

# Multicasting in Delay Tolerant Networks: A Social Network Perspective

Wei Gao, Qinghua Li, Bo Zhao and Guohong Cao  
Department of Computer Science and Engineering  
The Pennsylvania State University  
University Park, PA, 16802, USA  
{wxg139, qxl118, bzhao, gcao}@cse.psu.edu

## ABSTRACT

Node mobility and end-to-end disconnections in Delay Tolerant Networks (DTNs) greatly impair the effectiveness of data dissemination. Although social-based approaches can be used to address the problem, most existing solutions only focus on forwarding data to a single destination. In this paper, we are the first to study multicast in DTNs from the social network perspective. We study multicast in DTNs with single and multiple data items, investigate the essential difference between multicast and unicast in DTNs, and formulate relay selections for multicast as a unified knapsack problem by exploiting node centrality and social community structures. Extensive trace-driven simulations show that our approach has similar delivery ratio and delay to the Epidemic routing, but can significantly reduce the data forwarding cost measured by the number of relays used.

## Categories and Subject Descriptors

C.2.1 [Network and Architecture Design]: Wireless communication, Store and forward networks; C.4 [Performance of Systems]: Modeling techniques

## General Terms

Design, Algorithms, Performance

## Keywords

Multicast, Delay Tolerant Network, Social Network, Centrality, Community

## 1. INTRODUCTION

In Delay Tolerant Networks (DTNs) [6], mobile users contact each other opportunistically in corporate environments, such as conference sites and university campuses. Due to low node density and unpredictable node mobility, end-to-end connections are hard to maintain. Alternatively, node mobility is exploited to let mobile nodes physically carry data as relays, and forward data opportunistically upon contacts. The key problem is how to determine appropriate relay selection strategy and data forwarding criteria.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

MobiHoc '09, May 18-21, 2009, New Orleans, USA.

Copyright 2009 ACM 978-1-60558-083-9/08/05 ...\$5.00.

Recent trace-based study on campus wireless networks [8] shows that different nodes have heterogeneity in their contact patterns, and such heterogeneity validates the use of Social Network Analysis (SNA) for data forwarding in DTNs [17, 21]. There are two key concepts in SNA methods: (i) *Communities*, which are naturally formed according to people's social relations. Social communities are derived from the "small-world" phenomenon, which is first investigated by Milgram's experiment [16] in 1967, and is later formalized as a random graph problem in [23]. (ii) *Centrality*, which shows that some nodes in a community are the common acquaintances of other nodes and act as communication hubs. Since social relations among mobile users are more likely to be long-term characteristics and less volatile than node mobility, social-based forwarding schemes outperform traditional approaches based on oblivious heuristics [22] or mobility-based predictions [13, 25, 20].

The aforementioned work focuses on forwarding data to a single destination. Multicast, on the other hand, is more effective for data dissemination and multi-party communication, but is also more difficult to model and implement in opportunistic DTNs. Although there are some initial efforts on studying multicast in DTNs, they are limited to semantic multicast models [26] and multicast capacity analysis [12], and none of them considers multicast in DTNs from the social network perspective.

Due to DTN dynamics, deterministic data forwarding, either unicast and multicast, is only guaranteed in two cases: (i) the network is flooded, and (ii) the data forwarding process does not have time constraint. Neither of the two cases are practical in DTNs due to the inevitably high forwarding cost. Thus, a more practical solution is to maximize the data forwarding probability with a given time constraint. From such a probabilistic perspective, the *essential difference* between multicast and unicast in DTNs is that, a selected relay for multicast is expected to forward data to as many destinations as possible. The cumulative probability for a relay to forward data to multiple destinations therefore needs to be calculated, and such calculation may require global knowledge of social relations among nodes.

In this paper, we focus on improving the cost-effectiveness of multicast in DTNs by exploiting the two key concepts in Social Network Analysis, i.e., *centrality* and *communities*. We aim at minimizing the multicast cost, in terms of the number of relays used, given the required delivery ratio and time constraint. We first consider multicasting a single data item to the network, and then generalize the problem to multiple data items with node buffer constraints. Our detailed contributions are as follows:

- We develop analytical models for multicast relay selection using social network concepts.
- We formulate the relay selections for single-data and multiple-data multicast in DTNs as a unified knapsack problem.

- We provide deep insights into the problem of multicasting multiple data items, by investigating individual nodes’ awareness of their data forwarding probabilities to destinations.

Our approach is based on a weighted social network model for DTNs, such that the edges in the network contact graph are modeled as Poisson processes with pairwise node contact rates as the parameters. Based on this model, node centrality and social community structures are exploited for relay selections under the unified knapsack formulation.

The rest of this paper is organized as follows. Section 2 provides an overview about problem definitions, basic ideas and the weighted social network modeling. Based on such modeling, single-data and multiple-data multicast problems are studied in Sections 3 and 4, respectively. Section 5 evaluates the performance of our approach, Section 6 reviews existing work, and Section 7 concludes the paper.

## 2. OVERVIEW

### 2.1 Problem Definitions and Assumptions

We first focus on multicasting a single data item:

**PROBLEM 1. *Single-Data Multicast (SDM)***

$\{p, \mathbb{D}, T\}$ : To deliver a data item to a set  $\mathbb{D}$  of destinations, how to choose the minimum number of relays to achieve the delivery ratio  $p$  within the time constraint  $T$ ?

The SDM Problem is then generalized as follows:

**PROBLEM 2. *Multiple-Data Multicast (MDM)***

$\{p, \mathbb{D}_1, \dots, \mathbb{D}_n, s_1, \dots, s_n, T\}$ : To deliver a set of data items  $d_1, d_2, \dots, d_n$  with sizes  $s_1, \dots, s_n$ , from a data source to destination sets  $\mathbb{D}_1, \mathbb{D}_2, \dots, \mathbb{D}_n$ , how to choose the minimum number of relays to achieve the delivery ratio  $p$  within the time constraint  $T$ ?

In these problems, we assume that the selected relays can forward data when contacting other nodes. Suppose each node  $N_k$  has buffer  $B_k$ . For SDM, such buffer constraint is trivial because a node can simply refuse to receive the data if it does not have enough buffer. For MDM, since a node most likely cannot carry all the data items simultaneously due to the buffer constraint, and we should consider which data items to be carried by a selected relay.

We define the required delivery ratio  $p$  as the average ratio of data items being delivered to destinations. Such ratio is defined from a statistic perspective based on a large number of data items generated at the data source. For an arbitrary destination node, its delivery ratio is therefore equivalent to the average probability that it can receive a data item within time constraint  $T$ .

### 2.2 The Basic Approach

The basic idea of our approach is to develop *social-based metrics* based on the probabilities of nodes forwarding data to their destinations. Such metrics are developed based on social network concepts including centrality and social communities. Based on the social-based metrics, we formulate the relay selections in SDM and MDM uniformly as a knapsack problem:

$$\begin{aligned} \min \quad & \sum_{k=1}^n x_k \\ \text{s.t.} \quad & \sum_{k=1}^n w_k x_k \geq W \end{aligned} \tag{1}$$

where  $x_k \in \{0, 1\}$  indicates whether node  $N_k$  is selected as the relay, and the constraint indicates that the selected relays should satisfy the performance requirements in delivery ratio and delay.

The solution to this knapsack problem itself is trivial, as we can select the node with the maximum weight  $w_k$  round by round until the constraint is satisfied. Note that only best-effort solution is available if  $\sum_{k=1}^n w_k < W$ . Social-based metrics will be developed to calculate the weight  $w_k$  associated with each node  $N_k$  in the network, and the total required weight  $W$  is determined by the performance requirements. The rest of this paper therefore focuses on answering the following questions:

1. What are the appropriate social-based metrics for SDM and MDM, respectively?
2. How to calculate the weights  $w_k$  of individual nodes?
3. How can the source calculate the total required weight  $W$ ?

Generally speaking, the essential difference between SDM and MDM is on the required knowledge about node social relations for relay selections. For SDM, relay selection can be done based on the local knowledge of the data source about its contacted neighbors<sup>1</sup>, because the data source only multicasts a single data item, and does not need to distinguish the data forwarding probabilities to different destinations for relay selection. In Section 3, we use cumulative contact probability as the centrality metric to develop a centrality-based heuristic for SDM, and show that such heuristic is able to satisfy the given performance requirements.

On the other hand, for MDM the relays should be aware of their probabilities for forwarding each data item to the destinations. Such capability is called “*destination-awareness*” throughout this paper, and is required mainly due to the node buffer constraints. For example, suppose the data source  $S$  multicasts two data items  $d_1, d_2$  to the destination sets  $\mathbb{D}_1, \mathbb{D}_2$  respectively, and a selected relay  $R$  can only carry one data item. To maximize the delivery ratio,  $R$  should carry  $d_1$  if its probability of forwarding  $d_1$  to destinations in  $\mathbb{D}_1$  is higher. Otherwise,  $d_2$  should be carried by  $R$ . Thus, the source has to be destination-aware, which may require global knowledge about the social relations between the destinations and other nodes in the network. Since such global knowledge is hard to maintain in DTNs, in Section 4 we propose a community-based approach which only requires nodes to maintain destination-awareness about other nodes in the same social community.

**Table 1: Trace summary**

Trace	<i>Infocom</i>	<i>MIT Reality</i>
Network type	Bluetooth	Bluetooth
No. of devices	41	97
No. of internal contacts	22,459	54,667
Duration (days)	3	246
Granularity (secs)	120	300
Pairwise contact frequency (per day)	4.6	0.024

### 2.3 Experimental Traces

We use two experimental traces collected from realistic DTNs to validate our social network modeling, and to evaluate the performance of our multicast scheme. These traces record contacts among users in corporate environments carrying Bluetooth devices. The Bluetooth devices periodically discover their peers in the neighborhood and record contacts. We believe that the chosen traces cover a large diversity of environments, from university campuses (*MIT Reality*) to conference sites (*Infocom*), with experimental periods from a few days (*Infocom*) to several months (*MIT Reality*). The two traces are summarized in Table. 1.

<sup>1</sup>The nodes that have been directly contacted by the data source.

## 2.4 Social Network Modeling

In this paper, we model a weighted social network which differentiates the contact frequencies of different node pairs. In this model, the contact process of each node pair is formulated as a Poisson process, with the corresponding pairwise node contact rate as its parameter. Similar assumptions have been used in other existing works to analyze the multicast capacity [12] and content dissemination [11] in DTNs,

**Table 2: Acceptance ratio of  $\chi^2$  tests for *Infocom* Trace**

No. of test intervals	5	10	15	20	25
$\alpha=0.95$	74.71	85.02	87.14	91.03	94.71
$\alpha=0.75$	78.54	86.46	88.29	91.41	94.79
$\alpha=0.50$	82.83	87.44	89.59	91.88	94.87

**Table 3: Acceptance ratio of  $\chi^2$  tests for *MIT Reality* Trace**

No. of test intervals	5	10	15	20	25
$\alpha=0.95$	54.45	87.53	89.20	93.67	90.67
$\alpha=0.75$	60.74	88.18	90.50	93.78	91.21
$\alpha=0.50$	66.93	88.65	91.20	93.96	91.94

Although [12] and [11] have this assumption, they did not validate it experimentally. Next, we validate this social network model based on realistic trace analysis. As a prerequisite, letting the random variable  $X_{AB}(t)$  be the cumulative number of contacts of a node pair  $A$  and  $B$  at time  $t$ , we assume that any two contacts between  $A$  and  $B$  are independent from each other. Hence,  $X_{AB}(t)$  is a stochastic process with independent increments, i.e., for any  $0 \leq t_1 < t_2 < \dots < t_n$ ,  $X_{AB}(t_2) - X_{AB}(t_1)$ ,  $X_{AB}(t_3) - X_{AB}(t_2)$ , ...,  $X_{AB}(t_n) - X_{AB}(t_{n-1})$  are all independent random variables.

We observe that for most of the contacted node pairs in the *Infocom* and *MIT Reality* traces, the pairwise inter-contact time is exponentially distributed. To validate this result, we conduct  $\chi^2$  hypothesis test [7] on each contacted node pair, to test whether the hypothesis “the pairwise inter-contact time is exponentially distributed with parameter  $\bar{\lambda} = n / \sum_{i=1}^n T_i$ ”, where  $T_1, T_2, \dots, T_n$  are the inter-contact time samples, can be accepted. Since exponential distribution is continuous, in the  $\chi^2$  tests we divide the range of the sample values into several test intervals, and compare the sample frequencies with theoretical probabilities on each interval. The results of acceptance ratio on the *Infocom* and *MIT Reality* traces under different significance levels  $\alpha$  are listed in Tables 2 and 3. The results show that, when a enough number of test intervals ( $\geq 10$ ) is used, over 85% of the contacted node pairs pass the test.

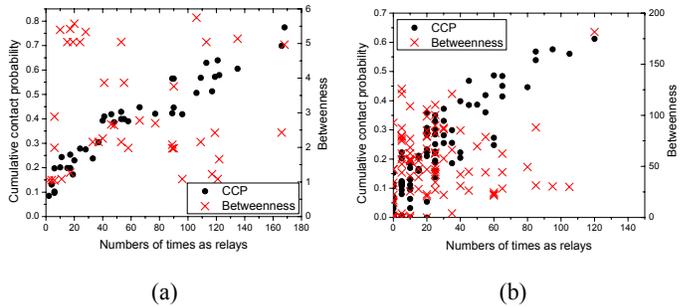
$X_{AB}(t)$  is therefore modeled as a homogeneous Poisson process. For any  $t > 0$ , the number of contacts  $X_{AB}(t + \Delta t) - X_{AB}(t)$  between nodes  $A$  and  $B$  during time  $\Delta t$  follows Poisson distribution  $P(\lambda_{AB}\Delta t)$ , i.e.,

$$P(X_{AB}(t + \Delta t) - X_{AB}(t) = k) = \frac{(\lambda_{AB}\Delta t)^k e^{-\lambda_{AB}\Delta t}}{k!}. \quad (2)$$

Then, the contact frequency between node pair  $\{A, B\}$  is indicated by the contact rate  $\lambda_{AB}$ , and its stochastic properties are represented by the Poisson process.

## 3. SINGLE-DATA MULTICAST

In this section, we develop a centrality-based heuristic for the SDM problem based on the local knowledge of the data source. The relays are selected among the contacted neighbors of the data source based on their centrality, to ensure that the required delivery ratio can be achieved within the time constraint.



**Figure 1: Node centrality values and the number of times for them to be counted as relays: (a) *Infocom*, (b) *MIT Reality***

## 3.1 Centrality Metric

Currently, the “betweenness” centrality metric is widely used in social-based data forwarding [4, 9]. Betweenness measures the extent to which a node lies on the shortest paths linking other nodes, and a node with higher betweenness has better capability of facilitating communication between other nodes. For utilization in distributed environments, the localized version of betweenness is proposed by Marsden [14] in the “ego-centric” network for each node, which only includes the contacted neighbors of that node.

Unfortunately, betweenness is defined and calculated based on the topology of network contact graph, and is not sufficient to analytically represent the probabilities for a node to contact others. To analytically model the relay selection process in SDM, we propose a new centrality metric based on the Poisson modeling of social networks, in the ego-centric network of a mobile node:

**DEFINITION 1.** Suppose there are totally  $N$  nodes in the network, and the contact rate between a node  $N_i$  to node  $N_j$  is  $\lambda_{ij}$ . The **cumulative contact probability (CCP)** of  $N_i$  is defined as

$$C_i = 1 - \frac{1}{N-1} \sum_{j=1, j \neq i}^N e^{-\lambda_{ij}T}. \quad (3)$$

$C_i$  indicates the average probability that a randomly chosen node in the network is contacted by  $N_i$  within time  $T$ . Since all the nodes in DTNs can exchange their centrality values upon contacts with each other, the data source knows the centrality values of all of its contacted neighbors when it selects relays.

To show the effectiveness of our centrality metric in characterizing the capability of a node to contact other nodes and to deliver data to destinations, we run 500 SDM scenarios with random data sources and destinations in both *Infocom* and *MIT Reality* traces, using Epidemic routing [22] to forward data. If a node delivers the data to a destination, the node is counted as a relay. The statistical results on the number of times for nodes to be counted as relays are shown in Figure 1. The CCP spots in the figures show the trend to form straight lines, which lead to the result that nodes with higher CCP values are more effective in delivering data to destinations. Comparatively, betweenness as the centrality metric is not effective enough, as the corresponding spots are scattered in wide ranges.

## 3.2 Relay Selection

The minimum number of relays is selected to satisfy the required delivery ratio  $p$  within time constraint  $T$ . We assume that the destinations are uniformly distributed in the network, so that each node other than the data source has equal chance to be a destination. To forward data to all the destinations, we ensure that all the nodes are contacted by the data source or the selected relays within  $T$ .

### 3.2.1 Relays In Contact

Suppose the data source  $S$  is in contact with a set of nodes  $\mathcal{R} = \{R_1, R_2, \dots, R_k\}$ , which can be selected as relays. The selected relays still need to contact the remaining  $N - k - 1$  nodes other than  $S$  and nodes in  $\mathcal{R}$ , to make sure that all the destinations are contacted. We assume that none of the destinations is in the set  $\mathcal{R}$ , otherwise  $S$  can trivially forward data to destinations immediately.

We define the random variable  $X_{ij}$  as

$$X_{ij} = \begin{cases} 1 & \text{If a node } N_j \text{ is contacted by } R_i \text{ within } T \\ 0 & \text{Otherwise} \end{cases}$$

then  $X_{ij}$  follows Bernoulli distribution with  $\bar{X}_{ij} = 1 - e^{-\lambda_{ij}T}$  where  $\lambda_{ij}$  is the contact rate between  $R_i$  and  $N_j$ .

When we randomly choose a node other than  $S$  and  $\mathcal{R}$  in the network, which needs to be contacted by the relays, the probability for each node to be chosen is  $\frac{1}{N-k-1}$ . Therefore, for a randomly chosen node  $N_j \notin \tilde{\mathcal{R}}$ , where  $\tilde{\mathcal{R}} = \{S\} \cup \mathcal{R}$ , the probability that  $N_j$  is not contacted by  $R_i$  within  $T$  is

$$\begin{aligned} p_i &= \frac{1}{N-k-1} \cdot \sum_{j=1, j \notin \tilde{\mathcal{R}}}^N (1 - X_{ij}) = \frac{1}{N-k-1} \sum_{j=1, j \notin \tilde{\mathcal{R}}}^N e^{-\lambda_{ij}T} \\ &= \frac{N-1}{N-k-1} \cdot (1 - C_i) - \frac{1}{N-k-1} \cdot \sum_{j \in \tilde{\mathcal{R}}, j \neq i} e^{-\lambda_{ij}T}. \end{aligned} \quad (4)$$

Since  $S$  is in contact with every  $R_i \in \mathcal{R}$ ,  $S$  can request each  $R_i$  to calculate  $p_i$  locally based on its centrality value  $C_i$ .

To ensure that the average delivery ratio is higher than  $p$ , the probability that a randomly chosen node  $N_j \notin \tilde{\mathcal{R}}$  is not contacted by the selected relays should be lower than  $1 - p$ , i.e.,

$$\prod_{i=1}^r p_i^{x_i} \leq 1 - p \quad (5)$$

where  $x_i \in \{0, 1\}$  indicates whether  $R_i$  is selected as relay. Such problem can be transformed to the unified knapsack formulation in Eq. (1) by taking logarithms on both sides of the inequality, where  $w_i = \log \frac{1}{p_i}$  and  $W = \log \frac{1}{1-p}$ .

### 3.2.2 Relays Not In Contact

When the data source selects relays, some contacted neighbors of the data source with higher centrality may not be in contact with the data source. To select relays among these nodes, we should also consider the time needed for them to contact the data source.

For a contacted neighbor  $R_i$  of the data source  $S$ , let the random variable  $T_1$  be the time for  $S$  to contact  $R_i$ , and  $T_2$  be the time for  $R_i$  to contact another node  $N_j$ , then the probability that  $S$  forwards data to  $N_j$  via  $R_i$  is  $P(T_1 + T_2 \leq T)$ . Assuming that the probability density functions (PDF) of  $T_1$  and  $T_2$  are  $f_1(t)$  and  $f_2(t)$  ( $t \geq 0$ ), respectively,  $P(T_1 + T_2 \leq T)$  is calculated through the convolution  $f_1(t) \otimes f_2(t)$  as

$$\begin{aligned} P(T_1 + T_2 \leq T) &= \int_0^T f_1(t) \otimes f_2(t) dt \\ &= \int_0^T dt \left( \int_0^t f_1(\tau) f_2(t - \tau) d\tau \right). \end{aligned} \quad (6)$$

In order to select relays using the node centrality values as weights, in our approach we exploit the following lower bound as an approximation to such probability.

**THEOREM 1.** For any fixed  $\tilde{T} \in (0, T)$ , we have

$$P(T_1 + T_2 \leq T) \geq (1 - e^{-\lambda_{Si}\tilde{T}})(1 - e^{-\lambda_{ij}(T-\tilde{T})}). \quad (7)$$

**PROOF.** The r.h.s. of Eq. (7) indicates the probability

$$P(T_1 \leq \tilde{T}) \cdot P(T_2 \leq T - \tilde{T}).$$

We define two sets  $\mathbb{T}_1$  and  $\mathbb{T}_2$  as

$$\begin{cases} \mathbb{T}_1 &= \{(T_1, T_2) | T_1 + T_2 \leq T\} \\ \mathbb{T}_2 &= \{(T_1, T_2) | T_1 \leq \tilde{T}, T_2 \leq T - \tilde{T}\} \end{cases}$$

It is obviously that for any element  $\{t_1, t_2\} \in \mathbb{T}_2$ , we also have  $\{t_1, t_2\} \in \mathbb{T}_1$ , which means  $\mathbb{T}_2 \subseteq \mathbb{T}_1$ . Therefore, we have

$$P(T_1 + T_2 \leq T) \geq P(T_1 \leq \tilde{T}) \cdot P(T_2 \leq T - \tilde{T})$$

which proves the theorem.  $\square$

From Theorem 1, for a randomly chosen node other than  $S$  and  $R_i$  in the network, the average probability  $\bar{P}(T_1 + T_2 \leq T)$  for the relay choice  $R_i$  has the similar lower bound:

$$\bar{P}(T_1 + T_2 \leq T) \geq (1 - e^{-\lambda_{Si}\tilde{T}}) \left(1 - \frac{1}{N-2} \sum_{j=1, j \notin \{S, i\}}^N e^{-\lambda_{ij}(T-\tilde{T})}\right).$$

Note that, since  $R_i$  is not in contact with  $S$ , in the above equation we also consider the contacts between  $R_i$  to another contacted neighbor, say  $R_j$ , of  $S$ . If  $R_j$  is in contact with  $S$ ,  $\bar{P}(T_1 + T_2 \leq T)$  is reduced and therefore provides a lower bound which also guarantees the required delivery ratio in later relay selection. If  $R_j$  is not in contact with  $S$ , we should consider the contacts between  $R_i$  and  $R_j$  because  $S$  contacts  $R_j$  opportunistically, too.

The tightest bound in Eq. (7) is achieved by  $\tilde{T}_0$  when  $\frac{\partial \bar{P}}{\partial \tilde{T}} \Big|_{\tilde{T}=\tilde{T}_0} = 0$ , and such equation can be proved to have only one solution in  $(0, T)$ . Such differential equation can be written as

$$\frac{\lambda_{Si} e^{-\lambda_{Si}\tilde{T}}}{1 - e^{-\lambda_{Si}\tilde{T}}} = \frac{\sum_{j=1, j \notin \{S, i\}}^N \lambda_{ij} e^{-\lambda_{ij}(T-\tilde{T})}}{N-2 - \sum_{j=1, j \notin \{S, i\}}^N e^{-\lambda_{ij}(T-\tilde{T})}} \quad (8)$$

which can be solved by Newton's method numerically.

For each node  $N_i$  in the network, upon contact with another node  $N_j$ , it calculates the optimal  $\tilde{T}_0$  for  $N_j$  from Eq. (8), based on their contact rate  $\lambda_{ij}$ .  $N_i$  then calculates

$$\tilde{C}_i = (1 - e^{-\lambda_{ij}\tilde{T}_0}) \left(1 - \frac{1}{N-2} \sum_{k=1, k \notin \{i, j\}}^N e^{-\lambda_{ik}(T-\tilde{T}_0)}\right)$$

and send  $\tilde{C}_i$  to  $N_j$  along with  $C_i$  described in Eq. (3). Since the pairwise inter-contact time is shown to be exponentially distributed, the pairwise node contact rate is going to be invariant over time, which makes it unnecessary to solve Eq. (8) and recompute  $\tilde{C}_i$  repetitively upon every contact. In relay selection, if  $R_i$  is not in contact with  $S$ ,  $\tilde{C}_i$  is used in calculating  $w_i$  instead of  $C_i$ , i.e.,

$$w_i = \log \frac{1}{1 - \tilde{C}_i}.$$

## 4. MULTIPLE-DATA MULTICAST

We exploit a community-based approach to solve the MDM Problem to which localized heuristic is not applicable due to the node buffer constraints and the subsequent requirement of destination-awareness, as discussed in Section 2.2. In our approach, a node

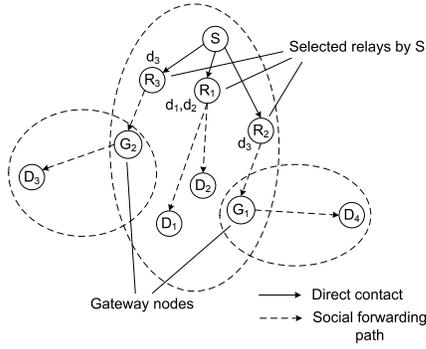


Figure 2: MDM data forwarding process

maintains its destination-awareness about other nodes in the same community. The data source selects relays based on destination-awareness, and places appropriate data items on each relay. Our relay selection scheme ensures that the average probability that a data item is delivered to its destinations by the selected relays is higher than the required  $p$ .

In practice, each data item is forwarded to the destinations by the selected relays according to their local knowledge about the destinations. Such process is illustrated in Figure 2, where  $S$  multicasts three data items  $d_1, d_2, d_3$  to destination sets  $\{D_1\}, \{D_2\}, \{D_3, D_4\}$ , respectively. The relay and data item selections are shown in the figure. For destinations  $D_1$  and  $D_2$  in the same community with relay  $R_1$ ,  $R_1$  forwards  $d_1$  and  $d_2$  to the destinations according to its local destination-awareness about  $D_1$  and  $D_2$ , in the form of social forwarding paths to be described in Section 4.1. For destinations  $D_3$  and  $D_4$  which reside in other communities, data forwarding is conducted through the “gateway” nodes  $G_1$  and  $G_2$ , which belong to multiple communities. Such hierarchical scheme limits inter-community data forwarding to the gateway nodes, and is therefore able to greatly reduce the forwarding cost.

#### 4.1 Social Forwarding Path

We first introduce the concept of social forwarding path.

**DEFINITION 2.** A  $k$ -hop social forwarding path  $P_{AB} = (V_P, E_P)$  between two nodes  $A$  and  $B$  consists of a node set  $V_P = \{A, N_1, N_2, \dots, N_{k-1}, B\}$  and an edge set  $E_P = \{e_1, e_2, \dots, e_k\}$  with edge weights  $\{\lambda_1, \lambda_2, \dots, \lambda_k\}$ . The path weight is the probability  $p_{AB}(T)$  that a data item is forwarded from  $A$  to  $B$  along  $P_{AB}$  within time  $T$ .

We now describe how to determine the weight of a social forwarding path. The inter-contact time  $X_i$  between nodes  $N_i$  and  $N_{i+1}$  follows an exponential distribution with PDF  $p_{X_i}(x) = \lambda_i e^{-\lambda_i x}$  ( $x \geq 0$ ). As a result, the total time needed to transfer data from  $A$  to  $B$  along  $P_{AB}$  is  $Y = \sum_{i=1}^k X_i$ . The PDF  $p_Y(x)$  can be calculated by convolutions on  $p_i(x)$  as

$$p_Y(x) = p_1(x) \otimes p_2(x) \dots \otimes p_k(x).$$

**THEOREM 2.** For a  $k$ -hop social forwarding path with edge weights  $\lambda_1, \lambda_2, \dots, \lambda_k$ , Let  $p_{X_i}(x) = \lambda_i e^{-\lambda_i x}$ ,  $p_Y(x)$  is expressed as

$$p_Y(x) = \sum_{i=1}^k C_i^{(k)} p_{X_i}(x)$$

where the coefficients  $C_i^{(k)} = \prod_{j=1, j \neq i}^k \frac{\lambda_j}{\lambda_j - \lambda_i}$ .

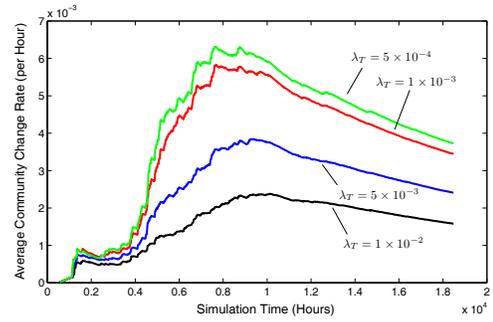


Figure 3: Rate of community changes

PROOF. This theorem is proved in the Appendix.  $\square$

From  $p_Y(x)$ , the weight of social forwarding path  $P_{AB}$  is

$$\begin{aligned} p_{AB}(T) &= P(Y < T) = \int_{-\infty}^T p_Y(x) dx \\ &= \sum_{i=1}^k (C_i^{(k)} \cdot \int_{-\infty}^T p_{X_i}(x) dx) = \sum_{i=1}^k C_i^{(k)} \cdot (1 - e^{-\lambda_i T}) \end{aligned} \quad (9)$$

and a node  $A$  maintains its destination-awareness to another node  $B$  in the form of the social forwarding path  $P_{AB}$ .

#### 4.2 Community-based Destination-Awareness

Each node maintains the “best” social forwarding path with the largest path weight to all the other nodes within the same community. To do this, we assume that each node in the network belongs to at least one social community.

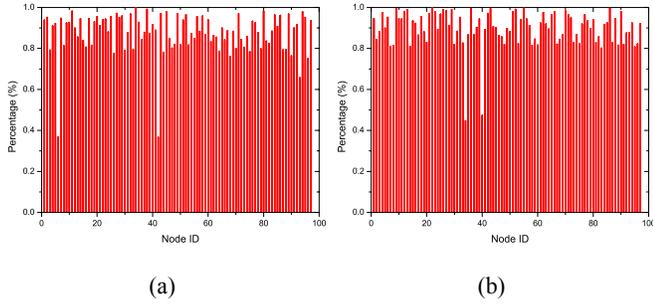
Table 4: Record of social forwarding path table

$D$	Hops $\{N_1, \dots, N_{k-1}, D\}$	edge weights $\{\lambda_1, \dots, \lambda_k\}$
-----	-----------------------------------	--

An efficient community detection mechanism is therefore needed. We use the  $k$ -clique community detection mechanism [18], because they are able to detect overlapping communities. A social community can be defined differently by the community detection mechanisms. According to Palla *et al.* [18], a  $k$ -clique community is defined as a union of all  $k$ -cliques (complete subgraphs of size  $k$ ) that can be reached from each other through a series of adjacent  $k$ -cliques. In the distributed implementation of the  $k$ -clique method for DTNs [10], each node first builds a familiar set containing its contact neighbors, based on specific admission criteria, then builds its local community by merging the familiar sets of other nodes. We adopt such method using the contact rate specified in Section 2.4 as the admission criterion.

Each node maintains a social forwarding path table for all the other nodes within the same community, and the record format for the path to a node  $D$  is shown in Table 4. If a node belongs to multiple communities, a separate table is maintained for each community. Initially, each node only has the information about its contacted neighbors. When a node  $A$  contacts another node  $B$ , they exchange and update their social forwarding path tables.

For a record of node  $C$  in  $B$ ’s social forwarding path table, if  $C$  has not been recorded at  $A$  and is in the same community with  $A$ ,  $A$  adds this record into its own table. Otherwise, if the path to  $C$  recorded by  $B$  has larger weight than that recorded by  $A$ ,  $A$  updates its local record about  $C$ . When updating a record,  $A$  also checks whether itself is on the path to prevent possible loops. If so,



**Figure 4: Coverage of destination-awareness: (a) 10 destinations in total, (b) 40 destinations in total**

$A$  extracts and stores the part  $\{N_1, \dots, N_{k-1}, D\}$  from the looped path  $\{N_0, \dots, N_s, A, N_1, \dots, N_{k-1}, D\}$ .

Compared with node mobility, the long-term social community relations among nodes are much more stable over time. This fact is validated by experimental results on the *MIT Reality* trace, where a community change is an operation adding a node to or deleting a node from a social community. Figure 3 shows that the community change rates are lower than  $7 \times 10^{-3}$  per hour. These results, together with the fact that contact frequencies between nodes in the same community are much higher than the network average level [9], ensure that the community-based destination-awareness can be maintained up-to-date and accurately at individual nodes.

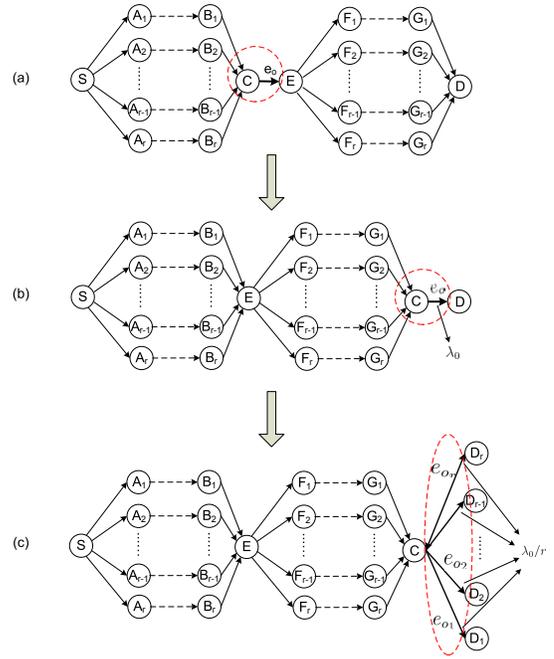
Figure 4 shows the coverage of the community-based destination-awareness from an experiment on 500 MDM scenarios with random sources and destinations. Such coverage is defined as the percentage of destinations that the data source is aware of. It is shown that the data source is aware of most of the destinations, and this result ensures that the data source has enough knowledge about the destinations to select relays effectively. The other destinations and the corresponding data items are considered as a separate SDM problem, and data relays are selected as described in Section 3.2.

### 4.3 Edge Splitting Process

Generally, the probabilities for various relays to forward a data item to the same destination are not independent, due to the possible overlap of their social forwarding paths on some common edges. For example, in Figure 5(a), paths from  $S$  to  $D$  via different relays  $A_1, \dots, A_r$  share the edge  $e_0 = (C, E)$  in common<sup>2</sup>. Such overlap makes it difficult to calculate the cumulative data forwarding probabilities for multiple relays, and we eliminate such overlap by exploiting an edge splitting process.

First, due to the commutativity of convolution, we are able to “move” the common edge  $e_0$  to the “end” of the paths, i.e., the last hop to the destination, without changing the weight of any involved path. This is illustrated in Figure 5(b). Since a node contacts each of its neighbors independently, node  $E$  shared by multiple paths does not affect the independence of the paths.

The edge splitting process is illustrated in Figure 5(c). For an edge  $e_0$  with rate  $\lambda_0$  being shared by  $r$  paths, we split  $e_0$  to  $r$  distinct edges, each of which has rate  $\lambda_0/r$ , and each of the  $r$  paths is allocated a splitted edge. Such process is equivalent to create a virtual destination node for each of the  $r$  paths. In Figure 5(c), letting the weight of the  $i$ -th path  $\{S, A_i, B_i, \dots, G_i, C, D_i\}$  after the edge splitting process be  $p_i$ , the cumulative probability for  $S$  to



**Figure 5: Edge splitting process**

send data to  $D$  within time  $T$  can be equivalently calculated as

$$1 - \prod_{i=1}^r (1 - p_i)$$

which is the probability that  $S$  sends data to at least one node among the virtual destination nodes  $D_i$ .

**THEOREM 3.** *Edge splitting process gives a lower bound to the original data forwarding probability from source to destination.*

**PROOF.** Let the random variable  $T_1$  be the time for  $S$  to send data to  $C$  with PDF  $f_1(t)$ , and  $T_2$  be the time for  $C$  to contact  $D$  in Figure 5(b), which has PDF  $f_2(t) = \lambda_0 e^{-\lambda_0 t}$ , we have

$$P(T_1 + T_2 \leq T) = P(\lambda_0) = \int_0^T f_1(t) \otimes f_2(t) dt.$$

Similarly, the probability for  $S$  to deliver data to any node  $D_i$  within time  $T$  in Figure 5(c) is equally  $P(\lambda_0/r)$ , and therefore the lower bound is written as

$$\prod_{i=1}^r (1 - p_i) \geq (1 - P(\frac{\lambda_0}{r}))^r \geq (1 - P(\lambda_0)). \quad (10)$$

The first inequality in Eq. (10) is obvious, because  $1 - P(\frac{\lambda_0}{r})$  is the probability for  $S$  to send data to a node  $D_i$  via all the  $r$  paths from  $S$  to  $C$  in Figure 5(c), and  $1 - p_i$  is that probability only via the  $i$ -th path. For the second inequality, let  $h(r) = (1 - P(\frac{\lambda_0}{r}))^r$ , it is equivalent to prove that  $\frac{\partial h(r)}{\partial r} \geq 0$  for all  $r > 1$ .

$$\begin{aligned} \frac{\partial h(r)}{\partial r} &= r \left(1 - P\left(\frac{\lambda_0}{r}\right)\right)^{r-1} \cdot \left(-\frac{\partial P\left(\frac{\lambda_0}{r}\right)}{\partial r}\right) \\ &= -r \left(1 - P\left(\frac{\lambda_0}{r}\right)\right)^{r-1} \cdot \frac{\partial P\left(\frac{\lambda_0}{r}\right)}{\partial \left(\frac{\lambda_0}{r}\right)} \cdot \frac{\partial \left(\frac{\lambda_0}{r}\right)}{\partial r} \\ &= \frac{\lambda_0}{r} \left(1 - P\left(\frac{\lambda_0}{r}\right)\right)^{r-1} \cdot \frac{\partial P(\lambda_0)}{\partial \lambda_0}. \end{aligned}$$

<sup>2</sup>The dashed lines in Figure 5(a) indicate multi-hop social forwarding paths.

It is easy to know that the contact rate  $\lambda_0$  between  $C$  and  $D$  is negatively proportional to the time needed for  $C$  to contact  $D$ , and is therefore positively proportional to the probability that  $S$  send data to  $D$  via  $C$ , i.e.,  $\frac{\partial P(\lambda_0)}{\partial \lambda_0} \geq 0$ . So we have  $\frac{\partial h(r)}{\partial r} \geq 0$  for all  $r > 1$ , and the theorem is proved.  $\square$

#### 4.4 Two-Stage Relay Selections

The data source  $S$  selects relays among its contacted neighbors, based on its knowledge about the destinations. Suppose at  $S$  there are data items  $d_1, \dots, d_n$  with sizes  $s_1, \dots, s_n$  and destination sets  $\mathbb{D}_1, \dots, \mathbb{D}_n$ , and  $S$  selects relays among nodes  $R_1, \dots, R_m$  with buffer sizes  $B_1, \dots, B_m$ . The relay selection problem is formulated at  $S$  as the following knapsack problem:

$$\begin{aligned} & \min \left| \{j \mid \sum_{i=1}^n x_{ij} > 0\} \right| \\ & s.t. \sum_{i=1}^n x_{ij} s_i \leq B_j, \text{ for } j = 1, \dots, m \\ & \quad \frac{1}{|\mathbb{D}_i|} \sum_{k \in \mathbb{D}_i} \prod_{j=1}^m (1 - x_{ij} p_{jk}) \leq (1 - p), \text{ for } i = 1, \dots, n \end{aligned}$$

where  $x_{ij} \in \{0, 1\}$  indicates that data item  $d_i$  is placed on relay  $R_j$ , and  $p_{jk}$  is the probability for  $S$  to send data to destination  $k$  via  $R_j$ , in the form of the weight of corresponding social forwarding path from  $S$  to  $k$  via  $R_j$ . The second set of constraints guarantees that, for each data item, the average probability that it is delivered to its destinations is higher than  $p$ .

Such problem is NP-hard due to the second set of constraints, and we instead propose an effective heuristic for relay selection consisting of two stages. First, the optimal data item selection for each  $R_j$  is calculated, based on the node buffer constraint of  $R_j$ . Such optimal data item selection leads to the maximal average probability that a data item is forwarded to its destinations via  $R_j$ . Second, relay selection is conducted using the optimized data forwarding probabilities in the first phase as node weights.

##### 4.4.1 Data Item Selection

On each node  $R_j$ , searching for the optimal data item selection leads to solving the following 0-1 knapsack problem:

$$\begin{aligned} & \max \sum_{i=1}^n \sum_{k \in \mathbb{D}_i} \frac{x_{ij} p_{jk}}{|\mathbb{D}_i|} \\ & s.t. \sum_{i=1}^n x_{ij} s_i \leq B_j \end{aligned}$$

where  $\frac{1}{|\mathbb{D}_i|} \sum_{k \in \mathbb{D}_i} p_{jk}$  is the average probability that data item  $d_i$  is delivered to its destinations from  $S$  via  $R_j$ .

Since  $B_j$  and  $s_i$  can be represented as integers (numbers of bytes), this knapsack problem can be solved in pseudo-polynomial time using a dynamic programming approach [15]. The solution to this problem maximizes the average probability that a data item is delivered to its destinations via  $R_j$  within  $T$ .

##### 4.4.2 Relay Selection

According to the optimal data item selection on  $R_j$ , we define the node weight  $w_j$  for  $R_j$  as

$$w_j = 1 - \frac{1}{n} \sum_{i=1}^n \sum_{k \in \mathbb{D}_i} \frac{x_{ij} p_{jk}}{|\mathbb{D}_i|}$$

which indicates the average probability that a data item cannot be delivered to its destinations by  $R_j$  within  $T$ . Such probability is minimized for  $R_j$  at the data item selection stage.

Then, the node weight  $w_j$  is used for relay selection, satisfying the following performance requirement

$$\prod_{j=1}^m w_j^{x_j} \leq 1 - p \quad (11)$$

where  $x_j \in \{0, 1\}$  indicates whether  $R_j$  is selected as the relay. Similar to the relay selection scheme in SDM, Eq. (11) can also be unified to the knapsack formulation in Eq. (1) by taking logarithms on both sides of the inequality.

## 5. PERFORMANCE EVALUATIONS

In this section, we compare the performance of our SDM and MDM schemes, with the flooding-based approach (Epidemic routing) [22] and the mobility-based approach (PROPHET) [13]. We also compare our SDM scheme with other social-based data forwarding schemes including SimBet [4] and BUBBLE Rap [9] to show the essential difference between multicast and unicast in DTNs. In Epidemic routing, each relay forwards the data to all the contacted nodes. In PROPHET, each data item is forwarded to the nodes with higher delivery predictability, and each of the specified destination is handled separately. We use the default PROPHET parameter settings recommended in [13].

The following metrics are used in our simulations. We only count *delivered destinations*, which are the destinations that have received the data. Each simulation is repeated 500 times with random data sources and destinations for statistical convergence.

1. **Delivery ratio**, the ratio of the number of delivered destinations to the total number of destinations.
2. **Actual delay**, the average delay for all the delivered destinations to receive the data.
3. **Average cost**, the average number of relays used for one delivered destination to receive a data item.

### 5.1 Performance of SDM

We use the *Infocom* trace with higher contact rates to evaluate our SDM scheme. In all simulations, we fix the required delivery ratio  $p = 80\%$ . The data source multicasts a data item to 10 randomly selected destinations, with various time constraints  $T$ . For PROPHET, multicast is handled as separate unicast processes for each destination. We also evaluate the performance of our MDM scheme to the SDM problem, by setting the number of data items to be 1. These results are indicated as ‘‘S-MDM’’ in Figure 6.

Due to the low rates of node contacts in DTNs, the selected relays may not be able to contact the destinations if the time constraints are short. As a result, the actual delivery ratio shown in Figure 6(a) is tightly related to the time constraint. When the time constraint becomes longer, such ratio increases dramatically because the selected relays have more chances to contact destinations, and the average delay increases accordingly.

The delivery ratio of our centrality-based SDM scheme is also limited by the time constraint, as the required delivery ratio  $p$  cannot be achieved when the time constraint is shorter than 14 hours. In such cases, it is most likely that all the available relays for the data source together cannot satisfy the performance requirements in Eq. (5), and the source therefore can only forward the data with best effort by selecting all the available relays. Nevertheless, under various time constraints, our SDM scheme shows only 5% degradation in delivery ratio and delay, compared with Epidemic, and outperforms PROPHET by 20%. Similar results are shown in Fig-

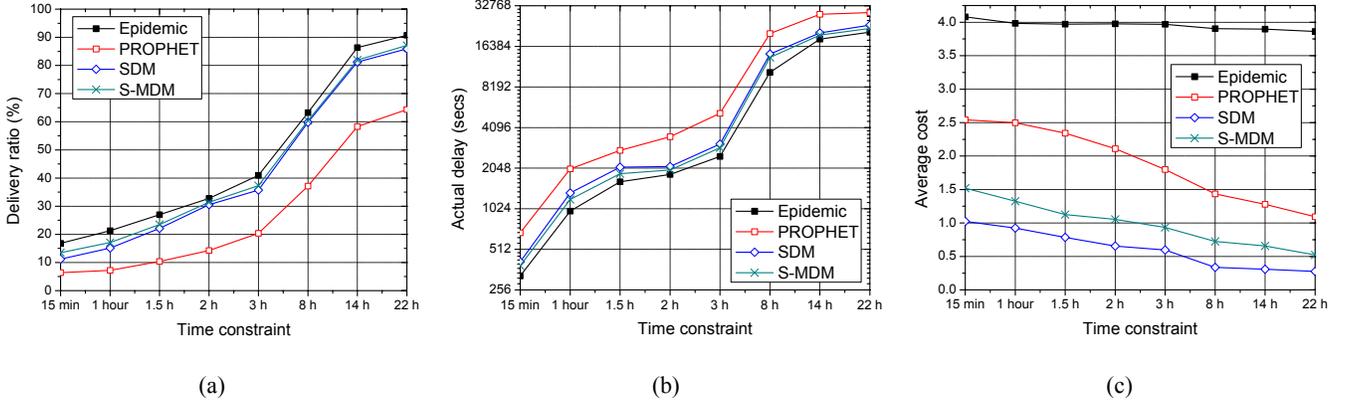


Figure 6: Performance of SDM on the *Infocom* trace: (a) Delivery ratio, (b) Actual delivery delay, (c) Average cost

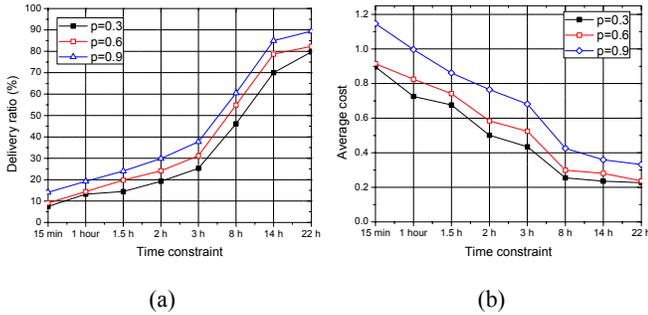


Figure 7: Delivery ratio and cost of SDM with different requirements on delivery ratio

ure 6(b) in terms of delay. The actual delay of SDM is 5% longer than that of Epidemic, but is over 10% shorter than PROPHET.

Figure 6(c) shows that our approach has much less cost than Epidemic and PROPHET. In Epidemic routing, since the data is flooded, the average cost remains at 4 relays for each delivered destination when the time constraint changes. When the time constraint increases from 15 minutes to 22 hours, the average cost of SDM is reduced from 1.0 relays to 0.25 relays. When the time constraint is 22 hours, the cost of SDM is only 25% of the cost of PROPHET, and 6.25% of the cost of Epidemic.

We also apply MDM (S-MDM in the figure) to the SDM problem by selecting relays based on the social forwarding paths to the destinations. Figures 6(a) and 6(b) show that the improvements on the delivery ratio and delay is less than 10%, but the average cost increases by 50%, because the data is going to be forwarded along the social forwarding paths. In practice, the average cost is even higher due to the maintenance of social forwarding paths and community structures. Therefore for the SDM problem, our centrality-based localized heuristic is cost-effective.

In Figure 7, the effects of different delivery requirements ( $p$ ) on SDM are investigated. When  $p$  is low, increasing  $p$  leads to a considerable improvement of the actual delivery ratio, as shown by Figure 7(a). When  $p$  increases from 0.3 to 0.6, the actual delivery ratio increases by 15%-25%. Such improvement becomes smaller when  $p$  is high. Correspondingly, higher  $p$  requires the data source to select more relays. Figure 7(b) shows that when  $p$  increases from 0.3 to 0.9, the average cost increases by 25%-30%.

## 5.2 Performance of MDM

We use the *MIT Reality* trace with larger network scale to evaluate our MDM scheme. We fix the required delivery ratio  $p = 60\%$ ,

and 5 data items are generated at the source node. The number of destinations for each data item is uniformly randomized in the range  $[3, 9]$ . Letting the total size of all the data items be  $S$ , the buffer size of each node is uniformly randomized in  $[\frac{1}{2}S, S]$ .

In Epidemic routing, each data item is flooded in the network. In PROPHET, each data item is forwarded to each destination as a separate unicast process. We also evaluate the performance of our SDM scheme to the MDM problem by multicasting each data item separately, and such results are indicated as “M-SDM” in Figure 8. In the three cases, data items are randomly selected at a relay when it cannot carry all the data items simultaneously.

Since the pairwise node contact rates in the *MIT Reality* trace is much lower (see Table 1), the time constraint and actual delay are much longer correspondingly. Similar to the results in SDM performance evaluation, the delivery ratio shown in Figure 8(a) is also highly related with the time constraint. Our MDM scheme can only achieve the required delivery ratio  $p = 60\%$  when the time constraint is longer than 6 weeks, but it keeps similar delivery ratio with Epidemic routing, and outperforms PROPHET over 100% when the time constraint is longer than 1 week. Meanwhile, the average cost of our approach is much lower than that of Epidemic and PROPHET, as shown in Figure 8(c). For the longest time constraint (6 months), the average cost of our approach is only 50% of that of PROPHET, and 11% of that of Epidemic.

The major difference between SDM and MDM is maintaining destination-awareness due to buffer constraints. When being used for the MDM problem, our SDM scheme selects data items for each relay at random, and therefore a selected relay may have low forwarding probabilities to the destinations of the data items it carries. Such random data item selection leads to reduction on the delivery ratio by 10%-20% and slight increase on the actual delay, as shown in Figures 8(a) and 8(b). The average cost is also 20%-25% higher in Figure 8(c), because the average number of destinations that a relay can deliver is smaller. Therefore, our MDM scheme shows the advantage of maintaining destination-awareness which helps select data items optimally for each relay.

## 5.3 Comparison with other social-based schemes

In this section, we compare the performance of our SDM scheme with other social-based forwarding schemes including SimBet [4] and BUBBLE Rap [9]. We fix the required delivery ratio of our SDM scheme as 80%, and 10 destinations are randomly selected. In SimBet and BUBBLE Rap, each destination is handled separately as unicast. We use the same parameter settings for community detection in BUBBLE Rap as in our MDM scheme.

SimBet calculates the betweenness and similarity metrics for

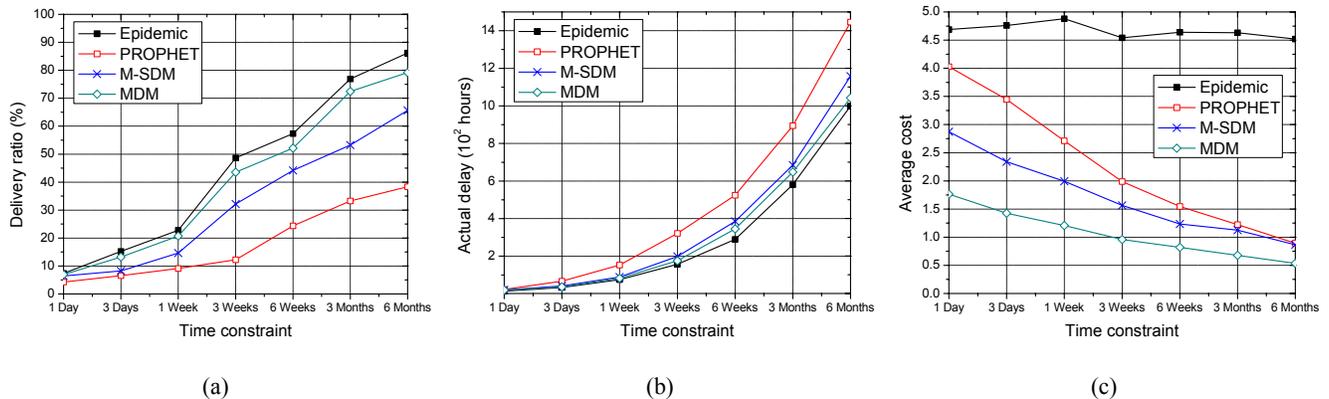


Figure 8: Performance of MDM on the *MIT Reality* trace: (a) Delivery ratio, (b) Actual delay, (c) Average cost

mobile nodes, and exchanges data between nodes based on their SimBet utilities. It does not consider contact frequencies between node pairs, and therefore leads to lower delivery ratio and longer delay, as shown in Figures 9(a) and 9(b). Comparatively, BUBBLE Rap benefits from the consideration of social community structure and hierarchical forwarding. It has similar delay to our SDM scheme, but still has a 20% lower delivery ratio. This is mainly because BUBBLE Rap also uses betweenness as centrality metric.

The essential difference between multicast and unicast in DTNs is illustrated in Figure 9(c). For multicast, a relay is required to deliver data to as many destinations as possible, and such requirement leads to different considerations in relay selection. Therefore, SimBet and BUBBLE Rap produces larger cost than our SDM scheme, because they select relays for each destination separately. The cost of BUBBLE Rap is smaller than that of SimBet because of exploiting hierarchical forwarding architecture, but its cost is still 50% higher than that of our SDM scheme. These results show that unicast schemes cannot be simply exploited for effective multicast.

## 6. RELATED WORK

In [22], the authors studied Epidemic routing in DTNs. The basic idea is to select all nodes in the network as relays. Some later work studied relay selection strategies based on node mobility patterns. For example, PROPHET [13] calculates the delivery predictability at each node by using encounter history, [25] employs some nodes with desirable mobility patterns as message ferries, [20, 3] analyze the performance of mobility-assisted schemes theoretically, and [5] provides a unified approach on mobility-based metrics. Some works make efforts on improving data forwarding performance by either determining the data delivery likelihood [2] or spraying data to relays waiting for contacts with destinations [19], which is similar with our SDM scheme. However, only simple heuristics are provided for selecting relays in these approaches.

Since node mobility patterns are highly volatile and hard to control, attempts on exploiting stable social network structure for data forwarding have emerged. Most social-based forwarding schemes exploit sociological centrality metrics [14] for relay selections. SimBet routing [4] uses ego-centric betweenness metric and forwards data to nodes with higher SimBet utility. Later work [9] considers node centrality in a hierarchical manner based on social community knowledge. Social communities are also investigated in a decentralized way [10] for publisher/subscriber applications [24].

The network contact graph in most social-based data forwarding schemes is considered as binary. Hui et al. [9] used cumulative contact length as edge weights for community detection in social networks, but did not exploit such weights for data forwarding.

The binary social network model considers node pairs with different contact frequencies as equivalent ones, and limits the performance of centrality-based data forwarding schemes because node centrality values do not really characterize the nodes' capabilities of contacting other nodes.

Some other work [1, 11] focus on modeling the content dissemination process in DTNs in an epidemic manner. In [1], efficient utility functions are developed for content dissemination, and [11] investigates optimal rate allocation schemes to maximize the data dissemination speed. However, the content dissemination processes are not oriented for specified destinations as in our work. As a result, research on content dissemination focuses on the optimal network design for improving dissemination speed, rather than relay selection schemes for better cost-effectiveness.

## 7. CONCLUSIONS

In this paper, we studied multicast in DTNs from the social network perspective, and exploited social network concepts, including centrality and social community, to improve the cost-effectiveness of multicast in DTNs. We investigated the essential difference between multicast and unicast in DTNs, and developed relay selection schemes considering the forwarding probabilities to multiple destinations simultaneously. Trace driven simulation results show that our approach achieves similar delivery ratio and delay to Epidemic routing, but significantly reduces the forwarding cost. We believe that this paper presents the first step in exploiting social network methods for efficient multi-party communication in DTNs. Further research can benefit from our results by developing specific applications based on the provided multicast architecture.

## Acknowledgment

We would like to thank Wenye Wang for her insightful comments and suggestions. We would also like to thank the anonymous reviewers for their helpful suggestions. This work was supported in part by the National Science Foundation under grant CNS-0721479.

## 8. REFERENCES

- [1] C. Boldrini, M. Conti, and A. Passarella. Modelling data dissemination in opportunistic networks. *Proc. ACM Workshop on Challenged Networks (CHANTS)*, 2008.
- [2] J. Burgess, B. Gallagher, D. Jensen, and B. Levine. Maxprop: Routing for vehicle-based disruption-tolerant networks. *Proc. INFOCOM*, 2006.
- [3] A. Chaintreau, P. Hui, J. Crowcroft, C. Diot, R. Gass, and J. Scott. Impact of Human Mobility on Opportunistic Forwarding Algorithms. *IEEE Trans. on Mobile Computing*, pages 606–620, 2007.

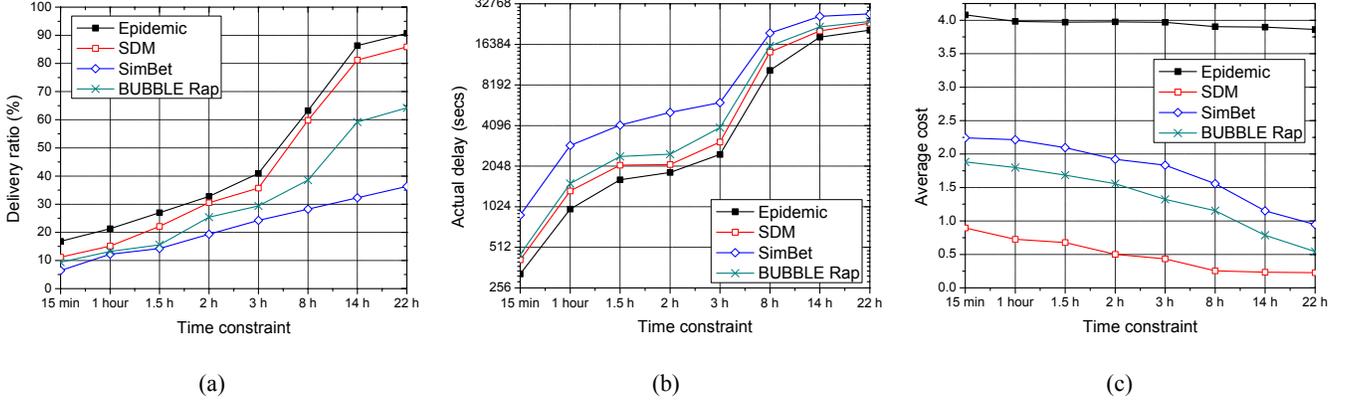


Figure 9: Comparison with other social-based schemes: (a) Delivery ratio, (b) Actual delay, (c) Average cost

[4] E. Daly and M. Haahr. Social network analysis for routing in disconnected delay-tolerant MANETs. *Proc. MobiHoc*, 2007.

[5] V. Erramilli, A. Chaintreau, M. Crovella, and C. Diot. Delegation Forwarding. *Proc. MobiHoc*, 2008.

[6] K. Fall. A delay-tolerant network architecture for challenged internets. *Proc. SIGCOMM*, pages 27–34, 2003.

[7] P. Greenwood and M. Nikulin. *A Guide to Chi-Squared Testing*. Wiley-Interscience, 1996.

[8] W. Hsu and A. Helmy. On Nodal Encounter Patterns in Wireless LAN Traces. *Proc. International Workshop On Wireless Network Measurement (WinMee)*, 2006.

[9] P. Hui, J. Crowcroft, and E. Yoneki. Bubble rap: social-based forwarding in delay tolerant networks. *Proc. MobiHoc*, pages 241–250, 2008.

[10] P. Hui, E. Yoneki, S. Chan, and J. Crowcroft. Distributed community detection in delay tolerant networks. *Proc. MobiArch*, 2007.

[11] S. Ioannidis, A. Chaintreau, and L. Massoulie. Optimal and scalable distribution of content updates over a mobile social network. *Proc. INFOCOM*, 2009.

[12] U. Lee, S.-Y. Oh, K.-W. Lee, and M. Gerla. Scalable multicast routing in delay tolerant networks. *Proc. ICNP*, 2008.

[13] A. Lindgren, A. Doria, and O. Schelen. Probabilistic routing in intermittently connected networks. *ACM SIGMOBILE CCR*, 7(3):19–20, 2003.

[14] P. Marsden. Egocentric and sociocentric measures of network centrality. *Social Networks*, 24(4):407–422, 2002.

[15] S. Martello and P. Toth. *Knapsack problems: algorithms and computer implementations*. John Wiley & Sons.

[16] S. Milgram. The small world problem. *Psychology Today*, 2(1):60–67, 1967.

[17] M. Motani, V. Srinivasan, and P. Nuggehalli. PeopleNet: engineering a wireless virtual social network. *Proc. MobiCom*, pages 243–257, 2005.

[18] G. Palla, I. Derényi, I. Farkas, and T. Vicsek. Uncovering the overlapping community structure of complex networks in nature and society. *Nature*, 435(7043):814–818, 2005.

[19] T. Spyropoulos, K. Psounis, and C. Raghavendra. Spray and wait: an efficient routing scheme for intermittently connected mobile networks. In *Proceedings of 2005 ACM SIGCOMM workshop on Delay-tolerant networking*, pages 252–259, 2005.

[20] T. Spyropoulos, K. Psounis, and C. Raghavendra. Performance analysis of mobility-assisted routing. *Proc. MobiHoc*, pages 49–60, 2006.

[21] V. Srinivasan, M. Motani, and W. Ooi. Analysis and implications of student contact patterns derived from campus schedules. *Proc. MobiCom*, pages 86–97, 2006.

[22] A. Vahdat and D. Becker. Epidemic routing for partially connected ad hoc networks. *Technical Report CS-200006*. Duke University, 2000.

[23] D. Watts and S. Strogatz. Collective dynamics of “small-world” networks. *Nature*, 393:440–442, 1998.

[24] E. Yoneki, P. Hui, S. Chan, and J. Crowcroft. A socio-aware overlay

for publish/subscribe communication in delay tolerant networks. *Proc. MSWiM*, pages 225–234, 2007.

[25] W. Zhao, M. Ammar, and E. Zegura. A message ferrying approach for data delivery in sparse mobile ad hoc networks. In *Proc. MobiHoc*, pages 187–198, 2004.

[26] W. Zhao, M. Ammar, and E. Zegura. Multicasting in delay tolerant networks: semantic models and routing algorithms. *Proc. SIGCOMM Workshop on Delay Tolerant Networking*, pages 268–275, 2005.

## APPENDIX

**Proof of Theorem 2:** We prove this theorem by induction. First, consider  $P_2 = p_1(x) \otimes p_2(x)$ , we have

$$\begin{aligned} P_2 &= \lambda_1 \lambda_2 \int_0^x e^{-(\lambda_1 - \lambda_2)t} e^{-\lambda_2 x} dt \\ &= \frac{\lambda_1 \lambda_2}{\lambda_2 - \lambda_1} (e^{-\lambda_1 x} - e^{-\lambda_2 x}) = -C_1^{(2)} p_1(x) - C_2^{(2)} p_2(x) \end{aligned}$$

This result can also be applied to two arbitrary functions  $p_i(x) \otimes p_j(x)$ . Suppose  $P_{k-1} = \sum_{i=1}^{k-1} C_i^{(k-1)} p_i(x)$ , then

$$\begin{aligned} P_k &= P_{k-1} \otimes p_k(x) = \sum_{i=1}^{k-1} C_i^{(k-1)} p_i(x) \otimes p_k(x) \\ &= \sum_{i=1}^{k-1} C_i^{(k-1)} \cdot \left( \frac{\lambda_k}{\lambda_k - \lambda_i} p_i(x) + \frac{\lambda_i}{\lambda_i - \lambda_k} p_k(x) \right) \end{aligned}$$

Consider that  $C_i^{(k)} = C_i^{(k-1)} \cdot \frac{\lambda_k}{\lambda_k - \lambda_i}$ , we have

$$P_k = \sum_{i=1}^{k-1} C_i^{(k)} p_i(x) + \sum_{i=1}^{k-1} C_i^{(k-1)} \frac{\lambda_i}{\lambda_i - \lambda_k} p_k(x) \quad (12)$$

For the second term in Eq. (12), we have

$$\begin{aligned} \sum_{i=1}^{k-1} C_i^{(k-1)} \frac{\lambda_i}{\lambda_i - \lambda_k} &= \sum_{i=1}^{k-1} \frac{\lambda_i}{\lambda_i - \lambda_k} \cdot \left( \prod_{j=1, j \neq i}^{k-1} \frac{\lambda_j}{\lambda_j - \lambda_i} \right) \\ &= \prod_{j=1}^{k-1} \lambda_j \cdot \sum_{i=1}^{k-1} \prod_{j=1, j \neq i}^{k-1} \frac{1}{\lambda_j - \lambda_i} = \prod_{j=1}^{k-1} \frac{\lambda_j}{\lambda_j - \lambda_k} = C_k^{(k)} \end{aligned}$$

Therefore, we have

$$P_k = \sum_{i=1}^{k-1} C_i^{(k)} p_i(x) + C_k^{(k)} p_k(x) = \sum_{i=1}^k C_i^{(k)} p_i(x)$$

which proves the theorem.