viewpoints.

The detection results of some images in the testing set are shown in Figure 7. The fractured bolts in the image were

Figure 8: Sample images of detected bolt delayed fracture considering the color and weather.

(a)                                                                                  (b)

FIGURE 9: Detected bolt delayed fracture on an actual bridge.

automatically detected (indicated with a solid red box) by the detector after the test images were inputted into the best detector. The detection process spends only 0.06 seconds for each input image ($416 \times 416$ resolution). It should be noted that the detection speed is affected by hardware limitations. In [41], the detection speed for a $416 \times 416$ resolution image is 0.029 seconds. Hence, this method can accomplish real-time autonomous damage detection when a camera is used in conjunction with a processor. The proposed method can facilitate the transition from manual inspection to automated inspection or monitoring carried out by fixed cameras, UAVs, or remote-controlled robots in the future.

The generalization ability of the trained detector was further demonstrated using the new images of bolts with different colors (black, red, gray, and blue) and covered with raindrops and the images of fractured bolts from an actual bridge (two $3024 \times 4032$-pixel images taken from a real long-span steel bridge in China). The detection results are shown in Figures 8 and 9, and the trained detector can correctly detect the damage from the new images. The results show that the trained detector does not overfit the two sample steel plates. It also demonstrates the practicality of the proposed method.

On the other hand, although the detector can detect the damage correctly, the predicted bounding boxes do not perfectly fit the fractured bolts. The minor errors may be induced by the limitation of the training set, such as the lack of images taken from actual bridges. In addition, a comprehensive analysis of the effectiveness of different augmentation techniques for this detection task needs a comprehensive image dataset. The images in the dataset should be collected from actual engineering. Thus, more actual images need to be collected and a larger image dataset will be established in the future to further analyze the effectiveness of different augmentation methods and how to use them in combination.

## 5. Conclusion

This paper presents a new, automated method to inspect fracture failures for bolts. The method is developed based upon the CNN-based object detection algorithm YOLOv3,

and the performance of the detector is improved by data augmentation. An image dataset was developed through image acquisition, image labeling, and data augmentation, and six YOLOv3 neural networks were trained using different augmented training sets to analyze the impact of different augmentation methods. The highest AP of the trained detectors is 89.14% when the IOU threshold equals to 0.5. The effectiveness of different data augmentation methods is evaluated by the increment of AP. The highest increment of AP on the testing set is achieved by perspective transformation augmentation. The detection speed of the trained detector achieved 0.06 seconds for each input image with $416 \times 416$ resolution. The generalization of the trained network and the practicality of the proposed method were validated using new images that were never used in the training and testing. The proposed method has the potential to enable safe, real-time, and autonomous detection of delayed fracture of high-strength bolts with high accuracy.

## Data Availability

The datasets, codes, and weight files 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

# References

[1] D. Pines and A. E. Aktan, "Status of structural health monitoring of long-span bridges in the United States," *Progress in Structural Engineering & Materials*, vol. 4, no. 4, pp. 372–380, 2002.

[2] J. Peng, L. Xiao, J. Zhang, C. S. Cai, and L. Wang, "Flexural behavior of corroded HPS beams," *Engineering Structures*, vol. 195, pp. 274–287, 2019.

[3] L. Xiao, J. Peng, J. Zhang, Y. Ma, and C. S. Cai, "Comparative assessment of mechanical properties of HPS between electrochemical corrosion and spray corrosion," *Construction and Building Materials*, vol. 237, article 117735, 2020.

[4] W. Wang, X. Wang, X. Hua, G. Song, and Z. Chen, "Vibration control of vortex-induced vibrations of a bridge deck by a single- side pounding tuned mass damper," *Engineering Structures*, vol. 173, pp. 61–75, 2018.

[5] X. Yin, G. Song, and Y. Liu, "Vibration suppression of wind/-traffic/bridge coupled system using multiple pounding tuned mass dampers (MPTMD)," *Sensors*, vol. 19, no. 5, p. 1133, 2019.

[6] J. Zhou, L. Huo, G. Song, and H. Li, "Deep learning-based visual inspection for the delayed brittle fracture of high-strength bolts in long-span steel bridges," in *International Conference on Image and Video Processing, and Artificial Intelligence*, p. 1132129, Shanghai, China, 2019.

[7] X. Ye, C. Dong, and T. Liu, "A review of machine vision-based structural health monitoring: methodologies and applications," *Journal of Sensors*, vol. 2016, no. 5, Article ID 7103039, 2016.

[8] Y. Tan and L. M. Zhang, "Computational methodologies for optimal sensor placement in structural health monitoring: a review," *Structural Health Monitoring-an International Journal*, vol. 19, no. 4, pp. 1287–1308, 2020.

[9] M. Abdulkarem, K. Samsudin, F. Z. Rokhani, and M. F. A Rasid, "Wireless sensor network for structural health monitoring: a contemporary review of technologies, challenges, and future direction," *Structural Health Monitoring-an International Journal*, vol. 19, no. 3, pp. 693–735, 2020.

[10] L. Huo, H. Cheng, Q. Kong, and X. Chen, "Bond-slip monitoring of concrete structures using smart sensors—a review," *Sensors*, vol. 19, no. 5, p. 1231, 2019.

[11] T. Wang, G. B. Song, S. P. Liu, Y. R. Li, and H. Xiao, "Review of bolted connection monitoring," *International Journal of Distributed Sensor Networks*, vol. 9, no. 12, 2013.

[12] Z. Zhao, P. Chen, E. Zhang, and G. Lu, "Health monitoring of bolt looseness in timber structures using PZT-enabled time-reversal method," *Journal of Sensors*, vol. 2019, 8 pages, 2019.

[13] T. Jiang, Q. Wu, L. Wang, L. Huo, and G. Song, "Monitoring of bolt looseness-induced damage in steel truss arch structure using piezoceramic transducers," *IEEE Sensors Journal*, vol. 18, no. 16, pp. 6677–6685, 2018.

[14] F. Wang, L. S. Huo, and G. Song, "A piezoelectric active sensing method for quantitative monitoring of bolt loosening using energy dissipation caused by tangential damping based on the fractal contact theory," *Smart Material Structures*, vol. 27, no. 1, article 015023, 2017.

[15] D. Chen, L. Huo, and G. Song, "EMI based multi-bolt looseness detection using series/parallel multi-sensing technique," *Smart Structures and Systems*, vol. 25, no. 4, pp. 423–432, 2020.

[16] Y. Liang, Q. Feng, D. Li, and S. Cai, "Loosening monitoring of a threaded pipe connection using the electro-mechanical impedance technique—experimental and numerical studies," *Sensors*, vol. 18, no. 11, p. 3699, 2018.

[17] E. E. Ungar, "The status of engineering knowledge concerning the damping of built-up structures," *Journal of Sound and Vibration*, vol. 26, no. 1, pp. 141–154, 1973.

[18] A. Milanese, P. Marzocca, J. M. Nichols, M. Seaver, and S. T. Trickey, "Modeling and detection of joint loosening using output-only broad-band vibration data," *Structural Health Monitoring-an International Journal*, vol. 7, no. 4, pp. 309–328, 2008.

[19] L. Huo, D. Chen, Q. Kong, H. Li, and G. Song, "Smart washer—a piezoceramic-based transducer to monitor looseness of bolted connection," *Smart Materials and Structures*, vol. 26, no. 2, article 025033, 2017.

[20] L. S. Huo, D. D. Chen, Y. B. Liang, H. N. Li, X. Feng, and G. B. Song, "Impedance based bolt pre-load monitoring using piezoceramic smart washer," *Smart Materials and Structures*, vol. 26, no. 5, p. 057004, 2017.

[21] N. Zhao, L. Huo, and G. Song, "A nonlinear ultrasonic method for real-time bolt looseness monitoring using PZT transducer–enabled vibro-acoustic modulation," *Journal of Intelligent Material Systems and Structures*, vol. 31, no. 3, pp. 364–376, 2020.

[22] Y. Hagihara, "Evaluation of delayed fracture characteristics of high-strength bolt steels by CSRT," *ISIJ International*, vol. 52, no. 2, pp. 292–297, 2012.

[23] E. Akiyama, "Evaluation of delayed fracture property of high strength bolt steels," *ISIJ International*, vol. 52, no. 2, pp. 307–315, 2012.

[24] B. F. Spencer, V. Hoskere, and Y. Narazaki, "Advances in computer vision-based civil infrastructure inspection and monitoring," *Engineering*, vol. 5, no. 2, pp. 199–222, 2019.

[25] C. Z. Dong and F. N. Catbas, "A review of computer vision-based structural health monitoring at local and global levels," *Structural Health Monitoring-an International Journal*, vol. 20, no. 2, pp. 692–743, 2021.

[26] J. H. Park, T. C. Huynh, S. H. Choi, and J. T. Kim, "Vision-based technique for bolt-loosening detection in wind turbine tower," *Wind and Structures*, vol. 21, no. 6, pp. 709–726, 2015.

[27] Y. J. Cha, K. You, and W. Choi, "Vision-based detection of loosened bolts using the Hough transform and support vector machines," *Automation in Construction*, vol. 71, pp. 181–188, 2016.

[28] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[29] B. Du, Y. He, Y. He, J. Duan, and Y. Zhang, "Intelligent classification of silicon photovoltaic cell defects based on eddy current thermography and convolution neural network," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 10, pp. 6242–6251, 2020.

[30] Y. Xu, M. Fu, Q. Wang et al., "Gliding vertex on the horizontal bounding box for multi-oriented object detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 4, pp. 1452–1459, 2021.

[31] X. H. Yuan, J. F. Shi, and L. C. Gu, "A review of deep learning methods for semantic segmentation of remote sensing imagery," *Expert Systems with Applications*, vol. 169, article 114417, 2021.

[32] T.-C. Huynh, J.-H. Park, H.-J. Jung, and J.-T. Kim, "Quasi-autonomous bolt-loosening detection method using vision-based deep learning and image processing," *Automation in Construction*, vol. 105, article 102844, 2019.

[33] X. F. Zhao, Y. Zhang, and N. N. Wang, "Bolt loosening angle detection technology using deep learning," *Structural Control & Health Monitoring*, vol. 26, no. 1, article e2292, 2019.

[34] C. Y. Wang, N. Wang, S. C. Ho, X. M. Chen, and G. B. Song, "Design of a new vision-based method for the bolts looseness detection in flange connections," *IEEE Transactions on Industrial Electronics*, vol. 67, no. 2, pp. 1366–1375, 2020.

[35] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 580–587, Columbus, USA, 2014.

[36] R. Girshick, "Fast R-CNN," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1440–1448, Santiago, USA, 2015.

[37] S. Q. Ren, K. M. He, R. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137–1149, 2017.

[38] W. Liu, D. Anguelov, D. Erhan et al., "SSD: single shot multibox detector," in *European Conference on Computer Vision*, pp. 21–37, Amsterdam, Netherland, 2016.

[39] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: unified, real-time object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 779–788, Seattle, USA, 2016.

[40] J. Redmon and A. Farhadi, "YOLO9000: better, faster, stronger," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7263–7271, Honolulu, USA, 2017.

[41] J. Redmon and A. Farhadi, "YOLOv3: an incremental improvement," 2018, http://arxiv.org/abs/1804.02767.

[42] A. Bochkovskiy, W. Chien-Yao, and H. Y. M. Liao, "YOLOv4: optimal speed and accuracy of object detection arXiv," 2020, http://arxiv.org/abs/2004.10934.

[43] Y. N. Tian, G. D. Yang, Z. Wang, H. Wang, E. Li, and Z. Z. Liang, "Apple detection during different growth stages in orchards using the improved YOLO-V3 model," *Computers and Electronics in Agriculture*, vol. 157, pp. 417–426, 2019.

[44] S. S. Kumar, M. Z. Wang, D. M. Abraham, M. R. Jahanshahi, T. Iseley, and J. C. P. Cheng, "Deep learning-based automated detection of sewer defects in CCTV videos," *Journal of Computing in Civil Engineering*, vol. 34, no. 1, article 04019047, 2020.

[45] D. M. Han, Q. G. Liu, and W. G. Fan, "A new image classification method using CNN transfer learning and web data augmentation," *Expert Systems with Applications*, vol. 95, pp. 43–56, 2018.

[46] C. Shorten and T. M. Khoshgoftaar, "A survey on image data augmentation for deep learning," *Journal of Big Data*, vol. 6, no. 1, p. 48, 2019.

[47] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.

[48] K. Wang, B. Fang, J. Y. Qian, S. Yang, X. Zhou, and J. Zhou, "Perspective transformation data augmentation for object detection," *IEEE Access*, vol. 8, pp. 4935–4943, 2020.

[49] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The Pascal visual object classes (VOC) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010.

*Research Article*

# Study on Smart Home Energy Management System Based on Artificial Intelligence

**Yunlong Ma,[1] Xiao Chen,[1] Liming Wang,[2] and Jianlan Yang [iD][3]**

[1]*State Grid Jiangsu Electric Power Co., Ltd., Nanjing 210024, China*
[2]*Jiangsu Frontier Electric Power Technology, Co., Ltd., Nanjing 211102, China*
[3]*College of Electrical & Information Engineering, Hunan University, Changsha 410082, China*

Correspondence should be addressed to Jianlan Yang; yangjianlan@hnu.edu.cn

With the increase of household electricity consumption and the introduction of distributed new energy sources, more attention has been paid to the issue of optimizing the cost of electricity purchase for household customers. An effective way to deal with these problems is through home energy management system (HEMS). In this paper, a model of home energy management is presented to optimize the home energy mix. The operation of home electricity consumption devices, distributed generation systems, and energy storage devices, as well as the charging and discharging of electric vehicles, are all considered. HEMS is a self-regulating system that can accommodate fluctuations in tariffs and home electricity consumption. The structure and the optimal scheduling algorithm of HEMS are introduced. The smart grid and demand response, smart home, new energy generation, energy storage, and other related technologies are discussed. Furthermore, the optimal scheduling of power consumption devices and energy sources in the HEMS and future development directions are explained and analyzed. A framework of HEMS is presented on the basis of advanced metering infrastructure (AMI). The framework adopts a local information management terminal as the core of data storage and scheduling in the home. Based on the timely purchase of electricity from the grid and the generation of electricity in combination with PV systems, an optimized simulation model for the scheduling of a new home energy management system is established. In addition, the application prospects of artificial intelligence in the HEMS are overviewed.

## 1. Introduction

Home energy management system (HEMS) is an intelligent network control system based on smart grid, smart home, and smart meters [1–3]. It integrates power generation, electricity consumption, and energy storage devices into a single system for management and control [4–6]. HEMS can improve the efficiency of household renewable energy and save electricity bills for customers [7, 8]. The traditional power market lacks interaction with customers, and the electricity tariff form is single, resulting in the insufficient supply of electricity during peak hours, as well as wasted electricity in low hours. Subsequently, the peak and off-peak tariff mechanism is introduced, which plays a role in guiding customers to adjust the time of electricity consumption [9]. However, it is less flexible and cannot reflect the real rela-tionship between electricity consumption and supply. More-over, HEMS can fully interact with the power grid to obtain accurate real-time price, cooperate with generation and load forecasting, perform an intelligent allocation of household energy, optimize the allocation of household load in the time dimension, achieve demand response on the customer side, relieve the pressure on the grid during peak hours, and improve the stability of grid [10]. HEMS is the minimal unit of smart grid, which is a new generation of information technologies such as Internet of Things, cloud computing, mobile Internet, and big data, combined with the household as a carrier to achieve a low-carbon, healthy, intelligent, comfortable, and safe family lifestyle [11, 12]. By combining distributed power technologies such as household photovoltaic and energy storage, it flexibly controls various household appliances and realizes an intelligent mode of

electricity and energy use. Currently, HEMS has been a hot research topic, and its optimization objectives contain the aspects such as economy, comfort, and load shedding.

Extensive research has been conducted to describe household electricity behavior and establish an intelligent model of household electricity, aiming for maximum peak load shedding and minimum electricity cost [13]. Alternatively, some studies consider the correlation between the use of home appliances and the optimization of household electricity behavior with the goal of minimizing electricity bills and maximizing comfort [14]. In addition to a variety of household appliances, there are scholars who investigate the impact of electric vehicles and energy storage devices in the optimization of smart homes, in order to propose a method of household energy that considers real-time control strategies for energy storage devices [15, 16]. Although the above studies coordinate the consideration of smart home energy management with the charging and discharge strategies of energy storage devices, there are very few studies concerned with the rational allocation methods.

In this paper, the structure of HEMS is introduced and the optimal scheduling algorithm of HEMS is analyzed; smart grid and demand response, smart home, new energy generation, and energy storage technologies are discussed; and an analysis of the optimal scheduling of power consumption devices and energy in the HEMS is discussed. Furthermore, a framework for an advanced measurement infrastructure (AMI) is presented for HEMS. Based on the timely purchase of electricity from the grid and the generation of electricity from PV systems, an optimized simulation model for the scheduling of a new HEMS is developed. The prospects for the application of artificial intelligence in the HEMS are also discussed.

## 2. Operating Principle of HEMS

*2.1. Structure of HEMS.* HEMS is a system for the residential user side, which is based on technologies such as AMI, intelligent collection, and intelligent interaction. It is a household area network with smart devices like smart meters, smart sockets/switches, smart appliances and smart interactive terminals in the home [17]. Moreover, it can support the access of distributed energy, electric vehicles, and other devices and uses the local information management terminal as a bridge for comprehensive management of user information and information interaction with the main station, thus realizing the bidirectional interaction between grid and user, energy management, and other functions [18, 19].

The bidirectional smart metering terminal is responsible for acquiring electricity generation and consumption information of the household. The mobile terminal supplies the function of interacting with users, which is responsible for acquiring electricity consumption settings of users and displaying household electricity consumption information. As the verification and control device of the HEMS, the local information management terminal is capable of communicating with the bidirectional smart meter and the mobile terminal, acquiring the necessary

electricity and setting data, and integrating with the weather, demand response, and other information acquired from the external network to invoke the localized forecasting module and scheduling module to achieve intelligent control of household electricity consumption. Particularly, the scheduling module considers the impact of distributed generation and energy storage access in order to find the optimal control result.

*2.2. AMI Architecture.* AMI is an open bidirectional communication platform, which is used to connect the system and power load and collect and manage grid data through electricity metering technology to achieve smart usage [20]. It provides customers with time-phased or instantaneous metering values, which improves the efficiency of equipment usage and supports the grid. AMI consists of four main components: smart meters, communication networks, measurement data management systems (MDMS), and home area network. AMI architecture is given in Figure 1.

MDMS is based on the main station and works in conjunction with the AMI Automatic Data Collection System to acquire and store metered values. After getting the data, validation, editing, and estimation are conducted through MDMS. It can provide the processed data to the required systems and ensure that the data stream from other systems is accurate and complete under communication disruptions and customer-side failures. By using the data provided by the MDMS, the utility can implement peak and off-peak tariffs, time-of-use tariffs, and a number of other complex billing methods.

The intelligence of smart meters is embodied in their programmable capability. Except for metering, smart meters also have functions such as compound rate metering, event recording, data storage, and bidirectional communication. As the foundation of HEMS, it offers data support for home energy dispatch and customer demand-side response.

Bidirectional communication network is the bridge between the company and customer, which is responsible for reading the data of smart meter at regular intervals and sending the demand response information to the customer. PLC, RF, GPRS, and McWiLL are the common communication methods.

Home network is used to connect the intelligent control terminal, the intelligent power consumption equipment, and the intelligent electricity meter [21]. The intelligent control terminal can acquire all the information on electricity consumption and equipment status and send the results of electricity dispatch to the electricity equipment. Wireless communication is often used in the household. The common wireless communication methods are ZigBee, Wi-Fi, etc. ZigBee has greater advantages in power consumption, cost, and networking whereas Wi-Fi has relatively fast speed and can be directly connected to the Internet [22–24]. It has a wide range of applications in mobile networking devices.

*2.3. HEMS Topology.* The home distributed PV/energy storage power generation system can be divided into two types: DC topology and AC topology, as shown in Figures 2(a) and

FIGURE 1: Structure of AMI architecture.



FIGURE 2: Structure of home distributed PV/battery system: (a) DC Topology; (b) AC topology.

2(b), respectively. The system consists of PV equipment, energy storage equipment, grid-connected inverter, and load. In this system, the photovoltaic panels are measured by a separate meter. The AC grid electricity consumption and the residual grid electricity are measured by a bidirectional meter. The appliance load can be monitored through a smart socket.

## 3. Home Energy Management Model

Household electrical appliances, in addition to room temperature heating and domestic hot water systems, can be divided into automatic of appliance (AOA) and manual of appliance (MOA). AOA refers to appliances that can be operated automatically without human intervention, such as washing machines and dishwashers. MOA means the devices that must be operated manually by the user, such as computers, TVs, and hoovers. Since MOAs are only suitable for manual switching, other electrical appliance strategies in the home are aimed at AOAs.

### 3.1. Photovoltaic Cell Model.
The output of power photovoltaic cell is a function of solar irradiance and temperature, and it can be obtained using the daily irradiance curve.

The output power can be expressed as follows:

$$\begin{cases} P_{\text{PV}}(t) = P_{\text{STC}} \dfrac{G(t)}{G_{\text{STC}}}[1 + k(T(t) - T_{\text{STC}})], \\ T(t) = T_{\text{air}}(t) + 0.0318G(t)(1 + 0.031T_{\text{air}}(t))(1 - 0.042V_{\text{W}}), \end{cases} \tag{1}$$

where $P_{\text{PV}}(t)$ is the photovoltaic output, $P_{\text{STC}}$ is the maximum output under standard test conditions, $G(t)$ is the current solar irradiance, $G_{\text{STC}}$ is the rated solar irradiance, $k$ is the temperature coefficient, $T(t)$ is the temperature of cell module at the current moment, $T_{\text{air}}(t)$ is the ambient temperature, $T_{\text{STC}}$ is the rated reference temperature, and $V_{\text{W}}$ is the current wind speed.

### 3.2. Battery Model.
The battery model mainly regards the state during the charging and discharging process. The remaining capacity of battery is expressed as

$$S_{\text{SOC}}(t + 1) = \frac{C_{\text{r}}}{C_{\text{N}}} \times 100\% = \begin{cases} S_{\text{SOC}}(t) + \eta_{\text{ch}}P_{\text{ch}}(t)\Delta t, \\ S_{\text{SOC}}(t) - \dfrac{P_{\text{dis}}(t)\Delta t}{\eta_{\text{dis}}}, \end{cases} \tag{2}$$

where $S_{\text{SOC}}(t+1)$ is the next charge state, $C_r$ is the actual charge capacity, $C_N$ is the nominal charge capacity; $S_{\text{SOC}}(t)$ is the current charge state, $\eta_{\text{ch}}$ is the battery charge efficiency; $\eta_{\text{dis}}$ is the battery discharge efficiency, $P_{\text{ch}}(t)$ is the current charge power, $P_{\text{dis}}(t)$ is the current discharge power, and $P_{dis}(t)$ is the charge and discharge time.

Additionally, the life of battery is related to the depth of discharge and the number of cycles, where the life consumption $D$ of lead-acid battery can be expressed as

$$D_i = \sum_{i=1}^{n} \frac{1}{a_1 + a_2 e^{-a_3\left(1-S_{\text{SOC}}^{(i)}\right)} + a_4 e^{-a_5\left(1-S_{\text{SOC}}^{(i)}\right)}}, \quad (3)$$

where $S_{\text{SOC}}^{(i)}$ is the state of charge when it is transferred from discharge to charge, which represents one discharge cycle.

The battery life parameter is obtained by fitting the number of cycle curves provided by its equipment manufacturer.

3.3. Load Model. The loads in the HEMS can be divided into 4 categories in accordance with their control level as follows:

(1) Temperature-controlled loads, which include air conditioners, water heaters, and refrigerators, with a certain degree of cooling or heat storage capacity

(2) Active controllable loads, including washing machines and rice cookers, with a fixed working cycle and a certain flexibility of use time

(3) Passive controllable loads, including lights and fans, which can be intelligently controlled but have inflexible operating hours

(4) Noncontrollable loads

3.3.1. Air Conditioner. Assume that the air conditioner is operating in cooling mode and the operating state is related to the room temperature setting. The air conditioner is energized when the room temperature is above the maximum value. As the temperature is below the minimum value, the air conditioner is disconnected. It maintains the original state if the temperature is within the set range. Its control model and the comfort index $K_{\text{AC},t}$ are shown as follows:

$$S_{\text{AC},t} = \begin{cases} 0 & T_{\text{AC},t} < T_{\text{AC,s}}, \\ 1 & T_{\text{AC},t} > T_{\text{AC,s}} + \Delta T_{\text{AC}}, \\ S_{\text{AC},t-1} & T_{\text{AC,s}} < T_{\text{AC},t} < T_{\text{AC,s}} + \Delta T_{\text{AC}}, \end{cases} \quad (4)$$

$$K_{\text{AC},t} = \frac{T_{\text{AC},t} - T_{\text{AC,s}}}{\Delta T_{\text{AC}}},$$

where $S_{\text{AC},t}$ is the state of air conditioning (the value of 0 means power off; the value of 1 means power on). $T_{\text{AC,s}}$ is the minimum setting temperature. $\Delta T_{\text{AC}}$ is the room temperature set range. $T_{\text{AC},t}$ is the room temperature at time $t$.

$K_{\text{AC},t}$ is the difference between the current room temperature and the minimum set value after standardization; the higher the room temperature, the greater the comfort index $K_{\text{AC},t}$; the lower the satisfaction of the user, thus the higher the power priority. During demand response, the power supply is controlled based on the priority of the air conditioner [25].

3.3.2. Water Heater. Water heater operation status is related to the water temperature setting. When the water temperature is above the maximum temperature $T_{\text{WH,s}}$, the water heater is disconnected; when it is below the minimum temperature, the water heater is powered on; when it is within the set range, it remains in the original state. The water heater control model and its comfort index $K_{\text{WH},t}$ are given as follows:

$$S_{\text{WH},t} = \begin{cases} 0 & T_{\text{WH},t} > T_{\text{WH,s}}, \\ 1 & T_{\text{WH},t} < T_{\text{WH,s}} - \Delta T_{\text{WH}}, \\ S_{\text{WH},t-1} & T_{\text{WH,s}} - \Delta T_{\text{WH}} < T_{\text{WH},t} < T_{\text{WH,s}}, \end{cases}$$

$$K_{\text{WH},t} = \frac{T_{\text{WH,s}} - T_{\text{WH},t}}{\Delta T_{\text{WH}}}, \quad (5)$$

where $S_{\text{WH},t}$ is the working state of water heater at time $t$ (the value of 0 means power off; the value of 1 means power on). $T_{\text{WH,s}}$ is the highest water temperature setting value. $\Delta T_{\text{WH}}$ is the water temperature setting range. $T_{\text{WH},t}$ is the water temperature at time $t$.

$K_{\text{WH},t}$ is the difference between the current water temperature and the highest set value after normalisation; the lower the water temperature, the greater the comfort index $K_{\text{WH},t}$; the lower the customer satisfaction, thus the higher the priority of electricity consumption. During demand response, the water heater is controlled based on its priority.

3.3.3. Electric Vehicles. It is assumed that the electric vehicle is plug-and-charge type. On the basis of its charging characteristics, the load demand is set as follows. The electric vehicle should be fully charged by the specified time [26]. For instance, if charging is assumed to start at 21:00, it is set to reach full charge at 04:00 on the next day. The electric vehicle control model is presented in Equation (6). The comfort index of electric vehicles is calculated in a different way from air conditioning and water heaters. It is specified that the comfort index tends to infinity when the electric vehicle is not expected to finish charging before the specified time; otherwise, the index is zero

$$S_{\text{EV},t} = \begin{cases} 0 & Q_t \geq Q_{\max}, \\ 1 & Q_t < Q_{\max}, \end{cases} \quad (6)$$

$$\begin{cases} K_{\text{EV},t} = 0 & Q_t > Q_{\min,t}, \\ K_{\text{EV},t} \longrightarrow \infty & Q_t \leq Q_{\min,t}, \end{cases} \quad (7)$$

where $S_{\text{EV},t}$ is the state of the EV at time $t$ (a value of 0 means disconnected; a value of 1 means energized), $Q_t$ is the charge of the EV at time $t$, $Q_{\max}$ is the maximum value of the

(a)

(b)

(c)

(d)

Figure 3: Schematic diagram of neural network model: (a) DAE; (b) DBN; (c) CNN. (d) LSTM.

battery state of charge (SOC), $Q_{\min,t}$ is the minimum value of the battery SOC at time $t$.

When $K_{\mathrm{EV},t}$ tends to infinity, it indicates that the EV cannot finish charging before the specified time, at which time its power consumption priority can be set to the highest. During demand response, the power supply state of the EV is controlled based on its priority level.

## 4. Artificial Intelligence and Its Application in HEMS

Artificial intelligence is a comprehensive discipline developed through the interplay of many disciplines such as mathematical logic, computer science, cybernetics, information theory, neurobiology, and linguistics. The main objective is to develop a theory of intelligent information processing and to design computer systems that can display certain behaviors approximating human intelligence.

4.1. Deep Learning. Deep learning was originally proposed by Hinton at the University of Toronto. Deep learning algorithms draw on the neural working mechanism of the brain, which is an extension and development on the traditional artificial neural network technology. Through increasing the number of hidden layers of artificial neural networks and proposing effective training methods, the gradient diffusion (GD) problem of neural network training has been solved, which effectively improves the feature extraction ability and classification ability of neural networks. According to the problems and tasks, different model structures and open-source technology platforms have been developed for deep learning techniques. The main deep learning models are deep autoencoder (DAE), deep belief networks (DBN), convolutional neural network (CNN), and long short-term memory (LSTM). A typical deep learning model structure is shown in Figure 3. The main open-source platforms are TensorFlow, Caffe, DMTK, SystemML, etc.

Figure 4: Technical architecture of knowledge graph.

The deep learning model has many parameters, a large training data scale, and a large amount of calculation, which consumes massive computing resources. Deep learning model parameters need to be debugged and optimized, such as network structure selection, neuron number setting, weight parameter initialization, learning rate adjustment, and minibatch control. In practice, it requires multiple trainings and constant exploration and experimentation, which further increases the demand for computing resources. With the increase of model depth and training data volume, the training acceleration method of the deep learning model becomes more and more important. Typical acceleration methods mainly include algorithm optimization, GPU acceleration, and computing cluster acceleration.

*4.2. Knowledge Graph.* Knowledge graph, as another important research direction in the field of artificial intelligence, is widely used in semantic search and automatic question answering. The knowledge graph usually organizes knowledge in the form of a network, describing the relationship between entities in the real world; each node represents an entity; and each edge represents the relationship between entities. After Google proposed the concept of knowledge graph, this form of network representation of knowledge has been widely recognized. The main research goal of knowledge graph is to propose knowledge from unstructured or semistructured information and carry out structured processing, automatic construction of knowledge base, knowledge reasoning, and so on. Knowledge representation is the basis of the research and application of knowledge graphs. The Word2Vec word representation model and toolkit found that there is a translation-invariant relationship in the word vector space, which makes representation learning gain widespread attention in the field of natural language processing. The TransE model expresses the relationship in the knowledge base as a translation vector



Figure 5: Structure of expert system.



Figure 6: Structure of Agent.

between entities, which has become the mainstream research method of knowledge representation today. The technical architecture of the knowledge graph is shown in Figure 4, including three parts: information extraction, knowledge fusion, and knowledge processing. Information extraction includes key technologies such as entity extraction, relationship extraction, and attribute extraction. Knowledge fusion includes entity disambiguation, coreference analysis, and

FIGURE 7: MAS system architecture: (a) centralized structure; (b) decentralized structure; (c) hybrid structure.

knowledge fusion. Knowledge processing includes knowledge reasoning, quality evaluation, and ontology extraction.

*4.3. Expert System.* The expert system was produced in the mid-1960s and is an important branch of artificial intelligence applications. Expert system is a computer program system that solves specific problems based on the knowledge of specialized fields. It can simulate the thinking activities of human experts to solve complex problems through reasoning and judgment like experts. A typical expert system is mainly composed of knowledge base, database, inference engine, and man-machine interface, and its structure is shown in Figure 5. There are many problems in the power system that need to be solved by expert planning, designers, dispatchers, etc. in related fields. Some rely on expert experience, and some integrate judgment based on experience with results obtained based on numerical analysis methods. Expert systems have become the most mature artificial intelligence technology used in power systems so far. The main application areas include power grid monitoring and fault diagnosis, power grid dispatching operation guidance, and fault recovery.

*4.4. Agent Technology.* Agent is an entity with high self-control capability that runs in a dynamic environment, and its structure is shown in Figure 6. From a software perspective, it is a computer program that communicates with the outside through a predefined protocol and is loosely coupled. Distributed intelligent solution is performed in a way. It is an entity that can work autonomously and has semantic interoperability and protocol interaction capabili-

ties. It is a distributed technology in the field of artificial intelligence. Due to the advantages of adaptability and openness, it has a good prospect in the new generation of dispatching automation system.

Agent encapsulates the tasks and goals to be completed in the target module and collects external data through the perception module. The information processing module makes corresponding decisions based on the data collected by the sensor. The communication module provides conditions for coordination between Agents. An independent rule library is set up in the Agent to provide choices for decision-making and improve the efficiency. Mobile agent server (MAS) achieves the goal of entire system by coordinating and controlling each agent. The architecture of MAS system can generally be divided into three types: centralized structure, decentralized structure, and hybrid structure, as shown in Figure 7.

## 5. Resident HEMS Application

On December 9, 2019, the first demonstration project for HEMS in Jiangsu was completed and put into operation in Huangzhuang Village in Jinhu County, Jiangsu Province. Jiangsu Electric Power Co., Ltd., of State Grid installed a set of ubiquitous Internet of Things devices such as energy controllers and household appliances in the demonstration area to realize the in-depth perception and precise adjustment of residential loads at the electrical level, allowing residents to interact friendly with the demand of power grid. By cooperating with the cloud master station, the energy controller can accurately predict load fluctuations in the station

FIGURE 8: User response to air conditioner.

area and effectively converge and regulate customer-side load resources without affecting the daily energy consumption. The temperature of residential air conditioner is adjusted through the energy control system. During the peak period of power consumption in the station area, the heating time of water heater would be adjusted to reduce the total load.

For instance, on a peak load forecasting day, the station area is in the second-level interval during the period from 19:00 to 20:40. According to the load forecasting result and the load coordinated control framework, the main adjustment potential of this period is the air conditioning load, and this period is selected through air conditioning adjustment. 10 users are selected, and their air conditioning temperatures are adjusted from the original 25°C to 23°C. The air condition response of user is shown in Figure 8.

## 6. Conclusion

(i) HEMS connects users and the grid. The smart terminal of HEMS enables to read, process, and display information such as household electricity, water, and faults, so as to guide users to use electricity reasonably and save energy. Users can realize remote monitoring of home appliances and achieve prepaid services through the Internet, mobile phones, etc.

(ii) Advanced sensing equipment can sense changes in the external environment in real time and communicate with humans in time. The artificial intelligence enables power equipment to calculate and fuse the sensed information to reach the corresponding conclusion and report it to the user. It can even analyze real-time information and historical data and propose long-term decision-making suggestions to provide reference for user services

(iii) Traditional artificial intelligence technologies such as expert systems, neural networks, fuzzy sets, and heuristic search algorithms have been widely used

in power systems. New-generation artificial intelligence technology is a breakthrough in distributed power and distributed energy storage. In response to the complex nonlinearity, uncertainty, and temporal and spatial differences brought by the high-proportion access of various new energy sources to the grid, effective solutions have been proposed

## Data Availability

The smart home energy management system data used to support the findings of this study were supplied by Yunlong Ma under license and so cannot be made freely available. Requests for access to these data should be made to Yunlong Ma (3397599241@qq.com).

## Conflicts of Interest

The authors declare no conflict of interest.

## Authors' Contributions

Y. Ma and X. Chen proposed the concepts and ideas. W. Li analyzed the results. J. Yang wrote this paper and revised the contents of this manuscript.

## Acknowledgments

## References

[1] S. Young-Sung and M. Kyeong-Deok, "Home energy management system based on power line communication," in *2010 Digest of Technical Papers International Conference on Consumer Electronics (ICCE)*, pp. 115-116, Las Vegas, NV, USA, Jan. 2010.

[2] Y. H. Lin and M. S. Tsai, "An advanced home energy management system facilitated by nonintrusive load monitoring with automated multiobjective power scheduling," *IEEE Transactions on Smart Grid*, vol. 6, no. 4, pp. 1839–1851, 2015.

[3] F. De Angelis, M. Boaro, D. Fuselli, S. Squartini, F. Piazza, and Q. Wei, "Optimal home energy management under dynamic electrical and thermal constraints," *IEEE Transactions on Industrial Informatics*, vol. 9, no. 3, pp. 1518–1527, 2013.

[4] Z. A. Khan, A. Khalid, N. Javaid, A. Haseeb, T. Saba, and M. Shafiq, "Exploiting nature-inspired-based artificial intelligence techniques for coordinated day-ahead scheduling to efficiently manage energy in smart grid," *IEEE Access*, vol. 7, pp. 140102–140125, 2019.

[5] J. H. Liu, *Study on the energy optimization scheduling model for home energy management system*, Hunan University, 2014.

[6] K. Zor, O. Timur, and A. Teke, "A state-of-the-art review of artificial intelligence techniques for short-term electric load forecasting," in *2017 6th International Youth Conference on Energy (IYCE)*, pp. 1–7, Budapest, Hungary, 2017.

[7] A. Marnerides, P. Smith, A. Schaeffer-Filho, and A. Mauthe, "Power consumption profiling using energy time-frequency

distributions in smart grids," *IEEE Communications Letters*, vol. 19, no. 1, pp. 46–49, 2015.

[8] A.-H. Mohsenian-Rad and A. Leon-Garcia, "Optimal residential load control with price prediction in real-time electricity pricing environments," *IEEE Transactions on Smart Grid*, vol. 1, no. 2, pp. 120–133, 2010.

[9] A. H. Mohsenian-Rad, V. W. S. Wong, J. Jatskevich, R. Schober, and A. Leon-Garcia, "Autonomous demand-side management based on game-theoretic energy consumption scheduling for the future smart grid," *IEEE Transactions on Smart Grid*, vol. 1, no. 3, pp. 320–331, 2010.

[10] C. O. Adika and L. Wang, "Autonomous appliance scheduling for household energy management," *IEEE Transactions on Smart Grid*, vol. 5, no. 2, pp. 673–682, 2014.

[11] A. Anvari-Moghaddam, H. Monsef, and A. Rahimi-Kian, "Optimal smart home energy management considering energy saving and a comfortable lifestyle," *IEEE Transactions on Smart Grid*, vol. 6, no. 1, pp. 324–332, 2015.

[12] L. Atzori, A. Iera, and G. Morabito, "The Internet of Things: a survey," *Computer Networks*, vol. 54, no. 15, pp. 2787–2805, 2010.

[13] G. K. F. Tso and K. K. W. Yau, "Predicting electricity energy consumption: a comparison of regression analysis, decision tree and neural networks," *Energy*, vol. 32, no. 9, pp. 1761–1768, 2007.

[14] M. Q. Raza and A. Khosravi, "A review on artificial intelligence based load demand forecasting techniques for smart grid and buildings," *Renewable and Sustainable Energy Reviews*, vol. 50, pp. 1352–1372, 2015.

[15] D. B. Richardson, "Electric vehicles and the electric grid: A review of modeling approaches, Impacts, and renewable energy integration," *Renewable and Sustainable Energy Reviews*, vol. 19, pp. 247–254, 2013.

[16] S. Han, S. Han, and K. Sezaki, "Development of an optimal vehicle-to-grid aggregator for frequency regulation," *IEEE Transactions on Smart Grid*, vol. 1, no. 1, pp. 65–72, 2010.

[17] J. Zheng, D. Gao, and L. Lin, "Smart meters in smart grid: an overview," in *2013 IEEE Green Technologies Conference (GreenTech)*, pp. 57–64, Denver, CO, USA, Apr. 2013.

[18] C. Zhou, K. Qian, M. Allan, and W. Zhou, "Modeling of the cost of EV battery wear due to V2G application in power systems," *IEEE Transactions on Energy Conversion*, vol. 26, no. 4, pp. 1041–1050, 2011.

[19] J. Donadee and M. D. Ilic, "Stochastic optimization of grid to vehicle frequency regulation capacity bids," *IEEE Transactions on Smart Grid*, vol. 5, no. 2, pp. 1061–1069, 2014.

[20] H. Farhangi, "The path of the smart grid," *IEEE Power and Energy Magazine*, vol. 8, no. 1, pp. 18–28, 2009.

[21] D.-M. Han and J.-H. Lim, "Design and implementation of smart home energy management systems based on ZigBee," *IEEE Transactions on Consumer Electronics*, vol. 56, no. 3, pp. 1417–1425, 2010.

[22] Y. G. Ha, "Dynamic integration of zigbee home networks into home gateways using OSGI service registry," *IEEE Transactions on Consumer Electronics*, vol. 55, no. 2, pp. 470–476, 2009.

[23] D. M. Han and J.-H. Lim, "Smart home energy management system using IEEE 802.15.4 and ZigBee," *IEEE Transactions on Consumer Electronics*, vol. 56, no. 3, pp. 1403–1410, 2010.

[24] C. Suh and Y. B. Ko, "Design and implementation of intelligent home control systems based on active sensor networks," *IEEE Transactions on Consumer Electronics*, vol. 54, no. 3, pp. 1177–1184, 2008.

[25] K. M. Tsui and S. C. Chan, "Demand response optimization for smart home scheduling under real-time pricing," *IEEE Transactions on Smart Grid*, vol. 3, no. 4, pp. 1812–1821, 2012.

[26] E. B. Iversen, J. M. Morales, and H. Madsen, "Optimal charging of an electric vehicle using a Markov decision process," *Applied Energy*, vol. 123, pp. 1–12, 2014.

*Research Article*

# DWCA-YOLOv5: An Improve Single Shot Detector for Safety Helmet Detection

**Zhang Jin** [1,2,3] **Peiqi Qu,**[1] **Cheng Sun,**[4] **Meng Luo,**[4] **Yan Gui,**[2] **Jianming Zhang** [2] **and Hong Liu** [1]

[1]*School of Information Science and Engineering, Hunan Normal University, Changsha 410081, China*
[2]*School of Computer and Communication Engineering, Changsha University of Science and Technology, Changsha 410114, China*
[3]*Key Laboratory of Industrial Control Technology, Zhejiang University, Hangzhou 310058, China*
[4]*School of Mathematics and Statistics, Hunan Normal University, Changsha 410081, China*

Correspondence should be addressed to Hong Liu; qpeggy@hunnu.edu.cn

Aiming at solving the problem that the detection methods used in the existing helmet detection research has low detection efficiency and the cumulative error influences accuracy, a new algorithm for improving YOLOv5 helmet wearing detection is proposed. First of all, we use the $K$-means++ algorithm to improve the size matching degree of the a priori anchor box; secondly, integrate the Depthwise Coordinate Attention (DWCA) mechanism in the backbone network, so that the network can learn the weight of each channel independently and enhance the information dissemination between features, thereby strengthening the network's ability to distinguish foreground and background. The experimental results show as follows: in the self-made safety helmet wearing detection dataset, the average accuracy rate reached 95.9%, the average accuracy of the helmet detection reached 96.5%, and the average accuracy of the worker's head detection reached 95.2%. Making a comparison with the YOLOv5 algorithm, our model has a 3% increase in the average accuracy of helmet detection, which is in line with the accuracy requirements of helmet wearing detection in complex construction scenarios.

## 1. Introduction

According to a series of statistical reports issued by the Ministry of Housing and Urban-Rural Development, compared with 934 accidents and 840 deaths in 2018, there were a total of 773 construction production safety accidents and 904 deaths across the country in 2019. The number of accidents and deaths increased by 5.31% and 7.62%. In general, the number of accidents in the construction industry is showing a gradual increase. In the literature [1], when studying the relationship between the use of safety protection equipment and the number of deaths in construction sites, it was found that 67.95% of the victims had not used or used safety protection (such as safety helmets and safety belts). Due to the weak awareness of safety protection of construction workers, the importance of wearing safety helmets is often ignored. At the construction site, manual supervision is usually used to

monitor whether workers wear safety helmets [2], which makes it impossible to manage all construction workers promptly on the construction site and to know the movement tracks of all construction workers. The use of automatic monitoring methods helps to monitor the construction personnel and confirm the specific conditions of all construction workers wearing helmets at the construction site, especially when the traditional monitoring methods are time-consuming and expensive, easy to detect errors, and are not enough to meet the safety of modern building construction management requirements. The use of automatic supervision of deep learning methods is conducive to supervising all construction personnel onsite.

Traditional object detection often uses an artificial selection of features and design and training classifiers based on specific detection objects. This method is highly subjective, complex in the design process, has poor generalization

ability, and has great limitations in engineering applications. In recent years, due to the fact that convolutional neural networks (CNN) do not use an artificial selection of features, they have gradually been sought after by scholars in the field of deep learning. The deep convolutional neural network has good comprehensive performance in the field of object detection. In 2014, Girshick et al. successfully proposed R-CNN [3], fast R-CNN [4], and faster R-CNN [5], which were verified in the PASCAL VOC2007dataset, respectively, and gradually improved the experimental effect. The method of extracting feature frames by these models gradually changes from selective search to regional proposal network (RPN), thus getting rid of the traditional manual feature extraction method. In 2015, Redmon and others proposed a one-stage object detection model YOLO [6], which abstracted the detection task as a regression problem for the first time, avoiding the cumbersome operation of dividing the detection task into two steps in the R-CNN series. In 2016, Liu et al. proposed the SSD [7] detection algorithm, which introduced a multiscale detection method, which can effectively detect groups of small targets. In 2017, Lin et al. proposed the RetinaNet [8] dense detector, which solves the problem of extreme foreground and background imbalance encountered during training by reshaping the standard entropy loss. In 2017, Redmon and others proposed the YOLOv2 [9] detection model, which selected a new basic model Darknet-19 to achieve end-to-end training. In 2018, Redmon et al. proposed YOLOV3 [10] based onYOLOV1 and YOLOV2. In this model, the FPN method was adopted to integrate three different sizes feature maps to accomplish detection tasks, which significantly improved the detection effect of small-size targets. In April 2020, Bochkovskiy proposed YOLOv4 [11], which uses PANet instead of FPN used in YOLOv3 as the path aggregation method; at the same time, the backbone network uses CSP Darknet53, which significantly enhances the detection accuracy of the network. In June 2020, Glenn proposed YOLov5 [12], which designed a new focus structure and added it to the backbone network to achieve a new benchmark for the perfect combination of speed and accuracy.

Because of the rapid rise of computer vision in the direction of object detection, more and more researchers are focusing on combining deep learning with practical application scenarios. For example, Chen et al. [13] improved the SSD model by adding an inception module before the prediction layer to achieve rapid and accurate detection of small vehicles. Tian et al. [14] used DenseNet to optimize the low-resolution layer in the feature layer of the YOLOv3 network and applied the improved YOLOv3 to the detection of anthrax lesions on the surface of orchard apples to achieve real-time detection. Dashun et al. [15] applied the improved RetinanNet network to the field of pedestrian detection and realized the rapid detection of multispectral pedestrians. Zhong et al. [16] used the LocNet positioning module to replace the boundary regression module to improve the faster R-CNN model and applied it to multidirectional text instance detection. Zhang et al. [17, 18] used the residual network (reset) in the prediction part to encode the input features of the image and chose to increase the deconvolu-

tion layer to change the MMDetection network model in the process of feature information decoding, to achieve a higher crowd in dense scenes. And it can be seen that deep learning has become a popular research direction, and it has become the mainstream field in combination with actual application scenarios.

Safety helmet detection is one of the application areas of object detection. So far, many researchers at home and abroad have conducted several related investigations on safety helmet detection. In 2013, Kelm [19] and others designed a mobile radio frequency identification (RFID) portal to check the compliance of construction workers wearing safety protective equipment. However, the recognition area of the radio frequency identification reader is limited. It is only recommended that the helmet be close to the worker, but it cannot be confirmed whether the helmet is worn correctly. In 2014, Liu [20] and others used a combination of support vector machines and skin color detection to achieve helmet detection. In 2016, Rubaiya [21] and Silva [22] and others combined the histogram of gradient (HOG) algorithm with the frequency domain-related information in the image for human detection and then used the circular Hough transform (CHT) to detect the helmet. In 2017, Li [23] and others used the vibe algorithm to locate the human body position, followed by the embossing algorithm to detect the worker's head and finally combined the HOG algorithm and SVM to realize the helmet wearing detection. In 2018, Wu et al. [24] used Hu moment invariant (HMI), color histogram (CH), and local binary pattern (LBP) to extract the characteristics of different color helmets and then constructed a hierarchical support vector machine (H-SVM) for safety cap wearing detection. Due to the complex environment, the detection accuracy of helmet wearing detection is low at this stage, which is quite different from the management requirements in actual building construction.

In this paper, two types of targets for construction workers wearing helmets and those not wearing helmets are the detection tasks, and more than 7,000 pictures are collected from the Internet for preprocessing to construct a helmet detection dataset. Select the YOLOv5 network model as the main body, and first, use the $k$-means++ algorithm to cluster the target anchor box to obtain a bounding box suitable for the target, so that the model can converge faster. Secondly, a new DWCA module is designed and integrated with the features of the backbone network to strengthen the attention to enhance the attention of the detection target and improve the ability to resist background interference. According to the final experimental results, the average detection accuracy (mAP) of the DWCA-YOLOv5 detection model has been significantly improved, and it can effectively detect the unsafe behavior of workers on the construction site not wearing helmets.

## 2. Related Work

*2.1. YOLOv5 Algorithm Principle.* YOLOv5 is a new-generation target detection network of the YOLO series. It is a product of continuous integration and innovation based on YOLOv3 and YOLOv4. Secondly, YOLOv5 has achieved

FIGURE 1: Network structure diagram of YOLO v5s.

better results in PASCAL VOC and COCO object detection tasks; so, this article uses the YOLOv5 detection network to detect the construction workers' helmet wearing.

The YOLOv5 object detection network official gave four network models: YOLOv5s, YOLOv5m, YOLOv5l, and YOLOv5x. The three models of YOLOv5m, YOLOv5l, and YOLOv5x are the products of continuous deepening and widening based on YOLOv5s. The YOLOv5 network structure is divided into four parts: input, backbone, neck, and prediction. YOLOv5 adds Mosaic data enhancement in the data input part; focus structure and CSP structure are used in the backbone; the FPN + PAN structure is added to the neck; the prediction part improves the bounding box loss function from CIOU_Loss to GIOU_Loss; YOLOv5 targets many in the postprocessing process of object detection. The screening of the target anchor frame adopts a weighted NMS operation.

Compared with YOLOv4, YOLOv5 has a new focus structure in the backbone network, which is mainly used for slicing operations. In the YOLOv5s network model, an ordinary image with a size of $3 \times 608 \times 608$ is input into the network, and after a focus slice operation, the feature map with a size of $12 \times 304 \times 304$ is converted, followed by the ordinary convolution operation of 32 convolution kernels. It is finally converted into a feature map with a size of $32 \times 304 \times 304$. Different from the YOLOv4 network model that only uses the CSP structure in the backbone network, the YOLOv5 network model has designed two new CSP structures. Taking the YOLOv5s network model as an example, the backbone network uses the CSP1_1 structure and the CSP1_3 structure, and the neck uses the CSP2_1 structure to enhance the feature fusion between networks. The network structure of YOLOv5s is shown in Figure 1.

*2.2. DWCA Moduel.* The traditional channel module is dedicated to constructing various channel importance weight functions. For example, SEnet [25] obtained a significant

effect improvement by calculating channel attention with the aid of a 2D global pool and with a small computational overhead. However, SENet only considers the encoding of information between channels and ignores the importance of position information, which is essential for capturing the structure of objects in vision tasks. Coordinate attention [26] has achieved significant performance improvement by encoding the interchannel relationship and long-term dependence. ECANet [27] proposed a method that does not take dimensionality reduction measures to achieve cross-channel local interaction and a method that automatically adapts to select one-dimensional ordinary convolution, thereby achieving performance improvement. CBAM [28] and BAM [29] reduce the channel input dimension of the tensor and secondly use convolution to calculate spatial attention to use position information. However, convolution can only capture local relationships, but not what is needed for modeling long-term dependence on visual tasks.

To solve the above problems, we designed a new attention mechanism based on previous work, which integrates the position information in the feature space into the channel attention, so that the network can participate in a larger area and at the same time avoid a lot of model parameters overhead. The structure diagram of DWCA mechanism is shown in Figure 2.

To reduce the lack of relevant location information caused by two-dimensional global sharing, we use two one-dimensional global aggregation operations to decompose the channel attention into two aggregated features along with the vertical and horizontal directions and then aggregate the obtained features into two independent directional perception features map. To promote the module to capture the remote spatial interaction with precise location information, this paper decomposes the global pooling according to formula (1) and transforms it into a one-to-one dimensional feature encoding operation. The specific operation process is

FIGURE 2: DWCA model network structure.

as follows: first, use a pooling kernel of size $(H, 1)$ or $(1, W)$ to encode the single dimension and horizontal and vertical coordinates of input $X$. Therefore, the $c$th channel of the output with a height of $h$ can be seen below, and the details are shown in formula (2).

$$Z_c = \frac{1}{H \times W} \sum_{i=1}^{H} \sum_{j=1}^{W} x_c(i, j). \tag{1}$$

In the formula, $Z_c$ is the output related to the $c$th channel, $H$ is the height of the input $X$, and $W$ is the width of the input $X$.

$$Z_c^h(h) = \frac{1}{W} \sum_{0 \leq i < W} x_c(h, i). \tag{2}$$

In the formula, $Z_c^h(h)$ is the specific output of the $c$th channel where the height is $h$, and $W$ is the width of the input $X$.

By analogy, the specific output of the $c$th channel with width $w$ can be seen below, see formula (3) for details.

$$Z_c^w(w) = \frac{1}{H} \sum_{0 \leq j < H} x_c(j, w). \tag{3}$$

In the formula, $Z_c^h(w)$ is the output of the $c$th channel at the width $w$, and $H$ is the height of the input $X$.

By extending the above two features to the transformation of the aggregation of the two spatial dimensions, the direction-aware feature map is obtained, followed by CONCAT operation, and then use the shared $1 \times 1$ conventional convolution transformation function to transform it, such as the formula (4) shown.

$$f = \delta\left(F1\left(\left[z^h, z^h\right]\right)\right). \tag{4}$$

In the formula, $f \in R^{C/r \times (H+W)}$ is an intermediate feature map, which encodes the spatially related information in the vertical and horizontal directions, $\delta$ is a nonlinear activation function, and $[\cdot, \cdot]$ represents the splicing operation along the spatial dimension.

Then, follow the spatial dimension, and $f$ is transformed into two independent tensors $f^h \in R^{C/r \times H}$ and $f^w \in R^{C/r \times W}$, using two effective depthwise separable convolution transforms $f^h$ and $f^w$ and then transforms the tensors $f^h$ and $f^w$ with the same number of channels into input $X$, as shown in formulae (5) and (6).

$$g^h = \sigma\left(F_h\left(f^h\right)\right), \tag{5}$$

$$g^w = \sigma\left(F_w\left(f^w\right)\right). \tag{6}$$

In the formula, $g^h$ and $g^w$ are the attention weights to be expanded, $\sigma$ is the Sigmoid function, and $r$ is the reduction ratio of the number of channels.

Finally, the entire DWCA module can be expressed as follows, see formula (7) for details:

$$y_c(i, j) = x_c(i, j) \times g_c^h(i) \times g_c^w(j), \tag{7}$$

### 2.3. Improve the YOLOv5 Algorithm

*2.3.1. K-Means++ for Target Frame Optimization.* Perform $K$-means dimensional clustering on the general target detection dataset COCO to obtain the initial a priori anchor frame parameters of YOLOv5. However, because the target types of the COCO dataset have 80 categories, the helmet detection types in this article only have two categories, which cannot be to meet the actual needs of helmet wearing detection, and the size of the a priori frame needs to be redesigned. Compared with the size of the anchor frame designed only relying on human prior knowledge, for the helmet wearing dataset, we select the $K$-means++ algorithm to perform multidimensional clustering on the marked target frame, resulting in different numbers and sizes. As far as possible, the accurate matching between the a priori anchor frame and the actual object is achieved, thereby further improving the accuracy of helmet detection. In the clustering process, the average intersection ratio (IoU) corresponding to the number of centers of different clusters is shown in Figure 3.

Observing Figure 3, we can get that when the number of prior anchor box clusters is 0 to 9, the average intersection ratio shows a rapid upward trend, but when the number of a priori anchor boxes is 9 to 12, the average intersection ratio

FIGURE 3: The number of centers of different clusters and the average intersection ratio.

increases gradually gentle. To balance the calculation accuracy and efficiency, 9 a priori anchor frames are finally selected and equally distributed to 3 prediction branches of different sizes. The determined a priori anchor frame sizes are normalized as shown in Table 1.

Table 2 illustrates the pros and cons of the model's performance. Among them, the clustering method of a priori anchor box is changed from $K$-means algorithm to $K$-means++ algorithm, mAP has a certain improvement, and the improved YOLOv5 algorithm changes due to the network structure and the detection accuracy. There is also a big improvement. At the same time, selecting the $K$-means++ clustering method and the improved YOLOv5 model is 3.2 percentage points higher than the original YOLOv5 algorithm. The average accuracy of the self-made helmet wearing detection dataset reaches 95.9%, which can accurately detect whether the construction personnel wears a hard hat.

### 2.3.2. DWCA Module Fusion Design.

For small target detection tasks, as the sum of the model network layers gradually increases, the feature information of small targets that can be collected gradually decreases. So, it is easy to cause the network model to false detection and miss detection of small targets. The DWCA module itself is to integrate the location information of the feature space with the channel features so that the network can grasp the "key points" of the target features during the training process. However, under specific circumstances, which position of the DWCA module to perform feature fusion in the network model is effective is still a question to be studied.

In this paper, the DWCA module is merged into different positions of the network model, and the detection results are studied. According to the structure of the YOLOv5s network model, this paper will integrate the DWCA module in the three areas of the backbone network, the neck, and the prediction module of YOLOv5s. Since the DWCA module is to enhance the relationship between channel information and channel information in the feature space, our embeds

TABLE 1: Prior anchor box scales.

| Feature map scale | Anchor box size | | |
| --- | --- | --- | --- |
| | Anchor 1 | Anchor 2 | Anchor 3 |
| Small scale | (11.09,18) | (21.5,30.8) | (30.8,43) |
| Middle scale | (38.1, 60) | (52.3, 73.6) | (63,103.3) |
| Large scale | (89.2, 135) | (120, 207.5) | (209.4, 324) |

TABLE 2: Effect evaluation of different models on the test set.

| Detection model | Clustering method | AP50/% | | mAP/% |
| --- | --- | --- | --- | --- |
| | | Hat | Person | |
| Original YOLOv5 | $K$-means | 93.3 | 91.7 | 92.7 |
| Original YOLOv5 | $K$-means++ | 94.4 | 92.8 | 93.6 |
| Improved YOLOv5 | $K$-means | 95.5 | 94.6 | 95.1 |
| Improved YOLOv5 | $K$-means++ | 96.5 | 95.2 | 95.9 |

the DWCA module into each feature fusion area in the above three parts, thereby generating three new types based on the YOLOv5s algorithm. Network model is as follows: DWCA-YOLOv5s-backbone, DWCA-YOLOv5s-neck, and DWCA-YOLOv5s-prediction. Figure 4 shows the specific location where the DWCA module is integrated into the network.

In Figure 3(a), the DWCA module is integrated at CSP1_3 (i. e., the feature fusion) in the backbone network of YOLOV5s. In Figure 3(b), the DWCA module is integrated behind the CONCAT layer on the neck of YOLOV5s. In Figure 3(c), the DWCA module is integrated, respectively, before the convolution of each prediction in YOLOV5s. Table 3 shows the experimental results of whether the DWCA module is integrated with three different positions.

By visualizing the output of the same channel of the three fusion-designed networks, as shown in Figure 5 (only the channel output of the same feature map is visualized), the experimental results show that, compared to fusing the DWCA module into the network neck and network

(a) DWCA-YOLOv5s-backbone



(b) DWCA-YOLOv5s-neck

Figure 4: Continued.

(c) DWCA-YOLOv5s-prediction

FIGURE 4: Three YOLOv5s modes embedded in the DWCA modules.

TABLE 3: Comparison of results of different detection models.

| Network model | P/% | R/% | Model parameters/M | mAP/% |
|---|---|---|---|---|
| YOLOv5s | 76.4 | 92.5 | 7.26 | 92.7 |
| DWCA-YOLOv5-backbone | 82.5 | 95.4 | 7.27 | 95.9 |
| DWCA-YOLOv5-neck | 70.9 | 93.7 | 7.26 | 91.6 |
| DWCA-YOLOv5-prediction | 72.5 | 92.8 | 7.27 | 92.4 |

prediction part, fusing DWCA into the backbone network can effectively strengthen the semantic information of the feature layer on the instance and pay more attention to the target hidden in the lower layer, which is easy to ignore. The texture information and contour information can effectively improve the network's attention to small targets.

## 3. Experimental Results and Analysis

### 3.1. Dataset Construction.
In the detection direction, the dataset required by experiments has always been an essential basic condition. The safety helmet dataset that has been open sourced is only SHWD (SafetyHelmetWearing-Dataset). In this dataset, the category label data of not wearing a helmet is mainly derived from the SCUT-HEAD dataset. The SCUT-HEAD dataset is used by students in classroom scenarios monitoring diagrams or photos taken, so the dataset is not a standard construction site scene dataset, which does not meet the detection requirements of actual building construction scenarios. To solve this problem, this article self-made a helmet wearing detection dataset in construction scenarios. The main process of constructing this dataset includes data collection, screening, and processing.

### 3.2. Data Collection.
The images required for the dataset in this article mainly come from the surveillance video framing of the construction site, self-collecting on the construction site, and Internet crawling. The collected data includes two types of pictures of workers wearing and not wearing helmets in different environments, different resolutions, and different construction sites. Multiple sets of interference pictures are added to the dataset, such as construction workers wearing baseball caps and safety helmets. Construction workers with hats placed on the table or in hand, construction workers wearing bamboo woven hats, etc., increase the diversity of the dataset, thereby enhancing the robustness of the network. The sample map of the dataset collected this time is shown in Figure 6.

### 3.3. Data Screening and Processing.
The pictures collected from the surveillance video of the construction site are divided into frames or crawled on the Internet. Many of the pictures do not contain the construction personnel as the research object. They can be regarded as background pictures and have no practical significance for the study of this article. The picture data is confirmed as the background is deleted. This paper conducts a preliminary screening of the collected image and selects the images that meet the requirements as the annotation dataset.

Preprocess the data, convert the images that meet the requirements into.jpg format, and use the labeling tool labellmg to manually label each image, and the construction personnel in the image i under wearing a helmet (hat) and not wearing a helmet (person) These two categories are labeled, as shown in Figure 7; after processing, a corresponding XML tag file is formed, which contains the four

| Original image | DWCA-YOLOv5s-backbone | DWCA-YOLOv5s-neck | DWCA-YOLOv5s-prediction |
|---|---|---|---|



Figure 5: Model heat map comparison.



(a) Put on the table

(b) Hand held

(c) Hand held

(d) Hand held

(e) Hand held

(f) Baseball cap

(g) baseball cap

(h) Normal sample

(i) Bamboo braided hat

(j) Police_cap

(k) Normal sample

(l) Sunhat

Figure 6: Safety helmet sample image.

FIGURE 7: Safety helmet wearing status mark.

TABLE 4: Dataset category allocation.

| Target category | Training set a target number | A test set the target number | Total number of labeled targets |
| --- | --- | --- | --- |
| Wearing helmets category | 81836 | 11316 | 93152 |
| Not wearing helmets category | 98187 | 12021 | 110208 |

TABLE 5: Experiment operating environment.

| Category | Entry | Version |
| --- | --- | --- |
| Hardware configuration | System | Ubuntu 18.04 |
| | GPU | GeForce RTX 2080 Ti |
| | CPU | AMD Ryzen 7 3800X 8-Core |
| Software configuration | Python version | 3.8 |
| | Deep learning framework | Pytorch |
| | CUDA | 10.0 |

TABLE 6: Comparison of experimental results of multiple detection algorithms.

| Detection model | AP50/% | | mAP/% |
| --- | --- | --- | --- |
| | Hat | Person | |
| Faster RCNN | 80.8 | 42.2 | 61.5 |
| SSD | 78.8 | 68.2 | 73.5 |
| YOLOv3 | 89.12 | 80.7 | 84.9 |
| YOLOv3 + SPP | 90.5 | 86.3 | 88.4 |
| YOLOv5m | 94.8 | 93.1 | 93.9 |
| YOLOv5l | 95.1 | 93.5 | 94.3 |
| YOLOv5x | 95.6 | 94.3 | 95.0 |
| YOLOv5s | 93.3 | 91.7 | 92.7 |
| Ours | 96.5 | 95.2 | 95.9 |

coordinates of the target in the frame and the given category (PASCAL VOC format).

The final dataset obtained in this paper has a total of 7076 images. Among them, the specific information of whether to wear a helmet in the dataset can be seen in Table 4. And the dataset contains a variety of construction scenes, which can more fully reflect the actual construction scenes. The final dataset is subdivided into training and validation in line with the 9 : 1 division ratio. The number of training set pictures in the final 7076 picture dataset is 6,370 pictures, and there are 706 pictures in the test set.

3.4. Experimental Environment. During the experiment, this article has high requirements for the configuration of the operating environment, and GPU acceleration is required for the experiment. Table 5 shows the configuration instructions for the experiment operating environment of this article. The model building, training, and result testing are all completed under the PyTorch framework, using the CUDA parallel computing architecture and at the same time integrating the cuDNN acceleration library into the PyTorch framework to accelerate computer computing capabilities.

(a) Detection of strong light construction scene



(b) Detection of construction scenes occluded by steel bars



(c) Detection of targets of different sizes



(d) Detection of long-distance construction scenes

FIGURE 8: Comparison of test results of model parts under different construction scenarios.

## 4. Result Analysis

*4.1. Evaluation Index.* In object detection, detection accuracy, and recall, the average accuracy rate (mAP) is the basic index to test the training model's overall stability and performance. This article also uses the above evaluation indicators to detect helmet wearing model performance that is evaluated.

Apply the above evaluation indicators to the stability test of the helmet detection model and then the detection results of whether the construction worker wears a helmet are compared. Among them, $TP_{hat}$ (true example), $FP_{hat}$ (false positive example), $TN_{hat}$ (true negative example), and $FN_{hat}$ (false negative example) are key indicators used to describe accuracy. Specifically, TP refers to the sum of workers who did not wear helmets and whose test results were correct within the monitoring range of the construction site. FP indicates the sum of workers who wear helmets but are mistakenly detected, TN indicates that results are completely correct, and FN indicates the sum of workers who did not wear helmets but were mistakenly

detected as wearing helmets. The calculation process of accuracy rate and recall rate is shown in formulae (8) and (9).

$$AP_{hat} = \frac{TN_{hat} + TP_{hat}}{TN_{hat} + TP_{hat} + FP_{hat}}, \tag{8}$$

$$Recall_{hat} = \frac{TP_{hat}}{FN_{hat} + TP_{hat}}, \tag{9}$$

$$Precision_{hat} = \frac{TP_{hat}}{TP_{hat} + FP_{hat}}. \tag{10}$$

$Precision_{hat}$ represents the ratio of real cases ($TP_{hat}$) to the sum of real cases and false real cases ($TP_{hat}+FP_{hat}$), and the sum of real cases and false real cases is the total number of helmets; $Recall_{hat}$ represents the sum of real cases($TP_{hat}$)and real cases and false counterexamples ($TP_{hat} + FN_{hat}$). The ratio of true cases and false counterexamples is the actual number of helmets.

AP refers to the average value of all precisions obtained under all possible recall rates. The average precision of the mean is the average of the AP value in all categories, and the calculation formula is shown in (3).

$$mAP = \frac{1}{C} \sum_{c \in C} AP(c). \tag{11}$$

*4.2. Result Analysis.* This article uses the YOLOv5 algorithm for helmet wearing detection. To verify that the algorithm proposed has better results, the same number of test sets is used under the same configuration conditions, and several popular object detection networks at this stage are used for comparative experiments: faster RCNN, SSD, and YOLOv3. Among them, SSD and YOLOv3 are single-stage detection algorithms, and faster RCNN is a two-stage detection algorithm. The experimental results are evaluated using two evaluation indicators AP50 and mAP. The experimental results are shown in Table 6.

Observing Table 6, we can know that the DWCA-YOLOv5 algorithm can significantly improve the accuracy of detecting whether a worker is wearing a helmet. The average accuracy of the DWCA-YOLOv5 algorithm in this paper can reach 96.2% for the construction personnel who wear the helmet correctly and 95.1% for the construction personnel who do not wear the helmet. mAP (mean average precision) can reach 95.7%. Compared with faster RCNN and SSD, our model detection results are better. Compared with YOLOv3 and YOLOv5, the algorithm in this paper has a certain improvement in AP50 and mAP. This shows that the DWCA-YOLOv5 algorithm has an excellent performance in the accuracy of detection and detection of helmet wearing, and it can ensure the accuracy of helmet detection in a complex construction environment.

In addition, to more intuitively see the detection gap between different algorithms, this paper also collected 158 pictures of the construction work site as a test set. In this test set, we use YOLOv5 and our model to test separately. Some of the detection results are shown in Figure 8 below.

From Figure 8, we can observe that the operator who wears the helmet correctly is marked with a red frame, and the operator who does not wear the helmet is marked with a light green frame. Figure 8(a) shows the detection in a strong light construction scene. In comparison, the detection accuracy of the original YOLOv5 algorithm is much lower than our algorithm; Figure 8(b) shows the detection of small targets in the construction scene where steel bars are shielded. After observation, the original model missed a construction worker wearing a helmet who was behind the steel bars; Figure 8(c) shows the detection of targets with different sizes. The target size in close range is larger, and the target size in distant range is smaller. Our model has detected all the targets, while the original model missed the small targets in distant range and mistakenly detected steel pipes as two construction workers wearing safety helmets; Figure 8(d) shows the detection of small targets in a long-distance construction scene. The comparison shows that the original YOLOv5 model has missed detection of long-distance small

helmets, and our model has a better detection effect. The original YOLOv5 model misses this situation. There are many inspections, but our model performs better. It can be seen from the detection comparison in the abovementioned various construction scenarios that the improved YOLOv5 model is better for the detection of safety helmets in a complex operating environment.

## 5. Conclusions

This paper proposes an improved YOLOv5 helmet wearing detection method. First, use the *K*-means++ method to perform dimensional clustering on the dataset of the self-made construction operation scene; secondly, so as to capture more detailed information, the DWCA attention mechanism is combined with the backbone network. According to the comparison of the final experimental results, our model can obtain high detection accuracy, which can meet the detection accuracy of helmets in the current complex operating environment. In the future, we will explore ways to keep the model detection accuracy as much as possible while reducing the weight of the model.

## Data Availability

All data included in this study are available upon request by contact with the corresponding author.

## Conflicts of Interest

The authors declare that there is no conflict of interest regarding the publication of this paper.

## Acknowledgments

## References

[1] X. Chang and X. M. Liu, "Fault tree analysis of unreasonably wearing helmets for builders," *Journal of Jilin Jianzhu University*, vol. 35, no. 6, pp. 65–69, 2018.

[2] Z. Y. Wang, *Design and Implementation of Detection System of Wearing Helmets Based on Intelligent Video Surveillance*, Beijing University of Posts and Telecommunications, Beijing, China, 2018.

[3] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *in Proceedings of the IEEE conference on Computer*

*Vision and Pattern Recognition*, pp. 580–587, Columbus, OH, USA, 2014.

[4] R. Girshick, "Fast R-CNN. computer science," in *in Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 1440–1448, Santiago, Chile, 2015.

[5] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal network," *IEEE Transactions on Pattern Analysis & Machine Intelligence*, vol. 39, 2017.

[6] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: unified, real-time object detection," in *in Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, 2016.

[7] W. Liu, D. Anguelov, and D. Erhan, "Single shot multibox detector," in *in Proceedings of the ECCV 2016: Computer vision ECCV 2016*, vol. 9905, pp. 21–37, Springer, Amsterdam, The Netherlands, 2016.

[8] T. Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar, "Focal loss for dense object detection," in *in Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, Venice, Ital, 2017.

[9] J. Redmon and A. Farhadi, "YOLO9000: better, faster, stronger," in *in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6517–6525, IEEE, Honolulu, HI, USA, 2017.

[10] J. Redmon and A. Farhadi, *YOLOv3: an incremental improvement*, 2018, http://arxiv.org/abs/1804.02767.

[11] A. Bochkovskiy, C. Y. Wang, and H. Y. M. Liao, *Yolov4: Optimal speed and accuracy of object detection*, 2020, https://arxiv.org/abs/2004.10934.

[12] G. JOCHER, *Yolov5*, Code repository, 2020, https://github.com/ultralytics/yolov5,2020.

[13] W. Chen, Y. Qiao, and Y. Li, "Inception-SSD: an improved single shot detector for vehicle detection," *Journal of Ambient Intelligence and Humanized Computing*, vol. 11, pp. 1–7, 2020.

[14] Y. Tian, G. Yang, Z. Wang, E. Li, and Z. Liang, "Detection of apple lesions in orchards based on deep learning methods of cyclegan and yolov3-dense," *Journal of Sensors*, vol. 2019, 13 pages, 2019.

[15] L. Zhang, L. Lin, X. Liang, and K. He, "Is faster R-CNN doing well for pedestrian detection?," in *in European conference on computer vision*, pp. 443–457, Cham, 2016.

[16] Z. Zhong, L. Sun, and Q. Huo, "Improved localization accuracy by LocNet for faster R-CNN based text detection in natural scene images," *Pattern Recognition*, vol. 96, p. 106986, 2019.

[17] J. Zhang, S. Chen, S. Tian, W. N. Gong, and G. S. Cai, "A crowd counting framework combining with crowd location," *Journal of Advanced Transportation*, vol. 2021, 14 pages, 2021.

[18] B. J. Cheng, J. Zhang, and Y. Wang, "Research on medical knowledge graph for stroke," *Journal of Healthcare Engineering*, vol. 2021, 10 pages, 2021.

[19] A. Kelm, L. Laußat, A. Meins-Becker et al., "Mobile passive radio frequency identification (RFID) portal for automated and rapid control of personal protective equipment (PPE) on construction sites," *Automation in Construction*, vol. 36, pp. 38–52, 2013.

[20] X. H. Liu and X. N. Ye, "Skin color detection and Hu moments in helmet recognition research," *Journal of East China University of Science and Technology (Nature Science Edition)*, vol. 40, no. 3, pp. 365–370, 2014.

[21] A. H. M. Rubaiyat, T. T. Toma, and M. Kalantari-Khandani, "Automatic detection of helmet uses for construction safety," in *in Proceedings of the 2016 IEEE ACM International Conference on Web Intelligence Workshops (WIW)*, ACM, Omaha, NE, USA, 2016.

[22] RRV E Silva, "Helmet detection on motorcyclists using image descriptors and classifiers," in *in 2014 27th SIBGRAPI Conference on Graphics, Patterns and Images*, pp. 141–148, 2014.

[23] Q. R. Li, *A Research and Implementation of Safety-Helmet Video Detection System Based on Human Body Recognition*, University of Electronic Science and Technology of China, Chengdu, China, 2017.

[24] H. Wu and J. Zhao, "An intelligent vision-based approach for helmet identification for work safety," *Computers in Industry*, vol. 100, pp. 267–277, 2018.

[25] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *in Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7132–7141, 2018.

[26] Q. Hou, D. Zhou, and J. Feng, "Coordinate attention for efficient mobile network design," in *in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13713–13722, 2021.

[27] Q W, "ECA-Net: efficient channel attention for deep convolutional neural networks," in *in CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Seattle, WA, USA, 2020.

[28] S. Woo, J. Park, J. Y. Lee, and I. S. Kweon, "Cbam: convolutional block module," in *in Proceedings of the European conference on computer vision (ECCV)*, pp. 3–19, 2018.

[29] J. Park, S. Woo, J. Y. Lee, and I. S. Kweon, *Bam: Bottleneck Module*, 2018, https://arxiv.org/abs/1807.06514.

*Research Article*

# Investigation on the Lightning Location and Warning System Using Artificial Intelligence

**Tianru Shi,[1] Danhui Hu,[1] Xiang Ren,[1] Zeqi Huang,[1] Yaodong Zhang,[1] and Jianlan Yang [1] [2]**

[1]*State Grid Hubei Electric Power Co., Ltd., Research Institute, Wuhan 430077, China*
[2]*College of Electrical & Information Engineering, Hunan University, Changsha 410082, China*

Correspondence should be addressed to Jianlan Yang; yangjianlan@hnu.edu.cn

An in-depth study on a lighting location system is conducted in this paper. Firstly, the history and application of this system are summarized. The overall structure is detailed, including the detection principle of the lightning location, the orientation method, the detection circuit, the method of discriminating cloud flash and ground lightning signal, the error analysis, the guideline for station deployment, the preprocessing of the central station, and the function and structure of data server and user interface. The development of a lightning monitoring system in China is presented, and the construction of a new generation of a lightning location system in the Hubei Province power grid is introduced. Through the collection of measured data, the performance of the lightning location system in the lightning accident inspection rate, lightning location, and lightning situation statistics are analyzed. Artificial intelligence algorithms are applied in the lightning warning system. The new system has a high predicting accuracy.

## 1. Introduction

Lightning is a high-intensity electromagnetic pulse phenomenon that frequently occurs in nature [1, 2]. As its impact is huge, it has received extensive attention from many industry fields, such as meteorology, aerospace, aviation, electric power, and petroleum. Among them, the power grid is susceptible to lightning due to its wide-area distribution and a geometric scale of thousands of kilometers [3]. It is estimated that the number of trips on high-voltage transmission lines caused by lightning accounts for 40% to 70% in China. Lightning is an important factor that seriously affects the safe operation of the power grid.

The observation of accurate lightning parameters is the basis for lightning protection [4–6]. The key to detect lightning is the lightning location. It refers to automatic detection equipment, which uses the characteristics of sound, light, and electromagnetic wave radiated by the lightning return strike to remotely measure the discharge parameters [7]. Several methods for detecting lighting have been proposed including acoustic, optical, and electromagnetic field methods [8–10]. The modern lightning location system started in 1976. Krider used a single-chip technique to successfully transform the original double-cathode oscilloscope lightning detector into an intelligent magnetic direction lightning location system, which effectively improved the accuracy of lightning angle measurement. In the early 1980s, the emergence and application of cloud-to-ground lightning waveform identification technology enabled the detection efficiency to reach up to 90%. Since then, all developed countries and regions in the world have begun to install lightning monitoring and location networks, e.g., the United States, Canada, Japan, France, and Germany. In the 1990s, due to the use of the global positioning system (GPS), lightning monitoring added GPS clocks based on a direction finding system to form a time difference direction hybrid system. Meanwhile, the use of digital signal processing (DSP) and integrated technology to perform correlation analysis and position processing on the waveform greatly improves the prediction performance. Currently, there are

more than 60 lightning location system networks worldwide that employ commercial instrumentation operating in the very low frequency/low frequency range.

The lightning location system has been widely used in the aerospace, disaster reduction, and prevention and power industries, especially in the global power system. Over 40 countries in the world have installed lightning monitoring systems. Over the past decades, with the development of science and technology and the continuous improvement of itself, the location accuracy and detection efficiency of the lightning location system have been greatly improved. The current lightning location system uses GPS satellite positioning technology, satellite communications, geographic information system (GIS), and other high-tech technologies to form a real-time dynamic multipurpose large-scale information system.

Extensive research has been done on the statistical analysis of lightning parameters [11, 12]. Chen et al. proposed a grid method based on the huge data accumulated by the lightning location systems and used data mining technology to analyze the temporal and spatial distribution of lightning [13]. An improved grid method using the two parameters of a grid area and observation range is developed to further improve the accuracy. Many scholars have analyzed the influence of region and climate change on lightning parameters [14–17]. The lightning parameters of crucial transmission line corridors are analyzed, reducing the error in the area where the line is located. The influence of lightning current amplitude probabilities on the trip rate of the transmission line is investigated, considering different topography and landforms.

## 2. Sensing Principle of the Lightning Location System

*2.1. Structure of the Lightning Location Station.* The detection station is composed of an electromagnetic field antenna, lightning waveform recognition and processing unit, high-precision crystal oscillator and GPS clock unit, communication, power supply, and protection unit [18]. It measures and outputs the characteristic quantities of ground-flash waves: the time, direction, and relative signal strength of each return strike, and sends the original measurement data to the central station in real time. Each part of the detection station has a unique function [19]. The GPS antenna is mainly used to receive a GPS synchronization signal. The electromagnetic antenna is composed of two vertical orthogonal frame antennas for receiving electromagnetic wave signals. The circuit structure of the detection station is shown in Figure 1.

The GPS clock unit is used to provide the required high-precision synchronization time signal. The lightning waveform delay processing circuit and the overrange timing circuit are specially designed to improve the accuracy of electromagnetic wave signal detection and lightning strike location. Meanwhile, a drift calibration is developed to avoid errors caused by the drift of the GPS clock crystal oscillator affected by temperature rise. These devices would improve both the detection efficiency and accuracy.

Each detection station of the lightning location system is equipped with a time difference clock, which is composed of a high-stability constant temperature crystal oscillator, a GPS antenna, and a clock board, as given in Figure 2. The clock consists of a highly stable crystal oscillator. GPS can receive a high-precision second pulse time signal and use this signal to correct the clock. The accuracy and reliability of the revised clock are greatly improved. The quality of the GPS receiving board and antenna is reliable, and the time error is less than $0.5\,\mu$s.

*2.2. Directional Location Principle of the Lightning Location System.* The lightning is accompanied by strong light, sound, and electromagnetic radiation. Among them, the most suitable signal for detecting in a relatively large range is electromagnetic radiation. The electromagnetic radiation of thunder and lightning mainly spreads along the earth surface through low frequency and very low frequency. The range is several hundred kilometers and sometimes can be wider, which is determined by the discharge energy. When extracting signals, the lightning location system activates multiple detection stations to measure the electromagnetic radiation generated by lightning, eliminate the signal of cloud flashes, and identify ground-to-ground flashes. The antenna can measure signals with a frequency ranging from 1 kHz to 1 MHz. Through the electronic circuit, the ground flashing signal is identified and the peak value of each return wave is sampled. The orientation method is the most widely used in the directional location principle.

It employs the magnetic field intensity to obtain the azimuth of the lightning strike point relative to the detection station. In order to detect the radiation waves of the ground flash magnetic field, as depicted in Figure 3, the two orthogonal antennas are in east-west and north-south directions, respectively. If a lightning strike occurs on A, the orthogonal antenna can receive two magnetic signals of different strengths. Assuming that the measured magnetic field strengths in the east-west and north-south directions are $H_{WE}$ and $H_{NS}$, respectively, the direction angle of the lightning strike point can be calculated as follows:

$$\tan \alpha = \frac{H_{NS}}{H_{WE}}. \tag{1}$$

The angles measured by two detection stations are shown in Figure 4. According to the angle relationship of the triangle, the azimuth of point A is expressed as

$$B = B_1 + \alpha_{1P}^B \left(1 + \eta_1^2\right)\left(1 - \frac{3}{2}\eta_1^2 \tan B\,\alpha_{1P}^B\right), \tag{2}$$

$$L = L_1 + \sin^{-1}\left[\frac{\sin \alpha_{1P} \sin \beta_{1p}}{\cos \left(\beta_1 + \alpha_{1P}^B\right)}\right]. \tag{3}$$

In Figure 4, $A$ is the location of the lightning strike, TDF1 and TDF2 are two different detection stations, the coordinates of TDF1 are $(B_1, L_1)$, the coordinates of TDF2
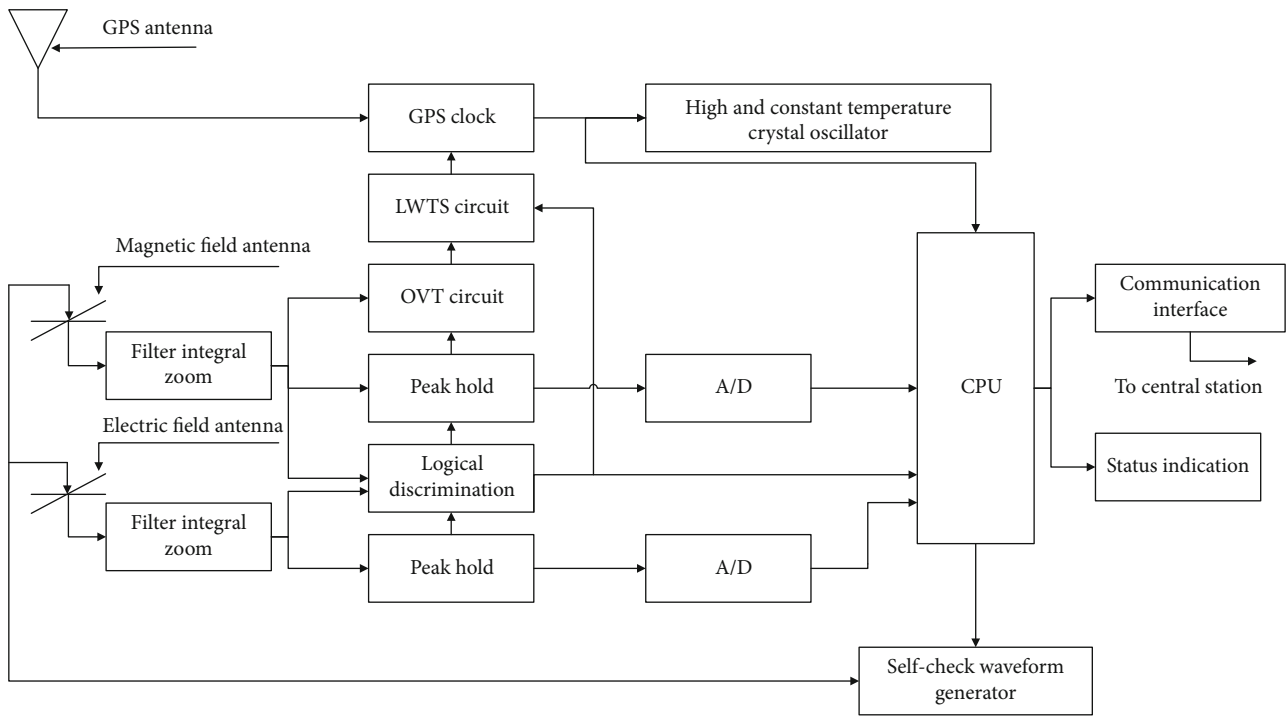
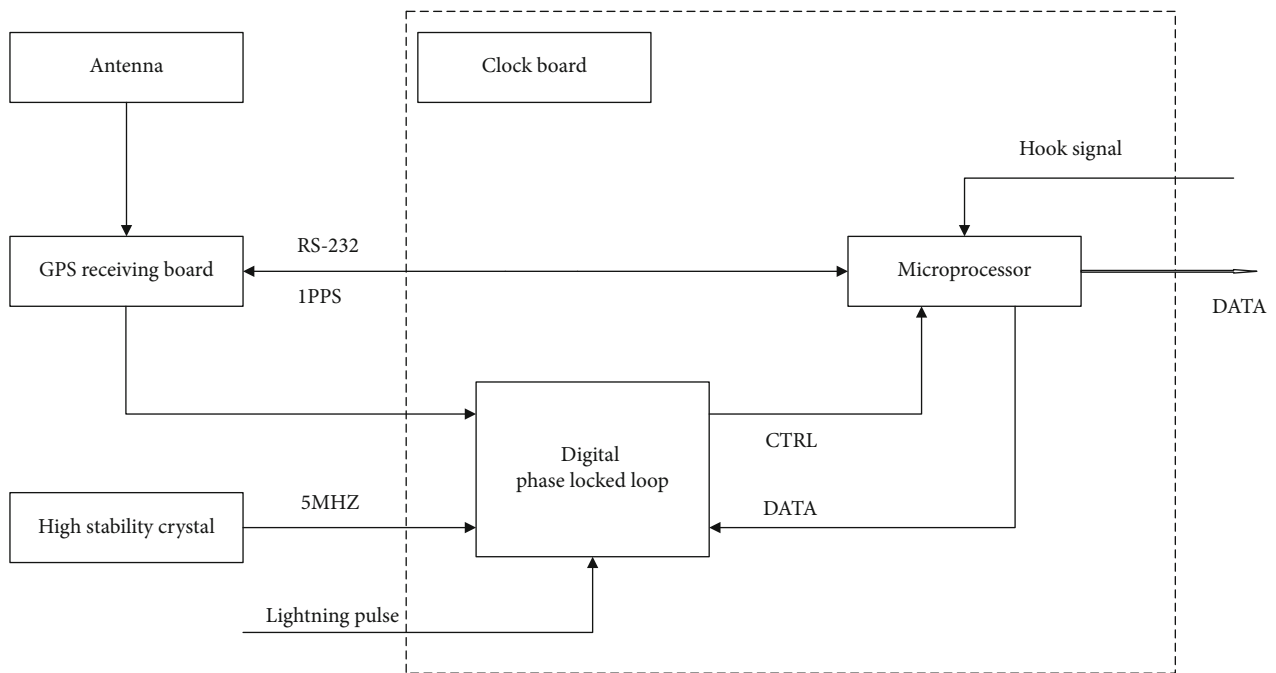FIGURE 1: Principle structure of the detection station circuit.
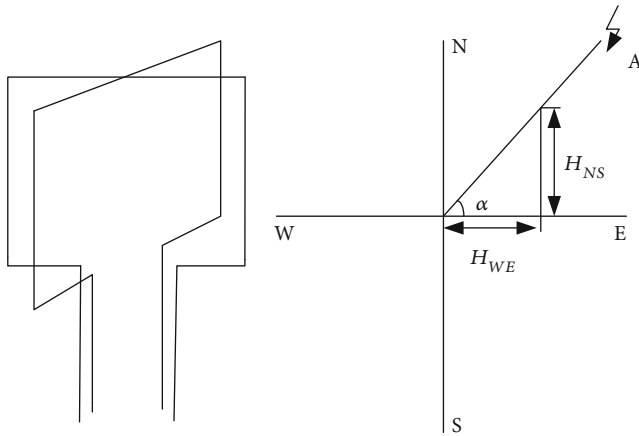


FIGURE 2: GPS clock unit.

Figure 3: Orientation method.



Figure 4: Obtaining the location of lightning strikes by the orientation method.

are $(B_2, L_2)$, $B$ is the latitude, $L$ is the longitude, and $\beta_{1P}$ and $\beta_{2p}$ are the azimuth angle.

*2.3. Structure of the Lightning Location System.* The lightning location system uses a browser/server mode to publish the lightning location information on the web, and its structure is given in Figure 5.

# 3. Application of the Lightning Location System in the Hubei Power Grid

*3.1. Hubei Power Grid Lightning Location System.* The lightning location system of the Hubei power grid was established in 1998, with six base stations located in Puqi, Huangshi, Xiaogan, Jingmen, Jingzhou, and Yichang, respectively. In 2000, three stations were added in Suizhou, Xiangfan, and Shiyan, and from 2002 to 2005, five other stations were built in Enshi, Qianjiang, Macheng, Shishou, and Badong. After 13 years of construction and operation, a lightning detection network consisting of 14 lightning detection stations has been developed. Based on the principles of time difference and direction location as well as modern communication technology, the automatic monitor with full real-time function in most areas in Hubei Province has been realized.

More than 1930 transmission lines of various voltage grades in Hubei Province were input into the system. It can be widely used in investigation of the line fault point, lightning parameter statistics, and lightning accident analysis. It reduces the loss of power failure caused by lightning strike and the labor intensity of searching the lightning strike point. The safety of the power system can be ensured.

*3.2. Frequency of Lightning Activity in the Hubei Power Grid.* From 2014 to 2018, the lightning activity was frequent in Hubei, and the average lightning density in the whole province was between 1.6 and 2.3 times/km². In 2018, lightning activity was the most intense, and the density reached 2.29 times/km². The lightning activity in Hubei in recent 5 years is listed in Table 1.

Taking 2018 as an example, there were 10 times of the lightning trip in the Hubei power grid for voltage levels of 500 kV and above, including once in March, once in April, thrice in June, quartic in July, and once in August. The most lightning trips were in June and July. The time distribution characteristics of lightning tripping are shown in Figure 6.

In 2018, overhead transmission lines of 500 kV and above were tripped 10 times due to lightning strikes, an increase of 6 times over the same period last year. The lightning density was 2.29 times/km², 1.44 times higher than that of the same period last year, which was the highest in the past five years. The results of ground lightning density and lightning trip times are shown in Figure 7. The number of lightning strikes increases with the rise of lightning density.

Compared with the lightning density map in 2017, the lightning activity is changed significantly in 2018. In 2017, the areas with high lightning density were mainly concentrated in the southern part of Jingzhou, Huangshi, Huanggang, and Xianning. In 2018, a large area of C1 level regions appeared in central Yichang, Jingzhou, Jingmen, Shiyan, and Wuhan, as shown in Figure 8.

220 kV Hubei Qiaoshun line fault analysis: on July 31, 2012, 220 kV Qiaoshun line was tripped and the reclose was successful. There were obvious discharge traces on both the left front and the right back of ground wire of the #017 pole. In the large side of the middle phase (phase C), the internal string porcelain insulators and connecting fittings had obvious discharge traces. The other poles passed ground and pole climbing inspection, and no abnormality was found. During the fault inspection, it was found that around 15:30-22:00 of July 31, strong lightning and heavy rain began to appear in the area. The lightning information query system is shown in Table 2: within 5 minutes before and after 21:21 on July 31, there were 4 lightning strikes along the second circuit of 220 kV Qiaoshun, and the lightning current amplitude was from -3.4 kA to -18.5 kA, distributed near #016~#020 towers.

After inspecting the transmission tower, the faulty section was located in the vegetable garden in Baijiawan Village, Xiangyang. The line right-of-way was in good condition, without tree barriers, external damage traces, industrial pollution sources, and fouling on insulator strings. There was no strong wind or abnormal local air flow during the fault, and possibilities caused by wind deviation, external damage, pollution, and tree barriers were ruled out, and it was judged

Figure 5: Structure of the lightning location system.

Table 1: Statistics of lightning parameters in Hubei from 2014 to 2018.

| Years | Thunder day | Density (times/km$^2$) | Positive number | Negative number | Positive density (times/km$^2$) | Negative density (times/km$^2$) |
|---|---|---|---|---|---|---|
| 2014 | 293 | 1.7 | 60419 | 257572 | 0.3 | 1.4 |
| 2015 | 320 | 2.1 | 93267 | 303948 | 0.5 | 1.6 |
| 2016 | 256 | 1.7 | 82627 | 236089 | 0.4 | 1.3 |
| 2017 | 213 | 1.6 | 59874 | 201418 | 0.3 | 1.3 |
| 2018 | 247 | 2.3 | 103148 | 321648 | 0.6 | 1.7 |



Figure 6: Lightning frequency and lightning tripping frequency of the Hubei power grid.

Figure 7: Tripping data of 500 kV and above voltage level lines in the Hubei power grid.

as a lightning fault, line fault scene photo as shown in Figure 9.

## 4. Application of Artificial Intelligence in the Lightning Warning System

*4.1. Principle of Lightning Warning.* The key of the lightning warning system is the thunderstorm forecast model, which can send out the lightning warning information in advance and effectively avoid the damage caused by lightning strike to the staff and equipment in the protected area. Thunderstorm prediction is based on a subjective and objective prediction algorithm. Subjective forecast takes observations from Doppler weather radar and combines them with other meteorological satellite cloud maps. The objective algorithms include radar echo or cloud image extrapolation and severe convective weather recognition [20]. However, the success rate of forecast and warning is quite limited. The advancement of big data and artificial intelligence technology, massive historical lightning data, and other meteorological monitoring data are processed and modeled through a deep learning method. As a result, a more accurate local lightning warning model is obtained. In practice, most warning methods are based on the data of the lightning location and atmospheric electric field. The electric field meter is used to determine the lightning probability by measuring the intensity and change trend of the atmospheric electric field, so that different alarm levels can be determined by setting the proper threshold. Under different alarm levels, contingency plans are made to implement actions such as stopping sending and receiving oil to achieve the purpose of active lightning protection.

*4.2. Lightning Warning Model.* The data of the lightning warning model comes from a 3D lightning locator and other meteorological observation data such as meteorological radar cloud image. The 3D lightning locator is a high-precision system for locating lightning strikes and lightning

within clouds. The average detection accuracy is about 300 m, and the detection efficiency is up to 95%. The lightning locator is used to collect the lightning location data within the region, including the time, the location (latitude and longitude information), the height of lightning from the ground, and other attribute information. The XGB (Extreme Gradient Boosting) algorithm is used to build a scoring model for the occurrence probability of lightning strike in the protected area to solve the problem of 0-2-hour approaching warning. The model structure is shown in Figure 10.

$A_n$ is the $n$ basic feature data. An XGB lightning prediction model is built based on the data, and the output is the probability of lightning strike in the surveillance area (the probability value between 0 and 1).

*4.3. Feature Extraction.* After obtaining the relevant lightning data, it is necessary to extract the data features. According to the periodicity, instantaneity, and mobility of lightning, the following targeted features were extracted from the collected data sources:

(1) Thunderstorm proximity: the distance between the thunderstorm cluster and the protection point

(2) Total number of lightning at close range

(3) Thunderstorm approaching speed in the protected area: the speed of the nearest thunderstorm group approaching the protected area

(4) The increasing trend of lightning strikes: the increasing trend of thunderstorm cluster energy in the presentation window

*4.4. Model Training and Evaluation.* After completing the selection of characteristic values required for model training, the next step is to use a characteristic variable to train the model and get the best parameters. XGB is used for data

Hubei
Cloud-to-ground lighting flash density (0.05 201801–09)

<1.313  1.313–2.06  2.06–3.018  3.018–4.255  >4.255

(a)



Hubei
Cloud-to-ground lighting flash density (0.05 201701–06)

<.184  .184–.375  .375–.683  .683–1.197  >1.197

(b)

FIGURE 8: Lightning density distribution in Hubei: (a) 2018 year and (b) 2017 year.

TABLE 2: Query results of the lightning location system of the Hubei power grid.

| Number | Time | Longitude | Latitude | Current (kA) | Reply |
|--------|------|-----------|----------|--------------|-------|
| 1 | 21:18:54 | 112.0387 | 32.0604 | -6.3 | 1 |
| 2 | 21:20:56 | 112.0521 | 32.0515 | -14.5 | 1 |
| 3 | 21:21:48 | 112.0269 | 32.0737 | -18.5 | 1 |
| 4 | 21:23:26 | 112.0403 | 32.0732 | -3.4 | -1 |

classification with machine learning integration, and the XGB algorithm structure is shown in Figure 11.

XGB is based on the gradient lifting decision tree (GDBT), which reduces the complexity of the model and avoids overfitting by adding a regularization term into the objective function. The objective function is

$$\text{Obj} = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k) + c, \tag{4}$$

FIGURE 9: Photo of the line fault scene.



FIGURE 10: Lighting warning model.

$$\Omega(f_k) = \gamma T + \frac{1}{2}\lambda\|\omega\|^2 = \gamma T + \frac{1}{2}\lambda\sum_{j=1}^{T}\omega_j^2, \qquad (5)$$

where $\hat{y}_i$ is the predicted value, $y_i$ is the actual value, and $l$ is the loss function, which shows the residual value between the predicted value and the real value. $\Omega(f_k)$ shows the complexity of the model. $f_k$ is the $k^{\text{th}}$ decision tree. $\gamma$ and $\lambda$ represent the penalty coefficient of the model. $T$ and $\omega$ are the number and the weight of leaves for the $k^{\text{th}}$ tree, respectively.

$c$ is a constant. It is relatively simple to solve the optimal solution for the general least square loss. However, when it is replaced by other loss functions [15], the solution process will become more complex. To solve this problem, the XGB algorithm performs second-order Taylor expansion on this basis. Assume that the $t^{\text{th}}$ loss function is defined as

$$\text{Obj}^{(t)} = \sum_i l\left(y\wedge_i^{(t-1)}, y_i + f_t(x_i)\right) + \Omega(f_t). \qquad (6)$$

The second-order Taylor expansion of formula (6) is carried out, and formula (7) is simplified by removing the constant term.

$$\text{Obj}^{(t)} = \sum_{i=1}^{n}\left[g_i f_t(x_i) + \frac{1}{2}h_i f_t(x_i)^2\right] + \Omega(f_t). \qquad (7)$$

Here,

$$\begin{cases} g_i = \partial_{\hat{y}_i(t-1)} l\left(y_i, y\wedge_i^{(t-1)}\right), \\ h_i = \partial^2_{\hat{y}_i(t-1)} l\left(y_i, y\wedge_i^{(t-1)}\right). \end{cases} \qquad (8)$$

FIGURE 11: XGB algorithm structure.

The objective function is

$$\mathrm{Obj}^{(t)} = \sum_{i=1}^{n} \left[ g_i f_t(x_i) + \frac{1}{2} h_i f_t(x_i)^2 \right] + \gamma \mathrm{T} + \|\omega\|^2$$

$$= \sum_{j=1}^{n} \left[ \left( \sum_{i \in I_j} g_i \right) \omega_j + \frac{1}{2} \left( \sum_{i \in I_i} h_i + \lambda \right) \omega_j^2 \right] + \gamma T, \quad (9)$$

where $I_j = \{i | q(x_i) = j\}$ represents the $j^{\text{th}}$ group of leaf nodes. At this time, the objective function is transformed into the problem of seeking the minimum of the element quadratic equation on $\omega_j$. Assume that the tree structure is fixed. The optimal weight of leaf node $j$ is in

$$\omega_j = \frac{G_i}{H_j + \lambda}. \quad (10)$$

Then, the objective function is expressed as

$$\mathrm{Obj}^* = -\frac{1}{2} \sum_{j=1}^{T} \frac{G_j^2}{H_j + \lambda} + \lambda T. \quad (11)$$

Here,

$$\begin{cases} G_j = \sum_{i \in I_j} g_i, \\ H_j = \sum_{i \in I_j} h_i. \end{cases} \quad (12)$$

Obj $*$ represents the structure score of the regression tree. The smaller the value, the better the structure. Since the structure of all the trees cannot be listed, the greedy algorithm is used for the division of subtrees. Each attempt is made to add a division point to the existing leaf node. The feasible division points are listed, and the division point with

the smallest objective function and the largest gain is selected [16]. The gain formula is given in

$$\text{Gain} = \frac{1}{2}\left[\frac{G_L{}^2}{H_L + \lambda} + \frac{G_R{}^2}{H_R + \lambda} - \frac{G_L + G_R}{H_L + H_R + \lambda}\right] - \gamma. \quad (13)$$

After the XGB model integrates several regression trees, the nodes of each tree are doing feature splitting. The number of times a feature is selected as a split feature can be used as the importance.

After the model training, the model evaluation index is the AUC (Area Under ROC Curve) value and ROC (Receiver Operating Characteristic) curve. Based on the data of the lightning location and other meteorological observation, the AUC value of the model reached 0.95 and the best performance was 1, indicating that the model has a good classification effect. AUC is the area of the ROC curve, which is used to evaluate the quality of the dichotomy system.

$$\text{AUC} = \int_{-\infty}^{+\infty} y(t)dx(t), \quad (14)$$

where $x$ and $y$ are the false-positive rate and true-positive rate, respectively, and also the horizontal and vertical coordinates of the ROC curve.

## 5. Conclusion

The conclusions are drawn as follows.

(i) The charge distribution of thunderstorm, the density and pressure level of air, and the topography and geological conditions are various in different climatic and geographical areas. In order to improve the accuracy and performance evaluation of the lightning location system, it is necessary to strengthen the observation of natural lightning in different areas

(ii) The cause of a typical lightning fault for a 500 kV transmission line is analyzed in this paper. Several lightning protection measures are introduced. Currently, the lightning protection methods available for transmission lines cannot completely eliminate the impact of lightning. In recent years, extreme weather occurred frequently and there are great uncertainties in lightning activity. It is necessary to further study climate change and lightning behavior

(iii) The lightning warning system can send out the lightning warning information in advance and effectively avoid the damage caused by lightning strike to the staff and equipment in the protected area. The use of artificial intelligence algorithms in the lightning warning system can improve the predicted accuracy

## Data Availability

The lightning sensing, location, and warning data used to support the findings of this study were supplied by Hubei Electric Power Company Research Institute under license and so cannot be made freely available. Requests for access to these data should be made to Tianru Shi (739601030@qq.com).

## Conflicts of Interest

The authors declare no conflict of interest.

## Authors' Contributions

T. Shi and D. Hu proposed the concepts and ideas. X. Ren analyzed the results. Y. Zhang and J. Yang wrote this paper and revised the contents of this manuscript.

## Acknowledgments

## References

[1] L. Cai, J. Wang, Q. Li, and Y. Fan, "The Foshan Total Lightning Location System in China and its initial result, presented at the 10th Asia-Pacific," in *Conf. Lightning Protection (APL 2017)*, Krabi, Thailand, 2017.

[2] L. Cai, X. Zou, J. Wang, Q. Li, M. Zhou, and Y. Fan, "The Foshan Total Lightning Location System in China and its initial operation results," *Atmosphere*, vol. 10, no. 3, p. 149, 2019.

[3] R. K. Said, U. Inan, and K. Cummins, "Long-range lightning geolocation using a VLF radio atmospheric waveform bank," *Journal of Geophysical Research*, vol. 115, no. D23, 2010.

[4] P. W. Casper and R. B. Bent, "Results of the LPATS USA national lightning detection and tracking system for the 1991 lightning season," in *21st International Conference on Lightning Protection*, Berlin 1992, 1992.

[5] V. Cooray, Ed., *Lightning Protection, IET Power and Energy Series 58*, The Institution of Engineering and Technology, London, UK, 2010.

[6] IEC technical committee-88, *IEC 61400-24 Wind Turbines-Part 24: Lightning Protection*, International Standard IEC, Geneva, CH, 2010.

[7] A. C. L. Lee, "An experimental study of the remote location of lightning flashes using a VLF arrival time difference technique," *Quarterly Journal of the Royal Meteorological Society*, vol. 112, no. 471, pp. 203–229, 1986.

[8] P. M. Lo, *A Simplified Model for Lightning Exposure of Wind Turbines, Master's Thesis*, McGill University, Montreal, 2008.

[9] V. Cooray, U. Kumar, F. Rachidi, and C. A. Nucci, "On the possible variation of the lightning striking distance as assumed in the IEC lightning protection standard as a function of structure height," *Electric Power Systems Research*, vol. 113, pp. 79–87, 2014.

[10] R. Rodrigues, V. Mendes, and J. Catalão, "Estimation of lightning vulnerability points on wind power plants using the

rolling sphere method," *Journal of Electrostatics*, vol. 67, no. 5, pp. 774–780, 2009.

[11] M. Becerra, M. Long, W. Schulz, and R. Thottappillil, "On the estimation of the lightning incidence to offshore wind farms," *Electric Power Systems Research*, vol. 157, pp. 211–226, 2018.

[12] N. Yang, Q. Zhang, W. Hou, and Y. Wen, "Analysis of the lightning-attractive radius for wind turbines considering the developing process of positive attachment leader," *Journal of Geophysical Research: Atmospheres*, vol. 122, no. 6, pp. 3481–3491, 2017.

[13] J.-H. Chen, Q. Zhang, W.-X. Feng, and Y.-H. Fang, "Lightning location system and lightning detection network of China power grid (in Chinese)," *High Voltage Eng*, vol. 34, pp. 425–431, 2008.

[14] J. Jerauld, V. A. Rakov, M. A. Uman et al., "An evaluation of the performance characteristics of the U.S. National Lightning Detection Network in Florida using rocket-triggered lightning," *Journal of Geophysical Research*, vol. 110, no. D19, article D19106, 2005.

[15] A. Mäkelä, *Thunderstorm Climatology and Lightning Location Applications in Northern Europe*, PhD Thesis, University of Helsinki, Finland, 2011.

[16] A. Mäkelä, T. J. Tuomi, and J. Haapalainen, "A decade of high-latitude lightning location: effects of the evolving location network in Finland," *Journal of Geophysical Research*, vol. 115, no. D21, 2010.

[17] A. T. Pessi, S. Businger, K. L. Cummins, N. W. S. Demetriades, M. Murphy, and B. Pifer, "Development of a long-range lightning detection network for the Pacific: construction, calibration, and performance," *Journal of Atmospheric and Oceanic Technology*, vol. 26, no. 2, pp. 145–166, 2009.

[18] W. Schulz and G. Diendorfer, *Detection Efficiency and Site Errors of Lightning Location Systems*, USA, International Lightning Detection Conference, Tucson Arizona, 1996.

[19] V. Cooray and R. E. Orville, "The effects of variation of current amplitude, current risetime, and return stroke velocity along the return stroke channel on the electromagnetic fields generated by return strokes," *Journal of Geophysical Research*, vol. 95, no. D11, pp. 18617–18630, 1990.

[20] W. Schulz and M. M. F. Saba, "First results of correlated lightning video images and electric field measurements in Austria," in *X International Symposium on Lightning Protection (SIPDA)*, Curitiba, Brazil, 2009.

Research Article

# An Optimization Method for the Layout of Soil Humidity Sensors Based on Compressed Sensing

**Yunsong Jia, Xueyun Tian, Xin Chen, and Xiang Li** [iD]

*College of Information and Electrical Engineering, China Agricultural University, Beijing 100083, China*

Correspondence should be addressed to Xiang Li; cqlixiang@cau.edu.cn

In the farmland Internet of Things, to achieve precise control of production, it is necessary to obtain more data support, which requires the deployment of many sensors, and this will inevitably bring about high investment and high-cost problems. This paper mainly studies the optimization of sensor placement in the agricultural field. Through compressed sensing and algorithm optimization, the number of sensors used is reduced and the cost is reduced on the premise of ensuring the effect. At present, there are many mature sensor layout optimization methods, but these methods will have incomplete parameters due to experimental conditions and environmental factors. They are more suitable for structural health monitoring and lack research in agricultural applications. Considering that the sensor layout optimization can be converted into the characteristics of image compression selection and the compression effect of the compressed sensing theory is better, therefore, this paper proposes a sensor layout optimization method based on compressed sensing. Due to the structural characteristics of the existing measurement matrix in the compressed sensing theory, the specific position distribution of the optimized sensor layout cannot be obtained directly. This paper improves the existing sparse random measurement matrix to determine the number of sensors required for a given region and the function of the specific location of each sensor. The experimental results show that soil moisture can be measured with a small error of 0.91 by using 1/3 of the original sensor number. The result of data reconstruction using 1/6 of the original sensor is average, and the average error is up to 1.68, which is suitable for the environment with small data fluctuation.

## 1. Introduction

Precise irrigation and irrigation automation are the inexorable trends of the development of modern agriculture. The accurate measurement of soil moisture content is the basis of precision irrigation and irrigation automation [1]. There are mainly 3 methods for soil moisture measurement, namely, manual soil sampling and drying, remote radar sensing, and acquisition by sensors. Manual soil sampling and drying are accurate but require the manual collection of soil samples at multiple locations. High in cost and weak in timeliness [2], this method is difficult to adapt to modern agricultural production. Remote radar sensing is the use of microwave radar on the satellite to measure the water on the soil surface. However, the measurement result is too coarse-grained [3] to guide fine agriculture production. Soil moisture sensors can quickly and accurately measure the

same point in the soil and are widely used in precision irrigation.

When soil humidity sensors are applied to measure the soil moisture content, the more soil humidity sensors are buried in a profile, the more precisely the soil moisture measurement will be [4]. However, as the number of sensors increases, so does the cost of production systems, and there is a contradiction between cost and data accuracy. To save the agricultural irrigation system cost and improve the efficiency of the state estimators, the paper proposes a method for soil humidity sensor layout based on compressed sensing, aiming at reducing the number of sensors as much as possible on the premise of accuracy.

Compressed sensing, also known as compressive sampling or sparse sampling, is a technique for finding sparse solutions of underdetermined linear systems. According to this theory, if the signal is sparse, it can be reconstructed

and recovered from sampling points much lower than the sampling theorem requires [5]. Compressed sensing is used in electronic engineering, especially signal processing, to obtain and reconstruct sparse or compressible signals. This method takes advantage of the sparsity of the signal. Compared with the Nyquist theory, this method can recover the entire desired signal from fewer measured values. Simply put, the process of data compression is completed during the sampling process. In the process of signal sampling, compressed sensing uses a few sampling points to achieve the same effect as full sampling. In the compressed sensing theory, the original signal can be accurately reconstructed under the condition of few measuring points. Based on this feature, an optimization method based on compressed sensing of soil humidity sensor layout is proposed in the paper. Firstly, the sensors are densely placed in farmland soil by which original soil humidity data is obtained, and then Fourier Transform is applied for sparse representation of the data; secondly, the sparse presentation data is operated through the improved sparse random measurement matrix and an observed value is obtained; finally, through the reconstruction algorithm, the reconstructed signals are obtained. Through comparison of the three kinds of sensor layout optimization strategies, it is found that 1/3 of the original number of sensors can measure the soil moisture with a minor error.

This paper theoretically ameliorates the existing sparse random measuring matrix, proposing a soil humidity sensor layout optimization method based on compressed sensing. The implemented functions are to quantify the required sensors in a given area and place them at specific locations. Effects are achieved that with fewer sensors, the whole farmland soil temperature distribution is measured, which reduces the costs effectively while increasing the efficiency of information processing of the system.

## 2. Related Work

In farmland IoT, strengthening the research on soil quality protection and management and realizing the intelligent management of farmland protection are the key to guarantee the safety of agricultural products [6]. In the field of farmland irrigation IoT, many scholars have been solving the problems on the hardware level. For example, Feng, to lower the irrigation water cost, achieved water saving by intelligent irrigation and raising irrigation water efficiency through embedded control technology [7]. Liu, on the possible time out and cross-restriction problems, proposed the farmland data collection mechanism to guarantee the reliable transmission of data [8]. Liu and Yang proposed a network management project of network topology management, location management, energy management, and fault management, referring to the features of node power and limited processing energy in sensor networks, to realize the remote management of the sensor monitoring network and the effective detection of the farmland environment for the users [9]. Singh and Saikia proposed an irrigation control system based on Arduino. The system collects and receives data through Arduino and uploads it to a designated interactive

website, on which the real-time soil status factor and the standard value of different factors required by crops are shown [10]. All these research studies ensure the instantaneity, accuracy, and reliability of the soil humidity data on the hardware level. It alleviates the noise folding phenomenon of compressed sensing [11], making compressed sensing a more effective choice in the optimization of agricultural sensors.

Sensor layout optimization plays an important role in different fields. In structural dynamics, a good sensor layout could recognize accurately the model parameter of a structure and ascertain the damage degree of the structure [12]. In direct kinematics, the location of the sensors will influence the calculation complexity, the accuracy of position sensing, and the reliability of the system [13]. In a network warning system, the position of sensors affects the effectiveness of warnings [14]. In thermology, placing a temperature sensor in an optimal position helps to accurately and in real time predict thermally induced deformation at a particular location [15].

Sensor layout optimization based on the model analysis method is a methodology by which the layout strategy is obtained through optimum analysis based on establishing a finite element model and setting an optimization target. The sensor layout optimization method based on the model analysis method was first proposed by Kammer as the effective influence method [16]; that is, a sensor placement program is to be obtained by maximizing the spatial independence and signal intensity of the target finite element model. Then, Heo et al. proposed kinetic energy to place the sensors [17]. Based on the EI method and KE method, Wu et al. proposed the effective independence driving-point residue method, improving the spatial independence and element strain energy of the above two methods [18]. Mukherjee et al. applied a reweighing method replacing repeated function simulation to estimate the expected influence value and proposed a mode analyzing method for sensor placement for nonlinear uncertain systems [19]. Modeling error was caused inevitably in the course of establishing a finite element model with the abovementioned sensor layout optimization method, and modeling tends to be trapped into the local optimum.

To avoid modeling error, Krause et al. proposed a sensor node placement method driven by data [20]. Guestrin et al. proposed to place the sensor based on rules of mutual information [21]. Xu and Choi used noise measurement and a mobile self-adaption anisotropy space-time Gaussian process, which enables nonparametric prediction toward a given space-time phenomenon [22]. The abovementioned methods presume that the space random process is Gaussian distribution; however, soil humidity in farmland does not completely follow Gaussian distribution, and thus, it does not work very well.

Compressed sensing is a technology that could be used to obtain and reconstruct sparse signals and does not depend on the Gaussian distribution of data. Put forward in 2006 by Donoho [23], now, compressed sensing has been extensively used in fields like wireless communication. Compressed sensing is extensively applied in sensor node information collection [24–26], but few studies have been made in the

sensor. In this paper, soil water distribution is treated as a two-dimensional image. The compressed sensory theory has a good effect and has been proved by a large amount of theoretical proofs in the field of image compression. Therefore, this method has general conditions.

## 3. Method

*3.1. Compressed Sensing Theory.* Compressed sensing of the signal mainly includes three steps [23]: the first step is sparse representation, which converts the original signal into a sparse signal on another dimension; the second step is to reduce the dimension of the observation matrix to minimize the information loss of the original signal; and the third step is to design a reconstruction algorithm to recover the $N$-dimensional original signal from the $m$-dimensional sampled signal ($m < n$). Figure 1 shows the compressed sensing framework.

*3.1.1. Sparse Representation.* Based on the signal sparse decomposition theory, $N$-dimensional discrete real value signal $x = (x_1, x_2, \cdots, x_n)$ could be denoted as a linear combination of a group of uncorrelated bases $\psi_i (i = 1, 2, 3, \cdots, N)$.

$$x = \sum_{i=1}^{N} \psi_i \alpha_i = \psi\alpha. \tag{1}$$

In the formula, $\psi = [\psi_1, \psi_2, \cdots, \psi_N]$ is the basis matrix $N \times N$. If there are only $K$ nonzero coefficients in $\alpha$, then $x$ is called the $K$ sparse signal in basis matrix $\psi$. If the conversion coefficient of the signal decays to zero exponentially with the order sorted, the signal is compressible.

*3.1.2. Measurement Matrix.* Suppose signal $x$, with length as $N$, is reflected by a group of unit vectors $\Phi = [\Phi_1, \Phi_2, \cdots, \Phi_N]$ and the measured value $y \in R^M (M \ll N)$ is obtained. This process could be shown as

$$y = \Phi x. \tag{2}$$

We put formula (1) into formula (2) and obtain

$$y = \Phi x = \Phi\psi\alpha = \Theta\alpha. \tag{3}$$

In the formula, $\Phi$ is the measurement matrix, while $\Theta = \Phi\psi$ is the sensing matrix; both are matrix $M \times N$. The sampled signal $y$ obtained is the linear combination of the column of matrix $\Theta$. The linear combination coefficient is that in the corresponding original signal $\alpha$.

Since measurement matrix dimension $M \ll N$, the process of solving formula (1) is pathological and it is impossible to obtain original signal $x$ directly from $y$. However, as $\alpha$ is sparse, the estimated signal $\hat{\alpha}$ could be obtained almost perfectly through the compressed sensing reconstruction algorithm by the known sensing matrix $\Theta$, and then the original signal could be approximated with $\hat{x} = \varphi\hat{\alpha}$.

To ensure that $K$ coefficients can be accurately recovered from $M$ measurements, that is, to ensure that the algorithm is convergent, the $\Theta$ in formula (3) must satisfy the restricted



FIGURE 1: The compressed sensing framework.

equidistance (RIP) criterion [5]; that is, for the matrix $\Theta$ of size $M \times N$ and $M \ll N$, if there is a constant $\delta_k \in (0, 1)$, make all submatrices $\Theta_k \in R^{M \times k}$ for any vector $s \in R^{|k|}$ and $\Theta$, that is,

$$(1 - \delta_k)\|s\|_2^2 \le \|\Theta_k s\|_2^2 \le (1 + \delta_k)\|s\|_2^2. \tag{4}$$

It is said to satisfy the $k$-bound isometric property ($K$-RIP).

*3.1.3. Reconstruction Algorithm.* When $\Theta$ satisfies the limited equidistance property, the known perceptual matrix $\Theta$ can be used to solve formula (3) through the $l_0$ norm.

$$\min \|\alpha\|_{l_0} \text{ s.t.} y = \Theta\hat{\alpha}. \tag{5}$$

Thus, the estimated signal $\hat{\alpha}$ is obtained. However, since the solution of Equation (5) is an NP-hard problem, literature shows that under certain conditions, the minimum norm of $l_1$ and the minimum norm of $l_0$ are equivalent, and the same solution can be obtained [23]. Then, Equation (5) can be transformed into an optimization problem of the minimum norm of $l_1$.

$$\min \|\alpha\|_{l_1} \text{ s.t.} y = \Theta\hat{\alpha}. \tag{6}$$

The original signal is then approximated by $\hat{x} = \varphi\hat{\alpha}$. However, the algorithm for solving the minimum norm of $l_1$ is slow. Therefore, new reconstruction algorithms such as OMP [27], CoSaMP [28], and GOMP [29] have been proposed and achieved good results.

*3.2. Algorithm Flow.* The algorithm is divided into three steps. In the first step, the data obtained by the soil moisture sensor is not sparse, so the partial Fourier Transform (Permute Fast Fourier Transform (PFFT)) is adopted for sparse representation. The second step is to optimize the selection method of the sparse random observation matrix as the observation matrix suitable for this study to obtain the number and optimal location of sensors. Third, OMP, GOMP, and CoSaMP reconstruction algorithms were used to reconstruct the compressed data and compared with the original data, and it was found that OMP was more accurate in calculating the distribution of soil moisture data.

*3.2.1. Introduction and Evaluation of the Original Observation Matrix.* The sparse random observation matrix construction method [30] first generates an $M \times N$ all-zero

matrix, and $M \ll N$ in the matrix $\Phi$ of each column vector, randomly selected $d$ positions; in the selected location assignment 1, $d$ values are generally $d \in \{4, 8, 10, 16\}$ and have little effect on the reconstruction results [31]. $M$ is the observed value, and the number of sensors is shown herein.

When $d = 4$, the expansion of the matrix multiplication is as follows:

$$Y = \Phi X = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_{M-1} \\ y_M \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 1 & \cdots & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & \cdots & 0 \end{bmatrix}$$

$$\cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ \vdots \\ x_{N-1} \\ x_N \end{bmatrix} = \begin{bmatrix} x_1 + x_2 \\ x_2 + x_3 \\ x_2 + x_4 \\ \vdots \\ x_5 + x_n \\ x_5 + \cdots \end{bmatrix},$$

$$(7)$$

where $Y$ is an $M \times 1$ matrix and $X$ is an $N \times 1$ matrix. The column vector selects four positions of 1 so that each row vector is likely to have many 1's, and the row position of the elements in the different row vectors is different. The result $Y$ may be related to all the elements in $X$, that is, $Y \notin X$.

The idea of this paper is to select a part of the sensors to collect data in the original sensor layout. That is, in the matrix $X$, select a part of the elements to form the $Y$ matrix to meet $Y \in X$. The above observation matrix cannot be satisfied, and it is necessary to optimize the observation matrix.

*3.2.2. Observation Matrix Optimization.* The sparse random observation matrix is changed twice. First, the randomly selected object is the row vector of the matrix $\Phi$. Second, only one position is selected in each row vector to assign it to 1. Make sure that the optimized observation matrix has only one element value per line.

The improved method for constructing a sparse random measurement matrix is as follows. Firstly, generate an identity matrix $\Phi \in \text{ones}^{N \times N}$. Secondly, randomly select $M$ row vectors from the generated matrix to form a matrix of $M \times N$. Since the identity matrix is an orthogonal matrix, the partial identity matrix of $M \times N$ size obtained after taking $M$ rows from it still has a strong noncorrelation and partial orthogonality. It satisfies the RIP theorem and ensures that the observation matrix will not combine two different sparse data mapped to the same collection.



FIGURE 2: Original sensor distribution.



FIGURE 3: Soil moisture surface.

When $d = 4$, after improving the measurement matrix, $Y = \Phi X$ corresponds to the expanded form of matrix multiplication as follows:

$$Y = \Phi X = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_{M-1} \\ y_M \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 0 & \cdots & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & \cdots & 0 \end{bmatrix}$$

$$\cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ \vdots \\ x_{N-1} \\ x_N \end{bmatrix} = \begin{bmatrix} x_2 \\ x_4 \\ x_3 \\ \vdots \\ x_N \\ x_6 \end{bmatrix}.$$

$$(8)$$

FIGURE 4: Humidity distribution at each sensor point.

TABLE 1: Comparison of different reconstruction algorithms.

| Soil type | Reconstruction algorithm | Observed value $M$ | Average relative error |
|---|---|---|---|
| Dry soil | OMP | 41 | 3% |
| | GOMP | 41 | 5% |
| | CoSaMP | 41 | 6% |
| Moist soil after irrigation | OMP | 39 | 9% |
| | GOMP | 39 | 11% |
| | CoSaMP | 39 | 12% |

$Y[m, 1] = \Phi[m, n] * X[n, 1]$, when $\Phi[m, n] = 1$ and $Y[m, 1] = X[n, 1]$. If the element 1 in $\Phi$ is in the $m$ rows and $n$ columns, the $m$th sensor position after sampling corresponds to the position of the $n$th sensor in $X$.

Compared with the gridded dense sensors, $M$ sensors should be selected based on $N$ sensors for sampling; that is, for the observation matrix $\Phi$ of $M \times N$, each row has an only one value of 1, and each column has at most one value of 1; that is, the $M$ sensor should be deployed at the element position with the $M$ row vector of 1.

## 4. Experiments

*4.1. Data Acquisition.* The soil humidity sensors numbered successively with $1, 2, \cdots, 64$ are placed evenly at 10 cm below the soil surface of a 40 m × 40 m farmland as shown in Figure 2. Experimental data is the soil moisture value measured by all sensors on June 1, 2015, solstice, and November 30, 2015.

Through data analysis, it is found that farmland soil moisture has strong spatial and temporal differences, so it is necessary to deploy more sensors to accurately monitor it.

Firstly, the soil moisture data of farmland have spatial differences. Figure 3 shows the soil moisture surface at 21:00 on May 31, 2015, obtained by the bilinear interpola-

tion method. The maximum value of moisture at point $(8, 4)$ is 32.91, and the minimum value at point $(7, 8)$ is 25.52, with a difference of 22.5%.

Secondly, the soil moisture data of farmland varies with time. Figure 4 shows the soil moisture curves of 10 sensors at 100 time points during June 1, 2015, solstice, and June 30, 2015. For sensor number 10, the maximum humidity at point $(10, 83)$ is 34.36, and the minimum humidity at point $(10, 1)$ is 17.5, with a difference of 49.07%.

*4.2. The Evaluation Indexes.* To analyze the performance index of sensor layout optimization in different compressive sampling conditions, an indicator of absolute error is provided in this paper. The formula is as follows:

$$\text{mean absolute error}: a = \frac{\sum_{i=1}^{n} |x_i - \hat{x}_i|}{n},$$
$$\text{average relative error}: r = \frac{\sum_{i=1}^{n} (|x_i - \hat{x}_i|/x_i)}{n}, \quad (9)$$

where $x_i$ represents the original soil moisture data, $\hat{x}_i$ represents the reconstructed soil moisture data, $i$ represents the number of the sensor in Figure 2, and $n$ represents the number of soil moisture data. The unit of $a$ is %, which indicates the relative moisture content value.

*4.3. Selection of the Reconstruction Algorithm.* In this experiment, the soil moisture distribution image composed of 64 sensor points is relatively simple; OMP, CoSaMP, GOMP, and other algorithms are suitable for a relatively simple image selected point compression. Therefore, OMP, GOMP, and CoSaMP algorithms are used for comparative analysis.

Compare the three kinds of reconstruction accuracy of the algorithm (OMP, GOMP, and CoSaMP). The data of soil moisture measured by all sensors in a set of dry soil (23:00 on May 28, 2015) and a set of irrigated soil (6:00 on July 15, 2015) were selected. The number of original sparse signals is $N = 64$, the observed values are $0 < M < 64$, and the

(a) Error-sparsity curve when $M = 10$



(b) Error-sparsity curve when $M = 20$



(c) Error-sparsity curve when $M = 40$

Figure 5: Error-sparsity variation curve.



(a) Position of 10 sensors



(b) Position of 20 sensors



(c) Position of 40 sensors

Figure 6: Error-sparsity variation curve.

sparsity is $K \in \{3, 7, 9, 13\}$. The observation matrix is an optimized sparse random measurement matrix. Each observed value was simulated 500 times to determine the probability of accurate recovery. OMP, GOMP, and CoSaMP algorithms were compared to determine the relationship between the recovered data and the observed value $M$ and the sparsity $K$ under a given sparsity. The experimental results are shown in Table 1.

For the two sets of data, under the same upper limit of residual error and the same observed value $M$, the average relative error of the OMP algorithm is only 3%-9%, which is better than that of the GOMP algorithm (5%-11%) and CoSaMP algorithm (6%-12%). Therefore, among the three reconstruction algorithms, the OMP reconstruction algorithm has a better effect. In the following experiments, the OMP algorithm is used to reconstruct data.

(a) Reconstruction effects of 10 sensors

(b) Reconstruction effects of 20 sensors

(c) Reconstruction effects of 40 sensors

FIGURE 7: Error-sparsity variation curve.

### 4.4. Sparsity Selection.

When original data is converted to a sparse vector through PFFT, the number of nonzero elements in the sparse vector is denoted as the sparsity degree $K$. For the determined reconstruction algorithm and $M$, the errors vary with the value of $K$. With the soil humidity value at 21 o'clock on July 30, 2015, under the conditions $M_1 = 10$, $M_2 = 20$, and $M_3 = 40$, respectively, we change the value of $K$ and take an iteration of 500 times for each sparsity to average the errors and observe the variation of errors with $K$. The simulation results are shown in Figure 5.

When the observed value $M = 10$ or $M = 20$ and $K = 3$, the mean absolute error is the minimum, while when $M = 40$, sparsity $K$ should be 13, to minimize the mean absolute error. Therefore, in the following experiments, we all adopted the optimal sparsity.

### 4.5. Reconstruction of Soil Water Spatial Distribution.

Soil humidity at 21 o'clock on July 30, 2015, is selected as experimental data. With the OMP reconstruction algorithm, under the circumstances of $M_1 = 10$ plus $K_1 = 3$, $M_2 = 20$ plus $K_2 = 3$, and $M_3 = 40$ plus $K_3 = 13$, respectively, multiple

TABLE 2: Comparison of reconstruction effects of different $M$ values.

| $M$ | Mean absolute error | Least absolute error | Maximum absolute error |
|-----|---------------------|----------------------|------------------------|
| 10  | 1.68                | 0.41                 | 3.71                   |
| 20  | 0.91                | 0.37                 | 2.12                   |
| 40  | 0.47                | 0.18                 | 1.14                   |

iterations are used to get a minimum of errors to determine the location of corresponding sensors. The locations under the 3 circumstances are shown in Figure 6.

Consider three cases, and the refactoring effect is shown in Figure 7. The abscissa represents 64 sensors, and the ordinate represents the soil moisture value. When $M = 40$, the predicted value almost coincides with the actual value. When $M = 20$, the conformance is also good, and the result is acceptable. However, when $M = 10$, the predicted value differs greatly from the actual value and cannot be used.

In the three cases, the distribution of mean absolute errors at all sampling moments is shown in Table 2.

(a) 10 sensors



(b) 20 sensors



(c) 40 sensors

FIGURE 8: Error distribution at different moments.

When $M = 10$, the mean absolute error of 80% of the sampling moments is concentrated below 1.5, the error of 12% of the sampling moments is between 1.5 and 2.5, and the error of 8% of the sampling moments is between 2.5 and 3.5. When $M = 20$, the mean absolute error of 85% of the sampling moments is below 1, the error of about 12% is between 1 and 1.5, and the error of only about 3% is above 1.5. When $M = 40$, the errors of all sampling moments are below 1.14. Therefore, considering the balance between sensor installation cost and reconstruction accuracy, it is recommended to deploy 20 sensor sampling points to obtain more accurate reconstruction results of soil water distribution.

## 5. Evaluation of the Effectiveness of Sensor Placement

It is limited to get the abovementioned sensor arrangement optimization at a moment. Therefore, the 250 time points from June 1, 2015, to November 30, 2015, are selected as the experimental data. The overall error distribution of the three strategies at different times is shown in Figure 8.



FIGURE 9: Recovery of abnormal points.

In Figure 8, the $r_{x,y}$ distributions in the three graphs are relatively concentrated, which proves that the strategies of optimization sensor placement in this experiment are

effective and feasible. Most of the $r_{x,y}$ in Figure 8(a) are concentrated below 1.5, and the $r_{x,y}$ of about 30 points is between 1.5 and 2.5, and the $r_{x,y}$ of about 20 points is between 2.5 and 3.5; the error of Figure 8(b) is mostly below 1, and the error of a few points is above 1; the $r_{x,y}$ of Figure 8(c) is below 1.134. It can be seen that the $r_{x,y}$ of 10 sensors is maximum, and the placement of 20 sensors and 40 sensors is similar. Considering the cost of the sensor, 20 sensors can be arranged to achieve more accurate measurements.

In Figure 8(a), there are some errors between [3, 3.5]. Choose two of them for detailed analysis, and get Figure 9.

Figure 9 is the 220th time in Figure 8(a). The error is 3.2325. The effect is generally and only part of the data close to the original data. The placement of the 10 sensors is too difficult to obtain all the features mainly due to significant changes in soil characteristics. However, it can be found that the amplitude of the reconstruction curve is stable, and the optimization strategy of 10 sensors is not suitable for the environment with high accuracy or obvious soil characteristic change.

## 6. Discussion

In summary, the overall error of arranging 20 sensors is close to that of arranging 40 sensors. Considering the cost issue, using 20 sensors can get a more accurate acquisition of soil moisture. Due to the obvious changes in soil moisture characteristics, the placement of 10 sensors is too small, resulting in large errors. By observing the data recovery of 10 sensors at a certain moment, it is found that the error is large at the maximum and minimum values, but the overall curve fluctuates stably. Therefore, the method of using only 10 sensors is suitable for situations where the data change range is not large.

## 7. Conclusions

Aiming at the problems of unreasonable sensor placement and high cost in the agricultural IoT, this paper proposes an optimization strategy of sensor placement based on the compressed sensing theory. By analyzing the soil moisture data at a certain moment and the optimization of the observation matrix, three optimal strategies of sensor layout were obtained and then verified at more time points. Through experiments, it is found that 1/3 of the original sensor can be used to measure soil moisture with a lower error. The purpose of obtaining more accurate data with fewer sensors is realized. The overall error of 20 sensors is close to that of 40 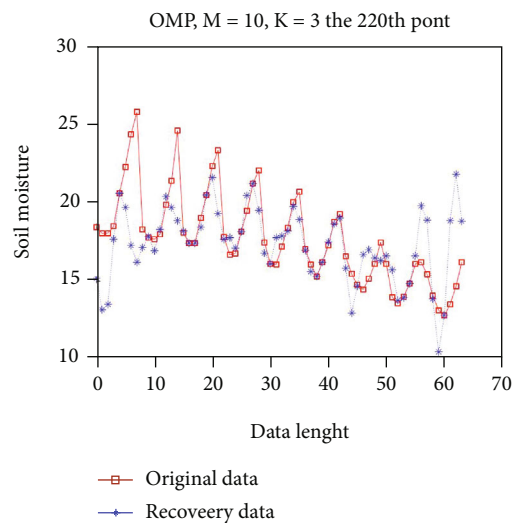sensors. Considering the cost, 20 sensors can be used to obtain soil moisture more accurately. Due to the obvious change of soil moisture characteristics, the 10 sensors were too few, resulting in a large error. When observing the data recovery situation at a certain moment, it was found that the error was large at the maximum and minimum values, but the overall fluctuation of the curve was stable, which was suitable for the situation with a small range of data changes.

When used, the moisture distribution of the mesh point can rely on manual multiple measurements, no need to install the sensor. After several measurements, the sensor deployment can be determined by this article, and only about 1/3 of the sensor can achieve a better effect, so the cost is relatively low.

The shortcomings of the experiment mainly include the following two aspects: (1) the sensor was numbered in one dimension, while the two-dimensional spatial correlation of soil moisture was ignored; and (2) soil moisture data at different moments were uniformly set to equal sparsity, resulting in large errors. These will be the focus of future research.

## Data Availability

The soil moisture 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 is no conflict of interest regarding the publication of this paper.

## Acknowledgments

## References

[1] Y. Huang, W. Han, L. Zhou, W. Liu, and J. Liu, "Farmer cognition on water-saving irrigation technology and its influencing factors analysis," *Transactions of the Chinese Society of Agricultural Engineering*, vol. 28, no. 18, pp. 113–120, 2012.

[2] W. Jiulin, Z. Chengping, and Y. Hua, "Design and research of water saving irrigation system based on LoRa," *Water Saving Irrigation*, vol. 12, pp. 104–111, 2017.

[3] M. H. Cosh, T. J. Jackson, R. Bindlish, and J. H. Prueger, "Watershed scale temporal and spatial stability of soil moisture and its role in validating satellite estimates," *Remote Sensing of Environment*, vol. 92, no. 4, pp. 427–435, 2004.

[4] Y. Wang, C. Wu, H. Liu, and M. Chong, "The influence of the depth and amount of soil moisture sensors on the accuracy of soil moisture content," *Water Saving Irrigation*, vol. 1, pp. 87–91, 2019.

[5] E. J. Candes, J. Romberg, and T. Tao, "Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information," *IEEE Transactions on Information Theory*, vol. 52, no. 3, pp. 489–509, 2006.

[6] Y. Zhou, X. U. Sheng, and Y. Cheng, "Preliminary study on the application of Internet of Things technique in farmland quality protection and management," *Agriculture Network Information*, 2016.

[7] Z. Feng, "Research on water-saving irrigation automatic control system based on Internet of Things," in *2011 International Conference on Electric Information and Control Engineering*, pp. 2541–2544, Wuhan, China, 2011.

[8] Y. Liu, "Data gathering mechanism for farmland wireless sensor network based on the Internet of Things," *Journal of Anhui Agricultural Sciences*, vol. 26, 2011.

[9] Y. Liu and W. Yang, "Management for farmland environment monitoring WSN based on the Internet of Things," *Chinese Agricultural Science Bulletin*, vol. 27, no. 30, pp. 297–302, 2011.

[10] P. Singh and S. Saikia, "Arduino-based smart irrigation using water flow sensor, soil moisture sensor, temperature sensor and ESP8266 WiFi module," in *2016 IEEE Region 10 Humanitarian Technology Conference (R10-HTC)*, Agra, India, 2016.

[11] S. Biwei, "Noise folding phenomenon and noise signal reconstruction in compressed sensing," Wuhan University, 2017.

[12] C. Papadimitriou and G. Lombaert, "The effect of prediction error correlation on optimal sensor placement in structural dynamics," *Mechanical Systems & Signal Processing*, vol. 28, no. 2, pp. 105–127, 2012.

[13] R. Stoughton and T. Arai, "Optimal sensor placement for forward kinematics evaluation of a 6-DOF parallel link manipulator," in *Proceedings IROS '91:IEEE/RSJ International Workshop on Intelligent Robots and Systems '91*, pp. 785–790, Osaka, Japan, 1991.

[14] J. Göbel and P. Trinius, "Towards optimal sensor placement strategies for early warning systems," Gesellschaft für Informatik, Bonn, 2010.

[15] R. Herzog and I. Riedel, "Sequentially optimal sensor placement in thermoelastic models for real time applications," *Optimization and Engineering*, vol. 16, no. 4, pp. 737–766, 2015.

[16] D. C. Kammer, "Sensor placement for on-orbit modal identification and correlation of large space structures," *Journal of Guidance, Control, and Dynamics*, vol. 14, no. 2, pp. 251–259, 1991.

[17] G. Heo, M. L. Wang, and D. Satpathi, "Optimal transducer placement for health monitoring of long span bridge," *Soil Dynamics and Earthquake Engineering*, vol. 16, no. 7-8, pp. 495–502, 1997.

[18] Z. Y. Wu, X. H. Jian, and Z. Bin, "Multi-objective optimal sensor placement methodology for vibration test," *Journal of Mechanical Strength*, vol. 6, pp. 888–892, 2008.

[19] R. Mukherjee, U. M. Diwekar, and A. Vaseashta, "Optimal sensor placement with mitigation strategy for water network systems under uncertainty," *Computers & Chemical Engineering*, vol. 103, pp. 91–102, 2017.

[20] A. Krause, C. Guestrin, A. Gupta, and J. Kleinberg, "Near optimal sensor placements: maximizing information while minimizing communication cost," in *Proceedings of the 5th international conference on Information processing in sensor networks*, Nashville Tennessee USA, 2006.

[21] C. Guestrin, A. Krause, and A. P. Singh, "Near-optimal sensor placements in Gaussian processes," in *Proceedings of the 22nd international conference on Machine learning*, New York, NY, United States, 2005.

[22] Y. Xu and J. Choi, "Adaptive sampling for learning Gaussian processes using mobile sensor networks," *Sensors*, vol. 11, no. 3, pp. 3051–3066, 2011.

[23] D. L. Donoho, "Compressed sensing," *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.

[24] M. A. Razzaque and S. Dobson, "Energy-efficient sensing in wireless sensor networks us-ing compressed sensing," *Sensors*, vol. 14, no. 2, pp. 2822–2859, 2014.

[25] L. Yu, D. Xiong, L. Guo, and J. Wang, "A compressed sensing-based wearable sensor network for quantitative assessment of stroke patients," *Sensors*, vol. 16, no. 2, p. 202, 2016.

[26] D. Wang, R. Xu, X. Hu, and W. Su, "Energy-efficient distributed compressed sensing data aggregation for cluster-based underwater acoustic sensor networks," *International Journal of Distributed Sensor Networks*, vol. 2016, 2016.

[27] J. A. Tropp and A. C. Gilbert, "Signal recovery from random measurements via orthogonal matching pursuit," *Information Theory IEEE Transactions on*, vol. 53, no. 12, pp. 4655–4666, 2007.

[28] D. Needell and J. A. Tropp, "CoSaMP: iterative signal recovery from incomplete and inaccurate samples," *Applied & Computational Harmonic Analysis*, vol. 26, no. 3, pp. 301–321, 2008.

[29] W. Jian, S. Kwon, and B. Shim, "Generalized orthogonal matching pursuit," *IEEE Transactions on Signal Processing*, vol. 60, no. 12, pp. 6202–6216, 2011.

[30] A. Gilbert and P. Indyk, "Sparse recovery using sparse matrices," *Proceedings of the IEEE*, vol. 98, no. 6, pp. 937–947, 2008.

[31] X. Li, "Research on measurement matrix based on compressed sensing," Beijing Jiao Tong University, 2010.

*Research Article*

# IoT Cloud-Based Framework for Face Spoofing Detection with Deep Multicolor Feature Learning Model

**Sajad Einy** [1,2] **Cemil Oz** [1] **and Yahya Dorostkar Navaei** [3]

[1]*Computer Engineering Department, Sakarya University, Turkey*
[2]*Application and Research Center for Advanced Studies, Istanbul Aydin University, Turkey*
[3]*Computer and Information Technology Engineering, Qazvin Branch, Islamic Azad University, Qazvin, Iran*

Correspondence should be addressed to Yahya Dorostkar Navaei; y.dorostkar@qiau.ac.ir

A face-based authentication system has become an important topic in various fields of IoT applications such as identity validation for social care, crime detection, ATM access, computer security, etc. However, these authentication systems are vulnerable to different attacks. Presentation attacks have become a clear threat for facial biometric-based authentication and security applications. To address this issue, we proposed a deep learning approach for face spoofing detection systems in IoT cloud-based environment. The deep learning approach extracted features from multicolor space to obtain more information from the input face image regarding luminance and chrominance data. These features are combined and selected by the Minimum Redundancy Maximum Relevance (mRMR) algorithm to provide an efficient and discriminate feature set. Finally, the extracted deep color-based features of the face image are used for face spoofing detection in a cloud environment. The proposed method achieves stable results with less training data compared to conventional deep learning methods. This advantage of the proposed approach reduces the time of processing in the training phase and optimizes resource management in storing training data on the cloud. The proposed system was tested and evaluated based on two challenging public access face spoofing databases, namely, Replay-Attack and ROSE-Youtu. The experimental results based on these databases showed that the proposed method achieved satisfactory results compared to the state-of-the-art methods based on an equal error rate (EER) of 0.2% and 3.8%, respectively, for the Replay-Attack and ROSE-Youtu databases.

## 1. Introduction

Nowadays, the Internet of Things (IoT) affects human lives in a wide range of technology from smart homes to smart cities. An enormous number of IoT devices are utilized for collecting and analyzing information for different reasons, such as healthcare, security, and management. According to the estimation of scientifics, around 90% of storing data would be useless [1]. Therefore, the researchers proposed [1] utilizing the edge devices in the architecture of applications or services for cloud computing. In this way, the data can be analyzed and filtered in edge devices and send more enhanced data for processing in the cloud. For example, the deployed sensors for traffic monitoring can be also utilized for fire detection with low-cost and low-performance devices. However, IoT-based systems are faced with different problems such as security threats from the Internet. For instance, let us consider an IoT-based healthcare application which contains critical information such as blood sugar level and blood pressure. The authentication system for data communication through wireless channels should be secured for protecting critical information of clients. Biometric authentication can be utilized for identifying a person in wireless communication. This authentication requires using personal attributes, such as speech, face, fingerprints, palmprint, gait, and iris [2]. This kind of authentication is based on a comparison between the physical aspect of the client that is collected with the help of different sensors and a copy that was stored. The physiological information of clients is more reliable when compared to knowledge-based or token-based methods because this information is unique and not shareable. For this reason, IoT-based cloud computing systems for authentication of clients applied their biometric information.

For instance, Kumari and Thangaraj [3] proposed a feature selection technique in biometric authentication using a cloud framework. In another similar study, Shakil et al. [4] proposed a biometric authentication system and data management application for security of healthcare data in the cloud. Also, Vidya and Chandra [5] proposed a multimodal biometric authentication system based on entropy-based local binary pattern feature description technique for cloud computing. Additionally, Masud et al. [6] proposed a deep learning-based approach for face recognition in IoT environments. Face recognition systems have achieved significant interest in many applications such as cell phones' and laptops' authentication or registration systems at places such as online exam centers and airports [1]. These kinds of security systems in the Big Data analytics platform are a topic of concern for real-time applications. Consider the scenario when a person is to be recognized in an airport for registration or a student is attending an online exam. In these scenarios and other similar conditions, the camera captures images of the face continuously and sends these data for processing in the cloud environment. Based on meaningful information of face image, a certain person can easily be identified. Nevertheless, these kinds of authentication and registration systems are vulnerable to different types of attacks. For improving the security of biometric authentication systems, various methods and models are proposed.

For example, Ali et al. [1] proposed a multimodal biometric authentication system using an encryption method for protecting the privacy of biometric information in the IoT-based cloud environment. In another study, Gomez-Barrero et al. [2] proposed a framework for the protection of the privacy of multibiometric templates with an encryption method. However, the aforementioned methods are designed for protection based on man-in-the-middle attacks in wireless communication. According to the literature, face spoofing attacks in IoT cloud environments are not discussed and studied yet. The main objective of this study is to present an IoT cloud-based framework for protecting client's information from face spoofing attacks. In a face spoofing attack, the intruder bypasses the authentication system by presenting a fake face of the victim. Due to this threat, robust and stable face Presentation Attack Detection (PAD) methods must be developed and designed. Face spoofing attacks may be classified into four main groups: print, display, replay, and mask attacks [7].

According to the types of sensors for detection of these kinds of attacks, different algorithms are proposed [9–11]. Generally, light field camera sensors are more popular compared to other sensors such as infrared and thermal ones [8] or multibiometric fusion systems [9] because this additional equipment increases the cost of authentication systems. In this case, many researchers investigate feature-based methods. These kinds of spoofing detection methods attempt to extract discriminative features to recognize the genuine user from a fake face. For example, in print, display, and mask attacks, facial liveness features such as lip movement, head movement, and eye blinking can help recognize spoofing attacks. Furthermore, detection of replay attacks is more challenging because they contain this kind of liveness feature [7]. In some cases, the intruder applies liveness features in a mask attack by cropping the lip and eye area from a mask, which shows that liveness features alone cannot detect spoofing attacks properly. Replay display and printed attack images contain some noise and defects because of recapturing of information by a camera. During recapturing of information, the fake face loses the high-frequency information by getting affected in terms of the texture and color information of images, and these features can help distinguish a genuine person and a recaptured face image. Especially in printing and displaying attacks, during recapturing of information, some defects and noises appear in the spoofing face image. These artifacts lead to inadequate color reproduction in comparison to real biometric samples [10]. RGB is the commonly employed color space for sensing and displaying color images on many devices. Nevertheless, this color space in image analysis is inadequate due to the high correlation between the red, green, and blue color components and incomplete separation of the luminance and chrominance information [11]. Therefore, a different color space may help extract discriminative features for extraction of liveness cues of skin tones for detection of live and fake images. Therefore, image texture analysis based on different color spaces has attracted the consideration of research areas in the field of face spoofing attacks [11, 12]. By the success of deep learning algorithms in the field of computer vision and multimedia analysis, deep texture analysis-based algorithms have been employed in face spoofing problems. Nevertheless, deep learning-based face spoofing detection algorithms are faced with some problems such as few numbers of spoofing data and lack of diversity of scenarios which make it difficult to train a deep network [13, 14]. Additionally, IoT-based authentication systems encountered several difficulties such as storing or processing in a real-time manner [6].

To address these problems, we presented a novel approach based on hybrid convolutional neural network (CNN) models on different color spaces for IoT-based cloud computing. The proposed deep learning approach utilized three pretrained models in different color spaces for extracting luminance and chrominance information which are useful in recognition of spoofing face images. Due to extracted robust and discriminative features from a single image, this proposed model can achieve satisfactory results with less training dataset. This advantage of the proposed approach helps to decrease the storing training data in cloud computing which tackles one of the major problems of cloud computing systems. To the best of our knowledge, for the first time, in this paper, an IoT security framework is proposed for face spoofing detection. Extensive experimental analysis was conducted based on two challenging public access spoofing databases with their predefined evaluation protocols for comparison of our proposed approach against state-of-the-art methods. These experimental results show that our proposed approach outperforms all existing deep-based methods among state-of-the-art methods based on benchmark databases. In addition, experimental results show that the proposed approach can achieve stable results with less training dataset compared to benchmark deep learning models.

In light of this information, the main contributions of this paper are presenting an IoT security framework for face spoofing detection which achieved significant results compared to the state of algorithms based on two public databases. Also, the proposed approach achieved stable results with less training dataset compared to benchmark deep learning models.

This paper is briefly organized as follows: In Section 2, short information about types of existing systems and related works on face spoofing methods are available. In Section 3, the methodology of the proposed approach is briefly presented. In Section 4, the experimental results and state-of-the-art algorithms with benchmark databases and protocols are presented. As the final section, conclusion statements are provided in Section 5.

## 2. Related Work for Face Spoofing Methods

Recently, a lot of face spoofing detection algorithms have been proposed [1–7], based on different cues and attacks. Based on our prior knowledge, the algorithms can be categorized into four different groups: texture analysis, motion analysis, image quality analysis, and hardware-based methods.

*2.1. Texture-Based Methods.* Face liveness detection algorithms based on texture analysis usually recognize the effects of illumination limitations of a printer or any other device during display, such as printing failures, blurring, and other effects. The RGB color space, as discussed in Section 1, cannot clearly present features regarding illumination and chrominance. In this case, a previous study [12] proposed a deep learning system based on the RGB, HSV, and YCbCr color spaces. In the paper, the CompactNet model was proposed as a layer-by-layer progressively generated color space. Additionally, features of spoofing databases are extracted by a pretrained feature extractor model. Researchers [11] proposed a color feature descriptor method based on different color spaces. In this method, information on the luminance and chrominance channels was extracted by a low-level feature descriptor. Due to the impact of a smaller number of databases in face spoofing detection on training deep learning methods and overfitting problems, researchers investigate the extraction of discriminative and deep features. For instance, a study [15] proposed a perturbation layer (low-level deep features) to extract the deep features of a convolutional neural network (CNN) for classification. Another study [16] presented an adaptive fusion of convolutional feature models to learn the features of face images, and a deep autoencoder was utilized for generating a face image to detect spoofing face images. Some authors [7] proposed a Spatial Pyramid Coding Microtexture (SPMT) feature extractor with a deep learning system for detection of liveness cues and employed the Single Shot Multibox Detector (SSD) as an end-to-end face spoofing detection model. Besides the aforementioned color-based deep learning methods, some methods presented local binary pattern- (LBP-) based feature descriptors for spoofing detection. For instance, a hybrid method was proposed [17] based on the Chromatic Cooccurrence of Local Binary Pattern (CCoLBP) and Ensemble Learning (EL) algorithms. In the case of reducing the param-

eters of CNN models and extraction of deep features, an end-to-end learnable LBP network was proposed [18]. A previous study [19] proposed an algorithm by integrating the LBP descriptor with a modified convolution neural network that extracted deep texture. For extraction of discriminative features of presentation attacks, the Extended Local Ternary Corelation Pattern (ELTCP) feature extraction method was proposed [20]. This feature descriptor with extraction of spatial information of an image in multiple directions achieved robust results on presentation attacks. In recent years, with increasing attention to 3D face spoofing attacks, several studies have been devoted to recognizing 3D mask attacks. For instance, the 3D wax face attacks [21] approach is proposed with a convolutional neural network based on the Residual Attention Network (RAN) for 3D face spoofing detection. In another similar study, a multichannel CNN [22] approach with a one-class Gaussian mixture model is proposed for the detection of 2D and 3D attacks. Another study [23] presented a shading-based 3D feature description method to extract discriminative and robust 3D features from the face image. In another study, researchers proposed [24] a face spoofing framework with the help of convolutional autoencoders for the detection of 3D mask attacks. Another study [10] investigated various factors of affection of acquisition conditions and devices with different resolutions on the generalization of color texture features for spoofing detection. In this light, another possibility seems to be analyzing image textures based on deep features from multiple color spaces, which is proposed in this paper. The experimental results show that our proposed algorithm is superior in color texture extraction and classification over state-of-the-art methods.

*2.2. Motion Analysis Algorithms.* Among texture recognition techniques, motion-based analysis also plays an important role in spoofing detection. For instance, a study [25] proposed a motion-based analysis approach based on rigid and nonrigid facial movements. The proposed system extracted motion cues such as face movement, lip movement, and hand shaking and classified them into natural and fake motions. In another study [8], an undirected conditional random field in video processing was proposed for the detection of eye blinking. Other researchers [26] proposed a dynamic mode decomposition pipeline with SVM and LBP. This algorithm extracted facial dynamic information in videos as an image sequence.

*2.3. Image Quality Analysis.* In spoofing attacks, the image quality is mostly reduced due to the image being reproduced. Based on this inability of devices, some methods have been proposed. For instance, in a previous study [27], an algorithm was proposed where real and fake face images were determined by analysis and comparison of both reflections taken from an LCD screen. In another study [28], it was posited that it is possible to differentiate a fake image from a real one by analyzing the noise signatures with the Fourier spectrum.

*2.4. Hardware-Based Analysis.* Researchers [29] proposed video-based stereo face antispoofing recognition systems. In this approach, for learning a dynamic disparity map, a CNN classifier with a disparity layer was proposed. In

another study [30], it was proposed to assign a light field to traditional HOG which was utilized for gathering texture information from 2D images and Light Field Histogram Of Gradient (LFHOG).

Apart from the mentioned proposed systems based on single cues, some methods have been proposed based on multicue approaches. For instance, another study [31] proposed a multicue face spoofing detection framework involving image quality analysis by employing the Shearlet method and motion analysis by utilizing the dense optical flow method. In this study, the extracted multicue features were fused and classified with a deep neural network.

# 3. Proposed IoT-Based Framework Face Spoofing Detection

The smart city framework contains multiple components such as smart devices, high-speed wireless networks, and cloud servers, as presented in Figure 1. The captured face images by IoT devices are analyzed and preprocessed with edges. The preprocessing section with edges and smart devices included Viola and Jone's [32] face detection algorithm for extracting face images and sending more enhanced data to optimize the resource of the cloud. Then, the captured faced images are continually sent to a cloud environment using wireless technology. In the cloud section, several Virtual Machines (VMs) work in a parallel mode. These VMs by employing a deep learning approach recognize spoofing attacks.

Before feeding the face image to the deep model for classification in cloud computing environments, RGB color space is transformed to the HSV and YCbCr color spaces. Three parallel pretrained models are utilized in the proposed deep learning approach. Based on the literature, because of the small number of data and lack of scenarios in controlled environments, it is quite hard to train CNN models from scratch and achieve a stable and high-performance model. In this case, we utilized the VGG-face [33] model in the RGB color space for face spoofing detection [14, 18]. In addition, the transformed images of the HSV and YCbCr color spaces are trained by the VGG16 [34] model individually on the cloud side. After fine-tuning models by a different color space, the features of the last fully connected layer which consists of 4096 features for each deep model are extracted. These features are combined and then selected by employing the Minimum Redundancy Maximum Relevance (mRMR) feature selection algorithm. These selected features are classified with the help of different classification algorithms such as linear regression (LR), Support Vector Machine (SVM), Linear Discriminant Analysis (LDA), and $K$ Nearest Neighborhood (KNN) for detection of the spoof image, as presented in Figure 2.

Suppose a scenario where a student wants to access an online exam. A smart device such as a smart phone or computer captures the student's facial image and sends this image to the cloud using 5G wireless technology. In the cloud server, by employing the face spoofing image database and deep learning method, a deep feature set of face image is extracted in three different color spaces. These combined feature sets contain various aliveness keys from face skin tones

which help to detect face spoofing in the online exam scenario. The proposed method is tested and evaluated based on two public access databases, namely, Replay-Attack and ROSE-Youtu. The Replay-Attack database is captured by a MacBook laptop webcam and the ROSE-Youtu database captured by Huawei, iPhone 5s, ZTE, and Hasee smart phone.

*3.1. Color Space Transform.* RGB is a common color space for many devices and sensors for displaying and sensing color images. Nevertheless, this color space is quite limited for analyzing images because of the high correlation of red, green, and blue colors and incomplete separation of the luminance and chrominance information.

In this case, for the detection of recapping artifacts in spoofing databases, different color spaces are utilized [12]. HSV and YCbCr in addition to RGB provide robust features to detect different liveness cues from face skin tones. Both the HSV and YCbCr color spaces provide color texture information such as the luminance and the chrominance components. In the HSV color space, the H and S define the hue and saturation dimensions for presenting the chrominance information, and V defines the value dimension for presenting the luminance information of images. The YCbCr space separates RGB into luminance (Y), Chrominance Blue (Cb), and Chrominance Red (Cr). The HSV and YCbCr spaces provide discriminative color-based texture from face skin tones in different spoofing attacks [11, 12]. Figure 3 presents different color spaces on the Replay-Attack database for both live and fake face images.

*3.2. Convolutional Neural Networks.* Convolutional neural networks (CNNs) are designed and developed to automatically learn the spatial hierarchies of features with the help of back propagation algorithms [35]. CNNs are designed based on multiple layers of neurons which mainly include multiple basic structural blocks such as the convolution, pooling, and fully connected (FC) layers. Each convolutional layer contains a set of filters whose sizes can be $3 \times 3$, $5 \times 5$, or $7 \times 7$ pixels. Therefore, each convolutional layer, by applying a filter, creates the input of the next layer [36]. The results of this convolution process are activation maps which contain local distinctive features. Based on Equation (1), the output of $Y_i^{(l-1)}$ of the $L$ layer contains $m_1^{(l)}$ feature maps with sizes of $m_2^{(l)} \times m_3^{(l)}$. In this equation, $B_i^{(l)}$ and $k_{I,j}^{(l)}$ represent, respectively, the basis matrix and the filter size for the $i$th feature map [37]:

$$Y_i^{(l)} = f\left( B_i^{(l)} + \sum_{j=1}^{m_i^{(l-1)}} k_{i,j}^{(l)} \times Y_j^{(l-1)} \right). \tag{1}$$

The pooling layer reduces the spatial size of the image to reduce the number of parameters and computations in the model. This layer operates on each feature map independently to keep the image features and information intact. Each pooling layer $L$ contains two main parameters as the spatial size of the filter $F^{(l)}$ and $S^{(l)}$ step. The input of the pooling layer is data

FIGURE 1: Proposed IoT-based framework for face spoofing detection.



FIGURE 2: The architecture of deep learning approach.

with the size of $m_1^{(l-1)} \times m_2^{(l-1)} \times m_3^{(l-1)}$, and the output volume of this layer is $m_1^{(l)} \times m_2^{(l)} \times m_3^{(l)}$. Equation (2) briefly presents the operation of the pooling layer:

$$\begin{cases} m_1^{(l)} = m_1^{(l-1)}, \\ m_2^{(l)} = \dfrac{m_2^{(l-1)} - F^{(l)}}{S^{(l)}} + 1, \\ m_3^{(l)} = \dfrac{m_3^{(l-1)} - F^{(l)}}{S^{(l)}} + 1. \end{cases} \quad (2)$$

The output of feature maps of the last convolutional or pooling layer is flattened in the layer named the fully connected layer. The FC layer transforms the output of previous layers into a one-dimensional feature vector, updates the

weights, and provides the latest possible values for each label [37]. These layers may be connected to a more fully connected layer which is also known as the dense layer. By employing a learning rate, every input is connected to every output. The features are extracted by the convolution layers, downsampled by the pooling layers, and mapped by the FC layer to the final output of the model. The last FC layer contains a number of nodes equal to the number of classes of classification images. Each FC layer is supported by a nonlinear function such as the ReLU function. Equation (3) presents the FC layer's processing steps by weights ($W$) and the $f(Z_i^{(l)})$ nonlinear function:

$$Y_i^{(l)} = f\left(Z_i^{(l)}\right) \text{ with } Z_i^{(l)} = \sum_{j=1}^{m_i^{(l-1)}} w_{i,j}^{(l)} \times y_j^{(l-1)}. \quad (3)$$

(a) RGB color space live

(b) YCbCr color space live

(c) HSV color space live

(d) RGB color space spoof

(e) YCbCr color space spoof

(f) HSV color space spoof

FIGURE 3: Different color spaces based on Replay-Attack databases.

*3.2.1. Pretrained Models.* To modify the pretrained experiment models for face spoofing recognition, the models were fine-tuned by spoofing databases. The binary classification was utilized for spoofing detection problems and changing the output of the classification layer to two classes of spoof and real face.

After modifying the SoftMax classification layer based on the spoofing database in the training phase, the VGG16 and VGG-face models were fine-tuned based on the spoofing database. The VGG-face model is one of the popular pretrained models for face recognition systems. This model was developed by the Oxford Visual Geometry Group [33]. The model was trained by 2.6 M to face images in the RGB color space, and the default size of an input image is 224 × 224 [18]. This model contains five max pooling, thirteen convolutional layers with the rectified linear unit (ReLU) function, and three fully connected layers, namely, FC6, FC7, and FC8. The last fully connected layer (FC8) modifies from 2622 (face image classes) to 2 classes of spoof and real. The architecture of the VGG-face model is a variant of VGG16, which is trained by face images, as presented in Table 1. In this approach, the fine-tuned VGG-face and VGG-16 models based on the face spoofing database are utilized as a deep feature extractor. The deep features are taken from FC7 (seventh fully connected layer), the last layer before the output layer. The activation values of this FC layer for all models are set as default values equal to 4096 (dimensional feature vectors) for the input images.

*3.3. Feature Selection.* The main purpose of the mRMR method is to select the subset of features which has the most correlation with the class and reduce irrelevant and redundancy features based on mutual information [38, 39]. Measurement of the mutual information of $I$ between two $x$ and $y$ attributes is defined based on

$$I(x, y) = \sum_{i,j} p\left(x_i, y_j\right) \log \frac{p\left(x_i, y_j\right)}{p(x_i)p\left(y_j\right)}, \tag{4}$$

TABLE 1: VGG16 architecture.

| Layer | Patch size/stride | Input size |
|---|---|---|
| Conv × 2 | 3 × 3/1 | 64 × 224 × 224 |
| Pool | 2 × 2 | 64 × 224 × 224 |
| Conv × 2 | 3 × 3/1 | 128 × 112 × 112 |
| Pool | 2 × 2 | 128 × 112 × 112 |
| Conv × 3 | 3 × 3/1 | 256 × 56 × 56 |
| Pool | 2 × 2 | 256 × 56 × 56 |
| Conv × 3 | 3 × 3/1 | 512 × 28 × 28 |
| Pool | 2 × 2 | 512 × 28 × 28 |
| Conv × 3 | 3 × 3/1 | 512 × 14 × 14 |
| Pool | 2 × 2 | 512 × 14 × 14 |
| FC | 25088 × 4096 | 25088 |
| FC | 4096 × 4096 | 4096 |

where $p(x_i)$ and $p(y_j)$ represent the marginal probabilities and $p(x_i, y_j)$ represents the joint probabilistic distribution. Let us define each property of the equation as $F_i$ in a $K$-size vector ($F_i = [F_{1i}, F_{2i}, F_{3i}, \cdots, F_{Ki}]$). In this case, the mutual information of the variables $(i, j)$ is defined as $I(F_i, F_j)$. In order to find the best features of the selected subset, Equations (5) and (6) must be satisfied. The minimum redundancy feature is presented in Equation (8), and the maximum relevance condition is presented in Equation (6):

$$\min W, W = \frac{1}{|s|^2} \sum_{F_i F_j} I(F_i, F_j), \tag{5}$$

$$\max V, V = \frac{1}{|s|} \sum_{F_i F_j} I(H, F_i), \tag{6}$$

where $H$ represents the class label and $s$ shows the number of features selected. The mRMR feature set is obtained by optimizing the combination of feature selection criteria, namely, Mutual Information Difference (MID) and Mutual Information Quotient (MIQ), which are presented in

$$\begin{cases} \text{MID} = \max{(v - w)}, \\ \text{MIQ} = \max{\left(\dfrac{v}{w}\right)}. \end{cases} \tag{7}$$

For optimizing the MID and MIQ conditions, it is required to combine them into a single criterion function [40], as shown in the following equation:

$$f_{\text{mRMR}}(X_i) = I(H, F_i) - \frac{1}{|s|} \sum_{F_i F_j} I(F_i, F_j), \tag{8}$$

where $I(H, F_i)$ measures the relevance feature to be added for the class and $1/|s|\sum_{F_i F_j} I(F_i, F_j)$ estimates the redundancy of features with respect to previously selected $s$ features. These selected features are classified with a linear regression classification algorithm for detection of face presentation attacks.

## 4. Experimental Results

The proposed method as shown in Figure 2 was compiled with an NVIDIA GeForce 4 GB graphics card (GPU). Other hardware details were Intel Core i5 3.6 GHz processor and 16 GB RAM. As presented in Table 2, these parameters were used with their default values. Additionally, the minibatch size was set as 32.

### 4.1. Experimental Databases

*4.1.1. The Replay-Attack Database [41].* The Replay-Attack database consists of 1300 videos of 2D face attacks under different conditions. This database contains three main subgroups for training, validation, and testing folders with names of training data, development data, and test data. Two main different lighting conditions in this database were named as controlled and adverse. The controlled scenario

data were collected under homogeneous backgrounds and with office lights turned on, and the adverse data were collected with more complex backgrounds and without office lights as presented in Figure 4.

*4.1.2. ROSE-Youtu Face Liveness Detection Dataset [42].* This database contains a large variety of illumination conditions, cameras with different resolutions, and types of attacks such as display, print, and mask attacks. The ROSE-Youtu database contains 4225 videos with 25 subjects, and each video duration average is around 10 seconds. The ROSE-Youtu database is divided into two subsets of training and testing. The first 10 indexed units are separate for training, and the rest of the videos belong to testing. The numbers of samples from this database are presented in Figure 5.

*4.2. Evaluation Metric.* To measure the performance of the models, accuracy (Acc), sensitivity (Se), specificity (Sp), precision (Pr), and $F$-score metrics derived from the confusion matrix were used, and the formulations of the metrics were as follows [43]:

$$\begin{cases} \text{Acc} = \dfrac{(\text{TP} + \text{TN})}{(\text{TF} + \text{FN}) + (\text{FP} + \text{TN})}, \\ \text{Se} = \dfrac{(\text{TP})}{(\text{TP} + \text{FN})}, \\ \text{Pr} = \dfrac{(\text{TP})}{(\text{TP} + \text{FP})}, \\ F\text{-score} = \dfrac{(2 \times \text{TP})}{(2 \times \text{TP} + \text{FP} + \text{FN})}. \end{cases} \tag{9}$$

To evaluate our new approach against state-of-the-art methods, we applied the formula of the Half Total Error Rate (HTER) in

$$\text{HTER} = \frac{\text{FRR}(\mathscr{K}, \mathscr{D}) + \text{FAR}(\mathscr{K}, \mathscr{D})}{2}, \tag{10}$$

where $\text{FRR}(\mathscr{K}, \mathscr{D})$ is a false rejection rate, $\mathscr{D}$ denotes the used database, and $\mathscr{K}$ is estimated on the equal error rate (EER). In this context $\text{FAR}(\mathscr{K}, \mathscr{D})$ stands for the False Acceptance Rate.

*4.3. Fine-Tuning VGG-Face Model for Face Spoofing Detection.* Our face spoofing recognition approach in the first steps was based on the VGG-face model. The VGG-face model is trained by a large database of face images. As presented in Figure 6, each convolution block contains the rectified linear unit (ReLU) function and a $3 \times 3$ kernel size. Also, each convolution block contains a max pooling layer with a kernel size of $2 \times 2$. Two FC layers are set with 4096 channels with the ReLU function and batch normalization. The last FC layer contains the ReLU function, batch normalization, and the SoftMax activation function where the output of this layer presents categorical distribution over face spoofing recognition labels.

TABLE 2: Parameter values of the proposed approach used in this study.

| Software | Optimization | Activation function | Momentum | Decay | Minibatch | Learning rate |
|---|---|---|---|---|---|---|
| Keras | Adam | ReLU | 0.9 | $1e-6$ | 32 | 0.01 |



(a) Live face



(b) Spoof face

FIGURE 4: Replay-Attack database samples for live and spoof images.



(a) Live face



(b) Spoof face

FIGURE 5: ROSE-Youtu face liveness detection samples for live and spoof face images.

The performance of the VGG-face model for face spoofing detection databases depends on the level of fine-tuning of the convolutional blocks. For this reason, in this test, we evaluated the effects of each pretrained convolutional block on the accuracy of the model [14]. Different models arranged based on the retrained and frozen levels of the parameters of the network with names of the A, B, C, and D models are presented in Figure 7. Five convolution blocks with the names of Conv1, Conv2, Conv3, Conv4, and Conv5 and two FC layers were trained based on the level of fine-tuning. For example, the first model (A) consisted of the Conv2-5 and FC layers, which means that the convolutional blocks from 2 up to 5 were trained based on new datasets, and the rest of the parameters of the model were frozen. In the same way, the models B, C, and D were, respectively, trained from the third, fourth, and fifth convolutional blocks with the fully connected layer.

Based on the experimental results presented in Figure 8, the best accuracy was for model A (Conv2-5 and FC layers) with 97.99% and 82%, respectively, for the Replay-Attack and ROSE-Youtu databases which were highlighted with gray shading. All models (A, B, C, and D) were trained based on the parameters presented in Table 2 and 1000 epochs. Additionally, for the classification of the images, the SoftMax classifier was utilized with two channels of live and spoof labels. As a result, for the Replay-Attack and ROSE-Youtu databases, model A stayed on the best accuracy, respectively, with (97%, 82%) compared to B (96%, 66%), C (96%, 76%), and D (92%, 66%). Based on these experimental results, it may be proven that, for spoofing detection based on the RGB color space, the optimum level of fine-tuning of the VGG-face model was the trained convolutional blocks numbered 2 up to 5 with two fully connected layers and by freezing the first convolutional block parameters.

FIGURE 6: Structure of VGG-face model [33].



FIGURE 7: Green shaded blocks are frozen and pretrained, and blue shaded blocks are retrained during the training process.

In this case, in the rest of the experimental results for the RGB color space, we utilized the same level of fine-tuning (model A) which stayed on the best accuracy rate for the VGG-face model for face spoofing detection. For training the deep models with the ROSE-Youtu database, we selected 70% of the data from the first 10 indexed samples of data for training, and the rest of these were used for validation. In this case, the training and validation data were totally separated. Because the ROSE-Youtu database contains data with different rotations such as 90 degrees clockwise and counterclockwise, the image data augmentation technique in the Keras library was utilized.

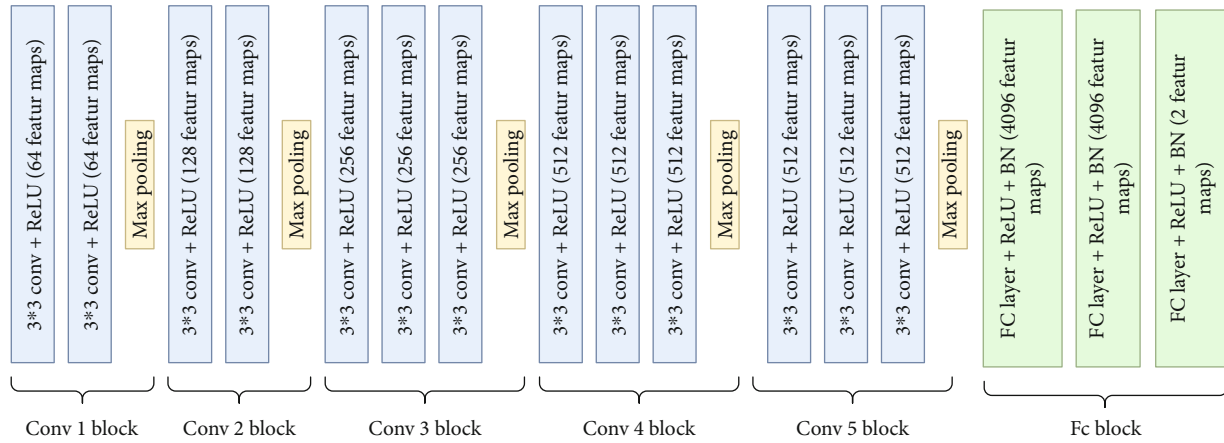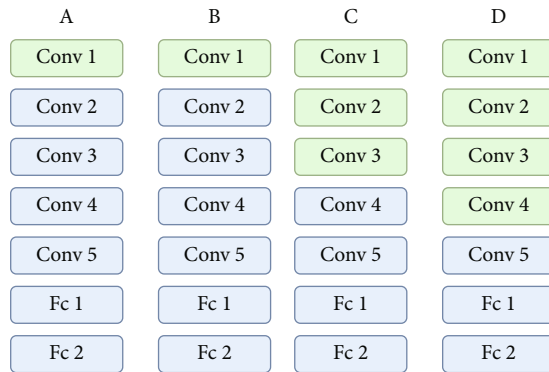*4.4. Color-Based Approach Model.* In this section, we explain the process of converting the color space from RGB to HSV and YCbCr. Furthermore, we evaluated three benchmark VGG models for finding the effects of each color space on the accuracy of classification. In this test, we utilized two VGG16 models and trained the entirety of each network with HSV and YCbCr color space images from spoofing datasets with the default window size. All models were trained based on the parameters presented in Table 2 and 1000 epochs.

Table 3 presents the experimental results on the HSV and YCbCr color spaces and the evaluation of fine-tuning of the entirety of the networks with these color spaces. According

to the results obtained, the HSV color space-based image in the Replay-Attack database achieved significant results compared to the YCbCr color space by improving 0.71% in accuracy. Nevertheless, in the ROSE-Youtu database, the YCbCr space provided better results compared to HSV by improving 7.59%. According to these results, it may be concluded that, for face spoofing recognition under different conditions such as illumination changes and displaying a high-resolution camera, both color spaces contain discriminative features which can help distinguish a live image from a fake face in different scenarios.

*4.5. Deep Feature Extraction.* In the second step of our experimental procedure, the features of the fully connected layer (FC7) of the pretrained VGG-face model based on the RGB color space were extracted, which included 4096 channels. The features extracted from this layer were classified with different typical classifiers such as SVM, LDA, and KNN. Moreover, these results were compared to the SoftMax classifier to evaluate the performance of the extracted deep features with other classification algorithms. Based on the experimental results shown in Table 4, the best results were for SVM and KNN in the Replay-Attack database with 98.93 (Acc), 98.50 (Se), 100 (Sp), 98.97 (Pr), and 98.93% (F-score) for both classification algorithms.

In the Replay-Attack database, the SoftMax classifier was placed on the fourth stage among the other classifiers based on the results. However, in the ROSE-Youtu database, the SoftMax classifier achieved significant results compared to the other classifiers with 82.84 (Acc), 97.42 (Se), 72.41 (Sp), 89.52 (Pr), and 88.00% (F-score).

*4.6. Feature Selection and Classification.* In this step, we utilized mRMR to reduce the size of the extracted features from three different models and select robust and discriminative feature sets. The size of the extracted features for each model was 4096, and by combining these three VGG models, the size increased to 12288 features. For finding the optimum dimension of feature sets, we analyzed different sizes of features with the help of mRMR feature selection as presented in Figures 9(a) and 9(b). Based on the results, the best feature
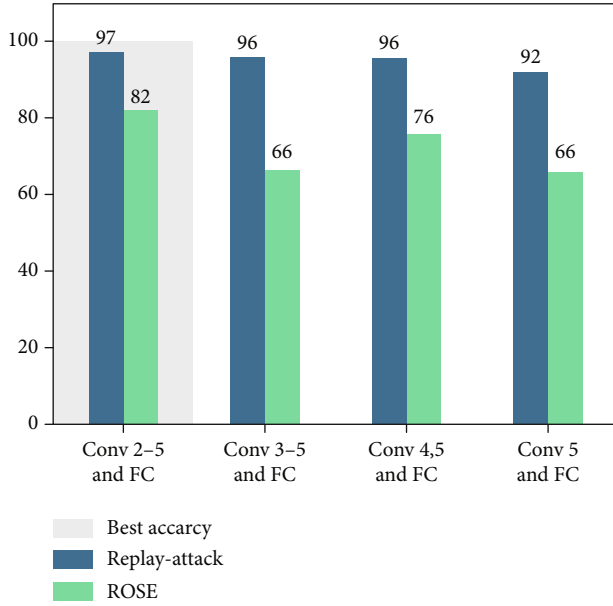
Figure 8: Accuracy of VGG-face model based on level of fine-tuning of the networks.

Table 3: Experimental results of fine-tuning pretrained VGG16 models with the HSV and YCbCr color spaces.

| Metrics (%) | HSV | | YCbCr | |
|---|---|---|---|---|
| | Replay-Attack | ROSE-Youtu | Replay-Attack | ROSE-Youtu |
| Acc | 99.46 | 71.94 | 98.75 | 79.53 |
| Se | 99.25 | 77.42 | 99.25 | 88.61 |
| Sp | 100 | 66.67 | 97.50 | 45.77 |
| Pr | 99.47 | 71.87 | 98.75 | 83.75 |
| F-score | 99.47 | 71.87 | 98.75 | 71.03 |

sizes for Replay-Attack were 400, 500, and 700, and those for the ROSE-Youtu database were 300, 500, and 700, respectively, for RGB, HSV, and YCbCr based on the LR classifier. In this case, the optimum feature size for covering both databases and all color spaces may be set to 1600 features. In continuation of this test, we analyzed the effects of the deep features of HSV color spaces on the improvement of accuracy rates. In this case, we combined extracted features from the FC7 layer of the pretrained VGG-face model (RGB) with the VGG16 model (HSV). The experimental results presented in Table 5 show that the accuracy of the face spoofing detection approach was improved drastically in the Replay-Attack database.

In this database, all evaluation metrics with the LR, SVM, and KNN classifiers stayed on significant rates with 99.82 (Acc), 99.75 (Se), 100 (Sp), 99.82 (Pr), and 99.82 (F-score) %. In the ROSE-Youtu database, also, all evaluation metrics were improved with four different classifiers, and the best results were obtained for the linear regression classifier by 95.98 (Acc), 99.00 (Se), 93.24 (Sp), 95.98 (Pr), and 95.98 (F-score) %. The experimental results in this table compared to Table 4 showed that HSV deep features improved the

effectiveness of detection of spoofing data. The comparison of two experimental results of Tables 4 and 5 showed that all evaluation metrics were improved by combining HSV deep features with VGG-face deep features, and these results were improved by 13.14 (Acc), 1.58 (Se), 20.83 (Sp), 6.46 (Pr), and 7.98 (F-score) based on the LR classifier in the ROSE-Youtu database.

In Table 6, the experimental results of the proposed deep model by applying the feature selection method are presented. After concatenation of three extracted features from different color spaces from the VGG models, mRMR feature selection was applied. As discussed in Section 3.3, the main reason for applying the mRMR algorithm was to reduce the irrelevant features and select robust and discriminative features. Figure 10 presents a visualization of the first four feature maps of each five convolutional blocks with the RGB, HSV, and YCbCr color spaces. According to the extracted features from each convolutional block and specifically the fifth convolutional block, it was obtained that combining features from each model with different color spaces includes redundant and irrelevant features which decrease the effectiveness of our proposed approach. Based on these results in Table 6, the extracted YCbCr features cannot improve the evaluation metrics in the replay-attack database. However, on the other hand, these features improved the effectiveness of recognition of spoofing data in the ROSE-Youtu database and increased the results by 1.18 (Acc), 2.91 (Sp), 2.25 (Pr), and 1.19 (F-score) based on the LR classifier. Additionally, the linear regression classifier stayed on the best results compared to SVM, KNN, and LDA.

To better present the results, we utilized ROC curve analysis for both experiment databases as shown in Figure 11. The ROC curve analysis showed that the proposed approach with the help of well-known pretrained models in the RGB, HSV, and YCbCr color spaces extracted discriminative features for the detection of spoofing face images. Based on these results, the LR classifiers stayed on the best AUC compared to the other mentioned classification algorithms by 0.995 and 1.00 for the ROSE-Youtu (Figure 11(b)) and Replay-Attack (Figure 11(a)) databases, respectively. In this case, we selected the LR classifier as the base classification algorithm for our proposed approach and employed this classification algorithm in the rest of the paper.

*4.7. Evaluation of Different Attacks.* For evaluation of our proposed approach in different scenarios of spoofing attacks and for finding the advantages and disadvantages of our proposed approach, we tested our deep learning approach on different attacks individually. Based on the experimental results on Replay-Attack (Table 6), it may be concluded that our proposed approach had satisfactory results in the replay, display, and print attacks which are presented in the Replay-Attack database. Furthermore, this approach achieved 97.16% accuracy in the ROSE-Youtu database, in which, for finding misclassification reasons, in this test, the spoofing scenarios were individually analyzed. We categorized the ROSE-Youtu database into five different groups such as the real, display and print, mask with cropping, and mask

TABLE 4: The classification results based on different classifiers and deep features of the VGG-face model on the Replay-Attack and ROSE-Youtu databases.

| Model | Database | Classification | Acc (%) | Se (%) | Sp (%) | Pr (%) | F-score (%) |
|---|---|---|---|---|---|---|---|
| VGG-face (RGB color space) | Replay-Attack database | SoftMax | 97.32 | 99.25 | 99.50 | 97.34 | 97.30 |
| | | SVM | 98.93 | 98.50 | 100 | 98.97 | 98.93 |
| | | LDA | 98.91 | 98.91 | 100 | 99.78 | 99.78 |
| | | KNN ($K = 1$) | 98.93 | 98.50 | 100 | 98.97 | 98.93 |
| | ROSE-Youtu | SoftMax | 82.84 | 97.42 | 72.41 | 89.52 | 88.00 |
| | | SVM | 78.38 | 59.75 | 90.03 | 78.46 | 77.65 |
| | | LDA | 70.30 | 50.13 | 82.91 | 69.61 | 69.39 |
| | | KNN ($K = 1$) | 78.38 | 59.75 | 90.03 | 78.46 | 77.65 |



(a) Replay-Attack database

(b) ROSE-Youtu database

FIGURE 9: Accuracy of LR classification based on different sizes of features.

TABLE 5: The classification results of the extracted features from RGB and HSV on the Replay-Attack and ROSE-Youtu databases.

| Model | Databases | Classification | Acc (%) | Se (%) | Sp (%) | Pr (%) | F-score (%) |
|---|---|---|---|---|---|---|---|
| VGG-face (RGB)+VGG16 (HSV) | Replay-Attack database | LR | 99.82 | 99.75 | 100 | 99.82 | 99.82 |
| | | SVM | 99.82 | 99.75 | 100 | 99.82 | 99.82 |
| | | LDA | 98.75 | 99.50 | 96.88 | 98.75 | 98.75 |
| | | KNN($K = 1$) | 99.82 | 99.75 | 100 | 99.82 | 99.82 |
| | ROSE-Youtu | LR | 95.98 | 99.00 | 93.24 | 95.98 | 95.98 |
| | | SVM | 95.98 | 97.51 | 94.59 | 96.04 | 95.98 |
| | | LDA | 83.34 | 77.11 | 92.79 | 85.97 | 85.22 |
| | | KNN($K = 1$) | 94.79 | 97.51 | 92.34 | 94.96 | 94.80 |

without cropping groups containing videos from persons as presented in Figure 12. Display and replay attacks are already tested in different conditions such as light change and shaking hands in experimental databases, namely, Replay-Attack. We set the displayed attack and print attack catego-ries together and labeled them as display. However, the main difference of the ROSE-Youtu database is mask attack in different conditions and scenarios which are not available in other experimental databases. Mask attack in the ROSE-Youtu database contains scenarios such as a mask with two

TABLE 6: The classification results of the extracted features from RGB and HSV and YCbCr.

| Model | Databases | Classification | Acc (%) | Se (%) | Sp (%) | Pr (%) | F-score (%) |
|---|---|---|---|---|---|---|---|
| VGG-face (RGB)+VGG16 (HSV)+VGG16 (YCbCr) | Replay-Attack database | LR | 99.82 | 99.75 | 100 | 99.82 | 99.82 |
| | | SVM | 99.82 | 99.75 | 100 | 99.82 | 99.82 |
| | | LDA | 98.75 | 99.50 | 96.88 | 98.75 | 98.75 |
| | | KNN($K = 1$) | 99.82 | 99.75 | 100 | 99.82 | 99.82 |
| | ROSE-Youtu | LR | 97.16 | 98.41 | 96.15 | 97.21 | 97.17 |
| | | SVM | 95.98 | 93.12 | 98.29 | 96.05 | 95.97 |
| | | LDA | 96.45 | 97.73 | 95.73 | 96.49 | 96.46 |
| | | KNN ($K = 1$) | 88.17 | 86.77 | 89.32 | 88.18 | 88.18 |



FIGURE 10: Extracted feature maps from each convolutional block.

eyes and mouth cropped out, mask without cropping, mask with the upper part cut in the middle, and mask with the lower part cut in the middle.

In this test, we categorized these mask attack scenarios into two main groups as a mask without cropping and mask with cropping. Based on the experimental results which are presented in Table 7, it appeared that the main advantage of the proposed approach was the detection of spoofing attacks such as display and print attacks. The accuracy of recognition of display and print attacks was 98.00%, which stayed on the highest value compared to other spoofing data. The second highest value of accuracy was for the mask with cropping attacks with 97.82% accuracy. The results for replay attacks were also compatible with 94.64% accuracy. On the other hand, the lowest results were for a mask without cropping with 92.59 (Acc), 96.81 (Se), 98.93 (Sp), 92.70 (Pr), and 92.33 (F-score) %. These results proved that the proposed approach has a significant accuracy in recognition of display and printed attack and compatible accuracy in a mask without cropping scenarios.

In continuation of this test, we utilized the scatter plot of the extracted features based on the attack groups and real videos. In this part, we selected one frame from each video from the test set and reduced the dimensions of the features with the help of Principal Component Analysis (PCA) from 1600 to 3 to obtain the $X$, $Y$, and $Z$ values for each image and present them in 3D scatter plots. As presented in Figure 13, it appeared that the mask without cropping and replay attack features were overlapped with real video frames. Furthermore, other spoofing attacks such as display and mask with cropping were clearly separated from real videos.

*4.8. Evaluation Efficiency of the Proposed Method in Cloud System.* As presented before, one of the main problems of cloud computing systems is the management of storing data and optimizing resources. For this reason, we proposed a deep learning approach that trains with fewer data and achieved significant results based on accuracy compared to existing models. For evaluating our approach, we train the model in four different types. First, the models are trained

(a) Replay-Attack database

(b) ROSE-Youtu database

Figure 11: ROC curve analysis based on different classifiers.



(a) Display and print attacks

(b) Mask with cropping

(c) Mask without cropping

(d) Replay attack

Figure 12: Categorization of spoofing attacks of the ROSE-Youtu database.

TABLE 7: Evaluation of different types of attacks on the ROSE-Youtu database.

| Database | Types of attacks | Acc (%) | Se (%) | Sp (%) | Pr (%) | F-score (%) |
|---|---|---|---|---|---|---|
| ROSE-Youtu | Mask without cropping | 92.59 | 96.81 | 98.93 | 92.70 | 92.33 |
| | Replay attack | 94.64 | 90.99 | 96.81 | 94.64 | 94.63 |
| | Mask with cropping | 97.82 | 95.83 | 98.89 | 97.83 | 97.82 |
| | Display and print attack | 98.00 | 96.46 | 98.93 | 98.00 | 98.00 |



(a) Replay-Attack database



● Mask without cropping     ● Display
● Replay     ● Real
● Mask with cropping

(b) ROSE-Youtu database

FIGURE 13: 3D and 2D scatter plots of features based on attacks.

on 10% of frames of each video and test on all frames. The second, third, and fourth modes of evaluation are in the same condition, such as 20, 30, and 40% of the frames for training and evaluating on all frames of test sets. These scenarios are tested on well-known deep learning models in RGB color space such as Inception V3 [44], InceptionResNetV2 [45], and VGG 19 [34]. These pretrained models on the ImageNet database are employed as a deep feature extractor. For fine-

(a) Replay-Attack database

(b) ROSE-Youtu database

FIGURE 14: Evaluation of different sizes of training data on the accuracy of classification.

TABLE 8: Comparison of the proposed approach against state-of-the-art algorithms based on the Replay-Attack database.

| Method | EER (%) | HTER (%) |
| --- | --- | --- |
| Motion+LBP [45] | 4.5 | 5.1 |
| DMD [26] | 3.8 | 5.3 |
| SURF color texture [10] | 1.2 | 4.2 |
| Color texture [11] | 0.4 | 2.8 |
| LBP net [18] | 0.6 | 1.3 |
| Color LBP [46] | 0.9 | 4.9 |
| Partial CNN [14] | 2.9 | 4.3 |
| CompactNet [12] | 0.8 | 0.7 |
| Dense optical flow+Shearlet [31] | 0.83 | 0.0 |
| Proposed method | 0.2 | 0.4 |

TABLE 9: Comparison of the proposed approach against state-of-the-art algorithms based on the ROSE-Youtu database.

| Method | EER (%) |
| --- | --- |
| Deep color-based feature [42] | 8.0 |
| SE-ResNet 18 [48] | 7.2 |
| 3D CNN [49] | 7.0 |
| Two-stage deep model [47] | 4.56 |
| Proposed method | 3.8 |

tuning of parameters in these models with the face spoofing database, we changed the SoftMax classification layer to two classes of spoof and real face. In addition, the small number of learning rates with 0.0001 is set for all models; besides, we employed Adam optimization, batch size 16, and 10000 epochs. Suppose we capture a one-minute video with 720 p resolution at 30 fps containing 1800 frames which are around

60 MB. Therefore, training the model with 10% of frames of each video not only reduced the size of data for training (around 6 MB) but also decreased the computation cost in the training phase.
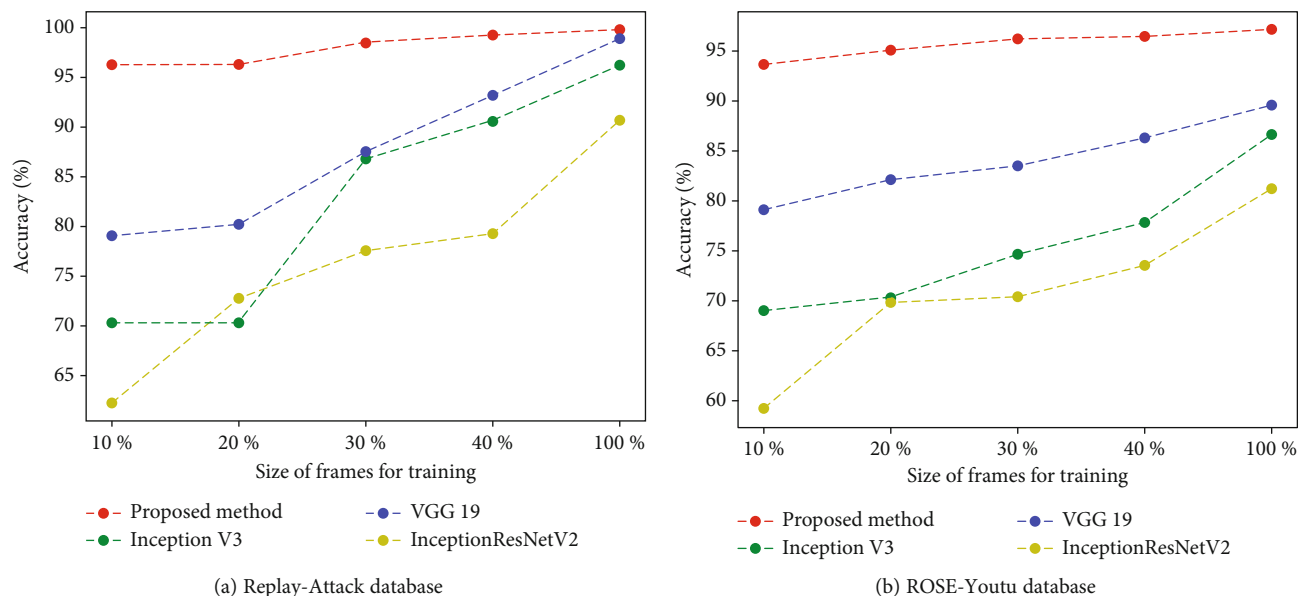
Based on experimental results presented in Figure 14, it appears that the proposed method achieved significant results in the detection of spoofing attacks with less training data compared to benchmark deep learning methods. The proposed method achieved the accuracy of classification with 96.3% in ten percent of frames of each video for training and testing on entire videos which this score is better than the results achieved by Inception V3, InceptionResNetV2, and VGG 19 with 70.32, 62.3, and 79.1, respectively, in the Replay-Attack database. The results of the proposed method are 96.3, 96.3, 98.5, 99.2, and 99.8% which are better than other experimented deep learning methods, respectively, for 10, 20, 30, 40, and 100% of frames of each video in the Replay-Attack database. In the same condition, in the ROSE-Youtu database, also, our proposed method stayed on the best results with 93.7, 95.1, 96.2, and 96.5 in 10, 20, 30, and 40 percent of frames of each video for training.

*4.9. Comparison of the Proposed Approach against State-of-the-Art Algorithms.* Table 8 provides a comparison between the proposed approach and state-of-the-art methods. The experimental results shown in Table 8 demonstrated the effectiveness of our extracted deep features in the Replay-Attack database.

We may observe that, among the state-of-the-art methods presented in this table, the best results were for deep learning-based methods like the LBP net [18] with 0.6 (EER) and 1.3 (HTER). The best HTER was for dense optical flow +Shearlet [31] with 0.0. Furthermore, our proposed method achieved 0.2 (EER), which was better than the multicue deep method proposed in a previous study [31] with a single cue (color texture analysis).

Table 9 provides a comparison between the proposed approach and state-of-the-art methods in the aspect of EER. According to these experimental results, it may be argued that our proposed approach is more applicable and stayed on the best EER (%) values in comparison to state-of-the-art methods in the Replay-Attack database. In the other benchmark public access database (ROSE-Youtu), our proposed approach also stayed on the best ERR (%) values. In this database, the best EER in state-of-the-art algorithms was for the two-stage deep model [47] approach over which our proposed approach improved the EER value by 0.76%. Based on these experimental results and comparison with state-of-the-art algorithms, it may be concluded that our proposed approach achieved robust and significant results for distinguishing fake faces from live faces with 0.2 and 3.8 for EER (%) in the replay-attack and ROSE-Youtu databases, respectively.

## 5. Conclusion

The IoT cloud-based framework for face spoofing detection is proposed and implemented in this study. The proposed system detects face spoofing attacks by applying the new deep learning framework. This approach can be used reliably in the cloud-based environment by storing less data which decreased both processing cost and size of data in the training phase. Moreover, the proposed multicolor deep feature-based approach outperformed the baseline methods on the Replay-Attack database, while achieving competitive results on the ROSE-Youtu database. The results obtained for the Replay-Attack and ROSE-Youtu databases proved that environmental factors and scenarios such as background changes, shaking hands, high-resolution camera, and illumination did not limit the effectiveness of our proposed approach. Furthermore, our proposed approach achieved satisfactory results in scenarios such as print, display, and replay attacks. In the case of mask attacks in different scenarios such as without cropping, with cropping, upper part cut, lower part cut, and mask with two eyes and mouth cropped out, the proposed approach presented compatible results. Furthermore, in mask without cropping attacks, the proposed approach achieved the lowest rate of accuracy (92.59%) compared to different attacks such as replay or print attacks. This inefficiency of the proposed approach in mask attack types makes us eager to solve this problem in future work. In future work, we will investigate adding depth information to our color-based deep features to improve the effectiveness of recognition of spoofing attacks in different mask scenarios in IoT cloud environments.

## Data Availability

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

## Conflicts of Interest

The authors declare no conflict of interest.

## References

[1] Z. Ali, M. S. Hossain, G. Muhammad, I. Ullah, H. Abachi, and A. Alamri, "Edge-centric multimodal authentication system using encrypted biometric templates," *Future Generation Computer Systems*, vol. 85, pp. 76–87, 2018.

[2] M. Gomez-Barrero, E. Maiorana, J. Galbally, P. Campisi, and J. Fierrez, "Multi-biometric template protection based on homomorphic encryption," *Pattern Recognition*, vol. 67, pp. 149–163, 2017.

[3] P. Kumari and P. Thangaraj, "A fast feature selection technique in multi modal biometrics using cloud framework," *Microprocessors and Microsystems*, vol. 79, p. 103277, 2020.

[4] K. A. Shakil, F. J. Zareen, M. Alam, and S. Jabin, "BAMHealthCloud: a biometric authentication and data management system for healthcare data in cloud," *Journal of King Saud University - Computer and Information Sciences*, vol. 32, no. 1, pp. 57–64, 2020.

[5] B. Sree Vidya and E. Chandra, "Entropy based local binary pattern (ELBP) feature extraction technique of multimodal biometrics as defence mechanism for cloud storage," *Alexandria Engineering Journal*, vol. 58, no. 1, pp. 103–114, 2019.

[6] M. Masud, G. Muhammad, H. Alhumyani et al., "Deep learning-based intelligent face recognition in IoT-cloud environment," *Computer Communications*, vol. 152, pp. 215–222, 2020.

[7] X. Song, X. Zhao, L. Fang, and T. Lin, "Discriminative representation combinations for accurate face spoofing detection," *Pattern Recognition*, vol. 85, pp. 220–231, 2019.

[8] G. Pan, L. Sun, Z. Wu, and S. Lao, "Eyeblink-based anti-spoofing in face recognition from a generic webcamera," in *2007 IEEE 11th International Conference on Computer Vision*, Rio de Janeiro, Brazil, 2007.

[9] A. Gumaei, R. Sammouda, A. M. S. Al-Salman, and A. Alsanad, "Anti-spoofing cloud-based multi-spectral biometric identification system for enterprise security and privacy-preservation," *Journal of Parallel and Distributed Computing*, vol. 124, pp. 27–40, 2019.

[10] Z. Boulkenafet, J. Komulainen, and A. Hadid, "On the generalization of color texture-based face anti-spoofing," *Image and Vision Computing*, vol. 77, pp. 1–9, 2018.

[11] Z. Boulkenafet, J. Komulainen, and A. Hadid, "Face spoofing detection using colour texture analysis," *IEEE Transactions on Information Forensics and Security*, vol. 11, no. 8, pp. 1818–1830, 2016.

[12] L. Li, Z. Xia, X. Jiang, F. Roli, and X. Feng, "CompactNet: learning a compact space for face presentation attack detection," *Neurocomputing*, vol. 409, pp. 191–207, 2020.

[13] D. T. Nguyen, T. D. Pham, N. R. Baek, and K. R. Park, "Combining deep and handcrafted image features for presentation attack detection in face recognition systems using visible-light camera sensors," *Sensors (Switzerland)*, vol. 18, no. 3, p. 699, 2018.

[14] L. Li, X. Feng, Z. Boulkenafet, Z. Xia, M. Li, and A. Hadid, "An original face anti-spoofing approach using partial convolutional neural network," in *2016 Sixth International Conference on Image Processing Theory, Tools and Applications (IPTA)*, pp. 16–21, Oulu, Finland, 2017.

[15] Y. A. U. Rehman, L. M. Po, and J. Komulainen, "Enhancing deep discriminative feature maps via perturbation for face presentation attack detection," *Image and Vision Computing*, vol. 94, p. 103858, 2020.

[16] Y. A. U. Rehman, L. M. Po, M. Liu, Z. Zou, W. Ou, and Y. Zhao, "Face liveness detection using convolutional-features fusion of real and deep network generated face images," *Journal of Visual Communication and Image Representation*, vol. 59, pp. 574–582, 2019.

[17] F. Peng, L. Qin, and M. Long, "Face presentation attack detection based on chromatic co-occurrence of local binary pattern and ensemble learning," *Journal of Visual Communication and Image Representation*, vol. 66, p. 102746, 2020.

[18] L. Li, X. Feng, Z. Xia, X. Jiang, and A. Hadid, "Face spoofing detection with local binary pattern network," *Journal of Visual Communication and Image Representation*, vol. 54, pp. 182–192, 2018.

[19] G. B. De Souza, S. Member, D. Felipe, R. G. Pires, A. N. Marana, and J. P. Papa, "Deep texture features for robust face spoofing detection," *IEEE Trans. Circuits Syst. II Express Briefs*, vol. 64, no. 12, pp. 1397–1401, 2017.

[20] R. J. Raghavendra and R. S. Kunte, "Extended local ternary co-relation pattern: a novel feature descriptor for face anti-spoofing," *J. Inf. Secur. Appl.*, vol. 52, p. 102482, 2020.

[21] S. Jia, C. Hu, X. Li, and Z. Xu, "Face spoofing detection under super-realistic 3D wax face attacks," *Pattern Recognition Letters*, vol. 145, pp. 103–109, 2021.

[22] A. George and S. Marcel, "Learning one class representations for face presentation attack detection using multi-channel convolutional neural networks," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 361–375, 2021.

[23] J. M. Di Martino, Q. Qiu, and G. Sapiro, "Rethinking shape from shading for spoofing detection," *IEEE Transactions on Image Processing*, vol. 30, pp. 1086–1099, 2021.

[24] S. Arora, M. P. S. Bhatia, and V. Mittal, "A robust framework for spoofing detection in faces using deep learning," *The Visual Computer*, vol. 13, 2021.

[25] T. Edmunds and A. Caplier, "Motion-based countermeasure against photo and video spoofing attacks in face recognition," *Journal of Visual Communication and Image Representation*, vol. 50, pp. 314–332, 2018.

[26] S. Tirunagari, N. Poh, D. Windridge, A. Iorliam, N. Suki, and A. T. S. Ho, "Detection of face spoofing using visual dynamics," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 4, pp. 762–777, 2015.

[27] D. Wen, H. Han, and A. K. Jain, "Face spoof detection with image distortion analysis," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 4, pp. 746–761, 2015.

[28] A. Pinto, H. Pedrini, W. R. Schwartz, and A. Rocha, "Face spoofing detection through visual codebooks of spectral temporal cubes," *IEEE Transactions on Image Processing*, vol. 24, no. 12, pp. 4726–4740, 2015.

[29] Y. A. U. Rehman, L. M. Po, and M. Liu, "SLNet: stereo face liveness detection via dynamic disparity-maps and convolutional neural network," *Expert Systems with Applications*, vol. 142, p. 113002, 2020.

[30] N. Erdogmus and S. Marcel, "Spoofing in 2D face recognition with 3D masks and anti-spoofing with Kinect," in *2013 IEEE Sixth International Conference on Biometrics: Theory, Applications and Systems (BTAS)*, Arlington, VA, USA, 2013.

[31] L. Feng, L.-M. Po, Y. Li et al., "Integration of image quality and motion cues for face anti-spoofing: a neural network approach," *Journal of Visual Communication and Image Representation*, vol. 38, pp. 451–460, 2016.

[32] Y.-Q. Wang, "An analysis of the Viola-Jones face detection algorithm," *Image Processing On Linec*, vol. 4, pp. 128–148, 2014.

[33] O. M. Parkhi, A. Vedaldi, and A. Zisserman, "Deep face recognition," in *Procedings of the British Machine Vision Conference 2015*, Swansea, UK, 2015.

[34] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*, pp. 1–14, Kuala Lumpur, Malaysia, 2015.

[35] N. Aloysius and M. Geetha, "A review on deep convolutional neural networks," in *2017 International Conference on Communication and Signal Processing (ICCSP)*, pp. 588–592, Chennai, India, 2018.

[36] Z. Wang, "Deep convolutional neural networks for image classification: a comprehensive review," *Neural Comput*, vol. 2733, pp. 2709–2733, 2017.

[37] M. Toğaçar, B. Ergen, and Z. Cömert, "A deep feature learning model for pneumonia detection applying a combination of mRMR feature selection and machine learning models," *Irbm*, vol. 1, pp. 1–11, 2020.

[38] H. Peng, F. Long, and C. Ding, "Feature selection based on mutual information: criteria of max-dependency, max-relevance, and min-redundancy," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 8, pp. 1–6, 2005.

[39] C. Ding and H. Peng, "Minimum redundancy feature selection from microarray gene expression data," *Proc. 2003 IEEE Bioinforma. Conf. CSB 2003*, vol. 3, no. 2, pp. 523–528, 2003.

[40] C. Perez, J. Tapia, P. Estévez, and C. Held, "Gender classification from face images using mutual information and feature fusion," *International Journal of Optomechatronics*, vol. 6, no. 1, pp. 92–119, 2012.

[41] I. Chingovska, A. Anjos, and S. Marcel, "On the effectiveness of local binary patterns in face anti-spoofing," in *Proc. Int. Conf. Biometrics Spec. Interes. Group, BIOSIG 2012*, pp. 1–7, 2012.

[42] Z. Yang, W. Chen, F. Wang, and B. Xu, "Unsupervised domain adaptation for face anti-spoofin," in *2018 24th International Conference on Pattern Recognition (ICPR)*, pp. 338–343, Beijing, China, 2018.

[43] D M W Powers and Ailab, "Evaluation: from precision, recall and F-factor to ROC, informedness, markedness & correlation," *J. Mach. Learn. Technol.*, vol. 2, pp. 37–63, 2007.

[44] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-ResNet and the impact of residual connections on learning," in *31st AAAI Conf. Artif. Intell. AAAI 2017*, pp. 4278–4284, 2017.

[45] J. Komulainen, A. Hadid, M. Pietikainen, A. Anjos, and S. Marcel, "Complementary countermeasures for detecting scenic face spoofing attacks," in *2013 International Conference on Biometrics (ICB)*, Madrid, Spain, 2013.

[46] Z. Boulkenafet, J. Komulainen, and A. Hadid, "Face anti-spoofing based on color texture analysis," in *2015 IEEE International Conference on Image Processing (ICIP)*, pp. 2636–2640, Quebec City, QC, Canada, 2015.

[47] M. M. Hasan, M. S. U. Yusuf, T. I. Rohan, and S. Roy, "Efficient two stage approach to detect face liveness : motion based and deep learning based," in *2019 4th International Conference on Electrical Information and Communication Technology (EICT)*, pp. 20–22, Khulna, Bangladesh, 2019.

[48] G. Wang, H. Han, S. Shan, and X. Chen, "Unsupervised adversarial domain adaptation for cross-domain face presentation attack detection," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. , 202156–69, 2021.

[49] H. Li, P. He, S. Wang, A. Rocha, X. Jiang, and A. C. Kot, "Learning generalized deep feature representation for face anti-spoofing," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 10, pp. 2639–2652, 2018.

*Research Article*

# BLNN: Multiscale Feature Fusion-Based Bilinear Fine-Grained Convolutional Neural Network for Image Classification of Wood Knot Defects

**Mingyu Gao,**[1] **Fei Wang** [ID],[2,3] **Peng Song** [ID],[4] **Junyan Liu,**[2,3] **and DaWei Qi**[1]

[1]*College of Science, Northeast Forestry University, Harbin 150040, China*
[2]*School of Mechatronics Engineering, Harbin Institute of Technology, Harbin 150001, China*
[3]*State Key Laboratory of Robotics and System, Harbin Institute of Technology, Harbin 150001, China*
[4]*School of Instrumentation Science and Engineering, Harbin Institute of Technology, Harbin 150001, China*

Correspondence should be addressed to Fei Wang; wangfeipublic@163.com

Wood defects are quickly identified from an optical image based on deep learning methodology, which effectively improves the wood utilization. The traditional neural network technique is unemployed for the wood defect detection of optical image used, which results from a long training time, low recognition accuracy, and nonautomatic extraction of defect image features. In this paper, a wood knot defect detection model (so-called BLNN) combined deep learning is reported. Two subnetworks composed of convolutional neural networks are trained by Pytorch. By using the feature extraction capabilities of the two subnetworks and combining the bilinear join operation, the fine-grained features of the image are obtained. The experimental results show that the accuracy has reached up 99.20%, and the training time is obviously reduced with the speed of defect detection about 0.0795 s/image. It indicates that BLNN has the ability to improve the accuracy of defect recognition and has a potential application in the detection of wood knot defects.

## 1. Introduction

Wood knot defect detection is an important link in evaluating wood quality, which ultimately affects the quality of wood products [1]. Rapid detection of knot defects on wood surface can effectively improve the qualified rate of wood products [2, 3]. Consequently, it is important to identify the defects of wood knots in a short time. Although manual recognition is accurate, it takes a lot of time to train the staff, and the recognition speed on the assembly line is very slow compared to machine recognition [4, 5]. With the development of artificial intelligence and computer vision technology, deep learning has potential significance in the application of wood knot defect classification [6–8].

In recent years, image recognition based on artificial neural network and image processing has been widely studied. In order to identify the target accurately, the first step is to extract image features. For example, a Hu invariant moment

feature extraction method combined with a BP (back propagation) neural network to classify wood knot defects was proposed by Qi and Mu [9]. The accuracy of this method for wood knot defect recognition is over 86%. In the same year, Khwaja et al. proposed a defect detection and classification method for wet-blue leather using artificial neural network (ANN). The features of several defects on leather were extracted by using grey level cooccurrence matrix (GLCM) and grey level run-length matrix (GLRLM). The acquired features are passed to the multilayer perceptron using the Levenberg-Marquardt (LM) algorithm. The accuracy of this model is 97.85% [10]. In 2021, Aditya et al. proposed a method based on statistical texture features in GLCM to classify leaf blight of four plants by selecting appropriate thresholds. The accuracy of this method can reach 74% under optimal conditions [11]. The above methods require manual feature extraction, and the recognition rate is not high. Consequently, a convolutional neural network (CNN) which can

(a)

(b)

(c)

(d)

FIGURE 1: Four types of wood knot defects: (a) dry knot, (b) edge knot, (c) leaf knot, and (d) sound knot.

automatically learn the target features is needed to replace the complex artificial defect feature extraction. In 2020, Zhang et al. proposed a CNN image recognition algorithm for supermarket shopping robots. This algorithm overcomes the problems of low accuracy and slow speed in image recognition. The experimental results show that the accuracy of the algorithm can reach more than 98%. It also verifies that the image recognition algorithm can be applied to supermarket shopping robots to meet the needs of competition [12]. In the same year, Liu et al. proposed an intangible cultural her-

itage image recognition model based on color feature extraction and CNN, with the recognition rate reaching 94.8% [13]. In 2021, a new method based on transfer learning and ResNet-34 convolutional neural network for recognizing wood knot defects was presented by Gao et al. The experimental results show that the classification accuracy of this method can reach 98.69% [14]. Although these methods are practical, their accuracy can still be improved, and they have less application in wood knot defect detection. In order to solve these problems, improve the accuracy and recognition

(a)



(b)



(c)

FIGURE 2: Continued.

(d)



(e)



(f)



(g)

FIGURE 2: Original images of wood knot defect and those created through data augmentation: (a) original image, (b) vertical mirroring, (c) rotation by 180˚, (d) horizontal mirroring, (e) adding Gaussian noise to image, (f) increasing the hue by 10, and (g) adding salt-and-pepper noise to the image.

speed of the model, and reduce the training time, a high-accuracy wood knot defect detection method based on convolutional neural network is required.

In this paper, a bilinear classification model based on feature fine-grained fusion strategy named BLNN was proposed to detect wood knot defects. This paper is arranged and structured as follows. Firstly, the dataset of wood knot defects is acquired and augmented. Then, the proposed BLNN model is introduced. Subsequently, the network is trained and tested

by using the dataset of wood knot defects. Finally, based on a benchmark dataset, the test results are compared and analyzed with other deep learning models.

## 2. Materials and Methodology

*2.1. Dataset Acquisition.* The dataset was downloaded from the website of the Computer Laboratory, Department of Electrical Engineering, University of Oulu [15–17], and

TABLE 1: Number of datasets.

| Wood knot defect | Before data augmentation | | | | After data augmentation | | | |
|---|---|---|---|---|---|---|---|---|
| | Training dataset | Validation dataset | Testing dataset | Original dataset | Training dataset | Validation dataset | Testing dataset | Total dataset |
| Dry knot | 41 | 14 | 14 | 69 | 291 | 96 | 96 | 483 |
| Edge knot | 39 | 13 | 13 | 65 | 273 | 91 | 91 | 455 |
| Leaf knot | 27 | 10 | 10 | 47 | 198 | 65 | 66 | 329 |
| Sound knot | 110 | 37 | 37 | 184 | 772 | 266 | 250 | 1288 |
| Total | 217 | 74 | 74 | 365 | 1534 | 518 | 503 | 2555 |



FIGURE 3: Structure of the proposed fusion network.

consists of 365 images with four types of spruce knot defects. These are dry knot, edge knot, leaf knot, and sound knot, respectively. Figure 1 shows the four types of wood knot defects in the dataset used in this paper.

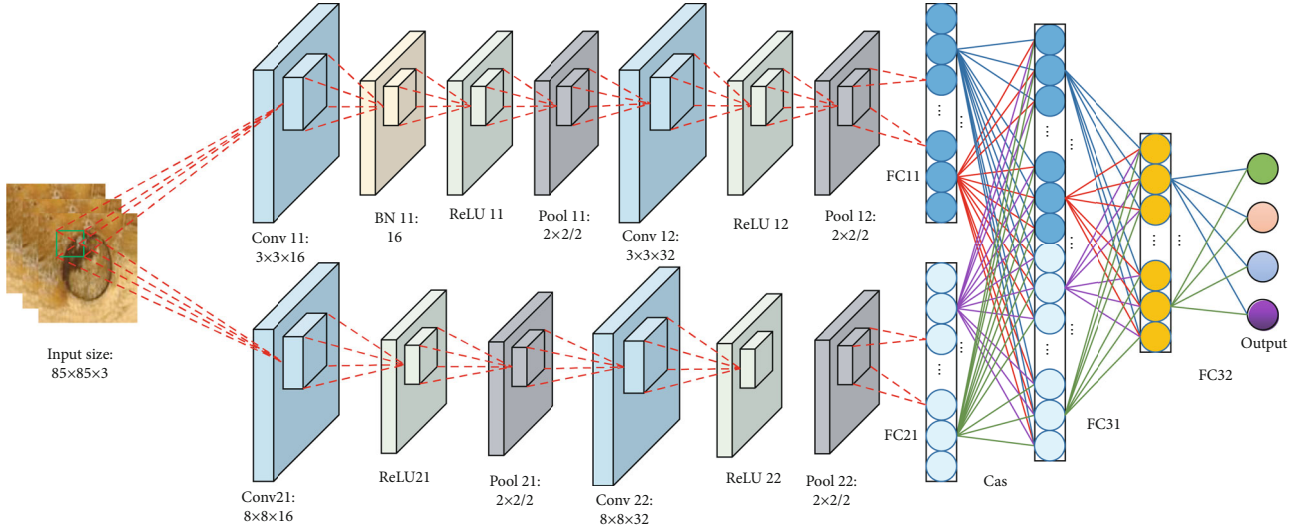*2.2. Image Preprocessing and Augmentation.* Deep learning networks have to be trained on massive datasets to achieve good performance [18]. Therefore, when the original dataset contains a limited number of images, data augmentation [19] is required to improve accuracy and prevent overfitting [20]. In this case, six methods are employed to augment the dataset, namely, vertical mirroring, rotation by 180°, horizontal mirroring, adding Gaussian noise, increasing the hue by 10, and adding salt-and-pepper noise. Consequently, the number of images was increased to seven times the original number. Due to more image augmentation, the learning ability of the network has increased. The data augmentation is shown in Figure 2. Table 1 lists the names and the number of images used for the experiments. Eventually, the dataset was randomly divided into a training set, a validation set, and a testing set in ratio of 3 : 1 : 1.

*2.3. Proposed Classification Model.* A CNN network called BLNN is proposed for fine-grained feature extraction [21–23] based on images, which consists of two different branching convolutional neural networks. Since the two

CNNs are different, they are used to extract features of different scales. These two features are confluence together to form a one-dimensional feature vector using the bilinear pooling operation [24, 25], and finally, the feature vector is classified using a classifier to obtain the recognized class. An overview of the proposed network architecture is shown in Figure 3. The parameters of BLNN are shown in Table 2.

*2.4. Multiscale Information Fusion Strategy.* The core of the BLNN lies in the fusion of two bilinear layer output vectors. According to this, a CNN-based fusion network structure is proposed to extract information about wood knot defects from different dimensions. BLNN can be expressed as follows:

$$B = (F_1, F_2, \mathrm{Fc}_{31}, \mathrm{Fc}_{32}), \qquad (1)$$

where $F_1$ and $F_2$ denote two feature extraction functions and $\mathrm{Fc}_{31}$ and $\mathrm{Fc}_{32}$ are the fully connected layers.

$$\begin{aligned} F_1 &= (C, B, R, P, \mathrm{Fc}_{11}), \\ F_2 &= (C, R, P, \mathrm{Fc}_{21}), \end{aligned} \qquad (2)$$

where $C$, $B$, $R$, $P$, $\mathrm{Fc}_{11}$, and $\mathrm{Fc}_{21}$ denote the convolutional layer [26], BatchNorm layer [27], ReLU activation function

TABLE 2: Parameters of BLNN layers.

| Layer | Type | Patch size | Kernel sum | Stride | Output size | Neuron sum |
|---|---|---|---|---|---|---|
| Input | Input | | | | $85 \times 85 \times 3$ | |
| Conv11 | Convolution | $3 \times 3$ | 16 | 1 | $83 \times 83 \times 16$ | $83 \times 83 \times 16$ |
| BN11 | BatchNorm | | | | $83 \times 83 \times 16$ | $83 \times 83 \times 16$ |
| ReLU11 | ReLU | | | | $83 \times 83 \times 16$ | $83 \times 83 \times 16$ |
| Pool11 | Avg-pooling | $2 \times 2$ | | 2 | $41 \times 41 \times 16$ | $41 \times 41 \times 16$ |
| Conv12 | Convolution | $3 \times 3$ | 32 | 1 | $39 \times 39 \times 32$ | $39 \times 39 \times 32$ |
| ReLU12 | ReLU | | | | $39 \times 39 \times 32$ | $39 \times 39 \times 32$ |
| Pool12 | Avg-pooling | $2 \times 2$ | | 2 | $19 \times 19 \times 32$ | $19 \times 19 \times 32$ |
| FC11 | Fully connected | $1 \times 1$ | 120 | | $1 \times 1 \times 120$ | $1 \times 1 \times 120$ |
| Conv21 | Convolution | $8 \times 8$ | 16 | 1 | $78 \times 78 \times 16$ | $78 \times 78 \times 16$ |
| ReLU21 | ReLU | | | | $78 \times 78 \times 16$ | $78 \times 78 \times 16$ |
| Pool21 | Avg-pooling | $2 \times 2$ | | 2 | $39 \times 39 \times 16$ | $39 \times 39 \times 16$ |
| Conv22 | Convolution | $8 \times 8$ | 32 | 1 | $32 \times 32 \times 32$ | $32 \times 32 \times 32$ |
| ReLU22 | ReLU | | | | $32 \times 32 \times 32$ | $32 \times 32 \times 32$ |
| Pool22 | Avg-pooling | $2 \times 2$ | | 2 | $16 \times 16 \times 32$ | $16 \times 16 \times 32$ |
| FC21 | Fully connected | $1 \times 1$ | 120 | | $1 \times 1 \times 120$ | $1 \times 1 \times 120$ |
| Cas | Cascade | | | | $1 \times 1 \times 240$ | $1 \times 1 \times 240$ |
| FC31 | Fully connected | $1 \times 1$ | 50 | | $1 \times 1 \times 50$ | $1 \times 1 \times 50$ |
| FC32 | Fully connected | $1 \times 1$ | 4 | | $1 \times 1 \times 4$ | $1 \times 1 \times 4$ |
| Output | Output | $1 \times 1$ | 4 | | $1 \times 1 \times 4$ | |



FIGURE 4: Local structure illustration of multiscale information fusion.

[28], pooling layer [29], and fully connected layers [30] of $F_1$ and $F_2$, respectively. First of all, the algorithm uses two branch networks named $F_1$ and $F_2$ to train the wood knot defect images, respectively. A smaller $3 \times 3$ convolutional kernel is used in $F_1$ to extract a rough feature; it can reduce the parameters. $F_2$ uses a larger $8 \times 8$ convolutional kernel

TABLE 3: Experimental environment.

| Hardware environment | | Software environment | |
| --- | --- | --- | --- |
| Memory | 128.00 GB | System | Windows 10 |
| CPU | Intel(R) Xeon(R) Bronze 3204 CPU @ 1.90 GHz (6 core) | Environment configuration | Pytorch-gpu 1.8.0 + Python 3.8.8 + cuda 11.1 + cudnn 8.0.5 |
| Graphics card | NVIDIA GeForce RTX 3090 (24 G) | | |

to extract features. The larger convolutional kernel can provide higher receptive field and extract more fine features. Therefore, $F_2$ is designed to capture the fine-grained characteristics [31] of wood knots. The fusion of two branches in the fully connected layer is shown in Figure 4.

After the first fully connected layer, vectors $x_1$ and $x_2$ with a dimension of $1 \times 120$ are obtained from the two branches, respectively (Figure 4). Then, $x_1$ and $x_2$ cascade to get $x_3$. Cascade fusion [32] is employed to superpose the two outputs, which can be expressed as follows:

$$f_{cas}(x_1, x_2), \tag{3}$$

where $x_1$ and $x_2$ are the outputs behind $Fc_{11}$ and $Fc_{21}$, respectively. Two vectors are cascaded and spliced along the vertical axis into one vector with a dimension of $1 \times 240$. Therefore, the vector $x_3$ contains all the eigenvectors computed by the two branches, which is computed from the image features of two different scales, and the features are represented more comprehensively. Next, a one-dimensional vector with a dimension of $1 \times 50$ is set after $x_3$, and finally, set the output of the fully connected layer to 4, indicating the category of classification.

*2.5. Loss Function and Optimizer.* The loss function is applied to evaluate the difference between the predicted and actual values of the model [33–35]. The smaller the difference, the smaller the cross-entropy. This study uses the cross-entropy loss function, which is expressed as follows:

$$L = -\sum_{i=1}^{n} p_i(x) \log [q_i(x)], \tag{4}$$

where $L$ represents the loss value of the sample and $p_i(x)$ and $q_i(x)$ represent the target output and the actual output, respectively. Cross-entropy overcomes the problem that weights and deviations are updated too slowly. When the error is large, the weight updates quickly, and when the error is small, the weight updates slowly.

The optimizer is used to update and compute the network parameters that affect the model training and output to approximate or reach the optimal value, thereupon then minimizing (or maximizing) the loss function [36]. In this case, the Adam optimizer is used. The Adam optimizer combines the advantages of AdaGrad [37] and RMSProp [38]. It takes the first-order moment estimation (i.e., the mean of the gradient) and second-order moment estimation (i.e., the uncentered variance of the gradient) of the gradient into

TABLE 4: Training hyperparameters.

| Related parameter | Value |
| --- | --- |
| Batch size | 128 |
| Learning rate | $1e-3$ |
| Epoch | 200 |
| Optimizer | Adam |
| Loss function | Cross-entropy |
| CUDA | Enable |
| CUDNN | Enable |

account and calculates the update step. Adam is simple to implement, is computationally efficient, and has low memory requirements, and the hyperparameters usually require no or little fine-tuning.

## 3. Experiment Results and Discussion

The experiment was performed on a Windows 10 64-bit PC equipped with an Intel(R) Xeon(R) Bronze 3204 CPU @ 1.90 GHz processor and 128 GB RAM. The deep learning programs were run on two NVIDIA GeForce RTX 3090 GPUs with 24 G RAM. The code is mainly implemented in Python, including data preprocessing and algorithm implementation. The deep learning framework is Pytorch. The experimental environment is shown in Table 3.

*3.1. Model Training.* In this study, the dataset is divided into a training set, a validation set, and a testing set, which contain 1534, 518, and 503 images, respectively. The hyperparameter setting for model training is shown in Table 4. The epoch, batch size, and learning rate are set to 200, 128, and $1e-3$ to make all models converge stably. The model training process is shown in Figure 5.

*3.1.1. The Training Results of the BLNN Model.* The accuracy and loss curves for the training and verification stages are shown in Figure 6, respectively.

Figure 6 shows that the model has trained 200 epochs; it can be seen that the training accuracy of the model remains stable after 50 epochs. Most of the fluctuations are between 0.95 and 1.00, and the loss decreases to around 0.2 to 0.35 with little fluctuation. After nearly 100 epochs, the loss of training phase decreased to about 0.2, but there are still fluctuations. The accuracy remained stable during the validation phase, most of which fluctuated between 0.95 and 1.00. Better classification results are obtained.
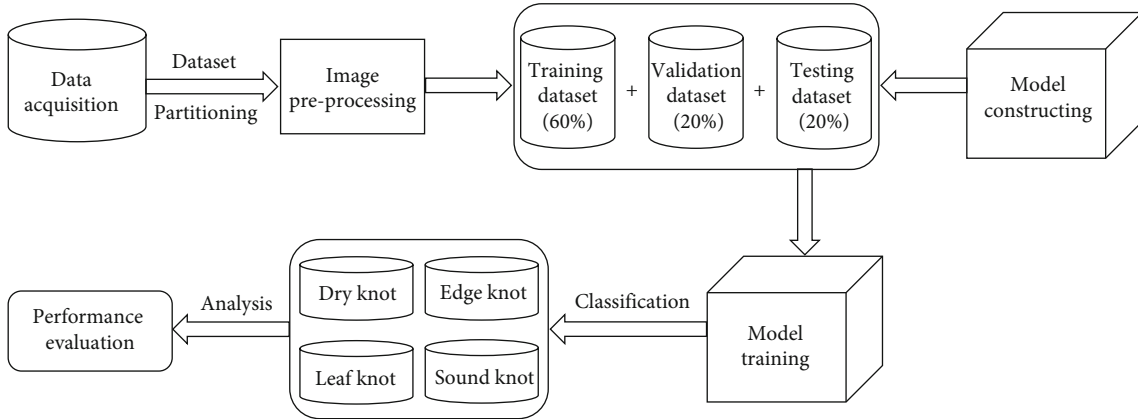
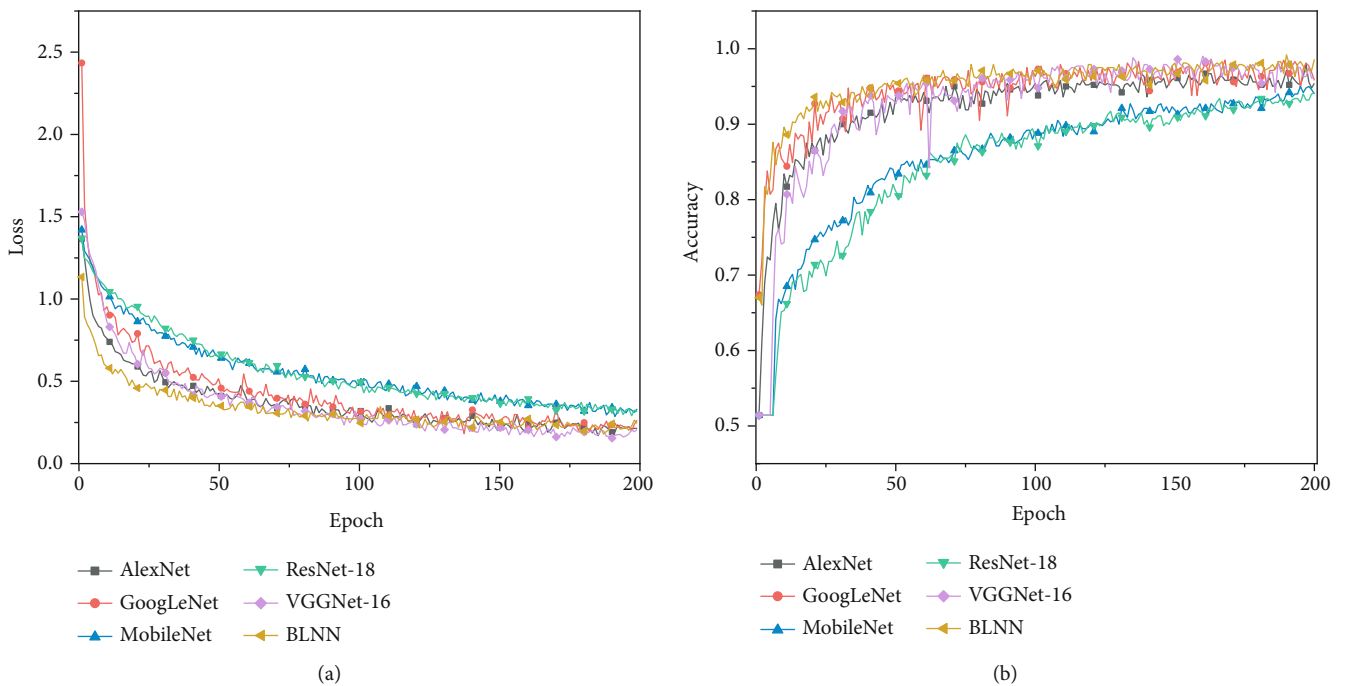FIGURE 5: Process flow diagram of wood knot defect detection.



FIGURE 6: The model was trained in the training dataset and validation dataset: (a) loss value; (b) accuracy value.

*3.1.2. Contrast Experiment.* The results of BLNN are compared with those of AlexNet, VGGNet-16, GoogLeNet, ResNet-18, and MobileNet-V2 to verify the effectiveness of the model. ResNet-18 achieves feature reuse by identity shortcut. Similar to ResNet, the fusion strategy of BLNN is to combine in-depth and shallow-depth features to obtain more detailed feature information. By comparing the performance of different network structures on the same wood knot defect dataset, the effectiveness and the superiority in identifying wood knot defects of BLNN are proved.

As shown in Figure 7, BLNN has a faster convergence rate than other models and finishes convergence at the 50th epoch. Consequently, a smaller epoch has the opportunity to be chosen to use in practice.

Five learning rates, 0.1, 0.01, 0.001, 0.0001, and 0.00001, were tested after establishing the BLNN model. The experimental results are shown in Table 5.

In Table 5, it is observed that when the learning rate is 0.1, the model does not converge effectively. The main reason is that an excessively large learning rate will cause the parameters of the model to oscillate beyond the valid range rapidly. When the learning rate has been reduced to 0.01, 0.001, and 0.0001, good results have been achieved, the error has been converged, and test accuracy has reached 94.43%, 99.20%, and 96.62%, respectively. When the learning rate continues to drop to 0.00001, the network convergence is very slow and the time to find the optimal value increases. At the same time, convergence may occur when entering the local extreme point, and no optimal value can be found. By continuously reducing the learning rate, it is found that the training results of different learning rates are different. Consequently, considering the accuracy and training time of the model, 0.001 is chosen as the initial learning rate to train the model.

(a)



(b)

FIGURE 7: Results in the training set of all the applied models.

TABLE 5: The comparison of results in different learning rates.

| Leaning rate | Number | Accuracy (%) |
|---|---|---|
| 0.1 | 250 | 49.70 |
| 0.01 | 475 | 94.43 |
| 0.001 | 499 | 99.20 |
| 0.0001 | 486 | 96.62 |
| 0.00001 | 436 | 86.68 |

The optimization algorithm is applied to find the optimal solution of the model. In this case, the Adam is employed and compared with SGD, AdaGrad, and Adax, as shown in Figure 8. The results show that the model with Adam has the fastest convergence speed and the highest accuracy. Table 6 shows the prediction results of the four optimization algorithms under the same condition. The results show that the accuracy of SGD, AdaGrad, Adamax, and Adam is 79.32%, 94.04%, 98.01%, and 99.20%, respectively. Consequently, considering the accuracy and training time of the model, Adam is chosen as the optimizer of the model.

*3.2. Evaluation Metrics.* To evaluate the performance of the BLNN, the precision ($P$), recall ($R$), $F1$ score ($F1$), and false alarm rate (FAR) were applied for the evaluation shown as follows:

$$P = \frac{TP}{TP + FP}, \tag{5a}$$

$$R = \frac{TP}{TP + FN}, \tag{5b}$$

$$FAR = \frac{FP}{FP + TN}, \tag{5c}$$

(a)



(b)

Figure 8: Results in the training and validation sets of all the applied optimizers.

Table 6: The comparison of results in different optimizers.

| Optimizer | Number | Accuracy (%) |
|-----------|--------|--------------|
| AdaGrad   | 473    | 94.04        |
| Adamax    | 493    | 98.01        |
| SGD       | 399    | 79.32        |
| Adam      | 499    | 99.20        |

$$F1 = 2 \frac{P \cdot R}{P + R}, \tag{5d}$$

where TP, FP, TN, and FN represent the true positive, false positive, true negative, and false negative.

### 3.3. Model Evaluation.

The performance of BLNN is evaluated in the task of wood knot defect classification. 503 wood knot defect images were used as testing dataset. The trained BLNN was compared with AlexNet, GoogLeNet, MobileNet, ResNet-18, and VGGNet-16, and the network was evaluated according to confusion matrix, precision, recall, $F1$ score, FAR, accuracy, training time, and detection time.

As shown in the confusion matrix in Figure 9, the accuracy of each category is described by comparing the actual category with the predicted category. The numerical

Figure 9: Confusion matrix in the testing dataset of all the applied models: (a) AlexNet, (b) GoogLeNet, (c) MobileNet, (d) ResNet-18, (e) VGGNet-16, and (f) BLNN.

distribution of confusion matrix shows that AlexNet and BLNN have better classification results. BLNN can recognize edge knot and 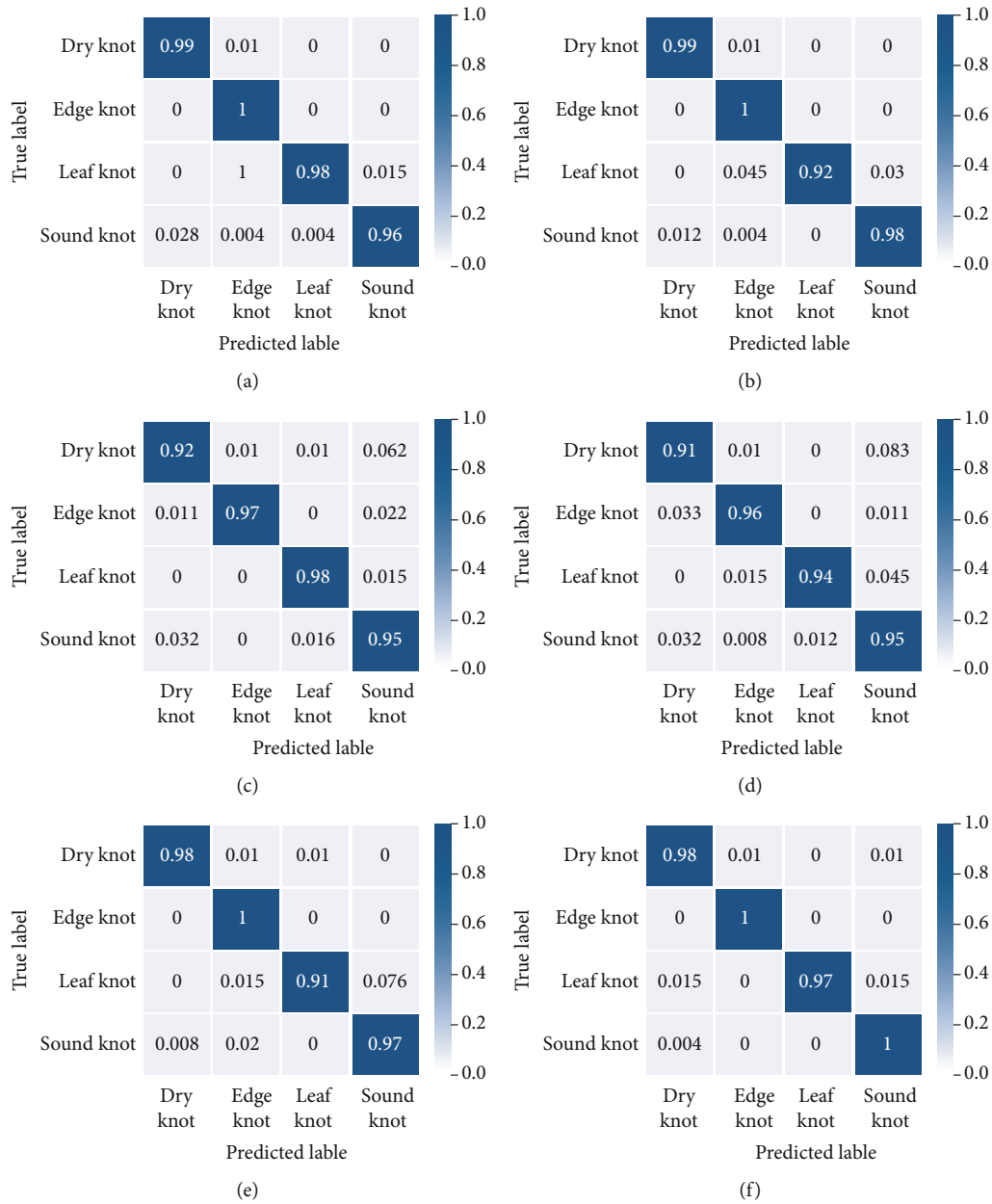sound knot up to 100%, and dry knot and leaf knot are slightly lower than AlexNet, which is the direction to improve in the future. However, as shown in Figure 10, BLNN has the highest overall recognition rate of knot defects, reaching 99.20%. Table 7 shows the training time and the detection time of all models for each wood image. It can be seen that BLNN has the shortest training time and the fastest detection speed in all models due to its fewer parameters and higher feature extraction ability.

Precision, recall, $F1$, and FAR of the four categories of wood knot defect images in the testing set are shown in

Figure 11. It can be seen that BLNN is superior to Mobile-Net-V2, ResNet-18, and VGGNet-16 in the classification of four wood knot defects. Compared with AlexNet and GoogLeNet, some of the BLNN metrics are slightly worse, but the gap is not big, which requires further improvement in the future. As shown in Figure 10 and Table 7, although BLNN is not always optimal in these models, BLNN has the highest accuracy and the fastest training time and detection speed, and it is easy to be built and embedded into other models because of its small parameters and computation, which makes it possible to identify wood knot defects. Compared with other models, BLNN has obvious advantages in accuracy and calculation, so it has more practical application value. An unexpected phenomenon is that MobileNet,
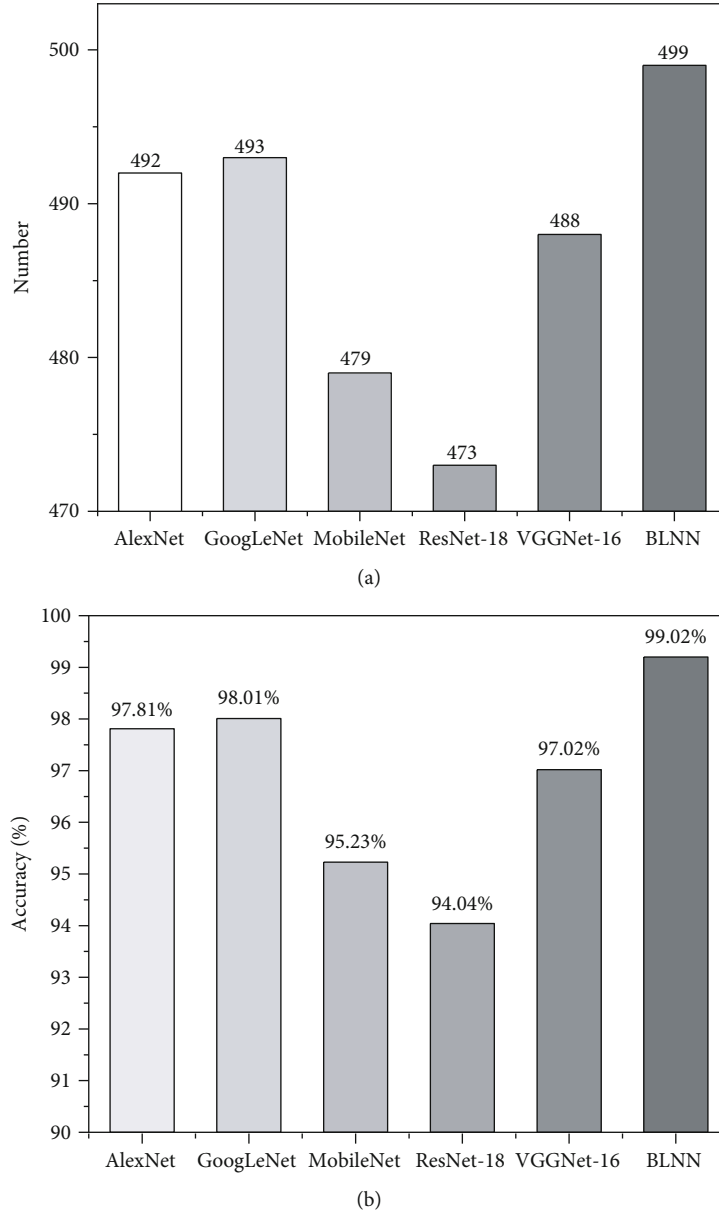
(a)



(b)

FIGURE 10: Prediction results of all the applied models in the testing dataset.

TABLE 7: Training time and detection time of all the applied methods.

| Method | Training time (min) | Detection time (s/image) |
| --- | --- | --- |
| AlexNet | 37.32 | 0.2744 |
| GoogLeNet | 44.27 | 0.3519 |
| MobileNet-V2 | 12.97 | 0.2425 |
| ResNet-18 | 15.95 | 0.4573 |
| VGGNet-16 | 36.88 | 1.9583 |
| BLNN | 11.22 | 0.0795 |

ResNet-18, and VGGNet-16 do not achieve the desired performance, especially ResNet which has the lowest recognition rate. Therefore, the network structure has a great impact on the training results.

As shown in Figure 3, BLNN consists of two single-branch networks. To verify the improvement of model performance by using two-branch networks, the upper and lower branches of BLNN are compared with BLNN, respectively. The results are shown in Figures 12 and 13.

From Figures 12 and 13, it can be seen that BLNN has the fastest convergence speed and highest accuracy in the three networks. In addition, the convergence speed of the upper branch network in the training set is faster than that of the lower branch network, and the performance of the lower branch network in the verification set is better than that of the upper branch network. As shown in Figure 13, BLNN has the best performance, the lower network has the second performance, and the upper network has the worst performance, because the upper network uses $3 \times 3$ convolutional kernel, the lower network uses $8 \times 8$ convolutional kernel,
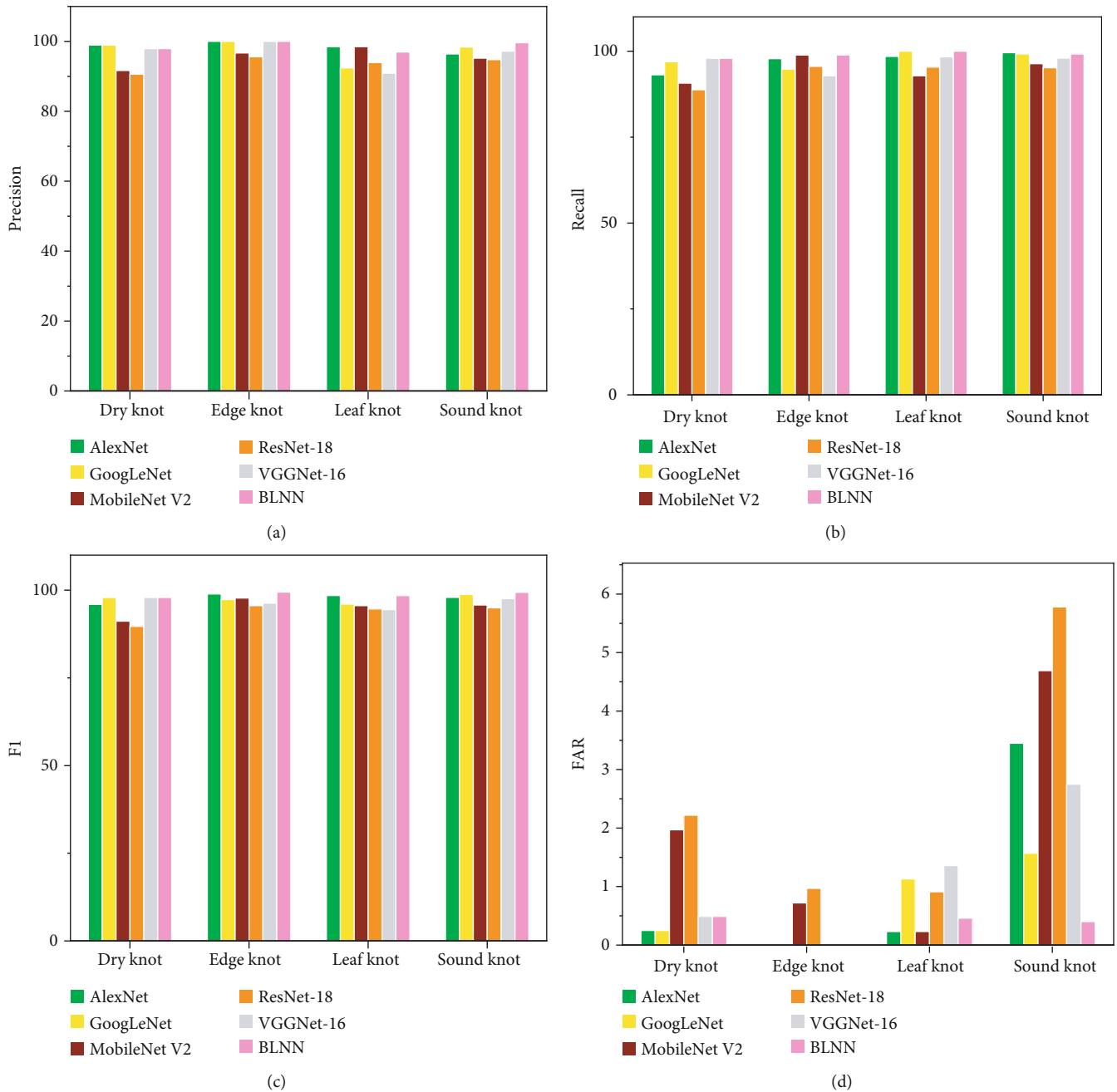
(a)

(b)

(c)

(d)

FIGURE 11: The evaluation index values of network.

and the lower network has a larger receptive field. Therefore, the bilinear structure of BLNN has better performance than that of single-branch networks.

As shown in Figure 3, BLNN has two single-branch networks. The upper and lower branch networks use different sizes of convolutional kernel; the upper branch network convolutional kernel is $3 \times 3$, and the lower branch network convolutional kernel is $8 \times 8$. To verify the effect of different convolutional kernel sizes on the model performance, we separately use BLNN (the upper branch network is $3 \times 3$, the lower branch network is $8 \times 8$) compared with two networks with $3 \times 3$ and $8 \times 8$; the results are shown in Figures 14 and 15.

From Figures 14 and 15, it can be seen that BLNN has the fastest convergence speed and highest accuracy in these three networks. In addition, the network with convolutional kernel size $3 \times 3$ in the training set converges faster than $8 \times 8$, and the network with convolutional kernel size $8 \times 8$ in the verification set performs better than $3 \times 3$. As shown in Figure 15, BLNN performs best, the network with convolutional kernel size $8 \times 8$ performs second, and the network with convolutional kernel size $3 \times 3$ performs worst. This is because networks with $8 \times 8$ convolutional kernel have a larger receptive field, but BLNN uses dual-branch networks with different sizes of convolutional kernel, smaller convolutional kernel ($3 \times 3$) for upper branch networks to extract local
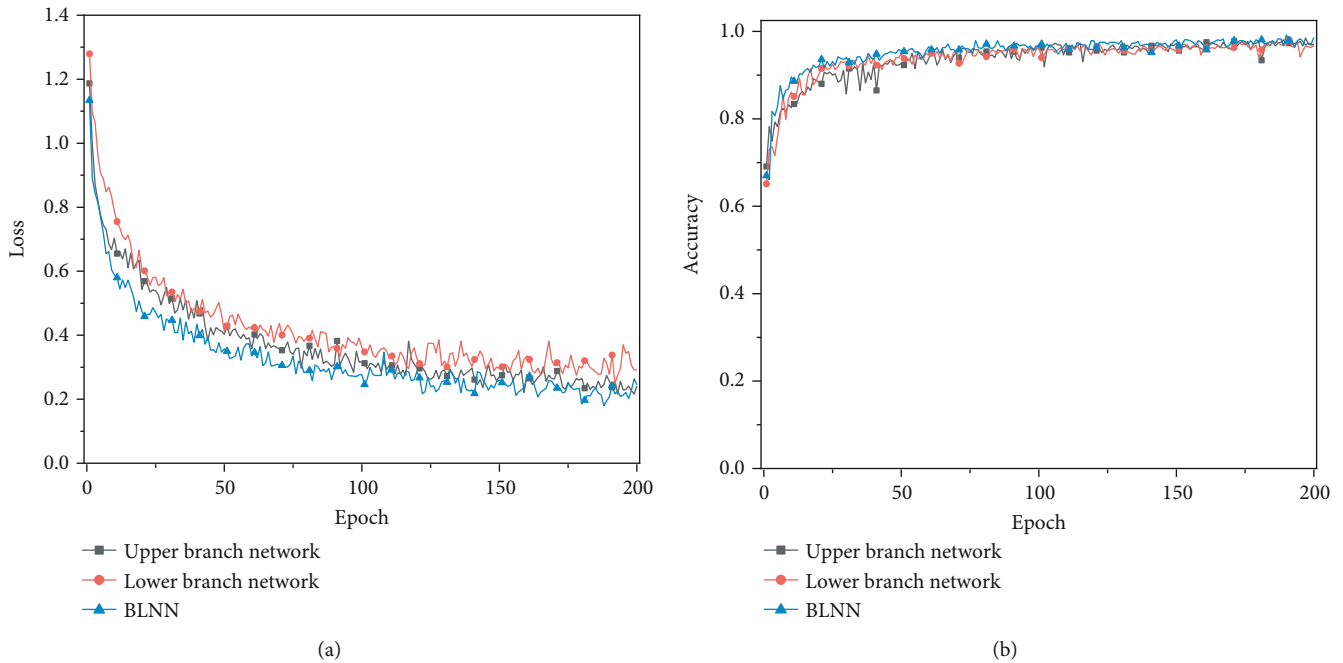
(a)

(b)

FIGURE 12: Results in the training and validation sets of BLNN and its component.



(a)

(b)

FIGURE 13: Prediction results of BLNN and its component in the testing dataset.

details and larger convolutional kernel (8 × 8) for lower branch networks to extract more comprehensive global information, and then, these two kinds of feature information are fused. More comprehensive information can be acquired, so the performance of BLNN is better than that of the other two networks with different convolutional kernels.

*3.4. Model Generalization.* In order to evaluate the generalization ability of BLNN, we tested the classification ability of BLNN on some boards. Green means correct recognition was used to mark in green and the wrong recognition was marked in grey in this case. Details of the identification such

as the name and probability of wood knot defects are displayed next to each label. Figure 16 shows four wood knot defects and the corresponding identification results.

It can be seen that most of the wood knot defects in the image are correctly identified. Some of the wood knot defects are similar in shape to other defects, and some of the wood defects are not trained, which makes the model appear to identify errors. In most cases, our method (BLNN) still has high accuracy. This indicates that BLNN has certain application value in practice.

As shown in Figure 16, since we only focus on the four defects of dry knot, edge knot, leaf knot, and sound knot

(a)

(b)

FIGURE 14: Results in the training and validation sets of BLNN and other convolutional kernel sizes.



(a)

(b)

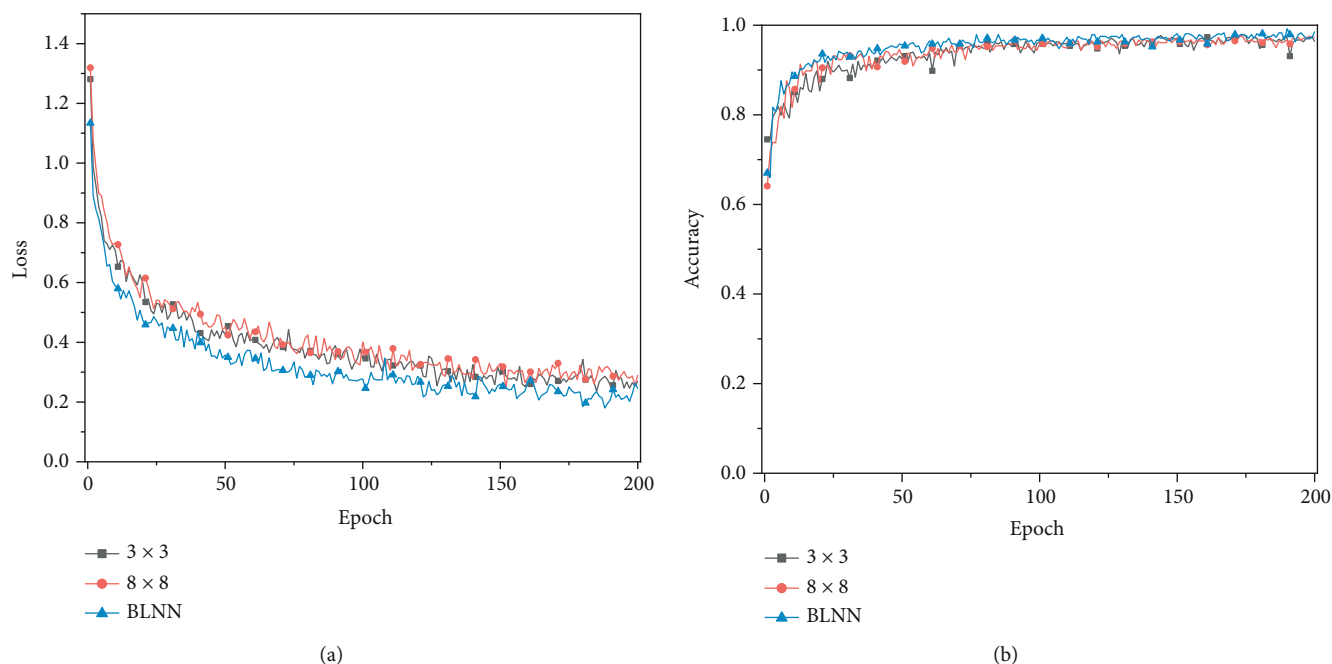FIGURE 15: Prediction results of BLNN and other convolutional kernel sizes in the testing dataset.

when training the network, it can be seen that there are some defects that have not been identified. This is one of our future research directions to increase the types of defect classification.

### 3.5. Discussion.
The effectiveness of BLNN can be discussed in two aspects.

### 3.5.1. Feasibility of Bilinear Network Structure.
Compared with single-branch network, BLNN has obvious advantages in accuracy and convergence speed, which proves that the classification ability of the network can be improved by extracting and fusing features from the bilinear network. This network extracts features from two parallel single-branch networks, which can make the extracted features more com-

prehensive. This is the key to improve the classification performance. Although classical network structures such as ResNet are generally single-branch networks, their features are relatively single. Bilinear network can extract more information than a single network.

### 3.5.2. Rationality of Using Different Convolutional Kernel Sizes.
Compared with other classical networks, BLNN has obvious advantages in accuracy and computation, which proves that the classification ability of networks can be improved by fusing local features (convolutional kernel size $3 \times 3$) and global features (convolutional kernel size $8 \times 8$) through a bilinear fusion structure. The network uses convolutional kernel with different sizes to extract multiscale

FIGURE 16: The generalization test of BLNN.

features from the same image, and this fine-grained information is the key to classification.

For the proposed BLNN network, the local and global features extracted by the convolutional layer are fused in the fully connected layer. In other words, it fuses all the features of different scales together through a fusion operation. Therefore, BLNN expands the number of features without generating many complex feature maps. In the fully connected layer, we improve the robustness and classification accuracy of the network by setting an appropriate number of neurons.

BLNN performs well in the classification of wood knot defects. However, performing network fusion operations in the fully connected layer may not be optimal for other tasks. This requires more research in the future.

## 4. Conclusion

In conclusion, a bilinear classification model based on feature fine-grained fusion strategy named BLNN was proposed in this case. The convolutional kernel size of the upper branch network of BLNN was set to $3 \times 3$, and the convolutional kernel size of the lower branch network was set to $8 \times 8$. Two different sizes of convolutional kernels were used to extract features at different scales, and feature fusion was used to classify the wood knot defects. 2052 images of wood knot defects were used for training after 200 training epochs. The experimental results show that the accuracy of BLNN reaches 99.20% during the testing phase. In addition, when wood knot defects are detected by this method, a large number of image preprocessing and manual feature extraction are not demanded, which greatly improves the recognition efficiency. The speed of defect detection is only *0.0795 s/image*, and the training time is reduced. This means that BLNN has potential application value in wood nondestructive testing and wood knot defect detection and provides a feasible solution for future wood knot defect identification. In addition, the experimental results also show that multiscale information fusion is effective to improve model performance through network fusion.

## Data Availability

The datasets, codes, and weight files 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

## References

[1] Y. Fang, L. Lin, H. Feng, Z. Lu, and G. Emms, "Review of the use of air-coupled ultrasonic technologies for nondestructive testing of wood and wood products," *Computers and Electronics in Agriculture*, vol. 137, pp. 79–87, 2017.

[2] W. Zhou, M. Fei, H. Zhou, and K. Li, "A sparse representation based fast detection method for surface defect detection of bottle caps," *Neurocomputing*, vol. 123, pp. 406–414, 2014.

[3] C. Todoroki, E. Lowell, and D. Dykstra, "Automated knot detection with visual post-processing of Douglas-fir veneer images," *Computers and Electronics in Agriculture*, vol. 70, no. 1, pp. 163–171, 2010.

[4] D. Yadav and A. Yadav, "A novel convolutional neural network based model for recognition and classification of apple leaf diseases," *Traitement du Signal*, vol. 37, no. 6, 2020.

[5] X. Zhu, M. Zhu, and H. Ren, "Method of plant leaf recognition based on improved deep convolutional neural network," *Cognitive Systems Research*, vol. 52, pp. 223–233, 2018.

[6] T. He, Y. Liu, Y. Yu, Q. Zhao, and Z. Hu, "Application of deep convolutional neural network on feature extraction and detection of wood defects," *Measurement*, vol. 152, article 107357, 2020.

[7] J. Shi, Z. Li, T. Zhu, D. Wang, and C. Ni, "Defect detection of industry wood veneer based on NAS and multi-channel mask R-CNN," *Sensors*, vol. 20, no. 16, p. 4398, 2020.

[8] Y. Huang, J. Jing, and Z. Wang, "Fabric defect segmentation method based on deep learning," *IEEE Transactions on Instrumentation and Measurement*, vol. 70, pp. 1–5, 2021.

[9] D. Qi and H. Mu, "Detection of wood defects types based on Hu invariant moments and BP neural network," *Journal of Southeast University*, vol. 43, pp. 63–66, 2013.

[10] K. Mohammed, S. K. S, and P. G, "Defective texture classification using optimized neural network structure," *Pattern Recognition Letters*, vol. 135, pp. 228–236, 2020.

[11] A. Sinha and R. Singh Shekhawat, "A novel image classification technique for spot and blight diseases in plant leaves," *The Imaging Science Journal*, vol. 5, pp. 1–5, 2021.

[12] X. Zhang, H. Lu, Q. Xu et al., "Image recognition of supermarket shopping robot based on CNN," in *2020 IEEE International Conference on Artificial Intelligence and Computer Applications (ICAICA)*, pp. 1363–1368, Dalian, China, 2020.

[13] E. Liu, "Research on image recognition of intangible cultural heritage based on CNN and wireless network," *EURASIP Journal on Wireless Communications and Networking*, vol. 2020, no. 1, 2020.

[14] M. Gao, D. Qi, H. Mu, and J. Chen, "A transfer residual neural network based on ResNet-34 for detection of wood knot defects," *Forests*, vol. 12, no. 2, p. 212, 2021.

[15] H. Kauppinen and O. Silven, Eds., "A color vision approach for grading lumber," in *Theory & Applications of Image Processing II-Selected Papers from the 9th Scandinavian Conference on Image Analysis*, pp. 367–379, Singapore, 1995.

[16] O. Silven and H. Kauppinen, "Recent developments in wood inspection," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 10, no. 1, pp. 83–95, 1996.

[17] H. Kauppinen and O. Silven, "The effect of illumination variations on color-based wood defect classification," in *Proceedings of the 13th International Conference on Pattern Recognition (13th ICPR)*, pp. 828–832, Vienna, Austria, August 1996.

[18] G. Folego, M. Weiler, R. Casseb, R. Pires, and A. Rocha, "Alzheimer's disease detection through whole-brain 3D-CNN MRI," *Frontiers in Bioengineering and Biotechnology*, vol. 8, 2020.

[19] A. El Bilali, A. Taleb, M. Bahlaoui, and Y. Brouziyne, "An integrated approach based on Gaussian noises-based data augmentation method and AdaBoost model to predict faecal coliforms in rivers with small dataset," *Journal of Hydrology*, vol. 29, article 126510, 2021.

[20] M. Monshi, J. Poon, V. Chung, and F. Monshi, "CovidXrayNet: optimizing data augmentation and CNN hyperparameters for improved COVID-19 detection from CXR," *Computers in Biology and Medicine*, vol. 133, article 104375, 2021.

[21] C. Liu, H. Ding, and X. Jiang, "Towards enhancing fine-grained details for image matting," in *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 385–393, Waikoloa, HI, USA, 2021.

[22] X. Chen and J. Lai, "Salient points driven pedestrian group retrieval with fine-grained representation," *Neurocomputing*, vol. 423, pp. 255–263, 2021.

[23] X. Chen and J. Lai, "Salient points driven pedestrian group retrieval with fine-grained representation," *Neurocomputing*, vol. 423, pp. 255–263, 2021.

[24] W. Wu and J. Yu, "An improved bilinear pooling method for image-based action recognition," in *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 8578–8583, Milan, Italy, 2021.

[25] X. Chen, X. Zheng, and X. Lu, "Bidirectional interaction network for person re-identification," *IEEE Transactions on Image Processing*, vol. 30, pp. 1935–1948, 2021.

[26] T. Pradhan, P. Kumar, and S. Pal, "CLAVER: an integrated framework of convolutional layer, bidirectional LSTM with attention mechanism based scholarly venue recommendation," *Information Sciences*, vol. 559, pp. 212–235, 2021.

[27] S. Gao, Q. Han, D. Li, P. Peng, M. Cheng, and P. Peng, "Representative batch normalization with feature calibration," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8669–8679, 2021.

[28] F. Laakmann and P. Petersen, "Efficient approximation of solutions of parametric linear transport equations by ReLU DNNs," *Advances in Computational Mathematics*, vol. 47, no. 1, pp. 1–32, 2021.

[29] C. Ren, J. Dulay, G. Rolwes, D. Pauli, N. Shakoor, and A. Stylianou, "Multi-resolution outlier pooling for sorghum classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2931–2939, 2021.

[30] P. Staszewski, M. Jaworski, J. Cao, and L. Rutkowski, "A new approach to descriptors generation for image retrieval by analyzing activations of deep neural network layers," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–8, 2021.

[31] H. Šimonová, B. Kucharczyková, V. Bílek, L. Malíková, P. Miarka, and M. Lipowczan, "Mechanical fracture and fatigue characteristics of fine-grained composite based on sodium hydroxide-activated slag cured under high relative humidity," *Applied Sciences*, vol. 11, no. 1, p. 259, 2021.

[32] X. Zhu, S. Ye, L. Zhao, and Z. Dai, "Hybrid attention cascade network for facial expression recognition," *Sensors*, vol. 21, no. 6, p. 2003, 2021.

[33] R. Ferdous, M. Arifeen, T. Eiko, and S. Al Mamun, "Performance analysis of different loss function in face detection architectures," in *Proceedings of International Conference on Trends in Computational and Cognitive Engineering*, pp. 659–669, Singapore, 2021.

[34] M. Shorfuzzaman and M. Hossain, "MetaCOVID: a Siamese neural network framework with contrastive loss for _n_ -shot diagnosis of COVID-19 patients," *Pattern Recognition*, vol. 113, article 107700, 2021.

[35] P. Negi, R. Marcus, A. Kipf et al., "Flow-Loss: learning cardinality estimates that matter," 2021, http://arxiv.org/abs/2101.04964.

[36] Z. Zhang, "Improved Adam optimizer for deep neural networks," in *2018 IEEE/ACM 26th International Symposium on Quality of Service (IWQoS)*, pp. 1-2, Banff, AB, Canada, 2018.

[37] R. Ward, X. Wu, and L. Bottou, "AdaGrad stepsizes: sharp convergence over nonconvex landscapes," in *International Conference on Machine Learning*, pp. 6677–6686, 2019.

[38] F. Zou, L. Shen, Z. Jie, W. Zhang, and W. Liu, "A sufficient condition for convergences of Adam and RMSProp," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11127–11135, Long Beach, CA, USA, 2019.

*Research Article*

# A Novel Deep Convolutional Neural Network Based on ResNet-18 and Transfer Learning for Detection of Wood Knot Defects

**Mingyu Gao,[1] Peng Song [ID],[2] Fei Wang,[3,4] Junyan Liu,[3,4] Andreas Mandelis,[5,6] and DaWei Qi[1]**

[1]College of Science, Northeast Forestry University, Harbin 150040, China
[2]School of Instrumention Science and Engineering, Harbin Institute of Technology, Harbin 150001, China
[3]School of Mechatronics Engineering, Harbin Institute of Technology, Harbin 150001, China
[4]State Key Laboratory of Robotics and System, Harbin Institute of Technology, Harbin 150001, China
[5]Center for Advanced Diffusion-Wave and Photoacoustic Technologies, University of Toronto, Toronto, Canada M5S 3G8
[6]Institute for Advanced Non-Destructive and Non-Invasive Diagnostic Technologies (IANDIT), University of Toronto, Toronto, Canada M5S 3G8

Correspondence should be addressed to Peng Song; songpeng@hit.edu.cn

Wood defects are quickly identified from an optical image based on deep learning methodology, which effectively improves wood utilization. Traditional neural network techniques have not yet been employed for wood defect detection due to long training time, low recognition accuracy, and nonautomatical extraction of defect image features. In this work, a model (so-called ReSENet-18) for wood knot defect detection that combined deep learning and transfer learning is proposed. The "squeeze-and-excitation" (SE) module is firstly embedded into the "residual basic block" structure for a "SE-Basic-Block" module construction. This model has the advantages of the features that are extracted in the channel dimension, and it is fused in multiscale with original features. Instantaneously, the fully connected layer is replaced with a global average pooling; consequently, the model parameters could be reduced effectively. The experimental results show that the accuracy has reached 99.02%, meanwhile the training time is also reduced. It shows that the proposed deep convolutional neural network based on ReSENet-18 combined with transfer learning can improve the accuracy of defect recognition and has a potential application in the detection of wood knot defects.

## 1. Introduction

Wood knot defect detection is an important part in the production of wood products and finally affects the quality of wood products. Rapid detection of wood knot defects on the surface of the wood can effectively improve the qualification rate of wood products. Consequently, it is important to quickly identify the wood knot defects in a short time [1–4]. Although the traditional manual recognition is widely used and accurate, it is still a subjective [5] and inefficient method to identify wood knot defects [6]. With the rapid development of digital image processing and computer vision, artificial intelligence technology can improve the recognition speed and accuracy at a certain extent [7–9]. Among them, deep learning is the most potential method in the field of artificial intelligence.

In recent years, wood knot defect recognition based on the artificial neural network and image analysis processing has been widely studied [10–15]. Because of its simple basic structure, the neural network can fit various data in theory. Because of this, large-scale neural network combination is needed. Due to the limitation of hardware, the current tools are not enough to run this complex network, resulting in its slow evolution. At present, with the development of robots and so on, the demand for computer vision technology based on CNN (convolutional neural network) is gradually increasing. Therefore, the neural network still has a great application value in the future. In the field of wood defect detection, the accurate recognition of wood needs to collect the defect image by camera or X-ray and then recognize it by image processing and artificial intelligence. In order to accurately identify the wood knot defects, image features must be
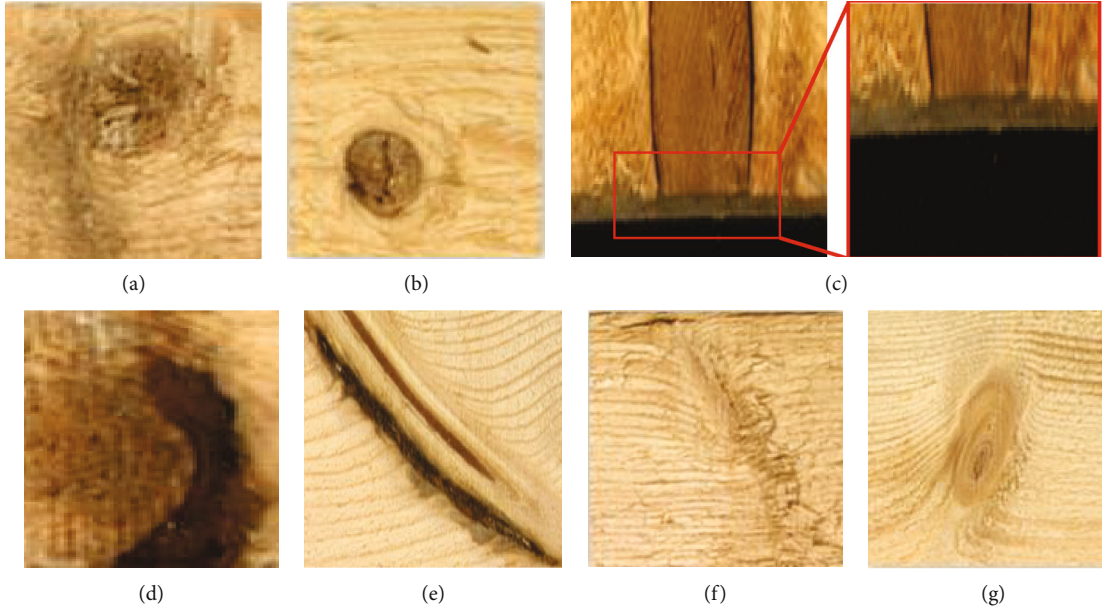
(a)  (b)  (c)

(d)  (e)  (f)  (g)

FIGURE 1: Seven types of wood knot defects: (a) decayed knot, (b) dry knot, (c) edge knot, (d) encased knot, (e) horn knot, (f) leaf knot, and (g) sound knot.

TABLE 1: Number of datasets.

| Wood knot | Before data augmentation | | | After data augmentation | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Training dataset | Validation dataset | Testing dataset | Original dataset | Training dataset | Validation dataset | Testing dataset | Total dataset |
| Decayed knot | 10 | 3 | 3 | 16 | 68 | 25 | 19 | 112 |
| Dry knot | 41 | 14 | 14 | 69 | 291 | 96 | 96 | 483 |
| Edge knot | 39 | 13 | 13 | 65 | 273 | 91 | 91 | 455 |
| Encased knot | 20 | 6 | 6 | 32 | 136 | 44 | 44 | 224 |
| Horn knot | 21 | 7 | 7 | 35 | 147 | 49 | 49 | 245 |
| Leaf knot | 27 | 10 | 10 | 47 | 198 | 65 | 66 | 329 |
| Sound knot | 110 | 37 | 37 | 184 | 772 | 266 | 250 | 1288 |
| Total | 268 | 90 | 90 | 448 | 1885 | 636 | 615 | 3136 |

extracted first. For example, Lin et al. in 2015 proposed a method to classify wood knot defects by combining the aspect ratio, grayscale, and variance feature extraction method of the back propagation (BP) network [10]. The accuracy of this method can reach 86.67%. In the same year, Mu et al. proposed a wood defect classification method by extracting the perimeter, area, aspect ratio, and mean grayscale value of the defect, combined with the radial basis function (RBF) neural network with the accuracy over 85% [11]. In 2019, Ji et al. proposed a wood defect classification method based on Hu moment invariant feature extraction and a combination of wavelet moment with BP network [12]. The crack identification accuracy of this method can reach 98%. However, due to the shape of flying knot scar and hole being similar, it is easy to induce a misclassification in some cases. Due to the quite unique shape of each wood knot defect, it is difficult and complex to identify the defect

by extracting the image features manually [13]. Therefore, a convolutional neural network (CNN) which can automatically learn the wood knot features is needed to replace the complex manual defect feature extraction. In 2019, Liu et al. proposed a CNN based on split-shuffle-residual (SSR) for real-time classification of rubber boards [14]. Comprehensive experiments show that the algorithm is superior than other classification methods and the latest deep learning classification network at that time has an accuracy of 94.86%, but there is still room for improvement. In 2021, a new method based on transfer learning and ResNet-34 convolutional neural network for recognizing wood knot defects was presented by Gao et al. [15]. The experimental results show that the classification accuracy of this method can reach 98.69%. Although both methods are practical, with the increase of network depth, the model parameters become more complex and the amount of calculation becomes

(a) Decayed knot

(b) Dry knot

(c) Edge knot

(d) Encased knot

(e) Horn knot

(f) Leaf knot

(g) Sound knot

FIGURE 2: Seven common wood knots and data augmentation of the dataset. Original images and those created through data augmentation: ① original image, ② vertical mirror, ③ rotated by 180, ④ horizontal mirror, ⑤ added Gaussian noise to image, ⑥ increased the hue by 10, and ⑦ added salt-and-pepper noise to image.



FIGURE 3: "Residual Basic-Block" structure of ResNet-18 acting as a building block for the network.

TABLE 2: The structure of ResNet-18.

| Layer name | Output size | 18-layer |
|---|---|---|
| Conv1 | $112 \times 112$ | $7 \times 7$, 64, stride 2 |
| | | $3 \times 3$ max pool, stride 2 |
| Conv2_x | $56 \times 56$ | $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$ |
| Conv3_x | $28 \times 28$ | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$ |
| Conv4_x | $14 \times 14$ | $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$ |
| Conv5_x | $7 \times 7$ | $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$ |

larger. To solve these problems and improve the accuracy of the model, a high accuracy wood knot defect detection method based on the convolutional neural network is required.

FIGURE 4: A squeeze-and-excitation block.



FIGURE 5: A "SE-Basic-Block" building block.



FIGURE 6: Comparison of the FC layer (a) and the GAP layer (b).

In this paper, a model based on the attention mechanism and deep transfer residual convolutional neural network structure named ReSENet-18 is proposed to detect wood knot defects. This paper is arranged and structured as follows. Firstly, the dataset of wood knot defects is acquired and preprocessed. Then, the proposed ReSENet-18 model is introduced. A squeeze-and-excitation-basic block (SE-Basic-Block) is added, and the fully connected layer is replaced by a global average pooling layer to adjust the network structure. At the same time, combined with the ideology of transfer learning, the ReSENet-18 network is pretrained on ImageNet. Subsequently, the network is trained and tested by using the dataset of wood knot defects. Finally, based on a benchmark dataset, the test results are compared and analyzed with other deep learning models.
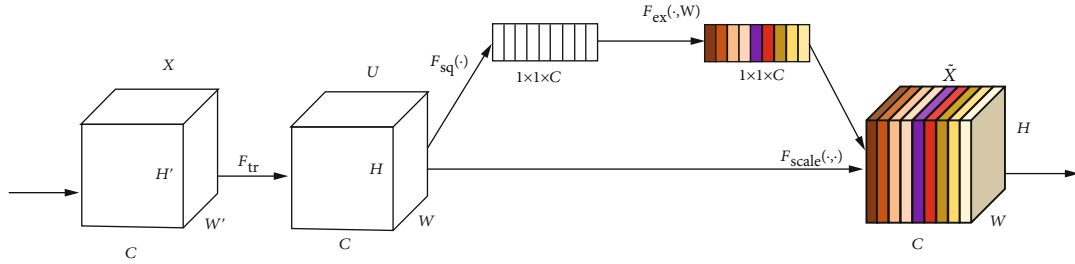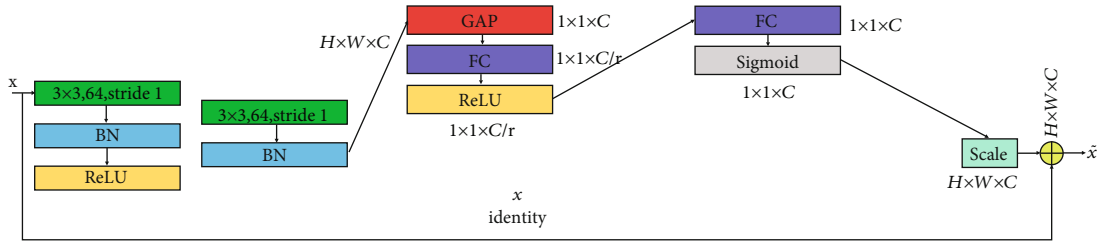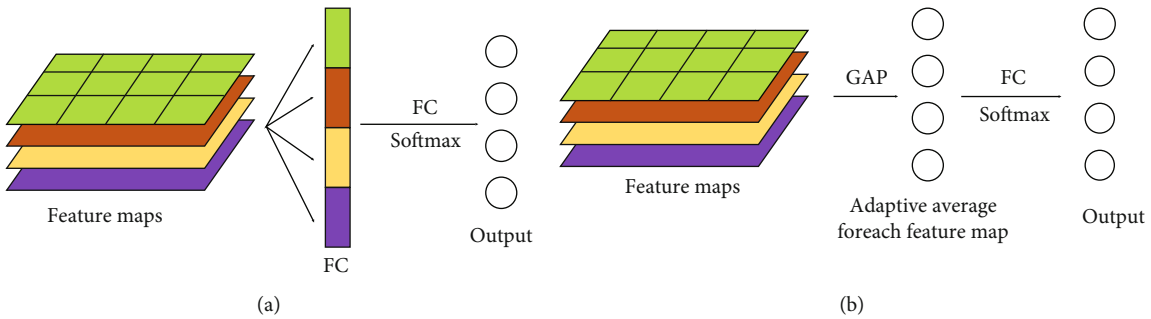
## 2. Image Processing and Methods

*2.1. Dataset.* In order to realize the classification and recognition of wood knot defects, firstly, the image information of 448 wood knot defects of spruce trees with seven kinds of

knot defects were collected on the website of Computer Laboratory of Department of Electrical Engineering, University of Oulu [16–18] (shown in Figure 1), and made them into a dataset that can simulate the actual use scene of ReSENet-18 model. Then, the preprocessing operations such as image scaling and adding noise were carried out to realize data augmentation. Finally, the dataset was divided into three parts: a training set, a verification set, and a testing set for training, verification, and testing.

*2.2. Data Preprocessing and Augmentation.* The dataset of wood knot defects with 448 images was divided into a training set, a verification set, and a testing set according to the ratio of $6:2:2$, which refers to 268 training images, 90 verification images, and 90 testing images, respectively (Table 1). The powerful generalization ability of the convolutional neural network is based on a large amount of data; thus, the model will induce the overfitting problem when the amount of data is not large enough which greatly limits the generalization ability [19–22]. Data augmentation technology [23, 24] was always used to expand the dataset of wood knot
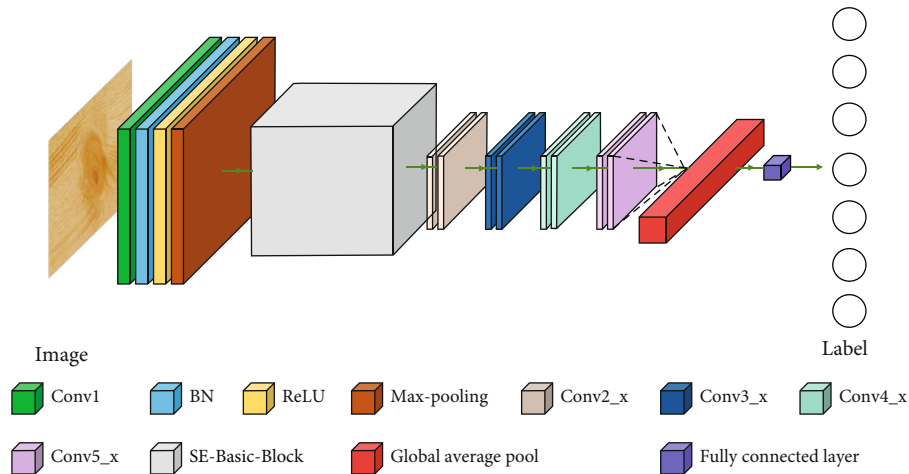
FIGURE 7: Architecture of the proposed ReSENet-18 model.

TABLE 3: Experimental environment.

| | Hardware environment | | Software environment |
|---|---|---|---|
| Memory | 16.00GB | System | Windows 10 |
| CPU | Intel Core i5-4210H 2.90GHz (2 core) | Environment configuration | Pytorch-gpu 1.0.0 + Python 3.7.3 + cuda 8.0 + cudnn7.1.3 |
| Graphics card | NVIDIA GeForce GTX 960 M(2G) | | |

TABLE 4: Training parameters.

| Related parameter | Value | Meaning |
|---|---|---|
| Batch size | 128 | Number of pictures per training |
| Learning rate | 1e-4 | Initial learning rate |
| Epoch | 200 | Training iteration times |
| CUDA | Enable | Computer unified device architecture |
| CUDNN | Enable | A GPU acceleration library for deep neural networks |

defects using color digital image processing technology to expand the data set and add it to the original image dataset; the problem of insufficient data can be easily solved.

The preprocessing of wood knot defect images was completed by simulating the change of angle, noise, and color of different tree species. In order to simulate these changes, the images of wood knots were mirrored horizontally, rotated by 180°, and mirrored vertically to simulate different angles of actual images. By the operation of increasing hue by 10, the color of different defects in actual image acquisition is simulated. At the same time, in order to simulate the noise that may appear in the process of image acquisition, appropriate Gaussian noise and salt-and-pepper noise are added to the image of wood knot defect to further enhance the dataset. In this work, the results of the data augmentation are shown in Figure 2. After data augmentation, the size of the dataset of wood knot defects is expanded from 448 images to 3136 images (7 times expansion). The number of training set, verification set, and testing set is 1885, 636, and 615, respectively, which can effectively reduce the overfitting phenomenon of the convolutional neural network during the training phase.
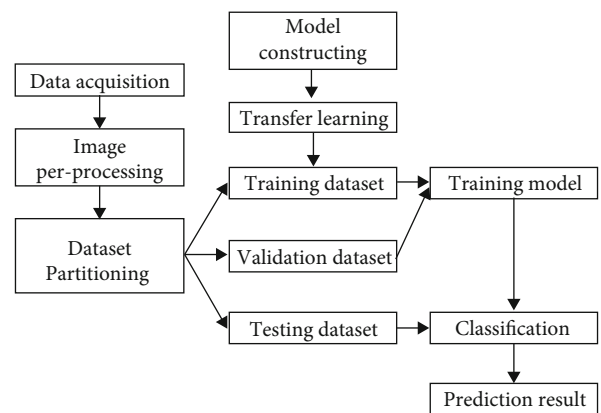


FIGURE 8: Wood knot defect detection process.

### 2.3. Improved Deep Convolutional Neural Network- (DCNN-) Based Method

*2.3.1. ResNet-18.* ResNet-18 consists of a convolutional layer and eight residual building blocks. A residual building block is the basic structure of the ResNet-18 network. The structure

FIGURE 9: The model was trained with a training dataset and validating datasets: (a) loss value and (b) accuracy value.





FIGURE 10: Confusion matrix of the model with 99.02% accuracy.

of the residual building block [25, 26] is shown in Figure 3. A kind of short-cut is used to skip the convolutional layer [27]. The input vector and the vector output through the convolutional layer can be added directly [28] and then output through the rectified linear unit (ReLU) activation function. This method can powerfully alleviate the problem of a vanishing gradient or exploding gradient caused by the increase of neural network depth and can eventually improve the recognition accuracy of wood knot defects.

The output of the residual building block is written as follows:

$$y = F(x) + x, \tag{1}$$

where $F$ presents the residual function and $x$ and $y$ stand for the input and output, respectively.

ResNet-18 consists of 17 convolutional layers, a max-pooling layer with the filter size of $3 \times 3$, and a fully con-

nected layer. A classical ResNet-18 model involves 33.16 million parameters, in which ReLU activation function and batch normalization (BN) are applied to the back of entire convolutional layers in "basic block." The structure of ResNet-18 is shown in Table 2 [27].

*2.3.2. SE-Basic-Block Module.* The SE-Basic-Block module has been used during the champion of ImageNet 2017 classification competition [29]. The structure is shown in Figure 4, which mainly includes squeeze and excitation [30]. The input image has the size of $W \times H \times C$, where $W$ and $H$ represent the width and height, respectively, and $C$ represents the number of channels. The structure of the SE module is uncomplicated and easy to implement. It can be easily embedded into the existing network framework. The SE module mainly studies the correlation between channels, which only increases a small amount of calculation but can achieve better results.

The attention mechanism of the SE module is mainly realized by multiplying the fully connected layer and input vector for feature fusion. Assume that the size of the input image is $H \times W \times C$, after passing through the global pooling layer and the fully connected layer, the input image is stretched to $1 \times 1 \times C$ and then multiplied with the original image to give weights to each channel. In the denoising task, each noise point is given weight, the low weight noise points are removed automatically, and the high weight noise points are retained. During this process, the network running efficiency can be improved, the parameters and computational cost can be reduced, and the recognition accuracy is improved [31]. As shown in Figure 4, by processing the feature map of convolutional, a one-dimensional vector with the same number of channels is obtained as the evaluation score of each channel [32], and then, the score is used for the corresponding channel to get the result.

The SE module can be embedded into the residual basic block of ResNet-18. Figure 5 shows the combined structure of the SE module and residual basic block module.

*2.3.3. Transfer Learning of DCNNs.* Due to the small size of the data in this experiment and a certain depth of the proposed network (ReSENet-18), it is easy to induce the overfitting problem in the training process, which leads to a poor recognition ability [33, 34]. In this case, the transfer learning is used to pretrain the deep learning model and then retrain for the wood knot defect detection task using the dataset in this study, which can make our model converge rapidly; thus, a lot of training time can be saved. The deep learning model includes a hierarchical architecture with various layers to learn the complex features of images with wood knot defects [35, 36]. Finally, all these layers are connected to the final fully connected layer classifier to obtain the final results. In the transfer learning, ResNet, VGG and AlexNet models have been trained in ImageNet [37], so that the better classification performance of wood knot defects can be achieved with less training time.

*2.3.4. Global Average Pooling.* The fully connected layer is usually used as a classifier of CNN, but too many parameters of the fully connected layer will increase the calculation amount of the network and thus slow down the training speed and also easily appear the overfitting problem [38]. Global average pooling (GAP) is a global average of all pixels in the feature map of each channel and obtains the output of each feature map [39–41]. GAP directly removes the features of black box in the fully connected layer and gives each channel practical significance; then, the vectors composed of these output features will be sent to the classifier for classification directly [42]. Figure 6 shows the comparison between the fully connected layer and the global average pool layer.

*2.3.5. Overall Architecture.* ReSENet-18 is a deep neural network based on the residual structure and attention mechanism (Figure 7). The main features of the architecture are described below.

ReSENet-18 network based on deep learning consists of a ResNet-18, a residual basic block and a squeeze-and-

TABLE 5: The evaluation index values of network.

| Classes | Model | *P* | *R* | *F1* | FAR |
|---|---|---|---|---|---|
| Decayed knot | LeNet-5 | — | 0% | — | 3.11% |
| | AlexNet | 72% | 94.74% | 81.82% | 0.17% |
| | VGGNet-16 | 90.91% | 52.63% | 66.67% | 1.5% |
| | GoogLeNet | 65.22% | 78.95% | 71.43% | 0.68% |
| | MobileNet V2 | 100% | 68.42% | 81.25% | 1.00% |
| | ReSENet-18 | 100% | 94.74% | 97.30% | 0.17% |
| Dry knot | LeNet-5 | 66.39% | 82.29% | 73.49% | 3.46% |
| | AlexNet | 96.81% | 94.79% | 95.79% | 0.97% |
| | VGGNet-16 | 90.48% | 98.96% | 94.53% | 0.20% |
| | GoogLeNet | 93.75% | 78.13% | 85.23% | 3.95% |
| | MobileNet V2 | 95.83% | 95.83% | 95.83% | 0.78% |
| | ReSENet-18 | 100% | 100% | 100% | 0% |
| Edge knot | LeNet-5 | 91.49% | 94.51% | 92.98% | 0.97% |
| | AlexNet | 94.74% | 98.90% | 96.78% | 0.19% |
| | VGGNet-16 | 94.68% | 97.80% | 96.21% | 0.39% |
| | GoogLeNet | 97.67% | 92.31% | 94.92% | 1.33% |
| | MobileNet V2 | 96.70% | 96.70% | 96.7% | 0.58% |
| | ReSENet-18 | 100% | 100% | 100% | 0% |
| Encased knot | LeNet-5 | 95.45% | 52.5% | 67.74% | 3.23% |
| | AlexNet | 100% | 90% | 94.74% | 0.70% |
| | VGGNet-16 | 100% | 90% | 94.74% | 0.70% |
| | GoogLeNet | 100% | 75% | 85.71% | 1.72% |
| | MobileNet V2 | 97.30% | 90% | 93.51% | 0.70% |
| | ReSENet-18 | 100% | 95% | 97.44% | 0.35% |
| Horn knot | LeNet-5 | 83.33% | 51.02% | 63.29% | 4.13% |
| | AlexNet | 97.96% | 97.96% | 97.96% | 0.18% |
| | VGGNet-16 | 88.89% | 97.96% | 93.20% | 0.18% |
| | GoogLeNet | 87.27% | 97.96% | 92.31% | 0.18% |
| | MobileNet V2 | 88% | 89.80% | 88.89% | 0.89% |
| | ReSENet-18 | 100% | 100% | 100% | 0% |
| Leaf knot | LeNet-5 | 70.89% | 84.85% | 77.24% | 1.88% |
| | AlexNet | 98.39% | 92.42% | 95.31% | 0.91% |
| | VGGNet-16 | 98.15% | 80.30% | 88.33% | 2.33% |
| | GoogLeNet | 98.25% | 84.85% | 91.06% | 1.81% |
| | MobileNet V2 | 90% | 95.45% | 92.64% | 0.55% |
| | ReSENet-18 | 97.01% | 98.48% | 97.74% | 0.18% |
| Sound knot | LeNet-5 | 85.39% | 91.2% | 88.20% | 6.40% |
| | AlexNet | 96.8% | 96.8% | 96.8% | 2.22% |
| | VGGNet-16 | 95.33% | 98% | 96.65% | 1.41% |
| | GoogLeNet | 88.93% | 99.6% | 93.96% | 0.30% |
| | MobileNet V2 | 96.06% | 97.6% | 96.82% | 1.68% |
| | ReSENet-18 | 99.20% | 99.20% | 99.2% | 0.55% |

excitation (SE) module. The ReSENet-18 model has 22 layers, including 8 parts: Conv1, SE-Basic-Block, Conv2_x, Conv3_x, Conv4_x, Conv5_x, a global average pooling layer, and a fully connected layer. The first part (Conv1) includes a convolutional layer, a batch normalization layer, a ReLU

TABLE 6: Accuracy of different models.

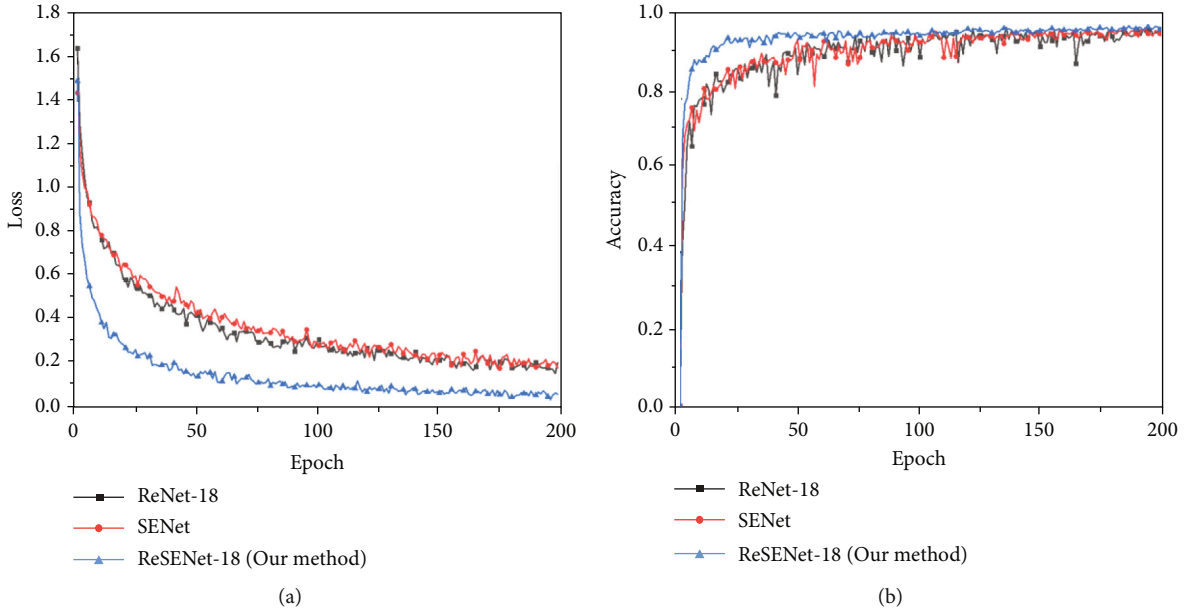| Different models | Accuracy | Different models | Accuracy |
| --- | --- | --- | --- |
| Add a logsoftmax classifier and a NLLLoss function | 14.08% | Add a ReLU and a fully connected layer | 14.08% |
| Add a softmax classifier | 14.40% | Without ReLU | 84.78% |
| Without batch normal | 96.07% | Add a convolutional layer, a BN layer, and a ReLU function | 96.24% |
| Without squeeze-and-excitation basic-block module | 98.53% | ReSENet-18 (our method) | 99.02% |



FIGURE 11: Training and prediction results of the ReSENet-18 model and its component models: (a) training loss value and (b) validation accuracy value.

activation function, and a max-pooling layer. The convolutional layer with the kernel size of $7 \times 7$, stride of 2, padding of 3, and the max-pooling layer with the kernel size of $3 \times 3$, stride of 2, padding of 1 were employed. Adding the max-pooling layer helps to reduce the dimensions and the parameters of model, to expand the receptive fields, and to retain important feature information. The second part (SE-Basic-Block) consists of a residual basic block and a squeeze-and-excitation (SE) module. There are two convolutional layers in the second part. The SE module was embedded into the residual basic block to form a SE-Basic-Block. The structure of the SE-Basic-Block module is shown in Figure 5. In the proposed SE-Basic-Block module, two convolutional layers with the kernel size of $3 \times 3$ and the stride of 1 were used. The first convolutional layer is followed by a BN layer and a ReLU activation function, while the second convolutional layer is only followed by a BN layer. As discussed above, the SE module mainly includes two parts. The first is squeeze, which makes the input image global average pooling; then, the feature map is compressed into a $1 \times 1 \times C$ vector. The second is excitation, which is composed of two fully connected layers and two activation functions (ReLU and Sigmoid). The input of the first fully connected layer is $1 \times 1 \times C$, and the output is $1 \times 1 \times C \times 1/r$, where $r$ is a scaling parameter which is used to reduce the number of chan-

nels so as to reduce the amount of calculation. The input of the second fully connected layer is $1 \times 1 \times C \times 1/r$, and the output is $1 \times 1 \times C$. In this paper, $r = 16$ is used. After getting the vector of $1 \times 1 \times C$, the initial feature map and the vector of $1 \times 1 \times C$ will be scaled. The size of the original feature map is $W \times H \times C$, the weight value of each channel output by the SE module is multiplied by the two-dimensional matrix of the corresponding channel of the original feature map, and the final output result is obtained. Parts three to six (Conv2_x, Conv3_x, Conv4_x, and Conv5_x) are shown in Figure 2. The seventh part (global average pool) uses AdaptiveAvgPool function, and the output size of this layer was set to $1 \times 1$. The eighth part (fully connected layer) is the classifier of ReSENet-18. Its output was set to 7, which corresponds to the types of datasets to train and classify.

ReSENet-18 takes RGB image with the random size as input, and then, the image is adjusted to $85 \times 85$ in batch. The input layer of ReSENet-18 is followed by a series of convolutional blocks and a subsampling layer. The CNN structure used in this paper is a variant of ResNet-18, and the feature extraction part of this network is similar to ResNet-18. We used 17 convolutional layers of ResNet-18 to self-study the features of input RGB images from low to high. With the deepening of convolutional layers, the resolution of feature map is reduced, and more abstract high-level

TABLE 7: Wood knots classification results of ReSENet-18 and its component networks.

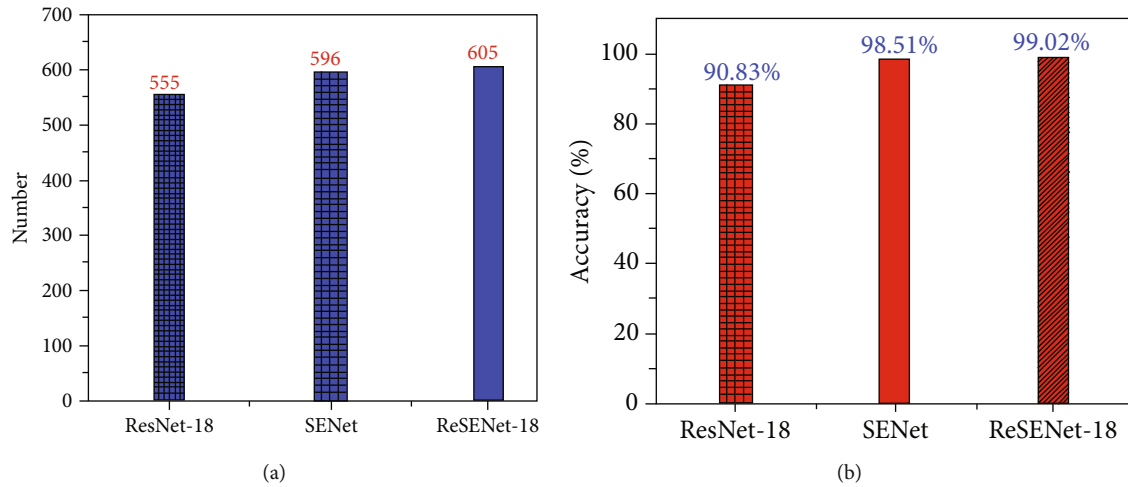| Method | Actual category | Predict category | | | | | | | Total |
| | | Decayed knot | Dry knot | Edge knot | Encased knot | Horn knot | Leaf knot | Sound knot | |
|---|---|---|---|---|---|---|---|---|---|
| ReSENet-18 (our method) | Decayed knot | 18 | 0 | 0 | 0 | 0 | 0 | 1 | **19** |
| | Dry knot | 0 | 96 | 0 | 0 | 0 | 0 | 0 | **96** |
| | Edge knot | 0 | 0 | 91 | 0 | 0 | 0 | 0 | **91** |
| | Encased knot | 0 | 0 | 2 | 38 | 0 | 0 | 0 | **40** |
| | Horn knot | 0 | 0 | 0 | 0 | 49 | 0 | 0 | **49** |
| | Leaf knot | 0 | 0 | 0 | 0 | 0 | 65 | 1 | **66** |
| | Sound knot | 0 | 0 | 0 | 0 | 0 | 2 | 248 | **250** |
| ResNet-18 | Decayed knot | 17 | 0 | 1 | 0 | 0 | 0 | 1 | **19** |
| | Dry knot | 4 | 86 | 1 | 2 | 0 | 0 | 3 | **96** |
| | Edge knot | 0 | 2 | 89 | 0 | 0 | 0 | 0 | **91** |
| | Encased knot | 2 | 2 | 0 | 33 | 0 | 0 | 3 | **40** |
| | Horn knot | 0 | 0 | 1 | 0 | 47 | 0 | 1 | **49** |
| | Leaf knot | 0 | 3 | 0 | 0 | 19 | 42 | 2 | **66** |
| | Sound knot | 4 | 5 | 0 | 0 | 0 | 0 | 241 | **250** |
| SENet | Decayed knot | 16 | 0 | 1 | 0 | 0 | 0 | 2 | **19** |
| | Dry knot | 0 | 95 | 1 | 0 | 0 | 0 | 0 | **96** |
| | Edge knot | 0 | 1 | 89 | 0 | 0 | 0 | 1 | **91** |
| | Encased knot | 0 | 1 | 1 | 36 | 0 | 0 | 2 | **40** |
| | Horn knot | 0 | 0 | 1 | 0 | 48 | 0 | 0 | **49** |
| | Leaf knot | 0 | 0 | 0 | 0 | 1 | 65 | 0 | **66** |
| | Sound knot | 0 | 1 | 0 | 0 | 0 | 2 | 247 | **250** |



FIGURE 12: Prediction results of the ReSENet-18 model and its component models.

features are extracted. Then, inspired by the success of SENet, we use a SE-Basic-Block module to improve the receptive field. The SE-Basic-Block was embedded between Conv1 and Conv2_x. To improve the sensitivity of the model to channel features, the original features are recalibrated in the channel dimension, so that the model can automatically learn the importance features of different channels. Then, the fully connected layer is replaced by the global average pooling layer to reduce the training parameters, to accelerate the convergence speed of the model, and to improve the classification accuracy of the model.

Finally, we use transfer learning in ReSENet-18 to maximize the collected data and prevent overfitting problem. During the training phase, the weights of 17 convolutional layers except SE-Basic-Block are initialized by the pretraining model of ResNet-18. In this paper, all the parameters would
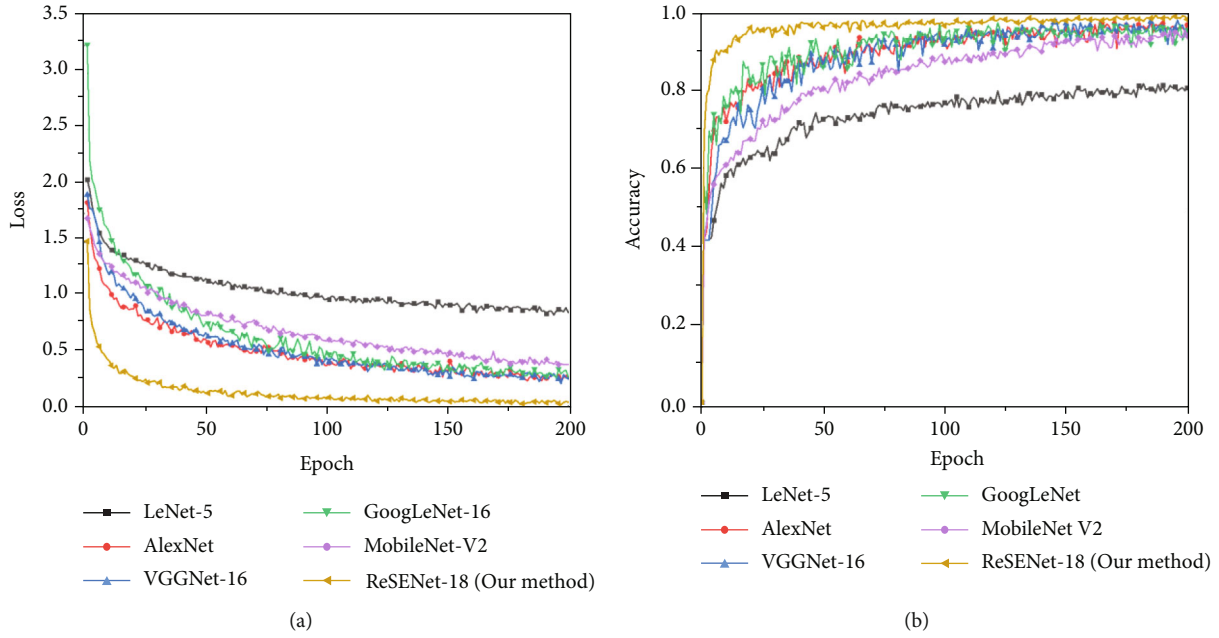
FIGURE 13: Training and prediction results of ReSENet-18 model and others classical CNN models: (a) training loss value and (b) validation accuracy value.

TABLE 8: Prediction results of ReSENet-18 model and other classical CNN models.

| Different models | Number | Accuracy | Parameters (million) |
|---|---|---|---|
| LeNet-5 | 495 | 81.01% | 0.06 |
| AlexNet | 586 | 95.91% | 60 |
| VGGNet-16 | 576 | 94.27% | 138 |
| GoogLeNet | 557 | 91.16% | 6.8 |
| MobileNet V2 | 580 | 94.93% | 3.5 |
| ReSENet-18 (our method) | 605 | 99.02% | 33 |

not be frozen. After loading the pretraining weights, the current dataset of wood knot defects was used to retrain the whole model, which can not only improve the accuracy and speed of training but also improve the recognition ability of the model in the current dataset of wood knot defects. This is very important for effective and stable feature learning.

*2.4. Training.* The proposed ReSENet-18 was used and trained on one GPU (GTX 960M 2G). The experimental environment is presented in Table 3. The parameter configuration is shown in Table 4. The model using the Adam optimization algorithm and the cross-entropy loss function was trained for 200 epochs, whose batch size is 128 and learning rate is 1e-4.

The flow diagram of the detection process of wood knot defects is shown in Figure 8. First, the images of knot defects were collected from logs. The original datasets were classified by experienced professionals according to the types of defects. Then, the datasets were divided into a training dataset, a verification dataset, and a testing dataset. Subsequently,

the proposed ReSENet-18 model was trained on the dataset of wood knot defects. Finally, the model is used to detect the defect types of each image in the testing dataset.

Figure 9 shows the process of training the model using the training and validation datasets. The best accuracy is 99.062%, the best loss is about 0.044, and the overall accuracy in the test phase is about 99.02%.

## 3. Experimental Results and Discussion

*3.1. Comparisons of Model Performance.* To evaluate the performance of the proposed model, the dataset was randomly divided and trained 10 times in our case. The classification accuracy of these 10 models is 99.02%, 98.20%, 98.20%, 98.20%, 98.20%, 98.36%, 97.71%, 97.87%, 99.02%, and 98.20%, respectively. The average classification accuracy of the 10 models is $98.30 \pm 0.16\%$ and the variance is 0.40%, which indicates a good stability. Taking the first model with the accuracy of 99.02% as an example, the confusion matrix is established by analyzing the predicted labels and true labels of the testing dataset, as shown in Figure 10. All the correct predictions are on the square of the diagonal.

Figure 10 shows that the recognition accuracy of the model for dry knot, edge knot, and horn knot are 100%. In the testing set, the total number of images is 611. The classification accuracies for decayed knot defect and encased knot defect are both 95%, which is due to the small number of decayed knot images and the quite shape difference of encased knot. For leaf knot defect, the classification accuracy is 98% due to the similarity of geometric features between the horn knot and the leaf knot which is easy to be mixed. The classification accuracy of sound knot is 99% which is due to the largest number of sound knot images.
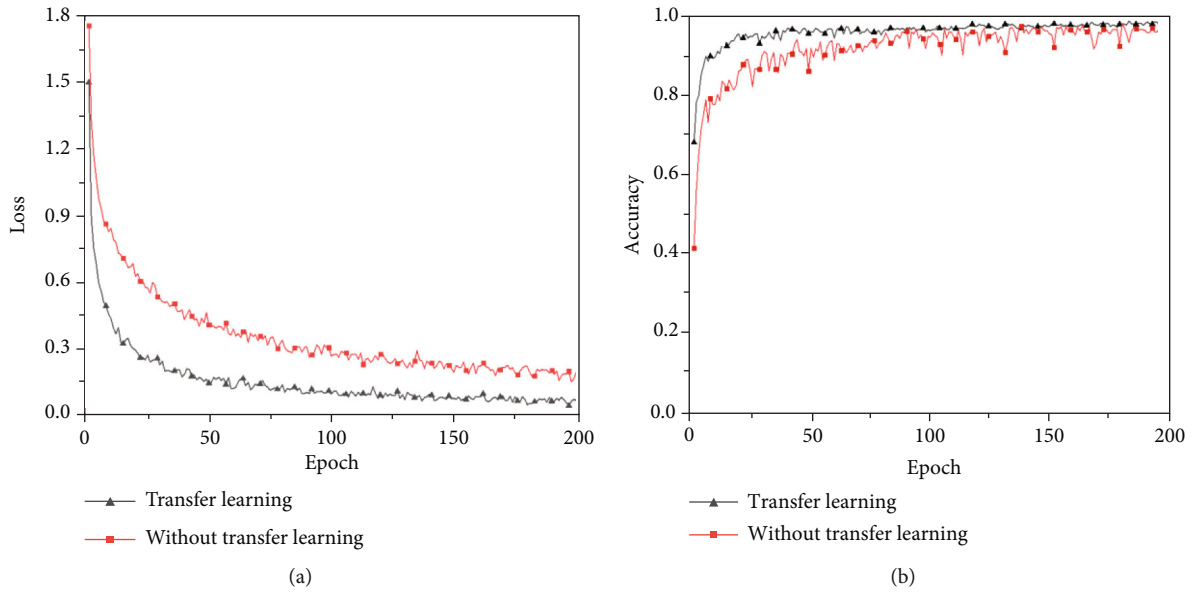
FIGURE 14: Training results of ReSENet-18 model with and without transfer learning: (a) training loss value and (b) validation accuracy value.
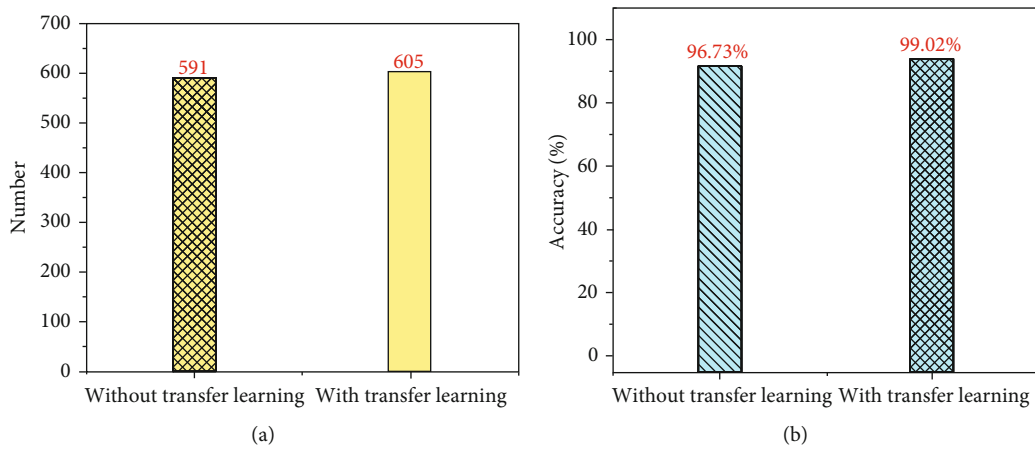


FIGURE 15: Prediction results of ReSENet-18 model with and without transfer learning.

To evaluate the performance of the ReSENet-18, the precision ($P$), recall ($R$), f1-score ($F1$), and false alarm rate (FAR) were applied for the evaluation shown as follows:

$$P = \frac{T_{ii}}{T_{ii} + T_{ij}}, \tag{2a}$$

$$R = \frac{T_{ii}}{T_{ii} + T_{ji}}, \tag{2b}$$

$$FAR = \frac{T_{ij}}{T_{ij} + T_{jj}}, \tag{2c}$$

$$F1 = 2\frac{P \cdot R}{P + R}, \tag{2d}$$

where $T_{ii}$, $T_{ij}$, $T_{ji}$, and $T_{jj}$ represent the confusion matrix components.

Table 5 shows the precision, recall, f1-score, and false alarm rate of ReSENet-18 for the seven types of wood knot defect and the other five models for comparison. It can be seen from Table 5 that the four indicators of ReSENet-18 are the best in the recognition of five knots (decayed knot, dry knot, edge knot, encased knot, and horn knot) compared with other five classical CNN models, i.e., LeNet-5, AlexNet, VGGNet-16, GoogLeNet, and MobileNet V2. In the recognition of leaf and sound knots, some indicators of ReSENet-18 are slightly worse than other models. For example, the precision of ResNet-18 is higher than that of LeNet-5 and Mobile-Net V2, but slightly worse than that of AlexNet, VGGNet-16, and GoogLeNet. Among the other three indicators ($R$, $F1$, and FAR), ReSENet-18 is still the best of the six models. In the recognition of sound knot, recall and false accept rate of ReSENet-18 are slightly worse than GoogLeNet and precision and f1-score are better than GoogLeNet. Compared with LeNet-5, AlexNet, VGGNet-16, and MobileNet V2, all the
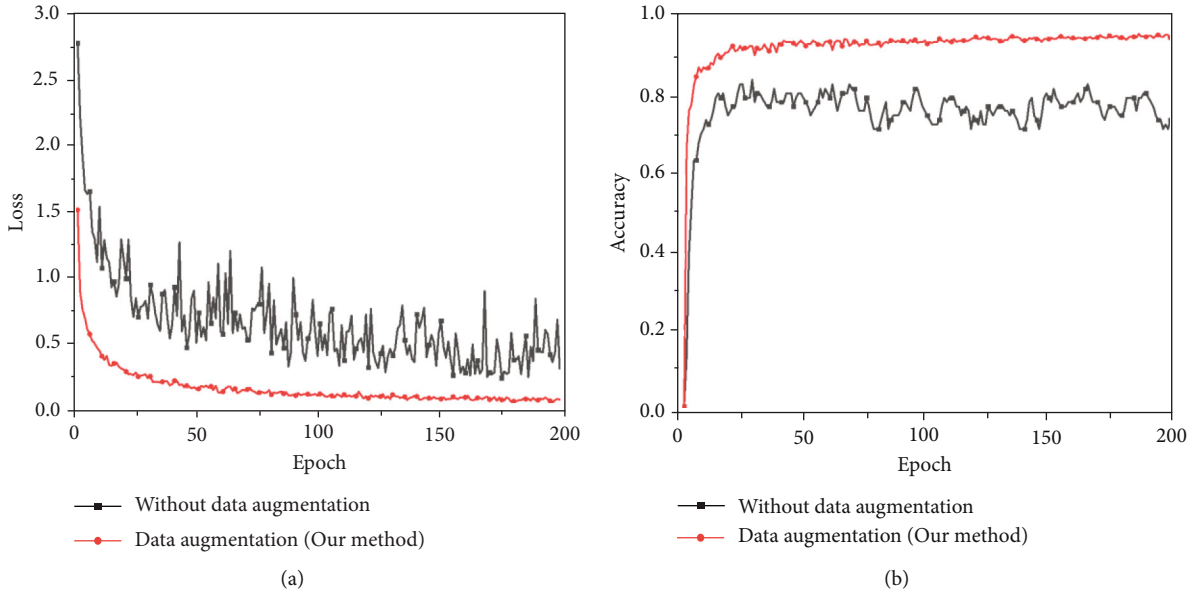
FIGURE 16: Training results of the ReSENet-18 model with and without data augmentation: (a) training loss value and (b) validation accuracy value.
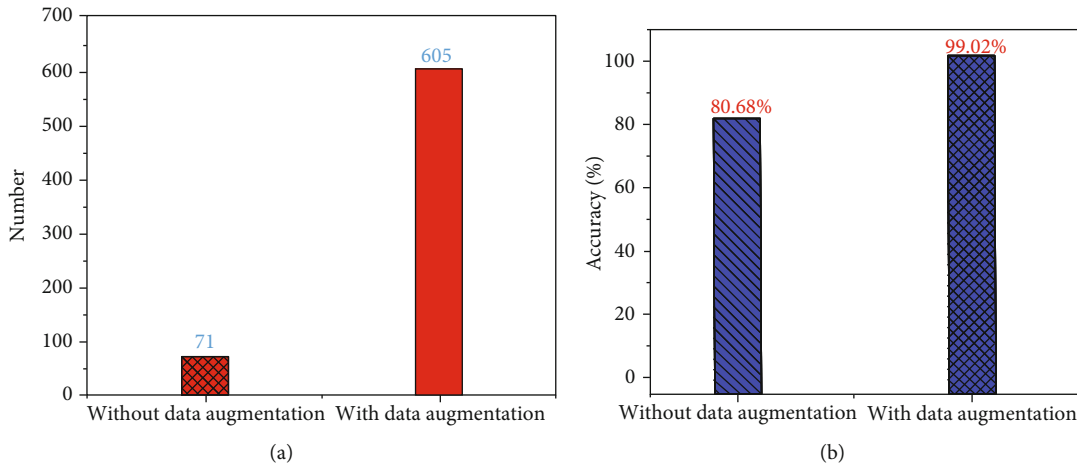


FIGURE 17: Prediction results of the ReSENet-18 model with and without data augmentation.

indicators of ReSENet-18 are better than them. Based on the above analysis, although there is still room for improvement in a few indicators, compared with the other five methods, it can be seen that ReSENet-18 still has a good performance in the identification of wood knot defects.

*3.2. Comparison of the Accuracy of Different Models.* The ReSENet-18 model was modified, and the test results are shown in Table 6. The accuracy of without ReLU, BN, or SE Basic-Block models is 84.78%, 96.07%, and 98.53%, respectively. The accuracy of adding a BN layer and a ReLU is 96.24%. The accuracy of adding a logsoftmax classifier and replacing Adam with NLLLoss, adding a ReLU and a fully connected layer, and adding a softmax classifier are 14.08%, 14.08%, and 14.40%, respectively. By comparing the performance of different structures of the ReSENet-18

model, the experimental results show that the ReSENet-18 model has the highest recognition accuracy.

*3.3. Convergence and Prediction Accuracy Analysis of ReSENet-18 Network Model*

*3.3.1. Comparison with ReSENet-18 Model and Its Component Models.* As demonstrated above, ReSENet-18 is composed of ResNet-18, SE module, and residual basic block. Therefore, to test the performance of the ReSENet-18 model, it was compared with ResNet-18 (composed of residual basic block) and SENet (composed of SE module). Figure 11 shows the loss curve and accuracy curve of ResNet-18, SENet, and ReSENet-18 during the training phase, which are trained by the wood knot defects dataset. One could learn from Figure 11 that our proposed method has the lowest loss value and highest accuracy compared with ResNet-18 and SENet.
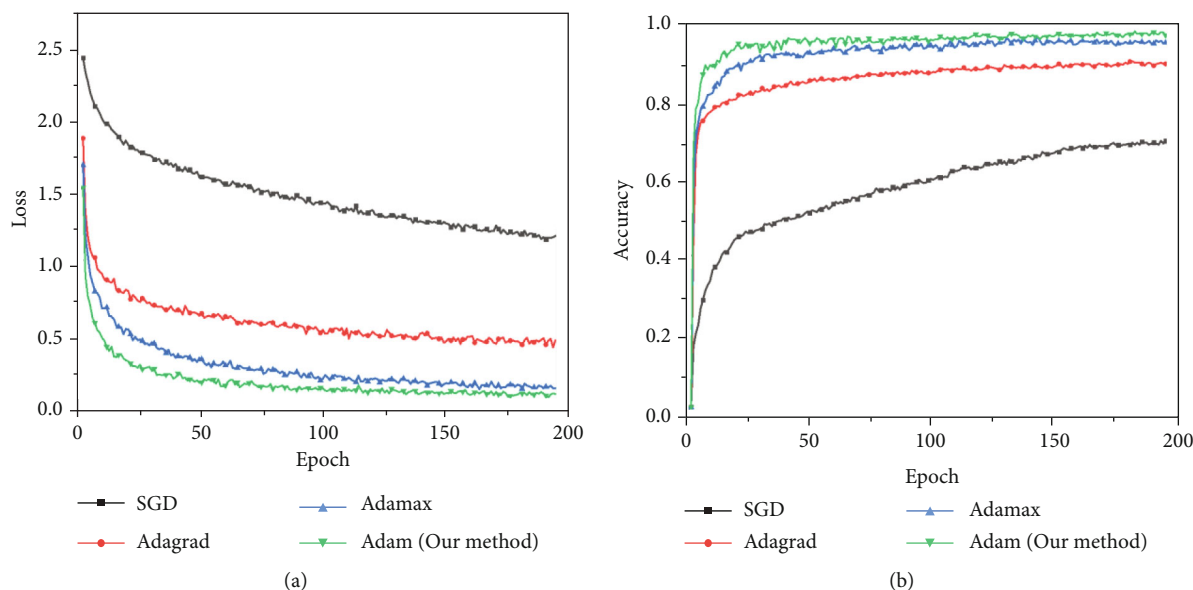
FIGURE 18: Training results of the ReSENet-18 model with different optimization algorithms: (a) training loss value and (b) validation accuracy value.

The classification results of wood knot defects are shown in Table 7, in which the italic entries represent the number of knot defects correctly identified by the corresponding model, and bold entries represent the total number of the wood knot defects. From Table 7 and Figure 12, ReSENet-18 has the best recognition effect among seven kinds of wood knot defects. Compared with other networks, the ResNet-18 network is relatively shallow in depth, so some degrees of underfitting phenomenon might appear, which leads to a low accuracy on the testing set. The SENet network has the largest number of layers among the three networks, but it can be seen from Table 7 and Figure 12 that the result of identification of SENet is not the best among them due to the increase of network layers and the appearance of overfitting phenomenon. The accuracy of the proposed model with lightweight of ResNet-18 in the testing dataset reaches 99.02% (Figure 12). At the same time, it can combine the features of channel into the network, which improves the feature extraction ability.

Based on the above analysis, ReSENet-18 has been proved to have the highest accuracy and fastest convergence speed in the wood knot defect dataset than other models.

*3.3.2. Comparison with Classical CNN Model.* Figure 13 shows the training of ReSENet-18 and other five CNN models which was mentioned in Section 3.1. These networks are trained through the dataset of wood knot defects. It can be seen that the ReSENet-18 network has the highest accuracy and the fastest convergence speed than other models on the wood knot defect dataset.

Table 8 compares the number and accuracy of the six network models on the testing dataset. The results show that the LeNet-5 model has the minimal training parameters, which may lead to the underfitting of the network which leads to the lowest accuracy. The parameters of GoogLeNet and MobileNet V2 models are slightly more than LeNet-5, but they are more complex than LeNet-5. It can be seen from Table 8 that

TABLE 9: Accuracy of different optimization algorithms.

| Different models | Number | Accuracy |
| --- | --- | --- |
| SGD algorithm | 429 | 70.21% |
| AdaGrad algorithm | 562 | 91.98% |
| AdaMax algorithm | 586 | 95.91% |
| Adam algorithm (our method) | 605 | 99.02% |

their accuracy is improved compared with that of LeNet-5. VGGNet-16 has the maximum parameters among the six models, and AlexNet follows. However, the accuracy is lower than that of ReSENet-18 even through the increase of parameters and longer training time. Compared with VGGNet-16, ReSENet-18 is a kind of lightweight network. At the same time, it can use SE-Basic-Block to weight and recalibrate features. It has stronger feature extraction ability and higher accuracy. It can be seen that, compared with other models, we have improved the performance of the ReSENet-18 model by adding appropriate parameters, while maintaining the robustness and efficiency of the model. Among the six models, the recognition accuracy of the ReSENet-18 model is the highest.

*3.3.3. Transfer Learning.* A pretraining model of ResNet-18 which includes 1.2 million color images and 1000 categories is used in this study. The weight of the pretrained model is taken as the initial weight of the dataset of wood knot defects. Figures 14 and 15 show the influence and prediction results of transfer learning on the classification accuracy and convergence speed of the ReSENet-18 model. It can be seen that the convergence speed and accuracy of the model have been improved after using transfer learning. The experimental results show that the accuracy of the model with transfer learning is 2.29% higher than that of the model without transfer learning on the testing dataset. Therefore, better convergence can be achieved using transfer learning.
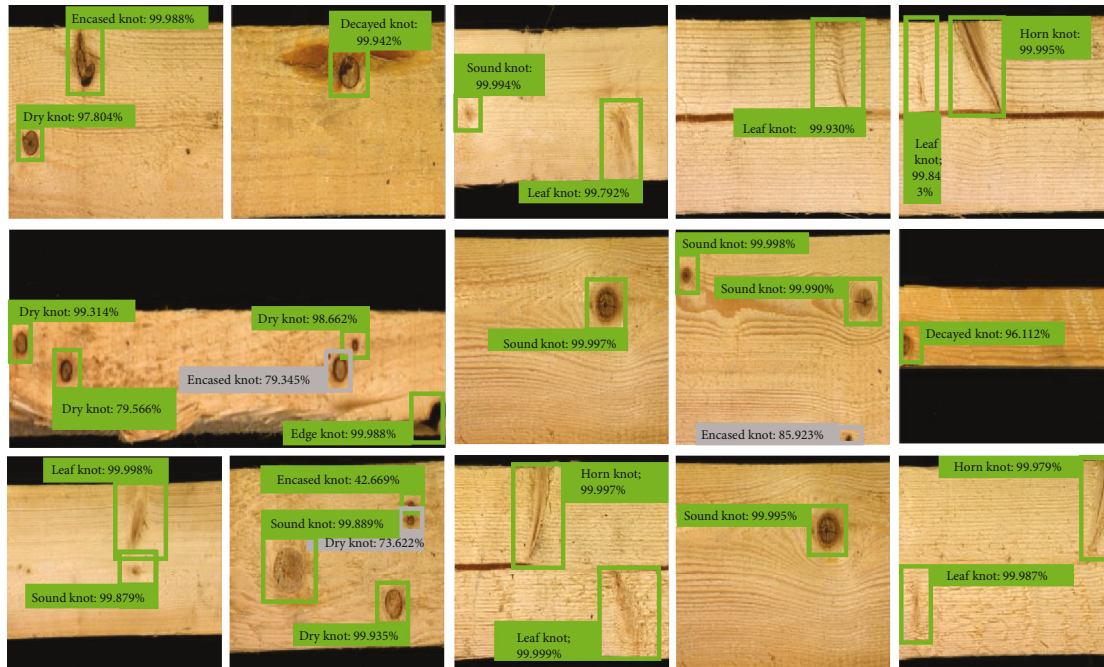
FIGURE 19: The recognition results of the ReSENet-18 model after training.

*3.3.4. Data Augmentation.* Figure 16 shows the convergence and recognition accuracy curves of ReSENet-18 with and without data augmentation. Under the same experimental conditions, ReSENet-18 was trained on 3136 images after data augmentation, and the final classification accuracy in testing dataset reached 99.02% while the accuracy was 80.68% before data augmentation, which is shown in Figure 17.

*3.3.5. Comparison of Optimization Algorithms.* The optimization algorithm has an important influence on the model performance. In this study, the Adam optimization algorithm is used and compared with SGD, AdaGrad, and AdaMax, which are shown in Figure 18. The results show that the model with Adam algorithm has fastest convergence speed. Table 9 shows the prediction results of these four optimization algorithms under the same environment. The results show that the accuracy in the testing phase is 70.21% for SGD algorithm, 91.98% for AdaGrad algorithm, 95.91% for AdaMax algorithm, and 99.02% for Adam algorithm. It can be seen that the ReSENet-18 model has the best training effect using the Adam optimization algorithm.

*3.3.6. Recognition Results of Different Kinds of Wood Knot Defects.* "Correct recognition" was used to mark in green, and the "Wrong recognition" was marked in grey in this study. Details of the identification, such as the name and probability of wood knot defects, are displayed next to each label. Figure 19 shows seven wood knot defects and the corresponding identification results.

It can be seen that most of the wood knot defects in the image were correctly identified. Due to the shape of some wood knot defects being similar to other defects, there is no clear feature to extract under this background to induce a few regions incorrectly identified. In addition, the shape of the defect is blurred due to the low resolution of some images, which also makes the extracted features different from those in the training set. In the most cases, our method (ReSENet-18) still has a high accuracy.

## 4. Conclusions

In conclusion, a novel convolutional neural network model ReSENet-18 is proposed. In the feature extraction part of the network, the SE module is embedded into the residual basic blocks to form SE-Basic-Block. The classifier of the network selects the global average pool to replace the fully connected layer after the convolutional layer at the end to speed up the convergence speed and reduce the model parameters. 2521 images of wood knot defects were used for training after 200 training epochs. Experimental results show that the accuracy of ReSENet-18 in the test phase reaches 99.02%, which is 8.19% higher than the classical ResNet-18 (90.83%). In addition, when various wood knot defects are detected by this method, a large amount of image preprocessing and manual feature extraction are not required, which greatly improves the recognition efficiency. This means that ReSENet-18 has a potential application in wood nondestructive testing and knot defect identification, and it provides a feasible solution for future wood knot defect identification.

## Data Availability

The datasets, codes, and weight files used to support the findings of this study are available from the corresponding author upon request.

## Conflicts of Interest

## Acknowledgments

## References

[1] R. Norlander, J. Grahn, and A. Maki, "Wooden knot detection using ConvNet transfer learning," in *Scandinavian Conference on Image Analysis*, pp. 263–274, Cham, 2015.

[2] W. Pölzleitner and G. Schwingshakl, "Real-time surface grading of profiled wooden boards," *Industrial Metrology*, vol. 2, no. 3-4, pp. 283–298, 1992.

[3] Z. Qiu, *A simple machine vision system for improving the edging and trimming operations performed in hardwood sawmills*, Doctoral dissertation, Virginia Tech, 1996.

[4] D. Schmoldt, P. Li, and A. Abbott, "Machine vision using artificial neural networks with local 3D neighborhoods," *Computers and Electronics in Agriculture*, vol. 16, no. 3, pp. 255–271, 1997.

[5] D. Yadav and A. K. Yadav, "A novel convolutional neural network based model for recognition and classification of apple leaf diseases," *Traitement du Signal*, vol. 37, no. 6, pp. 1093–1101, 2020.

[6] X. Zhu, M. Zhu, and H. Ren, "Method of plant leaf recognition based on improved deep convolutional neural network," *Cognitive Systems Research*, vol. 52, pp. 223–233, 2018.

[7] D. Patrício and R. Rieder, "Computer vision and artificial intelligence in precision agriculture for grain crops: a systematic review," *Computers and Electronics in Agriculture*, vol. 153, pp. 69–81, 2018.

[8] J. Lemley, S. Bazrafkan, and P. Corcoran, "Deep learning for consumer devices and services: pushing the limits for machine learning, artificial intelligence, and computer vision," *IEEE Consumer Electronics Magazine*, vol. 6, no. 2, pp. 48–56, 2017.

[9] T. He, Y. Liu, Y. Yu, Q. Zhao, and Z. Hu, "Application of deep convolutional neural network on feature extraction and detection of wood defects," *Measurement*, vol. 152, p. 107357, 2020.

[10] L. Wenshu, S. Lijun, and W. Jinzhuo, "Study on wood board defect detection based on artificial neural network," *The Open Automation and Control Systems Journal*, vol. 7, no. 1, 2015.

[11] H. Mu, M. Zhang, D. Qi, and H. Ni, "The application of RBF neural network in the wood defect detection," *International Journal of Hybrid Information Technology*, vol. 8, no. 2, pp. 41–50, 2015.

[12] X. Ji, H. Guo, and M. Hu, "Features extraction and classification of wood defect based on Hu invariant moment and wavelet moment and BP neural network," in *Proceedings of the 12th International Symposium on Visual Information Communication and Interaction*, pp. 1–5, Shanghai, China, 2019.

[13] T. He, Y. Liu, C. Xu, X. Zhou, Z. Hu, and J. Fan, "A fully convolutional neural network for wood defect location and identification," *IEEE Access*, vol. 7, pp. 123453–123462, 2019.

[14] S. Liu, W. Jiang, L. Wu, H. Wen, M. Liu, and Y. Wang, "Real-time classification of rubber wood boards using an SSR-based CNN," *IEEE Transactions on Instrumentation and Measurement*, vol. 69, no. 11, pp. 8725–8734, 2020.

[15] M. Gao, J. Chen, H. Mu, and D. Qi, "A transfer residual neural network based on ResNet-34 for detection of wood knot defects," *Forests*, vol. 12, no. 2, p. 212, 2021.

[16] H. Kauppinen and O. Silven, "A color vision approach for grading lumber," in *Theory & Applications of Image Processing II-Selected Papers from the 9th Scandinavian Conference on Image Analysis*, pp. 367–379, Singapore, 1995.

[17] O. Silven and H. Kauppinen, "Recent developments in wood inspection," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 10, no. 1, pp. 83–95, 1996.

[18] H. Kauppinen and O. Silven, "The effect of illumination variations on color-based wood defect classification," in *Proceedings of the 13th International Conference on Pattern Recognition (13th ICPR)*, pp. 828–832, Vienna, Austria, 1996.

[19] L. Jiang, Y. Wang, Z. Tang, Y. Miao, and S. Chen, "Casting defect detection in X-ray images using convolutional neural networks and attention-guided data augmentation," *Measurement*, vol. 170, p. 108736, 2021.

[20] S. Liu, H. Jiang, Z. Wu, and X. Li, "Rolling bearing fault diagnosis using variational autoencoding generative adversarial networks with deep regret analysis," *Measurement*, vol. 168, p. 108371, 2021.

[21] W. Zhang, X. Li, X. Jia, H. Ma, Z. Luo, and X. Li, "Machinery fault diagnosis with imbalanced data using deep generative adversarial networks," *Measurement*, vol. 152, p. 107377, 2020.

[22] H. Kim and M. Kang, "A comparison of methods to reduce overfitting in neural networks," *International journal of advanced smart convergence*, vol. 9, no. 2, pp. 173–178, 2020.

[23] N. Farda, J. Lai, J. Wang, P. Lee, J. Liu, and I. Hsieh, "Sanders classification of calcaneal fractures in CT images with deep learning and differential data augmentation techniques," *Injury*, vol. 52, no. 3, pp. 616–624, 2021.

[24] S. Han, S. Oh, and J. Jeong, "Bearing fault diagnosis based on multiscale convolutional neural network using data augmentation," *Journal of Sensors*, vol. 2021, 14 pages, 2021.

[25] H. Ali, S. Kabir, and G. Ullah, "Indoor scene recognition using ResNet-18," *International Journal of Research Publications*, vol. 69, no. 1, p. 7, 2021.

[26] D. Zhu, Y. Yang, W. Zhai, F. Ren, C. Cheng, and M. Huang, "Geosot grid remote sensing intelligent interpretation model based on fine-tuning ResNet-18: a case study of construction land," in *IGARSS 2020-2020 IEEE International Geoscience and Remote Sensing Symposium*, pp. 2535–2538, Waikoloa, HI, USA, 2020.

[27] D. Sarwinda, R. Paradisa, A. Bustamam, and P. Anggia, "Deep learning in image classification using residual network (ResNet) variants for detection of colorectal cancer," *Procedia Computer Science*, vol. 179, pp. 423–431, 2021.

[28] F. Zeng, W. Hu, G. He, and C. Yue, "Imbalanced Thangka image classification research based on the ResNet network," *Journal of Physics: Conference Series*, vol. 1748, no. 4, pp. 042–054, 2021.

[29] H. Ma, G. Han, L. Peng, L. Zhu, and J. Shu, "Rock thin sections identification based on improved squeeze-and-excitation

networks model," *Computers & Geosciences*, vol. 152, p. 104780, 2021.

[30] F. Qi, Z. Xie, Z. Tang, and H. Chen, "Related study based on Otsu watershed algorithm and new squeeze-and-excitation networks for segmentation and level classification of tea buds," *Neural Processing Letters*, vol. 19, pp. 1–5, 2021.

[31] Y. Ying, N. Zhang, P. Shan, L. Miao, P. Sun, and S. Peng, "PSigmoid: improving squeeze-and-excitation block with parametric sigmoid," *Applied Intelligence*, vol. 9, pp. 1–3, 2021.

[32] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7132–7141, Salt Lake City, Utah, USA, 2018.

[33] Y. Pathak, P. Shukla, A. Tiwari, S. Stalin, and S. Singh, *Deep Transfer Learning Based Classification Model for COVID-19 Disease*, Irbm, 2020.

[34] V. Chouhan, S. Singh, A. Khamparia et al., "A novel transfer learning based approach for pneumonia detection in chest X-ray images," *Applied Sciences*, vol. 10, no. 2, p. 559, 2020.

[35] M. Loey, G. Manogaran, M. Taha, and N. Khalifa, "A hybrid deep transfer learning model with machine learning methods for face mask detection in the era of the COVID-19 pandemic," *Measurement*, vol. 167, p. 108288, 2021.

[36] Z. Wu, H. Jiang, K. Zhao, and X. Li, "An adaptive deep transfer learning method for bearing fault diagnosis," *Measurement*, vol. 151, p. 107227, 2020.

[37] A. Abbas, M. Abdelsamea, and M. Gaber, "DeTrac: transfer learning of class decomposed medical images in convolutional neural networks," *IEEE Access*, vol. 8, pp. 74901–74913, 2020.

[38] Q. Zhang, C. Bai, Z. Liu et al., "A GPU-based residual network for medical image classification in smart medicine," *Information Sciences*, vol. 536, pp. 91–100, 2020.

[39] H. Lee, J. Park, and J. Hwang, "Channel attention module with multiscale grid average pooling for breast cancer segmentation in an ultrasound image," *IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control*, vol. 67, no. 7, pp. 1344–1353, 2020.

[40] Y. Li and K. Wang, "Modified convolutional neural network with global average pooling for intelligent fault diagnosis of industrial gearbox," *Eksploatacja i Niezawodność*, vol. 22, 2020.

[41] W. Li, K. Liu, L. Zhang, and F. Cheng, "Object detection based on an adaptive attention mechanism," *Scientific Reports*, vol. 10, no. 1, pp. 1–3, 2020.

[42] X. Zhang and X. Zhang, "Global learnable pooling with enhancing distinctive feature for image classification," *IEEE Access*, vol. 8, pp. 98539–98547, 2020.

*Research Article*

# Data-driven Learning Algorithm of Neural Fuzzy Based Hammerstein-Wiener System

**Feng Li [ID],[1] Yinsheng Luo,[1] Naibao He,[1] Ya Gu,[2] and Qingfeng Cao[3]**

[1]*College of Electrical and Information Engineering, Jiangsu University of Technology, Changzhou 213001, China*
[2]*School of Electrical Engineering and Automation, Changshu Institute of Technology, Changshu, 215500 Jiangsu, China*
[3]*College of Electrical, Energy and Power Engineering, Yangzhou University, Yangzhou 225127, China*

Correspondence should be addressed to Feng Li; lifeng@jsut.edu.cn

A novel data-driven learning approach of nonlinear system represented by neural fuzzy Hammerstein-Wiener model is presented. The Hammerstein-Wiener system has two static nonlinear blocks represented by two independent neural fuzzy models surrounding a dynamic linear block described by finite impulse response model. The multisignal theory is designed for employing Hammerstein-Wiener system to separate parameter learning issues. To begin with, the output nonlinearity parameters are learned utilizing separable signal with different amplitudes. Furthermore, correlation analysis method is implemented for estimating linear block parameters using separable signal inputs and outputs; thereby, the interference of process noise is effectively handled. Finally, multi-innovation learning technology is introduced to improve system learning accuracy, and then, multi-innovation extended stochastic gradient algorithm is obtained for optimizing input nonlinearity and noise model using multi-innovation technique and gradient search method. The simulation results display that presented data-driven learning approach has the availability of learning Hammerstein-Wiener system.

## 1. Introduction

The real industrial processes are almost nonlinear systems to some extent, and linear approximation means are usually unacceptable, and nonlinear models should be taken into account that they can present the nonlinearity successfully. For this, block-oriented nonlinear models which are composed of linear dynamic block and static nonlinear functions for instance Hammerstein model and Wiener model have been performed on account of their simple structures. The two nonlinear models can approximate nonlinear dynamics of many practical industrial processes applications [1–7].

In the past few years, many theoretical researchers and engineers have been performed for extensions of the Hammerstein and Wiener models to improve approximation capabilities of nonlinear systems for instance Hammerstein-Wiener system. In the existing literatures, a lot of optimization techniques have been developed to research the Hammerstein-

Wiener system [8–15]. For Hammerstein-Wiener system, the least square algorithm and blind identification method are put forward by Bai in [8, 9]. In literature [10], recursive parameter learning method are developed for a special nonlinear form described by a Hammerstein-Wiener nonlinear system including dead-zone input nonlinear function. Vörös [11] applied least square-based iterative technique to research Hammerstein-Wiener model parameters using measured input-output data. Xu et al. [12] used two extreme learning machine networks to approximate nonlinear functions of Hammerstein-Wiener system, and parameter estimation method of extreme learning machine-based Hammerstein-Wiener system is developed for large-scale complex nonlinear dynamic systems. The major drawback of the above-analyzed literatures is that the unmodeled dynamic or stochastic disturbances of the Hammerstein-Wiener process is not taken into account, which is an important factor for designing significant parameters learning algorithms [16, 17].

The stochastic gradient estimation methods have attracted much attention due to its less computational load in parameter learning. In recent years, the stochastic gradient-based algorithms have also been implemented to optimize Hammerstein-Wiener models corrupted by stochastic disturbances. For the Hammerstein-Wiener ARMAX system, Wang and Ding [18] investigated extended stochastic gradient estimation. Mansouri et al. [19] developed parameter estimation method through extended Kalman filter theory. Based on data filtering technique, a data filtering-based generalized extended stochastic gradient algorithm is derived for estimating Hammerstein-Wiener system parameters for improving computational efficiency in [20]. It is recognizable that two main problems should be considered in these proposed methods. On the one hand, methods mentioned above assume that unknown nonlinearities in systems are modeled by polynomial functions; if these nonlinearities are not polynomial functions or nonsmooth, methods mentioned do not converge [21]. On the other hand, the parameter cross-products of estimated system are included in learned system, thereby separating each block parameters from obtained parameter estimation of the cross-product terms is required, which increases computation load of learning algorithms [22].

Although many contributions in existing literatures have developed to learn nonlinear system represented by Hammerstein-Wiener model, the problem of stochastic disturbances is not fully considered. This paper focuses attention on a three-stage parameter learning approach of the Hammerstein-Wiener nonlinear systems with stochastic disturbances using multisignal data. In the first stage, the output nonlinearity are estimated depending on separable signal with different amplitudes. In phase two, correlation analysis method is implemented for estimating the linear dynamic block parameter according to one of separable signal inputs and outputs. In the third stage, in order to achieve a fast convergence rate of stochastic gradient algorithm, multi-innovation-based extended stochastic gradient scheme by expanding the scalar innovation to an innovation vector is used to learn parameters of input nonlinearity and noise model. The contributions of developed learning approach lies in:

(1) Multisignal theory is designed to employ the Hammerstein-Wiener system to separate parameter learning issues, thereby avoiding redundant parameters

(2) The unmeasurable problems of Hammerstein-Wiener system are well settled by using correlation analysis method

(3) The multi-innovation-based extended stochastic gradient scheme by expanding the scalar innovation to an innovation vector is used to achieve a fast convergence rate

The paper is organized as follows. Section 2 introduces problem statement of the neural fuzzy Hammerstein-Wiener system. Section 3 analyzes parameter learning based on multisignal data for the Hammerstein-Wiener systems with stochastic disturbances. Section 4 presents simulation cases of presented learning method. Lastly, the concluding remark is approached.

## 2. Preliminaries and Problem Statements

As described Figure 1, the nonlinear system represented by Hammerstein-Wiener model with disturbance is modeled by two neural fuzzy networks and finite impulse response model, which is formulated by

$$v(k) = f(u(k)), \tag{1}$$

$$x(k) = B(z)v(k), \tag{2}$$

$$w(k) = D(z)e(k), \tag{3}$$

$$z(k) = x(k) + w(k), \tag{4}$$

$$y(k) = g(z(k)), \tag{5}$$

where $u(k)$ and $y(k)$ denote input and output, $v(k)$ and $x(k)$ represent outputs of input nonlinearity and linear block, $e(k)$ indicates white noise sequence, $f(\cdot)$ shows input nonlinearity, $g(\cdot)$ is output nonlinearity, $B(z)$ is finite impulse response model with $B(z) = b_1 z^{-1} + \cdots + b_{n_b} z^{-n_b}$, and $D(z) = 1 + d_1 z^{-1} + \cdots + d_{n_d} z^{-n_d}$ is noise model.

For the given parameter $\varepsilon$, the establishment of the Hammerstein nonlinear system is to seek parameters satisfying the following conditions:

$$E\left(\widehat{f}(u(k)), \widehat{b}_1, \cdots, \widehat{b}_{n_b}, \widehat{d}_1, \cdots, \widehat{d}_{n_d}, \widehat{g}(\widehat{z}(k))\right)$$
$$= \frac{1}{2N} \sum_{k=1}^{N} [y(k) - \widehat{y}(k)]^2 \le \varepsilon,$$
$$\text{subject to} \qquad \widehat{v}(k) = \widehat{f}(u(k)), \tag{6}$$
$$\widehat{z}(k) = \widehat{B}(z)\widehat{v}(k) + \widehat{D}(z)e(k),$$
$$\widehat{y}(k) = \widehat{g}(\widehat{z}(k)),$$

where "∧" is estimate and $N$ represents length of measured data. From the perspective of easy analysis, the output nonlinearity is expressed by $\widehat{z}(k) = \widehat{g}^{-1}(y(k))$.

In this research, input nonlinear function and output nonlinear function are modeled using two independent neural fuzzy networks [23]. Figure 1 exhibits the neural fuzzy network, and its output is expressed as

$$\widehat{v}(k) = \widehat{f}(u(k)) = \sum_{l=1}^{L} \phi_l(u(k)) w_l, \tag{7}$$

where

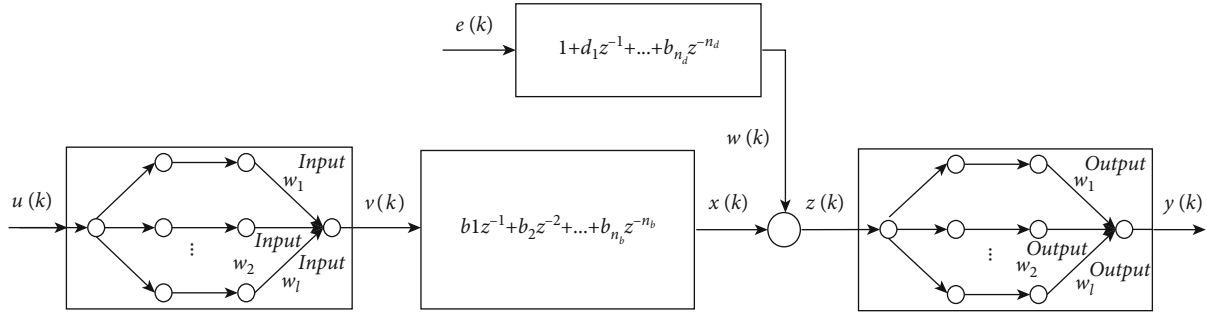$$\phi_l(u(k)) = \frac{\mu_l(u(k))}{\sum_{l=1}^{L} \mu_l(u(k))}, \tag{8}$$

FIGURE 1: Mathematical model of neural fuzzy Hammerstein-Wiener system.

where $\mu_l = \exp\left(-((u(k)-c_l)^2/\sigma_l^2)\right)$, $w_l$ represents weights of neural fuzzy network, $c_l$ and $\sigma_l$ are the center and width, and $L$ is the number of fuzzy rules.

Moreover, expressions of input-output nonlinear blocks are provided

$$\hat{v}(k) = \hat{f}(u(k)) = \sum_{l=1}^{L^{\text{input}}} \phi_l^{\text{input}}(u(k))w_l^{\text{input}}, \tag{9}$$

$$\hat{z}(k) = \hat{g}^{-1}(y(k)) = \sum_{l=1}^{L^{\text{output}}} \phi_l^{\text{output}}(y(k))w_l^{\text{output}}, \tag{10}$$

where "input" refers to input nonlinearity and "output" means output nonlinearity.

## 3. Learning Approach of Neural Fuzzy Hammerstein-Wiener System with Moving Average Noise

The tasks of parameter learning method are to estimate Hammerstein-Wiener system parameters, that is, two nonlinear blocks, linear block, and noise model. Previous research [24] pointed out that the separable signals are employed to realize separation identification of nonlinear block and linear block for the Hammerstein model. Inspired by this work, the separable signals are extended to present Hammerstein-Wiener system with unknown disturbance.

**Theorem 1.** *Considering a type of Hammerstein-Wiener system, when the separable signals are used as input signal, then the following expression maintains.*

$$R_{vu}(\tau) = b_0 R_u(\tau), \quad \forall \tau \in Z, \tag{11}$$

*where $R_u(\tau) = E(u(k)u(k-\tau))$ is the autocorrelation function, $R_{vu}(\tau) = E(v(k)u(k-\tau))$ is the cross-correlation function, and $b_0 = E(v(k)u(k))/E(u(k)u(k))$ is a constant.*

The proof can be done by referring to previous method in [23], hence omitted here.

According to Theorem 1, cross-correlation function $R_{vu}(\tau)$ is taken over by autocorrelation function $R_u(\tau)$ utilizing

separable signal. Therefore, the unknown variable $v(k)$ in Hammerstein-Wiener system is solved.

*3.1. Learning Parameters of Output Nonlinearity.* The parameters of output nonlinearity are computed using separable signals with multiple relation. In the light of description, the output nonlinearity is modeled by neural fuzzy network, thereby the centre $c_l^{\text{output}}$, the width $\sigma_l^{\text{output}}$, and the weights $w_l^{\text{output}}$ need to be estimated. The centre $c_l^{\text{output}}$ and width $\sigma_l^{\text{output}}$ are learned using previous cluster method [24]. Now, a crucial problem needs to be solved for learning parameters $w_l^{\text{output}}$.

Under the condition of two groups of separable signal with multiple relation, we can obtain following output nonlinearities:

$$z_1(k) = \hat{g}^{-1}(y_1(k)) = \sum_{l=1}^{L^{\text{output}}} \phi_l^{\text{output}}(y_1(k))w_l^{\text{output}}, \tag{12}$$

$$z_2(k) = \hat{g}^{-1}(y_2(k)) = \sum_{l=1}^{L^{\text{output}}} \phi_l^{\text{output}}(y_2(k))w_l^{\text{output}}. \tag{13}$$

From Equation (2) to Equation (4), we derive

$$z_1(k) = \sum_{j=1}^{n_b} b_j v_1(k-j) + \sum_{m=1}^{n_d} d_m e(k-m) + e(k), \tag{14}$$

$$z_2(k) = \sum_{j=1}^{n_b} b_j v_2(k-j) + \sum_{m=1}^{n_d} d_m e(k-m) + e(k). \tag{15}$$

Using $u_1(k-\tau)$ and $u_2(k-\tau)$ to multiply Equation (14) and Equation (15), respectively, the relation of correlation function are as below.

$$R_{z_1 u_1}(\tau) = \sum_{j=1}^{n_b} b_j R_{v_1 u_1}(\tau-j) + \sum_{m=1}^{n_d} d_m R_{eu_1}(\tau-m) + R_{eu_1}(\tau),$$

$$R_{z_2 u_2}(\tau) = \sum_{j=1}^{n_b} b_j R_{v_2 u_2}(\tau-j) + \sum_{m=1}^{n_d} d_m R_{eu_2}(\tau-m) + R_{eu_2}(\tau).$$

$$\tag{16}$$

It should be noted that $e(k)$ is uncorrelated with $u(k)$, thereby $R_{eu}(\tau) = 0$, and $R_{eu}(\tau - m) = 0$. Thus,

$$R_{z_1 u_1}(\tau) = \sum_{j=1}^{n_b} b_j R_{v_1 u_1}(\tau - j),$$

$$R_{z_2 u_2}(\tau) = \sum_{j=1}^{n_b} b_j R_{v_2 u_2}(\tau - j). \tag{17}$$

Using Theorem 1, we have

$$R_{z_1 u_1}(\tau) = \sum_{j=1}^{n_b} b_{01} b_j R_{u_1}(\tau - j), \tag{18}$$

$$R_{z_2 u_2}(\tau) = \sum_{j=1}^{n_b} b_{02} b_j R_{u_2}(\tau - j) = \lambda^2 \sum_{j=1}^{n_b} b_{02} b_j R_{u_1}(\tau - j), \tag{19}$$

where $b_{01} = E(v_1(k)u_1(k))/E(u_1(k)u_1(k))$ and $b_{02} = E(v_2(k)u_2(k))/E(u_2(k)u_2(k))$.

Using Equation (18) and Equation (19) acquires

$$R_{z_2 u_2}(\tau) = \beta R_{z_1 u_1}(\tau), \tag{20}$$

where $\beta = (\lambda^2 b_{01})/b_{02}$ and $\lambda = u_2/u_1$.

Equation (20) is given by

$$E(z_2(k)u_2(k - \tau)) = \beta E(z_1(k)u_1(k - \tau)). \tag{21}$$

According to Equation (12), Equation (13) and Equation (21) yields

$$\sum_{l=1}^{L^{\text{output}}} w_l^{\text{output}} E\left(\phi_l^{\text{output}}(y_2(k))u_2(k - \tau)\right)$$
$$= \beta \sum_{l=1}^{L^{\text{output}}} w_l^{\text{output}} E\left(\phi_l^{\text{output}}(y_1(k))u_1(k - \tau)\right). \tag{22}$$

Let $\phi_{l,i}^{\text{output}}(k) = \phi_l^{\text{output}} y_i(k)(l = 1, \cdots, L^{\text{output}}; i = 1, 2)$, we have

$$\sum_{l=1}^{L^{\text{output}}} w_l^{\text{output}} R_{\phi_{l,2}^{\text{output}} u_2}(\tau) = \beta \sum_{l=1}^{L^{\text{output}}} w_l^{\text{output}} R_{\phi_{l,1}^{\text{output}} u_1}(\tau). \tag{23}$$

Equation (23) is divided by $w_1^{\text{output}}$ gets

$$R_{\phi_{1,2}^{\text{output}} u_2}(\tau) - \beta R_{\phi_{1,1}^{\text{output}} u_1}(\tau)$$
$$= \frac{\beta \sum_{l=2}^{L^{\text{output}}} w_l^{\text{output}} R_{\phi_{l,1}^{\text{output}} u_1}(\tau) - \sum_{l=2}^{L^{\text{output}}} w_l^{\text{output}} R_{\phi_{l,2}^{\text{output}} u_2}(\tau)}{w_1^{\text{output}}}. \tag{24}$$

Let $\tilde{w}_l^{\text{output}} = w_l^{\text{output}}/w_1^{\text{output}}$, we have

$$R_{\phi_{1,2}^{\text{output}} u_2}(\tau) - \beta R_{\phi_{1,1}^{\text{output}} u_1}(\tau)$$
$$= \left[\sum_{l=2}^{L^{\text{output}}} \beta w_l^{\text{output}} R_{\phi_{l,1}^{\text{output}} u_1}(\tau) - R_{\phi_{l,2}^{\text{output}} u_2}(\tau)\right] \tilde{w}_l^{\text{output}}. \tag{25}$$

Assuming $\tau = 1, 2, \cdots, P(P \geq L^{\text{output}} - 1)$, and defining the following cost function according to Equation (20),

$$E\left(\tilde{w}_l^{\text{output}}\right) = \frac{1}{2} \sum_{\tau=1}^{P} \left[R_{z_2 u_2}(\tau) - \beta R_{z_1 u_1}(\tau)\right]^2. \tag{26}$$

Based on least square method, parameter $\theta$ is estimated

$$\widehat{\theta} = \left(X^{\text{T}} X\right)^{-1} X^{\text{T}} Y, \tag{27}$$

where $\widehat{\theta} = [\tilde{w}_2^{\text{output}}, \tilde{w}_3^{\text{output}}, \cdots, \tilde{w}_{L^{\text{output}}}^{\text{output}}]^{\text{T}}$ is estimation, and

$$X = [x_1, x_2, \cdots, x_{L^{\text{output}}-1}], x_{l-1} = \begin{bmatrix} \beta R_{\phi_{l,1}^{\text{output}} u_1}(1) - R_{\phi_{l,2}^{\text{output}} u_2}(1) \\ \beta R_{\phi_{l,1}^{\text{output}} u_1}(2) - R_{\phi_{l,2}^{\text{output}} u_2}(2) \\ \vdots \\ \beta R_{\phi_{l,1}^{\text{output}} u_1}(P) - R_{\phi_{l,2}^{\text{output}} u_2}(P) \end{bmatrix},$$

$$Y = \begin{bmatrix} R_{\phi_{1,2}^{\text{output}} u_2}(1) - \beta R_{\phi_{1,1}^{\text{output}} u_1}(1) \\ R_{\phi_{1,2}^{\text{output}} u_2}(2) - \beta R_{\phi_{1,1}^{\text{output}} u_1}(2) \\ \vdots \\ R_{\phi_{1,2}^{\text{output}} u_2}(P) - \beta R_{\phi_{1,1}^{\text{output}} u_1}(P) \end{bmatrix}. \tag{28}$$

The correlation functions $R_{\phi_{l,1}^{\text{output}} u_1}(\tau)$ and $R_{\phi_{l,2}^{\text{output}} u_1}(\tau)$ are given by

$$R_{\phi_{l,1}^{\text{output}} u_1}(\tau) = \frac{1}{N} \sum_{k=1}^{N} \sum_{l=2}^{L^{\text{output}}} \phi_{l,1}^{\text{output}}(y_1(k))u_1(k - \tau),$$

$$R_{\phi_{l,2}^{\text{output}} u_2}(\tau) = \frac{1}{N} \sum_{k=1}^{N} \sum_{l=2}^{L^{\text{output}}} \phi_{l,2}^{\text{output}}(y_2(k))u_2(k - \tau). \tag{29}$$

Taking the derivative of Equation (26) obtains

$$\beta = \frac{\sum_{\tau=1}^{P} R_{z_1 u_1}(\tau) R_{z_2 u_2}(\tau)}{\sum_{\tau=1}^{P} \left(R_{z_1 u_1}(\tau)\right)^2}. \tag{30}$$

*3.2. Learning Parameters of the Linear Block.* The measured input-output data of separable signal are implemented to optimize linear block relying on correlation analysis method.

Using Equation (4) gets

$$z_1(k) = \sum_{j=1}^{n_b} b_j v_1(k-j) + \sum_{m=1}^{n_d} d_m e(k-m) + e(k). \quad (31)$$

According to Equation (14)–Equation (18), we have

$$R_{z_1 u_1}(\tau) = \sum_{j=1}^{n_b} \tilde{b}_j R_{u_1}(\tau - j), \quad (32)$$

where $\tilde{b}_j = b_{01} b_j$, $b_{01} = E(v_1(k)u_1(k))/E(u_1(k)u_1(k))$.

Using Equation (32) gets

$$R = \theta_1 \psi, \quad (33)$$

where

$$\theta_1 = \left[\tilde{b}_1, \tilde{b}_2, \cdots, \tilde{b}_{n_b}\right], R = \left[R_{z_1 u_1}(1), R_{z_1 u_1}(2), \cdots, R_{z_1 u_1}(P)\right], \psi$$

$$= \begin{bmatrix} R_{u_1}(0) & R_{u_1}(1) & R_{u_1}(2) & \cdots & R_{u_1}(P-1) \\ 0 & R_{u_1}(0) & R_{u_1}(1) & \cdots & R_{u_1}(P-2) \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & 0 & \cdots & R_{u_1}(P-n_b) \end{bmatrix}. \quad (34)$$

Defining the following criterion function:

$$E(\theta_1) = \|R - \theta_1 \psi\|^2. \quad (35)$$

Taking derivative of Equation (43) obtains

$$\frac{\partial E(\theta_1)}{\partial \theta_1} = \frac{\partial \left[(R - \theta_1 \psi)^{\mathrm{T}}(R - \theta_1 \psi)\right]}{\partial \theta_1} = 2\theta_1 \psi \psi^{\mathrm{T}} - 2R\psi. \quad (36)$$

Let $\partial E(\theta_1)/\partial \theta_1 = 0$, we get

$$\theta_1 \psi \psi^{\mathrm{T}} = R\psi^{\mathrm{T}}. \quad (37)$$

Equation (37) is multiplied by $(\psi \psi^{\mathrm{T}})^{-1}$ achieves

$$\widehat{\theta}_1 = R\psi^{\mathrm{T}}(\psi \psi^{\mathrm{T}})^{-1}. \quad (38)$$

The correlation functions $R_{z_1 u_1}(\tau)$ and $R_{u_1}(\tau)$ are presented by

$$R_{z_1 u_1}(\tau) = \frac{1}{N} \sum_{k=1}^{N} \sum_{l=2}^{L^{\mathrm{output}}} \phi_{l,1}^{\mathrm{output}}(y_1(k)) u_1(k-\tau),$$
$$\quad (39)$$
$$R_{u_1}(\tau) = \frac{1}{N} \sum_{k=1}^{N} u_1(k) u_1(k-\tau).$$

### 3.3. Learning Parameters of Input Nonlinearity and Noise Model.
Based on the measured data of random signals, parameters of input nonlinearity and noise model, that is, $c_l^{\mathrm{input}}$, $\sigma_l^{\mathrm{input}}$, $w_l^{\mathrm{input}}$, and $d_m$, are learned. Parameters $c_l^{\mathrm{input}}$ and $\sigma_l^{\mathrm{input}}$ of input neural fuzzy network are learned using cluster method. Therefore, we need to learn parameters $w_l^{\mathrm{input}}$ and $d_m$.

Using Equation (1)–Equation (4) and Equation (9) gets

$$z(k) = \sum_{j=1}^{n_b} \sum_{l=1}^{L^{\mathrm{input}}} b_j \phi_l(u(k)) w_l^{\mathrm{input}} + \sum_{m=1}^{n_d} d_m e(k-m) + e(k). \quad (40)$$

For convenience, the above equation is described as below:

$$z(k) = \varphi^{\mathrm{T}}(k)\theta_2 + e(k), \quad (41)$$

where

$$\theta_2 = [\theta_s, \theta_e]^{\mathrm{T}}, \theta_s = \left[b_1 \tilde{w}_2^{\mathrm{input}}, b_1 \tilde{w}_3^{\mathrm{input}}, \cdots, b_1 \tilde{w}_{L^{\mathrm{input}}}^{\mathrm{input}}, \cdots, b_{n_b} \tilde{w}_2^{\mathrm{input}}, \cdots, b_{n_b} \tilde{w}_{L^{\mathrm{input}}}^{\mathrm{input}}\right]^{\mathrm{T}}, \quad (42)$$
$$\theta_e = [d_1, d_2, \cdots, d_{n_d}]^{\mathrm{T}}, \varphi_s(k) = [\phi_1(u(k-1)), \cdots, \phi_{L^{\mathrm{input}}}(u(k-1)), \cdots,$$
$$\phi_1(u(k-n_b)), \cdots, \phi_{L^{\mathrm{input}}}(u(k-n_b))]^{\mathrm{T}}, \varphi(k) = [\varphi_s(k), \varphi_e(k)]^{\mathrm{T}}, \varphi_e(k)$$
$$= [e(k-1), \cdots, e(k-n_d)]^{\mathrm{T}}.$$

The quadratic cost function is defined as

$$J(\theta_2) = \sum_{k=1}^{N} \|z(k) - \varphi^{\mathrm{T}}(k)\theta_2\|^2. \quad (43)$$

Based on negative search theory, minimizing $J(\theta_2)$ draws

$$\widehat{\theta}_2(k) = \widehat{\theta}_2(k-1) - \frac{1}{2r(k)} \mathrm{grad}\left[J\left(\widehat{\theta}_2(k-1)\right)\right]$$
$$= \widehat{\theta}_2(k-1) + \frac{\varphi(k)}{r(k)}\left[z(k) - \varphi^{\mathrm{T}}(k)\widehat{\theta}_2(k-1)\right], \quad (44)$$

$$r(k) = r(k-1) + \|\varphi(k)\|^2. \quad (45)$$

It is worth emphasizing that the algorithm in Equation (44) and Equation (45) is not carried out due to unknown noise terms $e(k)$ in $\varphi(k)$. In order to solve this issue, a feasible method is to use noise estimation, that is, replacing unmeasurable noise terms $e(k)$ with corresponding estimates $\widehat{e}(k)$.

The estimate $\widehat{e}(k)$ is expressed as

$$\widehat{e}(k) = z(k) - \widehat{\varphi}^{\mathrm{T}}(k)\widehat{\theta}_2(k), \quad (46)$$

where

$$\widehat{\varphi}(k) = [\varphi_s(k), \widehat{\varphi}_e(k)]^{\mathrm{T}}, \widehat{\varphi}_e(k)$$
$$= [\widehat{e}(k-1), \cdots, \widehat{e}(k-n_d)]^{\mathrm{T}}, \widehat{\theta}_2 = \left[\widehat{\theta}_s, \widehat{\theta}_e\right]^{\mathrm{T}}. \quad (47)$$
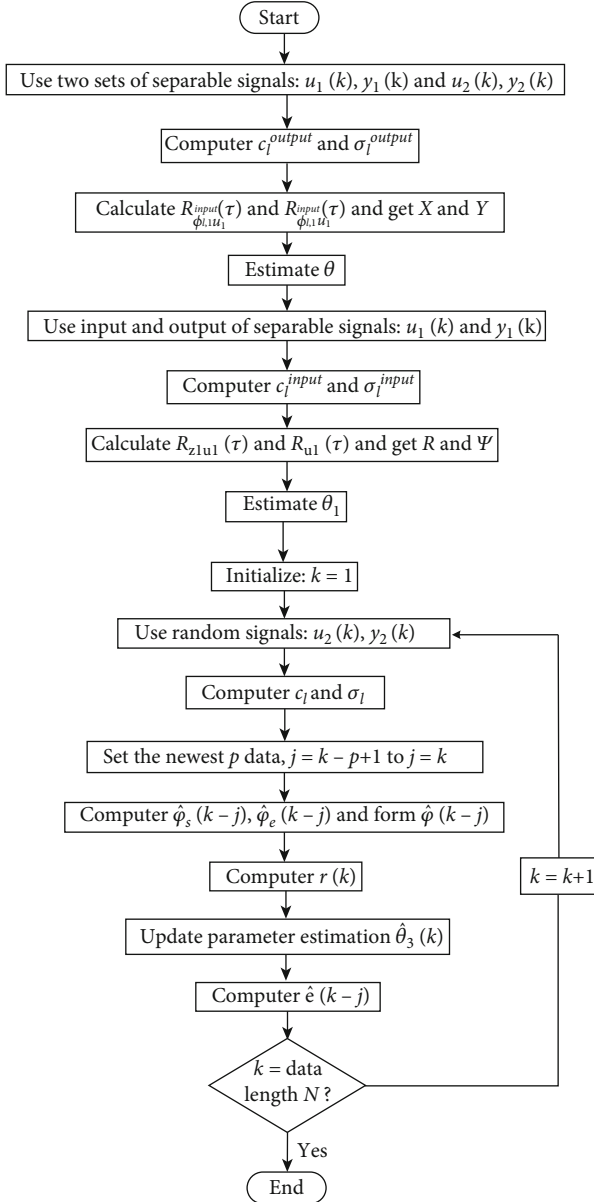
Figure 2: The flowchart of developed learning method.

As a consequence, the following algorithm is obtain:

$$
\begin{aligned}
&\widehat{\theta}_2(k) = \widehat{\theta}_2(k-1) + \frac{\widehat{\varphi}(k)}{r(k)}\left[z(k) - \widehat{\varphi}^{\mathrm{T}}(k)\widehat{\theta}_2(k-1)\right], \\
&r(k) = r(k-1) + \|\widehat{\varphi}(k)\|^2, \\
&\widehat{\varphi}(k) = [\varphi_s(k), \widehat{\varphi}_e(k)]^{\mathrm{T}}, \ \widehat{\theta}_2 = \left[\widehat{\theta}_s, \widehat{\theta}_e\right]^{\mathrm{T}}, \\
&\widehat{\varphi}_e(k) = [\widehat{e}(k-1), \cdots, \widehat{e}(k-n_d)]^{\mathrm{T}}, \\
&\widehat{e}(k) = z(k) - \widehat{\varphi}^{\mathrm{T}}(k)\widehat{\theta}_2(k).
\end{aligned}
\tag{48}
$$

As is known to all, stochastic gradient algorithm has poor convergence rate. To improve convergence rate, an effective method is to use multi-innovation learning theory by expanding the scalar innovation to an innovation vector
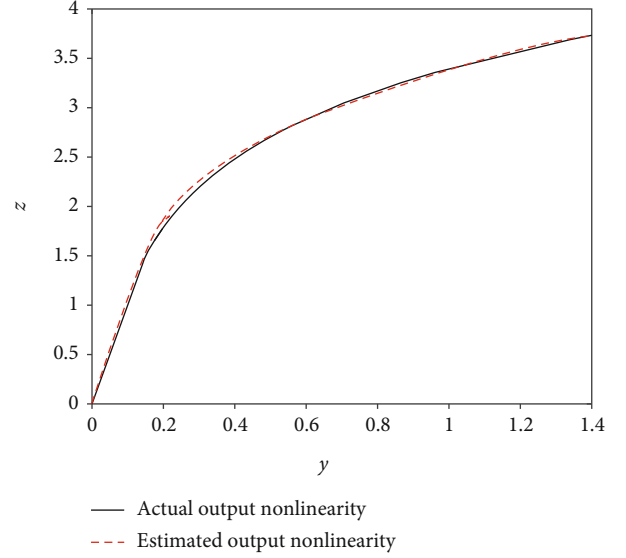


Figure 3: The estimation of the output nonlinearity.

[25], which uses not only the current data but also past data at each recursive computation.

Set the length of $p$ from $t = k - p + 1$ to $t = k$ and define cost function as below.

$$
J(\theta_2) = \sum_{t=0}^{p-1}\left\|z(k-t) - \varphi^{\mathrm{T}}(k-t)\theta_2\right\|^2.
\tag{49}
$$

Using stochastic gradient and minimizing $J(\theta_2)$ gets

$$
\widehat{\theta}_2(k) = \widehat{\theta}_2(k-1) + \frac{1}{r(k)}\sum_{t=0}^{p-1}\varphi(k-t)\left[z(k-t) - \varphi^{\mathrm{T}}(k-t)\widehat{\theta}_2(k-1)\right],
\tag{50}
$$

where $p$ is innovation length.

It is similar to extend stochastic gradient method, replacing unknown variables $\varphi(k-t)$ in Equation (50) by their estimates $\widehat{\varphi}(k-t)$, and then, the following approach combining multi-innovation theory with stochastic gradient technique is accomplished:

$$
\begin{aligned}
&\widehat{\theta}_2(k) = \widehat{\theta}_2(k-1) + \frac{1}{r(k)}\sum_{t=0}^{p-1}\widehat{\varphi}(k-t) \\
&\qquad \cdot \left[z(k-t) - \widehat{\varphi}^{\mathrm{T}}(k-t)\widehat{\theta}_2(k-1)\right], \\
&r(k) = r(k-1) + \sum_{t=0}^{p-1}\|\widehat{\varphi}(k-t)\|^2, \\
&\widehat{\theta}_2 = \left[\widehat{\theta}_s, \widehat{\theta}_e\right]^{\mathrm{T}}, \quad \widehat{\varphi}(k) = [\varphi_s(k), \widehat{\varphi}_e(k)]^{\mathrm{T}}, \\
&\widehat{\varphi}_e(k) = [\widehat{e}(k-1), \cdots, \widehat{e}(k-n_d)]^{\mathrm{T}}, \\
&\widehat{e}(k) = z(k) - \widehat{\varphi}^{\mathrm{T}}(k)\widehat{\theta}_2(k).
\end{aligned}
\tag{51}
$$

$\delta_{ns} = 12.05\%$

(a)

$\delta_{ns} = 33.73\%$

(c)

$\delta_{ns} = 22.83\%$

(b)

$\delta_{ns} = 42.19\%$
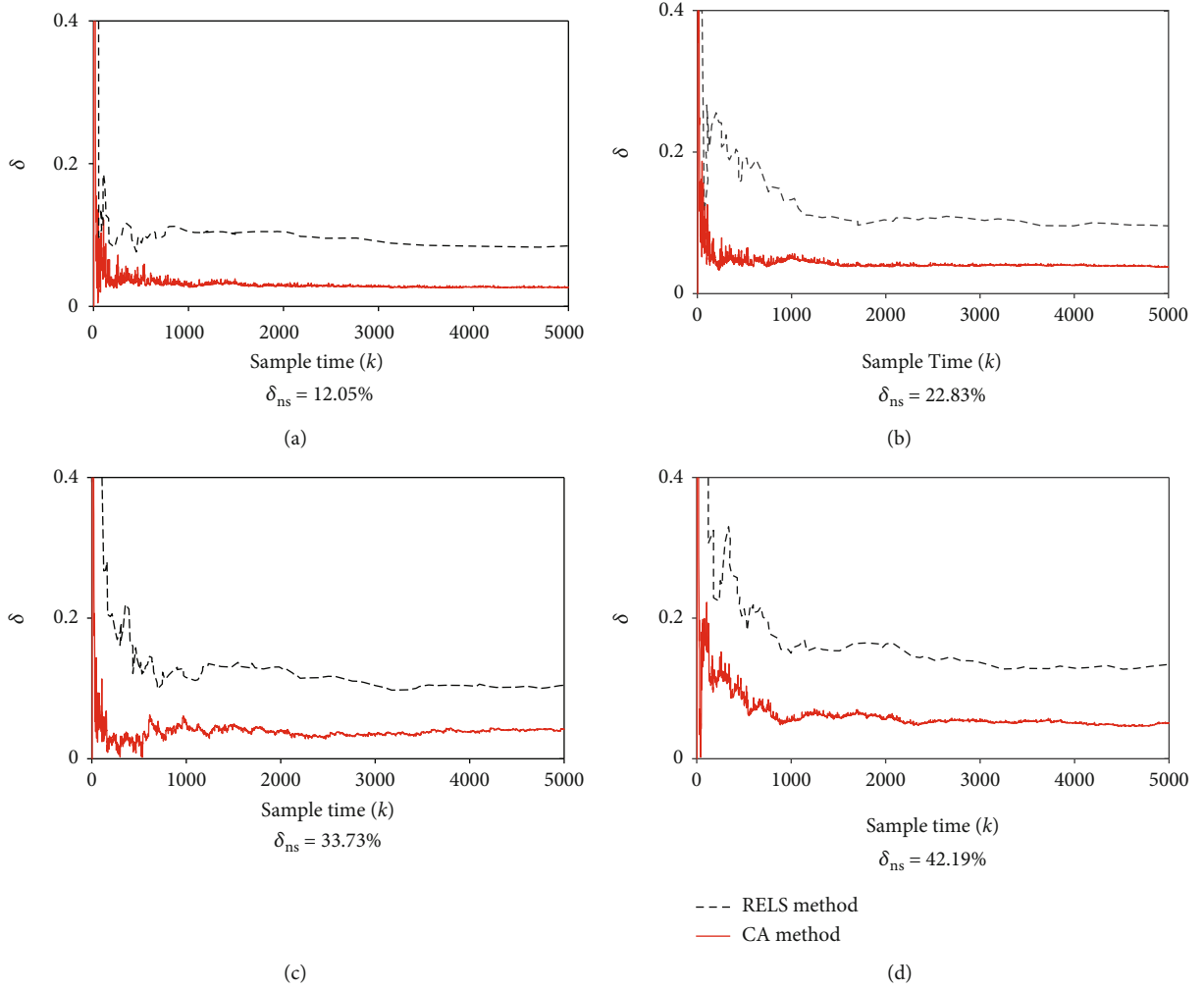
- - - RELS method
—— CA method

(d)

FIGURE 4: Error comparisons using CA method and RELS method.

From the above analysis, the flowchart of developed data-driven learning method is shown in Figure 2.

*Remark 2.* The proposed three-stage parameter learning approach estimates independently each block parameters of identified Hammerstein-Wiener system using designed multisignals, which avoids the redundant parameters of the system. In contrast, other algorithms like blind parameter identification method [9], extended stochastic gradient identification algorithm [18], and modified bias-eliminating least square algorithm [26] estimate parameters in the product term form, and they need another algorithms such as singular value decomposition method and average method to separate the hybrid parameters. Therefore, the computational complexity of these approaches increases.

## 4. Numerical Examples

For the developed learning approach, two kinds of multisignals are designed, and numerical cases of nonlinear system represented by Hammerstein-Wiener model with disturbance are applied into certificating the availability.

*4.1. Numerical Example 1.* The Hammerstein-Wiener system corrupted by noise is concerned with, where the input nonlinearity is polynomial.

$$
\begin{aligned}
v(k) &= 0.98u(k) + 0.2u(k)^2, \\
x(k) &= 0.2v(k-1) + 0.5v(k-2), \\
z(k) &= x(k) + w(k), \\
w(k) &= e(k) + 0.5e(k-1), \\
y(k) &= \begin{cases} 0.1z(k) & z(k) \le 1.5, \\ 0.15\exp\,(z(k)-1.5) & z(k) > 1.5, \end{cases}
\end{aligned}
\tag{52}
$$

where $e(k)$ is stochastic white noise.

Define the noise-to-signal ratios as $\delta_{ns} = \sqrt{\mathrm{var}\,[w(k)]/\mathrm{var}\,[x(k)]} \times 100\%$ and parameter estimation error $\delta = \|\widehat{\theta}_1(k) - \theta_1\|/\|\theta_1\|$ at sample time $k$.

The designed multisignal data consist of two sets of Gaussian signals and random signals, including Gaussian signals with mean value of 0 and variance of 1, the mean
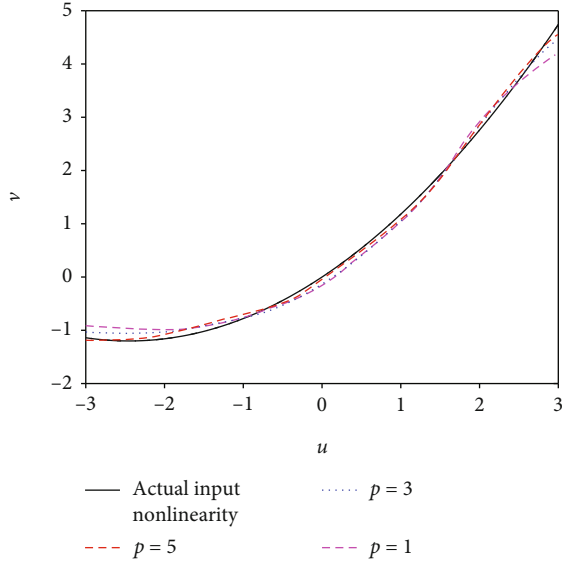
FIGURE 5: Approximation of the input nonlinearity using different innovation length.
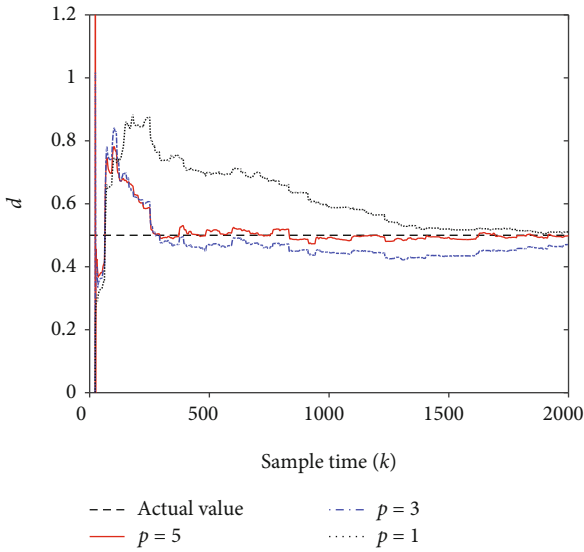


FIGURE 6: Estimate of noise model.



FIGURE 7: Estimate of output nonlinearity in example 2.

value of Gaussian signal is 0, and the variance is 0.5, and the range of random signal is -3 to 3.

To begin with, parameters of output nonlinear block are learned with the aid of collected input-output data of two sets of Gaussian signals using least square method. Set the parameters as below: $S_0 = 0.99$, $\rho = 1$, and $\lambda = 0.01$. The estimation of output nonlinearity is depicted in Figure 3. From Figure 3, neural fuzzy networks can well approximate the output nonlinearity by means of developed parameter learning approach.

Moreover, based on input-output data of Gaussian signals with variance of 1, the CA (correlation analysis) algorithm and RELS (recursive extended least square) algorithm [20] are implemented for optimizing linear block. Figure 4 shows error comparisons using CA method and RELS method of different noise-to-signal ratios. The CA algorithm
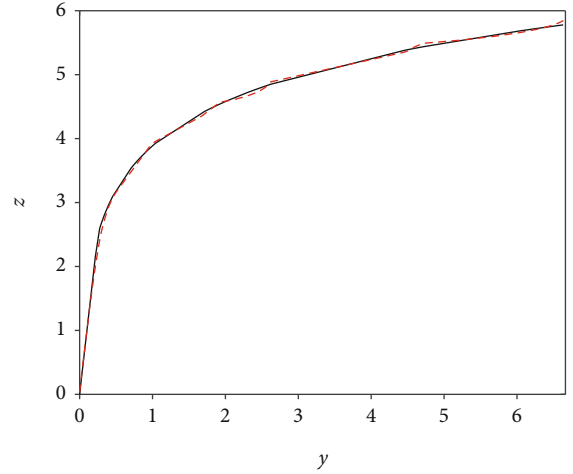
uses cross-covariance function between input and output variables and auto covariance function of input variables to learn the model parameters, which can effectively handle noise interference and improve learning accuracy. From Figure 4, with noise-to-signal ratio increases, the CA algorithm has higher precision than RELS method.

Finally, on the basis of measured input-output data of random signals, parameters of input nonlinearity and noise model are learned adopting $S_0 = 0.9$, $\lambda = 0.01$, and $\rho = 1$. Figure 5 displays the approximation of the input nonlinearity with different innovation length. Figure 6 shows estimate of noise model.

According to Figure 5, it is evident that presented learning method can effectively model input nonlinearity and obtain small approximation error with $p$ increases. According to Figure 6, with the increase of $p$, the noise model estimate is closer to real value. As a consequence, the introduction of innovation length in developed algorithm can obtain fast convergence rate. This demonstrates that presented three-stage method can accurately learn the Hammerstein-Wiener system.

*4.2. Numerical Example 2.* In view of a class of Hammerstein-Wiener system with disturbance whose input nonlinearity is discontinuous function:

$$v(k) = \begin{cases} 2 - \cos(3u(k)) - \exp(-u(k)) & u(k) \leq 3.15, \\ 3 & u(k) > 3.15, \end{cases}$$

$$x(k) = 0.9v(k-1) + 0.6v(k-2) + 0.3v(k-3) + 0.1v(k-4),$$

$$z(k) = x(k) + w(k),$$

$$w(k) = e(k) + 0.5e(k-1),$$

$$y(k) = \begin{cases} 0.25 \exp(z(k) - 2.5) & z(k) > 2.5, \\ 0.1z(k) & z(k) \leq 2.5, \end{cases}$$

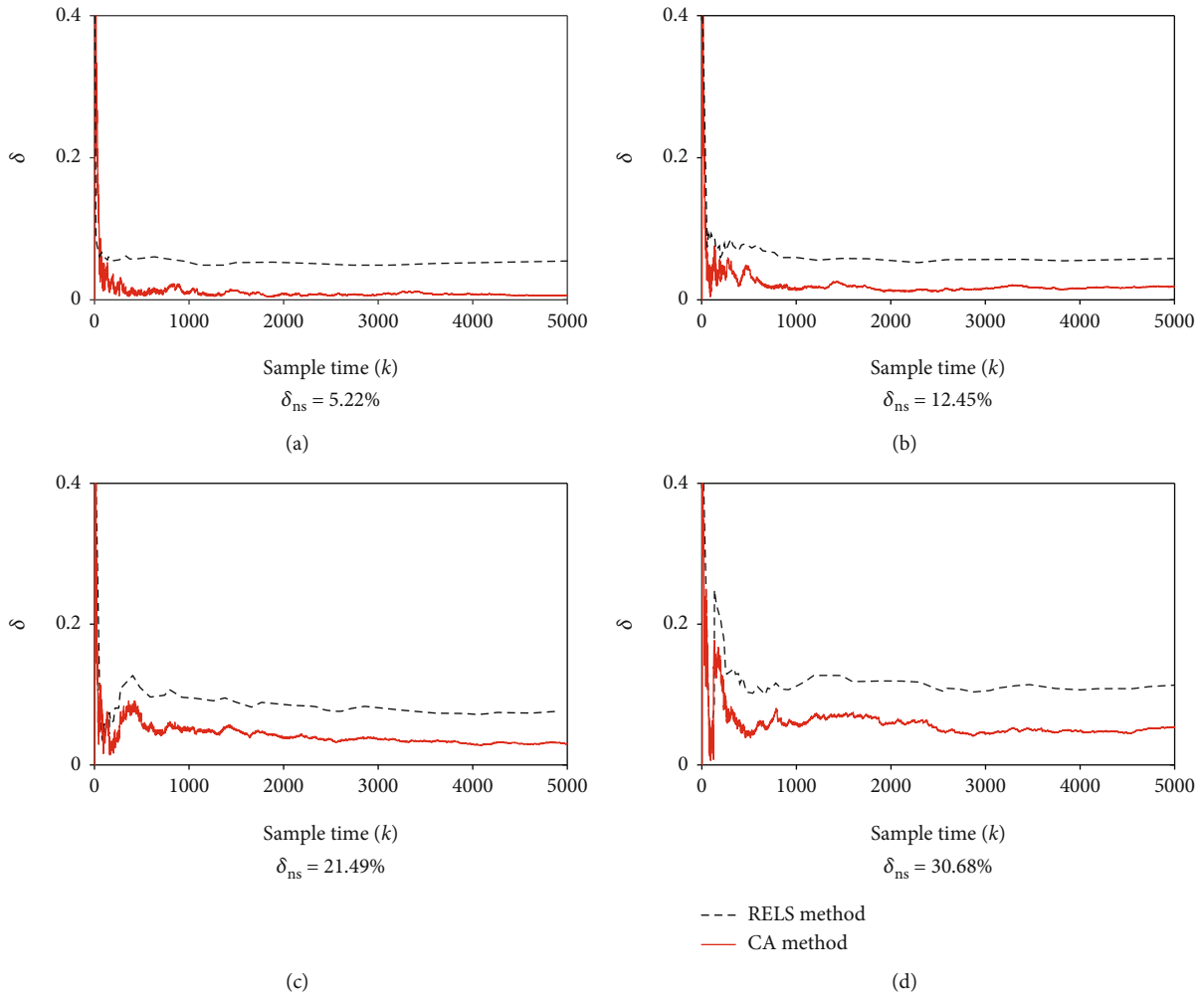$$(53)$$

where $e(k)$ is noise sequence.

FIGURE 8: Error comparisons using two methods with different noise-to-signal ratios.

The designed multisignal data consist of two sets of binary signals and random signals, including the amplitudes of the two binary signals are 2 and 4, respectively, and interval of random signal is 0 to 5.

The parameters of output nonlinearity are calculated with the aid of collected input-output data of two sets of binary signals using least square method. Set the parameters as below: $S_0 = 0.99$, $\rho = 1$, and $\lambda = 0$. The estimation of output nonlinearity is described in Figure 7. From Figure 7, the neural fuzzy networks can well approximate the output nonlinearity with the help of developed parameter learning approach.

In addition, using data of binary signals whose amplitude is 4, the CA algorithm and RELS algorithm are used. Figure 8 gives error comparisons using two methods in presence of different noise-to-signal ratios. The CA method can effectively deal with the process noise disturbance, so it achieves good parameters learning results. As can be evidently seen from Figure 8, the CA method can more effectively obtain linear block parameters and have better robustness than RELS method.

Lastly, on the basis of measured input-output data of random signals, parameters of input nonlinearity and noise model are learned adopting $S_0 = 0.92$, $\lambda = 0.01$, and

$\rho = 1$. Figure 9 displays the approximation of the input nonlinearity with different innovation length. Figure 10 lists estimate of moving average noise model for different innovation length.

Multi-innovation learning theory is combined with stochastic gradient technique to jointly improve convergence rate by expanding the scalar innovation to an innovation vector. According to Figure 9, it is recognizable that presented learning method can effectively model input nonlinearity and obtain small approximation error with $p$ increases. According to Figure 10, the noise mode estimate is closer to real value with larger innovation length.

Remark 3. For more complex Hammerstein-Wiener system with unknown disturbance in example 2, its input nonlinearity is a discontinuous function; the learning accuracy of parameter learning method proposed is reduced. In addition, it is a common knowledge that convergence rate of stochastic gradient algorithm is poor, the parameter estimation results fluctuate greatly owing to the less information in data used. With data length increases, more data information are used in parameter learning; thus, the fluctuation decreases gradually.
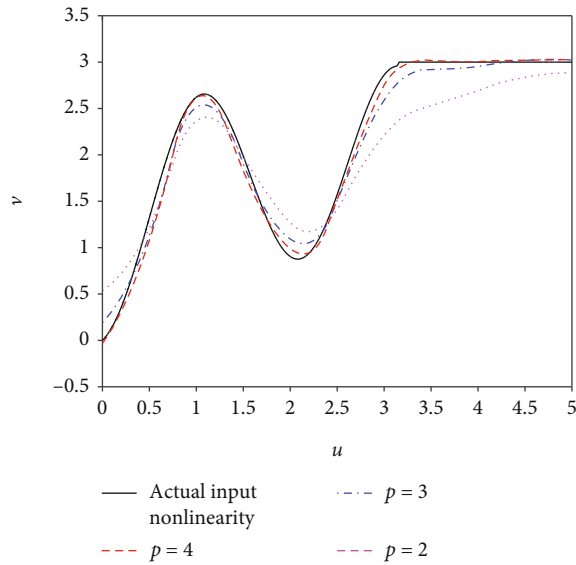
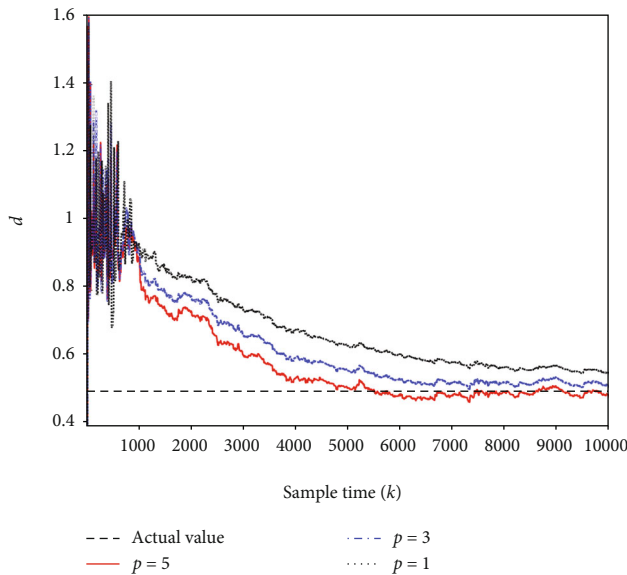FIGURE 9: Approximation of input nonlinearity with different innovation length.



FIGURE 10: Estimate noise model in example 2.

## 5. Conclusions

A novel data-driven learning approach of the nonlinear system represented by neural fuzzy Hammerstein-Wiener model with stochastic disturbances is presented. The idea of the developed method is to model the structure of two nonlinear blocks and a linear block at first and then learn the unknown parameters of each block. In the process of Hammerstein-Wiener model modeling, two nonlinear functions are approximated utilizing two neural fuzzy networks, the linear block is modeled applying impulse response model, and stochastic disturbances are described by means of moving average noise.

Multisignal theory is designed for implementing the Hammerstein-Wiener system to segregate parameter learning issues of each block, simplifying parameter learning process. Firstly, the output nonlinear block parameters are learned utilizing separable signal with different amplitudes. Secondly, the correlation analysis algorithm is used; thereby, the interference of process disturbance is effectively settled when estimating linear block. In the end, multi-innovation learning technique is combined with stochastic gradient theory to jointly learn parameters of input nonlinearity and noise model by expanding the scalar innovation to an innovation vector; the convergence rate of the system is improved. This demonstrates the availability of presented Hammerstein-Wiener system with stochastic disturbances using developed learning method.

## 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

## References

[1] F. Li, L. Jia, D. Peng, and C. Han, "Neuro-fuzzy based identification method for Hammerstein output error model with colored noise," *Neurocomputing*, vol. 244, pp. 90–101, 2017.

[2] J. C. Jeng and Y. W. Lin, "Data-driven nonlinear control design using virtual reference feedback tuning based on the block-oriented modeling of nonlinear systems," *Industrial and Engineering Chemistry Research*, vol. 57, no. 22, pp. 7583–7599, 2018.

[3] S. I. Biagiola and J. L. Figueroa, "Identification of uncertain MIMO Wiener and Hammerstein models," *Computers and Chemical Engineering*, vol. 35, no. 12, pp. 2867–2875, 2011.

[4] M. Lawrynczuk, "Nonlinear predictive control of dynamic systems represented by Wiener-Hammerstein models," *Nonlinear Dynamics*, vol. 86, no. 2, pp. 1193–1214, 2016.

[5] O. Naeem and A. E. M. Huesman, "Non-linear model approximation and reduction by new input-state Hammerstein block structure," *Computers and Chemical Engineering*, vol. 35, no. 5, pp. 758–773, 2011.

[6] F. le, I. Markovsky, C. Freeman, and E. Rogers, "Recursive identification of Hammerstein systems with application to electrically stimulated muscle," *Control Engineering Practice*, vol. 20, no. 4, pp. 386–396, 2012.

[7] J. G. Smith, S. Kamat, and K. P. Madhavan, "Modeling of PH process using wavenet based Hammerstein model," *Journal of Process Control*, vol. 17, no. 6, pp. 551–561, 2007.

[8] E. W. Bai, "An optimal two-stage identification algorithm for Hammerstein-Wiener nonlinear systems," *Automatica*, vol. 34, no. 3, pp. 333–338, 1998.

[9] E. W. Bai, "A blind approach to the Hammerstein-Wiener model identification," *Automatica*, vol. 38, no. 6, pp. 967–979, 2002.

[10] F. Yu, Z. Mao, and M. Jia, "Recursive identification for Hammerstein-Wiener systems with dead-zone input nonlinearity," *Journal of Process Control*, vol. 23, no. 8, pp. 1108–1115, 2013.

[11] J. Vörös, "Iterative identification of nonlinear dynamic systems with output backlash using three-block cascade models," *Nonlinear Dynamics*, vol. 79, no. 3, pp. 2187–2195, 2015.

[12] K. K. Xu, H. D. Yang, and C. J. Zhu, "A novel extreme learning machine-based Hammerstein-Wiener model for complex nonlinear industrial processes," *Neurocomputing*, vol. 358, pp. 246–254, 2019.

[13] F. Li, K. Yao, B. Li, and L. Jia, "A novel learning algorithm of the neuro-fuzzy based Hammerstein-Wiener model corrupted by process noise," *Journal of the Franklin Institute*, vol. 358, no. 3, pp. 2115–2137, 2021.

[14] F. Li, L. Chen, S. Wo, S. Li, and Q. Cao, "Modeling and parameter learning method for the Hammerstein-Wiener model with disturbance," *Measurement and Control*, vol. 53, no. 5-6, pp. 971–982, 2020.

[15] F. Li and L. Jia, "Parameter estimation of Hammerstein-Wiener nonlinear system with noise using special test signals," *Neurocomputing*, vol. 344, pp. 37–48, 2019.

[16] J. Wang, H. W. Wang, and H. Gu, "A novel recursive subspace identification approach of closed-loop systems," *Mathematical and Computer Modelling of Dynamical Systems*, vol. 19, no. 6, pp. 526–539, 2013.

[17] A. Mahmoudi, M. Karimi, and H. Amindavar, "Parameter estimation of autoregressive signals in presence of colored AR(1) noise as a quadratic eigenvalue problem," *Signal Processing*, vol. 92, no. 4, pp. 1151–1156, 2012.

[18] D. Wang and F. Ding, "Extended stochastic gradient identification algorithms for Hammerstein-Wiener ARMAX systems," *Computers and Mathematics with Applications*, vol. 56, no. 12, pp. 3157–3164, 2008.

[19] M. Mansouri, H. Tolouei, and M. A. Shoorehdeli, "Identification of Hammerstein-Wiener ARMAX systems using extended kalman filter," in *2011 Chinese Control and Decision Conference (CCDC)*, pp. 1110–1114, Mianyang, China, May 2011.

[20] D. Wang and F. Ding, "Recursive least squares algorithm and gradient algorithm for Hammerstein-Wiener systems using the data filtering," *Nonlinear Dynamics*, vol. 84, no. 2, pp. 1045–1053, 2016.

[21] Z. Q. Lang, "On identification of the controlled plants described by the Hammerstein system," *IEEE Transactions on Automatic Control*, vol. 39, no. 3, pp. 569–573, 1994.

[22] F. Ding and T. Chen, "Identification of Hammerstein nonlinear ARMAX systems," *Automatica*, vol. 41, no. 9, pp. 1479–1489, 2005.

[23] F. Li and L. Jia, "Correlation analysis-based error compensation recursive least-square identification method for the Hammerstein model," *Journal of Statistical Computation and Simulation*, vol. 88, no. 1, pp. 56–74, 2018.

[24] F. Li, J. Li, and D. Peng, "Identification method of neuro-fuzzy-based Hammerstein model with coloured noise," *IET Control Theory and Applications*, vol. 11, no. 17, pp. 3026–3037, 2017.

[25] X. Wang and F. Ding, "Modelling and multi-innovation parameter identification for Hammerstein nonlinear state space systems using the filtering technique," *Mathematical and Computer Modelling of Dynamical Systems*, vol. 22, no. 2, pp. 113–140, 2016.

[26] Z. Wang, Y. Wang, and Z. Ji, "A novel two-stage estimation algorithm for nonlinear Hammerstein-Wiener systems from noisy input and output data," *Journal of the Franklin Institute*, vol. 354, no. 4, pp. 1937–1944, 2017.

*Research Article*

# The Effect of Piezoelectric Fiber Rosette Configurations on Lamb Wave Direction Detection for Damage Localization

**Shuai Jiang** [1,2] **Yiping Shen** [1] **Songlai Wang** [1] **Yanfeng Peng** [1] **and Yi Liu** [3]

[1]*Hunan Provincial Key Laboratory of Health Maintenance for Mechanical Equipment, Hunan University of Science and Technology, Xiangtan 411201, China*
[2]*Hunan Railway Professional Technology College, Zhuzhou 412000, China*
[3]*Zhuzhou National Innovation Railway Technology Co., Ltd., Zhuzhou 412000, China*

Correspondence should be addressed to Yiping Shen; yiping1011@163.com

Piezoelectric fiber rosettes respond to the directivity characteristics of Lamb waves, and therefore, are useful in detecting the Lamb wave propagation direction. Considering material damage as a secondary wave source, two piezoelectric fiber rosettes are arranged to measure the scattered wave propagation directions for damage localization. The influences of various rosette configurations, i.e., $45°$-rectangular, $135°$-rectangular, $60°$-delta, and $120°$-delta, on the estimation accuracy of the propagation direction are investigated in this paper. The response of the piezoelectric fiber to the $A_0$ mode Lamb wave under narrowband tone-burst excitation is theoretically derived. Experimental tests and piezoelectric coupling simulations are performed to obtain the Lamb wave signal of each fiber. The matching pursuit (MP) algorithm is applied to extract the weak damage-related wave packet by using Hann-windowed narrowband excitation as an atom. The Lamb wave propagation directions are estimated based on the error function. The accuracies of the directions with 4 types of rosette configurations are compared, and their error sources are discussed. The results show that the accuracy of the $135°$-rectangular configuration is relatively satisfactory, and the errors depend on the size and location of each fiber in the rosette. The proposed damage localization method is validated by experimental tests. The predicted locations are close to the actual damage location. The research results are significant for piezoelectric fiber rosette design and optimization and damage location without wave speed or time-of-flight information in complex or irregular structures.

## 1. Introduction

Lamb-wave-based damage detection in plate-like structures draws increasing attention as Lamb waves can travel a long distance even in materials with low attenuation and are highly susceptible to small damage along a propagation path [1, 2]. For isotropic plates, damage can be located after detecting the scattered Lamb wave signal by at least three sensors and applying conventional time-of-flight triangulation. However, the approach requires a priori knowledge of the Lamb wave velocity in a plate to translate arrival time measurements into damage locations. This requirement is a fundamental limitation for complex or irregularly shaped structures.

An alternative damage localization technique is to apply a rosette-like directional sensor to predict the wave direction [3–6]. These directional sensors are similar with well-known electrical resistance strain gage rosette constructed of three gage grids in a certain configuration, which are generally used to resolve the principal strain directions. Thus, the wave propagation direction can be evaluated from the principal strain direction when this direction coincides with the principal strain direction for isotropic plates. Consequently, damage can be located from the point of the intersection of two wave propagation directions obtained from two directional sensors. Two directional sensor types have been proposed: the first sensor is based on fiber optics [7], and the second sensor is based on different piezoelectric elements,

e.g., macrofiber composites (MFCs) [8, 9], metal-core piezo-electric fibers [10], rectangular piezoelectric sheets [11], and round piezoelectric fibers [12]. The directivities of the piezo-electric element responses to Lamb waves have been well explored, which enables the use of rosettes for wave direction evaluation.

Kundu et al. [4] proposed a technique that performs acoustic source localization by acquiring and analyzing the signal data at several sensors in L-shaped clusters. The wave propagation direction (of the group velocity) is determined by the time difference of the arrival of waves at each sensor of a cluster. Their method was subsequently extended by Yin et al. [13, 14], and different Z-shaped clusters were introduced to decrease the number of required sensors. This method can works for an anisotropic plate despite the wave direction does not coincide with the principal strain direction [15]. However, it is difficult to extract the small differences in the time of arrival (TOA) from the weak noise signal of each sensor.

Three-element rectangular or delta rosettes are preferable in application where the principal strains are unknown [16, 17]. Rosette configurations allow us to determine the principal strain direction of the Lamb wave. The 120° and 60° delta configurations are applied in Refs. [7–10] to locate acoustic sources and 45° rectangular rosette in Refs. [11, 12]. The high directivity of those piezoelectric elements to sensing Lamb waves contributes to the principal strain direction evaluation. The delta configuration has the advantage that it enables a somewhat simpler estimation when the directivity of the rosette sensor is approximated by a cosine-squared function, as in Refs. [7, 10]. However, the presumed directivity function is not adapted to high-frequency excitation conditions. In conclusion, the rosette configuration is an important parameter for the accuracy of Lamb wave direction estimation but is still not sufficiently characterized.

The primary focus of this paper is the piezoelectric fiber rosette configurations; these rosettes are used to determine the damage location based on the measured scattered wave propagation directions. Four types of well-established rectangular and delta rosette configurations of conventional electrical strain gauges are discussed, i.e., the 45° and 135° rectangular configurations and the 60° and the 120° delta configurations. Considering the damping effect, the directivity response of the piezoelectric fiber to $A_0$ mode Lamb wave is theoretically derived, which allows the Lamb wave propagation direction to be evaluated. As the damage is small and can be treated as a secondary wave source, the damage location is determined by the intersection of the scattered wave propagation directions with two rosettes. Coupled finite element analysis and experimental tests are performed to demonstrate the accuracies of Lamb wave propagation direction estimations with various rosette configurations. The matching pursuit (MP) algorithm is applied to extract the incident and the scattered wave signal from the measured noisy Lamb wave signals by using Hann-windowed narrowband excitation as a so-called atom. Error analyses for Lamb wave direction estimations are discussed in terms of various rosette configurations.

The performance of the rosettes for damage localization is validated through artificial damage manufactured on the specimen.

## 2. Damage Localization Method with Directional Piezoelectric Fiber Rosettes

Assume that piezoelectric fiber is well bonded to the top surface of a plate with thickness of $2h$, as shown in Figure 1. The angle between the wave propagation direction $x'$ and the lengthwise direction of piezoelectric fiber is defined as $\theta$. The response voltage amplitude depends on the angle $\theta$ between the wave propagation direction $x'$ and the lengthwise direction of piezoelectric fiber. The directional response of piezoelectric fiber under narrowband tone-burst excitation has been theoretically deduced in our previous work [12]. To investigate the effect of piezoelectric fiber rosette configurations on Lamb wave direction detection, the lengths and the actual positions of piezoelectric fibers are considered. Therefore, the wave attenuation factors, including geometry spreading and material damping, are considered in this paper. The in-plate displacement of the flexural Lamb wave can be written as [8]

$$u_{x'}|_{x'=L} = Bk\sqrt{\frac{r_a}{L}}\left(\frac{\sinh az}{\cosh ah} - \frac{2ab}{k^2+b^2}\cdot\frac{\sinh bz}{\cosh bh}\right) \\ \times e^{-k_d(L-r_a)+i(kL-\omega t-\pi/2)}, \tag{1}$$

where $L$ is the distance between the Lamb wave source point and the center of the piezoelectric fiber, $B$ is an arbitrary constant, $r_a$ is the radius of the actuator, $k_d$ is the attenuation factor in the material, $h$ is half of the plate thickness, and the parameters $a$ and $b$ are defined as

$$a = \sqrt{k^2 - \frac{\omega^2}{c_L^2}}, b = \sqrt{k^2 - \frac{\omega^2}{c_T^2}}, k = \frac{2\pi}{\lambda}, \tag{2}$$

where $k$, $\lambda$, and $\omega$ are the wavenumber, wavelength, and circular central frequency, respectively, and $c_L$ and $c_T$ are the longitudinal velocity and transversal velocity, respectively, in the plate.

Similarly, the piezoelectric fiber's response to the flexural Lamb wave is expressed as [12]

$$\bar{V} = \frac{Ed_{33}\lambda}{\pi e_{33}}\sqrt{\frac{r_a}{L}}e^{-k_d L}S(\theta)\times\bar{\varepsilon}_{x'x'}, \tag{3}$$

where $\bar{\varepsilon}_{x'x'}$ is the amplitude of the in-plane strain under excitation and $S(\theta)$ is the sensitivity factor

$$S(\theta) = \cos\theta\sin\left(\frac{\pi l\cos\theta}{\lambda}\right). \tag{4}$$

Hann-windowed narrowband excitation signals are typically employed in applications to excite a Lamb wave, which is expressed as
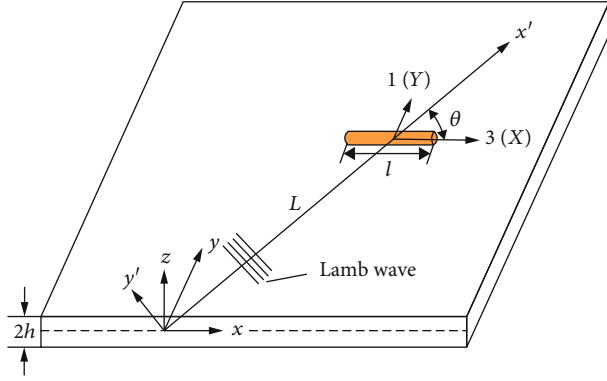
FIGURE 1: Directional response of a piezoelectric fiber to Lamb wave.

$$A(t) = A_a \left[ H(t) - H\left( t - \frac{2\pi n}{\omega_c} \right) \right] \times \sin (\omega_c t) \left( 1 - \cos \left( \frac{\omega_c t}{n} \right) \right), \quad (5)$$

where $A_a$ is the excitation amplitude, $n$ is the cycle number of the tune-burst excitation signal, and $\omega_c$ is the central circular frequency.

The response to the narrowband excitation can be expressed as [12]

$$U(t) = \int_{\omega_c - \Delta\omega/2}^{\omega_c + \Delta\omega/2} A(\omega) \bar{V}(\omega) e^{-i\omega t} d\omega$$
$$= \int_{\omega_c - \Delta\omega/2}^{\omega_c + \Delta\omega/2} \frac{E d_{33} \lambda}{\pi e_{33}} \sqrt{\frac{r_p}{L}} e^{-k_d L} S(\theta) \bar{\varepsilon}_{x'x'} A(\omega) e^{-i\omega t} d\omega, \quad (6)$$

where $\Delta\omega = 4\omega_c/n$ denotes the frequency bandwidth.

From Equation (6), the time-domain response to the narrowband excitation depends on the angle $\theta$ and the distance $L$. Considering Equation (6), the response voltage is just equal to the excitation with a time shift, a phase variation, and additional amplitude attenuation. The response voltages can be represented by their Hilbert envelope for disregarding the effect of the phase change with in the wave packet [18], and the scattered wave peak of the energy envelopes is introduced to quantify the response amplitude, which is expressed as

$$\tilde{U} = |U(t) + iH[U(t)]|_{\text{peak}}, \quad (7)$$

where $H[U(t)]$ denotes the Hilbert transform (HT) of the wave signal.

Let three piezoelectric fibers $A$, $B$, and $C$ are arranged in an arbitrary rosette configuration, as shown in Figure 2(a), and $\alpha_i (i = A, B, C)$ denote the angle between the $i^{\text{th}}$ piezoelectric fiber and the referenced piezoelectric fiber $A$. It is assumed that $\alpha_A$ is 0. According to Equation (7), the voltage response of the three piezoelectric fibers can be expressed as

$$\tilde{U}_i = \tilde{U}_{\text{max}} \cos (\theta - \alpha_i) \sin \left( \frac{\pi l \cos (\theta - \alpha_i)}{\lambda} \right) \quad i = A, B, C, \quad (8)$$

where $\tilde{U}_{i \max}$ is the maximum voltage response of the piezoelectric fiber, which is parallel to the Lamb wave propagation direction.

In this paper, the actual sum is still applied to normalize the response. The normalized amplitudes of piezoelectric fibers can be expressed as

$$T_i = \frac{3\tilde{U}_i}{\sum_{i=1}^3 \tilde{U}_i} = \frac{3 \cos (\theta + \alpha_i) \sin (\pi l \cos (\theta + \alpha_i)/\lambda)}{\sum_{i=1}^3 \cos(\theta + \alpha_i) \sin (\pi l \cos (\theta + \alpha_i)/\lambda)}. \quad (9)$$

The angle $\theta$ can be evaluated by the error between the experimental normalized voltage amplitude and the theoretical normalized voltage amplitude using the numerical computation method. The error is defined as [12]

$$e\left( \hat{\theta} \right) = \sqrt{\frac{1}{3} \sum_{i=1}^3 (T_i - T_i(\theta\wedge))^2}, \quad (10)$$

where $\hat{\theta}$ is the estimation of the Lamb wave propagation direction for theoretical calculation.

According to Equation (10), the error value $e(\hat{\theta})$ will be 0 when $\hat{\theta} = \theta$. In practical application, $e(\hat{\theta})$ is impossible to be 0 because of the unavoidable measurement error. Therefore, the estimation of $\theta$ is assumed to be $\hat{\theta}$ when $e(\hat{\theta})$ trends to the minimum. The detailed discussion is presented in our previous work [12]. Considering the damage as a secondary wave source, the damage location can be evaluated by the intersection of the scattered wave propagation directions by two rosettes, as shown in Figure 2(b). The damage location $(x, y)$ can be determined by the scattered wave propagation directions $\theta_1$ and $\theta_2$ according to

$$\begin{cases} x = \dfrac{y_2 - y_1 + x_1 \tan \theta_1 - x_2 \tan \theta_2}{\tan \theta_1 - \tan \theta_2}, \\ y = (x - x_1) \tan \theta_1 + y_1, \end{cases} \quad (11)$$

where $(x_1, y_1)$ and $(x_2, y_2)$ are the coordinates of the layout origin of the corresponding rosette.

Considering the dimension and the configuration of the piezoelectric rosette, the actual direction angle is different from the theoretical angle of the rosette. In theory, the voltage responses of the piezoelectric fibers $B$ and $C$ are referenced to the piezoelectric fiber $A$ as the base, and the distance error $\Delta L_i$ and the angle error $\Delta\theta_i$ are included in Equation (8). Figure 3 shows the distance error $\Delta L_B$ and the angle error $\Delta\theta_B$ of the piezoelectric fiber $B$. Therefore, Equation (8) is rewritten as

$$\tilde{U}_i = \tilde{U}_{i \max} \cos (\theta - \alpha_i + \Delta\theta_i) \sin \left( \frac{\pi l \cos (\theta - \alpha_i + \Delta\theta_i)}{\lambda} \right), \quad (12)$$

where $U_{i \max}$ is the maximum voltage response of the piezoelectric fiber which is located in parallel with the Lamb wave propagation direction at the distance $L + \Delta L_i$, while $U_{i \max}$

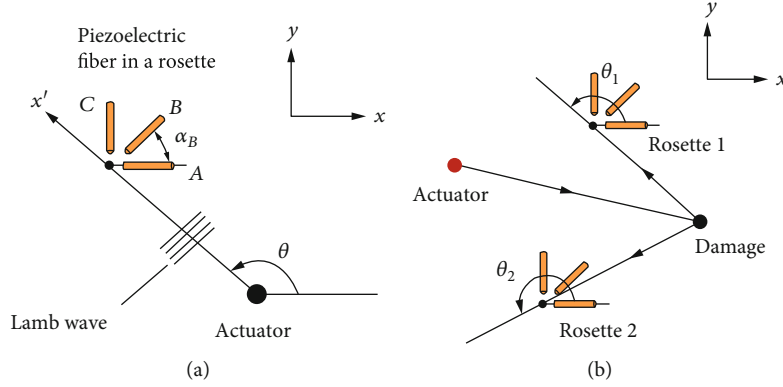(a)                                                                    (b)

FIGURE 2: Damage location with the estimated scattered wave directions. (a) The Lamb wave propagation direction. (b) Damage location with two rosettes.
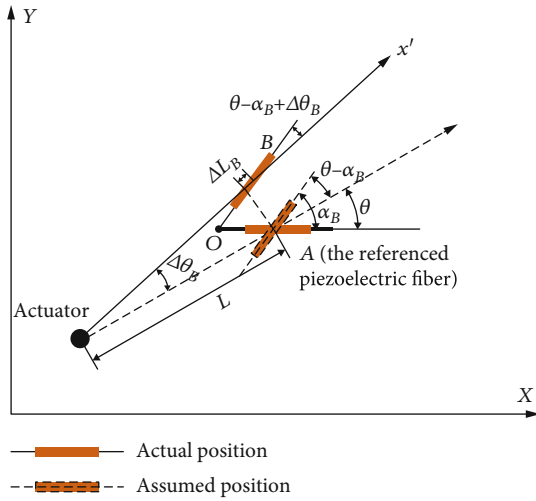


FIGURE 3: The error of the direction estimation with the rosette.

used in Equation (8) is the response of the piezoelectric fiber located at the distance $L$.

Similar to conventional electrical strain gauges, piezoelectric fiber rosettes can be arranged with different configurations. Four types of rosette configurations are discussed in this paper, as shown in Figure 4. Three piezoelectric fibers are numbered in the counterclockwise direction. $O$ is the layout origin for rosette sensor placement. Figure 4(a) shows the 45° rectangular configuration denoted by RC1 in this paper, and the configuration angles of the three piezoelectric fibers are 0°, 45°, and 90°, respectively. Figure 4(b) shows the 135° rectangular configuration denoted RC2, and the configuration angles of the three piezoelectric fibers are 0°, 135°, and 90°, respectively. Figure 4(c) shows the 60° delta configuration denoted by DC1, and the configuration angles of the three piezoelectric fibers are 0°, 120°, and 240°, respectively. Figure 4(d) shows the 120° delta configuration denoted by DC2, and the configuration angles are the same with DC1. In this paper, the piezoelectric fiber with a length of 10 mm is cut from a round piezoelectric fiber with a length of 150 mm and a diameter of 0.8 mm, produced by Smart Material Corp. The piezoelectric material is PZT SP505 (Navy type II). The electrodes at two ends of a piezoelectric fiber

are covered with silver paint. The signal wires are wired to the electrodes of the piezoelectric fibers and connected to the signal collector.

## 3. Experiment Test and Simulation Analysis

The accuracies of the Lamb wave direction estimations using the four different rosette configurations are compared by experimental tests and finite element simulations with coupled-field elements. The damage localization method is validated by experimental tests with artificial damage.

*3.1. Experimental Test Setup.* Rectangular aluminum plate specimens are employed with dimensions of 1 m × 1 m and a thickness of 1 mm for both simulation analysis and experimental testing. The density is 2730kg/m³, the elastic modulus is 68.9GPa, and Poisson's ratio is 0.33. A piezoelectric wafer is applied to excite the Lamb wave in the plate. This wafer is made from PZT8 material with a radius of 10 mm and a thickness of 0.8 mm. The piezoelectric wafer is manufactured by Smart Material Corp., USA. A 5-cycle narrowband tone-burst signal modulated by the Hamming window is employed to excite a Lamb wave in the plate. An EPA-10 power amplifier, which is produced by Piezo System Inc., USA, is applied to amplify the excitation narrowband signals. An 80 V peak-to-peak amplification excitation is applied to the actuator. An NI PXle-6361 platform is applied to collect the response output voltage of the piezoelectric sensors.

The piezoelectric fiber rosettes are directly connected to the platform to acquire the electric signals generated through the piezoelectric coupling between the strain field and the electric field. Four piezoelectric fiber rosette configurations are applied to estimate the incident Lamb wave propagation direction, as shown in Figure 5(a). The actuator is placed at an angle of 32° with four rosettes, and the distance between the actuator and the sensors is 300 mm. For the damage location experimental test, a hole with a diameter of 20 mm is manufactured in the aluminum plate. The arrangement of the actuator, rosettes, and damage is shown in Figure 4(b). The central frequencies of the excitation waves are 20 kHz, 40 kHz, and 60 kHz. The sampling frequency is 2 MHz. The measured signals of the piezoelectric sensors are averaged
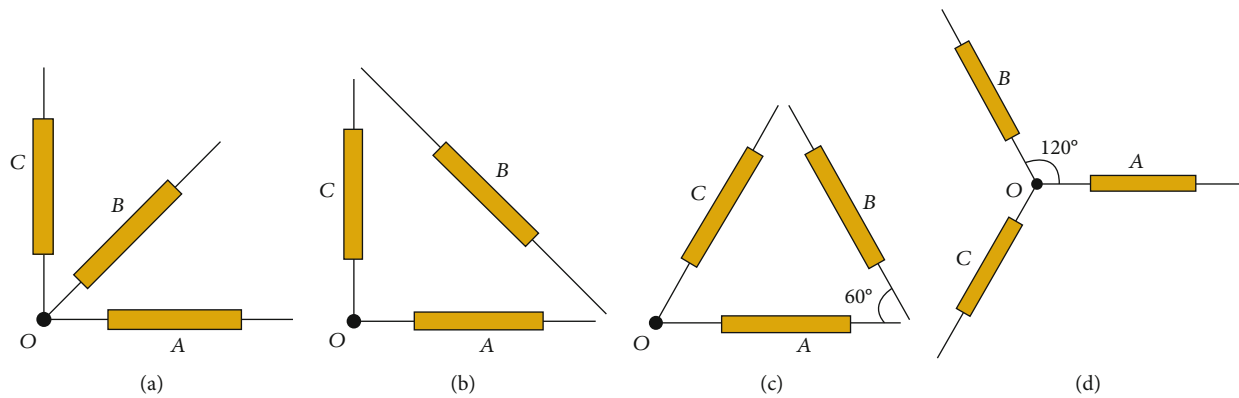
FIGURE 4: Four types of piezoelectric fiber rosette configuration: (a) rectangular (RC1); (b) rectangular (RC2); (c) delta (DC1); and (d) delta (DC2).
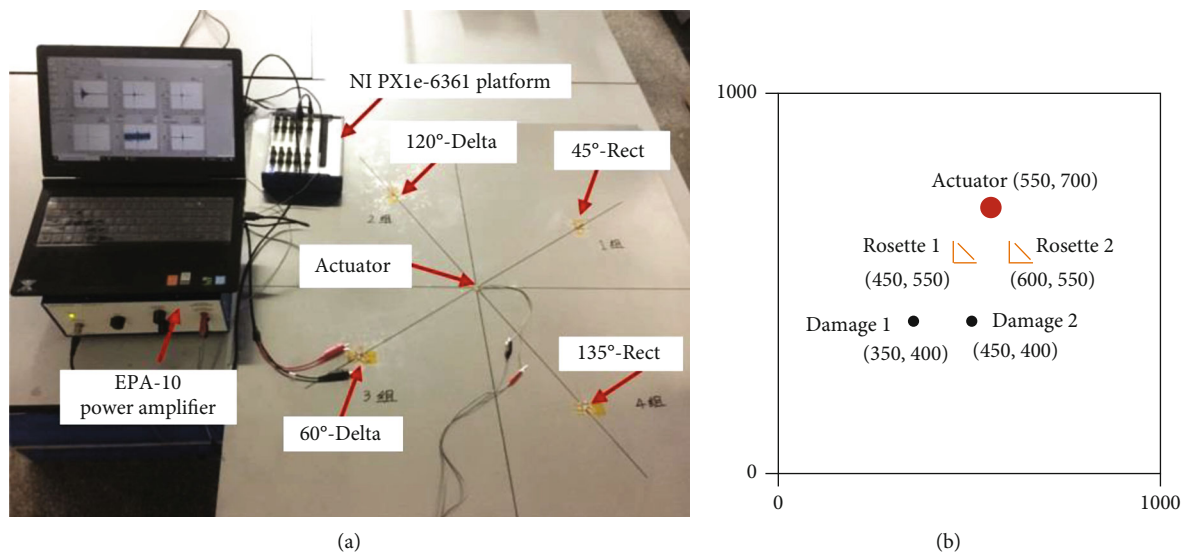


FIGURE 5: Experimental test setup. (a) Four rosette configurations. (b) Damage location.

100 times. An NI LabVIEW program was written to generate excitation waves and acquire the sensor data. A MATLAB program was written for Lamb wave signal analysis, Lamb wave propagation direction estimation, and damage location computation.

*3.2. Coupled Finite Element Analysis.* Mechanical-electrical coupled finite element analysis, which takes into account both the piezoelectric wafer actuator and the piezoelectric fibers, is performed. The corresponding material physical parameters are listed in Table 1.

The commercial finite element software ANSYS is employed for this analysis. SOLID5 coupled-field elements are applied to simulate the piezoelectric effects of the piezoelectric wafer and piezoelectric fibers. SHELL181 elements are employed to model the aluminum plate specimen. The mesh size of the finite element model is 1 mm which is smaller than one-twentieth of the Lamb wavelength at 60 kHz to ensure the accuracy of the analysis results. The time step is set to $0.5 \mu s$. Both the mesh size and the time step set satisfy the criteria of transient dynamic analysis [19]. A

total of 1,033,000 elements are employed. The voltage DOFs of the nodes located on two surfaces of the piezoelectric wafer and two ends of each piezoelectric fiber are coupled to only one master node to simulate their electrodes, as shown in Figure 6.

The excitation voltage of the narrowband tone-burst signal is applied to the upper electrode of the piezoelectric wafer. The output voltages of the three piezoelectric fibers in the rosettes are analyzed to calculate the Lamb wave propagation direction. To verify the effects of the rosette configurations on the accuracies of the direction estimation, the first arrival flexural Lamb waves are employed to extract the response amplitudes of each piezoelectric fiber. The resultant voltage outputs of the piezoelectric fiber $C$ in RC2 at a central frequency of 40 kHz are plotted in Figure 7, which is located in Rosette 1 to measure the scattered wave from damage 1. The first arrival wave of the simulation resultant wave shows agreement with the experimental measured signal, and an amplitude difference is caused by the measurement noise. Note that the obscurity of the scattered waves is attributed to their weakness. Considering the overlapped wave packets

Table 1: Material physical parameters of the piezoelectric wafer and piezoelectric fibers.

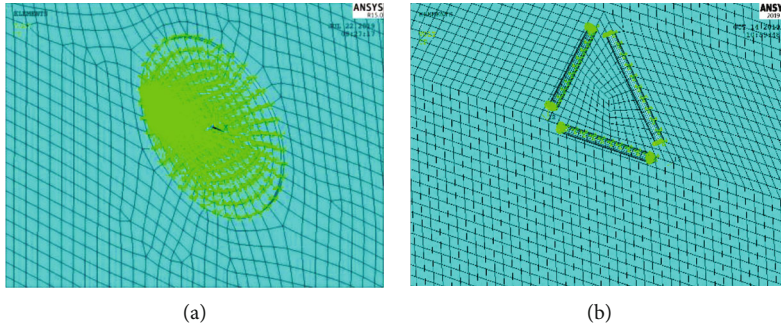| Parameter | Piezoelectric wafer | Piezoelectric fibers |
|---|---|---|
| Density (kg/m$^3$) | 7600 | 7850 |
| Relative dielectric constant | 900 | 1850 |
| Piezoelectric constant $d_{33}$ ($\times 10^{-12}$C/N) | 225 | 440 |
| Piezoelectric constant $d_{31}$ ($\times 10^{-12}$C/N) | 97 | 185 |
| Elastic compliance constant $s_{11}^E$ ($\times 10^{-12}$m/N) | 11.20 | 18.50 |
| Elastic compliance constant $s_{33}^E$ ($\times 10^{-12}$C/N) | 13.36 | 20.70 |
| Electromechanical coupling factors $k_p$ | 0.6 | 0.62 |
| Electromechanical coupling factors $k_{33}$ | 0.7 | 0.72 |



(a)　(b)

Figure 6: Coupled analysis model. (a) Local mesh of piezoelectric wafer. (b) Local mesh of RD1.
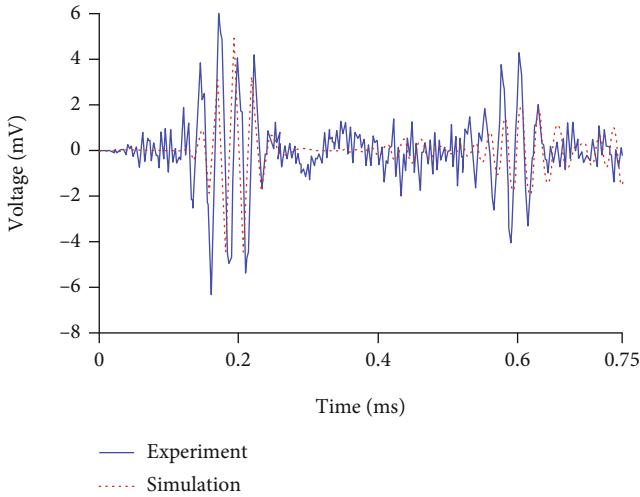


Figure 7: Resultant voltage outputs of the piezoelectric fiber $C$ in RC2 at 40 kHz.

and the measurement noise, extraction of the damage-related wave packets from the measured signal is difficult. Therefore, considering a suitable signal decomposition method is essential to effectively separate the individual wave packets and obtain the exact required wave for further scattered wave direction calculation.

### 3.3. Signal Decomposition Based on the MP Algorithm.

More advanced signal processing techniques are required to accurately separate the weak scattered wave packet from the noisy overlapped signal in an application. Sparse reconstruction has attracted a substantial amount of attention in ultrasonic guided wave-based damage detection [18, 20]. The MP algorithm introduced by Mallat and Zhifeng [21] is one of the most extensively applied algorithms for sparse signal representation. MP is an iterative greedy algorithm that computes an accurate solution for a signal in terms of the linear combinations of predefined atoms that construct an overcomplete dictionary. The algorithm of MP is described as follows.

*Step 1.* Construct a dictionary $\mathbf{D}$

$$\mathbf{D} = \left\{ g_{\gamma_1}, g_{\gamma_2}, \cdots, g_{\gamma_j}, \cdots g_{\gamma_J} \right\}, \tag{13}$$

where $\gamma_j$ is the $j^{\text{th}}$ parameter set of possible parameter combinations and $g_{\gamma_j}$ is the atom determined by the $j^{\text{th}}$ parameter set.

*Step 2.* Initialize the iteration number $\kappa = 1$ and set the measured signal $f$ to the residual $R^\kappa$.

*Step 3.* Search for the dictionary atom $g_{\gamma_\kappa}$ that best resembles the measured signal $f$, which is achieved by solving the optimization problem

$$g_{\gamma_\kappa} = \arg \max_{g_{\gamma_j} \in \mathbf{D}} \left| \left\langle R^\kappa, g_\gamma \right\rangle \right|. \tag{14}$$

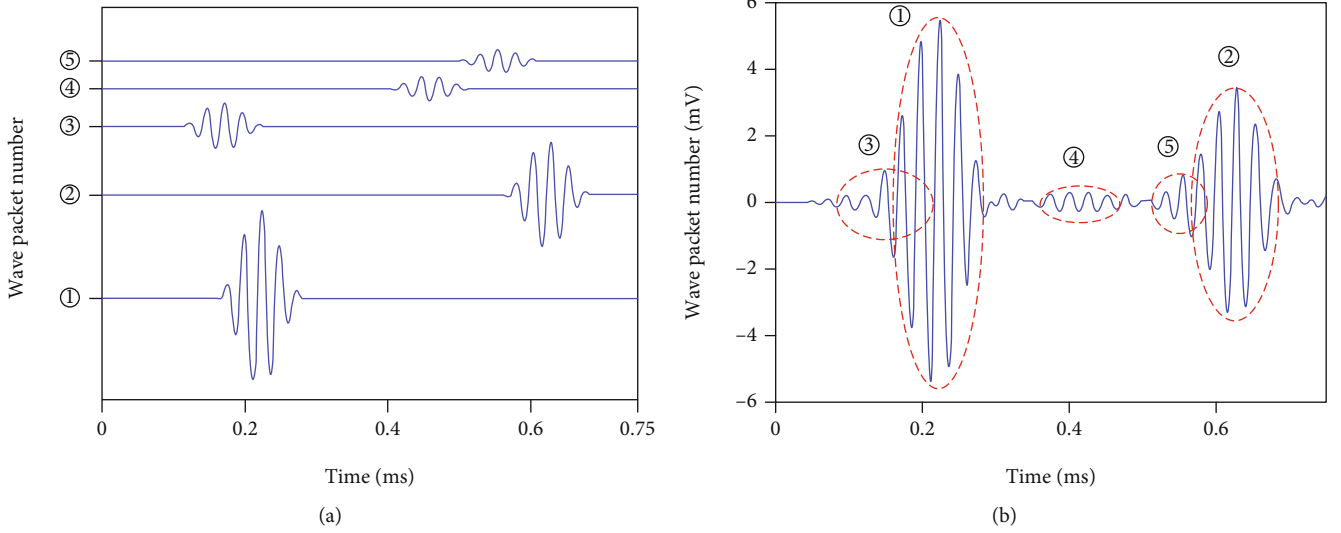*Step 4.* Calculate the amplitude of the chosen atom to the measured signal

(a)

(b)

FIGURE 8: Wave packets decomposed after five iterations using the MP algorithm. (a) MP decomposition results and (b) reconstructed signal.

$$c_\kappa = \left\langle R^\kappa, \tilde{g}_{\gamma_\kappa} \right\rangle, \tag{15}$$

where $\tilde{g}_{\gamma_\kappa} = g_{\gamma_\kappa}/\langle g_{\gamma_\kappa}, g_{\gamma_\kappa}\rangle$.

*Step 5.* Calculate a new residue by subtracting the chosen atom and adjusting the amplitude

$$R^{\kappa+1} = R^\kappa - c_\kappa g_{\gamma_\kappa}. \tag{16}$$

*Step 6.* Perform the next iteration with the residual signal $R^{\kappa+1}$, i.e., return to Step 3, until the energy of the residual signal becomes sufficiently small. After $K$ iterations, MP decomposes the signal into

$$f_{recon} = \sum_{i=1}^{K} c_\kappa g_{\gamma_k} + R^{K+1}. \tag{17}$$

To quickly and effectively decompose the measured signal, many possible atom functions are employed to construct an overcomplete dictionary, such as the Gabor atom [22], Chirplet atom [23], and an atom based on Hann-windowed narrowband excitation [18, 24, 25]. The atom dictionary should take into account the actual problem. Equations (1) and (6) indicate that the measured signal of each piezoelectric fiber is equal to the excitation with a time shift, a phase variation, and additional amplitude attenuation. Therefore, the atom in the overcomplete waveform dictionary is defined as [18]

$$g_\gamma(t) = \left[ H(t - \tau) - H\left(t - \tau - \frac{2\pi n}{\omega_c}\right) \right] \times \sin\left(\omega_c(t - \tau) + \varphi\right)$$
$$\cdot \left(1 - \cos\left(\frac{\omega_c(t - \tau)}{n}\right)\right), \tag{18}$$

where $\tau$ is the time delay and $\varphi$ is the shifted phase.

The finite set of parameters $\tau$, $\phi$, and $\omega_c$ for the dictionary should be discretized uniformly for the measured Lamb wave signal [26]. $\tau$ is discretized as $nT_s$, where $n$ and $T_s$ denote the sample length and the sampling rate, respectively. $\omega_c$ is the excitation central frequency. Many optimization algorithms can be applied to determine the parameter set, such as the genetic algorithm (GA) [27] and the artificial bee colony algorithm [28]. The GA is employed in this paper to obtain the global optimal solution in a continuous parameter space.

The damage-scattered Lamb wave signal with an excitation frequency 40 kHz, as shown in Figure 8, is decomposed after five iterations by applying the proposed MP algorithm, which is based on GA optimization. Figure 8(a) represents the individual wave packet after decomposing the measured signal, and Figure 8(b) is the reconstructed signal. Wave packets ① and ③ are the first-arrival A0 and S0 waves, respectively, ④ is the damage-scattered wave, and ② and ⑤ are bounced back from the edge. The dictionary is based on Hann-windowed narrowband excitation as an atom, which can match the individual wave packet and the weak scattered wave packet. The MP method has the advantage of excellent noise robustness.

## 4. Results and Discussion

*4.1. Lamb Wave Direction Estimation Results and Error Analysis.* Lamb wave direction estimation results of the four types of rosette configurations are listed in Table 2. Comparing the results of four rosette configurations from simulation and experiment signals, ANSYS results indicate that the estimation errors increased as the excitation frequency increased; the rosette in RC2 shows the best estimation accuracy, which is less than 4%; DC2 shows the largest error of 15%; and RC1 and DC1 have an equivalent accuracy of less than 8%. Experimental test results show larger errors than the ANSYS results. RC2 and DC1 have equivalent accuracies, which are less than 8%; RC1 presents a larger error than 10%, and DC2 shows the largest error of 30%. The differences

TABLE 2: Lamb wave direction estimation of the four types of rosette configurations.

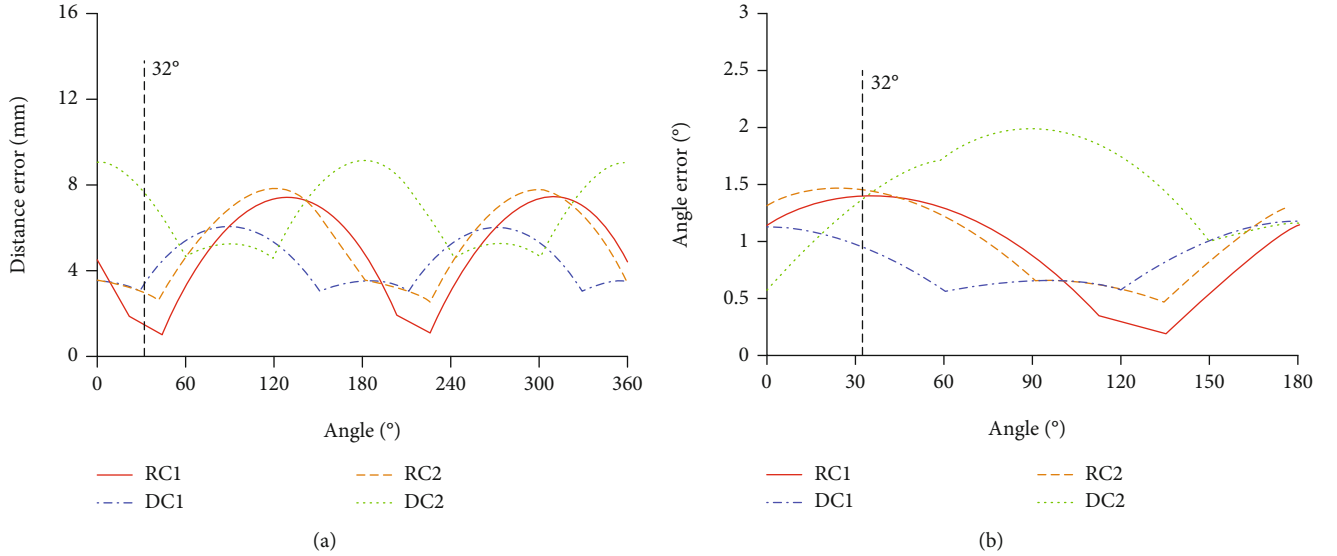| Frequency (kHz) | Results | RC1 | | RC2 | | DC1 | | DC2 | |
|---|---|---|---|---|---|---|---|---|---|
| | | /° | /%E | /° | /%E | /° | /%E | /° | /%E |
| 20 | ANSYS | 32.3 | 0.8 | 32.0 | -0.2 | 30.2 | -5.6 | 33.3 | 3.9 |
| 40 | ANSYS | 33.3 | 3.9 | 33.1 | 3.3 | 32.6 | 1.8 | 36.4 | 13.8 |
| 60 | ANSYS | 30.3 | -5.5 | 32.3 | 0.8 | 34.8 | 7.8 | 36.8 | 15.0 |
| 20 | Test | 37.8 | 18.1 | 32.4 | 1.3 | 33.7 | 5.2 | 42.4 | 32.5 |
| 40 | Test | 36.2 | 13.1 | 31.4 | -1.9 | 33.1 | 3.3 | 42.9 | 33.9 |
| 60 | Test | 34.0 | 6.3 | 30.0 | -6.3 | 34.6 | 8.1 | 41.9 | 30.8 |



(a)

(b)

FIGURE 9: The distance errors and the angle errors of various rosette configurations. (a) Distance error and (b) Angle error.

TABLE 3: The experimental results of damage localization.

| Damage number | Rosette 1 | | | Rosette 2 | | | Location | |
|---|---|---|---|---|---|---|---|---|
| | Actual (°) | Test (°) | %E | Actual (°) | Test (°) | %E | Actual (mm) | Test (mm) |
| Damage 1 | 56.4 | 59.0 | 4.6 | 31.0 | 25.2 | -18.8 | (350, 400) | (391, 451) |
| Damage 2 | 108.4 | 101.4 | -6.5 | 56.3 | 57.2 | 1.5 | (450, 400) | (486, 373) |

between the simulation results and the experimental results are caused by measurement errors and noise.

As presented in Equation (12), the attenuation caused by different distances can generate differences in the wave amplitude detected by each piezoelectric fiber in different rosette configurations, and the actual angle deviates from the ideal value. This deviation is the error in the Lamb wave direction estimation method. As the scatter wave direction is varied, the distance and angle errors also vary. To compare the accuracies of four types of rosette configurations, the average errors $(|\Delta L_B| + |\Delta L_C|)/2$ and $(|\Delta\theta_B| + |\Delta\theta_C|)/2$ of different rosette configurations are plotted in Figure 9. The distance and angle error ranges of DC2 are larger than those of the other rosette configurations; those of RC1 and RC2 are equivalent, and those of DC1 are slightly better than those of RC1 and RC2. The largest distance error in the direction of 32° is associated with DC2, which shows poor performance in estimating the Lamb wave propagation direction, as listed

in Table 2. RC2 shows better performance than the other rosette configurations due to its small distance and angle error.

*4.2. Damage Localization Results.* Considering its strong performance, RC2 is employed to perform damage localization tests. Figure 5(b) shows the arrangement of two rosettes and the actuator. The excitation central frequency is 40 kHz. Considering the damage as a secondary wave source actuator, the corresponding damage location can be estimated with the scattered signals. The related scattered wave propagation directions of the two rosettes and the predicted damage locations are listed in Table 3. The intersection point of two direction lines provides the predicted damage location. The predicted and actual damage locations are shown in Figure 10. The predicted locations are not located far from the actual locations. Note that when the damage is on the path between the actuator and the rosette, the damage-
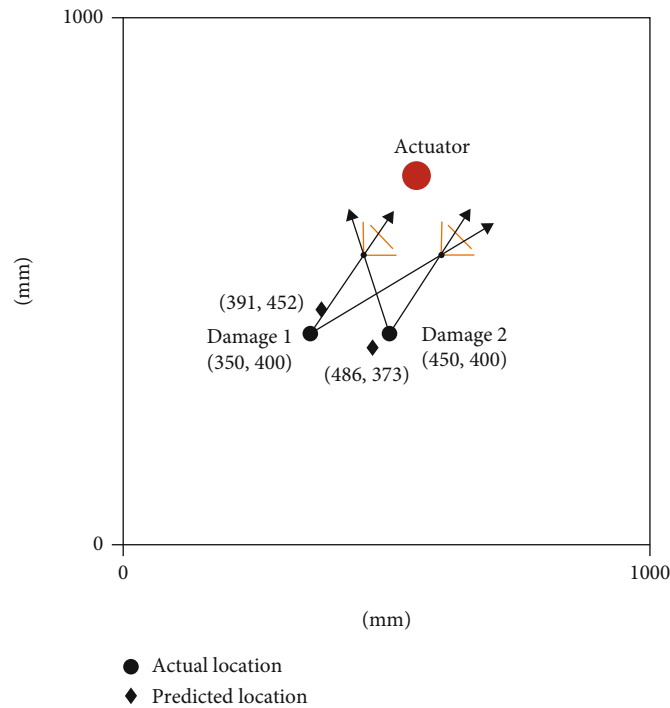
FIGURE 10: The predicted and actual damage locations.

transmitted wave packets are applied to estimate the scattered wave propagation direction. When the damage is located on or near the line between two rosettes, more rosettes are required to solve for the damage location [10]. In the future, the prediction process can be improved by using transmitted and reflected wave packets from additional rosettes for wave direction estimation.

## 5. Conclusion

This paper focuses on a damage localization method by using two piezoelectric fiber rosettes to measure the scattered Lamb wave propagation direction. The advantage of this method is that wave speed or time-of-flight information is not needed. The effects of various piezoelectric fiber rosette configurations, i.e., 45°-rectangular, 135°-rectangular, 60°-delta, and 120°-delta configurations, on the accuracies of Lamb wave propagation direction estimation are investigated. Mechanical-electric coupled finite element analyses and experimental tests are performed. The MP algorithm that is based on GA optimization by using Hann-window excitation as an atom is proposed to extract the weak damage-related wave packet. The rosette in the 135°-rectangular configuration shows satisfactory performance in determining the wave direction, but the 120°-delta configuration suffers from poor accuracy. Error analyses are performed by analyzing the distance and the angle error of each piezoelectric fiber, which deviates from the theoretical assumption. Considering damage as a secondary wave source, the damage location is determined by the intersection of two scattered wave propagation directions with two rosettes. The proposed damage localization method is validated by experimental tests, and the predicted locations are close to the actual damage locations. Future work will focus on improving the

damage localization by using transmitted and reflected wave packets from a larger number of rosettes.

## Data Availability

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

## Conflicts of Interest

The authors declare no conflict of interest.

## Authors' Contributions

S.J. contributed to the methodology, experimental test, and writing—original. Y.S. supervised this research and helped in data analysis and modification. S.W. was responsible for the data analysis and the writing—review and editing. Y.P. and Y.L. were responsible for the simulation and validation. All authors have read and agreed to the published version of the manuscript.

## Acknowledgments

# References

[1] X. Qing, W. Li, Y. Wang, and H. Sun, "Piezoelectric transducer-based structural health monitoring for aircraft applications," *Sensors*, vol. 19, no. 3, p. 545, 2019.

[2] X. Li, Z. Yang, and X. Chen, "Quantitative damage detection and sparse sensor array optimization of carbon fiber reinforced resin composite laminates for wind turbine blade structural health monitoring," *Sensors*, vol. 14, no. 4, pp. 7312–7331, 2014.

[3] T. Stepinski, M. Mańka, and A. Martowicz, "Interdigital lamb wave transducers for applications in structural health monitoring," *NDT & E International*, vol. 86, pp. 199–210, 2017.

[4] T. Kundu, H. Nakatani, and N. Takeda, "Acoustic source localization in anisotropic plates," *Ultrasonics*, vol. 52, no. 6, pp. 740–746, 2012.

[5] X. Lin, G. Chen, J. Li, F. Lu, S. Huang, and X. Cheng, "Investigation of acoustic emission source localization performance on the plate structure using piezoelectric fiber composites," *Sensors and Actuators A: Physical*, vol. 282, pp. 9–16, 2018.

[6] V. Giurgiutiu, *SHM of Aerospace Composites–Challenges and Opportunities*, CAMX Conference Proceedings, Dallas, TX, USA, 2015.

[7] D. C. Betz, G. Thursby, B. Culshaw, and W. J. Staszewski, "Lamb wave detection and source location using fiber Bragg gratin rosettes," in *Smart Structures and Materials 2003: Smart Sensor Technology and Measurement Systems*, vol. 5050, pp. 117–128, San Diego, CA, USA, July 2003.

[8] H. M. Matt and F. L. di Scalea, "Macro-fiber composite piezoelectric rosettes for acoustic source location in complex structures," *Smart Materials and Structures*, vol. 16, no. 4, pp. 1489–1499, 2007.

[9] S. Salamone, I. Bartoli, P. di Leo et al., "High-velocity impact location on aircraft panels using macro-fiber composite piezoelectric rosettes," *Journal of Intelligent Material Systems and Structures*, vol. 21, no. 9, pp. 887–896, 2010.

[10] C. Zhang, J. Qiu, H. Ji, and S. Shan, "An imaging method for impact localization using metal-core piezoelectric fiber rosettes," *Journal of Intelligent Material Systems and Structures*, vol. 26, no. 16, pp. 2205–2215, 2015.

[11] P. Zhao, D. Pisani, and C. S. Lynch, "Piezoelectric strain sensor/actuator rosettes," *Smart Materials and Structures*, vol. 20, no. 10, p. 102002, 2011.

[12] S. Wang, W. Wu, Y. Shen, H. Li, and B. Tang, "Lamb wave directional sensing with piezoelectric fiber rosette in structure health monitoring," *Shock and Vibration*, vol. 2019, Article ID 6189290, 12 pages, 2019.

[13] S. Yin, Z. Cui, and T. Kundu, "Acoustic source localization in anisotropic plates with "Z" shaped sensor clusters," *Ultrasonics*, vol. 84, pp. 34–37, 2018.

[14] N. Sen and T. Kundu, "Acoustic source localization in a highly anisotropic plate with unknown orientation of its axes of symmetry and material properties with numerical verification," *Ultrasonics*, vol. 100, article 105977, 2020.

[15] J. Zhao, J. Qiu, H. Ji, and N. Hu, "Four vectors of Lamb waves in composites: semianalysis and numerical simulation," *Journal of Intelligent Material Systems and Structures*, vol. 24, no. 16, pp. 1985–1994, 2013.

[16] V. Micro-Measurements, *Strain Gage Selection: Criteria, Procedures, Recommendations*, Technical Note. Vishay Precision Group, Inc. TN-5052007, 2007.

[17] D. A. Drake, R. W. Sullivan, and J. C. Wilson, "Distributed strain sensing from different optical fiber configurations," *Inventions*, vol. 3, no. 4, p. 67, 2018.

[18] C. Xu, Z. Yang, S. Tian, and X. Chen, "Lamb wave inspection for composite laminates using a combined method of sparse reconstruction and delay-and-sum," *Composite Structures*, vol. 223, p. 110973, 2019.

[19] Y. Shen and V. Giurgiutiu, "Combined analytical FEM approach for efficient simulation of Lamb wave damage detection," *Ultrasonics*, vol. 69, pp. 116–128, 2016.

[20] W. Wang, Y. Bao, W. Zhou, and H. Li, "Sparse representation for Lamb-wave-based damage detection using a dictionary algorithm," *Ultrasonics*, vol. 87, pp. 48–58, 2018.

[21] S. G. Mallat and Z. Zhifeng, "Matching pursuits with time-frequency dictionaries," *IEEE Transactions on Signal Processing*, vol. 41, no. 12, pp. 3397–3415, 1993.

[22] J. C. Hong, K. H. Sun, and Y. Y. Kim, "The matching pursuit approach based on the modulated Gaussian pulse for efficient guided-wave damage inspection," *Smart Materials and Structures*, vol. 14, no. 4, pp. 548–560, 2005.

[23] A. Raghavan and C. E. S. Cesnik, "Guided-wave signal processing using chirplet matching pursuits and mode correlation for structural health monitoring," *Smart Materials and Structures*, vol. 16, no. 2, pp. 355–366, 2007.

[24] H. W. Kim and F. G. Yuan, "Enhanced damage imaging of a metallic plate using matching pursuit algorithm with multiple wavepaths," *Ultrasonics*, vol. 89, pp. 84–101, 2018.

[25] Y. Xu, M. Luo, Q. Liu, G. du, and G. Song, "PZT transducer array enabled pipeline defect locating based on time-reversal method and matching pursuit de-noising," *Smart Materials and Structures*, vol. 28, no. 7, article 075019, 2019.

[26] Y. Lu and J. E. Michaels, "Numerical implementation of matching pursuit for the analysis of complex ultrasonic signals," *IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control*, vol. 55, no. 1, pp. 173–182, 2008.

[27] Y. Liu, C. Shen, Y. Wang, and F. Sun, "Guided wave NDT signal recognition with orthogonal matching pursuit based on modified evolutionary programming," *AASRI Procedia*, vol. 3, pp. 43–48, 2012.

[28] A. L. Qi, G. M. Zhang, M. Dong, H. W. Ma, and D. M. Harvey, "An artificial bee colony optimization based matching pursuit approach for ultrasonic echo estimation," *Ultrasonics*, vol. 88, pp. 1–8, 2018.