

Computer Networks

MLM-WR: An Effective Workers Recruitment Scheme Based on Swarm Intelligence for Mobile Crowd Sensing --Manuscript Draft--

Manuscript Number:	
Article Type:	VSI:SI DCOSS 2022
Section/Category:	
Keywords:	Mobile Crowd Sensing Truthful workers discovery Sensing difference discovery Task assignment Swarm Intelligence
Corresponding Author:	Anfeng Liu Central South University Changsha, China
First Author:	Jiaheng Lu
Order of Authors:	Jiaheng Lu Zhenzhe Qu Anfeng Liu Shaobo Zhang Neal N. Xiong Tian Wang
Manuscript Region of Origin:	Asia Pacific
Abstract:	Mobile Crowd Sensing (MCS) is a powerful paradigm of sensing systems, which collects data from workers to build applications. It is vital to identify truthful workers, then match sensing data tasks and truthful workers to collect data. We apply matching theory to improve the quality of the MCS system from two aspects, i.e., one is truthful workers discovery, and the other is sensing difference discovery. In the truthful workers discovery phase, we establish an approach to obtain workers' credibility by comparing their report data with Ground Truth Data (GTD) and Sub-GTD. In the sensing difference discovery phase, the workers' sensing quality of different attributes is obtained by comparing Attribute-GTD and Sub-Attribute-GTD. And workers' sensing quality of different locations is attained by the integration of absolute and relative sensing location preference. Combining the two aspects, an effective workers recruitment scheme based on swarm intelligence for MCS, Multi-attribute and Local Matching based Workers Recruitment (MLM-WR), is proposed. MLM-WR utilizes the PSO algorithm to dynamically assign workers. We finally carry out extensive evaluations, where results demonstrate the superiority of our approach over state-of-the-art approaches.
Suggested Reviewers:	Huan Zhou China Three Gorges University houhuan117@gmail.com Yongxuan Lai Xiamen University laiyx@xmu.edu.cn
Opposed Reviewers:	

Dear Editor in Chief, Editors and Reviewers,

We would like to submit the enclosed manuscript entitled "**MLM-WR: An Effective Workers Recruitment Scheme Based on Swarm Intelligence for Mobile Crowd Sensing**", which we wish to be considered for publication in the Computer Networks journal.

The main contributions of this manuscript include:

(1) A truthful workers discovery approach is proposed to evaluate the credibility of workers by comparing their report data with the GTD and Sub-GTD. First, the UAV is dispatched and a small amount of data is collected as GTD. Then, we propose a complete workers credibility calculation and inference method. In the proposed approach, the credibility of workers is continuously adjusted based on the results of the comparison of their submitted data with GTD. Through such comparison, the credibility of truthful workers increases. When the credibility rises to a certain threshold, the workers are considered truthful, at that time, their reported data can act as Sub-GTD so to speed up the trust identification of workers. By utilizing the data collected by the UAV, this approach can accurately identify truthful workers and exclude malicious workers. Moreover, Sub-GTD makes the speed of identifying different types of workers very fast.

(2) We propose a sensing difference discovery approach, in which the variability of the workers' sensing quality for different attributes of data, and different locations are obtained by comparing Attribute-GTD and Location-GTD. Then, A-GTD and L-GTD are used as criteria to compare the accuracy of workers on the attribute and location dimensions. Based on the results of the comparison, the workers' preferences in each dimension are adjusted, and the data submitted by those workers who have reached a threshold in attributes and locations are used as the sub-A-GTD and sub-L-GTD of the attributes and locations. The location preference is divided into absolute location preference and relative location preference to adapt to a wider range of scenarios so that more workers' attributes and location credibility can be calculated and inferred.

(3) By utilizing swarm intelligence, we convert the MCS from the traditional IoT system into the AIoT system. Combining truthful workers discovery and sensing difference discovery, an effective workers recruitment scheme based on swarm intelligence for Mobile Crowd Sensing Systems, MLM-WR, is proposed to select truthful workers matching sensing data tasks in terms of attributes and locations, which optimizes the quality of sensing data and saves costs. We finally carry out extensive evaluations, where results demonstrate the superiority of our approach over the state-of-the-art approaches. In terms of workers' recognition rate, data quality improvement, and cost-saving, we have improved by 36.12%, 27.10%, and 40.19% respectively.

We state that:

- a. This manuscript is the authors' original work and has not been published nor has it been submitted simultaneously elsewhere.
- b. All authors have checked the manuscript and have agreed to the submission.

We greatly appreciate your consideration of our manuscript and look forward to receiving comments from reviewers!

With best wishes,

Sincerely,

All authors

MLM-WR: An Effective Workers Recruitment Scheme Based on Swarm Intelligence for Mobile Crowd Sensing

Jiaheng Lu ^a, Zhenzhe Qu ^a, Anfeng Liu^a, Shaobo Zhang ^b, Neal N. Xiong^{c,*}, Tian Wang^d

^a School of Computer Science and Engineering, Central South University, Changsha 410083 China

^b School of Computer Science and Engineering, Hunan University of Science and Technology, Xiangtan 411201 China

^c Department of Computer Science and Mathematics, Sul Ross State University, Alpine, TX 79830, USA

^d Artificial Intelligence and Future Networks, Beijing Normal University & UIC, Zhuhai, Guangdong 519087, China

ARTICLE

INFO

ABSTRACT

Article history:

Received

Received in revised form

Accepted

Available online

Keywords:

Mobile Crowd Sensing

Truthful workers discovery

Sensing difference discovery

Task assignment

Swarm Intelligence

Mobile Crowd Sensing (MCS) is a powerful paradigm of sensing systems, which collects data from workers to build applications. Recruitment of workers is one of the most important issues for constructing high-quality applications. It is vital to identify the truthful workers first, then match the sensing data task and truthful workers to collect data. Based on these observations, we apply matching theory to improve the quality of the whole MCS system from two aspects, i.e., one is truthful workers discovery, and the other is sensing difference discovery. In the truthful workers discovery phase, we establish a truthful workers discovery approach to obtain workers' credibility by comparing their report data with the Ground Truth Data (GTD) and Sub-GTD. In the sensing difference discovery phase, the workers' sensing quality of different attributes is accurately obtained by comparing Attribute-GTD and Sub-Attribute-GTD. And workers' sensing quality of different locations is attained by the integration of absolute and relative sensing location preference. Combining truthful workers discovery and sensing difference discovery, an effective workers recruitment scheme based on swarm intelligence for Mobile Crowd Sensing, Multi-attribute and Local Matching based Workers Recruitment (MLM-WR), is proposed to match truthful workers with sensing data tasks in terms of attributes and locations at high quality and low cost. Our MLM-WR scheme utilizes the PSO algorithm to assign truthful workers to the most suitable places. The sensing location credibility and sensing attribute credibility which are obtained from the sensing difference discovery are used to calculate the fitness value of PSO. Meanwhile, MLM-WR dynamically sets the number of recruited workers according to the number of detected truthful workers which can make the most of limited resources. We finally carry out extensive evaluations, where results demonstrate the superiority of our approach over the state-of-the-art approaches.

1. Introduction

In traditional sensing techniques such as Wireless Sensor Networks (WSNs), distributed sensors are leveraged to obtain real-world conditions [1]. However, traditional commercial sensor network techniques have never been successfully deployed in the real world due to several reasons, such as insufficient node coverage, high installation cost, and lack of scalability [2]. Mobile Crowd Sensing (MCS) [3] has emerged as a new networked sensing paradigm that uses humans as sensors to report the states of the physical world [4]. With the prevalence of sensor-rich smartphones carried by humans (namely workers), billions of smartphones can be employed to sense and collect data in an on-demand style [5]. Smartphones are typically equipped with infrared sensors, temperature sensors, light sensors, etc., which are contact-free and can be

utilized to monitor different environments [6]. Workers can sense numerous data from different environments, then the data are reported to the Data Process Center (DPC) in the MCS. And a large number of data-based applications have been constructed by MCS such as NoiseTube [7], Ear-phone [8]. In the MCS, the platform publishes data collection tasks whose information mainly includes the time and place of data collection and quality requirements for the multi-attribute of the collected data [9], [10], [11]. After workers are informed of the sensing data task, they evaluate their ability to sense data without contact and choose their suitable tasks to launch a request to the platform to participate in the sensing task, platform recruits some workers for data sensing [12], [13], [14].

The data submitted by workers directly affects the quality of the application constructed by the platform [15], [16], [17]. Thus, it is required that the platform recruits truthful and high-quality workers for data collection [18], [19]. Sensing data usually consumes computation, storage, communication resources, and time resources. The higher the sensing quality,

* Corresponding author. Tel.: +86 731 8879628.

E-mail address: jiahengl@csu.edu.cn (J. H. Lu);

zhenzhequ@csu.edu.cn (Z. Z. Qu); afengliu@csu.edu.cn (A. F. Liu);

shaobozhang@hnust.edu.cn (S. B. Zhang); xiongnai@bnu.edu.cn (N. Xiong);

tianwang@bnu.edu.cn (T. Wang).

the higher the demand on workers, and the more effort and time they need to spend. Besides, some tasks need workers to move to the specified location to sense, thus the cost is higher [20], [21], [22]. So, the platform needs to give some reward to the workers to compensate for the cost of completing the task [23], [24], [25]. However, there are some low-credible and even malicious workers, who want to pay as little cost as possible to earn the maximum reward. For this reason, these workers often do not sense data but fabricate and report fake data to the platform which pays the least cost and gets a large reward [26], [27]. Malicious workers not only want to get big rewards but also use the reported malicious data to attack the platform and make the system suffer more than false data [28], [29], [30]. For example, navigation and weather forecasts based on faulty data can disorient users and even kill them in harsh environments which have happened a lot. It is an urgent issue to recruit truthful workers to report high-quality data to construct high-quality applications [18], [19], [31]. The most decisive element in the quality of data is to ensure authenticity of data, so there have been several studies on truth data discovery [32], [33], [34]. These studies are divided into 2 main categories. One category is the computational approach to obtain truthful data from the obtained data, which is mainly a kind of method based on mathematical calculations such as Mean, Median, Weighted average, et al [23], [30]. The main feature of this kind of method is to recruit n workers to sense the same sensing object at the same time. Then, the Mean, Median, and Weighted average of the data of these n workers are estimated truth data [23], [30]. The basic idea behind these approaches is that most of the workers in the MCS are trustworthy and the reported data are true, thus if the n data are averaged, the effect of false data can be reduced [23], [30]. It can be seen that in this type of method, the truthful data is not known, and the accuracy of the result is unclear too. Also, these methods are costly and recruit n workers to do a sensing task, which increases the cost by a factor of n . This type of method, truth discovery, does not identify the credibility of workers but uses the majority rule to obtain values close to the truth data. It is a type of method that collects data and then computes them to obtain the truth data. Another type of approach is to evaluate the credibility of workers before collecting data, and then, recruit truthful workers to sense data [34], [35]. Since credibility is a stable property of workers, the data submitted by truthful workers is more authentic [34], [35]. This approach is theoretically efficient as it requires only 1 worker per task to be recruited. Moreover, if the platform is accurate in recognizing the trustworthiness of workers, it can achieve a very high degree of accuracy in obtaining the truth data. Therefore, the difficulty of this approach is to identify the credibility of workers. Although there are many studies on workers recruitment approach to obtain high-quality data. However, this paper concludes that the current research is still in three areas: truthful workers discovery, sensing difference discovery, and task assignment which deserve further research.

(1) Truthful workers discovery. From the above, we can see that it is a good approach to identify the credibility of workers, and then, select truthful workers to do contact-free sensing data tasks. The key to this approach is how to identify the credibility of workers. In such an approach, it is assumed that the system has identified some truthful workers [10], and then, the

collected data is clustered, and those workers that are in the same class as the truthful workers are credible, while those that are far from the class are not trusted. However, there are several problems here. One is that this clustering method cannot be used in many cases because sensing data has a strong correlation with time, location, and workers' sensing ability. Therefore, there are serious shortcomings in the method of classifying data from different times and locations by mixing them together. For example, in the task of sensing natural phenomena such as temperature and humidity, worker 1 is in the hot spring area (called the hot zone), so its temperature and humidity are very high, and its truth data value is 50, 95. Malicious worker 2 is also in the hot zone, and malicious worker 3 is in the normal zone. The truth data value of the temperature and humidity in the normal zone is 8, 35. Worker 1 truthfully reports his or her sensing data, its value is 51, 96. while worker 2 reports false data, its value is 7.5, 34, and the reported value of worker 3 is 15, 50. Because the number of known truthful workers is small and the majority of the areas are in the normal zone, the known truthful workers in the system are all in the normal zone. Let the data obtained from these known truthful workers which are the same as truthful data of the normal zone be called Ground Truth Data (GTD). After using the classification method, worker 1 is classified as an untrustworthy worker because its value is very different from the GTD. Malicious worker 2 reports false data that deviates significantly from the actual one, but its value is very close to the GTD value and is classified as trustworthy instead. Malicious worker 3 reports data with a much smaller distance from GTD than worker 1. Therefore, it seems that worker 3 is more credible than worker 1. Thus, classification without considering the temporal and spatial properties of the data will give wrong results [36], [37], [38]. Second, this kind of method is based on the assumption that the data collected by a fraction of truthful workers act as GTD [28], [33], [34]. In this way, the data collected by workers can be compared with GTD to obtain workers' credibility. In practice, it is a great challenge to obtain these initial truthful workers, and there is no good solution yet. In addition, even if there are initially some truthful workers, we cannot assume that the credibility of workers will remain unchanged and their data always act as GTD. In fact, the credibility of workers changes frequently, which will lead to the problem of how to continuously ensure that GTD is available after the system runs for a while. There is a fundamental problem of truth discovery, that is, how to obtain and maintain Ground Truth Data (GTD) which has not been well addressed in previous studies.

(2) Sensing difference discovery: In the real-world distributed sensing task of Mobile Crowding Sensing, the sensing task is distributed in different locations, and has multi-attributes such as sound, image, et al [9]. Each attribute has a different metric to represent the quality of sensing data. For workers, since there is a difference in the ability of workers to acquire data from different locations and different attributes, that is, workers' performance of contact-free sensing tasks on different attributes and locations varies a lot [9]. For example, since the quality of sensing different attributes of data fluctuates greatly among cell phones from different manufacturers, e.g., the cell phone from company A is good at capturing images, especially night scenes, while the cell phone from company B is particularly good at sensing sound. Since workers have cell

phones from different manufacturers, their sensing quality varies in different attributes [9]. Some workers are equipped with particularly accurate instruments for observing specific physical phenomena so that they have high sensing quality, while most workers are equipped with only generic sensing devices and have average sensing quality. This kind of variability is due to the different sensing devices held by the workers. In addition, different workers have different skill levels in contact-free sensing data, which reflects the differences in the sensed data. For example, younger workers are better at taking pictures and thus submit data with a higher quality of image attributes, while older workers may submit data with lower quality of image attributes. Some workers are good at taking landscape photos with water and mountains, so they have high sensing quality in locations with mountains and water. Some workers are especially good at taking photos indoors, while their sensing quality is low in outdoor locations. It can be seen that even if the workers are credible, the quality of workers' sensing data shows great variability in different attributes and locations [9]. Therefore, to improve the data quality of the MCS system, we should identify the difference in workers' sensing quality of different attributes and locations. However, the phenomenon of workers sensing difference has not been well investigated in previous studies which lead to poor data collection quality.

(3) Task assignment: In the field of MCS, many task assignment schemes have been proposed [15], [16], [25]. The basic idea is to assign one task to multiple workers and then aggregate their answers using mechanisms like majority voting [30]. However, that kind of method usually costs a lot which wastes valuable resources, and does harm to the environment. A task conducted by several workers consumes several times more energy than that of a task conducted only by only one truthful worker. With the development of smart cities and the increment of smart city applications, it is important to hire the fewest workers for the highest quality tasks to realize low-power sensing. Even if the sensing difference of workers has been calculated, there is no guarantee that all workers will apply to the platform to do sensing tasks. Some workers who can sense high-quality may be not available from time to time [14]. Besides, the location-based character makes the problem that workers are restricted to tasks' spatial constraints [14], so workers are normally only willing to do specific tasks. For these reasons, the task assignment has to be dynamically well-designed. The emergence of Artificial Intelligence (AI) paves the evolution of the IoT into the AI of Things (AIoT) which brings a chance for the traditional Mobile Crowding Sensing to make the task assignment more intelligent. Nevertheless, this kind of study has not been fully studied.

In summary, the key elements to improving data quality are the following three points: (1) Truthful workers discovery, the platform needs to discover truthful workers and then select truthful workers for data collection. (2) The platform not only needs to select truthful workers for data collection but also needs to perform sensing difference discovery to find the differences of workers on sensing tasks in terms of different locations and attributes. (3) Finally, based on this difference, the preferred attributes of the sensing task should also be well-matched with workers who have a high sensing preference for this attribute to obtain high-quality data. However, truth data or

truthful workers discovery is already a very difficult problem. It is even more challenging to discover the preferences of different workers for locations and attributes at a finer granularity.

So, in this paper, an effective workers recruitment scheme based on swarm intelligence for Mobile Crowding Sensing, Multi-attribute and Local Matching based Workers Recruitment (MLM-WR), is proposed to improve the data quality and save energy consumption. The main innovations of this paper are as follows.

(1) We establish a truthful workers discovery approach to evaluate the credibility of workers by comparing their report data with the Ground Truth Data (GTD) and Sub-GTD. First, the Unmanned Aerial Vehicle (UAV) is dispatched and a small amount of data is collected as GTD. Then, we propose a complete calculation of workers' credibility method. In the proposed approach, the credibility of workers is continuously adjusted based on the results of the comparison of their submitted data with GTD. Through such comparison, the credibility of truthful workers increases. When the credibility rises to a certain threshold, the workers are considered truthful, at that time, their reported data can act as Sub-GTD so to speed up the credibility identification of workers.

(2) We propose a sensing difference discovery approach, in which the variability of the workers' sensing quality for different attributes of data, and different locations are obtained by comparing Attribute-GTD (A-GTD) and Location-GTD (L-GTD). A-GTD and L-GTD refer to the data with high accuracy on an attribute or location which mainly come from the UAV respectively. Then, A-GTD and L-GTD are used as criteria to compare the accuracy of workers on the attribute and location dimensions. Based on the results of the comparison, the workers' preferences in each dimension are adjusted, and the data submitted by those workers who have reached a threshold in attributes and locations are used as the sub-A-GTD and sub-L-GTD of the attributes and locations. The location preference is divided into absolute location preference and relative location preference to adapt to a wider range of scenarios so that more workers' attributes and location credibility can be calculated and inferred.

(3) By utilizing swarm intelligence, we convert the MCS from the traditional IoT system into the AIoT system. Combined truthful workers discovery and sensing difference discovery, an effective workers recruitment scheme based on swarm intelligence for Mobile Crowding Sensing, MLM-WR, is proposed to select truthful workers matching sensing data tasks in terms of attributes and locations, which optimizes the quality of sensing data and saves costs. We finally carry out extensive evaluations, where results demonstrate the superiority of our approach over the state-of-the-art approaches. In terms of workers' recognition rate, data quality improvement, and cost-saving, we have improved by 36.12%, 27.10%, and 40.19% respectively.

The rest of this paper is organized as follows: The related works are introduced in Section 2. In Section 3, the system model and problem statement are presented. The design of the MLM-WR approach is proposed in Section 4. Then, the performance analysis of the MLM-WR approach is presented in Section 5. Finally, section 6 provides conclusions.

2. System Model

2.1 Network Model

Truth discovery is an important research issue in MCS which is crucial for ensuring the quality of applications and has attracted the attention of many researchers [19], [29]-[34]. It is also often studied in conjunction with security [39], [40], [41], privacy, etc. In this paper, we focus on a pure truth data discovery study to enable the platform of MCS to obtain authentic, high-quality data [9], [10], [23]. In this paper, we divide the current research on truth data discovery into two categories. One is the research that does not need to get ground truth data for truth data discovery [23], [28], [30], [33]. One is truth data discovery studies that require comparison with ground truth data (GTD) [9], [34], in which GTD is obtained and then, the data reported by workers is compared with GTD to determine whether the data is authentic and whether the workers are truthful.

The first category of studies that perform truth data discovery without obtaining GTD is [23], [28], [30], [33]. The main methods for this type of computed estimate truth data are the Mean, Median, Weighted mean, and Major voting methods [23], [30]. Each of them is described below. (1) First, the Mean method is a simple and easy-to-understand method [30]. In this approach, the basic idea of finding the estimate truth data is that most of the workers in the network are truthful and only a few of them are malicious, so if n workers are recruited to execute the same task and report n data [30]. These n data are then averaged, and since the proportion of false data among these n data is relatively small under normal circumstances, the average value is more accurate than false data. And the larger the n , the closer the average value is to the truth data [30]. However, if multiple malicious workers jointly attack, there will be more malicious data in these n data, and thus the average value is far from the truth data. (2) Median method. The method is to take the median from n data as the truth data [30]. This method is potentially more accurate than the Mean method. Because, in the Mean method, if the value of false data is very different from the truth data. Then, only one false data can make the estimate truth data obtained by the mean method to be far from the truth data. The Median method can avoid this malicious situation because if the false data only works when it is far from the truth data, the probability that the false data will be at the two ends of the truth data and become the median among the n data is low. So, the probability that the estimate truth data is false data in this method is also very low [30]. (3) Weighted average method. The idea of the Weighted average method proposed by Sun et al [23] is as follows: Generally, the data observed on the same object obey normal distribution, and the data located in the center of the normal distribution are closer to the true value, so these data should be given a large weight, while the data farther from the normal center are farther from the truth data, so they are given a smaller weight. Finally, these data are weighted and averaged to obtain the best truth data. This method is a combination of the Mean method and the Median method [23]. The weighted sum reflects the characteristics of the Mean method so that each data contributes to the objective truth data. The large weight given to the data at the center of the normal distribution reflects the greater weight given to the median. However, the shortcoming of this approach is that it assumes

that the data are served from a normal distribution, which is valid when n is large, but in the MCS, the value of n can be not large, and theoretically only 1 truth data is sufficient, and it costs n times as much to obtain n data. In the case where n is small, if there is false data, or if malicious workers join together and dominate, the distribution is either not normal, or the distribution of the service is centered on false data. In this case, the optimal truth data is wrong. (4) Major voting. Major voting has two types of methods [30], one is the majority principle, i.e., the majority of the n data is the estimate truth data. This is especially applicable when the values of the observed objects are discrete. For example, to determine whether there is a specified person in the graph, or to determine which type of flower is in the graph. This method is also subject to malicious worker joint attacks. Another type of method used in Major voting is to ask experts to evaluate the data, the data considered to be true by most experts are used as estimate truth data. This type of method is more expensive, and the speed of identifying data is very slow, which makes it difficult to be suitable for large-scale data collection like MCS systems. Moreover, it is difficult for experts to achieve fine-grained identification [30]. Ye et al [10] proposed a method called the Mean and Median Check (MMC) method. The MMC method is an iterative truth-finding method that combines truth detection, removal of false data, truth detection, and removal of false data again [10]. They first show that if there are data in a dataset that are far from the mean and median values, then there must be incorrect data in the dataset. Thus, the key is to remove false data from the dataset [10]. In the MMC method, the data set is first averaged, then the data that deviate far from the estimate truth data are excluded, then the remaining data are averaged, and then the data that deviate far from estimate truth data is removed. The above process is repeated and the final data is the truth data [10].

In the above methods, we do not check the credibility of workers. We calculate the truth value from n data directly after acquiring the data set. In fact, this kind of method can also be done in conjunction with the credibility of workers, and the result will be better. The specific method is: The higher the credibility of workers, the closer the report data is to the truth data because credibility represents a stable characteristic of workers. Therefore, under such circumstances, if a worker reports with truth data, the credibility will be increased, and if not, the credibility will be decreased. In this way, only workers with high credibility can be selected to report data. This reduces the number of workers needed to report data, thus reducing the cost and making the obtained data more accurate.

The above methods for truth data discovery without GTD are a class of unsupervised methods that have been extensively studied [23], [28], [30], [33]. The advantage of these methods is that they do not require GTD because obtaining GTD is a labor-intensive task and is difficult due to the temporal and spatial correlation of the data [42], [43]. However, these methods require collecting n data for the same object. The cost of this is n times. Thus, if the number of acquired GTD is small, especially with the current development of UAV, which makes it easier and cheaper to obtain GTD, the cost is much smaller than that of these unsupervised methods which require n times the cost. And more importantly, the unsupervised method without GTD is very vulnerable to joint attacks by malicious workers. Since there is no GTD for comparison, the

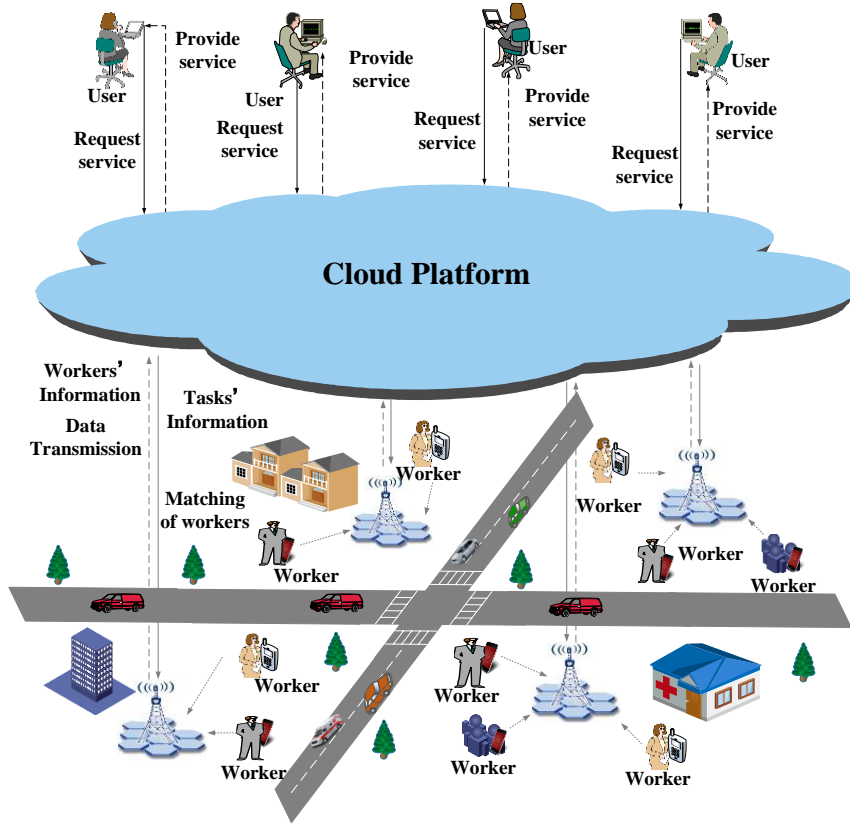


Fig. 1. The MCS system framework

system does not know whether the obtained data is truth data or not, and the accuracy is unstable, which makes it difficult to be applied to applications that have strict requirements.

A class of studies that requires obtaining GTD to perform truth data discovery [23], [28], [30], [33]. Du et al. argue that there are M tasks, among which there are M_g tasks which are golden tasks whose truths are already known (i.e., GTD), and M_0 tasks whose truth data are not known and are needed for discovery [34]. Then, some of the workers execute golden tasks and report their sensing data. The entity is the physical object to be sensed [34]. The main idea of Du et al [34] is that different workers have different accuracy in sensing data for different entities, e.g., some workers may be good at identifying some types of entities but worse at others. Therefore, if the entities are divided into classes, then the workers are classified based on the accuracy of the recognition of these classes. Then, combining the clustering results of entities and workers, a co-clustering reliability matrix representing the accuracy of different classes of workers for different classes of entities is obtained [34]. Finally, based on the co-clustering reliability matrix, good results are achieved by selecting workers with high accuracy on the entities.

The work of Tian et al [9] has similarities with the above work which is also the most similar work to this paper. Du et al [34] argue that there is a difference in sensing accuracy between different workers sensing different entities, so it is necessary to match the workers with the entities to make the sensing accuracy high. The work of Tian et al [9] suggests that the sensing data of the entities have multidimensional attributes and different workers have different sensing accuracy for different attributes, therefore, it is necessary to match the appropriate

workers to sense different tasks. However, in their study, their method only performs unsupervised computation [9] to calculate the relative sensing accuracy of different workers for different attributes. GTD based truth data discovery has an accurate GTD to compare for truth data identification, thus, its accuracy is guaranteed. However, in practice, it is very difficult to obtain ground truth data. The difficulties exist in the following aspects: (1) GTD is difficult to obtain because most of the data in MCS are spatial-temporal related, and the data are even difficult to be measured repeatedly. For example, the simplest observation is the temperature. The temperature varies from region to region and varies with time. Therefore, its value can only be measured repeatedly within a certain period. When the time and location change, the value cannot be reproduced. (2) The cost of collecting GTD is high. The cost of collecting GTD is higher than that of recruiting workers, which requires sending truthful instruments or workers to collect GTD. And the amount of collection should not be too large. Otherwise, the cost is higher. However, in many studies, the question of how to obtain GTD is avoided, and such studies seem to produce beautiful results, but it is difficult to apply them in practice because they avoid the crucial issue of GTD acquisition. Some studies assume the existence of GTD, and some studies assume that there are some known credible workers so that the data collected by these credible workers can act as GTD. However, there are problems with this assumption. One is that the credible workers are either obtained by other external means or compared with GTD to determine whether they are credible. In many cases, it is not feasible to obtain them through external means, and if the basis of the study is based on external forces, it is not applicable under many circumstances. If GTD is used

to identify credible workers, it comes back to the question of how to obtain GTD. Second, even if some workers are initially known to be truthful, the credibility of workers is not constant. It will also change with time and scenario. Therefore, this part of truthful workers can only last for a while, which is far from enough for such long-running MCS applications. The key to the problem is back to the starting point, how to identify and maintain a certain number of truthful workers for truth discovery. On the critical issue of effective GTD acquisition, we also propose an effective method which can be found in Ref. [35]. Our approach is to send UAVs to the sampling point to sample data as truth data [35], then the data reported by the workers are checked for comparison, if the data reported by the workers are consistent with the data reported by the UAV, then the data reported by the workers are considered to be truthful. Thus the method of truthful data discovery and detection based on truthful data is first proposed.

3. System Model and Problem Statement

3.1 Network Model

The network model used in this paper is similar to most of the studies that can be found in Ref. [9, 34]. Mobile Crowd Sensing generally consists of 3 components:

(1) Workers, i.e., various types of sensing devices that can sense the surrounding environment, these sensing devices are mainly mobile sensing devices on the vehicle devices, and various sensor nodes deployed in the network, but the largest and most abundant of them are various types of handheld mobile phones. The number of cell phones has exceeded 5 billion and is growing by hundreds of thousands every year, and their functions are very rich, i.e., they can sense traditional physical phenomena such as temperature, humidity, location, etc. They can also sense multimedia data such as images, sounds, videos, etc.

(2) Cloud platform. Cloud platform issues sensing tasks, which include: the type of data to be sensed, the time, the location, the quality of the sensing requirements, and the reward for data samples. After the workers have been informed of the sensing tasks, they decide whether to participate in the tasks or not according to their situations. Then the platform chooses workers from the participants and assigns them to do sensing tasks. The assigned workers report their sensing data to the cloud platform. Once the cloud platform has received enough data, it processes the data and constructs applications. A wide variety of applications has emerged in recent years, including CrowdAtlas [7], CrowdPark [8], and Parknet [9].

(3) Users. Users request services such as querying location services, querying noise distribution, etc. from the applications platform.

As shown in Fig. 1, the interaction between these three components in the Mobile Crowd Sensing system is as follows: Applications platform issues sensing tasks at some reward for reporting data samples. Then workers apply to the platform to do sensing tasks. The platform matches workers with tasks according to the tasks' information and workers' information. The assigned workers sense data and report data to the platform to get rewards to compensate for their effort in sensing data. After the cloud platform has obtained the data, it processes and computes the data and constructs an application to provide

services to users. The users get the service by paying a fee to the application platform.

3.2 Problem Statement

The main problem in this paper is to design a Multi-attribute and Local Matching based Workers Recruitment (MLM-WR) scheme to improve the data quality at a low cost for the MCS system. The following issues need to be addressed in the design of the MLM-WR approach.

(1) The first goal is to recruit truthful workers to report data to ensure that high-quality data for the MCS system is obtained. To characterize the ability of the proposed method to identify truthful workers, this paper defines the proportion of truthful workers among recruited workers as \mathcal{P}_c to reflect the accuracy of the proposed method to identify truthful workers.

$$\mathcal{P}_c = \frac{w_c}{w_{total}}. \quad (1)$$

Where w_{total} denotes the total number of recruited workers while the proposed method is being used, and w_c denotes the number of truthful workers among them. It can be seen that if \mathcal{P}_c is close to 1, which means that the method can effectively identify the truthful workers and the ratio of truthful workers among selected workers is high. Thus, the objective of the proposed strategy is $\max(\mathcal{P}_c)$.

t_p is the number of workers predicted to be truthful and actually truthful, f_n is the number of workers predicted not to be truthful and actually truthful, f_p is the number of workers predicted to be truthful and actually not truthful.

$$Precision = \frac{t_p}{t_p + f_p}. \quad (2)$$

$$Recall = \frac{t_p}{t_p + f_n}. \quad (3)$$

$$F1 - score = 2 * \frac{Precision * Recall}{Precision + Recall}. \quad (4)$$

The F1-score is also used to evaluate the algorithm.

(2) The second research goal is to obtain high-quality multidimensional data. The data discussed in this paper is multi-dimensional attribute data, and the one-dimensional attribute data can be considered a special case of multi-dimensional attribute data. Although the platform tries to recruit truthful workers, the quality of data collected by even truthful workers for different attributes and different locations can vary greatly. Therefore, the quality of the data collected by the platform can be optimized when the workers are well matched with the collected data and locations.

$Q_k^{(i)}$ is used to denote the i -th dimensional true quality of the k -th data collection task. $Q_{k,j}^{(i)}$ denotes the i -th dimensional data quality of the k -th task collected by the j -th worker. Thus, the difference between the data quality reported by worker j for the k -th task and the true quality of the data in the i -th dimension is shown below.

$$\mathfrak{Z}_{k,j}^{(i)} = \frac{|Q_k^{(i)} - Q_{k,j}^{(i)}|}{Q_k^{(i)}}. \quad (5)$$

And the difference between the data quality reported by workers j and the truthful data quality of task k in all dimensions is :

$$\mathfrak{Z}_{k,j} = \sum_{i=1}^n \mathfrak{Z}_{k,j}^{(i)}. \quad (6)$$

For task k , the platform recruits s workers to collect data,

Table 1 Parameter Description

Parameter	Meaning
$\mathcal{D}_{m,n}^{t,k}$	The observation of the n -th attribute for the m -th location offered by the worker k at time t
$U_{m,n}^t$	The observation of the n -th attribute for the m -th location offered by UAV at time t
$Z_{k,n}^t$	The sensing attribute preference of n -th attribute of worker k at time t
$H_{k,m}^t$	The absolute sensing location preference of m -th location of worker k at time t
$b_{k,m}^t$	The relative sensing location preference of m -th location of worker k at time t
$C_{i,t}$	The credibility of worker i at time t
W_i	The worker i in the network
E_m	The number of attributes that the platform focuses on at location m
R^t	The set of Autonomy Supervised Sensing Location Preference Sequences at time t
L^t	The set of Absolute Sensing Location Preference Sequence at time t
G^t	The set of sensing location preference partial order relationship at time t
G_m^t	The sensing location preference partial order relationship of all workers on location m at time t .
W	The set of all workers' numbers in the network
Z	The set of all sets of sensing preferences of workers in the network for different attributes respectively at time t
C^t	The credibility set at time t
Z_n^t	The set of sensing preferences of workers in the network for attribute n at time t
$a_{k,n}^t$	The sensing attribute credibility of worker k for attribute n at time t .
$e_{k,m}^t$	The sensing location credibility of worker k at moment t for location m

thus, the average data quality difference of the data collected by the platform in task k is:

$$\tilde{\mathcal{S}}_k = \sum_{i=1}^s \mathcal{S}_{k,i} / s. \quad (7)$$

Let there be m tasks in the platform, thus the total average data quality difference of the platform is as follows:

$$\tilde{\mathcal{S}} = \sum_{k=1}^m \tilde{\mathcal{S}}_k / m. \quad (8)$$

Obviously, the higher the quality of the data collected by workers who are recruited by the platform, the smaller the difference between the data collected by the workers and the truthful data, i.e., $\min(\tilde{\mathcal{S}})$.

(3) In the MCS system, there is often a cost to the system to obtain high-quality data, so the proposed strategy is best to optimize the performance of the system with minimum cost. For example, in some schemes, redundant data is collected to make the data close to the truthful data, and there is a cost for collecting redundant data. Let the average cost of the system for workers to submit one data is c . Thus, theoretically, the cost of collecting m data is mc . However, in some strategies, the average redundant data per packet is k . The total cost is kmc . In this paper, the cost of the system has two parts: one is the cost of the UAV to collect data. The other part is the cost of recruiting workers. Thus, the total cost is:

$$\mathcal{C} = \mathcal{C}_{UAV} + \mathcal{C}_w. \quad (9)$$

Where \mathcal{C}_{UAV} is the cost of the UAV, \mathcal{C}_w is the cost of recruiting workers. Therefore, the goal of the strategy designed in terms of system cost is to reduce the cost, i.e. $\min(\mathcal{C})$.

In conclusion, the problem to be solved is the following three objectives: First, the strategy has the ability to effectively identify truthful workers, so as to select truthful workers for data collection and reduce the possibility of being attacked by malicious nodes; Second, the collected data is of high quality; Third, the cost of the strategy is low. The objective of the research problem can be expressed as follows:

$$\begin{cases} \text{Max}(\mathcal{P}_c) = \text{Max}(w_c / w_{total}); \\ \text{Min}(\tilde{\mathcal{S}}) = \text{Min}\left(\sum_{k=1}^m \tilde{\mathcal{S}}_k / m\right); \\ \text{Min}(\mathcal{C}) = \text{Min}(\mathcal{C}_{UAV} + \mathcal{C}_w). \end{cases} \quad (10)$$

To facilitate the reader's understanding of this paper, Table 1 gives a list of some of the main parameters used in this paper.

4. Our Proposed MLM-WR Scheme

4.1 Research Motivation

The research motivation of the Multi-attribute and Local Matching based Workers Recruitment (MLM-WR) scheme is as follows. The MLM-WR scheme consists of three main components: (1) Identification and calculation of workers' credibility. (2) Discovery and computation of workers' preferences on sensing data attributes. (3) Preference discovery and computation of workers' preferences over sensing locations. The first one is also known as truthful workers discovery, while the second and the third one are both sensing difference discovery.

First of all, the truthful workers discovery. As mentioned earlier, it costs workers to sense data so some low-credible, malicious workers submit false data or malicious data and cause damage to the platform. There are two main types of methods to obtain truthful data, one is to recruit multiple workers to obtain the truthful value by using methods like Average, Median, and Weighted average. The other category is to select truthful workers to sense the data because the data sensed by the truthful workers is true. However, in the first category, the platform is uncertain whether it has obtained truthful data because of the lack of ground truth data. The key issue in the approach about truthful workers is how to determine the credibility of the workers which has not been well studied in the past.

In this paper, we propose a method to obtain the credibility of workers by comparing deterministic ground truth data. The main points of this method are as follows: the UAV has become very common and low-cost nowadays, and can quickly acquire data from sensing locations, so in this paper, the platform sends a UAV to collect data from multiple locations. For example, UAV collects data from z locations, $\ell = \{l_1, l_2, l_3, \dots, l_z\}$. Then, take the GTD collected by UAV as the starting point. Let the data collected by the UAV at the location l_i be $d_{u,i}$, and the set of data collected by the UAV be $\mathcal{U} = \{d_{u,1}, d_{u,2}, d_{u,3}, \dots, d_{u,z}\}$. Meanwhile, the system recruits a total of x_i workers to sense the data at the location l_i , and the set of its workers is $S_i = \{w_{i,1}, w_{i,2}, w_{i,3}, \dots, w_{i,x_i}\}$. The set of these workers' report data is $D_i = \{d_{i,1}, d_{i,2}, d_{i,3}, \dots, d_{i,x_i}\}$. Specifically as shown in Fig. 2, the UAV collects data at seven

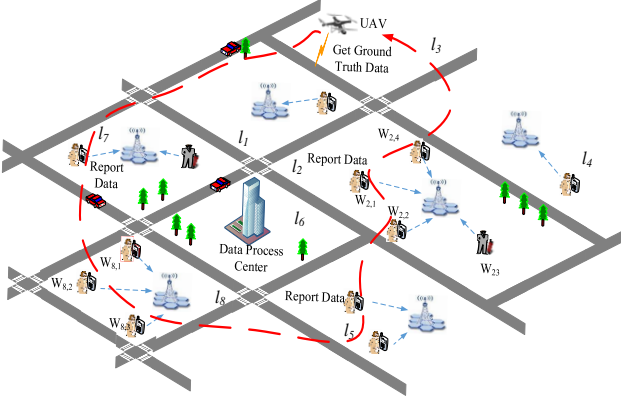


Fig. 2. The illustration of the MLM-WR scheme

locations, numbered $\{l_1, l_2, l_3, \dots, l_7\}$. At location l_2 the system recruited 4 workers: $w_{2,1}, w_{2,2}, w_{2,3}, w_{2,4}$ to report the data. Since the data collected by UAV are true and act as GTD, the data submitted by workers are compared with the data $d_{u,2}$ from UAV, and if the error is within the allowed error, then the data submitted by workers is considered to be true. On the contrary, if the error is beyond the allowed error, the data submitted by workers is considered to be untrue. Then, the credibility is increased or decreased based on the truthfulness of the data submitted by the workers, and the credibility evaluation is given to the workers. In this way, it is possible to evaluate $\sum_{i=1}^N x_i$ workers in one UAV data collection process.

However, there are shortcomings in only comparing the data collected by UAV: First, the number of workers that can be verified is very small compared to the number of workers recruited by the whole system. For example, in Fig. 2, there are only 7 locations for UAV data collection, and if the average number of workers per location is 5, the number of data that can be verified in one UAV data verification process is 35. Compared with the large-scale data collection of tens of thousands or millions, the proportion is very low and the proportion of workers that can be verified is low, which makes it difficult to obtain a sufficient number of truthful workers to collect data. Moreover, since the speed of identifying the credibility of workers by this method is very low, and if the workers in the network are dynamic, the identification speed may not be able to keep up with the speed of change of workers, which makes this method fail. In addition, among the above methods, GTD collected by UAV is costly and difficult to be extended. To address the above situation, this paper proposes a method to expand the scope of comparing to extend the calculation of the credibility of workers, then the speed of identifying truthful workers will be faster and the number of identified truthful workers will be enlarged. The approach is that workers will be identified as truthful by the above UAV method if the credibility of the workers exceeds a high threshold. Since the data submitted by the truthful workers are considered to be true.

In this paper, these truthful workers' data act as Sub-GTD to check the credibility of other workers. As mentioned above, let n identified truthful workers go to n locations. Let a total of x_j workers at location l_j have submitted data, whose set is $S_j = \{w_{j,1}, w_{j,2}, w_{j,3}, \dots, w_{j,x_j}\}$. The set of these workers' report data is $D_j = \{d_{j,1}, d_{j,2}, d_{j,3}, \dots, d_{j,x_j}\}$. There is one truthful worker $w_{8,1}$ in location l_8 as shown in Fig. 2, and there are

also workers $w_{8,2}, w_{8,3}$ submitting data on location l_8 . Thus, the data submitted by the worker $w_{8,1}$ is used as a Sub-GTD to check the credibility of workers

$w_{8,2}, w_{8,3}$. It is possible to check the credibility of $\sum_{j=1}^k x_j$ workers. Obviously, this approach is able to scale up the identification of truthful workers and accelerate the speed of the identification of truthful workers.

This is illustrated below by giving a concrete example. Let the UAV collect data from m locations, and each location has an average of y workers so that the UAV can detect my workers at a time. Assume that after k rounds of testing, ϖmy truthful workers can be detected in each subsequent round, and once ϖmy truthful workers are available, the data collected by these truthful workers will be used as Sub-GTD in credibility evaluation. Let the probability that a truthful worker is selected by the system be λ_c . Thus, if the current number of truthful workers is N , the number of selected truthful workers is $\lambda_c N$. And $\lambda_c N$ truthful workers travel to $\lambda_c N$ locations, $\lambda_c N$ truthful workers can detect $\lambda_c N y$ workers' credibility. And also after k rounds, $\varpi \lambda_c N y$ truthful workers are obtained. And in each subsequent round, as the number of truthful workers increases from N to $N + \varpi \lambda_c N y$, there will be a faster increase in the number of truthful workers that can be detected later. And the faster increase in the number of truthful workers will evaluate more workers. This leads to a snowball-like cumulative effect of the detected truthful workers.

The above part mainly identifies the credibility of workers, and truthful workers mainly guarantee that the data they submit are their own sensed data and not false or malicious data. However, truthful workers do not necessarily guarantee high sensing quality. The main reason is that different workers have preferences in terms of sensing attributes and sensing locations. The sensing attribute preference means that workers show sensing differences in different attributes, i.e., workers show good sensing quality in some attributes and mediocre sensing quality in other attributes. Different data have different attributes (or have the same attributes, but the application values certain attributes). Therefore, the platform should be able to identify the preference of the workers in terms of the quality of the collected attributes, and then select the workers who happen to have good sensing quality in attributes that applications and data require to sense data so that the collected data can be matched with the application requirements to achieve good results. But this issue has not been well studied in the past. There is no good method for this issue in previous research. This paper proposes an effective approach to discover the preference of workers' attributes. The specific method is divided into two steps: Step 1 needs to identify which workers have high sensing quality in which data attributes; Step 2 selects the workers with high sensing quality in demand attributes to sense the data. For the first step that identifying which workers have high sensing quality on which attributes, this paper adopts the following approach: First, the data quality collected by UAV is the highest, that is, the data quality collected by UAV is high in all dimensions of data attributes. Thus, it can be used as a criterion to check the quality of workers' data on attributes. If the difference in quality between workers and UAV in dimension i is less than a certain error, the workers can be considered to have a very high sensing quality in dimension i . On the contrary, workers' sensing quality in dimension i is low.

The calculation of sensing preference on attributes is somewhat similar to the previous method of workers' credibility calculation. In the beginning, workers' sensing preference for each dimension is an initial value. Then, as workers are compared with the UAV data, the sensing preference for each attribute is obtained similarly to the previous calculation of credibility. Again after some time, if the preference of the worker for attribute i is greater than a higher threshold value, then this worker can be considered to have a high sensing attribute preference on this attribute. Then, the other workers' sensing attribute preference in attribute i is checked against the workers with high sensing attribute preference in attribute i . In this way, after the above process, the sensing attribute preference of workers can be obtained.

The calculation of workers' sensing location preferences is given below. The calculation of workers' preferences over locations is equivalent to combining the calculation of credibility with the calculation of sensing attribute preference, but with significant differences. First, the data sensed by the UAV have the highest quality at each location. Thus, they can act as a criterion to check the workers' sensing preference for locations. The workers' sensing preference for location l_i is the deviation of workers' submitted data from the UAV's data in attributes of each dimension. The sensing preference obtained above is the exact, absolute sensing location preference. After the above process, if the absolute sensing location preference of the worker k at the location l_i is greater than the threshold value after multiple calculations, the worker is considered to have a better absolute sensing location preference for l_i . So to improve the data quality, when the platform recruits workers to sense the data of a location, it selects the workers with high absolute sensing preferences for the location from those who are willing to participate in sensing data. However, since the UAV data can cover only a small part of the locations, it is necessary to obtain as many workers' sensing preference for locations as possible. If the absolute sensing location preference of a worker for location l_i is greater than a high threshold, it is assumed that the worker can sense the location l_i accurately. Then, the workers with high absolute sensing location preference for the location l_i can check and update the absolute sensing location preference of other workers. However, because of the large number of locations and the large number of total workers who apply to execute sensing tasks, there may be many locations where the UAV or workers with high absolute sensing location preference for these locations cannot be dispatched; Or the UAV has not checked any worker on location l_i , so there is a situation where all workers have not been updated on their absolute sensing preference for location l_i . The workers performing the task at these locations need to be unsupervised to check their sensing preferences for these locations.

To achieve the above goal, we propose the concept of relative sensing location preference. Relative sensing location preference means that we need to compute a relative ranking relationship between workers' sensing quality at the location l_i . Such a relative quality ranking of workers' sensing quality on location l_i was created because the platform only selects the workers who have a high sensing preference for this location relative to other workers from workers who are willing to participate in the perception tasks, it is not necessary to obtain

Table 2 Summary of terminologies

Notation	Description
attribute	An attribute to describe the location, e.g., "Humidity"
location	A place with several attributes. e.g., " l_1 "
entry	An attribute value of a location. e.g., "Humidity of l_1 "

all workers' absolute sensing location preference. Thus, workers can be selected as long as they are ranked in a partial order with respect to sensing preference over this location. Therefore, we adopt the method of calculating the relative partial order: If there are workers who can get absolute sensing preference in the location l_i , the partial order relationship about sensing location preference of location l_i between them can be calculated by comparing their absolute sensing preference for the same location. And as for most workers who have not got absolute sensing location preference, we use the indirect inference method. When there are no workers with a high absolute sensing preference for location l_i sense data in l_i and UAV is also not in the location l_i , the partial order relationship about sensing location preference of l_i of workers in the location l_i are obtained by utilizing an unsupervised algorithm.

Then, these partial order relationships of sensing location preference are combined to form a comprehensive partial order relationship of sensing location preference, which guides the platform in recruiting workers to select workers with good sensing location preference for sensing data. So the platform can obtain high-quality data.

After the above calculation, we can get the credibility of workers, the sensing attribute preference of workers, and the sensing preference of location. With the above three values, the platform can first select truthful workers for perception tasks, thus ensuring that the obtained data is true. In fact, when multiple truthful workers are willing to participate in perception tasks, the platform can select the workers with a high preference for sensing nodes and data, so that it can obtain high-quality data. The following paper discusses how to get the credibility, sensing attribute preference, and sensing location preference of the workers. The way how the platform selects the workers with matching sensing attribute preference and sensing location preference to obtain high-quality data will also be discussed.

4.2 The Design of Truth Finding Mechanism

A. Credibility Evaluation Model in MLM-WR

Because there are some malicious workers in the MCS system who want to get the maximum reward at the least cost, the false data of these malicious workers can cause great loss to the platform, so it is necessary to identify the truthful workers.

Assuming that at time t , UAV collects a total of m data from m locations, each data containing n dimensional attributes. Each worker goes to one location at one time to collect data, and there are g workers collecting data from one of these m locations. Thus, we evaluate the credibility of these g workers based on these m data. For each worker, an initial value of credibility is given at the beginning and compared with the UAV data at each time. Because each time UAV collects data from different locations, the credibility of different workers can be updated each time. The Reward Mechanism is utilized to update the credibility when the calculated error is small, the

Punishment Mechanism is utilized when the calculated error is significant, and the Timeout Mechanism is used when the worker is not checked at this time.

Use U_i^t to denote the data collected by UAV on location i at time t . Because the UAV is sent by the platform, the data collected by the UAV are credible and trustworthy, so the data of the UAV act as GTD. Use $U_{i,j}^t$ to denote the value of the j -th attribute of the data collected by UAV on location i at time t . And $\mathcal{D}_i^{t,k}$ denotes the data collected by the worker k on location i at time t . Use $\mathcal{D}_{i,j}^{t,k}$ to denote the value of attribute i of data collected by the worker k on location i at time t . $S_{i,j}^k$ denotes the ratio of deviation from $U_{i,j}^t$ and $\mathcal{D}_{i,j}^{t,k}$ to $U_{i,j}^t$, and For attribute j of location i , if worker k has collected data at location i and UAV focuses the attribute j of location i , then $S_{i,j}^k$ will be calculated with Formula (11).

$$S_{i,j}^k = \frac{|\mathcal{D}_{i,j}^{t,k} - U_{i,j}^t|}{U_{i,j}^t}. \quad (11)$$

If $S_{i,j}^k < \theta$, let $\mathcal{X}_{i,j}^k = 1$; else let $\mathcal{X}_{i,j}^k = 0$.

As to attribute j that is not of interest to the platform at location i , let $\mathcal{X}_{i,j}^k = 0$. $A_{i,k}$ represents the number of right entries collected by the k -th worker at location i , the calculation of which is as Formula (12).

$$A_{i,k} = \sum_{j=1}^n \mathcal{X}_{i,j}^k. \quad (12)$$

The next time, UAV collects some other data, and this time the number of collected data and the location may not be the same as the last time, and the number of workers compared is also different. As a result, some workers are updated and some are not. The credibility of workers who are not evaluated will also decay with time.

Suppose that worker k is compared with the data of the UAV to evaluate the performance of worker k . E_i represents the total number of attributes of the data of location i that the platform is concerned about, which is calculated by the Formula (13). If the attribute is of concern to the platform, then UAV has collected data for this attribute, let $E_{i,j} = 1$; else let $E_{i,j} = 0$. Interaction right σ_r expresses the effect of the interaction between the worker and the UAV, calculated as the ratio of $A_{i,k}$ and E_i .

$$E_i = \sum_{j=1}^n E_{i,j}. \quad (13)$$

$$\sigma_r = \frac{A_{i,k}}{E_i}. \quad (14)$$

Assuming that at time t , worker 1, worker 2, and worker 3 collect data of l_1 , they collect data of l_2 at time $t+1$ and collect data of l_3 at time $t+2$. The data they collected are shown in Table 3 and the corresponding truth data are shown in Table 4. All of the attributes are concerned with the platform.

For example, Table 3 shows that worker 1 owns $A_{2,1}$ which is calculated as 1 by Equations (11) and (12), when $\theta = 0.001$, E_2 is calculated by Equation (13), σ_r is calculated by Equation (14). So, we have $E_2 = 3$ and $\sigma_r = \frac{2}{3}$:

$$E_2 = 3.$$

$$\sigma_r = \frac{2}{3}.$$

(1) Reward Mechanism

$C_{k,t}$ means the credibility of worker k at time t , and $C_{k,t-1}$ represents the credibility of worker k at time $t-1$.

Table 3 Workers' data of l_1, l_2, l_3

Location	Temperature	Humidity	Noise
(a)Worker1			
l_1	1	4	7
l_2	4	5	8
l_3	2	6	9
(b)Worker2			
l_1	1	4	10
l_2	2	15	8
l_3	3	6	10
(c)Worker3			
l_1	1	4	13
l_2	2	6	8
l_3	3	8	9

Table 4 Ground truth of l_1, l_2, l_3

Location	Temperature	Humidity	Noise
l_1	1	4	7
l_2	2	5	8
l_3	3	6	9

Assuming that there are K workers in the scenario, the set of credibility of K workers at time t is $C^t = \{C_{1,t}, \dots, C_{K,t}\}$. $C_{k,0}$ means the initial value of credibility of worker k . A threshold σ_s is defined to decide whether to perform the Reward Mechanism or Punishment Mechanism. If σ_r is greater than σ_s , then worker k is considered to be truthful and the Reward Mechanism is performed. The credibility will be increased by using the Formula (15). In Formula (15), ρ is a variable greater than 1, and ρ is used to control the speed of updating credibility. The credibility can be increased by a maximum of $\frac{1}{\rho}$ per update (when $\rho = 1$).

$$C_{k,t} = C_{k,t-1} + \frac{\sigma_r}{\rho} * (1 - C_{k,t-1}). \quad (15)$$

Theorem 1: The Reward Mechanism (RM) is used when a worker is evaluated to be truthful, and if the worker is evaluated to be truthful each time, the Reward Mechanism will be run repeatedly. The credibility of the worker that repeatedly applies the Reward Mechanism will gradually converge over time to 1.

Proof: Assuming that all entries of worker i being evaluated at each time are considered to be credible, i.e., $\sigma_r = 1$, the following inequality (16) is obtained. For the last three steps, since $\left(\frac{\rho-1}{\rho}\right)^t$ equals to 0 when t tends to be infinity, we finally obtain that $C_{i,t}$ converges to 1.

$$\begin{aligned} \lim_{t \rightarrow \infty} C_{k,t} &= C_{k,t-1} + \frac{\sigma_r}{\rho} * (1 - C_{k,t-1}) \\ &\leq \lim_{t \rightarrow \infty} \left[C_{k,t-1} + \frac{1}{\rho} * (1 - C_{k,t-1}) \right] \\ &= \lim_{t \rightarrow \infty} \left[\frac{\rho-1}{\rho} C_{k,t-1} + \frac{1}{\rho} \right] \\ &= \lim_{t \rightarrow \infty} \left[\frac{\rho-1}{\rho} \left(\frac{\rho-1}{\rho} C_{k,t-2} + \frac{1}{\rho} \right) + \frac{1}{\rho} \right] \\ &= \lim_{t \rightarrow \infty} \left[\left(\frac{\rho-1}{\rho} \right)^t C_{k,0} + \frac{1}{\rho} \sum_{i=0}^{t-1} \left(\frac{\rho-1}{\rho} \right)^i \right] \\ &= \lim_{t \rightarrow \infty} \left[0 + \frac{1}{\rho} \frac{1 - \left(\frac{\rho-1}{\rho} \right)^t}{1 - \frac{\rho-1}{\rho}} \right] = 1. \end{aligned} \quad (16)$$

According to the proof, we can see that when the historical credibility of workers is higher, the update will be slower, eventually, the credibility will converge. Theorem 1 shows that even in an extreme situation, the workers' credibility will not increase endlessly. The highest value of credibility of workers will be controlled.

(2) Punishment Mechanism

When σ_r is smaller than σ_s , worker k is considered to be incredible, so the credibility of workers k needs to be reduced by using the Formula (17).

$$C_{k,t} = C_{k,t-1} + (\sigma_r - 1) * \left(\frac{C_{k,t-1}}{\rho} \right). \quad (17)$$

Theorem 2: When the data submitted by a worker is determined to be untrustworthy, the Punishment Mechanism (PM) is applied. As time t is updated, the worker that is always been evaluated to be malicious will be continuously applied the Punishment Mechanism, and the credibility of this worker will be reduced to 0.

Proof: Assuming that all the entries given by the worker are judged to be untrustworthy at each time, i.e., $\sigma_r = 0$, the maximum update is performed. As shown in inequality (18), since $\left(\frac{\rho-1}{\rho}\right)^t$ turns out to be 0 as t tends to be infinity, the credibility eventually converges to 0.

$$\begin{aligned} \lim_{t \rightarrow \infty} C_{k,t} &= \lim_{t \rightarrow \infty} \left[C_{k,t-1} + (\sigma_r - 1) * \left(\frac{C_{k,t-1}}{\rho} \right) \right] \\ &\geq \lim_{t \rightarrow \infty} \left[C_{k,t-1} - 1 * \left(\frac{C_{k,t-1}}{\rho} \right) \right] \\ &= \lim_{t \rightarrow \infty} \left[\frac{\rho-1}{\rho} C_{k,t-1} \right] \\ &= \lim_{t \rightarrow \infty} \left[\left(\frac{\rho-1}{\rho} \right)^2 C_{k,t-2} \right] \\ &= \lim_{t \rightarrow \infty} \left[\left(\frac{\rho-1}{\rho} \right)^t C_{i,0} \right] = 0. \end{aligned} \quad (18)$$

Theorem 2 shows that the minimal value of the credibility of workers will be controlled even in an extreme situation, the workers' credibility will not decrease endlessly. According to the proof, we can see that when the σ_r of workers is lower, the update will be faster. To keep high credibility, workers have to keep submitting high-quality data.

(3) Timeout Mechanism

Because some workers may submit accurate data at the beginning, but they become very inactive later on. The value of selecting such workers is very small for us, so we should not only update the credibility of workers who submit data but also those who do not submit data. For a worker who hardly submits data, the credibility of the worker will be continuously reduced by a certain amount. This calculation is shown in Formula (19) and Formula (20). When the data is not submitted at time t , the following process is performed, where γ controls the speed of updating. Compared with the RM and PM, the Timeout Mechanism (TM) has a much smaller change in the credibility update. This is because the data quality is more important than the frequency of submitting data. The damage caused by malicious data is more serious than that caused by non-submission. Therefore, the value of γ is set to be small. When the credibility of a worker is higher, the rate of decline will be slower.

$$C_{k,t} = \text{Max}(C_{\min}, C_{k,t-1} - \partial). \quad (19)$$

$$\partial = \gamma(1 - C_{k,t-1}). \quad (20)$$

After repeating the above process for some time, if the credibility of some workers is greater than ν , then these workers are considered to be credible. If these workers collected data at locations l_1, l_2, l_3 , and other workers also collected data at one of these locations. Then the data of the truthful workers are used as Sub-GTD to update the credibility of other workers by utilizing Vice Reward Mechanism or the Vice Punishment Mechanism. The update process is shown below.

When the data of truthful workers are used as Sub-GTD to evaluate other workers, $\mathcal{D}_{i,j}^{t,b}$ denotes the value of the j -th attribute of truthful worker b on location i at time t . $P_{i,j}^k$ denotes the ratio of deviation between the value of the j -th attribute of worker k on location i at time t and Sub-GTD to Sub-GTD. $A_{i,k}^c$ denotes the number of wrong entries collected by worker k at location i under Sub-GTD evaluation. If the ordinary worker k has collected data at location i and attribute j of location i is concerned by the platform, then $P_{i,j}^k$ will be calculated by Formula (21).

$$P_{i,j}^k = \frac{|\mathcal{D}_{i,j}^{t,k} - \mathcal{D}_{i,j}^{t,b}|}{\mathcal{D}_{i,j}^{t,b}}. \quad (21)$$

If $P_{i,j}^k < \theta^c$, let $y_{i,j}^k = 1$; else let $y_{i,j}^k = 0$.

$A_{i,k}^c$ represents the number of wrong entries collected by the k -th worker at location i when Sub-GTD evaluates worker k , the calculation of which is as Formula (22). And the interaction right σ_r^c expresses the effect of the interaction between the worker and the Sub-UAV, calculated as the ratio of $A_{i,k}^c$ and E_i .

$$A_{i,k}^c = \sum_{j=1}^n y_{i,j}^k. \quad (22)$$

$$\sigma_r^c = \frac{A_{i,k}^c}{E_i}. \quad (23)$$

(4) Vice Reward Mechanism

A threshold σ_s^c is defined to determine whether the Vice Reward Mechanism or the Vice Punishment Mechanism is applied to ordinary workers. When σ_r^c is greater than σ_s^c , it means that this time GTD considers worker k to be credible. The Vice reward mechanism will be used to update the credibility of worker k , which is shown in Formula (24). In Formula (24), $C_{b,t-1}$ is multiplied after the updated term, which shows that when the credibility of the worker whose data act as Sub-GTD is higher, the change of the update is larger, and when $C_{b,t-1}$ equals to 1, which means the worker is completely credible and can be treated as a UAV. Then the mechanism of Sub-GTD is the same as the mechanism of GTD. ρ^c is a variable greater than 1, which is used to control the speed of updating credibility.

$$C_{k,t} = C_{k,t-1} + \frac{\sigma_r^c}{\rho^c} * (1 - C_{k,t-1}) * C_{b,t-1}. \quad (24)$$

(5) Vice Punishment Mechanism

When σ_r^c is smaller than σ_s^c , it means that at this time the Sub-GTD considers worker k to be untrustworthy. Formula (25) is used to reduce the credibility of worker k .

$$C_{k,t} = C_{k,t-1} + (\sigma_r^c - 1) * \left(\frac{C_{k,t-1}}{\rho^c} \right) * C_{b,t-1}. \quad (25)$$

When both Sub-GTD and GTD are involved in the evaluation, the Timeout Mechanism is applied if there are no workers evaluated by sub-GTD or GTD.

When there is no Sub-GTD and only the data collected by

UAV can be used as GTD, the Timeout Mechanism is applied to the workers that are not evaluated by GTD.

B. SAPE in MLM-WR

In the previous section, we gave a solution for the credibility of workers, which indicates whether the workers are credible at the macro level. However, different workers have different performances on multidimensional data in different dimensions. Therefore, this paper proposes the concept of sensing attribute preference and the Sensing Attribute Preference Evaluation Model (SAPE) to make the final collected data more accurate by evaluating the sensing preference of workers for different attributes.

Example 1. According to the data of l_1 collected by three workers shown in Table 3, when $\theta = 0.001$, time t is the first round, $C_{1,0} = C_{2,0} = C_{3,0} = 0.5$, $\rho = 10$, $\sigma_s = 1/2$, the credibility of workers is calculated by Equations (11), (12), (13), (14), (15), (17):

Worker 1: $C_{1,3} = 0.608$;

Worker 2: $C_{2,3} = 0.593$;

Worker 3: $C_{3,3} = 0.593$.

The result shows that the credibility of worker 1 is the highest, but the accuracy of the temperature of worker 1 is the lowest. So, in order to get the truth, the temperature data of worker 1 can not be referred to. $\mathcal{D}_{m,n}^{t,k}$ represents the n -th attribute of the data collected by the i -th worker at time t , at the m -th location. $U_{m,n}^t$ represents the value of the n -th attribute of the m -th location collected by UAV at time t , which is also named Attribute-GTD. Attribute-GTD is used to help establish the sensing attribute preference of workers. Suppose there are K workers in the network, $Z_{k,n}^t$ denotes the sensing preference of worker k for attribute n at time t . $Z_n^t = \{Z_{1,n}^t, Z_{2,n}^t, \dots, Z_{K,n}^t\}$, Z_n^t denotes the set of sensing preferences of workers in the network for attribute n at time t . The evaluation of sensing attribute preference is given below.

If the attribute n of location m is of interest to the task, then use Formula (26) to calculate $Y_{m,n}^k$.

$$Y_{m,n}^k = \frac{|\mathcal{D}_{m,n}^{t,k} - U_{m,n}^t|}{U_{m,n}^t}. \quad (26)$$

If $Y_{m,n}^k < \vartheta$, let $P_n = 1 - Y_{m,n}^k$;
else if $\vartheta \leq Y_{m,n}^k < 1$, $P_n = 1 - Y_{m,n}^k$;
else $1 \leq Y_{m,n}^k$, $P_n = 0$.

(1) Attribute Reward Mechanism

$Z_{k,n}^0$ represents the initial value of the sensing preference of worker k for attribute n . $Z_{k,n}^t$ represents the value of the sensing preference of worker k for attribute n at time t . A threshold Ω is defined to determine whether to perform the Attribute Reward Mechanism (ARM) or the Attribute Punishment Mechanism (APM). If $Y_{m,n}^k$ is less than Ω , then A-GTD considers worker k to be credible for attribute n . So Formula (27) will be used to increase the preference of attribute n . δ is a variable greater than 1. δ is used to control the speed of updating the sensing preference of each attribute. Each sensing attribute preference can be increased by a maximum of $\frac{1}{\delta}$ per the update (when $P_n = 1$).

$$Z_{k,n}^t = Z_{k,n}^{t-1} + \frac{P_n}{\delta} (1 - Z_{k,n}^{t-1}). \quad (27)$$

(2) Attribute Punishment Mechanism

If $\Omega < Y_{m,n}^k$, it means that worker k is not credible for attribute n . It is necessary to use Formula (28) to reduce the

sensing preference of worker k for attribute n .

$$Z_{k,n}^t = Z_{k,n}^{t-1} + (P_n - 1) * \left(\frac{Z_{k,n}^{t-1}}{\delta} \right). \quad (28)$$

(3) Attribute Timeout Mechanism

For inactive workers, the Attribute Timeout Mechanism (ATM) is used to slowly reduce the sensing attribute preference for attribute n of the workers, the calculation of which is as Formula (29) and Formula (30).

$$Z_{k,n}^t = \text{Max}(Z_{\min}, Z_{k,n}^{t-1} - \varphi). \quad (29)$$

$$\varphi = \mu(1 - Z_{k,n}^{t-1}). \quad (30)$$

After repeating the above process for a while, when sensing preference $Z_{k,n}^t$ of worker k for the attribute n is higher than the threshold δ , the worker is considered to have high sensing quality for attribute n . This workers' data of attribute n can act as Sub-Attribute-GTD for attribute n .

If there are multiple workers at location m with high sensing quality for attribute n at time t , and UAV does not go to location m at this time, the data of worker b with the highest sensing preference for attribute n at location m at this time is selected as the Sub-Attribute-GTD of attribute n at location m . Then worker b will evaluate attribute n of other workers at location m , and update the sensing preference of attribute n of other workers at location m by utilizing Vice Attribute Reward Mechanism or Vice Attribute Punishment Mechanism.

$Z_{b,n}^t$ represents the sensing attribute preference of attribute n of worker b who acts as Sub-Attribute-GTD of attribute n at time t , and $\mathcal{D}_{m,n}^{b,t}$ represents the value of attribute n collected by worker b at location m at time t . When worker b collects data at location m , and the attribute n at location m is of interest to the platform, the sensing attribute preference of attribute n of worker k can be updated who also collects data at location m at this time. Then Formula (31) will be used to calculate $q_{m,n}^k$.

$$q_{m,n}^k = \frac{|\mathcal{D}_{m,n}^{t,k} - \mathcal{D}_{m,n}^{b,t}|}{\mathcal{D}_{m,n}^{b,t}}. \quad (31)$$

If $q_{m,n}^k < \Omega^c$, let $P_n^c = 1 - q_{m,n}^k$;
else if $\Omega^c \leq q_{m,n}^k < 1$, $P_n^c = 1 - q_{m,n}^k$;
else $1 \leq q_{m,n}^k$, $P_n^c = 0$.

(4) Vice Attribute Reward Mechanism

A threshold Ω^c is defined to determine whether to perform the Vice Attribute Reward Mechanism or the Vice Attribute Punishment mechanism. When $q_{m,n}^k < \Omega^c$, it means that the attribute n of worker k is considered to be credible by Sub-Attribute-GTD. So the sensing attribute preference of worker k on attribute n will be updated by the Vice Attribute Reward Mechanism, which is calculated as Formula (32), in Formula (32), δ^c is used to control the speed of the update.

$$Z_{k,n}^t = Z_{k,n}^{t-1} + \frac{P_n^c}{\delta^c} * (1 - Z_{k,n}^{t-1}) * Z_{b,n}^t. \quad (32)$$

(5) Vice Attribute Punishment Mechanism

When $q_{m,n}^k > \Omega^c$, it means that the attribute n of worker k is considered to be untrustworthy by the Sub-Attribute-GTD. The minimum value of P_n^c is 0 because it is necessary to ensure that $Z_{k,n}^t$ converges to 0.

Then Vice Attribute Punishment Mechanism (VPRM) will be utilized by using Formula (33) to reduce the sensing

preference of attribute n .

$$Z_{k,n}^t = Z_{k,n}^{t-1} + (P_n^c - 1) * \left(\frac{Z_{k,n}^{t-1}}{\delta^c} \right) * Z_{b,n}^t. \quad (33)$$

When both Sub-A-GTD and A-GTD of attribute n are involved in the evaluation, the Attribute Timeout Mechanism is applied to worker k if the attribute n of worker k has not been evaluated by Sub-A-GTD or A-GTD.

When there is no Sub-A-GTD of attribute n and only A-GTD, the Attribute Timeout Mechanism is applied to worker k if attribute n of worker k has not been evaluated by A-GTD.

C. SLPE in MLM-WR

The above section evaluates the credibility and sensing attribute preference of the worker. Although the change in workers' personal quality and the variation of the reliability of sensing different attributes are considered, it is not taken into account that even if the data collected by the worker at each location contains the same attributes, such as temperature, humidity, etc., the accuracy of the data collected by workers in different locations for the same attributes may vary. Therefore, this paper proposes the concept of sensing location preference and the Sensing Location Preference Evaluation Model (SLPE). Workers need to be evaluated according to locations as well.

Assuming that at time t , worker 4, worker 5, and worker 6 collect data of l_4 , they collect data of l_5 at time $t+1$ and collect data of l_6 at time $t+2$. The data they collected are shown in Table 5 and the corresponding truth data are shown in Table 6. All of the attributes are concerned with the platform.

Example 2. According to the data of l_4 collected by three workers shown in Table 5, when $\theta = 0.001$, time t is the first round, $C_{4,0} = C_{5,0} = C_{6,0} = 0.5$, $\rho = 10$, $\sigma_s = 1/2$, the credibility of workers is calculated by Equations (11), (12), (13), (14), (15), (17):

$$\text{Worker 4: } C_{4,3} = 0.552;$$

$$\text{Worker 5: } C_{5,3} = 0.541;$$

$$\text{Worker 6: } C_{6,3} = 0.541.$$

The result shows that Worker 5, and Worker 6 have the same credibility and Worker 4 has the highest credibility. Although Worker 5, and Worker 6 have the same credibility, Worker 5 has a better sensing ability of l_5 , and Worker 6 has a better sensing ability of l_4 . And Worker 4 cannot sense l_4 well even though the credibility of Worker 4 is the highest. So, it is necessary to evaluate the sensing preference of workers for different locations in detail.

(1) Calculation of absolute sensing location preference

Assuming that there are a total of K workers in the network, and since the UAV is dispatched by the platform, the data collected by the UAV must be truthful. The data collected by UAV is used as Location-GTD (L-GTD) to compare the data submitted by workers who also collect data at the same location, and then the absolute sensing location preference will be updated. $H_{k,m}^t$ denotes the absolute sensing location preference of worker k for location m at time t .

If UAV and worker k collect data at location m at time t , and the attribute n of location m is of interest to the task, $O_{m,n}^k$ will be calculated by Formula (34).

$$O_{m,n}^k = \frac{|\mathcal{D}_{m,n}^{t,k} - U_{m,n}^t|}{U_{m,n}^t}. \quad (34)$$

If $O_{m,n}^k < \aleph$, let $O_{m,n}^k = 1$; else $O_{m,n}^k = 0$.

Table 5 Workers' data of l_4, l_5, l_6

Location	Temperature	Humidity	Noise
(a) Worker 4			
l_4	17	83	27
l_5	16	60	30
l_6	17	70	29
(b) Worker 5			
l_4	15	80	29
l_5	16	60	28
l_6	15	67	29
(c) Worker 6			
l_4	15	80	27
l_5	16	60	30
l_6	20	73	29

Table 6 Ground truth of l_4, l_5, l_6

Location	Temperature	Humidity	Noise
l_4	15	80	27
l_5	16	60	28
l_6	17	70	29

If the task of location m does not focus on the attribute n , let $O_{m,n}^k = 0$.

$F_{m,k}$ represents the number of entries in location m collected by worker k that are evaluated as correct by Location-GTD. $F_{m,k}$ is calculated by the Formula (35). The interaction right η_m expresses the effect of the interaction between the worker and the UAV, calculated as the ratio of $F_{m,k}$ and E_m .

$$F_{m,k} = \sum_{n=1}^N O_{m,n}^k. \quad (35)$$

$$\eta_m = \frac{F_{m,k}}{E_m}. \quad (36)$$

1) Location Reward Mechanism

$H_{k,m}^0$ represents the initial value of the absolute sensing preference of worker k for location m . A threshold η_s is defined to decide whether to perform the Location Reward Mechanism (LRM) or the Location Punishment mechanism (LPM).

If η_m is greater than η_s , then L-GTD considers that worker k is credible in terms of location m . And Formula (37) is used to increase the absolute sensing preference of worker k for location m . In Formula (37), τ is a variable greater than 1, and τ is used to control the speed of updating sensing location preference. The maximum value of increment in each evaluation is $\frac{1}{\tau}$ (When $\eta_m = 1$).

$$H_{k,m}^t = H_{k,m}^{t-1} + \frac{\eta_m}{\tau} * (1 - H_{k,m}^{t-1}). \quad (37)$$

2) Location Punishment Mechanism

When $\eta_m < \eta_s$, it means that worker k is not credible in terms of location m . Then Formula (38) will be utilized to reduce the absolute sensing location preference of worker k for location m .

$$H_{k,m}^t = H_{k,m}^{t-1} + (\eta_m - 1) * \left(\frac{H_{k,m}^{t-1}}{\tau} \right). \quad (38)$$

3) Location Timeout Mechanism

When the data of worker k for location m cannot be evaluated at time t , the absolute sensing preference of worker k for location m needs to be reduced slowly by using the Formulas (39), (40). And φ is relatively small to help control to the speed of decline of the absolute sensing preference.

$$H_{k,m}^t = \text{Max}(H_{k,m}^0, H_{k,m}^{t-1} - \Psi). \quad (39)$$

$$\Psi = \varphi(1 - H_{k,m}^{t-1}). \quad (40)$$

After repeating the above process for some time, some workers with a high absolute sensing preference for location m are obtained. When the absolute preference of a worker for location m is higher than the threshold Γ , then this worker is considered to have high sensing quality of location m . If the UAV does not collect data at location m at time t , and there are multiple workers with absolute sensing preference for location m above the threshold Γ at location m at time t , the data of location m of worker b who has the highest absolute sensing preference for location m is selected as the Sub-L-GTD for location m , and other workers who also collect data at location m will be evaluated by worker b by utilizing Vice Location Reward Mechanism or Vice Location Punishment Mechanism.

The value of attribute n of location m collected by worker b at time t is $\mathcal{D}_{m,n}^{b,t}$, and the absolute sensing preference for location m is $H_{k,b}^t$. $R_{m,n}^k$ denotes the ratio of the difference between entry from ordinary worker k and Sub-L-GTD respectively.

$$R_{m,n}^k = \frac{|\mathcal{D}_{m,n}^{t,k} - \mathcal{D}_{m,n}^{b,t}|}{\mathcal{D}_{m,n}^{b,t}}. \quad (41)$$

If $R_{m,n}^k < \aleph^c$, let $R_{m,n}^k = 1$;

Else $R_{m,n}^k = 0$;

If the attribute n is not of interest to the task at location m , let $R_{m,n}^k = 0$.

$F_{m,k}^c$ represents the number of entries in location m collected by worker k that are evaluated as correct by Sub-Location-GTD.

$$F_{m,k}^c = \sum_{n=1}^N R_{m,n}^k. \quad (42)$$

$$\eta_m^c = \frac{F_{m,k}^c}{E_m}. \quad (43)$$

4) Vice Location Reward Mechanism

A threshold η_s^c is used to decide whether to apply Vice Location Reward Mechanism or Vice Location Punishment Mechanism. When η_m^c is greater than η_s^c , it means that worker k is considered to be credible in terms of location m at this time by Location-GTD, and Formula (44) is used to increase the absolute sensing preference of worker k for location m . In Formula (44), τ^c is a variable greater than 1. Use τ^c to control the speed of updating absolute sensing location preference.

$$H_{k,m}^t = H_{k,m}^{t-1} + \frac{\eta_m^c}{\tau^c} * (1 - H_{k,m}^{t-1}) * H_{b,m}^t. \quad (44)$$

5) Vice Location Punishment Mechanism

If η_m^c is smaller than the threshold η_s^c , it is considered that worker k is not credible in terms of location m . So the absolute sensing location preference will be reduced by Formula (45).

$$H_{k,m}^t = H_{k,m}^{t-1} + (\eta_m^c - 1) * \left(\frac{H_{k,m}^{t-1}}{\tau^c}\right) * H_{b,m}^t. \quad (45)$$

When both Sub-L-GTD and L-GTD are involved in the evaluation, the Location Timeout Mechanism will be applied to worker k if the data of location m of worker k is not evaluated by either Sub-L-GTD or L-GTD.

When there is no sub-L-GTD and only L-GTD, the Location Timeout Mechanism will be applied to worker k if the data of location m of worker k is not evaluated by L-GTD.

(2) Calculation of relative sensing location preference

For most of the workers, the collected data cannot be evaluated by Location-GTD or Sub-Location-GTD, so the concept of relative sensing location preference is proposed here.

The relative sensing location preference is divided into two parts, one is the Absolute Sensing Location Preference Sequence (ASLPS), and the other is the Autonomy Supervised Sensing Location Preference Sequence (ASSLPS).

1) The setup of ASLPS

The ASLPS refers to the sequence constructed in order according to the value of the absolute sensing location preference of workers whose absolute sensing location preference has been updated except by the Location Timeout Mechanism.

$H_{k,m}^t$ is the absolute sensing location preference of worker k for location m at round t , $H_m^t = \{H_{1,m}^t, \dots, H_{K,m}^t\}$, H_m^t is the set of absolute sensing location preference of all K workers in the network for location m at round t .

The higher the absolute sensing location preference for location m , the higher the quality of the sensing data for location m . Copy H_m^t to H'_m , then the workers whose sensing quality of location m has not been evaluated by L-GTD or Sub-Location-GTD are removed from H'_m . Then k_t remaining elements in H'_m are sorted in descending order, and the ASLPS for location m of remaining k_t workers are obtained as $L_m^t = \{L_{m,1}^t, \dots, L_{m,k_t}^t\}$. Each element in the sequence is the worker's number. If k is in the j -th place in L_m^t , then $L_{m,j}^t = k$, which means worker k has the j -th largest value of absolute sensing location preference for location m among all workers. The earlier the worker's number appears in the sequence, the better absolute sensing location preference for location m the worker has.

2) Calculation of RSLQW

The Relative Sensing Location Quality Weights (RSLQW) mean the relative sensing quality of location m among workers who collect data at location m . The higher the relative sensing location quality weight of a worker for location m , the better the sensing quality of this worker for location m relative to other workers at location m . ASSLPS means the partial order relationship of sensing quality of different workers for location m obtained by relative perceived quality weights. Therefore, to obtain the ASSLPS, we need to establish the RSLQW first.

Supposing there are K_m workers collecting data at location m . Calculating the RSLQW of location m is an unsupervised

computational process that requires establishing the initial relative truthful data of location m . $\mathcal{T}_{m,n}^{(*)}$ represents the relative truthful value of the attribute n of location m , it is established when calculating RSLQW of location m . $\mathcal{T}_m^{(*)}$ represents the relative truthful data of location m . Entries of relative truthful data of location m are established by using the Formula (46).

$$\mathcal{T}_{m,n}^{(*)} = \frac{\sum_{k=1}^{K_m} \mathcal{D}_{m,n}^{t,k}}{K_m}. \quad (46)$$

Because the higher the RSLQW of workers for location m , the closer the data given by workers for this location should be to the truthful data. As a result, the more accurate the relative perceived quality weights are assigned, the smaller the weighted sum of the deviations between the data collected by K_m workers and $\mathcal{T}_m^{(*)}$ is. Therefore, we have the objective function (47) to minimize the weighted sum of deviations between the data given by K_m workers and relative truthful data. $b_{k,m}^t$ denotes the Relative Sensing Location Quality Weights of worker k for location m at time t . The higher the $b_{k,m}^t$, the better the sensing quality of worker k relative to other workers. The set of Relative Sensing Location Quality Weights of location m is $B_m^t = \{b_{1,m}^t, b_{2,m}^t, \dots, b_{K_m,m}^t\}$. Assuming that location m has N attributes, and the set of data sensed by K_m workers at time t for location m is $\mathcal{D}_m^t = \{\mathcal{D}_{m,1}^{t,1}, \dots, \mathcal{D}_{m,N}^{t,1}, \dots, \mathcal{D}_{m,1}^{t,K_m}, \dots, \mathcal{D}_{m,N}^{t,K_m}\}$, the set of data of all locations at time t is $\mathcal{D}^t = \{\mathcal{D}_1^t, \dots, \mathcal{D}_M^t\}$.

$$\min \sum_{k=1}^{K_m} b_{k,m}^t * \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|, \quad (47)$$

$$S.T. f_c(b_{k,m}^t) = 1.$$

After establishing the objective function, the Relative Sensing Location Quality Weights of location m are calculated iteratively by the following two steps.

First step: Update RSLQW of location m .

The calculation of the RSLQW of location m is based on the current $\mathcal{T}_m^{(*)}$ and the data of location m \mathcal{D}_m^t provided by workers.

The calculation process is as follows: Assuming that the distribution of RSLQW of location m satisfies the exponential distribution, Equation (48) can be obtained.

$$f_c(b_{k,m}^t) = \sum_{k=1}^{K_m} e^{(-b_{k,m}^t)} = 1. \quad (48)$$

Theorem 3. Assuming that the true value is static, the optimal solution for the location trust degree of workers is given by:

$$b_{k,m}^t = \log \left(\frac{\sum_{k=1}^{K_m} \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|}{\sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|} \right).$$

Proof. Supposing that $v_{k,m}^t = e^{(-b_{k,m}^t)}$, The objective function (47) can be translated into the objective function (49):

$$\min \sum_{k=1}^{K_m} -\log(v_{k,m}^t) \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|,$$

$$S.T. \sum_{k=1}^{K_m} v_{k,m}^t = 1. \quad (49)$$

Since the objective function (49) is a linear combination of negative logarithm functions, and the constraint is linear to $v_{k,m}^t$. So this objective function is convex and can be solved by the Lagrange multiplier method.

$$L(\{v_{k,m}^t, \lambda\}) = \sum_{k=1}^K -\log(v_{k,m}^t) * \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})| + \lambda \left(\sum_{k=1}^{K_m} v_{k,m}^t - 1 \right). \quad (50)$$

λ is a Lagrange multiplier. Let the partial derivative of Lagrangian with respect to all $v_{k,m}^t$ be 0 in order to obtain the optimal value.

$$\begin{cases} -\frac{1}{v_{1,m}^t} \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,1})| + \lambda = 0, \\ -\frac{1}{v_{2,m}^t} \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,2})| + \lambda = 0, \\ \dots \\ -\frac{1}{v_{K_m,m}^t} \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,K_m})| + \lambda = 0. \end{cases}$$

After summing the equations above, we can have:

$$\sum_{k=1}^{K_m} -\frac{1}{v_{k,m}^t} \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})| + \lambda \sum_{k=1}^{K_m} 1 = 0, \quad (51)$$

$$\lambda \sum_{k=1}^{K_m} v_{k,m}^t = \sum_{k=1}^{K_m} \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|.$$

Because of the constraint $\sum_{k=1}^{K_m} v_{k,m}^t = 1$ in the objective function (47), the Equation (51) becomes:

$$\lambda = \sum_{k=1}^{K_m} \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|. \quad (52)$$

Supposing that the value of k is static, Equation (51) becomes:

$$\lambda v_{k,m}^t = \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|. \quad (53)$$

Combine Equations (53) and (52), we can get:

$$v_{k,m}^t = \frac{\sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|}{\sum_{k=1}^{K_m} \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|}.$$

Because $v_{k,m}^t = e^{(-b_{k,m}^t)}$, the optimal solution is:

$$b_{k,m}^t = \log \left(\frac{\sum_{k=1}^{K_m} \sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|}{\sum_{n=1}^N |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|} \right). \quad (54)$$

Second step: Update the relative truthful data of location m .

In order to minimize the objective function (47), it is necessary to update the relative truthful data after the calculation of Relative Sensing Location Quality Weights

$$\mathcal{T}_m^{(*)} \leftarrow \arg \min f(B_m^t, \mathcal{D}_m^t, \mathcal{T}_m^{(*)}).$$

Algorithm 1 RSLQW Estimation Algorithm

Input: Observation from K_m workers who collect data in location m : $\mathcal{D}_m^t = \{\mathcal{D}_{m,1}^{t,1}, \dots, \mathcal{D}_{m,N}^{t,1}, \dots, \mathcal{D}_{m,1}^{t,K_m}, \dots, \mathcal{D}_{m,N}^{t,K_m}\}$

Output: Relative relative sensing location quality weights set $B_m^t = \{b_{1,m}^t, b_{2,m}^t, \dots, b_{K_m,m}^t\}$ for location m

```

1: For  $n \leftarrow 1$  to  $N$  do
2:    $\mathcal{T}_{m,n}^{(*)} = \frac{\sum_{k=1}^{K_m} \mathcal{D}_{m,n}^{t,k}}{K_m}$ 
3: End for
4: Repeat
5:   Arrange  $B_m^t = \{b_{1,m}^t, \dots, b_{K_m,m}^t\}$  based on e. q. (54)
6:   For  $k \leftarrow 1$  to  $K_m$  do
7:     For  $n \leftarrow 1$  to  $N$  do
8:       Update relative truthful data  $\mathcal{T}_{m,n}^{(*)}$ 
         based on e. q. (56)
9:   End for
10: End for
11: Until Convergence criterion is satisfied;
12: Return  $B_m^t$ 
  
```

The calculation process for this step is as follows:

As to objective function (47), when n and m are fixed and a new set of Relative Sensing Location Quality Weights B_m^t has been determined for location m , since higher Relative Sensing Location Quality Weights indicate better sensing quality, in order to optimize the quality of sensing data, the objective function (55) is set up.

$$\mathcal{T}_{m,n}^{(*)} \leftarrow \arg \min \sum_{k=1}^{K_m} b_{k,m}^t * |(\mathcal{T}_{m,n}^{(*)} - \mathcal{D}_{m,n}^{t,k})|. \quad (55)$$

The weighted average method is used to minimize the objective function (55), then the new relative truthful data of location m is obtained.

$$\mathcal{T}_{m,n}^{(*)} = \frac{\sum_{k=1}^{K_m} b_{k,m}^t * \mathcal{D}_{m,n}^{t,k}}{\sum_{k=1}^{K_m} b_{k,m}^t}. \quad (56)$$

$\mathcal{T}_{m,n}^{(*)x}$ is the $\mathcal{T}_{m,n}^{(*)}$ from the x -th iteration, Formula (56) is used to compute the difference of relative truthful data between the two consecutive iterations.

If $\frac{|\mathcal{T}_{m,n}^{(*)x} - \mathcal{T}_{m,n}^{(*)x-1}|}{\mathcal{T}_{m,n}^{(*)x-1}} < 3$, let $\mathcal{H}_{m,n} = 1$; else let $\mathcal{H}_{m,n} = 0$.

$$d_m = \frac{\sum_{n=1}^N \mathcal{H}_{m,n} * E_{m,n}}{E_m}. \quad (57)$$

If d_m is less than the threshold value \mathcal{A} , which means that the difference between two consecutive iterations is small, and the iteration will stop. When the number of iterations exceeds the specified value ψ , the iteration will also stop.

Example 6. In the first iteration, Relative Sensing Location Quality Weights of l_4 of workers in Table 5 are: $b_{4,4}^1 = 0.847$, $b_{5,4}^1 = 1.135$, $b_{6,4}^1 = 1.386$. Then the relative truthful data of temperature on l_4 is updated as follows:

$$\frac{0.847 * 17 + 1.135 * 15 + 1.386 * 15}{0.847 + 1.135 + 1.386} = 15.503.$$

The calculation of the set $B_m^t = \{b_{1,m}^t, b_{2,m}^t, \dots, b_{K_m,m}^t\}$ is shown in Algorithm 1.

3) Set up of ASSLPS

Suppose that there are K_m^t workers collected at location m at time t . The set of K_m^t workers' number is $X^t =$

$\{x_1^t, x_2^t, \dots, x_{K_m^t}^t\}$, and we already know the Relative Sensing Location Quality Weights $B_m^t = \{b_{x_1^t,m}^t, b_{x_2^t,m}^t, \dots, b_{x_{K_m^t}^t,m}^t\}$.

After sorting the Relative Sensing Location Quality Weights B_m^t in descending order, the Autonomy Supervised Sensing Location Preference Sequence $R_m^t = \{r_{m,1}^t, \dots, r_{m,K_m^t}^t\}$ can be obtained from the worker number corresponding to the elements in the sorted sequence. For example, if $b_{k,m}^t$ is the 2nd largest in B_m^t , then worker k is in the 2nd in R_m^t , which means $r_{m,2}^t = k$. $R_m^t(2)$ also means $r_{m,2}^t$. If the worker's number k appears earlier in the sequence R_m^t , it means that the sensing quality of worker k for location m is better than other workers who also collect data on location m at time t . When it is not the first round, since the workers who sense data at location m in round t may not be the same with workers who collect data in round $t-1$, it is necessary to combine the ASSLPS of round $t-1$ with the ASSLPS of round t . Then the mixed sequence will be the ASSLPS of round t .

Next, the Merge Sequences Algorithm is proposed to merge two sequences. The first sequence $O_m^t = \{y_{m,1}^t, \dots, y_{m,x_t}^t\}$ containing x_t elements, and the second sequence $O_m^{t-1} = \{y_{m,1}^{t-1}, \dots, y_{m,x_{t-1}}^{t-1}\}$ containing x_{t-1} elements. The first sequence O_m^t has higher priority. N_m^t denotes the sequence after combining and synthesizing, $O_m^t \cap O_m^{t-1} = I_m^t$. The operations under different conditions are described below. Because O_m^t has a higher priority, the principle of operations under different conditions is keeping elements in the O_m^t and adding as many elements of O_m^{t-1} to O_m^t as possible.

Condition 1 operation:

If $O_m^t == NULL$, $O_m^{t-1} != NULL$:

Do $N_m^t = O_m^{t-1}$.

Condition 2 operation:

If $O_m^{t-1} == NULL$, $O_m^t != NULL$:

Do $N_m^t = O_m^t$.

Condition 3 operation:

If $O_m^{t-1} == NULL$, $O_m^t == NULL$:

Do $N_m^t = NULL$

Condition 4 operation:

If $O_m^t != NULL$, $O_m^t \cap O_m^{t-1} = k$, k in O_m^t is the element $y_{m,1}^t$ and in O_m^{t-1} is the element $y_{m,w}^{t-1}$:

Do $N_m^t = \{y_{m,1}^{t-1}, \dots, y_{m,1}^t, y_{m,2}^t, \dots, y_{m,x_t}^t\}$, that is, a segment is copied from O_m^{t-1} which starts from $y_{m,1}^{t-1}$ to $y_{m,w}^{t-1}$ and does not contain $y_{m,w}^{t-1}$, then the segment is added to the top of O_m^t .

Condition 5 operation:

If $O_m^t != NULL$, $O_m^t \cap O_m^{t-1} = k$, the number of elements which I_m^t contains is 1. k in O_m^t is the element y_{m,x_t}^t and in O_m^{t-1} is the element $y_{m,w}^{t-1}$:

Do $N_m^t = \{y_{m,1}^t, \dots, y_{m,x_t}^t, \dots, y_{m,x_{t-1}}^{t-1}\}$, that is, a segment is copied from O_m^{t-1} which starts from $y_{m,w}^{t-1}$ to $y_{m,x_{t-1}}^{t-1}$, then the segment is appended to the end of O_m^t . Fig. 3 shows the process of the Condition 5 operation. The first sequence O_m^t is {5,3,2,9}, the second sequence O_m^{t-1} is {7,9,6,4,5}, the same element of the two sequences is 9. Because 9 is the last element in the O_m^t , combining the segment which consists of elements behind 9 in the O_m^{t-1} and O_m^t can generate N_m^t .

Condition 6 operation:

If $O_m^t != NULL$, $O_m^t \cap O_m^{t-1} = I_m^t$, the number of elements which I_m^t contains is greater than 1, $y_{m,x_t}^t, y_{m,1}^t$ are both in

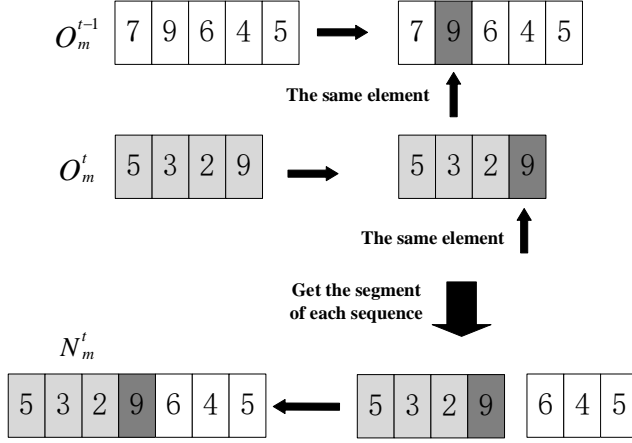


Fig. 3. The process of the Condition 5 operation

I_m^t :

Do $N_m^t = \{y_{m,1}^{t-1}, \dots, y_{m,1}^t, \dots, y_{m,x_t}^t, \dots, y_{m,x_t-1}^{t-1}\}$, that is, the segments from O_m^{t-1} are inserted at the top and end of O_m^t respectively.

Condition 7 operation:

If the number of elements which I_m^t contains is greater than 1, I_m^t contains $y_{m,1}^t$ but does not contain y_{m,x_t}^t :

Do $N_m^t = \{y_{m,1}^{t-1}, \dots, y_{m,1}^t, \dots, y_{m,x_t}^t\}$, that is, the segment of O_m^{t-1} is only added to the top of O_m^t .

Condition 8 operation:

If the number of elements which I_m^t contains is greater than 1, I_m^t contains y_{m,x_t}^t but does not contain $y_{m,1}^t$:

Do $N_m^t = \{y_{m,1}^t, \dots, y_{m,x_t}^t, \dots, y_{m,x_t-1}^{t-1}\}$, that is, the segment of O_m^{t-1} is only added to the top of O_m^t .

Condition 9 operation:

After judging the conditions above, if the value of two elements $y_{m,k}^t, y_{m,k+1}^t$ both exist in O_m^{t-1} , O_m^{t-1} has the segment $\{y_{m,k}^t, y_{m,j}^{t-1}, \dots, y_{m,j+u}^{t-1}, y_{m,k+1}^t\}$, and there exist two adjacent elements $y_{m,k}^t, y_{m,k+1}^t$ in O_m^t :

Remove the same elements of N_m^t and the segment $\{y_{m,j}^{t-1}, \dots, y_{m,j+u}^{t-1}\}$ from the segment $\{y_{m,j}^{t-1}, \dots, y_{m,j+u}^{t-1}\}$, assuming the $N_m^t = \{y_{m,1}^t, \dots, y_{m,x_t}^t\}$, then insert the segment into the N_m^t . we can get the renewed sequence $N_m^t: N_m^t = \{y_{m,1}^t, \dots, y_{m,k}^t, y_{m,q}^{t-1}, \dots, y_{m,e}^{t-1}, y_{m,k+1}^t, \dots, y_{m,x_t}^t\}$.

This case is to find two elements that both sequences contain and they are consecutive in O_m^t , then N_m^t is expanded by inserting a segment of O_m^{t-1} in the middle of two adjacent elements. And because operations of conditions 1-6 do not insert any element into the inner part of O_m^t , the adjacent elements of O_m^t will also be adjacent in N_m^t . Before inserting the segment, the removal of the same elements in two sequences is to make sure the new N_m^t will not contain duplicate elements.

The Autonomy Supervised Sensing Location Preference Sequence at time t R_m^t has a higher priority than the Autonomy Supervised Sensing Location Preference Sequence at time $t-1$ R_m^{t-1} . Input R_m^t and R_m^{t-1} into Algorithm 2, and the output is the updated R_m^t .

Fig. 4 shows the process of the Merge Sequence Algorithm. The first sequence O_m^t is $\{5, 7, 8, 12, 13, 11, 2\}$, the second sequence O_m^{t-1} is $\{3, 4, 5, 13, 6, 7, 11, 9, 2, 1, 10\}$, and the set I_m^t which consists of the same elements of two sequences is $\{5, 1, 11, 2\}$. Since the number of elements which I_m^t contains is greater than 1, and the first element and last element of O_m^t are both in I_m^t , the **Condition 6 operation** is utilized. Then a

Algorithm 2 Merge Sequences Algorithm

Input: Sequence $O_m^t = \{y_{m,1}^t, \dots, y_{m,x_t}^t\}$, Sequence $O_m^{t-1} = \{y_{m,1}^{t-1}, \dots, y_{m,x_t-1}^{t-1}\}$, O_m^t has a higher priority

Output: The result N_m^t

```

1: If  $O_m^t == \text{NULL}$  and  $O_m^{t-1} == \text{NULL}$ 
2:    $N_m^t = O_m^{t-1}$ 
3: End If
4: If  $O_m^t = \text{NULL}$  and  $O_m^{t-1} \neq \text{NULL}$ 
5:    $N_m^t = O_m^{t-1}$ 
6: End If
7: If  $O_m^t \neq \text{NULL}$  and  $O_m^{t-1} \neq \text{NULL}$ 
8:    $I_m^t = O_m^t \cap O_m^{t-1}$ 
9: End If
10: If  $I_m^t == \{y_{m,1}^t\}$ 
11:    $N_m^t = \{y_{m,1}^{t-1}, \dots, y_{m,1}^t, y_{m,2}^t, \dots, y_{m,x_t}^t\}$ 
12: End If
13: If  $I_m^t == \{y_{m,x_t}^t\}$ 
14:    $N_m^t = \{y_{m,1}^t, \dots, y_{m,x_t}^t, \dots, y_{m,x_t-1}^{t-1}\}$ 
15: End If
16: If  $\text{len}(I_m^t) > 1$ 
17:   If  $y_{m,1}^t$  in  $I_m^t$  and  $y_{m,x_t}^t$  in  $I_m^t$ 
18:      $N_m^t = \{y_{m,1}^{t-1}, \dots, y_{m,1}^t, \dots, y_{m,x_t}^t, \dots, y_{m,x_t-1}^{t-1}\}$ 
19:   End If
20:   If  $y_{m,1}^t$  not in  $I_m^t$  and  $y_{m,x_t}^t$  in  $I_m^t$ 
21:      $N_m^t = \{y_{m,1}^t, \dots, y_{m,x_t}^t, \dots, y_{m,x_t-1}^{t-1}\}$ 
22:   End If
23:   If  $y_{m,1}^t$  in  $I_m^t$  and  $y_{m,x_t}^t$  not in  $I_m^t$ 
24:      $N_m^t = \{y_{m,1}^{t-1}, \dots, y_{m,1}^t, \dots, y_{m,x_t}^t\}$ 
25:   End If
26:   For  $i \leftarrow 1$  to  $\text{len}(O_m^{t-1}) - 1$  do
27:     If  $O_m^t(i)$  and  $O_m^t(i+1)$  in  $I_m^t$ 
28:       do condition 9 operation
29:     End If
30:   End For
31: End If
32: return  $N_m^t$ 

```

new sequence N_m^t , $\{3, 4, 5, 7, 8, 12, 13, 11, 2, 1, 10\}$, is created. After that, traverse the O_m^t to find out if there exist some adjacent elements of I_m^t . First, element 5 and element 7 are adjacent in the sequence O_m^t , and O_m^{t-1} has the segment $\{5, 13, 6, 7\}$, so the **Condition 9 operation** is utilized. After elements that already exist in N_m^t are removed from the segment $\{5, 13, 6, 7\}$, segment $\{5, 13, 6, 7\}$ converts to segment $\{6\}$, then the segment $\{6\}$ is inserted into elements 5 and 7, sequence N_m^t converts to $\{3, 4, 5, 6, 7, 8, 12, 13, 11, 2, 1, 10\}$. Then keep traversing, because element 11 and element 2 are adjacent in the sequence O_m^t , element 11 and element 2 are in I_m^t and O_m^{t-1} exists segment $\{11, 9, 2\}$, N_m^t can be renewed. After elements that already exist in N_m^t are removed from the segment $\{11, 9, 2\}$, segment $\{11, 9, 2\}$ converts to segment $\{9\}$, then the segment $\{9\}$ is inserted into elements 11 and 2, sequence N_m^t converts to $\{3, 4, 5, 6, 7, 8, 12, 13, 11, 9, 2, 1, 10\}$.

4) The integration of ASSLPS and ASLPS

After updating the ASSLPS R_m^t , it is necessary to integrate R_m^t with ASLPS I_m^t .

The Merge Sequences Algorithm is utilized to integrate them and the result is G_m^t . G_m^t represents the sensing location preference partial order relationship of all workers on location m at time t . The set of workers' numbers in G_m^t is defined as

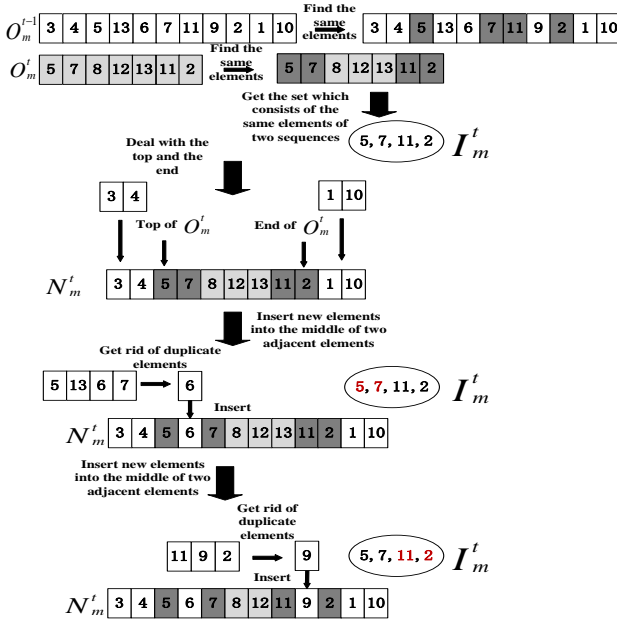


Fig. 4. The process of the Merge Sequence Algorithm

Algorithm 3 Integrated Merge Algorithm

Input: ASSLPS set at time t : $R^t = \{R_1^t, R_2^t, \dots, R_M^t\}$,
 ASSLPS set at time $t-1$: $R^{t-1} = \{R_1^{t-1}, R_2^{t-1}, \dots, R_M^{t-1}\}$
 ASLPS set: $L^t = \{L_1^t, L_2^t, \dots, L_M^t\}$, L_m^t has a higher
 priority than R_m^t , set of all workers' numbers W .

Output: The set of SLPPOR $G^t = \{G_1^t, G_2^t, \dots, G_M^t\}$

```

1: For  $m \leftarrow 1$  to  $M$  do
2:   If  $t > 1$ 
3:     Apply Merge Sequences Algorithm to merge
        $R_m^{t-1}$  and  $R_m^t$ ,  $R_m^t$  has a higher priority,
       then get  $N_m^t$  and assign  $N_m^t$  to  $R_m^t$ 
4:     If  $L_m^t = \text{NULL}$  and  $R_m^t = \text{NULL}$ 
5:        $G_m^t = L_m^t$ 
6:     End If
7:     If  $L_m^t = \text{NULL}$  and  $R_m^t \neq \text{NULL}$ 
8:        $G_m^t = R_m^t$ 
9:     End If
10:    If  $L_m^t \neq \text{NULL}$  and  $R_m^t \neq \text{NULL}$ 
11:       $G_m^t = \text{NULL}$ 
12:    End If
13:    If  $L_m^t \neq \text{NULL}$  and  $R_m^t \neq \text{NULL}$ 
14:      Apply Merge Sequences Algorithm to merge
         $L_m^t$  and  $R_m^t$ ,  $L_m^t$  has a higher priority,
        then get  $N_m^t$  and assign  $N_m^t$  to  $G_m^t$ 
15:    End If
16:    Get the set  $Q$ 
16:     $W' = W - Q$ 
17:    Shuffle  $W'$  to get  $w_m^t$ , then link  $w_m^t$  and  $G_m^t$ 
18:  End For
19: Return  $G^t$ 

```

$Q = \{i, j, \dots, q\}$, and the set of all workers' numbers in the network is defined as $W = \{1, 2, \dots, K\}$. So the set of workers' numbers whose sensing ability of location m has not been evaluated is $W' = W - Q$. After shuffling the set W' , we assign the shuffled set to a sequence $w_m^t = \{a, s, \dots, d\}$, which is added at the end of G_m^t , so that G_m^t contains all workers in the network, $G_m^t = \{i, j, \dots, q, a, s, \dots, d\}$. Sequence w_m^t is added at the end of G_m^t because of the uncertainty of workers whose numbers in w_m^t are large. Their priority needs to be

lower than that of workers who have been evaluated. The set of Autonomy Supervised Sensing Location Preference Sequences for M places is $R^t = \{R_1^t, R_2^t, \dots, R_M^t\}$, and the set of the Absolute Sensing Location Preference Sequence is $L^t = \{L_1^t, L_2^t, \dots, L_M^t\}$, the set of Sensing Location Preference Partial Order Relationship (SLPPOR) is $G^t = \{G_1^t, G_2^t, \dots, G_M^t\}$.

Algorithm 3 completely shows how to obtain the Sensing Location Preference Partial Order Relationship (SLPPOR) after obtaining the Absolute Sensing Location Preference Sequence (ASLPS) of location m and Autonomy Supervised Sensing Location Preference Sequence (ASSLPS) of location m .

4.3 The Task Assignment Model in MLM-WR

With Section 4.2, we obtain the set of sensing location preference partial order relationship $G^t = \{G_1^t, G_2^t, \dots, G_M^t\}$ for M location, set of sensing attribute preference $Z = \{Z_1^t, \dots, Z_N^t\}$ for N attributes, and the credibility set $C^t = \{C_{1,t}, \dots, C_{K,t}\}$ at time t . As to the K workers in the network, assuming that there are $K' = K * \phi$ workers applying to the platform to perform tasks in time t . Because of the limited resources, we need to make the most efficient recruitment and assignment of workers for sensing tasks.

The sequence of tasks is $p = \{p_1, p_2, \dots, p_M\}$. Each location has N attributes and the set of locations is $\ell = \{\ell_1, \ell_2, \dots, \ell_M\}$. Each worker submits multiple locations that he/she is willing to go to perform the task, and each worker is willing to go to different locations. A worker submits data for a location at a cost of c . A worker can only go to one location at a time to perform the task. The process of task assignment is as follows.

(1) Filter malicious workers

Because malicious data can cause bad consequences to the platform, it is necessary to exclude the identified malicious workers before recruiting workers. θ is a threshold to judge whether a worker is malicious or not. When the credibility of a worker is less than θ , the worker is considered to be a malicious worker. K'' workers are left after the malicious workers are excluded.

(2) Arrange the locations where the UAV senses data

If there are Y identified truthful workers among remaining K'' workers, then other workers are identified as ordinary workers. The number of locations collected by UAV is U , and the set of locations where Y workers are willing to go is ℓ' . The set of locations $\ell'' = \ell - \ell'$ means the locations where no truthful workers are willing to go. If the number of elements in ℓ'' is greater than U , then U elements of ℓ'' are randomly selected as locations for the UAV to sense data. The set of locations where UAV to sense data is ℓ^U , and the set of locations where truthful workers sense data is ℓ^C , then $\ell^C = \ell'$; if the number of elements in ℓ'' is less than U , locations in ℓ'' are assigned to the UAV. And UAV can also go to $U - \text{len}(\ell'')$ locations to sense data. The $U - \text{len}(\ell'')$ locations are randomly selected from ℓ' . Then ℓ^U is formed. The set of locations where truthful workers go is $\ell^C = \ell' - \ell^U + \ell''$.

(3) Arrange truthful workers

After arranging the locations where the UAV senses data, the set ℓ^C is determined. Assuming ℓ^C contains J elements. Therefore, not all locations need to be sensed by truthful workers. The sequence of tasks that need truthful workers can be determined from ℓ^C as $p_c = \{p_{\ell_1}, p_{\ell_2}, \dots, p_{\ell_J}\}$. ℓ_i denotes the location corresponding to the i -th task of the task

sequence of truthful workers, e.g. $\mathcal{L}_i = m$, which means that the i -th task $p_{\mathcal{L}_i}$ of the task sequence of truthful workers is to go to location m to execute the task p_m .

Because of the limited resources, it is necessary to allocate these Y workers to maximize efficiency. After the set of sensing attribute preference $Z_n^t = \{Z_{1,n}^t, \dots, Z_{K,n}^t\}$ is sorted in descending order, the sequence Z'_n is obtained by the worker number corresponding to the element of the sorted set Z_n^t . Sequence Z'_n consists of workers' numbers. If k shows earlier in the sequence, it means worker k has a higher value of sensing attribute preference of attribute n , that is, $Z_{k,n}^t$ is larger. $s_{k,n}^t$ is the position of worker k 's number in Z'_n . The smaller $s_{k,n}^t$ is, number k shows earlier in the sequence. Small $s_{k,n}^t$ means good sensing attribute quality of worker k for attribute n . The value range of $s_{k,n}^t$ is $[1, K]$, $s_n^t = \{s_{1,n}^t, \dots, s_{K,n}^t\}$. $a_{k,n}^t$ denotes the sensing attribute credibility of worker k for attribute n . When worker k has a smaller position in the sequence Z'_n , the sensing attribute credibility of attribute n is better, so $a_{k,n}^t$ is higher. Formula (58) is used to establish the sensing attribute credibility.

$$a_{k,n}^t = K - s_{k,n}^t + 1. \quad (58)$$

The set of sensing attribute credibility of attribute n is $a_n^t = \{a_{1,n}^t, \dots, a_{K,n}^t\}$, and the sensing attribute credibility set of workers with N attributes is $a^t = \{a_1^t, \dots, a_N^t\}$.

$g_{k,m}^t$ represents the rank of worker k 's number in G_m^t , the value range of $g_{k,m}^t$ is $[1, K]$, the smaller the rank of k in G_m^t , the higher the sensing location credibility of worker k for location m . $e_{k,m}^t$ is defined as the sensing location credibility of worker k at moment t for location m , which can be calculated by the Formula (59). The set of sensing location credibility of location m is $e_m^t = \{e_{1,m}^t, \dots, e_{K,m}^t\}$, and the set of sensing location credibility of workers at M locations is $e^t = \{e_1^t, \dots, e_M^t\}$.

$$e_{k,m}^t = K - g_{k,m}^t + 1. \quad (59)$$

T_m^t denotes the execution degree of the task p_m at time t . It is calculated by the Formula (60). Formula (60) reflects that the execution degree of the task p_m is high if worker k has high sensing location credibility for location m and high sensing attribute credibility for the attributes that task p_m focuses on.

$$T_m^t = \sum_{n=1}^N e_{k,m}^t * a_{k,n}^t * E_{mn}. \quad (60)$$

In order to obtain the optimal solution, the algorithm based on Particle Swarm Optimization (PSO) is used. PSO is used to simulate a school of fish or a flock of birds that moves in a group and can profit from the experience of all other members [39]. In other words, while a bird is flying and searching randomly for food, for instance, all birds in the flock can share their discovery and help the entire flock get the best hunt. The PSO algorithm relies on the exchange of information among individuals, and each solution candidate is called a particle. Swarm means all particles. The position of a particle represents the current solution of this particle. Two factors affect the particle's status: its position and velocity.

The number of truthful workers that can be assigned is Y ,

J is the number of tasks that truthful workers are willing to perform, and P is the size of the PSO population. $PSO[i]$ represents the position of the i -th particle, which is a J -dimensional vector. Since the workers who apply for each task are not the same, each element of $PSO[i][j]$ sets its unique numbering system of worker number. x_j represents the number of workers applying for the task $p_{\mathcal{L}_j}$. $PSO[i][j]$ encodes x_j workers of task $p_{\mathcal{L}_j}$, so the number of $PSO[i][j]$ ranges from $\{1, 2, 3 \dots x_j\}$, e.g. $PSO[i][j] = 3$ means task $p_{\mathcal{L}_j}$ is matched to the 3rd worker under $p_{\mathcal{L}_j}$'s numbering method.

$PSO[i]$ initialization starts from the first task in the task sequence. One worker is randomly selected from workers who are willing to perform $p_{\mathcal{L}_1}$ and have not been assigned yet. Then the worker is assigned to perform this task. If task $p_{\mathcal{L}_j}$ cannot select a worker willing to participate and not assigned, then $PSO[i][j] = -1$, which means no worker will be assigned to this task.

The fitness function is used to evaluate the quality of particles. The fitness function needs to consider the sensing quality of all tasks. Equation (61) is used as the fitness function. The aim of Equation (61) is to calculate the sum of the execution degree of all tasks. If the T_{total} is higher, which means the quality of collected data is better. So this algorithm is aimed to get the highest T_{total} which is calculated by Equation (61). $f[i]$ is the result value of the i -th particle according to the fitness function. G_b is the index of the current global optimal particle, $P_b[i]$ is the historical optimal solution of the i -th particle, $P_{b,f}[i]$ is the historical best fitness of the i -th particle, and P_g represents the historical optimal solution of all particles, and the historical optimal (highest) fitness of all particles is represented by f_g .

$$T_{total} = \sum_{i=1}^J T_{\mathcal{L}_i}^t. \quad (61)$$

For the i -th particle, we can obtain the following optimization equation.

$$V[i] = \omega * V[i] + q_1 * r_1 (P_b[i] - PSO[i]) + q_2 * r_2 (P_g - PSO[i]). \quad (62)$$

$$PSO[i] = PSO[i] + V[i]. \quad (63)$$

Where r_1 and r_2 are recalculated at each update, both are random numbers in the range of $(0, 1)$. q_1 and q_2 are positive constant parameters, they are used to control the maximum step size the particle can do. ω is the inertia weight, ω means the degree of influence of the past velocity on the present velocity. If the inertia weight is small it means a fine-tuning of the present position. If the inertia weight is large, it encourages global exploration.

After being updated by Formulas (62) and (63), the elements of each velocity vector will be rounded down, because there is a high probability that the value of elements in the position vector is not an integer. But the decimal number is meaningless, so it is necessary to round down the number. And if a number is less than -1, it is converted to -1. If $PSO[i][j]$ is greater than the x_j , it will be converted to x_j .

For example, assuming that the position vector is $\{1, 3, 2, 2\}$, and the velocity vector is $\{1.82, 1.94, 2.51, -3.67\}$. After

adding them together, we get a new vector equal to $\{2.82, 4.94, 4.51, -1.67\}$, which is then converted to $\{2, 4, 4, -1\}$. It means that p_{L_1} is assigned to the 2nd worker under the numbering method of p_{L_1} . p_{L_2} is assigned to the 4th worker under the numbering method of p_{L_2} , p_{L_3} is assigned to the 3rd worker under the numbering method of p_{L_3} , and no worker is assigned to p_{L_4} .

After the $PSO[i]$ is calculated, it is necessary to check whether the duplicate worker is scheduled. For the worker k scheduled more than once, the task execution degree of each location assigned to worker k will be calculated. Then the location with the highest task execution degree will still be assigned to worker k , other locations will not assign workers. So the corresponding values for other locations in $PSO[i]$ will become -1. For example, a position vector is $\{1, 2, 1, 2\}$. The first worker under the numbering method of p_{L_1} and the second worker under the numbering method of p_{L_4} both mean the worker k . The execution degree of p_{L_1} is 100 and the execution degree of p_{L_4} is 120, so we arrange for the worker k to do p_{L_4} , and the position vector becomes $\{-1, 2, 1, 2\}$.

$h[i]$ denotes the number of elements after removing -1 from the i th particle, it is also called the number of scheduled workers of particle i . $h_b[i]$ represents the largest number of historically arranged workers of particle i .

h_g represents the largest number of historically arranged workers of all particles. The more locations are assigned to truthful workers, the better the particle will be, so $h[i]$ will be calculated before $f[i]$ when judging the particle i .

When two particles are compared, for example, particle i and particle j are compared. First, compare $h[i]$ and $h[j]$. if $h[i]$ is larger, then it means that particle i is better, if $h[i] = h[j]$, then compare fitness, if $f[i] > f[j]$, then it is considered that particle i is better. When comparing two particles, the priority is to compare the number of elements after removing 1, and the second priority is to compare the value calculated by the fitness function. So for all the particles of a generation, the best particle is the one with the largest number of arranged workers first, and the one with the highest fitness among particles whose number of arranged workers is the largest. The task assignment algorithm for credible workers is shown in Algorithm 4.

(4) Arrangement of ordinary workers.

Ordinary workers means workers who are not judged to be truthful or malicious. There are two situations with different operations:

1) When there are no truthful workers applying.

$W_{p_m}^t = \{W_z, W_b, \dots, W_q\}$ is the set of workers who have not been arranged at time t and are willing to perform the task p_m . α is a threshold to arrange ordinary workers when no truthful worker or UAV is going to perform the task p_m . If the number of elements in $W_{p_m}^t$ is greater than threshold α , α workers are randomly selected from $W_{p_m}^t$ to be dispatched to execute the task p_m ; if the number of elements in $W_{p_m}^t$ is less than α , then all workers of $W_{p_m}^t$ are dispatched to execute the task p_m .

2) When there are truthful workers applying.

Assuming that there is no truthful worker or UAV that can execute the task p_m . $W_{p_m}^t = \{W_z, W_b, \dots, W_q\}$ is the set of workers who have not been arranged at time t and are willing to perform the task p_m . α is a threshold to arrange ordinary workers when no truthful worker or UAV is going to perform

Algorithm 4 Task Assignment of Truthful Workers

Input: Sensing attribute credibility set $a^t = \{a_1^t, \dots, a_N^t\}$, sensing location credibility set $e^t = \{e_1^t, \dots, e_M^t\}$
Output: Task assignment solution

```

1: For  $i \leftarrow 1$  to  $P$  do //Initialization
2:   For  $j \leftarrow 1$  to  $J$  do
3:     Initialize  $PSO[i][j]$  by randomly choosing ID from
        $\{1, 2, 3 \dots x_j\}$  and the ID
       corresponding worker is not assigned
4:   End For
5: End for
6: For each particle  $i$  in the population do
7:   Initialize  $V[i]$  randomly
8:   Evaluate  $f[i]$  according to equation (61)
9:   Evaluate  $h[i]$  by get rid of -1 in  $PSO[i]$ 
10:   $P_{b_f}[i] \leftarrow f[i]$ 
11:   $h_b[i] \leftarrow h[i]$ 
12:   $P_b[i] \leftarrow PSO[i]$ 
13: End for
14: Initialize  $G_b$  with the index of the particle which
    arrange most workers and gets highest fitness
15:  $P_g \leftarrow PSO[G_b]$ 
16:  $f_g \leftarrow f[G_b]$ 
17:  $h_g \leftarrow h[G_b]$ 
18: Repeat //Update
19:   For each particle  $i$  in the population do
20:     Update  $PSO[i]$  according to equations (62), (63)
21:     Evaluate fitness[i] according to equation (61)
22:     check repeatedly selected workers and adjust PSO
23:   End For
24:   Find  $G_b$  according to  $h[G_b] \geq h[i], \forall i \leq P$ 
    and get the highest fitness among the particles
    which get the highest  $h[i]$ 
25:   if  $h_g \leq h[G_b]$  and  $f_g \leq f[G_b]$  do
26:      $P_g \leftarrow PSO[G_b]$ 
27:      $h_g \leftarrow h[G_b]$ 
28:   End if
29:   For each particle  $i$  in the population do
30:     if  $h[i] \geq h_b[i]$  and  $f[i] \geq P_{b_f}[i]$  do
31:        $P_b[i] \leftarrow PSO[i]$ 
32:        $P_{b_f}[i] \leftarrow fitness[i]$ 
33:        $h_b[i] \leftarrow h[i]$ 
34:     End if
35:   End for
36: Until a number of generations;
37: return  $P_g$ 

```

the task p_m . If the number of elements in $W_{p_m}^t$ is greater than threshold α , α workers with the highest credibility are selected from $W_{p_m}^t$ to be dispatched to execute the task p_m ; if the number of elements in $W_{p_m}^t$ is less than α , then all workers of $W_{p_m}^t$ are dispatched to execute the task p_m .

Assuming that there is a truthful worker or UAV has been arranged to execute the task p_m . β is a threshold to arrange ordinary workers when a truthful worker or UAV has been arranged to perform the task p_m . If the number of elements in $W_{p_m}^t$ is less than β , then all the workers in $W_{p_m}^t$ are arranged to execute the task p_m . If the number of elements in $W_{p_m}^t$ is greater than β , β workers with the highest credibility are selected from $W_{p_m}^t$ to be dispatched to execute the task p_m ;

In order to save the cost of sending workers to perform tasks,

as the number of identified truthful workers increases, the value of β will be reduced.

There are M locations in the network. If the number of workers that can be dispatched is less than $\theta_1 * M$, let $\beta = \beta_1$; Else if the number of workers that can be dispatched is less than $\theta_2 * M$, let $\beta = \beta_2$; Else if the number of workers that can be dispatched is less than $\theta_3 * M$, let $\beta = \beta_3$; Else if the number of workers that can be dispatched is less than $\theta_4 * M$, let $\beta = \beta_4$; Else let $\beta = 0$. The adjustment of the number of dispatched workers is divided into 5 stages, $\theta_1, \theta_2, \theta_3, \theta_4$ are all variables greater than 1. If there are enough truthful workers to dispatch, there is no need to dispatch ordinary workers to places where truthful workers can be dispatched.

5. Performance Analysis

5.1 Experiment Setup

The MLM-WR scheme is implemented by Python. And the dataset The Diabetes Dataset [44] is used in our simulation experiments, which has 768 objects and 6912 observations. Each object has 9 attributes. Since these objects have different attributes, they can be utilized to simulate data from different locations.

Four different experiment scenarios are set up: (1) There are 100 workers in the network, 10 objects are selected as locations from The Diabetes Dataset and the values of 5 of the attributes are kept as the baseline data. UAV visits 2 locations each time. The platform wants to collect data on these 10 locations. (2) There are 150 workers in the network, 15 objects are selected as locations from The Diabetes Dataset and the values of 5 of the attributes are kept as the baseline data. UAV visits 3 locations each time. The platform wants to collect data on these 15 locations. (3) There are 200 workers in the network, 20 objects are selected as locations from The Diabetes Dataset and the values of 5 of the attributes are kept as the baseline data. UAV visits 2 locations each time. The platform wants to collect data on these 20 locations. (4) There are 2000 workers in the network, 200 objects are selected as locations from The Diabetes Dataset and the values of 5 of the attributes are kept as the baseline data. UAV visits 20 locations each time. The platform wants to collect data on these 200 locations.

For each experimental scenario, 20 rounds of data are generated to simulate the real-life location data which are constantly changing, with each entry in each round fluctuating by 20% from the baseline data. For different locations, the platform focuses on different attributes. The result of each experimental scenario is the mean of 20 times the experiments' results.

The worker's attribute deviation is u , the location deviation is v . Assuming the truthful value of an entry is T^r , then the entry value of the generated worker is $T^r * (1 + u) * (1 + v)$. u, v has a 10% probability of being negative. The proportion of malicious workers is 20%, that of ordinary workers is 20%, and that of truthful workers is 60%. The range of values of u, v for malicious workers is $[[0.15, 0.2]]$. The range of values of u, v for ordinary workers is $[[0.1, 0.15]]$, and the range of values of u, v for truthful workers is $[[0.05, 0.1]]$.

In the Credibility Evaluation Model, the appropriate parameters are set as: When data of location m is evaluated,

$\sigma_s = E_m/2, \sigma_s^c = E_m/2, \rho = 4, v = 0.7, \theta = 0.3, \gamma = 0.003, \theta = 0.25, \theta^c = 0.25$.

In the SAPE, $\delta = 4, \Omega = 0.25, \delta^c = 4, \Omega^c = 0.25, \mu = 0.003, \delta = 0.7$.

In the SLPE, $\eta_s = E_m/2, \eta_s^c = E_m/2$ when data of location m is evaluated. $\tau = 4, \tau^c = 4, \kappa = 0.25, \kappa = 0.25, \Gamma = 0.6, \varphi = 0.003$.

In the Task Assignment Model, $q_1 = 3, q_2 = 3, \omega = 0.9$. In the first round, the task assignment is that each location is assigned 5 workers who are randomly selected. After the first round, the set of parameters is: $\theta_1 = 1.8, \beta_1 = 4, \theta_2 = 3.0, \beta_2 = 3, \theta_3 = 4.2, \beta_3 = 2, \theta_4 = 5.0, \beta_4 = 1$.

For comparison, we chose two reference schemes. The first is TAFR and TFGR [9], which show state-of-the-art performance in estimating workers' reliability as to workers who collect data from the same locations. The TFAR calculates the attribute reliability based on the data collected by the worker, and TFGR calculates the location reliability based on the data collected by the worker. Both methods have only the process of calculating credibility without the process of task assignment. In order to make the comparison fair, TFAR and TFGR are combined. The original TFAR algorithm calculates the attribute reliability among workers who collect data from the same location. However, in reality, workers who perform the task at each location are different, so the attribute reliability calculated for each location is cross-mixed.

The principle of mixing is that in the sequence of attribute reliability obtained from l_1, l_2, \dots, l_m , the first worker in each sequence from l_1 to l_m is formed into an integrated attribute credibility sequence. Then these workers are deleted from the original sequence. After that, starting from l_1 to l_m , the first unarranged worker of each sequence is formed into a new sequence at the end of the integrated attribute credibility sequence, then these workers are deleted from the original sequence. Repeat doing the operations mentioned above, if all the workers of sequence l_i are deleted, no more workers in sequence l_i will be integrated. The repetition will stop until all the workers in the attribute reliability sequences of l_1, l_2, \dots, l_m are deleted. For example, if the attribute reliability sequence of location 1 is worker 1, worker 2, worker 3, and the attribute reliability sequence of location 2 is worker 4, worker 5, worker 6, then the integrated attribute credibility sequence is worker 1, worker 4, worker 2, worker 5, worker 3, worker 6. As to TFGR, the location credibility sequences of workers at each location of each round are combined by using the Integrated Merge Algorithm proposed in this paper. The part of the Absolute Sensing Location Preference Sequences is moved while using the Integrated Merge Algorithm to obtain the integrated location credibility sequence.

The final task assignment is based on the integrated attribute credibility sequence and integrated location credibility sequence which are mentioned above. The task assignment utilizes Algorithm 4 proposed in this paper.

In the first round, the TFAR+TFGR method randomly sends 5 workers to each location to sense data. After the first round, besides the UAV or workers assigned by the PSO algorithm, about 4 additional workers are assigned to each location.

The second scheme is the Random method which is mentioned in the reference [45]. This scheme does not identify the credibility of workers or identify the attribute credibility and

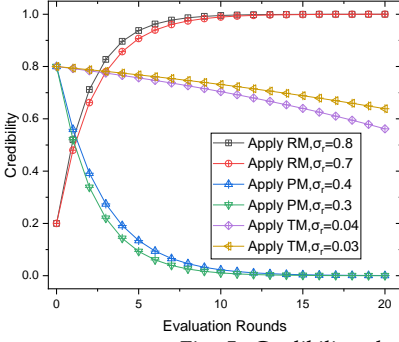


Fig. 5. Credibility changes in rounds when $\rho = 2$

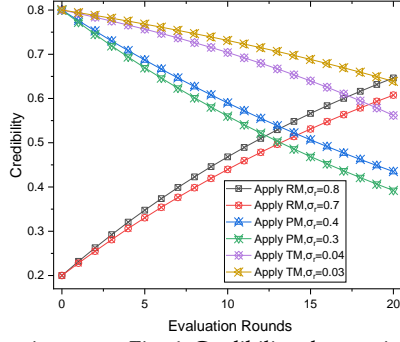


Fig. 6. Credibility changes in rounds when $\rho = 20$

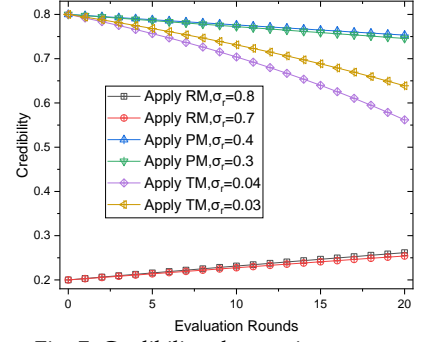


Fig. 7. Credibility changes in rounds when $\rho = 200$

location credibility. This scheme randomly dispatches 5 workers per location per round in the experiments.

5.2 Experiment Result

The strategy of credibility evaluation is to calculate σ_r according to the data submitted by workers. The Reward Mechanism (RM) is utilized when σ_r is greater than σ_s , and the Punishment Mechanism (PM) is used if σ_r is less than σ_s . The Timeout Mechanism (TM) is used for workers whose data have not been assessed by GTD or Sub-GTD. Three mechanisms can make a big difference in credibility between workers, then truthful workers or malicious workers can be easily identified.

Fig. 5 reflects the change of credibility by continuously using three assessment mechanisms. When σ_r exceeds σ_s , as the update time increases, the credibility continues increasing. Although the initial credibility is only 0.2, it eventually converges to 1. In addition, the larger the σ_r , the more obvious the increment in credibility. When σ_r is smaller than σ_s , the

Punishment Mechanism is utilized. Although the initial credibility is high, the credibility gradually decreases and eventually converges to 0. The curves shown in the graph are consistent with Theorem 1 and Theorem 2; σ_r affects the decrease of credibility. When σ_r is smaller, the decrease in credibility will be lower because the smaller σ_r means lower data quality. As to workers who hardly submit data, a penalty will also be applied, then the Timeout Mechanism will be used. But the credibility decreases much slower than workers who submit malicious data. The initial credibility is 0.8, which means that the previous performance is good, but because the data of the worker is always not been evaluated, the credibility keeps decreasing. This means that the previous merit is erased. It can be seen that these evaluation mechanisms not only evaluate the accuracy of the data submitted by workers but also evaluate the interaction frequency of workers. Fig. 6 shows the curves when the weight $\rho = 20$, and Fig. 7 shows the curves when the weight $\rho = 200$. As the value of ρ becomes larger, the update of credibility becomes slower.

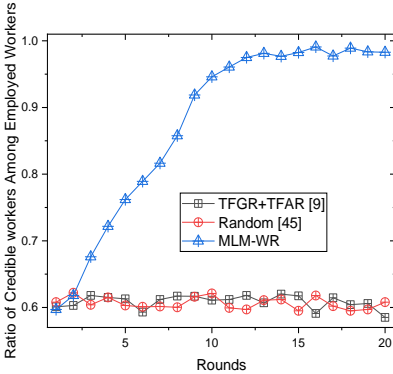


Fig. 8. Ratio of Credible Workers Among Employed Workers (when the number of workers = 100)

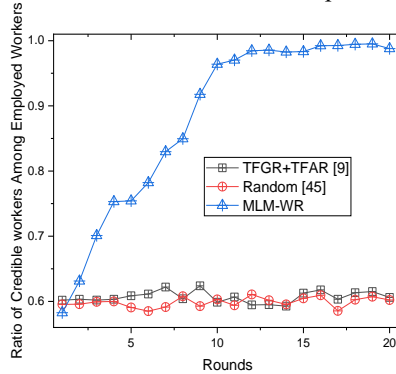


Fig. 9. Ratio of Credible Workers Among Employed Workers (when the number of workers = 150)

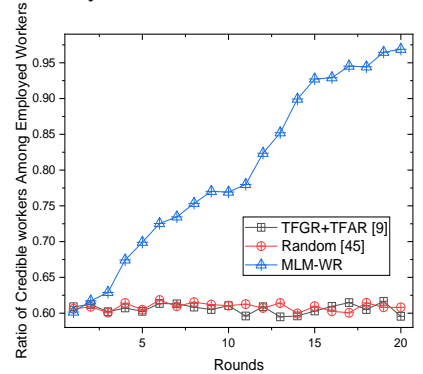


Fig. 10. Ratio of Credible Workers Among Employed Workers (when the number of workers = 200)

Fig. 8 shows the change of the ratio of credible workers among employed workers when there are 100 workers in the network. Because the TFAR+TFGR method and the Random method do not filter the workers by excluding the malicious ones after the credibility evaluation. Both methods can not identify truthful workers, so the ratio of truthful workers among the employed workers does not increase with the increase of rounds. In contrast, with the detection of truthful workers by GTD and Sub-GTD in MLM-WR, the ratio of truthful workers among employed workers keeps increasing. Fig. 9 shows the change of the ratio of credible workers when

the number of workers in the network is 150, and Fig. 10 shows the curves when 200 workers are in the network.

Fig. 11 shows the number of detected truthful workers with the increase of rounds. The number of workers evaluated in each round by GTD is constant, and as the number of rounds increases, the number of detected truthful workers whose data can act as Sub-GTD increases. So the curve of the number of detected truthful workers will show an exponential increase at first. The flattening out in the later part is a reflection of the decrease in the number of workers hired. After the number of

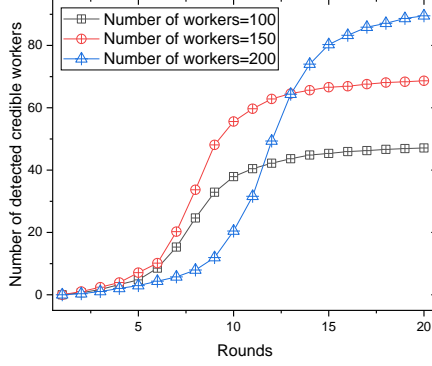


Fig. 11. Number of detected truthful workers

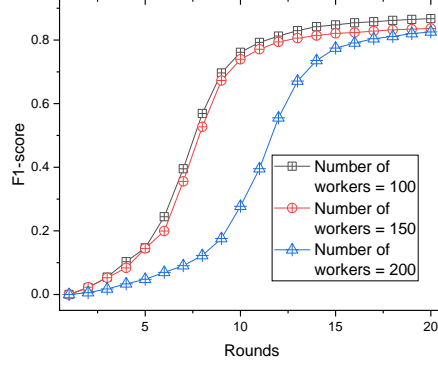


Fig. 12. F1-score

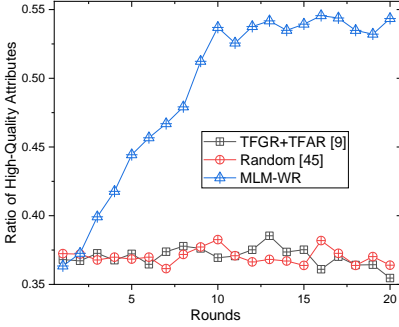


Fig. 13. Rate of High-Quality Attributes (when the number of workers = 100)

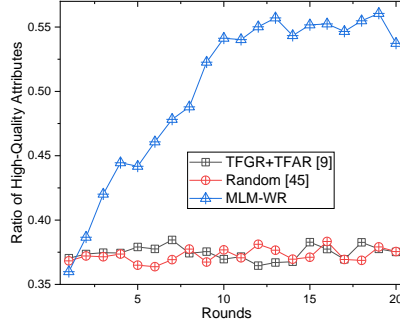


Fig. 14. Rate of High-Quality Attributes (when the number of workers = 150)

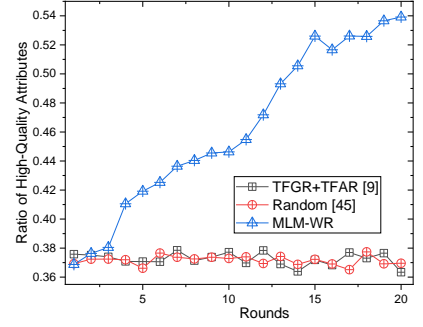


Fig. 15. Rate of High-Quality Attributes (when the number of workers = 200)

detected truthful workers reaches a certain value, the Task Assignment Model will decrease the number of hired workers, so the growth rate of the number of detected truthful workers will decrease.

Fig. 12 shows the change of the F1-score value with the increase of rounds. Because the higher the ratio of the number of locations detected by UAV to the number of all locations, the higher the proportion of workers can be evaluated. If more workers can be evaluated, the F1-score can increase faster. In Fig. 11, when the number of workers is 200 the ratio of the number of locations detected by UAV to the number of all locations is the lowest, so the speed of detecting truthful is the slowest. As a result, the F1-score increases more slowly than that of the other two cases.

Figs. 13-15 show the change in the rate of high-quality attributes. High-quality attributes are the sensing attributes of

workers whose attribute deviation \mathcal{U} range is $[[0.05, 0.75]]$. The rate of high-quality attributes is the ratio of the number of high-quality attributes among employed workers to the number of employed workers multiplied by the number of attribute categories. Because even if the truthful workers are identified, there are differences in the sensing attribute quality among truthful workers, the truthful workers with better sensing attribute quality should be selected as the hired workers. In MLM-WR, because the high-quality attributes of truthful workers will be gradually identified as the number of rounds increases, the rate of high-quality attributes will also increase. The TFAR+TFGR method does not integrate the historical result of evaluating sensing attribute quality, so there is no significant increase in the rate of high-quality attributes. High-quality locations are the sensing locations of workers with the location deviation \mathcal{U} ranging from $[[0.05, 0.75]]$.

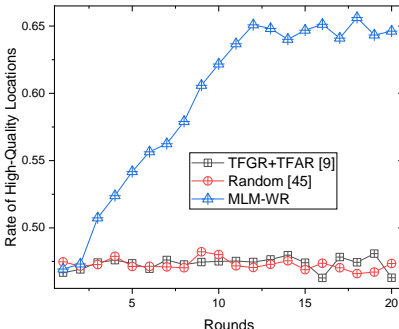


Fig. 16. Rate of High-Quality Locations (when the number of workers = 100)

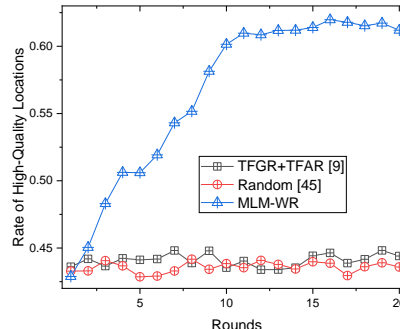


Fig. 17. Rate of High-Quality Locations (when the number of workers = 150)

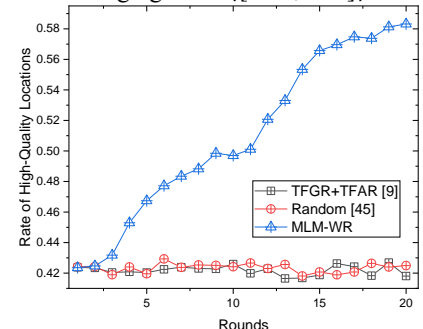


Fig. 18. Rate of High-Quality Locations (when the number of workers = 200)

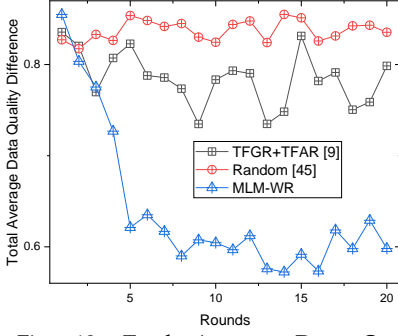


Fig. 19. Total Average Data Quality Difference (The number of workers = 100)

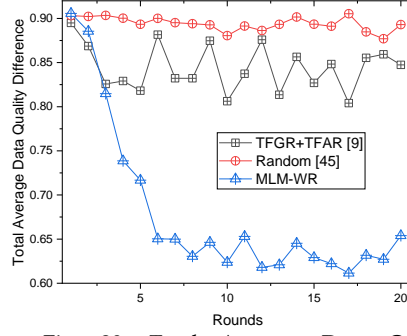


Fig. 20. Total Average Data Quality Difference (The number of workers = 150)

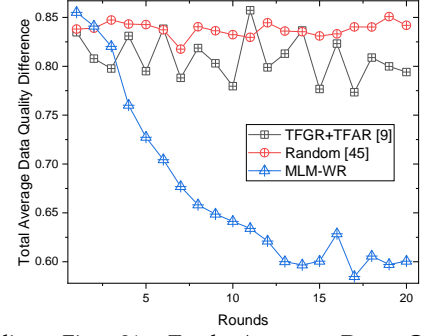


Fig. 21. Total Average Data Quality Difference (The number of workers = 200)

The MLM-WR method combines two aspects for sensing location preference calculation, one is the absolute sensing location preference and the other is the relative sensing location preference. The calculation of absolute sensing location preference utilizes the UAV like the calculation of sensing attribute preference, so the overall trend is similar. However, because of the existence of relative sensing location preference, the Sensing Location Preference Evaluation Model is more comprehensive than Sensing Attribute Preference Evaluation Model. So the curve is smoother when the high-quality locations of workers are detected than high-quality attributes of workers are detected. The rate of high-quality locations is the ratio of the number of high-quality locations to the number of workers dispatched multiplied by the number of locations. Figs. 16-18 show the change of rate of high-quality locations in the network, which shows similar but smoother curves than the curves in Figs. 13-15. The rate of high-quality locations increases until the number of high-quality locations are almost all selected.

Fig. 19 shows the change in the total average data quality difference when there are 100 workers in the network. The total average data quality difference is calculated without considering the locations visited by the UAV because the platform uses the data sensed by the UAV as the result data for the location visited by UAV, and the data collected by the UAV is truthful. As to the locations visited by truthful workers assigned by the platform, only the data sensed by truthful workers at this location are considered because the platform takes the data of truthful workers as the result data for locations

where they go. As to locations that are not assigned to the UAV or truthful workers, all sensing data are considered. For one thing, MLM-WR improves the quality of employed workers because it continuously identifies truthful workers. For another thing, it continuously identifies high-quality locations and high-quality attributes. So it can accurately dispatch workers with high sensing preference for the designated locations and attributes to perform tasks. As a result, MLM-WR can significantly reduce the total average data quality difference. In contrast, the TAGR+TFAR method does not improve the quality of workers as the number of rounds increases because there is no identification of truthful workers. But its identification of attribute and location trustworthiness makes the total average data quality difference smaller than that of the Random method. In the 20th round, the total average data quality difference of MLM-WR is 28.41% less than that of the Random method and 25.10% less than that of the TFGR+TFAR method.

Fig. 20 shows the curves when there are 150 workers in the network. In the 20th round, the total average data quality difference of MLM-WR is reduced by 26.77% compared to the Random method and 22.83% compared to the TFGR+TFAR method. Fig. 21 shows the curves when there are 200 workers in the network. At the 20 rounds, the total average data quality difference of MLM-WR is 28.67% less than that of the Random method and 24.37% less than that of the TFGR+TFAR method. The trend of curves is similar to the case of 150 workers because the ratio of the number of locations detected by UAV to the number of locations is the same;

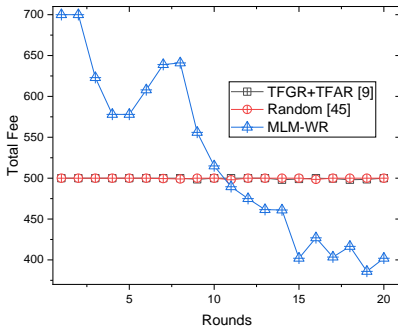


Fig. 22. Total Fee (when the number of workers = 100)

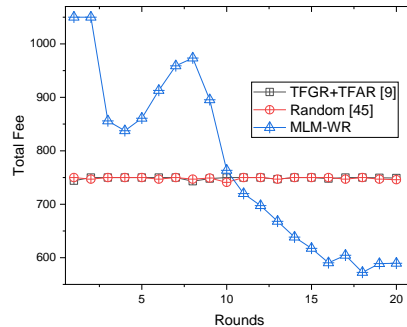


Fig. 23. Total Fee (when the number of workers = 150)

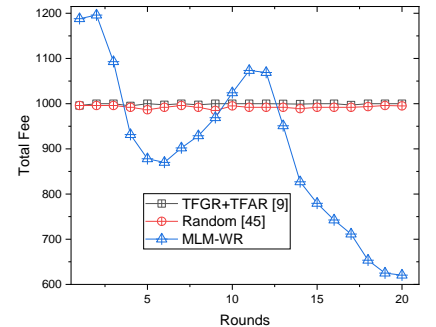


Fig. 24. Total Fee (when the number of workers = 200)

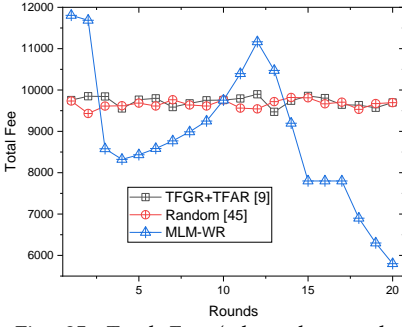


Fig. 25. Total Fee (when the number of workers = 2000)

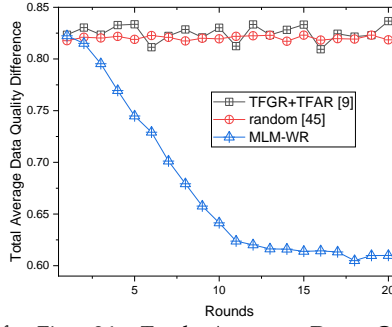


Fig. 26. Total Average Data Quality Difference (when the number of workers = 2000)

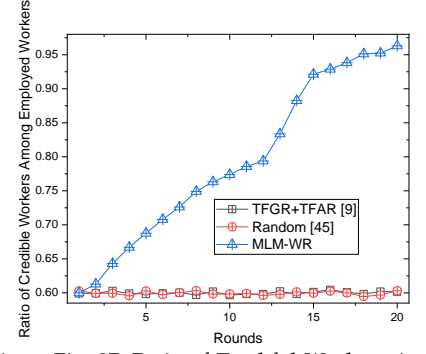


Fig. 27. Ratio of Truthful Workers Among Employed Workers (when the number of workers = 2000)

However, in the case of 200 workers, the ratio of the number of locations detected by UAV to the number of locations decreases, so the total average data quality difference decreases more slowly than the two figures on the left. TFAR+TFGR works relatively well when the number of locations is small because this method evaluates the relative credibility of locations and attributes among workers performing tasks at each location, but it lacks the absolute sensing location credibility as the sequence with high priority to mix the sensing location credibility of workers evaluated at different locations. The sensing attribute credibility of this method has the same problem. The increase in the number of locations and workers increases the randomness of mixing, so the total average data quality difference of this method is best when there are 100 workers in the network. And the total average data quality difference gradually becomes similar to the random method as the number of workers increases.

Figs. 22-24 show the change in the total fee. The cost of sending a UAV to execute one sensing task is 100, and the cost of sending a worker to execute one sensing task is 10. At the beginning of the MLM-WR method, because no truthful workers are identified, 5 workers are sent to each location, and the cost of employment is at the highest level. After a few truthful workers are identified, the number of hired workers decreases, and the cost decreases because locations, where no truthful workers or the UAV go to sense data, are arranged for 3 workers in each place. As the number of truthful workers increases, the cost goes up again because when the number of truthful workers is small, 4 additional workers need to be dispatched to each location where UAV or truthful workers sense data. As the number of identified truthful workers becomes larger, the number of additional workers that need to be dispatched to places where UAV or truthful workers sense data decreases, so the cost continues to decrease. Although the cost of MLM-WR is initially higher than that of the Random method and the TFAR+TFGR method, after a few rounds the decrease of the total fee of MLM-WR is significant.

To confirm the practicality of MLM-WR in real-world scenarios, the algorithm was executed with a network of 2,000 workers. Fig. 25 shows the total fee when there are 2000 workers in the network, and at the 20th round, MLM-WR saves 40.15% compared to the Random method and saves 40.19% compared to the TFAR+TFGR method. Fig. 26 shows the total average data quality difference when there are 2000 workers in

the network. TFAR+TFGR is effective in reducing the total average data quality difference compared to the Random method when the number of workers is small. However, when the number of workers and locations is large, because of the lack of data sensed by UAV as a criterion, the combination of sensing attribute credibility and sensing location credibility obtained from different locations is very random. Therefore there is no significant difference in the total average data quality difference between the Random method and the TFAR+TFGR. Nevertheless, MLM-WR is still able to work, and in the 20th round, the total average data quality difference of MLM-WR is 25.48% less than that of the Random method and 27.10% less than that of the TFGR+TFAR method. The percentage of truthful workers in MLM-WR is 96.30%, and the percentage of trusted workers in TFAR+TFGR is 60.18%. Fig. 27 shows the ratio of truthful workers when 2000 workers are in the network. In the 20th round, the ratio of truthful workers in MLM-WR is 96.30%, in Random is 60.31%, and in TFAR+TFGR is 60.18%. The result of MLM-WR improves by 36.12% compared to TFAR+TFGR.

5. Conclusion and Future Work

To serve future heterogeneous services and applications supported by generation multiple access, numerous scattered workers with computing power need to be identified accurately. It is necessary to exclude malicious workers and then efficiently schedule workers to complete tasks. Therefore in this paper, we propose a Multi-attribute and Local Matching based Workers Recruitment (MLM-WR) approach to improve the data quality for MCS, which realizes accurate identification of truthful workers, and accurate evaluation of workers' sensing attribute preference and sensing location preference. MLM-WR also provides a way to assign tasks with the best data quality at the least cost. In the Credibility Evaluation Model, the credibility is updated by comparing the data submitted by workers with GTD and Sub-GTD. Different mechanisms are utilized to update the credibility of workers. After identifying truthful workers, we obtain sensing attribute preference by comparing Attribute-GTD and Sub-Attribute-GTD. Absolute sensing location preference is calculated by comparing Location-GTD and Sub-Location-GTD. Then the sensing location preference is acquired by combining the Absolute Sensing Location

Preference Sequence and the Autonomy Supervised Sensing Location Preference Sequence. Finally, an algorithm is built to find the best task assignment with the consideration of total fee and data quality. Our strategy is based on Swarm Intelligence, which also combines sensing location credibility and sensing attribute credibility. Experimental results demonstrate that, compared with the previous evaluation methods, the MLM-WR can make a more accurate estimation of sensing attribute credibility and sensing location credibility. It also offers a state-of-art task assignment approach. In future works, we will consider the relationship between attributes and the relationship among locations to improve the data quality of tasks. Also, it is important to improve the incentive mechanism which helps to avoid the low participation rate. The task assignment model will consider the malicious competition among workers. For example, some workers may want to be hired by putting out the below-cost price which will inhibit the willingness of other workers to participate in the sensing tasks. So, it is necessary to avoid this kind of behavior existing. While choosing workers to do sensing tasks, we will combine the credibility of workers and the price of workers to make the whole MCS system more feasible.

Acknowledgment

This work was supported by the National Natural Science Foundation of China under Grant No. 62072475

Reference

- [1] K. Langendoen, N. Reijers. Distributed localization in wireless sensor networks: a quantitative comparison. *Computer Networks*, 43 (4) (2003) 499-518.
- [2] B. Guo, Z. Yu, X. Zhou and D. Zhang, From participatory sensing to Mobile Crowd Sensing, *IEEE International Conference on Pervasive Computing and Communication Workshops (PERCOM WORKSHOPS)*, 2014, pp. 593-598.
- [3] J. Wang, J. Tang, G. Xue, et al. Towards energy-efficient task scheduling on smartphones in mobile crowd sensing systems. *Computer Networks* 115 (2017) 100-109.
- [4] Y. Wang, Y. Gao, Y. Li, et al. A worker-selection incentive mechanism for optimizing platform-centric mobile crowdsourcing systems. *Computer Networks* 171 (2020) 107144.
- [5] J. Guo, H. Wang, W. Liu, et al. A Lightweight Verifiable Trust based Data Collection Approach for Sensor-Cloud Systems, *Journal of Systems Architecture*, 119 (2021) 102219.
- [6] T. Liu, Y. Zhu, L. Huang. TGBA: A two-phase group buying based auction mechanism for recruiting workers in mobile crowd sensing. *Computer Networks*, 149 (2019) 56-75.
- [7] N. Maisonneuve, M. Stevens, M. E. Niessen, and L. Steels. Noisetube: Measuring and mapping noise pollution with mobile phones, *Information Technologies in Environmental Engineering*, (2009) 215-228.
- [8] R. K. Rana, C. T. Chou, S. S. Kanhere, N. Bulusu, and W. Hu. Ear-phone: an end-to-end participatory urban noise mapping system, in *Proceedings of the 9th ACM/IEEE International Conference on Information Processing in Sensor Networks*, Stockholm, Sweden, April, pp. (2010) 105-116.
- [9] H. Tian, W. Sheng, H. Shen, C. Wang. Truth finding by reliability estimation on inconsistent entities for heterogeneous data sets, *Knowledge-Based Systems*, 187 (2020) 104828.
- [10] S. Ye, J. Wang, H. Fan, Z. Zhang. Probabilistic model for truth discovery with mean and median check framework, *Knowledge-Based Systems*, 233 (2021) 107482.
- [11] Y. Ren, W. Liu, A. Liu, T. Wang, A. Li. A privacy-protected intelligent crowdsourcing application of IoT based on the Reinforcement Learning, *Future generation computer systems*, 127 (2022) 56-69.
- [12] X. Zhu, Y. Luo, A. Liu, W. Tang, M. Z. A. Bhuiyan. A Deep Learning-Based Mobile Crowdsensing Scheme by Predicting Vehicle Mobility, *IEEE Transactions on Intelligent Transportation Systems*, 22 (7) (2021) 4648-4659.
- [13] Y. Qu, S. Tang, C. Dong, et al. Posted Pricing for Chance Constrained Robust Crowdsensing, *IEEE Transactions on Mobile Computing*, 19 (1) (2020) 188-199.
- [14] J. Goncalves, S. Hosio, J. Rogstadius, et al. Motivating participation and improving quality of contribution in ubiquitous crowdsourcing. *Computer Networks*, 90 (2015) 34-48.
- [15] I. Boutsis and V. Kalogeraki, On Task Assignment for Real-Time Reliable Crowdsourcing, *2014 IEEE 34th International Conference on Distributed Computing Systems*, 2014, pp. 1-10.
- [16] D. R. Karger, S. Oh, and D. Shah. Efficient crowdsourcing for multi-class labeling, *Proceedings of the ACM SIGMETRICS/international conference on Measurement and modeling of computer systems*. 2013.
- [17] T. Luo, J. Huang, S. S. Kanhere, J. Zhang and S. K. Das, Improving IoT Data Quality in Mobile Crowd Sensing: A Cross Validation Approach, *IEEE Internet of Things Journal*, 6 (3) (2019) 5651-5664.
- [18] Y. Wen, J. Shi, Q. Zhang, et al. Quality-Driven Auction-Based Incentive Mechanism for Mobile Crowd Sensing, *IEEE Transactions on Vehicular Technology*, 64 (9) (2014) 4203-4214.
- [19] X. Liu, J. Fu, Y. Chen, et al. Trust-Aware sensing Quality estimation for team Crowdsourcing in social IoT. *Computer Networks* 184 (2021) 107695.
- [20] K. Ota, M. Dong, J. Gui and A. Liu, QUOIN: Incentive Mechanisms for Crowd Sensing Networks, *IEEE Network Magazine*, 32 (2) (2018) 114-119.
- [21] C. Meng, W. Jiang, Y. Li, et al. Truth discovery on crowd sensing of correlated entities, in *Proceedings of the 13th acm conference on embedded networked sensor systems*. 2015, pp. 169-182.
- [22] Y. Liu, A. Liu, T. Wang, X. Liu, N. N. Xiong. An intelligent incentive mechanism for coverage of data collection in cognitive Internet of Things, *Future Generation Computer Systems*, 100 (2019) 701-714.
- [23] P. Sun, Z. Wang, L. Wu, Y. Feng, X. Pang, et al. Towards Personalized Privacy-Preserving Incentive for Truth Discovery in Mobile Crowdsensing Systems, *IEEE Transactions on Mobile Computing*, 21 (1) (2020) 352-365.
- [24] J. Xu, C. Guan, H. Dai, D. Yang, L. Xu and J. Kai, Incentive Mechanisms for Spatio-Temporal Tasks in Mobile Crowdsensing, in *2019 IEEE 16th International Conference on Mobile Ad Hoc and Sensor Systems (MASS)*, 2019, pp. 55-63.
- [25] C.-J. Ho, S. Jabbari, and J. W. Vaughan. Adaptive task assignment for crowdsourced classification. *International Conference on Machine Learning*. PMLR, 2013.
- [26] H. Jin, L. Su, and K. Nahrstedt. Theseus: Incentivizing truth discovery in mobile crowd sensing systems, in *Proceedings of the 18th ACM International Symposium on Mobile Ad Hoc Networking and Computing*. 2017.
- [27] H. Cai, Y. Zhu, and Z. Feng. A truthful incentive mechanism for mobile crowd sensing with location-Sensitive weighted tasks. *Computer Networks*, 132 (2018) 1-14.

- [28] P. Inkeaw, P. Udomwong, J. Chaijaruwanich. Density based semi-automatic labeling on multi-feature representations for ground truth generation: Application to handwritten character recognition, *Knowledge-Based Systems*, 220 (2021) 106953.
- [29] T. Li, A. Liu, N. N. Xiong, et al. A trustworthiness-based vehicular recruitment scheme for information collections in distributed networked systems, *Information sciences*, 545 (2021) 65-81.
- [30] Y. Zheng, G. Li, Y. Li, et al. Truth inference in crowdsourcing: Is the problem solved?, in *Proceedings of the VLDB Endowment*, 2017, pp. 541-552.
- [31] C. Zhao, S. Yang and J. A. McCann, On the Data Quality in Privacy-Preserving Mobile Crowdsensing Systems with Untruthful Reporting, *IEEE Transactions on Mobile Computing*, 20 (2) (2021) 647-661.
- [32] X. Yin, J. Han, and P. S. Yu. Truth discovery with multiple conflicting information providers on the web, in *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '07)*. Association for Computing Machinery, New York, NY, USA, 2007, pp. 1048-1052.
- [33] Xiu Susie Fang, Quan Z. Sheng, Xianzhi Wang, Wei Emma Zhang, Anne H. H. Ngu, and Jian Yang. 2020. From appearance to essence: comparing truth discovery methods without using ground truth. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11 (6) (2020) 1-24.
- [34] Y. Du, YE. Sun, H. Huang, et al. Bayesian Co-Clustering Truth Discovery for Mobile Crowd Sensing Systems, *IEEE Transactions on Industrial Informatics*, 16 (2) (2019) 1045-1057.
- [35] J. Guo, A. Liu, K. Ota, M. Dong, X. Deng and N. N. Xiong, ITCN: An Intelligent Trust Collaboration Network System in IoT, *IEEE Transactions on Network Science and Engineering*, 9 (1) (2022) 203-218.
- [36] J. Picaut, A. Boumchich, E. Bocher, et al. A Smartphone-Based Crowd-Sourced Database for Environmental Noise Assessment, *International Journal of Environmental Research and Public Health*, 18 (15) (2021) 7777.
- [37] W. Fourati, B. Friedrich. A method for using crowd-sourced trajectories to construct control-independent fundamental diagrams at signalized links, *Transportation research part C: emerging technologies*, 130 (2021) 103270.
- [38] Y. Ren, Y. Liu, N. Zhang, et al. Minimum-Cost Mobile Crowdsourcing with QoS Guarantee Using Matrix Completion Technique, *Pervasive and Mobile Computing*, 49 (2018) 23-44.
- [39] X. Kan, Y. Fan, Z. Fang, et al. A novel IoT network intrusion detection approach based on adaptive particle swarm optimization convolutional neural network, *Information Sciences*, 568 (2021), 147-162.
- [40] T. Huang, Y. Chen, B. Yao, et al. Adversarial attacks on deep-learning-based radar range profile target recognition, *Information Sciences*, 531 (2020) 159-176.
- [41] M. Huang, A. Liu, N. N. Xiong and J. Wu, A UAV-Assisted Ubiquitous Trust Communication System in 5G and Beyond Networks, *IEEE Journal on Selected Areas in Communications*, 39 (11) (2021) 3444-3458.
- [42] S. Olariu, S. Mohrehkesh, X. Wang, et al. On aggregating information in actor networks, *ACM SIGMOBILE Mobile Computing and Communications Review*, 18 (1) (2014) 85-96.
- [43] R. K. Ganti, F. Ye and H. Lei, Mobile crowdsensing: current state and future challenges, *IEEE Communications Magazine*, 49 (11) (2021) 32-39.
- [44] <https://www.kaggle.com/datasets/mathchi/diabetes-data-set>
- [45] C.-J. Ho, J. Vaughan. Online task assignment in crowdsourcing markets, in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2012, pp. 45-51.



Jiaheng Lu is with the School of Computer Science and Engineering, Central South University, China. His research interests include wireless sensor network. E-mail: jiahenglucsu.edu.cn.



Zhenzhe Qu received his master degree in the school of software, Central South University, China, in 2019. He is currently pursuing his Ph.D. degree in the school of Computer Science and Engineering, Central South University, China. His research interests include edge computing. E-mail: zhenzheQu@csu.edu.cn.



Anfeng Liu received the M.Sc. and Ph.D. degrees from Central South University, China, in 2002 and 2005, respectively, both in computer science. He is currently a professor of the School of Computer Science and Engineering, Central South University, China. His major research interests include Wireless Sensor Networks, Internet of Things, Information Security, Edge Computing and Crowdsourcing. Dr. Liu has published 4 books and over 100 international journal and conference papers, among which there are more than 30 ESI highly-cited papers. He was a recipient of the First Prize of Scientific Research Achievement of Colleges from the Ministry of Education of China in 2016. E-mail: afengliu@mail.csu.edu.cn.



Shaobo Zhang received the B.Sc. and M.Sc. degree in computer science both from Hunan University of Science and Technology, Xiangtan, China, in 2003 and 2009 respectively, and the Ph.D. degree in computer science from Central South University, Changsha, China, in 2017. He is currently an associate professor at School of Computer Science and Engineering of the Hunan University of Science and Technology, China. His research interests include privacy and security issues in social networks and cloud computing. E-mail: shaobozhang@hnust.edu.cn.



Neal N. Xiong (S'05-M'08-SM'12) is current an Associate Professor (4 year credits) at Department of Computer Science and Mathematics, Sul Ross State University, Alpine, TX 79830, USA. He received his both PhD degrees in Wuhan University (2007, about sensor system engineering), and Japan Advanced Institute of Science and Technology (2008, about dependable

communication networks), respectively. Before he attended Sul Ross State University, he worked in Georgia State University, Northeastern State University, and Colorado Technical University (**full professor about 5 years**) about 15 years. His research interests include Cloud Computing, Security and Dependability, Parallel and Distributed Computing, Networks, and Optimization Theory. Dr. Xiong published over 200 IEEE journal papers and over 200 international conference papers. Some of his works were published in IEEE JSAC, IEEE or ACM transactions, ACM Sigcomm workshop, IEEE INFOCOM, ICDCS, and IPDPS. He is serving as an Editor-in-Chief, Associate editor or Editor member for over 10 international journals (including Associate Editor for IEEE Tran. on Systems, Man & Cybernetics: Systems, IEEE Tran. on Network Science and Engineering, Information Science). Dr. Xiong is the Chair of “Trusted Cloud Computing” Task Force, IEEE Computational Intelligence Society (CIS), and he is a Senior member of IEEE Computer Society from 2012, E-mail: xionгнаixue@gmail.com.



Tian Wang received his BSc and MSc degrees in Computer Science from Central South University in 2004 and 2007. He received his Ph.D. degree at the City University of Hong Kong in 2011. Currently, he is a professor at the Artificial Intelligence and Future Networks, Beijing Normal University & UIC, China. His research interests include the internet of things, edge computing, and mobile computing. E-mail: tianwang@bnu.edu.cn.



Jiaheng Lu is with the School of Computer Science and Engineering, Central South University, China. His research interests include wireless sensor network. E-mail: jiahenglu@csu.edu.cn.



Zhenzhe Qu received his master degree in the school of software, Central South University, China, in 2019. He is currently pursuing his Ph.D. degree in the school of Computer Science and Engineering, Central South University, China. His research interests include edge computing. E-mail: zhenzheQu@csu.edu.cn.



Anfeng Liu received the M.Sc. and Ph.D. degrees from Central South University, China, in 2002 and 2005, respectively, both in computer science. He is currently a professor of the School of Computer Science and Engineering, Central South University, China. His major research interests include Wireless Sensor Networks, Internet of Things, Information Security, Edge Computing and Crowdsourcing. Dr. Liu has published 4 books and over 100 international journal and conference papers, among which there are more than 30 ESI highly-cited papers. He was a recipient of the First Prize of Scientific Research Achievement of Colleges from the Ministry of Education of China in 2016. E-mail: afengliu@mail.csu.edu.cn.



Shaobo Zhang received the B.Sc. and M.Sc. degree in computer science both from Hunan University of Science and Technology, Xiangtan, China, in 2003 and 2009 respectively, and the Ph.D. degree in computer science from Central South University, Changsha, China, in 2017. He is currently an associate professor at School of Computer Science and Engineering of the Hunan University of Science and Technology, China. His research interests include privacy and security issues in social networks and cloud computing. E-mail: shaobozhang@hnust.edu.cn.



Neal N. Xiong (S'05–M'08–SM'12) is current an Associate Professor (4 year credits) at Department of Computer Science and Mathematics, Sul Ross State University, Alpine, TX 79830, USA. He received his both PhD degrees in Wuhan University (2007, about sensor system engineering), and Japan Advanced Institute of Science and Technology (2008, about dependable communication networks), respectively. Before he attended Sul Ross State University, he worked in Georgia State University, Northeastern State University, and Colorado Technical University (full professor about 5 years) about 15 years. His research interests include Cloud Computing, Security and Dependability, Parallel and Distributed Computing, Networks, and Optimization Theory. Dr. Xiong published over 200 IEEE journal papers and over 200 international conference papers. Some of his works were published in IEEE JSAC, IEEE or ACM transactions, ACM Sigcomm workshop, IEEE INFOCOM, ICDCS, and IPDPS. He is serving as an Editor-in-Chief, Associate editor or Editor member for over 10 international journals (including Associate Editor for IEEE Tran. on Systems, Man & Cybernetics: Systems, IEEE Tran. on Network Science and Engineering, Information Science). Dr. Xiong

is the Chair of “Trusted Cloud Computing” Task Force, IEEE Computational Intelligence Society (CIS), and he is a Senior member of IEEE Computer Society from 2012, E-mail: xionгнаixue@gmail.com.



Tian Wang received his BSc and MSc degrees in Computer Science from Central South University in 2004 and 2007. He received his Ph.D. degree at the City University of Hong Kong in 2011. Currently, he is a professor at the Artificial Intelligence and Future Networks, Beijing Normal University & UIC, China. His research interests include the internet of things, edge computing, and mobile computing. E-mail: tianwang@bnu.edu.cn.



Jiaheng Lu



Zhenzhe Qu



Anfeng Liu



Shaobo Zhang



Neal N. Xiong



Tian Wang

Conflict of interest statement

We declare that we have no financial and personal relationships with other people or organizations that can inappropriately influence our work, there is no professional or other personal interest of any nature or kind in any product, service and/or company that could be construed as influencing the position presented in, or the review of, the manuscript entitled, “MLM-WR: An Effective workers recruitment Scheme based on Swarm Intelligence for Mobile Crowd Sensing”.