

## Research Article

# Replication-Based Data Dissemination in Connected Internet of Vehicles

**Xiying Fan** <sup>1,2</sup> **Chuanhe Huang** <sup>1,2</sup> **Junyu Zhu**<sup>3</sup> and **Bin Fu**<sup>4</sup>

<sup>1</sup>School of Computer Science, Wuhan University, China

<sup>2</sup>Collaborative Innovation Center of Geospatial Technology, China

<sup>3</sup>Research Center for Computer and Microelectronics Industry Development, MIIT (China Software Testing Center), China

<sup>4</sup>Department of Computer Science, The University of Texas Rio Grande Valley, Edinburg, USA

Correspondence should be addressed to Chuanhe Huang; [huangch@whu.edu.cn](mailto:huangch@whu.edu.cn)

Received 24 October 2018; Accepted 12 March 2019; Published 4 April 2019

Guest Editor: Zaobo He

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

Due to the dynamically changing topology of Internet of Vehicles (IoV), it is a challenging issue to achieve efficient data dissemination in IoV. This paper considers strongly connected IoV with a number of heterogenous vehicular nodes to disseminate information and studies distributed replication-based data dissemination algorithms to improve the performance of data dissemination. Accordingly, two data replication algorithms, a deterministic algorithm and a distributed randomised algorithm, are proposed. In the proposed algorithms, the number of message copies spread in the network is limited and the network will be balanced after a series of average operations among the nodes. The number of communication stages needed for network balance shows the complexity of network convergence as well as network convergence speed. It is proved that the network can achieve a balanced status after a finite number of communication stages. Meanwhile, the upper and lower bounds of the time complexity are derived when the distributed randomised algorithm is applied. Detailed mathematical results show that the network can be balanced quickly in complete graph; thus highly efficient data dissemination can be guaranteed in dense IoV. Simulation results present that the proposed randomised algorithm outperforms the present schemes in terms of transmissions and dissemination delay.

## 1. Introduction

As a promising branch of Internet of Things (IoT), Internet of Vehicles (IoV) mainly improves traffic efficiency and assists road safety through wireless communication technologies [1]. Interconnected by means of vehicle-to-vehicle (V2V) and vehicle-to-infrastructure (V2I) communications, IoV could provide data services including road safety (such as collision warning and smart traffic management), entertainment demand services (such as advertisements and online videos), and location-based services (such as interest points and path optimization). Thus IoV plays a vital role in accident warning, traffic management, and user entertainment services [2].

To enhance on-road transportation safety and efficiency, efficient data dissemination which can enable high-rate communications and rapid data dissemination, is essential for

applications in IoV. Data replication can improve dissemination performance effectively, as all the vehicles involved in data dissemination help disseminate a certain quantity of message copies. Therefore, the process of information dissemination could be expedited and the dissemination delay could be reduced [3].

Characterised by decentralised control, emerging applications in IoV are confronted with problems, such as efficient cooperation among vehicles and network consensus [4]. Adapting to dynamically changing network, a type of algorithm based on distributed averaging, gossip algorithm, attracts lots of interest [5]. Through a series of communications, the participants could have the same value or reach the common state. However, gossip-based algorithms might lead to a significant waste of network resources (network capacity, bandwidth, and computing resources)

by transmitting redundant information. Similarly, although dynamic data replication can accelerate data dissemination in distributed ad hoc networks, replication-based methods could also meet a variety of problems; for instance, the high density may result in longer communication delay, which causes network resources wasting and scalability issues. Towards data replication, Spyropoulos et al. [6] proposed to disseminate a limited number of replicas; however they did not consider available network capacity and bandwidth. RAPID [7] solved the problem by taking data utilities into account to determine how the replication should carry out. Additionally, traditional replication-based dissemination algorithms could lead to high communication overhead as well as congestions and sometimes even broadcast storm by passing around redundant information. Considering the mentioned problems, the quantity of data replicas spread in the area should be controlled.

To accelerate information dissemination, every vehicle could carry a number of data replicas. In this way, the computational burden can be distributed among the vehicles and the network load balancing can be achieved. Accordingly, a concept of network balance is proposed.

In this study, we mainly investigate data dissemination in dense IoV, which can be abstracted as complete graph by graph theory. In the situation of complete graph, we assume that every two vehicular nodes are within each other's communication range. As a tentative study, the conference paper [8] focuses on data dissemination in the context of complete graph.

Additionally, since nodes in the network can have different capabilities in terms of computation or processing due to their heterogeneity, it would be better to carry an appropriate number of data replicas according to the vehicles' own capabilities rather than an approximately equal number of data replicas as [9]. Dissemination strategies will be adjusted according to different capabilities of nodes. However, most previous work studies homogeneous vehicles in vehicular communication. Therefore, this study considers heterogeneous vehicles such that data replication strategy should be determined by the capabilities of the vehicles.

To achieve data dissemination to a target area with reduced dissemination delay and consumed resources, a deterministic algorithm and a distributed randomised algorithm based on data replication are proposed for dense vehicular scenarios. In the proposed framework, different types of vehicles have different dissemination capabilities. Each vehicular node is allocated with a corresponding value to indicate the quantity of replicas that the node can spread. Every node selects one of its neighbours to exchange data depending on the proposed algorithms, and then the pair of nodes take proportional average operations, such that the values of the vehicular nodes could be updated. By iterating the operations among the nodes, the network can reach a balanced status; that is, the network converges to a consensus. To prove the efficiency of the algorithms, we evaluate the convergence complexity by calculating the average operations needed for network balance. Detailed theoretical analysis of convergence complexity is provided.

To summarise, the current study presents the following key contributions.

(1) We consider heterogeneous vehicles with different capabilities and propose a deterministic algorithm and a distributed randomised algorithm for dense scenarios in IoV, by utilising data replication to enhance data dissemination. In the algorithms, the quantity of data copies is bounded while a network balanced status can be achieved.

(2) Theoretical analysis is presented to illustrate the number of stages needed when the network achieves a balanced status in the deterministic algorithm. The upper and lower bounds of the distributed randomised algorithm are also derived. Simulation results show the effectiveness of the distributed randomised algorithm.

The remainder of the paper is structured as below. Section 2 introduces related data dissemination schemes in vehicular networks as well as the average consensus problem. Section 3 describes the system framework. Section 4 presents a deterministic algorithm and a distributed randomised algorithm for complete graph. Section 5 gives the upper and lower bounds of the proposed randomised algorithm. Section 6 evaluates the performance of the distributed randomised algorithm. Finally, Section 7 summarises the study and Section 8 presents the future prospect.

## 2. Related Work

This section mainly introduces some information dissemination schemes in IoV. Meanwhile, related work on average consensus problem is discussed.

*2.1. Data Dissemination in Vehicular Networks.* As multiple data replicas can be forwarded to an area of interest, many works have studied replication-based data dissemination schemes and a variety of dissemination strategies have been developed [6, 7]. As a simple data dissemination scheme, while flooding has the merits of high dissemination speed and wide coverage, it could cause serious broadcast storm. Towards the problem, improvements have been made by Torres et al. [10] to adapt to various traffic scenarios.

In the routing mechanism developed by [11], the amount of data spread in the target area mainly depended on the distance from source to the base station within its communication range. Xing et al. [12] proposed a framework of utility maximisation problem for multimedia dissemination and obtained the closed form of the network utility. Wu et al. [13] aimed to fully utilise the available network capacity and presented a distributed data replication scheme. Shen et al. [14] designed a data dissemination framework to schedule data transmission with maximum dissemination utility and took advantage of the space-time network coding to improve the network efficiency. To minimise the dissemination delay to a desired number of receivers, Yan et al. [15] converted the problem to processor scheduling problem and proposed heuristics to solve the problem. Xiang et al. [16] quantified different classes of data preferences and designed a safety data dissemination protocol. Zhao et al. [17] incorporated link quality and diversity as the sender metric, based on which

an efficient selection mechanism for bulk data dissemination was proposed. Chen et al. [18] studied the relation between content replication and RSU deployment and developed a cooperative replication scheme. Given a set of tasks to be executed in vehicular clouds, Jiang et al. [19] proposed the balanced-task-assignment (BETA) policy to minimise the probability of deadline violation. The authors in [20] focused on data dissemination in IoV with social characteristic and applied the property in the design of dissemination strategies. Fan et al. [9] considered vehicles with the same capability while Lin et al. [21] studied resource allocation in vehicular cloud computing systems with heterogeneous vehicles and proposed a semi-Markov based architecture to achieve optimal resource allocation. Ghorai et al. [22] considered that the obstacles might affect radio propagation and then proposed a forwarding node selection algorithm based on fuzzy logic. Ding et al. [23] studied the cooperation in group vehicular interactions and presented a dynamic member public goods game model and a greedy based neighbour selection scheme towards the high density vehicular networks.

*2.2. Average Consensus Problem in Wireless Networks.* As the average consensus problem attracts lots of interest in research areas, such as wireless networks, many researches have been done to address the problem [24]. Boyd et al. [25] studied randomised gossip algorithms. They developed a distributed subgradient method to improve the speed of gossip algorithms and designed a framework that could be applied to analyse distribute algorithms in different scenarios. The consensus studied by Fagnani et al. [26] could be achieved at some point. Different from average preserving algorithms, this consensus point might not be the same as the average by initial states. As bidirectional communication among agents was not necessary, the studied algorithms could be applied to more settings. To describe gossip periodic sequences in an undirect graph, Yu et al. [27] used transfer function of the node. Chen et al. [28] utilised probabilistic grouping in the proposed distributed random grouping algorithm to converge to the sums. Therefore, the impact of dynamically changing topology would be alleviated. Aysal et al. [29] developed a novel gossiping algorithm for deriving the average values that could simplify the process of random gossiping and described the conditions to guarantee the network convergence. To study network consensus in strongly connected networks, Wu et al. [30] presented a gossip-based algorithm and showed it could quickly reach consensus as well as reducing the consumed transmissions. Franceschelli et al. [31] studied the execution time of heterogeneous tasks in an undirected graph. They proposed randomised interaction algorithm based on gossip to let the nodes cooperatively complete the tasks to minimise the task execution time. Nedić et al. [32] studied the characteristics of weighted-averaging dynamic for network consensus.

The mentioned literature talks about the reliability and efficiency of data dissemination as well as network consensus rather than both of the problems. Also, as few of the previous works study the heterogeneity of vehicles; this study considers a scenario that a number of vehicles with different

capabilities exist, according to which the replication strategy is determined. In summary, we aim to design replication-based dissemination schemes to facilitate data dissemination in heterogeneous vehicular networks while the network convergence rate is investigated.

### 3. System Framework

*3.1. Network Architecture.* The proposed network architecture is shown in Figure 1. As it is shown, a source vehicle carries a message and aims to disseminate the message to the area that is indicated by the circle. The message dissemination is completed by pure vehicle-to-vehicle communication. In the network, two vehicular nodes would update their own values after an average operation until the network consensus is reached.

Different from previous settings, this study considers heterogeneous vehicles with different capabilities such that the number of replicas assigned to each vehicle should be determined according to the vehicle's capability. For example, there are three types of vehicles that are classified as red vehicles, yellow vehicles, and black vehicles. Assume that red vehicle could carry 100 replicas while yellow and black ones could carry 200 and 300, respectively. Assign each type of vehicles a parameter to indicate the maximum number of replicas the vehicles can spread. When two vehicles communicate, they exchange the replicas not by simply averaging their values; instead, the average operations are taken according to the vehicles' capabilities. For example, in Figure 1, let  $n_R$  and  $n_Y$  denote the values of the red vehicle and the yellow one, respectively. The following operation will be taken when they communicate with each other,

$$\begin{aligned} n'_R &= (n_R + n_Y) \times \frac{100}{100 + 200} \\ n'_Y &= (n_R + n_Y) \times \frac{200}{100 + 200} \end{aligned} \quad (1)$$

where  $n'_R$  and  $n'_Y$  denote the new values of the red vehicle and the yellow one, respectively.

*3.2. Time Model.* In the proposed architecture, it is allowed to let pairs of independent nodes contact and exchange information in parallel. We apply the synchronous time model [25]. As in the synchronous time model, the nodes can communicate simultaneously, while it only allows one node to communicate in each time slot in the asynchronous time model.

*3.3. Assumptions and Definitions.* This part presents some related assumptions and definitions for bounded number of data dissemination.

*Assumption 1.* The vehicular network of our concern is described as an undirected graph  $G(V, E)$ . Assume that every two vehicular nodes can exchange information with each other. Consider dense IoV, such as the scenario when road congestions happen or parking lots with many parked cars.

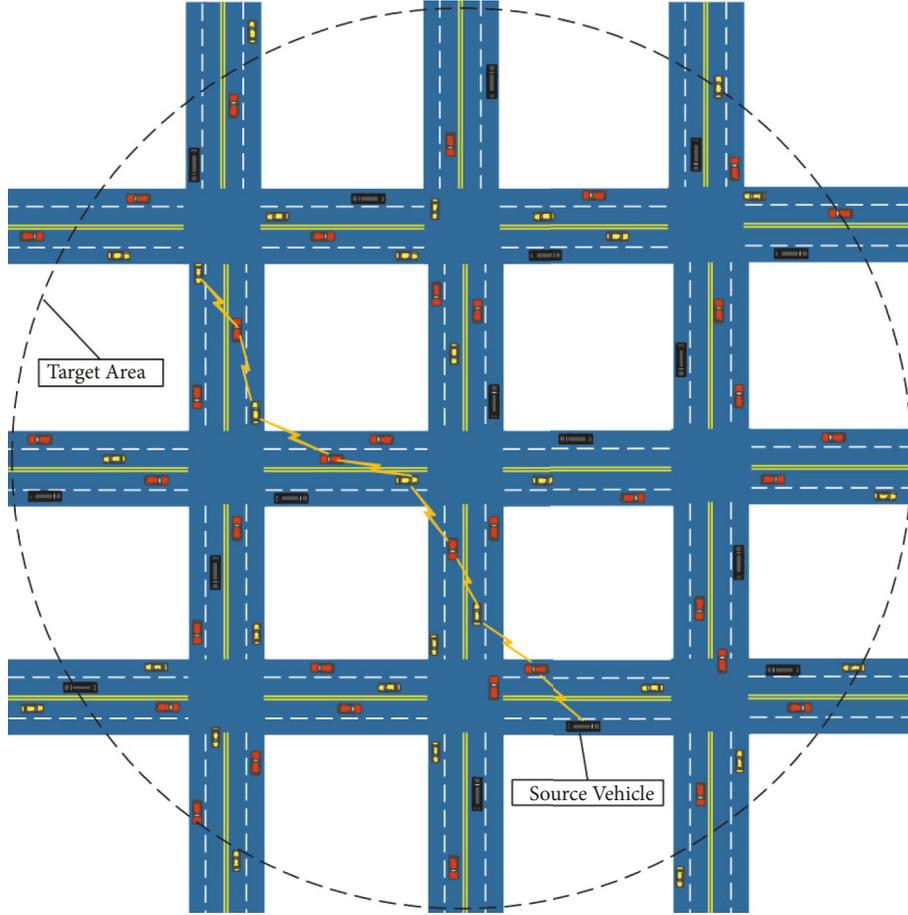


FIGURE 1: Data dissemination area.

According to graph theory, this type of network topology can be abstracted as complete graph. Every vehicle node owns a value to indicate the number of replicas it could spread. Let  $n_i$  denote the value of node  $i$ . Accordingly, graph  $G(V, E)$  becomes a weighted graph.

*Assumption 2.* As it is stated that the number of message replicas is limited to a value, we let parameter  $n$  denote the maximum quantity of replicas. Assume there are  $k$  types of vehicles in the system, e.g.,  $type_1, type_2, \dots, type_k$ . The corresponding capability of the vehicles are indicated as  $N_{type_1}, N_{type_2}, \dots, N_{type_k}$ . If the vehicle of  $type_i$  (value  $n_i$ ) and the vehicle of  $type_j$  (value  $n_j$ ) meet, the operation should be taken based on proportion,  $N_{type_i}$  or  $N_{type_j}/(N_{type_i} + N_{type_j})$ , such that the new values should be  $n'_i = (n_i + n_j) \times N_{type_i}/(N_{type_i} + N_{type_j})$ , and  $n'_j = (n_i + n_j) \times N_{type_j}/(N_{type_i} + N_{type_j})$ . In a system with homogeneous vehicles,  $n'_i = n'_j = (n_i + n_j)/2$ .

To calculate the communication stages needed to obtain network balance, the following lemmas and definitions are presented.

*Definition 3.* The nodes in the weighted graph are associated with corresponding nonnegative numbers. We say that an  $\epsilon$ -balanced status is achieved among the nodes in the graph with the following conditions met.

- (i) For any node, the number of message replicas is not smaller than 1, that is,  $n_i \geq 1$ .
- (ii) For any pair of nodes with  $n_i, n_j > 0$ ,  $|n_i - n_j| \leq \epsilon$ , where  $\epsilon > 0$ .
- (iii) If  $n_i \geq 2$  and  $n_j = 0$ , no edge should exist between the two nodes with values  $n_i$  and  $n_j$ .

**Lemma 4.** Let  $a, b$ , and  $c, d$  be real numbers satisfying the condition that  $a + b$  is equal to  $c + d$ . Then the following results can be obtained (1)  $(a^2 + b^2) - (c^2 + d^2) = 2(b - d)(b - c)$ , and (2)  $(a^2 + b^2) - (c^2 + d^2) \geq 0$  if the inequality  $a \leq c \leq d \leq b$  holds.

Lemma 4 is very easy to obtain and will be used to represent the change of potential.

*Definition 5.* Assume that  $\mathbb{R}$  denotes a list containing real numbers and  $\mathbb{N}$  denotes a list containing nonnegative integers. We present the following definitions.

- (i) A real average function  $A(.,.)$  is a mapping  $\mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R} \times \mathbb{R}$ , such that for two numbers  $a \leq b$ ,  $A(a, b) = ((a + b)/2, (a + b)/2)$  if  $a + b \geq 2$ , or  $A(a, b) = (a, b)$  if  $a + b < 2$ .
- (ii) An integer average function  $A(.,.)$  is a mapping  $\mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N} \times \mathbb{N}$  such that for two numbers  $a \leq b$ ,  $A(a, b) = (k, k)$  if  $a + b = 2k \geq 2$ ,  $A(a, b) = (k, k + 1)$  if  $a + b = 2k + 1 \geq 2$ , or  $A(a, b) = (a, b)$  if  $a + b < 2$ .
- (iii) As to list  $L : a_1, a_2, \dots, a_m$ , potential of list  $L$  is defined as  $P(L) = a_1^2 + a_2^2 + \dots + a_m^2$ .
- (iv) Assume  $A(a, b) = (c, d)$ ; we have  $S_A(\langle a, b \rangle) = 2(b - d)(b - c)$ , which could be indicating a small piece of length  $b - d$  from the bar of length  $b$  to go down by  $(b - c)$ . Function  $S_A(.)$  presents the potential change after an average operation (see Lemma 4).
- (v) Assume  $L : a_1, a_2, \dots, a_m$  is converted into a new list  $L' : a'_1, a'_2, \dots, a'_m$  after taking average operations. Let  $H$  denote all the tuples which involve in average operations and then fulfill the list transformation. Then we have  $S(H) = \sum_{(a,b) \in H} S_A(a, b) = P(L) - P(L')$  to represent the sum of product.

*Definition 6.* A communication stage represents an average operation that happens in the connected graph. Only the nodes linked by the independent edges could exchange information and take average operations. Pairs of nodes connected by different independent edges can communicate with each other at the same time.

Through iterative average operations among the nodes, we will achieve an  $\epsilon$ -balanced status. The quantity of stages needed for  $\epsilon$ -balance will be analysed to reflect the complexity of network convergence.

## 4. Algorithm Design

We abstract the strongly connected network topology as complete graph. For data dissemination in complete graph, we propose a deterministic algorithm and analyse the stages needed for network balance in Section 4.1. Then, we present a distributed randomised algorithm in Section 4.2. Upper and lower bounds of the randomised algorithm are derived through detailed theoretical analysis in Section 5.

*4.1. Deterministic Algorithm.* Here, a deterministic algorithm is proposed for complete connected graph. The algorithm is described as Algorithm 1.

In the proposed deterministic algorithm, the nodes perform operations in a deterministic manner to achieve network balance. The first step is to initialise the input graph, assume the source node has a value  $n_1$ , all other nodes have value zero. Following the parameter initialisation, the deterministic algorithm is executed in two ways, which is determined by the number of nodes with value at least two

and nodes with value zero. The operations will be iterated until the network reaches a consensus. The flowchart of the proposed deterministic algorithm is shown in Figure 2, to give a clear description of how the operations are done in a deterministic way. In the flowchart,  $t_1$  denotes the number of nodes with value at least two while  $t_2$  denotes the number of nodes with value zero. Finally, we have Theorem 7 to show the consumed communication stages.

**Theorem 7.** For complete connected graph, an algorithm exists such that after  $O(\log(n/\epsilon))$  stages of real average operations, the network can reach an  $\epsilon$ -balanced status. The algorithm is shown as Algorithm 1.

*Proof.* It is easy to see that there are at most  $O(\log n)$  stages for steps (6) - (10). The upper bound for steps (11) - (13) is given below.

Let  $n_i$  be the parameter value of node  $i$ . In general, assume that  $n_1 \geq n_2 \geq \dots \geq n_m$ . Let  $n_i$  and  $n_{m-i+1}$  take average.

Assume that after one stage, pair  $n_i$  and  $n_{m-i+1}$  has the largest average  $d_i = (n_i + n_{m-i+1})/2$ , and pair  $n_j$  and  $n_{m-j+1}$  has the least average  $d_j = (n_j + n_{m-j+1})/2$ . We assume that  $i \neq j$ .

It is noted that either  $(n_i - n_j)/2 \leq 0$  or  $(n_{m-i+1} - n_{m-j+1})/2 \leq 0$ .

$$\begin{aligned}
 d_i - d_j &= \frac{n_i + n_{m-i+1}}{2} - \frac{n_j + n_{m-j+1}}{2} \\
 &= \frac{(n_i - n_j) + (n_{m-i+1} - n_{m-j+1})}{2} \\
 &\leq \max \left( \left| \frac{n_i - n_j}{2} \right|, \left| \frac{(n_{m-i+1} - n_{m-j+1})}{2} \right| \right) \\
 &\leq \frac{n_1 - n_m}{2}.
 \end{aligned} \tag{2}$$

After  $t$  stages, the difference of the nodes is at most  $(n_1 - n_m)/2^t$ , such that after  $O(\log(n/\epsilon))$  stages, an  $\epsilon$ -balanced status can be achieved.  $\square$

*4.2. Distributed Randomised Algorithm.* The following part develops a distributed randomised algorithm (DRA), shown as Algorithm 2.

The flowchart of the proposed deterministic algorithm is shown in Figure 3, to give a clear description of how the operations are done in a random manner.

During the process of initialisation, related parameters as well as the settings (including the values of nodes, the maximum number of replicas, and number of vehicles of different types) in the network are set. Then, we have to determine how the replication strategy is carried out. For each node in sending status, it randomly selects a neighbour node to send the communication request. The receiving node would choose a node with the largest gap to take specific average operations according to the capabilities of vehicles. The algorithm would execute an iterative procedure until the network reaches consensus. Finally, the algorithm returns

**Input:** Graph  $G$ .  
**Output:** Graph  $G'$ , communication stages  $a$ .

- (1) Initialisation;
- (2)  $a = 1$ ;
- (3)  $k$  nodes, values with  $n_1, n_2, \dots, n_k$ ;
- (4) Stage  $a$  (step (5) - step (13)):
- (5) sort  $n_1, n_2, \dots, n_k$  in descending order;
- (6) **if** the nodes with value no smaller than two are more than the ones with value zero **then**
- (7)     let the latter ones take operations with the former ones;
- (8) **else**
- (9)     take operations the other way around.
- (10) **end if**
- (11) **if** there is no node with value at least two or with value zero **then**
- (12)     Let  $n_i$  and  $n_{m-i+1}$  take average for  $i = 1, 2, \dots, \lfloor m/2 \rfloor$ ;
- (13) **end if**
- (14) Enter into the next stage,  $a = a + 1$ ;
- (15) End of Algorithm.

ALGORITHM 1: Deterministic Algorithm.

**Input:** graph  $G$ ;  
**Output:** graph  $G'$ , parameter  $a$ .

- (1) Initialisation;
- (2) Let  $a = 0$ ;
- (3) Let  $i$  indicate a vehicular node;
- (4) Let  $n_i$  indicate the distribution task of node  $i$ ;
- (5) Let  $n'_i$  indicate the new distribution task of node  $i$ ;
- (6) **repeat**
- (7)      $a = a + 1$ ;
- (8)     Each vehicular node at sending status randomly chooses a vehicular node within its neighbourhood, then sends the communication request;
- (9)     Each vehicular node at receiving status selects the node from the received request if the gap between the two nodes is the largest;
- (10)     Take pairwise proportional average operations for the corresponding pairs of nodes;
- (11)     Let  $i$  and  $j$  indicate the vehicles who take average with each other;
- (12)     New values of node  $i$  and  $j$  are updated as  $n'_i = (n_i + n_j) \times N_{type_i} / (N_{type_i} + N_{type_j})$ , and  $n'_j = (n_i + n_j) \times N_{type_j} / (N_{type_i} + N_{type_j})$ , respectively;
- (13) **until**  $(|n_i - n_j| \leq \epsilon)$
- (14) End of Algorithm.

ALGORITHM 2: Distributed randomised algorithm.

graph  $G'$  and the number of average stages when the network reaches the  $\epsilon$ -balanced status.

To analyse the effectiveness of the randomised algorithm, we derive several important theorems based the famous Chernoff bounds [33]. The new theoretical results are shown as Theorems 8 and 9 and Corollary 10. The proof makes our entire study self-contained.

**Theorem 8** (see [8]). *Let  $X_1, \dots, X_n$  be  $n$  independent random 0-1 variables, where  $X_i$  takes 1 with probability at least  $p$  for  $i = 1, \dots, n$ . Let  $X = \sum_{i=1}^n X_i$ , and  $\mu = E[X]$ . Then for any  $\delta > 0$ ,  $\Pr(X < (1 - \delta)\mu) < e^{-(1/2)\delta^2 \mu}$ .*

**Theorem 9** (see [8]). *Let  $X_1, \dots, X_n$  be  $n$  independent random 0-1 variables, where  $X_i$  takes 1 with probability at most  $p$  for*

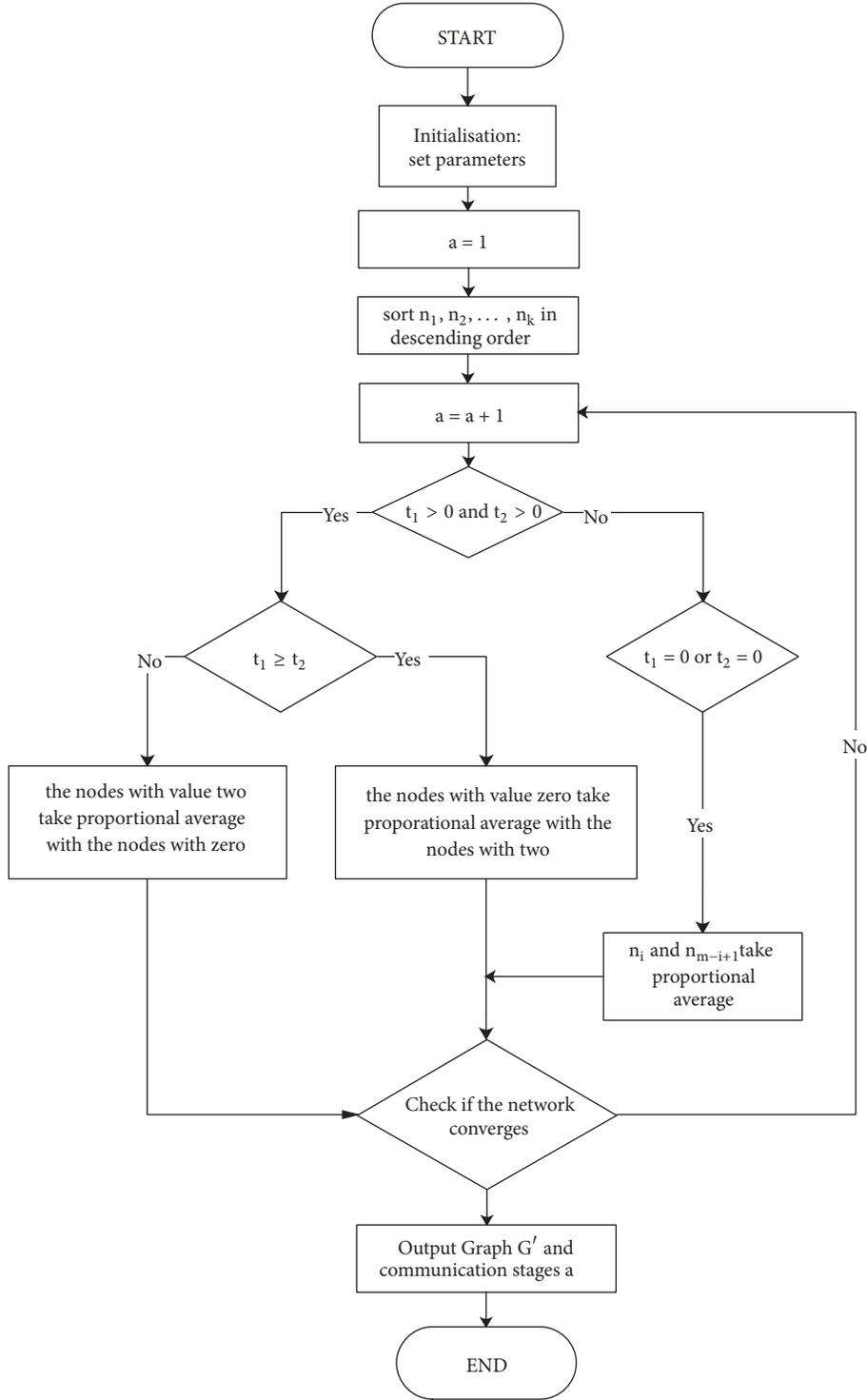


FIGURE 2: Flowchart of proposed deterministic algorithm.

$i = 1, \dots, n$ . Let  $X = \sum_{i=1}^n X_i$ . Then for any  $\delta > 0$ ,  $\Pr(X > (1 + \delta)pn) < [e^\delta / (1 + \delta)^{(1+\delta)}]^{pn}$ .

**Corollary 10** (see [34]). Let  $X_1, \dots, X_n$  be  $n$  independent random 0-1 with  $X$  being the sum of  $X_i, i = 1, \dots, n$ .

(1) If  $X_i$  takes 1 with probability at most  $p$  for  $i = 1, \dots, n$ , then for any  $1/3 > \epsilon > 0$ ,  $\Pr(X > pn + \epsilon n) < e^{-(1/3)\epsilon n^2}$ .

(2) If  $X_i$  takes 1 with probability at least  $p$  for  $i = 1, \dots, n$ , then for any  $\epsilon > 0$ ,  $\Pr(X < pn - \epsilon n) < e^{-(1/2)\epsilon n^2}$ .

## 5. Performance Analysis

In each time step, every node becomes active with probability  $1/2$  independently. Consider a node  $i$  that is active; let  $d(i)$

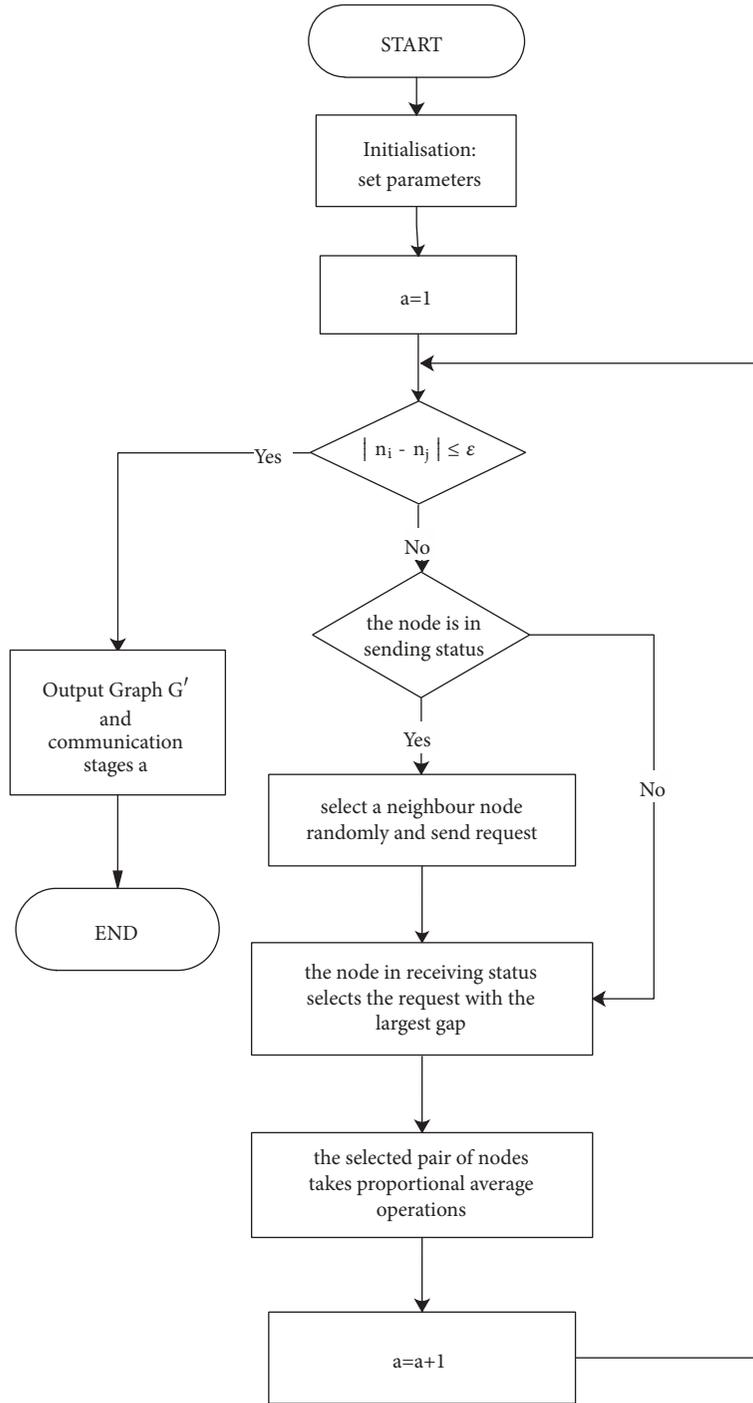


FIGURE 3: Flowchart of proposed randomised algorithm.

be its degree; that is, the number of its neighbours. Node  $i$  selects at most one of its neighbours to contact and take the proportional average operation. Each neighbour has an equal probability to be selected, i.e.,  $1/d(i)$ . The active nodes may receive more than one contact request and they would select one of the contact requests with the largest gap. To represent the averaging time of the randomised algorithm, we have the following theorem referring to Boyd et al. [25]. Here, the averaging time means the smallest time it takes a value to be  $\epsilon$ -close to the average value in the system.

**Theorem 11** (see [25]). *The averaging time of the distributed randomised algorithm described above is given as*

$$\frac{0.5 \log(1/\epsilon)}{\log(1/\lambda)} \leq T(\epsilon) \leq \frac{3 \log(1/\epsilon)}{\log(1/\lambda)}, \quad (3)$$

where  $\lambda = (1/2)(1 + \lambda_2(P))$ .

In the following part, an upper bound and a lower bound of the communication stages consumed for network balance are derived for the proposed randomised algorithm.

**5.1. Upper Bound.** Before we present the upper bound of the proposed randomised algorithm, the following concepts need to be clarified.

**Definition 12.** For two integers  $a$  and  $b$ ; the average operation generates two new integers  $(a', b')$ , with the value of  $a'$  being equal to  $\lfloor (a + b)/2 \rfloor$ , and  $b'$  is equal to the remaining value.

**Definition 13.** Let  $L = a_1, \dots, a_k$  denote a set of real numbers.  $gap(L)$  is defined as  $\max_{1 \leq i, j \leq k} |a_i - a_j|$ .

**Definition 14.** Let  $\alpha > 0$ , and  $K = a_1, \dots, a_k$  denote a list of real numbers. Assume that  $K$  is transformed into  $K' = a'_1, \dots, a'_k$  after a series of communication stages. We regard  $K'$  as an  $\alpha$ -shrink compared with  $K$  when  $gap(K')$  is within a factor  $(1 - \alpha)$  of  $gap(K)$ .

Lemma 15 is derived to show how the gap of a list of numbers shrinks via the specific average operations.

**Lemma 15** (see [8]). *Let  $r(\cdot)$  be a function from  $S \rightarrow S$  that  $r(x)$  generates a random element in  $S$ . Assume  $A$  and  $B$  are two subsets of  $S$  satisfying  $|A| \leq |B|$ , and  $R(A) = \{x : x \in A, r(x) \in B\}$ ,  $H(A) = \{r(x) : x \in A, r(x) \in B\}$ . Then with a probability at most*

$$g(\epsilon)^{|A||B|/|S|} + ((1 - \gamma))^{(2\gamma-1)(1-\epsilon) \cdot |B|/|S| \cdot |A|}, \quad (4)$$

we have

$$|H(A)| \leq (1 - \gamma)(1 - \epsilon) \cdot \frac{|B|}{|S|} \cdot |A|, \quad (5)$$

where  $\gamma$  is a constant in  $(0, 1)$ . Furthermore, if  $|B| \geq \delta|S|$  for some fixed  $\delta \in (0, 1)$  then the failure probability is at most  $2(1 - a)^{|A|}$  for some fixed  $a \in (0, 1)$ .

*Proof.* Let  $m$  denote the number of elements in  $R(A)$  and let  $n$  denote the number of elements in  $B$ . In subset  $A$ , with probability  $|B|/|S|$ , each element sends its corresponding request to an element in  $B$ . Combining with Chernoff bound, the inequality  $m < (1 - \epsilon) \cdot |B|/|S| \cdot |A|$  holds with a small probability

$$\zeta_1 \leq g(\epsilon)^{|A||B|/|S|}. \quad (6)$$

Assume  $\gamma \in (0, 1)$  and  $\epsilon(1 - \gamma) \leq 1$ . The probability that  $|H(A)| \leq (1 - \gamma)m$  is

$$\zeta_2 \leq \binom{n}{(1 - \gamma)m} \cdot \left( \frac{(1 - \gamma)m}{n} \right)^m \quad (7)$$

$$\leq \frac{n^{(1-\gamma)m} e^{(1-\gamma)m}}{((1 - \gamma)m)^{(1-\gamma)m}} \cdot \left( \frac{(1 - \gamma)m}{n} \right)^m \quad (8)$$

$$\leq \left( \frac{(1 - \gamma)m}{n} \right)^{(2\gamma-1)m} \quad (9)$$

$$\leq ((1 - \gamma))^{(2\gamma-1)m}. \quad (10)$$

Combining inequalities (6) and (10), the failure probability is at most  $\zeta_1 + \zeta_2 \leq g(\epsilon)^{|A||B|/|S|} + ((1 - \gamma))^{(2\gamma-1)(1-\epsilon) \cdot |B|/|S| \cdot |A|}$ . Thus the lemma is proved.  $\square$

**Lemma 16.** *Let  $S$  be the list of all  $m$  elements that will take average operations. For some fixed  $\alpha > 0$ , with the failure probability not larger than  $1/(\log m)^3$ , the following conclusions hold.*

- (1) After  $O(\log m)$  stages of average operations, there is an  $\alpha$ -shrink.
- (2) After  $O(\log m)$  stages of integer average operations, there is an  $\alpha$ -shrink if  $gap(L)$  is at least  $H$  for some  $H$  to be large enough.

*Proof.* Let  $h = \max S - \min S$ ,  $a = \min\{S\}$ , and  $b = \max\{S\}$ ; the median is  $(a + b)/2$ , which is also equal to  $a + h/2$ .  $A$  and  $B$  denote the sets of elements greater than and not larger than the median  $(a + h/2)$  of  $\min\{S\}$  and  $\max\{S\}$ , respectively. In general, assume  $|A|$  is not larger than  $|B|$ . Let  $A_0 = A$ ,  $B_0 = B$ ,  $S_0 = S$ , and  $j = 0$ . Three periods of communication stages will be discussed in the following part.

If  $|A_j| \geq (\log m)/(\log \log m)^5$ , enter Period 1 below, or otherwise, enter Period 3.

**Period 1.**  $O(1)$  communication phases will be performed as follows.

Use  $A_{i+1}$  to represent a set of elements  $a$ , which satisfies one of the following conditions:

- (1)  $a \in A_i$ ;  $a$  does not participate in any average operation;
- (2)  $a$  is one of the elements generated by averaging elements  $c$  and  $d$ , where  $c$  and  $d$  belong to set  $A_i$ .

Use  $B_{i+1}$  to represent a set of elements  $a$ , which satisfies one of the following conditions:

- (1)  $a \in B_i$ ;  $a$  does not participate in any average operation;
- (2)  $a$  is one of the elements generated by averaging elements  $c$  and  $d$ , where  $c$  and  $d$  belong to set  $B_i$ .

Assume that  $S_{j+1}^{(1)} = A_{j+1} \cup B_{j+1}$ .

Then, select three constants  $\tau_1 > \tau_2 > \tau_3 > 0$  with the relation  $\tau_1 = 2\tau_2 = 4\tau_3$ , and constants  $0 < \gamma_1, \gamma_2 < 1$ . For analysis,  $\gamma_1$  is set to 0.05, and  $\gamma_2$  is set to 0.1. It is assumed that  $\epsilon \leq \gamma_1$ . According to Lemma 15, for constant  $\beta \in (0, 1)$ , the failure probability of  $|A_{i+1}| \leq (1 - \beta)|A_i|$  is not larger than  $2(1 - a)^{|A_i|} \leq 2(1 - a)^{(\log m)/(\log \log m)^5}$ . Choose an integer  $i$  that makes  $|A_i| \leq (1 - \beta)^i |A_0| \leq \gamma_1 |A_0|$  hold. After  $i$  stages, we have the following inequalities.

$$|A_i| \leq (1 - \beta)^i |A_0| \leq \gamma_1 |S_0|, \quad (11)$$

$$|B_i| \geq (1 - \gamma_1) |S_0| \text{ and}, \quad (12)$$

$$\max\{B_i\} \leq a + (1 - \tau_1)h \text{ with } \tau_1 \in (0, 1). \quad (13)$$

Its failure probability is at most  $2i(1 - a)^{(\log m)/(\log \log m)^5}$ .

$D_3$  is represented as the set of elements that belong to  $(a + (1 - \tau_3)h, a + h)$  in set  $S_j^{(1)}$ . If  $|D_3| \geq (\log m)/(\log \log m)^5$ , the process will execute Period 2, or otherwise, Period 3 will be executed.

Period 2.  $O(\log m)$  communication phases will be performed as follows.

$S_{i_1}^{(1)}$  is defined to represent the final set generated by the set  $S$  from period 1 after a series of communication operations.  $S_0^{(2)}$  is defined as  $S_{i_1}^{(1)}$ .  $S_j^{(2)}$  is used to represent the set of elements via  $j$  phases in this period.

$D_{j,3}$  is defined to represent the list of elements with values in the range of  $(a + (1 - \tau_3)h, a + h]$  after  $j$  stages, while  $D'_{j,3}$  indicates the elements in the range of  $[a, a + (1 - \tau_3)h]$ . Similarly,  $D_{j,2}$  and  $D'_{j,2}$  represent the elements in  $(a + (1 - \tau_2)h, a + h]$  and  $[a, a + (1 - \tau_2)h]$  in  $S_j^{(2)}$  after  $j$  stages, respectively. As it is stated,  $|D_{j,3}| \geq (\log m)/(\log \log m)^5$  always holds in Period 2.

The number of elements that belong to  $[a + (1 - \tau_1)h, a + h]$  is at most  $\gamma_1 m$ , thus they can help up to  $\gamma_1 m$  elements (the value of these elements is at most  $a + (1 - \tau_1)h$ ) increase to at least  $a + (1 - \tau_2)h$ . The reason is that each element with a maximum value  $a + h$  can contribute up to  $\tau_2 h$ . For the  $\gamma_1 m$  elements with values greater than  $a + (1 - \tau_1)h$ , the contribution of these elements is at most  $\gamma_1 m \cdot \tau_2$ . Thus, the quantity of elements in the range of  $[a, a + (1 - \tau_2)h]$  is at least  $(1 - \gamma_2)m$ . According to Lemma 15, after each phase of operation, the number of elements in set  $D_3$  would be decreased in a fixed rate, as each element in  $D_3$  has a greater probability to be able to do operations with the elements in set  $D'_{j,2}$ .

Let  $P_1$  be the probability that the number of elements in  $D_{j,3}$  that select elements in  $D_{j,2}$  to take average is less than  $(\gamma_2 + \epsilon)|D_{j,3}|$ , and let  $P_2$  be the probability that the number of elements in  $D_{j,3}$  that take average with the ones in  $D'_{j,2}$  is less than  $(1 - \gamma)(1 - \epsilon) \cdot (|D'_{j,2}|/|S|)|D_{j,3}|$ . According to Theorem 9 and Lemma 15, we have  $P_1 \leq g(\epsilon)^{|D_{j,3}|}$  and  $P_2 \leq g(\epsilon)^{|D_{j,3}|/2} + 2(1 - a)^{|D_{j,3}|}$ . Thus, the elements in  $D_{j+1,3}$  will be reduced.

Accordingly, with a failure probability no smaller than  $2(1 - a)^{|D_{j,3}|} + g(\epsilon)^{|D_{j,3}|/2}$ , that is,  $o(1/(\log m)^4)$ , the number of elements in  $D_{j,3}$  is decreased by not smaller than

$$\begin{aligned} & (1 - \gamma_1)(1 - \epsilon) \cdot \frac{|D'_{j,2}|}{|S|} |D_{j,3}| - (\gamma_2 + \epsilon) |D_{j,3}| \\ & \geq (1 - \gamma_1)(1 - \epsilon)(1 - \gamma_1) |D_{j,3}| - (\gamma_2 + \epsilon) |D_{j,3}| \quad (14) \\ & \geq ((1 - \gamma_1)(1 - \epsilon - \gamma_1) - (2\gamma_1 - \gamma_1)) |D_{j,3}| \\ & \geq (1 - 9\gamma_1) |D_{j,3}|. \end{aligned}$$

It can be seen that  $O(\log m)$  communication stages are needed to reach stage  $j$ . Once  $|D_{j,3}| < (\log m)/(\log \log m)^5$ , enter Period 3.

It is easy to verify the case for all integer averages when the initial gap is big enough. We just need to set up the gap such that  $\tau_2 \text{gap}(K) \geq \tau_2 H \geq 3$ . There is a  $\tau_2 \text{gap}(K)$  gap from the new list to the old list  $K$ .

Period 3.  $O((\log \log m)^2)$  communication phases will be performed as follows.

The set produced from  $S_2^{(2)}$  after corresponding operations is denoted as  $S_{i_2}^{(2)}$ . The equalities  $|D_{i_2,3}| < (\log m)/(\log \log m)^5$  and  $|D'_{i_2,2}| \geq (1 - \gamma_2)|S|$  hold.

Assume that with a very small failure probability, the elements in  $D'_{i_2,2}$  in sending status would select elements from  $B$ . Meanwhile, with a probability not larger than  $3/4$ , the elements fail to do average with those in  $D'_{i_2,2}$ . If an element has sent the requests for  $\mu$  times, then the probability that it would not select an element from  $D'_{i_2,2}$  is not larger than  $(3/4)^\mu$ . The probability that an element has no more than  $t/4$  stages to send requests is not larger than  $g(1/4)^{t/2}$ .

Thus, the probability that an element in  $D_{i_2,3}$  would not do average with that in  $D'_{i_2,2}$  is no greater than  $p_3 = g(1/4)^{t/2} + (3/4)^{t/4}$ . The probability that two elements in  $D_{i_2,3}$  would do average with the same element is no greater than  $O((\log m)^3/m^2)$ .

Consequently, the probability that the number of stages  $t$  in Period 3 is selected to be  $(\log \log m)^2$  will not be larger than  $|A^*|p_3 \leq 1/(\log m)^4$ .

To summarise from the above analysis, the failure probability is not greater than  $1/(\log m)^3$ .  $\square$

*Definition 17.* We treat the communication that includes  $c$  stages as  $\alpha$ -successful when there is an  $\alpha$  shrink in aspect of gap. For further analysis, let parameter  $\delta$  indicate the probability when it fails to achieve a shrink within  $\alpha$ .

The following Lemma 18 is derived in our previous work [8]; thus proof of the lemma is omitted here. Based on this lemma, we present Theorem 19.

**Lemma 18** (see [8]).  *$c$  indicates a parameter. Partition the communication stages into groups with each containing  $c$  stages, which are denoted as  $G_1, G_2, \dots, G_k$ . Then there are  $k$  independent 0, 1 random variables  $r_i$  for each group  $G_i$  such that*

- (1)  $\Pr(G_i \text{ is } \alpha\text{-successful}) \geq \Pr(r_i = 1)$
- (2)  $\Pr(r_i = 1) \geq 1 - \delta$ .
- (3)  $\Pr(\text{there are at least } t \text{ } G_i \text{ to be } \alpha\text{-successful}) \geq \Pr(r_1 + r_2 + \dots + r_k \geq t)$ .

We use  $m$  to indicate the quantity of vehicles. The corresponding value of node  $i$  is  $n_i^*$  when the network converges.

**Theorem 19.** *The following conclusions will be achieved after  $O((\log n)(\log m)/\epsilon)$  stages by applying the proposed randomised algorithm.*

- (1) If  $m \geq (1 + \epsilon)n$ , then each  $n_i^*$  is either 1 or 0
- (2) If  $m \leq (1 - \epsilon)n$ , then each  $n_i \geq 1$  and has a difference not greater than two between every two of them
- (3) If  $(1 - \epsilon)n < m < (1 + \epsilon)n$ , then  $0 \leq n_i \leq 2$ , and the numbers of value zero and two are both no greater than  $\epsilon m$ .

*Proof.* Applying Lemma 18, and Chernoff bound, it can be known that a  $O(1)$  gap could be achieved after  $O((\log n)(\log m))$

stages. And apply Lemma 16  $O((\log n)(\log m))$  stages if the difference is bounded by a fixed integer. The following cases are presented if the gap is  $O(1)$ .

- (i)  $m \geq (1 + \epsilon)n$ : the quantities of vehicles with value zero and two are at least  $\epsilon n$  in the same condition. We bound the largest elements with a constant. In this case, the number of items of 0 and 1 is at least  $(1 + \epsilon)n/2$ .
- (ii)  $(1 - \epsilon)n < m < (1 + \epsilon)n$ : when the number of vehicles with value greater than zero is at least  $(1 - \epsilon)n$ , the vehicles with value three will be removed via  $O(\log n)$  stages. Thus,  $0 \leq n_i \leq 2$  and no more than  $\epsilon n$  vehicles have the value zero, and no more than  $\epsilon n$  with two.
- (iii)  $m < (1 - \epsilon)n$ : after  $O(\log n)$  stages, not less than  $\gamma n$  elements are at least two as to constant  $\gamma > 0$ .

When the number of vehicles with value two is not less than  $\gamma n$ , the number of zero would be decreased by  $\theta n$  for constant  $\theta$ . And after  $O(\log n)$  stages, all the zero items will be removed. When two vehicles are with gap of at least two encounters, their values will be updated such that the gap will not be greater than 1.

Assume  $S$  indicates the set of the rest of the items, and  $a = \max\{x : x \in S\}$ , and  $b = \min\{x : x \in S\}$ .  $\|x\|_S$  is defined as the value of  $x$ , and it is assumed that  $\|a\|_S \leq \|b\|_S$ . Additionally, let  $gap(S) > 2$  and  $|\{x : x \in S \text{ and } x \geq a + 2\}| \geq \gamma k$ . The number of  $a$  that selects value not smaller than  $a + 2$  is  $\delta m$ . In this way, value  $a$  should disappear after  $O(\log m)$  stages.

□

To show the efficiency of real average operations, Theorem 20 is proved.

**Theorem 20.** Assume  $m$  denotes the quantity of vehicles. Then the randomised algorithm takes  $O((\log n)(\log m)/\epsilon)$  stages to enter into an  $\epsilon$ -balanced status.

## 5.2. Lower Bound

**Theorem 21.** For a complete connected graph,  $\Omega(\log n)$  stages are needed to reach an  $\epsilon$ -balanced network status.

*Proof.* Referring to the proposed randomised algorithm, it may only double the number of elements via each phase. Consequently, the number of communication stages consumed is  $\Omega(\log n)$  if the quantity of vehicular nodes  $m$  is more than the number of replicas  $n$ . □

## 6. Simulation and Analysis

We conduct our simulations by using NS-3 simulator. Consider the complexity of performance evaluation; we use the OpenStreetmap [35] to extract an area for simulation. The size of the selected area is set to 2000 m  $\times$  2000 m, the satellite map of which is shown as Figure 4(a). Meanwhile, we use

TABLE 1: Simulation settings.

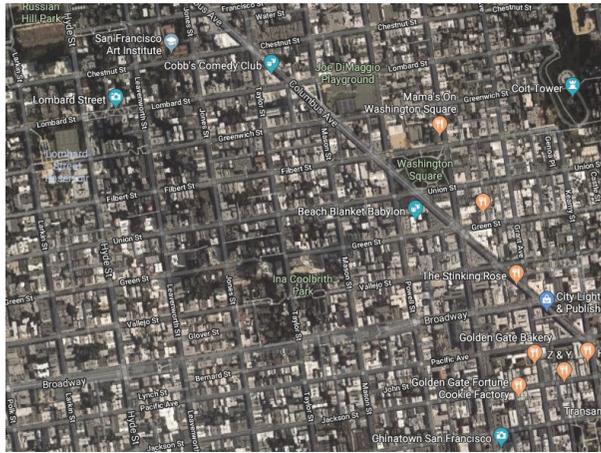
Parameters	Settings
Size of simulation area	2000m $\times$ 2000m
Simulation time	1 hour
MAC Protocol	IEEE 802.11p
Packet size	512 bytes
Vehicle communication range	300m
Vehicle velocity	30 - 60 km/h
Number of vehicles	600 - 800
Number of data replicas	400 - 800
Shadowing model	Lognormal Shadowing model
Path loss model	Two-ray
SNR threshold	4 dBm

SUMO [36] to transform the extracted area into a simplified road network presented as Figure 4(b). We also use SUMO to generate the movement trajectories of vehicles, which will be input in the simulator to describe the vehicles' movement.

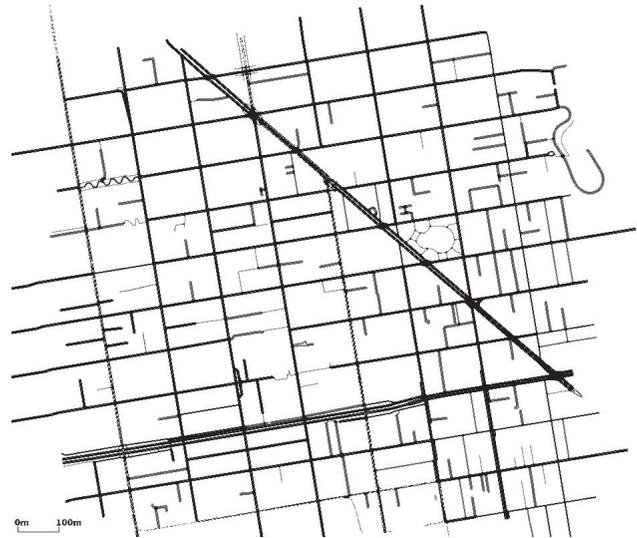
In the simulations, the number of vehicles is varied from 600 to 800 with a step length 50 to reflect a dense vehicular environment. The vehicle velocity is varied from 30 to 60 km/h that follows a normal distribution. The number of message copies is varied from 400 to 800 with a step length 100. The packet size is set to 512 bytes. The transmission range of vehicles is set to 300 m. As to the communication protocol, IEEE 802.11p is adopted to guarantee the reliability of information transmission. The two-ray path loss model is applied in the simulation as the model can calculate both the direct path and the ground reflection path. The signal-to-noise ratio (SNR) threshold is set to 4 dBm. Considering the impact of obstacles on wireless signal in urban environment as well as vehicular mobility, we apply the Lognormal Shadowing model that is suitable for the proposed scenarios. The list of simulation parameters is shown in Table 1.

**6.1. Compared Algorithms.** To evaluate the performance of the proposed replication-based randomised algorithm, we compared it with several data dissemination algorithms in vehicular networks, which are described below.

- (i) Constrained Capacity Replication (CCR): CCR is a distributed algorithm which can assist the vehicles to select data replication strategy autonomously according to the current network capacity.
- (ii) DOVE: DOVE controls the number of receivers in data dissemination and transforms the problem to the processor scheduling problem by utilising road layout and traffic information.
- (iii) EDDA: EDDA considers the common urban and highway scenarios. It selects all the independent pairs of nodes firstly, and then takes average operations between the corresponding nodes. The average operations will run iteratively until the network converges.
- (iv) Our proposed randomised algorithm: the proposed randomised algorithm considers the heterogeneous



(a) Google map



(b) Road topology layout

FIGURE 4: Selected area.

properties of vehicles while controlling the number of data replicas. Each node in sending status randomly selects one of its neighbours to send a contact request while the receiving node selects a request with the largest gap to take the proportional average operation.

6.2. *Performance Metrics.* We evaluate the performance of these algorithms according to the following metrics.

- (i) Number of communication stages: it indicates the average operations for the network to be balanced, which can represent the communication overhead of data dissemination to a certain extent. Meanwhile, it can also reflect the number of data transmissions for network balance as well as the network convergence complexity.
- (ii) Dissemination delay: it presents the consumed time to obtain network balance, which can be utilised to measure effectiveness of the algorithms.
- (iii) Packet delivery ratio: it can directly indicate how many vehicles could receive the replicas, as well as reflecting the dissemination performance.
- (iv) Throughput: it is used to evaluate the proposed randomised algorithm when the capabilities of vehicles are considered.

6.3. *Impact of Number of Vehicular Nodes.* To evaluate the performance of the proposed randomised algorithm with different network densities, we vary the number of vehicles from 600 to 800 to represent increasing network size.

The experimental performance analysis for the consumed communication stages is depicted in Figure 5. The proposed randomised algorithm outperforms other compared data dissemination schemes when the number of vehicles increases

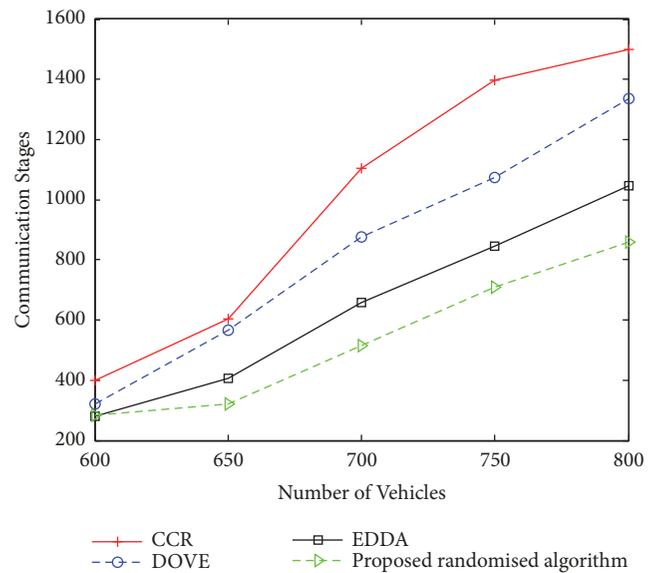


FIGURE 5: Communication stages comparison when the number of vehicles varies.

from 600 to 800, which means it needs fewer average operations to achieve network balance. CCR mainly considers network capacity to determine the replication limit while DOVE utilises road layout and traffic information to reach a desired number of vehicular receivers and minimise the dissemination delay. Our proposed randomised algorithm benefits from strong connectivity of the dense vehicular scenarios such that it can obtain network convergence with fewer communication stages than other schemes. As EDDA is developed for scenarios with normal urban density and highway, it is slightly inferior to the randomised algorithm. With the increasing number of vehicular nodes, the consumed communication stages grow for all the compared algorithms.

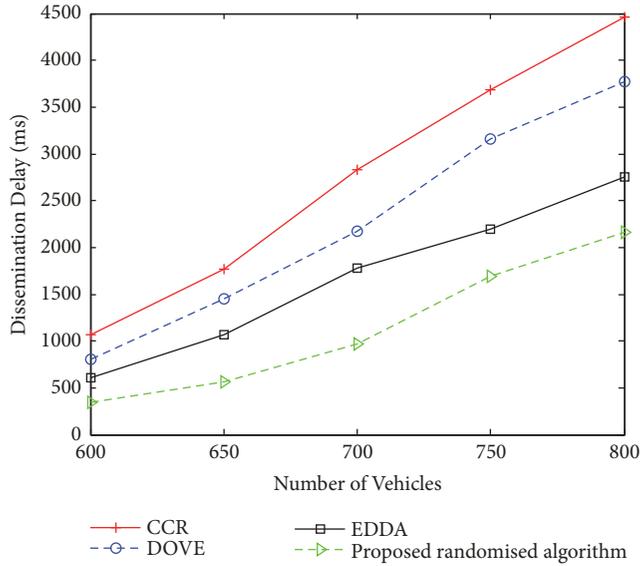


FIGURE 6: Dissemination delay comparison when the number of vehicles varies.

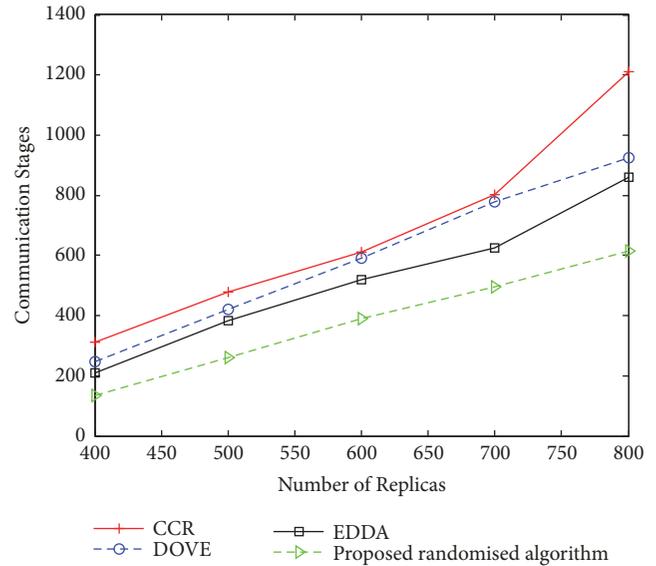


FIGURE 8: Communication stages comparison when the number of replicas varies.

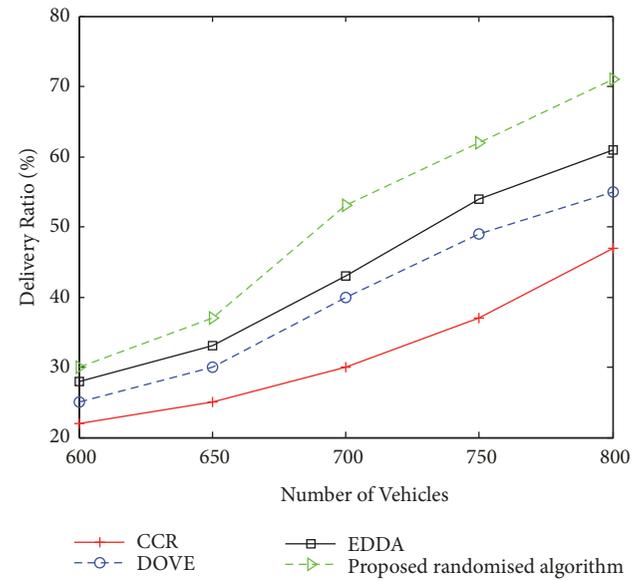


FIGURE 7: Delivery ratio comparison when the number of vehicles varies.

The variation of dissemination delay of the compared algorithms when the number of participating vehicles increases is shown in Figure 6. As is shown in the figure, when the traffic densities become higher, the dissemination delay increases for all the compared algorithms as it needs more time to fulfill data dissemination. The proposed randomised algorithm takes less time to complete data dissemination and realise network consensus compared to the other three schemes. This is because it considers the different capabilities of the vehicles in data dissemination and thus constructs a better replication strategy which could reduce data dissemination delay.

Figure 7 shows the performance of packet delivery ratio of the compared algorithms when the number of nodes involved

in data dissemination varies. In terms of delivery ratio, our proposed randomised algorithm presents an improvement compared with EDDA, DOVE, and CCR. Also, the delivery ratio of all the algorithms increases when the number of nodes increases from 600 to 800. More vehicles cooperatively participate in data dissemination such that the network connectivity would be enhanced. Accordingly, the successful packet transmissions will be improved by frequent vehicle communication instead of transmission failures caused by fewer forwarding vehicular nodes.

**6.4. Impact of Number of Data Replicas.** By varying the number of data replicas, we evaluate the impact of number of allowed data replicas on consumed communication stages and dissemination delay of the compared algorithms.

Figure 8 shows the changing trend of our proposed randomised algorithm and other compared algorithms with increasing number of data replicas. As to communication stages, the randomised algorithm performs fewer operations than EDDA profiting from the strongly connected network property of dense networks, while DOVE and CCR both need more stages for network balance. Additionally, as there are more data replicas to be disseminated in the network, it can be seen that more average operations are required to achieve network convergence. As a result, the four compared algorithms would consume increasing communication stages when more data replicas are spreading to the dissemination area.

The change of data dissemination delay under different number of data replicas is shown in Figure 9. As observed from the figure, the dissemination delay of all the compared algorithms increases when the number of replicas increases with a step length 100. The reason mainly lies in that more message replicas would take more communication stages for the network to be balanced, which leads to higher

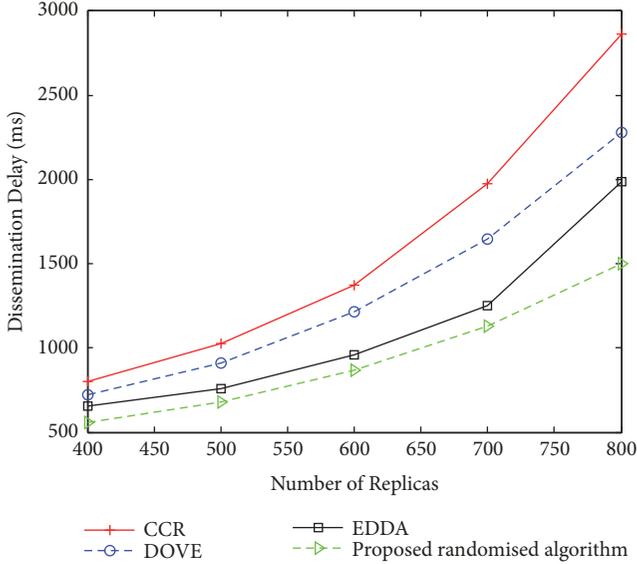


FIGURE 9: Dissemination delay comparison when the number of replicas varies.

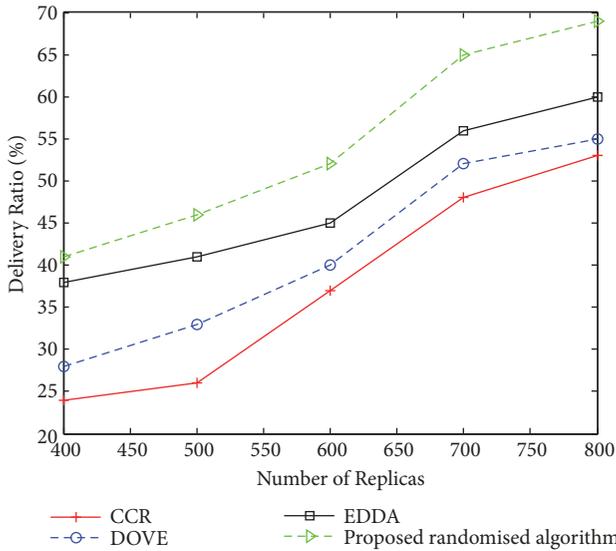


FIGURE 10: Delivery ratio comparison when the number of replicas varies.

dissemination delay. Meanwhile, data dissemination would be accelerated when the randomised algorithm is applied to the scenario as the algorithm selects pairs of nodes with the largest gap and adjusts how to do average operations according to the properties of vehicles. This is why the randomised algorithm outperforms EDDA, DOVE, and CCR.

Figure 10 depicts packet delivery ratio of the four compared algorithms varying the number of data replicas. Other than taking advantage of heterogeneous network and random average operations, the proposed randomised algorithm also benefits from better network connectivity to obtain higher delivery ratio than other algorithms. EDDA controls the number of replicas in arbitrarily connected network while

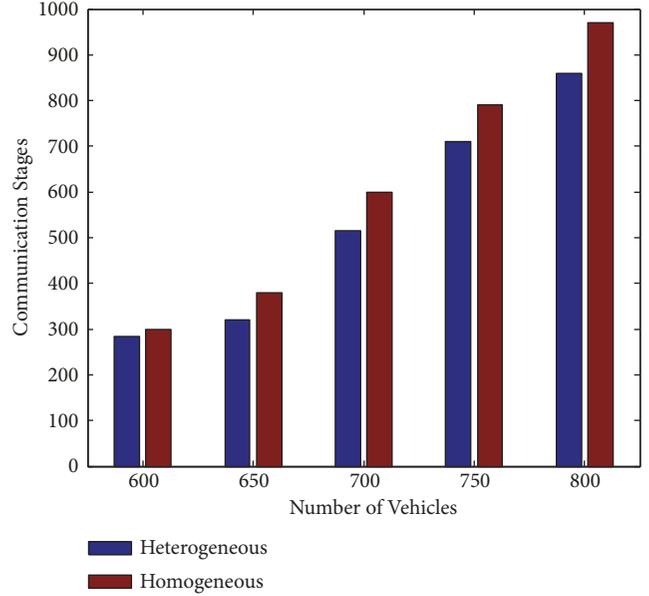


FIGURE 11: Communication stages comparison between the cases of homogeneous vehicles and heterogeneous vehicles.

CCR focuses on network capacity and DOVE wants to minimise the delay; however, they all only consider the vehicles with the same capability. The growth of delivery ratio is similar to Figure 9 that with the number of replicas changing from 400 to 800, the ratio increases for all the compared solutions. The performance comparison verifies that our proposed randomised algorithm can efficiently improve the performance of data dissemination and expedite the network convergence.

6.5. Comparing Different Versions of the Proposed Randomised Algorithm. We evaluate the communication stages and throughput of the proposed randomised algorithm in dense traffic scenario, in the case of two conditions; that is, the vehicles are homogeneous and heterogeneous. The number of vehicular nodes is varied from 600 to 800 with a step length 50. The comparison results are shown in Figures 11 and 12, respectively. In Figure 11, we can see that the number of communication stages consumed when the vehicles have the same capability is larger than the one consumed when the different capabilities of vehicles are considered. Figure 12 presents the total data sent during data dissemination. In the figure, we can see that the randomised algorithm with different capabilities achieves a 10% improvement compared with the algorithm with homogeneous vehicles. The comparison shows that the consideration of heterogeneous vehicular networks could reduce the consumed communication overhead while improving the system throughput to some extent.

## 7. Conclusion

To enhance data dissemination in dense vehicular scenarios, we propose a network architecture, considering heterogeneous network composed of vehicles with different capabilities. By studying data replication and network consensus

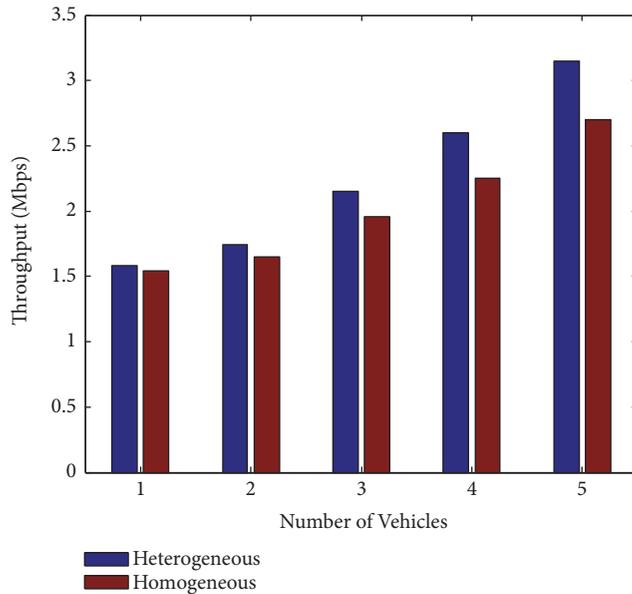


FIGURE 12: Throughput comparison between the cases of homogeneous vehicles and heterogeneous vehicles.

properties, two replication-based algorithms including a deterministic algorithm and a distributed randomised algorithm are designed. In the proposed algorithms, the vehicles take proportional average operations depending on their own capabilities. The operations will be iterated until the network converges. Mathematical analysis is derived to evaluate the complexity of network convergence. An upper bound and a lower bound of the randomised algorithm are analysed in detail. Simulation results show that the proposed randomised algorithm can reduce data dissemination delay and improve communication overhead.

## 8. Prospective Directions

In this study, we consider scenarios with relative ideal communication situations, which means that link interruption and other interferences are not considered. Future work should incorporate the factors that will affect data dissemination. Also, as a future prospect, we intend to enrich data dissemination mechanisms which can be adapted to complicated scenarios in IoT. Solutions that apply infrastructures such as unmanned aerial vehicles (UAVs) in cooperative networks could also be seen as a promising prospect to enhance network performance. The effectiveness of the replication-based algorithms in more IoT application scenarios needs further verification.

## Data Availability

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

## Disclosure

An earlier conference version of this paper [8] has been presented in 12th International Conference on WASA.

## Conflicts of Interest

The authors declare that they have no conflicts of interest.

## Acknowledgments

This work is supported by the National Science Foundation of China (no. 61772385, no. 61373040, and no. 61572370) and autonomous electric vehicle driving ability and safety evaluation technology and system development (SQ2018YFB010236-04).

## References

- [1] X. Wang, Z. Ning, X. Hu et al., "A city-wide real-time traffic management system: enabling crowdsensing in social internet of vehicles," *IEEE Communications Magazine*, vol. 56, no. 9, pp. 19–25, 2018.
- [2] R. Ghebleh, "A comparative classification of information dissemination approaches in vehicular ad hoc networks from distinctive viewpoints: A survey," *Computer Networks*, vol. 131, pp. 15–37, 2018.
- [3] Y. Li, D. Jin, P. Hui, and S. Chen, "Contact-aware data replication in roadside unit aided vehicular delay tolerant networks," *IEEE Transactions on Mobile Computing*, vol. 15, no. 2, pp. 306–321, 2016.
- [4] G. Wang, Z. Wang, and J. Wu, "A local average broadcast gossip algorithm for fast global consensus over graphs," *Journal of Parallel and Distributed Computing*, vol. 109, pp. 301–309, 2017.
- [5] R. Azimi and H. Sajedi, "A decentralized gossip based approach for data clustering in peer-to-peer networks," *Journal of Parallel and Distributed Computing*, vol. 119, pp. 64–80, 2018.
- [6] T. Spyropoulos, K. Psounis, and C. S. Raghavendra, "Efficient routing in intermittently connected mobile networks: the single-copy case," *IEEE/ACM Transactions on Networking*, vol. 16, no. 1, pp. 63–76, 2008.
- [7] A. Balasubramanian, B. N. Levine, and A. Venkataramani, "Replication routing in DTNs: a resource allocation approach," *IEEE/ACM Transactions on Networking*, vol. 18, no. 2, pp. 596–609, 2010.
- [8] J. Zhu, C. Huang, X. Fan, and B. Fu, "An efficient distributed randomized data replication algorithm in VANETs," in *Wireless Algorithms, Systems, and Applications*, pp. 369–380, Springer, Guilin, China, 2017.
- [9] X. Fan, C. Huang, J. Zhu, and B. Fu, "R-DRA: a replication-based distributed randomized algorithm for data dissemination in connected vehicular networks," *Wireless Networks*, pp. 1–16, 2018.
- [10] A. Torres, C. T. Calafate, J.-C. Cano, P. Manzoni, and Y. Ji, "Evaluation of flooding schemes for real-time video transmission in VANETs," *Ad Hoc Networks*, vol. 24, pp. 3–20, 2015.
- [11] A. Takahashi, H. Nishiyama, N. Kato, K. Nakahira, and T. Sugiyama, "Replication control for ensuring reliability of convergecast message delivery in infrastructure-aided DTNs," *IEEE Transactions on Vehicular Technology*, vol. 63, no. 7, pp. 3223–3231, 2014.
- [12] M. Xing, J. He, and L. Cai, "Utility maximization for multimedia data dissemination in large-scale VANETs," *IEEE Transactions on Mobile Computing*, vol. 16, no. 4, pp. 1188–1198, 2017.

- [13] Y. Wu, Y. Zhu, H. Zhu, and B. Li, "CCR: Capacity-constrained replication for data delivery in vehicular networks," in *Proceedings of the IEEE Conference on Computer Communications (INFOCOM '13)*, pp. 2580–2588, Turin, Italy, April 2013.
- [14] X. Shen, X. Cheng, L. Yang, R. Zhang, and B. Jiao, "Data dissemination in VANETs: a scheduling approach," *IEEE Transactions on Intelligent Transportation Systems*, vol. 15, no. 5, pp. 2213–2223, 2014.
- [15] T. Yan, W. Zhang, and G. Wang, "DOVE: Data dissemination to a desired number of receivers in VANET," *IEEE Transactions on Vehicular Technology*, vol. 63, no. 4, pp. 1903–1916, 2014.
- [16] Q. Xiang, X. Chen, L. Kong, L. Rao, and X. Liu, "Data preference matters: a new perspective of safety data dissemination in vehicular ad hoc networks," in *Proceedings of the 34th IEEE Annual Conference on Computer Communications and Networks (IEEE INFOCOM '15)*, pp. 1149–1157, Kowloon, Hang Kong, May 2015.
- [17] Z. Zhao, W. Dong, J. Bu, T. Gu, and G. Min, "Accurate and generic sender selection for bulk data dissemination in low-power wireless networks," *IEEE/ACM Transactions on Networking*, vol. 25, no. 2, pp. 948–959, 2017.
- [18] F. Chen, D. Zhang, J. Zhang et al., "Distribution-aware cache replication for cooperative road side units in VANETs," *Peer-to-Peer Networking and Applications*, vol. 11, no. 5, pp. 1075–1084, 2018.
- [19] Z. Jiang, S. Zhou, X. Guo, and Z. Niu, "Task replication for deadline-constrained vehicular cloud computing: optimal policy, performance analysis, and implications on road traffic," *IEEE Internet of Things Journal*, vol. 5, no. 1, pp. 93–107, 2018.
- [20] P.-Y. Chen, S.-M. Cheng, and M.-H. Sung, "Analysis of data dissemination and control in social internet of vehicles," *IEEE Internet of Things Journal*, vol. 5, no. 4, pp. 2466–2477, 2018.
- [21] C. Lin, D. Deng, and C. Yao, "Resource allocation in vehicular cloud computing systems with heterogeneous vehicles and roadside units," *IEEE Internet of Things Journal*, vol. 5, no. 5, pp. 3692–3700, 2018.
- [22] C. Ghorai and I. Banerjee, "A robust forwarding node selection mechanism for efficient communication in urban VANETs," *Vehicular Communications*, vol. 14, pp. 109–121, 2018.
- [23] Q. Ding, X. Zeng, X. Zhang, and D. K. Sung, "A public goods game theory-based approach to cooperation in VANETs under a high vehicle density condition," *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–11.
- [24] G. Shi, B. Li, M. Johansson, and K. H. Johansson, "Finite-time convergent gossiping," *IEEE/ACM Transactions on Networking*, vol. 24, no. 5, pp. 2782–2794, 2016.
- [25] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, "Randomized gossip algorithms," *IEEE Transactions on Information Theory*, vol. 52, no. 6, pp. 2508–2530, 2006.
- [26] F. Fagnani and S. Zampieri, "Randomized consensus algorithms over large scale networks," *IEEE Journal on Selected Areas in Communications*, vol. 26, no. 4, pp. 634–649, 2008.
- [27] C. Yu, B. D. Anderson, S. Mou, J. Liu, F. He, and A. S. Morse, "Distributed averaging using periodic gossiping," *Institute of Electrical and Electronics Engineers Transactions on Automatic Control*, vol. 62, no. 8, pp. 4282–4289, 2017.
- [28] J.-Y. Chen, G. Pandurangan, and D. Xu, "Robust computation of aggregates in wireless sensor networks: Distributed randomized algorithms and analysis," *IEEE Transactions on Parallel and Distributed Systems*, vol. 17, no. 9, pp. 987–1000, 2006.
- [29] T. C. Aysal, M. E. Yildiz, A. D. Sarwate, and A. Scaglione, "Broadcast gossip algorithms for consensus," *IEEE Transactions on Signal Processing*, vol. 57, no. 7, pp. 2748–2761, 2009.
- [30] S. Wu and M. G. Rabbat, "Broadcast gossip algorithms for consensus on strongly connected digraphs," *IEEE Transactions on Signal Processing*, vol. 61, no. 16, pp. 3959–3971, 2013.
- [31] M. Franceschelli, A. Giua, and C. Seatzu, "Gossip based asynchronous and randomized distributed task assignment with guaranteed performance on heterogeneous networks," *Nonlinear Analysis: Hybrid Systems*, vol. 26, pp. 292–306, 2017.
- [32] A. Nedic and J. Liu, "On convergence rate of weighted-averaging dynamics for consensus problems," *Institute of Electrical and Electronics Engineers Transactions on Automatic Control*, vol. 62, no. 2, pp. 766–781, 2017.
- [33] R. Motwani and P. Raghavan, "Randomized algorithms," *Acm Computing Surveys*, vol. 26, pp. 48–50, 1995.
- [34] M. Li, B. Ma, and L. Wang, "On the closest string and substring problems," *Journal of the ACM*, vol. 49, no. 2, pp. 157–171, 2002.
- [35] M. Haklay and P. Weber, "Openstreetmap: user-generated street maps," *IEEE Pervasive Computing*, vol. 7, no. 4, pp. 12–18, 2008.
- [36] SUMO-Simulation of Urban Mobility, <http://sumo.sourceforge.net>.



**Hindawi**

Submit your manuscripts at  
[www.hindawi.com](http://www.hindawi.com)

