

# Network Size Estimation in Opportunistic Mobile Networks: The Mark-Recapture Method

Bambang Soelistijanto\* and Geraldev Manoah  
Department of Informatics, Sanata Dharma University  
b.soelistijanto@usd.ac.id, udevmanoah@gmail.com

Received: July 25, 2020; Accepted: August 14, 2021; Published: September 30, 2021

## Abstract

This article addresses the issues of counting the number of nodes in opportunistic mobile networks (OMNs). The global knowledge of network size is commonly required to design optimal routing algorithms in OMNs. However, due to the inherent characteristic of long transfer delay, node counting in such intermittently-connected networks is a challenging task. In this paper, we propose the Mark-Recapture method to estimate the number of nodes in a network. In ecology, the statistical technique has been widely used to predict the population sizes of animals in open areas. The scheme initially samples nodes in the network, and an estimate of the network size is then calculated based on this partial knowledge of the network. Through extensive simulations driven by random movement and realistic mobility models, we show that the proposed method is able to produce a good estimate of network size within a relatively short duration of time. Finally, by tweaking Epidemic routing with the local estimate of network size, we can reduce the delivery cost of this flooding strategy without significantly degrading the overall network delivery performances.

**Keywords:** network size, node counting, Mark-Recapture, opportunistic mobile networks

## 1 Introduction

Nowadays, opportunistic mobile networks (OMNs) [16] have received much attention by industry and research community. These networks are an extension of mobile ad-hoc networks (MANETs), and are an instance of delay tolerant networks (DTNs). While MANETs require end-to-end paths between sources and destinations to enable message transfer, OMNs are capable of performing communication despite the absence of stable paths between any pair of nodes. In MANETs node movement is considered as a potential disruption, but in OMNs data transfer is performed by opportunistic communication, leading to a higher delay than that of MANETs. Data dissemination in OMNs is thus delay-tolerant in nature. Some realizations of OMNs exist, including emergency scenarios and natural disasters [21], military operations [15], and social-based networks [28]. The widely use of mobile wireless devices, such as smart phones, gadgets, and laptops, is the main factor in the proliferation of these systems.

In OMNs, searching for optimal paths between a pair of nodes is a non-trivial task. Since the stable paths between any pair of nodes rarely exist at all the time, conventional routing algorithms proposed for MANETs would fail in this setting. This imposes a new model for routing in OMNs, the *store-carry-forward* paradigm [2]. This suggests that a message is stored and carried by relay nodes, and finally is forwarded when the destination is encountered. In this regard, choosing good relays for message transfers is indeed crucial in OMNs. A bulk of researches in OMNs have focused on developing effective routing

---

*Journal of Wireless Mobile Networks, Ubiquitous Computing, and Dependable Applications (JoWUA)*, 12(3):29-46, Sept. 2021  
DOI:10.22667/JOWUA.2021.09.30.029

\*Corresponding author: Department of Informatics, Sanata Dharma University, Yogyakarta, 55282, Indonesia, Tel: +62-819-317-27372

protocols. To achieve this goal, the algorithms typically require complete information of the current network states. In practice, however, this global knowledge is commonly unavailable to all the network nodes. To improve the delivery performance, several algorithms opt to increase message redundancy in the network. Naive approaches (e.g., Epidemic routing [26]) forward a message replica to each contacted node, so that the copies are quickly dispersed over the network. This oblivious forwarding assumes unlimited node resources, but this is hard to achieve in practice. On the other hand, some algorithms (e.g., adaptive Spray-Wait [6]) attempt to reduce the number of message replicas by capping the message replication at a maximum value. To this aim, the protocol at each node needs to know the number of nodes in the network. However, estimating this global parameter in a decentralized manner is a difficult task in OMNs, due to the highly dynamic topology changes and long transfer delays.

In the present work, we focus on the particular case of node counting in OMNs (the global statistic of the total number of nodes in a network is also referred to as *network size*). To date, distributed node counting has attracted interest from researchers, since a local estimate of network size is often very useful for building applications that are adaptive and robust. For example, the population algorithm in [5] uses a random sample to estimate the size of a large network and its communities; a crowd counting system in [7] estimates crowd sizes and densities for city administration and disaster management; a data dissemination protocol in [4] predicts the network size for limiting message redundancy. In the literature, several distributed computing algorithms have been proposed in the area of global information collection and estimation in opportunistic networks. In addition, majority of them are modifications of data aggregation schemes proposed for well-connected networks (e.g., [17]). Even though Aggregation provides accurate estimates in the conventional networks, but it suffers from a number of difficulties in the context of OMNs as follows [14]: first, the delay time to converge to the actual network size is very long in such delay-tolerant networks; second, node failures will significantly degrade the performance of Aggregation. As an alternative to Aggregation, several distributed estimation algorithms for OMNs (e.g., [4, 23, 1]) are developed based on statistical sampling techniques.

In this paper, we propose the Mark-Recapture method [8], a statistical technique used to estimate the number of nodes in an OMN. This technique has been widely used in ecology to predict the population sizes of animals or fishes in forests or seas, respectively. In the area of communication networks, the method has been utilized to estimate the network size in peer-to-peer (P2P) networks as well as multicast networks [1]. To the best of our knowledge, however, this paper is the first work that applies Mark-Recapture to perform node counting in OMNs. In addition, most of the existing works in distributed node counting in OMNs only consider a simple random *i.i.d* model when designing and evaluating the algorithms. In fact, such model may not be realistic to describe real human mobility cases [3]. In this paper, we investigate the proposed algorithm under both random movement and realistic mobility scenarios. The underlying node mobility contributes to node mixing, and in turn to the spreading of data. Consequently, the most important questions we answer in this paper are:

- How does node mobility impact the performance of the Mark-Recapture counting algorithm in OMNs?
- Can Mark-Recapture outperform Aggregation in OMNs in terms of estimation accuracy and convergence time?
- Can a local estimate of network size improve the delivery performance of Epidemic routing [26] in OMNs?

The main contributions of this paper are:

- We present a distributed counting algorithm based on Mark-Recapture [8] to estimate the number of nodes in an OMN.

- We evaluate the proposed algorithm via extensive simulations driven by random movement and real-life mobility models.
- We identify the performance improvement of Mark-Recapture compared to Aggregation in terms of estimation accuracy and convergence time.
- Using local estimates of network size, we improve the delivery cost performance of Epidemic routing without significantly degrading the overall network delivery performances.

The remainder of the paper is organized as follows. In Section 2, we introduce the related works and position our work concerning the state-of-the-arts in the area of node counting in OMNs. The problem description and the proposed distributed counting algorithm based on Mark-Recapture are presented in Section 3. In Section 4, we evaluate the estimation accuracy and convergence time of the scheme in OMNs through simulations under random movement and realistic mobility scenarios. Subsequently, we compare the performance of Mark-Recapture with that of Aggregation in terms of estimation accuracy and convergence time. Finally, we investigate the delivery performance improvement of Epidemic routing with local estimates of network size. We conclude the paper and present directions for future work in the last section.

## 2 Related Work

In this section, we review some state-of-the-art node counting algorithms for OMNs to indicate our motivations and contributions. We discuss the existing works that are related to our work in the following three categories.

### 2.1 Distributed algorithms for node counting

We can broadly distinguish two classes of methods for node counting in OMNs. Techniques of the first type are based on *data aggregation algorithms*, while those of the second type are based on *statistical sampling algorithms*.

#### 2.1.1 Aggregation algorithms

To date, Aggregation has played an important role in modern distributed systems [18]. It can perform the evaluation of global properties of the systems in a decentralized way. Moreover, network size is a typical system-wide property required by algorithms in many contexts. Jelasity *et al.*[17] proposed a distributed gossip-based aggregation algorithm for large dynamic networks. In this algorithm, each node periodically chooses one node among the neighbours, and afterwards the pair of nodes exchange and update their local estimates to assure quick convergence to the desired aggregate value. Since the scheme was developed under the assumption of stable links, it will not work properly in the context of opportunistic communication. In OMNs, links between nodes are created by sporadic contacts, occurring when they come in direct radio range. Consequently, the list neighbours is often not known in advance, and there is no neighbour sampling before links are established between the node and its neighbours.

Guerrieri *et al.*[14] introduce a set of node counting strategies based on Aggregation for OMNs, namely pairwise average and population protocols. The former is a class of gossip protocols. At the beginning, one node (called an initiator) stores a value equals to “1” and all the remaining nodes stores “0”. At every contact, the nodes exchange their current values and update the stored value as the average of its value and the peer’s. Eventually, the algorithm converges to  $1/N$  and the number of the network nodes is achieved as the inverse of the estimate. In contrast, the population protocols use tokens to

calculate network size. At the initial run, each node is allocated a single token. At each contact, two nodes toss a fair coin and the one winning the ballot collects the whole peer's tokens. At the end, tokens gather on a node that has the accurate estimation of the network's size. However, randomly choosing nodes to collect tokens during node contacts may result in suboptimal performance: it leads to long convergence time and low estimation accuracy. To cope with these issues, Ning *et al.* [24] therefore propose a new technique that incorporates effective contact probability into counting process. On the other hand, the works in [7, 20] apply a different strategy based on Aggregation of node states. When two nodes come into contact, they exchange the state sets and each node then establishes a union set containing the elements of both its own set and the peer's. In the end, all nodes converge to have a set including the ids of all nodes in the network, and the network size is determined by the cardinality of the set.

However, all the abovementioned Aggregation schemes suffer from common problems in OMNs, namely long convergence time and estimation accuracy sensitive to node failures. To deal with these issues, we propose a node counting algorithm based on a statistical sampling technique, Mark-Recapture [8]. In the following, we discuss the existing works that are developed based on statistical sampling techniques.

### 2.1.2 Statistical sampling algorithms

Statistical sampling methods produce a prediction of the system's global properties based on the statistics attained from uniformly random samples. Sample-Collide [23] is proposed to calculate peer counting in overlay networks. The work in [4] applies a sampling technique based on Taxi-Problem (also known as Racing-Car Problem) to predict the number of active nodes in an OMN. In general, Taxi-Problem works as follows: one (also called initiator) wishes to estimate the number of taxis currently operating on the streets of a city. The taxis are numbered consecutively from one to some unknown number  $N$ . The initiator observes and records the ids (=serial numbers) of all taxis that have passed in a given time interval. In addition, this scheme assumes that each taxi is equally likely to pass the initiator at any given time. Using the sampling data, an unbiased minimum variance estimator (UMVE) is finally computed as the best estimate of the total taxis in the given city. As shown in [4], the counting algorithm based on Taxi-Problem can work properly in OMNs to give a good estimate of network size. Despite the elegance of this technique, however, the effectiveness of Taxi-Problem in OMNs strictly depends on two conditions as follows: first, all nodes are consecutively numbered from 1 to  $N$ ; second, the probability of encounter between any pair of nodes is uniformly distributed in the network. As opposed to Taxi-Problem, our proposed algorithm is relaxed from these constraints: it does not require the nodes either to be successively numbered or to have a homogeneous contact pattern.

## 2.2 Node mobility considerations

In the literatures, majority of the node counting algorithms proposed for OMNs are developed under the assumption of a simple random *i.i.d* model. In the class of Aggregation, for example, the gossip-based pair-wise average method [14] suggests that in each pair-wise contact nodes exchange their current values and store the new value as the average of their present values. Given that all nodes have an equal opportunity to meet any other node in the network, the algorithms of all the nodes eventually converge to a single values of actual network size. Moreover, the work of crowd counting in [7] proposes a fully decentralized Aggregation to calculate an accurate estimate of the crowd size. During a node contact, two nodes exchange their state sets containing the identities of the nodes already seen before. By assuming that all nodes (individuals) in the crowd follow a random walk (RW) mobility, they finally converge to have a set that includes the ids of all nodes in the crowd and the crowd size is then determined by the

cardinality of the set.

As similar to Aggregation, most of the existing works of node counting based on statistical sampling methods also rely on the assumption of random mobility. For instance, the work based on Taxi-Problem [4] strictly requires that the node contact pattern should be homogenous, so that the probability of any node encountering the initiator will be equal. In fact, however, real-life mobility deviates from the assumption of random *i.i.d.* mobility [3]. In [20], Li *et al.* study the effect of node mobility on data collection and node counting in OMNs. However, their investigation is still based on a homogeneous mobility pattern, where each node randomly selects the destination and speed: the destination follows a uniform distribution, but the speed follows a Gaussian distribution with the mean is constant, but the standard deviation varies during the experiment. Our proposed algorithm, however, is investigated under both random movement and realistic mobility scenarios. For the latter case, we use real human mobility models, which intrinsically possess a heterogeneous contact pattern [27], where a few nodes (called hub nodes) have many contacts with others, but majority of the network nodes only have few ones.

### 2.3 Applications of node counting

With the more powerful mobile wireless devices nowadays, it is not required to offload the processing to an edge server or a cloud computing service. In mobile computing, a computational task is executed independently in each node (mobile device), and by using communication all nodes share their individual outcomes and ultimately arrive at a convergence result. One of the typical tasks in distributed computing is calculating network size (i.e., the number of nodes in the network). This information can then be used as input by other applications or protocols. Some applications of node counting are as follows: network size is used for building and maintaining the distributed hash table in P2P networks [1]; in [22] the statistic is exploited in wireless mobile ad hoc networks (MANETs) to set up a quorum of a membership service; UrbanCount [7] applies a fully distributed crowd counting protocol to estimate crowd size during open-air events or rush hours for city administration; in [4] the knowledge of network size is required to optimize the performance of a routing algorithm in OMNs by minimizing the delivery cost. In this paper, we use local estimates of network size to improve the delivery cost performance of a flooding-based algorithm, Epidemic routing [26], by capping the message replicas to be a half of the network size.

## 3 The Mark-Recapture Distributed Estimation Scheme

In this section, we propose a novel strategy of distributed estimation based on the Mark-Recapture technique to predict the number of nodes in an OMN. We initially introduce the basic scheme of Mark-Recapture widely used in ecology. We then discuss the system model and problem description and finally propose the Mark-Recapture distributed estimation algorithm for OMNs.

### 3.1 The Basic Mark-Recapture Method

Wildlife managers commonly use the Mark-Recapture technique [8] (also called the Lincoln-Petersen method) to estimate the population size of animals or fishes in forests or seas before hunting or fishing seasons, respectively. The scheme comprises a single marking episode (also called a capture episode) and a single recapture episode. It initially starts with taking a sample of individuals in a natural population, marking and then sending them back to the original population and finally recapturing some of them as a basis for predicting the population size at the time of initial marking. The basic principle of the algorithm is that if a sample of the population is marked in some way, returned them to the original population, and after fully dispersed in the population a second sample (also called a recapture sample) is taken, the ratio

of total marked individuals ( $m$ ) to sample size ( $C$ ) in the recapture sample will be equal to the ratio of total marked individuals in the initial sample ( $M$ ) to the population size ( $N$ ). That is,

$$\frac{m(\text{total\_marked\_in\_recapture})}{C(\text{recapture\_size})} = \frac{M(\text{total\_marked\_initially})}{N(\text{population\_size})} \quad (1)$$

By rearrangement (1), we can calculate the estimate of the population's size at the time of initial marking, as

$$\hat{N} = \frac{MC}{m} \quad (2)$$

However, the accuracy of Mark-Recapture relies on several assumptions as follows:

- The population size should be constant during the period between the initial marking episode and the recapture episode.
- The probability of all individuals being captured should be the same during both the episodes.
- There must be sufficient time between the capture and recapture periods to allow all the marked individuals to be randomly mixed all over the population.
- The marked individuals should not lose their marks between the two periods.

### 3.2 System Model and Problem Description

We consider an opportunistic mobile network, where the nodes move independently in a given area and communicate to the peers wirelessly. Communication occurs when nodes come into contact within their radio ranges. Our study is based on several assumptions as follows:

- There are  $N$  mobile nodes in the network.
- Nodes participate equally in the counting process.
- Nodes do not provide fake information to others.
- Nodes do not stop operations or abruptly leave the network all the time.
- Any node can initiate a counting process whenever it needs to know the network size.

The purpose of this study is to make a prediction on the number of nodes in an OMN with high accuracy and a low delay. This particularly becomes a complex task in OMNs, since the node contacts are unpredictable and are limited in terms of time and bandwidth. Furthermore, this paper considers node counting in a closed system. In this setting, the number of nodes is fixed but unknown and needs to be predicted. Different scenarios may allow nodes to enter and leave the area (called an open system

initiator-id	total marks	seq. number	TTL
--------------	-------------	-------------	-----

Figure 1: The marking message structure

**Algorithm 1:** The Mark-Recapture Distributed Algorithm

```

/* initialization */
totalMarks ← M;
ttl ← TTL;
round ← 0;
estimates ← 0;
seqNum ← 0;
/* initial marking phase */
begin
  initiatorID ← InitiatorSerialNumber;
  /* starts a new counting round by creating a marking message s */
  if initiator then
    createMarkingMessage(s.initiatorID, s.totalMarks, s.seqNum ++, s.ttl);
    round++;
  end
  /* marking the encountered nodes until TotalMarks=1 */
  if contactedNode.marked=false and s.totalMarks > 1 then
    sendMarkingMessage(s.totalMarks = ⌊M/2⌋);
    updateMarkingMessage(s.totalMarks = ⌈M/2⌋);
  end
end
/* recapture phase */
begin
  recapture ← 0;
  markedNode ← 0;
  /* when contact occurring with a node B */
  if initiator and s.ttl > 0 then
    if ¬recapture.contains(B) then recapture.add(B);
    if B.marked = true then markedNode++;
  end
end
/* the counting round terminates */
begin
  if initiator and s.ttl = 0 then
    estimate=calculateEstimate();/* using(2) */
    estimates(round).add(estimate);
  end
  /* calculate final estimate */
  finalEstimate=avg(estimates());
end

```

with node churn). However, as previously demonstrated in [7], node counting in an open system is more challenging, and providing an accurate count is not trivial. Therefore, we restrict the discussion in this paper to the case of closed systems, and all kinds of opportunistic network applications that meet the requirements of Mark-Recapture mentioned above can use our proposed algorithm for estimating network size.

### 3.3 The Proposed Algorithm

We now discuss the proposed distributed estimation algorithm based on Mark-Recapture for OMNs. We divide the algorithm into two phases: *initial marking* and *recapture*. For node marking, we firstly define a marking message (Fig.1) as a small (control) message containing a number of variables: *initiator-id*, *total-marks*, *sequence-number*, and *TTL*. Initiator-id represents the identity of a node that initiates counting; total-marks is the maximum number of nodes that can be marked during the marking session; sequence-number is the unique identity of a marking message, incrementing by one for each new counting initiation; finally, TTL is the time-to-life of a marking message which directly represents the duration of a single counting round.

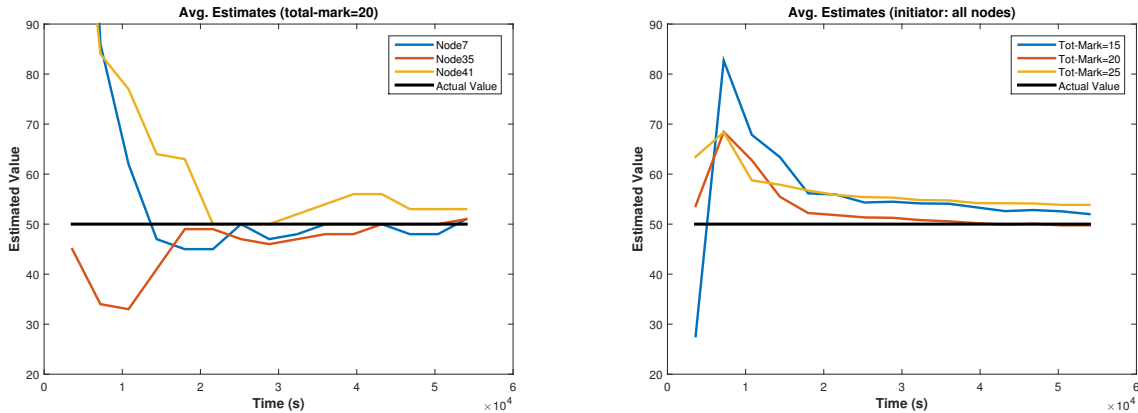
The marking episode starts when a node (called an *initiator*) initiates a counting process by creating a new marking message. When the initiator encounters another node, it marks the contacted node by sending a copy of the message to the peer. Moreover, we assume that the marked nodes do not drop the marking message before the message TTL expires. For the marking process, we have two possible strategies: first, only the initiator itself can mark the encountered nodes; indeed, this strategy is simple but takes a long time to completely perform node marking in an OMN. To speed up the process, the second strategy, we call it *binary-marking*, allows the already marked nodes to help the initiator to perform node marking: the initiator initially starts with a marking message with total-marks is set to  $M$  marks; when any node  $A$  (either the initiator or the marked node) that has total-marks  $m > 1$  encounters another node  $B$  that has not yet been marked,  $A$  then forwards the copy of the marking message to  $B$  with total-marks  $\lfloor m/2 \rfloor$  and keeps  $\lceil m/2 \rceil$  for itself; if the total-marks is left with only one mark, the node terminates marking other contacted nodes; particularly, when this case happens in the initiator, the algorithm subsequently switches the marking phase to a recapture phase (in this algorithm, we assume that only the initiator itself is able to perform the recapture process).

Before commencing a recapture episode, the initiator must wait for some time to allow all the marked nodes to be randomly dispersed over the network. During the recapture period, at each contact the initiator records the id of the encountered node and then categorizes it into a marked or unmarked node: if the contacted node has the marking message with the id-initiator matches with the id of the initiator, the initiator then increments the marked-node counter. When the message TTL expires, the recapture episode finishes, and in turn the current counting round completely ends. In future, the initiator can launch another counting round by initially creating a new marking message with a unique sequence number (that is, the algorithm increments the sequence number by one for each new marking message creation). At the end of each counting round, the initiator computes the total number of nodes in the network using (2). The algorithm eventually returns the final estimate of network size as an average of the estimates obtained from the all previous counting rounds. We depict the pseudo-code of the Mark-Recapture distributed algorithm in Alg. 1.

## 4 Performance Evaluation of the Proposed Algorithm

In this section, we evaluate the performance of the Mark-Recapture distributed algorithm in OMNs. Initially, we conduct extensive simulations to investigate the estimation accuracy and convergence time of the proposed algorithm. Subsequently, we examine the performance improvement of Epidemic routing with local estimates of network size. In this study, we use the ONE simulator [19], a discrete-event simulator for delay-tolerant networks. For simulation's mobility scenarios, we consider both random movement and real human mobility models. For the former case, we use the Random-Walk (RW) model packaged along with the ONE simulator. For the latter one, we consider two real human contact data traces, namely Huggle [25] and Reality [11], which represent the short-term and long-term human mobility models, respectively. The Huggle trace captured the activities of 41 participants during the 2005



Figure 2: Node counting in the random mobility scenario ( $N=50$ )

Infocomm conference lasted for 3 days in Miami, USA. However, Reality trace logged the activities of 97 students and staffs at the MIT campus during one academic year. The study was actually performed around 10 months.

#### 4.1 The Estimation Accuracy and Convergence Time of the Mark-Recapture Distributed Algorithm

In this section, we discuss the accuracy and convergence time of Mark-Recapture in estimating the network size of an OMN. We initially consider the random movement scenario. In Fig. 2, we show the simulation results of Mark-Recapture that estimates the number of nodes in an OMN in the random case. In this setting, the total nodes in the network ( $N$ ) is 50 nodes, the node mobility speed ( $v$ ) is 1.5-2.5 m/s, and the simulation area is 5000 X 5000 m<sup>2</sup>. We randomly choose nodes in the network as initiators (e.g., node ids 7, 35 and 41) and subsequently depict the counting results of these nodes with respect to simulation time in Fig. 2 (left) for total-marks=20. We see that the average estimates of all the given nodes eventually converge to the actual network size ( $N=50$  nodes) at nearly the same time ( $\approx 18,000$  sec or 5 hours). In Random-Walk, the probability of node contact is identically, independently distributed (*i.i.d*) in the network. All nodes therefore have the same probability of being captured in both the initial marking and recapture periods. As a result, as shown in Fig. 2 (left) the counting algorithms of all the nodes show similar counting performances, in terms of accuracy and convergence time. Afterwards, in Fig. 2 (right) we describe the effect of total-mark values on the algorithm's performance when all the network nodes simultaneously perform node counting. We notice that for network with  $N=50$ , total-marks=20 gives the best performance among the others in terms of both accuracy and convergence time. Nevertheless, the performance differences among the given total-marks are insignificant in this random scenario, and all of them are eventually able to converge to the actual network size at a slightly different time. Finally, in Fig. 3 we compare the performance of Mark-Recapture with that of an Aggregation scheme, i.e. the Pair-Wise Average method [14] (hereafter, we call it *PW-Avg* for short), in the random scenario (the brief discussion of how *PW-Avg* works is given in Section 2.1.1). We again randomly select nodes in the network as initiators and then run the simulations by successively applying both the algorithms on the nodes. In Fig. 3 (left), we depict the counting performance of Mark-Recapture (total-marks=20) compared with that of *PW-Avg* in the random scenario for  $N=50$ . Even though the estimates of both the algorithms in the given nodes can eventually reach the actual network size, the estimates of Mark-Recapture nodes converge at a shorter delay time. The use of sampling strategy in Mark-Recapture effectively produces a good estimate of network size within a relatively short duration of time, while *PW-Avg* requires more time to enable the

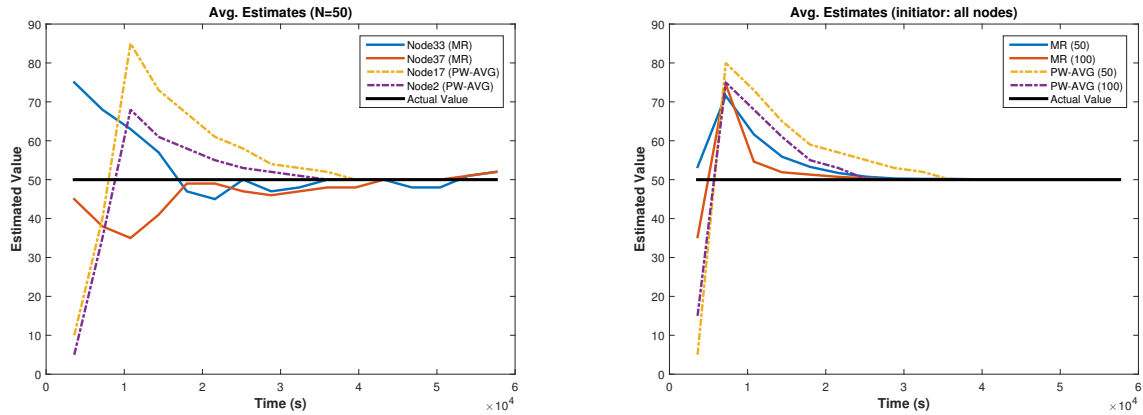


Figure 3: Mark-Recapture (MR) vs. Pair-Wise Average (PW-AVG) in the random mobility scenario

initiator to meet more nodes before having a proper result. Subsequently, in Fig. 3 (right) we show the effect of node density (e.g.,  $N=50$  and  $100$ ) on both the algorithms’ performances when all the network nodes perform node counting simultaneously. It is obvious that the increase of node density can reduce the convergence times of both the algorithms. Moreover, within the same node density, Mark-Recapture again outperforms PW-Avg in terms of convergence time.

We now discuss the performance evaluation of the Mark-Recapture algorithm in real human mobility scenarios. Typically, individuals move to places or meet other people to fulfill their social needs, and social (contact) graphs are commonly used to describe their social relationships. The authors of [29] investigated several real human contact datasets and confirmed that human mobility typically possesses a non-random contact pattern, where a few nodes (individuals) have contacts (or relations) with many others, but majority of nodes only have few ones. The nodes having a large number of contacts with others are therefore socially very popular in the networks (these popular nodes also called *hub nodes* in social network analysis, SNA). Firstly, we consider the short-term contact traces, the Huggle dataset [25]. In this scenario, we deliberately choose 3 nodes as counting initiators, namely node ids 21, 28, and 34, which represent the most-popular node, moderate-popular node, and the least-popular node, respectively, in Huggle. We then depict the Mark-Recapture performances on these nodes in Fig. 4 (left) for total-marks=10. We notice that node 21 (the most popular node) can accurately estimate the

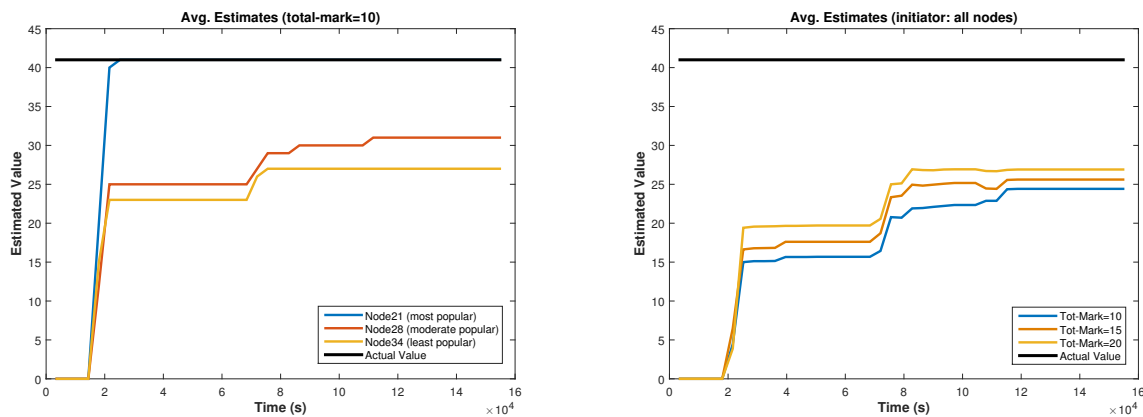


Figure 4: Node counting in the Huggle mobility scenario ( $N=41$ )

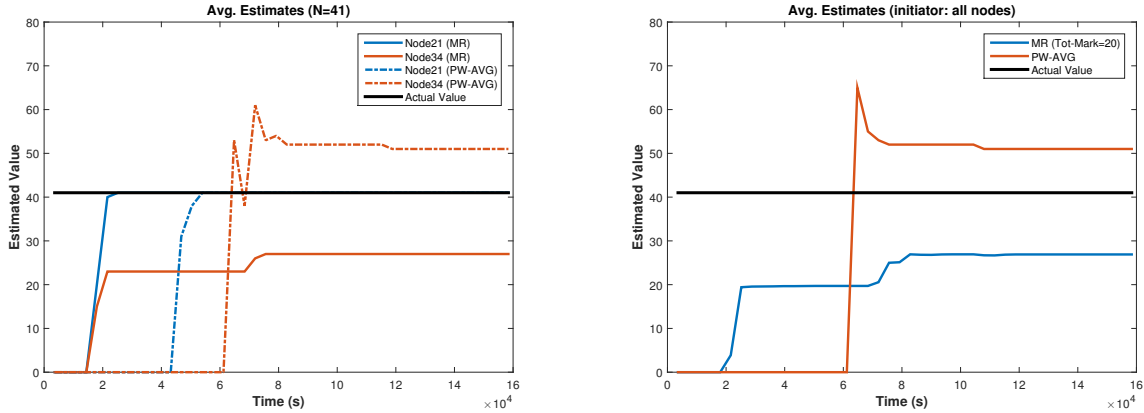


Figure 5: Mark-Recapture (MR) vs Pair-Wise Average (PW-AVG) in the Haggles mobility scenario ( $N=41$ )

network size in a relatively short time. In contrast, the less popular nodes (node 28 and 34) fail to predict the network size (i.e., the estimates of these nodes never converge to the real network size throughout the simulation time). As opposed to random movement, human mobility possesses a heterogeneous contact pattern. Consequently, the most popular node can perform the recapture process properly (i.e., it can meet the marked nodes with same probability), leading to accurately estimate the network size. In the less popular nodes, however, when the marking process can be assisted by other nodes (since our scheme uses the binary marking scheme), the nodes cannot (re)capture the marked nodes with the same probability; in turn, this results in inaccurate estimates of network size. Furthermore, in Fig. 4 (right) we depict the average estimates of network size for several total-marks when all the network nodes simultaneously initiate counting processes. We notice that the average estimates of all the network nodes are significantly below the actual network size. Due to the inherent characteristic of non-random contact, only a few popular nodes (as initiators) can produce a good estimate of network size, while most of the network nodes (i.e., less popular nodes) fail to do this. Consequently, as shown in Fig. 4 (right), the whole counting processes initiated by all the Haggles nodes result in the average estimates below the actual network size.

Lastly, we compare the performance of Mark-Recapture with that of PW-Avg in the Haggles scenario. We again choose two nodes in Haggles as before, namely node 21 and 34, that represent the most popular node and the least popular node, respectively, and then apply both the algorithms on the nodes successively. In Fig. 5 (left), we show the counting performance of Mark-Recapture (total-marks=20) compared with that of PW-Avg in the given nodes in Haggles ( $N=41$ ). It is obvious that both the algorithms in the most popular node (node 21) can work properly (i.e. the estimates of the node eventually converge to the actual network size). Moreover, Mark-Recapture can converge in a shorter time compared to PW-Avg in this popular node. In contrast, both the algorithms in the least popular node fail to produce a good estimate of network size. As described above, in Mark-Recapture the least popular node suffers from a recapture issue, as it cannot recapture the marked nodes with the same probability. Similarly, in PW-Avg the least popular node has only few contacts with others, therefore it cannot update the counting value properly, resulting in an incorrect estimate of network size. Furthermore, in Fig. 5 (right) we describe the counting performance of Mark-Recapture (total-marks=20) compared with that of PW-Avg when all the network nodes initiate counting processes simultaneously. Since majority of nodes in Haggles are less popular nodes (due to the inherent characteristic of heterogeneous contact in real human mobility), we then see in Fig. 5 (right) that both Mark-Recapture and PW-Avg are unsuccessful to predict the network

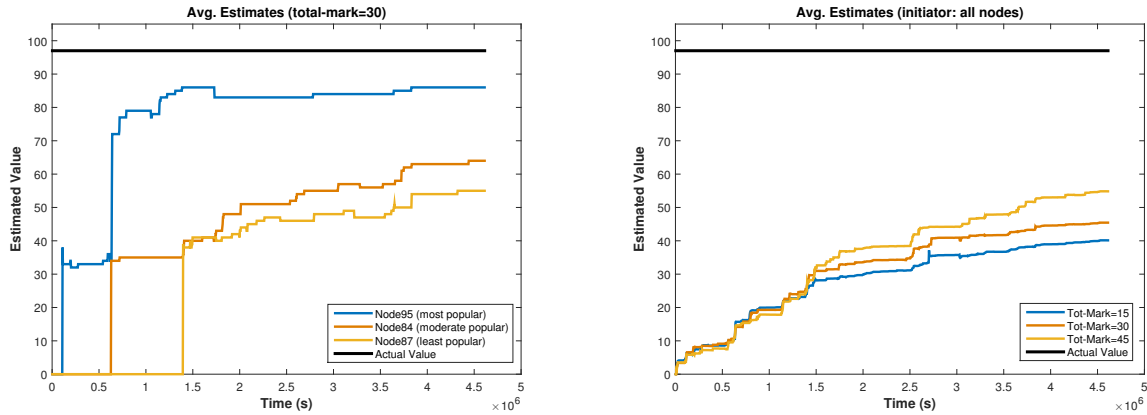


Figure 6: Node counting in the Reality mobility scenario ( $N=97$ )

size. Indeed, only a few nodes (i.e. popular nodes) in Huggle are able to effectively estimate the network size.

The final mobility scenario we consider is the Reality dataset [11], which is captured the long-term human mobility traces. We purposely select 3 nodes as initiators, namely node ids 95, 84, and 87, representing the most-popular node, moderate-popular node, and the least-popular node, respectively, and next apply Mark-Recapture on these nodes. In Fig. 6 (left), we illustrate the counting performances of Mark-Recapture in these nodes for total-marks=30. As similar to Huggle, we again see that the counting algorithm in the most popular node (node 95) outperforms those of the less popular nodes (node 84 and 87). The estimate of the most popular node is able to nearly approach the actual network size in a relatively short time. In contrast, the algorithms in the less popular nodes are ineffective to estimate the network size (i.e. the average estimates of these nodes never converge to the actual network size throughout the simulation time). As in Huggle, less popular nodes in Reality suffer from the recapture issue, since they cannot (re)capture the marked nodes with the same probability. Furthermore, in Fig. 6 (right) we depict the average estimates of network size when all nodes in Reality simultaneously initiate counting processes for several total-mark values. Since majority of the Reality nodes are less popular nodes (due to the heterogeneous contact pattern) and the less popular nodes are typically not able to perform node counting properly (due to the recapture issue), as shown in Fig. 6 (right) the average estimates produced by all the Reality nodes therefore are far below the actual network size for all the given total-mark values.

Finally, we compare the performance of Mark-Recapture with that of PW-Avg in Reality. We again choose node 95 and 87 representing the most popular and the least popular nodes, respectively, and then apply both the algorithms on these nodes consecutively. In Fig. 7 (left), we illustrate the counting performances of Mark-Recapture (total-marks=45) and PW-Avg on both the nodes. It is clear that both the algorithms can provide a good estimate when they are applied on the most popular node (node 95). In the least popular node, however, both the algorithms fail to accurately predict the network size and their estimates never converge to the real network size during all the simulation time. Furthermore, in Fig. 7 (right) we show the counting performance of Mark-Recapture (total-marks=45) compared with that of PW-Average when all the network nodes initiate counting processes concurrently. As similar to Huggle, we again see that majority of the Reality nodes (i.e., less popular nodes) fail to produce an accurate estimate of network size in both the algorithms, resulting in the average estimates of all the nodes are far from the actual network size. Actually, only a small number of (popular) nodes can contribute a correct result in the average estimates in Fig.7 (right).

To sum up, the counting performance of the Mark-Recapture distributed algorithm in OMNs is optimal, in terms of accuracy and convergence time, when all nodes move in a random manner in the area. In this case, all nodes are able to perform node counting properly, leading to accurately estimate the network size. Furthermore, Mark-Recapture can outperform an Aggregation scheme, namely PW-Avg, in terms of convergence time in the random scenario. However, both the algorithms suffer from a common problem in the real human mobility cases, where only a few (popular) nodes are able to carry out node counting appropriately, while majority of nodes work ineffectively. Meanwhile, in some specific cases of real human mobility, such as in disaster, military and crowd scenarios, human movement is typically modelled as a random process. For instance, the work of UrbanCount [7] used a City-Square model [9] to describe the movement of people in a city square. This movement model is actually an improvement of Random-Waypoint. Based on that study, we therefore believe that Mark-Recapture still can be utilized in these specific cases of human mobility. On the other hand, mobile crowd sensing applications [13] use a client-server paradigm, where mobile devices sensing and sending data to a server (in the cloud) that further processes the data and distributes the result to users who need the result. Using such client-server architecture, we suggest that Mark-Recapture can be utilized for counting the number of nodes in human-based opportunistic networks, where only a small number of (popular or most active) nodes sampling and reporting data to a central server, and the server eventually provide the final result to nodes requesting the information.

## 4.2 The Performance Improvement of Epidemic Routing with Local Estimate of Network Size

In this section, we discuss the application of node counting in data dissemination in OMNs. We exploit a local estimate of network size obtained from the Mark-Recapture algorithm to improve the delivery cost performance of Epidemic routing [26]. In Epidemic routing, a node forwards message copies to all the neighbours within the radio range so that the copies are quickly disseminated all over the network. This oblivious forwarding achieves near-optimal in terms of delivery latency when the node resources are assumed to be unlimited. In practice, however, Epidemic routing tends to quickly deplete the node resources, such as power and storage, and eventually greatly reduces the network delivery performance. We therefore improve Epidemic routing by tweaking it based on the observation in [30] as follows: *only the estimate of the number of nodes in the network ( $N$ ) is required to tune the number of copies ( $L$ ), and*

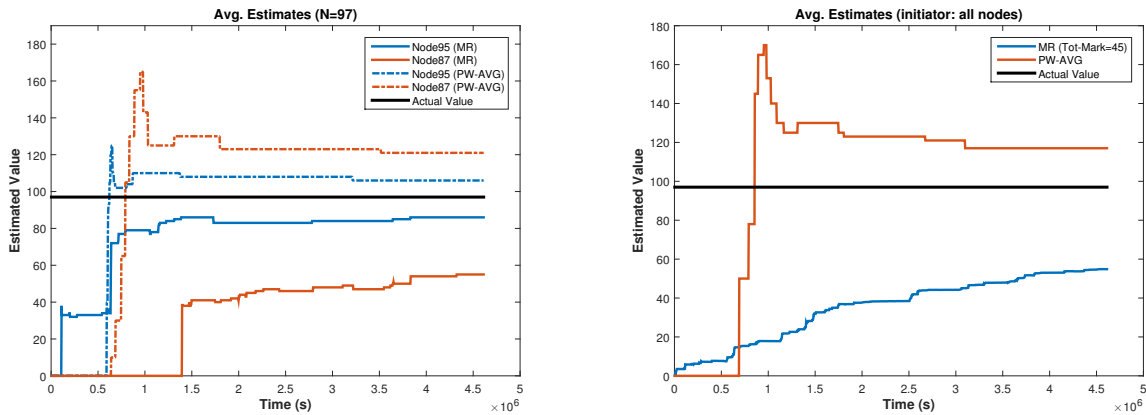


Figure 7: Mark-Recapture (MR) vs. Pair-Wise Average (PW-AVG) in the Reality mobility scenario ( $N=97$ )

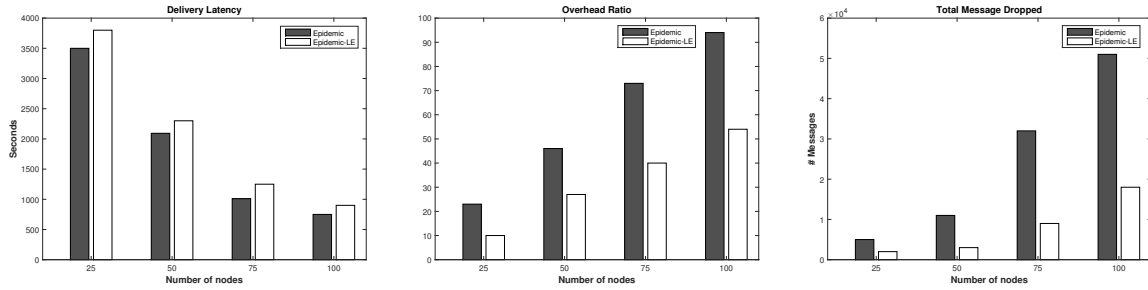


Figure 8: Performance comparison of Epidemic and Epidemic-LE in the random scenario

*Epidemic routing with  $L = N/2$  can achieve an optimal delivery delay with minimum resource overhead.* In order to incorporate the estimate of network size discussed so far in Epidemic routing, we associate a variable with each message, namely *total-copies* ( $L$ ) denoting the total number of message copies that can be forwarded by the source node and other nodes receiving a copy to  $L$  distinct relay nodes. When  $L$  copies have been spread, Epidemic routing stops to forward and lets each relay carrying a copy to perform direct transmission to the destination. Furthermore, in this experiment we set  $L$  to be a half of the local estimate of network size ( $\hat{N}$ ). We eventually compare Epidemic routing with a local estimate of network size (hereafter, we call it *Epidemic-LE*) to conventional Epidemic routing (hereafter, we just call it *Epidemic*) for three performance evaluations, namely delivery latency, overhead ratio, and total message dropped. We do not show the delivery ratio results since Epidemic-LE is able to achieve the delivery ratio as high as that of Epidemic in all scenarios.

We firstly discuss the performance of Epidemic-LE compared with that of Epidemic in the random scenario. In Fig. 8 we show the performance comparison of these routing schemes in terms of the given evaluation measures for different network sizes. In this random case, all nodes independently initiate a counting process to attain a local estimate of network size ( $\hat{N}$ ), and subsequently create a new message with total copies  $L = \hat{N}/2$  and it is sent to a randomly chosen destination. In the simulation, we set the message generation interval to be 5-10 minutes with the simulation time is 12 hours. For other simulation settings, we use the same settings used in the earlier experiment for the random case. As shown in Fig. 8, Epidemic-LE outperforms Epidemic in terms of overhead ratio and total message dropped for all the given network sizes. With the delivery ratio performance is almost the same between the two routing schemes, Epidemic-LE is able to reduce the copy redundancy in the network (indicated by the lower overhead ratio), leading to efficiently use the node resources, e.g. buffer or storage (showed by the significant decrease of total message dropped of Epidemic-LE compared to Epidemic's). Nevertheless, this reduced resource overhead of Epidemic-LE comes at a price, as the delivery latency slightly increases

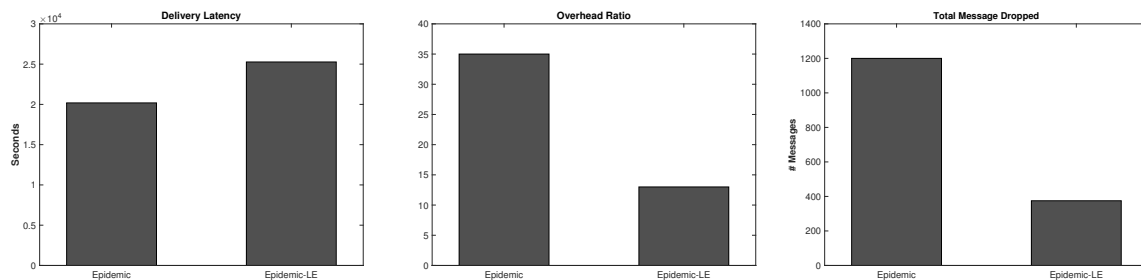


Figure 9: Performance comparison of Epidemic and Epidemic-LE in the Haggie scenario

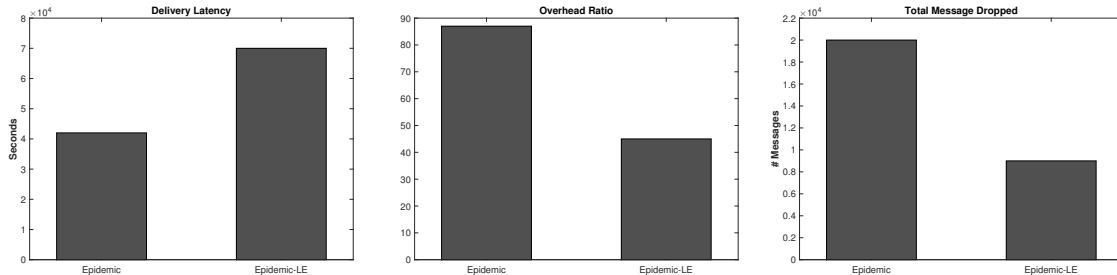


Figure 10: Performance comparison of Epidemic and Epidemic-LE in the Reality scenario

beyond that of Epidemic for all the given total number nodes in the network.

We next discuss the performance improvement of Epidemic-LE in the real human mobility scenarios, namely Huggle and Reality. In contrast to the random case, in these real-life cases only popular nodes can initiate a counting process (as we have described previously, the counting algorithms of less popular nodes fail to produce a good estimate of network size). In consequence, we assume that less popular nodes have ways to learn about the network size from the popular nodes (for example, using a simple flooding data dissemination algorithm or using a client-server architecture as in [13]). Subsequently, all the network nodes randomly create a new message with total copies  $L$  is set to a half of local estimate of network size, and the message is then sent to a randomly chosen destination. In Fig. 9 and 10, we show the simulation results of Epidemic and Epidemic-LE in Huggle and Reality, respectively, for the given performance metrics. In the simulations, we set the message generation interval to be 5-10 minutes with the simulation time is 3 days for Huggle, and the message generation interval to be 20-30 minutes with the simulation time is 3 months for Reality. From both the figures, we notice that by capping the total message copies distributed in the network at maximum  $L = \hat{N}/2$  replicas, Epidemic-LE can significantly reduce both the overhead ratio and total message dropped below those of Epidemic, while keeping the delivery ratio as high as Epidemic in both Huggle and Reality. However, as in the random case, we again see a trade-off between resource efficiency and delivery latency performance: the efficient use of node resources of Epidemic-LE increases the delivery delay beyond that of Epidemic in both the real human mobility scenarios. In addition, the increase of delivery delay is more obvious in Reality. Given that OMNs are a class of delay-tolerant networks (DTNs), this increase in delivery latency is not regarded substantial; instead, the reduction of node resource consumption, reflected in the improved overhead ratio and total message dropped, represents a significant improvement in the network's performance.

## 5 Conclusion

We have presented the Mark-Recapture distributed algorithm, a novel node counting technique targeted at accurately estimate the total number of active nodes in an OMN with a low delay. We have demonstrated that the algorithm achieves high accuracy and a low convergence time in estimating network size in the random *i.i.d.* movement case. In addition, Mark-Recapture can outperform a gossip-based Pair-Wise Average scheme in terms of convergence time in this random scenario. However, in the real human mobility scenarios, only the algorithm in popular nodes (both Mark-Recapture and PW-Avg) can produce an accurate estimate of network size in a relatively short delay time, while majority of nodes (i.e., less popular nodes) are ineffective to perform node counting.

After this, we improved Epidemic routing by incorporating a local estimate of network size to the routing scheme to reduce message redundancy in the network. We showed that Epidemic with LE (local estimates) can achieve delivery ratio as high as conventional Epidemic routing, but at a lower overhead

ratio and total message dropped in both the random and real-life scenarios. Nevertheless, this efficient delivery of Epidemic-LE slightly increases the delivery latency beyond that of (conventional) Epidemic routing.

Finally, for future works we can identify two points. First, we have shown that Mark-Recapture cannot work appropriately in less popular nodes in the real human mobility scenarios. Consequently, these nodes should rely on popular nodes to learn an accurate information of network size. In future, we therefore need to study a method to efficiently distribute the counting results of the popular nodes to all nodes in the network, such as a publish-subscribe scheme [12] or a client-server model [10]. Second, even though this paper only considers a closed system, we also need to take into account a more realistic scenario, i.e. an open system with churn, where nodes are allowed to enter and leave the area during the experiment. We believe that our proposed algorithm should be improved to accommodate this complex system.

## Acknowledgments

This work has been partially supported by the Sanata Dharma University's research grant 2019.

## References

- [1] N. Accettura, G. Neglia, and L. A. Grieco. The capture-recapture approach for population estimation in computer networks. *Computer Networks*, 89:107–122, October 2015.
- [2] C. Y. Aung, I. W.-H. Ho, and P. H. J. Chong. Store-carry-cooperative forward routing with information epidemics control for data delivery in opportunistic networks. *IEEE Access*, 5:6608–6625, April 2017.
- [3] H. Barbosa, M. Barthelemy, G. Ghoshal, C. R. James, M. Lenormand, T. Louail, R. Menezes, J. J. Ramasco, F. Simini, and M. Tomasini. Human mobility: Models and applications. *Physics Reports*, 734(6):1–74, March 2018.
- [4] S. Batabyal and P. Bhaumik. Estimators for global information in mobile opportunistic network. In *Proc. of the 7th IEEE International Conference on Advanced Networks and Telecommunications Systems (ANTS'13), Kattankulathur, India*, pages 1–6. IEEE, December 2013.
- [5] L. Chen, A. Karbasi, and F. W. Crawford. Estimating the size of a large network and its communities from a random sample. In *Proc. of the 30th International Conference on Neural Information Processing Systems (NIPS'16), Barcelona, Spain*, pages 3080—3088. Curran Associates Inc., December 2016.
- [6] J. Cui, S. Cao, Y. Chang, L. Wu, D. Liu, and Y. Yang. An adaptive spray and wait routing algorithm based on quality of node in delay tolerant network. *IEEE Access*, 7:35274–35286, March 2019.
- [7] P. Danielis, S. T. Kouyoumdjieva, and G. Karlsson. Urbancount: Mobile crowd counting in urban environments. In *Proc. of the 8th IEEE Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON'17), Vancouver, British Columbia, Canada*, pages 640–648. IEEE, October 2017.
- [8] J. Darroch. The multiple-recapture census i. estimation of a closed population. *Biometrika*, 45(3/4):343–359, December 1958.
- [9] M. S. Desta, E. Hyttiä, J. Ott, and J. Kangasharju. Characterizing content sharing properties for mobile users in open city squares. In *Proc. of the 10th Annual Conference on Wireless On-demand Network Systems and Services (WONS'13), Banff, Alberta, Canada*, pages 147–154. IEEE, March 2013.
- [10] L. Duan, T. Kubo, K. Sugiyama, J. Huang, T. Hasegawa, and J. Walrand. Motivating smartphone collaboration in data acquisition and distributed computing. *IEEE Transactions on Mobile Computing*, 13(10):2320–2333, October 2014.
- [11] N. Eagle and A. S. Pentland. Reality mining: Sensing complex social systems. *Personal and Ubiquitous Computing*, 10(4):255—268, May 2006.



- [12] C. Fang, H. Yao, Z. Wang, W. Wu, X. Jin, and F. R. Yu. A survey of mobile information-centric networking: Research issues and challenges. *IEEE Communications Surveys & Tutorials*, 20(3):2353–2371, February 2018.
- [13] R. K. Ganti, F. Ye, and H. Lei. Mobile crowdsensing: current state and future challenges. *IEEE Communications Magazine*, 49(11):32–39, November 2011.
- [14] A. Guerrieri, I. Carreras, F. De Pellegrini, A. Montresor, and D. Miorandi. Distributed estimation of global parameters in delay-tolerant networks. In *Proc of the 10th IEEE International Symposium on a World of Wireless, Mobile and Multimedia Networks Workshops (WoWMoM'09), Kos, Greece*, pages 1–7. IEEE, June 2009.
- [15] Y. He, X. Tang, R. Zhang, X. Du, D. Zhou, and M. Guizani. A course-aware opportunistic routing protocol for fanets. *IEEE Access*, 7:144303–144312, October 2019.
- [16] J. Hu, L.-L. Yang, H. V. Poor, and L. Hanzo. Bridging the social and wireless networking divide: Information dissemination in integrated cellular and opportunistic networks. *IEEE Access*, 3:1809–1848, September 2015.
- [17] M. Jelasity, A. Montresor, and O. Babaoglu. Gossip-based aggregation in large dynamic networks. *ACM Transactions on Computer Systems*, 23(3):219–252, August 2005.
- [18] P. Jesus, C. Baquero, and P. S. Almeida. A survey of distributed data aggregation algorithms. *IEEE Communications Surveys & Tutorials*, 17(1):381–404, September 2014.
- [19] A. Keränen, J. Ott, and T. Kärkkäinen. The one simulator for dtn protocol evaluation. In *Proc. of the 2nd International Conference on Simulation Tools and Techniques (Simutools'09), Rome, Italy*, pages 1–10. ICST, March 2009.
- [20] T. Li, S. T. Kouyoumdjieva, G. Karlsson, and P. Hui. Data collection and node counting by opportunistic communication. In *Proc. of the 18th IFIP Networking Conference (IFIP'19), Warsaw, Poland*, pages 1–9. IEEE, May 2019.
- [21] Z. Lu, G. Cao, and T. La Porta. Teamphone: Networking smartphones for disaster recovery. *IEEE Transactions on Mobile Computing*, 16(12):3554–3567, April 2017.
- [22] S. Manaseer and I. Alhabash. Number of node estimation in mobile ad hoc networks. *International Journal of Interactive Mobile Technologies (IJIM)*, 11(6):65–72, November 2017.
- [23] L. Massoulié, E. Le Merrer, A.-M. Kermarrec, and A. Ganesh. Peer counting and sampling in overlay networks: Random walk methods. In *Proc. of the 25th Annual ACM Symposium on Principles of Distributed Computing (PODC'06), Denver, Colorado, USA*, pages 123–132. ACM, July 2006.
- [24] T. Ning, Z. Yang, and H. Wu. Counting in delay-tolerant mobile networks. In *Proc. of the 2010 IEEE International Conference on Communications (ICC'10), Cape Town, South Africa*, pages 1–5. IEEE, May 2010.
- [25] J. Scott, R. Gass, J. Crowcroft, P. Hui, C. Diot, and A. Chaintreau. CRAWDAD dataset cambridge/haggle (v. 2009-05-29). <https://crawdad.org/cambridge/haggle/20090529/imote> [Online; accessed on September 15, 2021], May 2009.
- [26] A. Vahdat and D. Becker. Epidemic routing for partially-connected ad hoc networks. Technical Report CS-2000-06, Duke University, 2000.
- [27] E. M. Volz, J. C. Miller, A. Galvani, and L. Ancel Meyers. Effects of heterogeneous and clustered contact patterns on infectious disease dynamics. *PLOS Computational Biology*, 7(7):1–13, June 2011.
- [28] E. Wang, Y. Yang, J. Wu, and W. Liu. Phone-to-phone communication utilizing wifi hotspot in energy-constrained pocket switched networks. *IEEE Transactions on Vehicular Technology*, 65(10):8578–8590, January 2016.
- [29] E. Yoneki, P. Hui, and J. Crowcroft. Distinct types of hubs in human dynamic networks. In *Proc. of the 1st Workshop on Social Network Systems (SocialNets'08), Glasgow, Scotland*, pages 7–12. ACM, April 2008.
- [30] X. Zhang, G. Neglia, J. Kurose, and D. Towsley. Performance modeling of epidemic routing. *Computer Networks: The International Journal of Computer and Telecommunications Networking*, 51(10):2867–2891, July 2007.

## Author Biography



**Bambang Soelistijanto** is currently a senior lecturer with the Department of Informatics, Sanata Dharma University, Indonesia. He received his bachelor's degree and a M.Sc. degree in Electrical Engineering from Gadjah Mada University, Indonesia and Delft University of Technology, the Netherlands, respectively. In 2014, he obtained the PhD degree in Electrical Engineering from the University of Surrey, England. His expertise is on mobile communication, and his main research interests include opportunistic networks, mobile social networks, and distributed computing.



**Geraldev Manoah** was a bachelor student at the Department of Informatics, Sanata Dharma University, Indonesia. Currently, he is an IT staff in a senior high school, Palangka Raya, Indonesia. His research interests include opportunistic networks and Internet technologies.