# Social Network Analysis for Information Flow in Disconnected Delay-Tolerant MANETs

Elizabeth M. Daly and Mads Haahr, Member, IEEE

Abstract—Message delivery in sparse mobile ad hoc networks (MANETs) is difficult due to the fact that the network graph is rarely (if ever) connected. A key challenge is to find a route that can provide good delivery performance and low end-to-end delay in a disconnected network graph where nodes may move freely. We cast this challenge as an information flow problem in a social network. This paper presents social network analysis metrics that may be used to support a novel and practical forwarding solution to provide efficient message delivery in disconnected delay-tolerant MANETs. These metrics are based on social analysis of a node's past interactions and consists of three locally evaluated components: a node's "betweenness" centrality (calculated using ego networks), a node's social "similarity" to the destination node, and a node's tie strength relationship with the destination node. We present simulations using three real trace data sets to demonstrate that by combining these metrics delivery performance may be achieved close to Epidemic Routing but with significantly reduced overhead. Additionally, we show improved performance when compared to PRoPHET Routing.

Index Terms—Delay- and disruption-tolerant networks, MANETs, sparse networks, ego networks, social network analysis.

## **1** INTRODUCTION

mobile ad hoc network (MANET) is a dynamic  ${f A}$  wireless network with or without fixed infrastructure. Nodes may move freely and organize themselves arbitrarily [9]. Sparse MANETs are a class of ad hoc networks in which the node population is sparse, and the contacts between the nodes in the network are infrequent. As a result, the network graph is rarely, if ever, connected and message delivery must be delay tolerant. Traditional MANET routing protocols such as AODV [45], DSR [27], DSDV [46], and LAR [29] make the assumption that the network graph is fully connected and fail to route messages if there is not a complete route from source to destination at the time of sending. One solution to overcome this issue is to exploit node mobility in order to carry messages physically between disconnected parts of the network. These schemes are sometimes referred to as mobility-assisted routing that employ the store-carry-and-forward model. Mobility-assisted routing consists of each node independently making forwarding decisions that take place when two nodes meet. A message gets forwarded to encountered nodes until it reaches its destination.

Current research supports the observation that encounters between nodes in real environments do not occur randomly [24] and that nodes do not have an equal probability of encountering a set of nodes. In fact, one study by Hsu and Helmy observed that nodes never encountered more than 50 percent of the overall population [23]. As a consequence, not all nodes are equally likely to encounter each other, and nodes need to assess the probability that they will encounter the destination node. Additionally, Hsu and Helmy performed an analysis on real-world encounters based on network traffic traces of different university campus wireless networks [22]. Their analysis found that node encounters are sufficient to build a connected relationship graph, which is a small-world graph. Therefore, social analysis techniques are promising for estimating the social structure of node encounters in a number of classes of disconnected delay-tolerant MANETs (DDTMs).

Social networks exhibit the small-world phenomenon, which comes from the observation that individuals are often linked by a short chain of acquaintances. The classic example is Milgrams' 1967 experiment, where 60 letters were sent to various people located in Nebraska to be delivered to a stockbroker located in Boston [40]. The letters could only be forwarded to someone whom the current letter holder knew by first name and who was assumed to be more likely than the current holder to know the person to whom the letters were addressed. The results showed that the median chain length of intermediate letter holders was approximately 6, giving rise to the notion of "six degrees of separation." Milgram's experiment showed that the characteristic path length in the real world can be short. Of particular interest, however, is that the participants did not send on the letters to the next participant randomly but sent the letter to a person they perceived might be a good carrier for the message based on their own local information. In order to harness the benefits of small-world networks for the purposes of message delivery, a mechanism for intelligently selecting good carriers based on local information must be explored. In this paper, we propose the use of social network analysis techniques in order to exploit the underlying social structure in order to provide information flow from source to destination in a DDTM, which extends on the authors' previous work [10].

<sup>•</sup> The authors are with the Distributed Systems Group, Computer Science Department, O'Reilly Institute, Trinity College Dublin, Dublin 2, Ireland. E-mail: {elizabeth.daly, Mads.Haahr}@cs.tcd.ie.

Manuscript received 12 Oct. 2007; revised 13 June 2008; accepted 21 Oct. 2008; published online 6 Nov. 2008.

For information on obtaining reprints of this article, please send e-mail to: tmc@computer.org, and reference IEEECS Log Number TMC-2007-10-0309. Digital Object Identifier no. 10.1109/TMC.2008.161.

The remainder of this paper is organized as follows: Section 2 reviews related work in the area of message delivery in disconnected networks. Section 3 examines network theory that may be applied to social networks along with social network analysis techniques. Section 4 discusses SimBetTS, a sample routing protocol, which applies these techniques for routing in DDTMs. Section 5 evaluates the performance of the protocol along with a performance comparison between SimBetTS Routing and Epidemic Routing [51] and the PRoPHET Routing protocol [36] using three real trace data sets from the Haggle project [8], [24]. We conclude in Section 6.

## 2 RELATED WORK

A number of projects attempt to enable message delivery by using a virtual backbone with nodes carrying the data through disconnected parts of the network [15], [47]. The Data MULE project uses mobile nodes to collect data from sensors, which is then delivered to a base station [47]. The Data MULEs are assumed to have sufficient buffer space to hold all data until they pass a base station. The approach is similar to the technique used in [2], [15], and [17]. These projects study opportunistic forwarding of information from mobile nodes to a fixed destination. However, they do not consider opportunistic forwarding between the mobile nodes.

"Active" schemes go further in using nodes to deliver data by assuming control or influence over node movements. Li and Rus [32] explore message delivery where nodes can be instructed to move in order to transmit messages in the most efficient manner. The message ferrying project [54] proposes proactively changing the motion of nodes in order to meet a known "message ferry" to help deliver data. Both assume control over node movements and, in the case of message ferries, knowledge of the paths to be taken by these message ferry nodes.

Other work utilizes a time-dependent network graph in order to efficiently route messages. Jain et al. [26] assume knowledge of connectivity patterns where exact timing information of contacts is known and then modifies Dijkstra's algorithm to compute the cost edges and routes accordingly. Merugu et al. [39] and Handorean et al. [21] likewise make the assumption of detailed knowledge of node future movements. This information is time dependent and routes are computed over the time-varying paths available. However, if nodes do not move in a predictable manner or are delayed, then the path is broken. Additionally, if a path to the destination is not available using the time-dependent graph, the message is flooded.

Epidemic Routing [51] provides message delivery in disconnected environments where no assumptions are made with regard to control over node movements or knowledge of the network's future topology. Each host maintains a buffer containing messages. Upon meeting, the two nodes exchange summary vectors to determine which messages held by the other have not been seen before. They then initiate a transfer of new messages. In this way, messages are propagated throughout the network. This method guarantees delivery if a route is available but is expensive in terms of resources since the network is essentially flooded. Attempts to reduce the number of copies of the message are explored in [44] and [49]. Ni et al. [44] take a simple approach to reduce the overhead of flooding by only forwarding a copy with some probability p < 1, which is essentially randomized flooding. The Spray-and-Wait solution presented by Spyropoulos et al. [49] assigns a replication number to a message and distributes message copies to a number carrying nodes and then waits until a carrying node meets the destination.

A number of solutions employ some form of "probability to deliver" metric in order to further reduce the overhead associated with Epidemic Routing by preferentially routing to nodes deemed most likely to deliver. These metrics are based on either contact history, location information, or utility metrics. Burgess et al. [7] transmit messages to encountered nodes in the order of probability for delivery, which is based on contact information. However, if the connection lasts long enough, all messages are transmitted, thus turning into standard Epidemic Routing. PRoPHET Routing [36] is also probability based, using past encounters to predict the probability of meeting a node again, nodes that are encountered frequently have an increased probability whereas older contacts are degraded over time. Additionally, the transitive nature of encounters is exploited where nodes exchange encounter probabilities and the probability of indirectly encountering the destination node is evaluated. Similarly, Khelil et al. [28] and Tan et al. [50] define probability based on node encounters in order to calculate the cost of the route. In other work, Dubois-Ferriere et al. [11] and Grossglauser and Vetterli [20] use the so-called "time elapsed since last encounter" or the "last encounter age" to route messages to destinations. In order to route a message to a destination, the message is forwarded to the neighbor who encountered the destination more recently than the source and other neighbors.

Lebrun et al. [30] propose a location-based routing scheme that uses the trajectories of mobile nodes to predict their future distance to the destination and passes messages to nodes that are moving in the direction of the destination. Leguay et al. [31] present a virtual coordinate system where the node coordinates are composed of a set of probabilities, each representing the chance that a node will be found in a specific location. This information is then used to compute the best available route. Similarly, Ghosh et al. propose exploiting the fact that nodes tend to move between a small set of locations, which they refer to as "hubs" [16]. A list of "hubs" specific to each user's movement profile is assumed to be available to each node on the network in the form of a "probabilistic orbit," which defines the probability with which a given node will visit a given hub. Messages destined for a specific node are routed toward one of these user-specific "hubs."

Musolesi et al. [41] introduce a generic method that uses Kalman filters to combine multiple dimensions of a node's connectivity context in order to make routing decisions. Messages are passed from one node to a node with a higher "delivery metric." The messages for unknown destinations are forwarded to the "most mobile" node available. Spyropoulos et al. [48] use a combination of random walk and utility-based forwarding. Random walk is used until a node with a sufficiently high utility metric is found after which the utility metric is used to route to the destination node. More recently, Hui and Crowcroft [25] investigated assigning labels to nodes identifying group membership. Messages are only forwarded to nodes in the same group as the destination node.

Our work is distinct in that the SimBetTS Routing metric is comprised of both a node's centrality and its social similarity. Consequently, if the destination node is unknown to the sending node or its contacts, the message is routed to a structurally more central node where the potential of finding a suitable carrier is dramatically increased. We will show that SimBetTS Routing improves upon encounter-based strategies where direct or indirect encounters may not be available.

## **3** Social Networks for Information Flow

In a disconnected environment, data must be forwarded using node encounters in order to deliver data to a destination. The problem of message delivery in disconnected delay-tolerant networks can be modeled as the flow of information over a dynamic network graph with timevarying links. This section reviews network theory that may be applied to social networks along with social network analysis techniques. These techniques have yet to be applied to the context of routing in DDTMs. Social network analysis is the study of relationships between entities and on the patterns and implications of these relationships. Graphs may be used to represent the relational structure of social networks in a natural manner. Each of the nodes may be represented by a vertex of a graph. Relationships between nodes may be represented as edges of the graph.

## 3.1 Network Centrality for Information Flow

Centrality in graph theory and network analysis is a quantification of the relative importance of a vertex within the graph (for example, how important a person is within a social network). The centrality of a node in a network is a measure of the structural importance of the node; typically, a central node has a stronger capability of connecting other network members. There are several ways to measure centrality. Three widely used centrality measures are Freeman's degree, closeness, and betweenness measures [13], [14].

"Degree" centrality is measured as the number of direct ties that involve a given node [14]. A node with high degree centrality maintains contacts with numerous other network nodes. Such nodes can be seen as popular nodes with large numbers of links to others. As such, a central node occupies a structural position (network location) that may act as a conduit for information exchange. In contrast, peripheral nodes maintain few or no relations and thus are located at the margins of the network. Degree centrality for a given node  $p_i$ , where  $a(p_i, p_k) = 1$  if a direct link exists between  $p_i$ and  $p_k$ , is calculated as

$$C_D(p_i) = \sum_{k=1}^{N} a(p_i, p_k).$$
 (1)

"Closeness" centrality measures the reciprocal of the mean geodesic distance  $d(p_i, p_k)$ , which is the shortest path between a node  $p_i$  and all other reachable nodes [14]. Closeness centrality can be regarded as a measure of how long it will take information to spread from a given node to other nodes in the network [43]. Closeness centrality for a given node, where N is the number of reachable nodes in the network, is calculated as

$$C_C(p_i) = \frac{N-1}{\sum_{k=1}^N d(p_i, p_k)}.$$
(2)

"Betweenness" centrality measures the extent to which a node lies on the geodesic paths linking other nodes [13], [14]. Betweenness centrality can be regarded as a measure of the extent to which a node has control over information flowing between others [43]. A node with a high betweenness centrality has a capacity to facilitate interactions between nodes it links. In our case, it can be regarded as how much a node can facilitate communication to other nodes in the network. Betweenness centrality, where  $g_{jk}$  is the total number of geodesic paths linking  $p_j$  and  $p_k$ , and  $g_{jk}(p_i)$  is the number of those geodesic paths that include  $p_i$ , is calculated as

$$C_B(p_i) = \sum_{j=1}^{N} \sum_{k=1}^{j-1} \frac{g_{jk}(p_i)}{g_{jk}}.$$
(3)

Borgatti analyzes centrality measures for flow processes in network graphs [5]. A number of different flow processes are considered, such as package delivery, gossip, and infection. He then analyzes each centrality measure in order to evaluate the appropriateness of each measure for different flow processes. His analysis showed that betweenness centrality and closeness centrality were the most appropriate metrics for message transfer that can be modeled as a package delivery.

Freeman's centrality metrics are based on analysis of a complete and bounded network, which is sometimes referred to as a sociocentric network. These metrics become difficult to evaluate in networks with a large node population as they require complete knowledge of the network topology. For this reason, the concept of "ego networks" has been introduced. Ego networks can be defined as a network consisting of a single actor (ego) together with the actors they are connected to (alters) and all the links among those alters. Consequently, ego network analysis can be performed locally by individual nodes without complete knowledge of the entire network. Marsden introduces centrality measures calculated using ego networks and compares these to Freeman's centrality measures of a sociocentric network [37]. Degree centrality can easily be measured for an ego network where it is a simple count of the number of contacts. Closeness centrality is uninformative in an ego network, since by definition an ego network only considers nodes directly related to the ego node; consequently by definition, the hop distance from the ego node to all other nodes in the ego network is 1. On the other hand, betweenness centrality in ego networks has shown to be quite a good measure when compared to that of the sociocentric measure. Marsden calculates the egocentric and the sociocentric betweenness centrality for the network shown in Fig. 1.



Fig. 1. Bank wiring room network [25].

The betweenness centrality  $C_B(p_i)$  based on the egocentric measures does not correspond perfectly to that based on sociocentric measures. However, it can be seen that the ranking of nodes based on the two types of betweenness is identical in this network. This means that two nodes may compare their own locally calculated betweenness value, and the node with the higher betweenness value can be determined. In effect, the betweenness value captures the extent to which a node connects nodes that are themselves not directly connected. For example, in the network shown in Fig. 1, w9 has no connection with w4. The node with the highest betweenness value connected to w9 is w7, so if a message is forwarded to w7, the message can then be forwarded to w5, which has a direct connection with w4. In this way, betweenness centrality may be used to forward messages in a network. Marsden compared sociocentric and egocentric betweenness for 15 other sample networks and found that the two values correlate well in all scenarios. This correlation is also supported by Everett and Borgatti [12].

Routing based on betweenness centrality provides a mechanism for information to flow from source to destination in a social network. However, routing based on centrality alone presents a number of drawbacks. Yan et al. analyzed routing in complex networks and found that routing based on centrality alone causes central nodes to suffer severe traffic congestion as the number of accumulated packets increases with time, because the capacities of the nodes for delivering packets are limited [53]. Additionally, centrality does not take into account the time-varying nature of the links in the network and the availability of a link. In terms of information flow, a link that is available is one that is "activated" for information flow.

#### 3.2 Strong Ties for Information Flow

The previous section's discussion of information flow based on centrality measures does not take into account the strength of the links between nodes. In terms of graph theory, where the links in the network are time varying, a link to a central node may not be highly available. Brown and Reingen explored information flow in word-of-mouth networks and observe that it is unlikely that each contact representing potential sources of information has an equal probability of being activated for the flow of information [6]. They hypothesize that tie strength is a good measure of whether a tie will be activated, since strong ties are typically more readily available and result in more frequent interactions through which the transfer of information may arise. In a network where a person's contacts consisted of both strong and weak tie contacts, Brown and Reingen found that strong ties were more likely to be activated for information flow when compared to weak ties.

Tie strength is a quantifiable property that characterizes the link between two nodes. The notion of tie strength was first introduced by Granovetter in 1973. Granovetter suggested that the strength of a relationship is dependent on four components: the frequency of contact, the length or history of the relationship, contact duration, and the number of transactions. Granovetter defined tie strength as "the amount of time, the emotional intensity, the intimacy (mutual confiding), and the reciprocal services, which characterize a tie" [18]. Marsden and Campbell extended upon these measures and also proposed a measure based on the depth of a relationship referred to as the "multiple social context" indicator. Lin et al. proposed using the recency of a contact to measure tie strength [34]. The tie strength indicators are defined as follows:

**Frequency.** Granovetter observes that "the more frequently persons interact with one another, the stronger their sentiments of friendship for one another are apt to be" [18]. This metric was also explored in [3], [4], [18], [35], and [38].

**Intimacy/Closeness.** This metric corresponds to Granovetter's definition of the time invested into a social contact as a measure for a social tie [4], [18], [38]. A contact with which a great deal of time has been spent can be deemed an important contact.

Long period of time (longevity). This metric corresponds to Granovetter's definition of the time commitment into a social contact as a measure for a social tie [4], [18], [38]. A contact with which a person has interacted over a longer period of time may be more important than a newly formed contact.

**Reciprocity.** Reciprocity is based on the notion that a valuable contact is one that is reciprocated and seen by both members of the relationship to exist. Granovetter discusses the social example with the absence of a substantial relationship, for example, a "nodding" relationship between people living on the same street [4], [18]. He observes that this sort of relationship may be useful to distinguish from the absence of any relationship.

**Recency.** Important contacts should have interacted with a user recently [34]. This relates to Granovetter's amount of time component and investing in the relationship, where a strong relationship needs investment of time to maintain the intimacy.

Multiple social context. Marsden and Campbell discuss using the breadth of topics discussed by friends as a measure to represent the intimacy of a contact [4], [38].

Mutual confiding (trust). This indicator can be used as a measure of trust in a contact [18], [38].

Routing based on tie strength in network terms is routing based on the most available links. A combination of the tie strength indicators can be used for information flow to determine which contact has the strongest social relationship to a destination. In this manner, messages can be forwarded through links possessing the strongest relationship, as a link representing a strong relationship more likely will be activated for information flow than a weak link with no relationship with the destination. These social measures lend themselves well to a disconnected network by providing a local view of the network graph as they are based solely on observed link events and require no global knowledge of the network.

However, Granovetter argued the utility of using weak ties for information flow in social networks [18]. He emphasized that weak ties lead to information dissemination *between* groups. He introduced the concept of "bridges," observing that

information can reach a larger number of people, and traverse a greater social distance when passed through weak ties rather than strong ties ... those who are weakly tied are more likely to move in circles different from our own and will thus have access to information different from that which we receive [18].

Consequently, it is important to identify contacts that may act as potential bridges. Betweenness centrality is a mechanism for identifying such bridges. Granovetter differentiates between the usefulness of weak and strong ties, "weak ties provide people with access to information and resources beyond those available in their own social circle; but strong ties have greater motivation to be of assistance and are typically more easily available." As a result, routing based on a combination of strong ties and identified bridges is a promising trade-off between the two solutions.

#### 3.3 Tie Predictors

Marsden and Campbell distinguished between indicators and predictors [38]. Tie strength evaluates already existing connections whereas predictors use information from the past to predict likely future connections. Granovetter argues that strong tie networks exhibit a tendency toward transitivity, meaning that there is a heightened probability of two people being acquainted, if they have one or more other acquaintances in common [18]. In literature, this phenomenon is called "clustering." Watts and Strogatz showed that real-world networks exhibit strong clustering or network transitivity [52]. A network is said to show "clustering" if the probability of two nodes being connected by a link is higher when the nodes in question have a common neighbor.

Newman demonstrated this by analyzing the time evolution of scientific collaborations and observing that the use of examining neighbors, in this case coauthors of authors, could help predict future collaborations [42]. From this analysis, Newman determined that the probability of two individuals collaborating increases as the number m of their previous mutual coauthors increases. A pair of scientists who have five mutual previous collaborators, for instance, is about twice as likely to collaborate as a pair with only two, and about 200 times as likely as a pair with none. Additionally, Newman determined that the probability of collaboration increases with the number of times one has collaborated before, which shows that past collaborations are a good indicator of future ones.

Liben-Nowell and Kleinberg explored this theory by the following common neighbor metric in order to predict future collaborations on an author database by assigning a score to the possible collaboration [33]. The score of a future collaboration score(x, y) between authors x and y, where

N(x) and N(y) are the set of neighbors of authors x and y, respectively, is calculated by

$$score(x, y) = |N(x) \cap N(y)|.$$
 (4)

Their results strongly supported this argument and showed that links were predicted by a factor of up to 47 improvement compared to that of random prediction. The *common neighbor* measure in (4) measures purely the similarity between two entities. Liben-Nowell also explored using Jaccard's coefficient, which attempts to take into account not just similarity but also dissimilarity. The Jaccard coefficient is defined as the size of the intersection divided by the size of the union of the sample sets:

$$score(x,y) = \frac{|N(x) \cap N(y)|}{|N(x) \cup N(y)|}.$$
(5)

Adamic and Adar performed an analysis to predict relationships between individuals by analyzing user homepages on the World Wide Web (WWW) [1]. The authors computed features of the pages and also took into account the incoming and outgoing links of the page and defined the similarity between two pages by counting the number of common features, assigning greater importance to rare features than frequently seen features. In the case of neighbors, Liben-Nowell utilized this metric, which refines a simple count of neighbors by weighting rarer neighbors more heavily than common neighbors. The probability, where N(z) is the number of neighbors held by z, is then given by

$$P(x,y) = \sum_{z \in N(x) \cap N(y)} \frac{1}{\log|N(z)|}.$$
 (6)

All three metrics performed well compared with random prediction. Liben-Nowell explored a number of different ranking techniques based on information retrieval research but generally found that common neighbor, the Jaccard's coefficient, and the Adamic and Adar technique were sufficient and performed equally as well, if not better than the other techniques.

Centrality and tie strength are based on the analysis of a static network graph whose link availability is time varying. However, in the case of DDTMs, the network graph is not static; it evolves over time. Tie predictors can be used in order to predict the evolution of the graph and evaluate the probability of future links occurring. Tie predictors may be used not only to reinforce already existing contacts but to anticipate contacts that may evolve over time.

## 4 ROUTING BASED ON SOCIAL METRICS

We propose that information flow in a network graph whose links are time varying can be achieved using a combination of centrality, strong ties, and tie prediction. Social networks may consist of a number of highly disconnected cliques where neither the source node nor any of its contacts has any direct relationship with the destination. In this case, relying on strong ties would prove futile, and therefore, weak ties to more connected nodes may be exploited.

Centrality has shown to be useful for path finding in static networks, however, the limitation of link capacities causes congestion. Additionally, centrality does not account for the time-varying nature of link availability. Tie strength may be used to overcome this problem by identifying links that have a higher probability of availability. Tie strength evaluates existing links in a time-varying network but does not account for the dynamic evolution of the network over time. Tie predictors may be used to aid in predicting future links that may arise. As such, we propose that the combination of centrality, tie strengths, and tie predictors are highly useful in routing based on local information when the underlying network exhibits a social structure. This combined metric will be referred to as the SimBetTS utility. When two nodes meet, they exchange a list of encountered nodes, this list is used to locally calculate the betweenness utility, the similarity utility, and the tie strength utility. Each node then examines the messages it is carrying and computes the SimBetTS utility of each message destination. Messages are then exchanged where the message is forwarded to the node holding the highest SimBetTS utility for the message destination node. The remainder of this section describes the calculation of the betweenness utility, the similarity utility, and the tie strength utility and how these metrics are combined to calculate the SimBetTS utility.

## 4.1 Betweenness Calculation

Betweenness centrality is calculated using an ego network representation of the nodes with which the ego node has come into contact. When two nodes meet, they exchange a list of contacts. A contact is defined as a node that has been directly encountered by the node. The received contact list is used to update each node's local ego network. Mathematically, node contacts can be represented by an adjacency matrix A, which is an  $n \times n$  symmetric matrix, where n is the number of contacts a given node has encountered (for a worked example, see [10]). The adjacency matrix has elements:

$$A_{i,j} = \begin{cases} 1, & \text{if there is a contact between } i \text{ and } j, \\ 0, & \text{otherwise.} \end{cases}$$
(7)

Contacts are considered to be bidirectional, so if a contact exists between *i* and *j*, then there is also a contact between *j* and *i*. The betweenness centrality is calculated by computing the number of nodes that are indirectly connected through the ego node. The betweenness centrality of the ego node is the sum of the reciprocals of the entries of A', where A' is equal to  $A^2[1 - A]_{i,j}$  [12] where *i*, *j* are the row and column matrix entries, respectively. A node's betweenness utility is given by

$$Bet = \sum \frac{1}{A'_{i,i}}.$$
(8)

Since the matrix is symmetric, only the nonzero entries above the diagonal need to be considered. When a new node is encountered, the new node sends a list of nodes it has encountered. The ego node makes a new entry in the  $n \times n$  matrix. As an ego network only considers the contacts between nodes that the ego has directly encountered, only the entries for contacts in common between the

ego node and the newly encountered node are inserted into the matrix.

#### 4.2 Similarity Calculation

Node similarity is calculated using the same  $n \times n$  matrix discussed in Section 4.1. The number of common neighbors between the current node i and destination node j can be calculated as the sum of the total overlapping contacts as represented in the  $n \times n$  matrix (for a worked example, see [10]). This only allows for the calculation of similarity for nodes that have been met directly, but during the exchange of the node's contact list, information can be obtained with regard to nodes that have yet to be encountered. As discussed in Section 3.3, the number of common neighbors may be used for ranking known contacts but also for predicting future contacts. Hence, a list of indirect encounters is maintained in a separate  $n \times m$  matrix, where n is the number of nodes that have been met directly and mis the number of nodes that have not directly been encountered but may be indirectly accessible through a direct contact. The similarity calculation, where  $N_n$  and  $N_e$ are the set of contacts held by node n and e, respectively, is given as follows:

$$Sim(n,e) = \left| N_n \bigcap N_e \right|. \tag{9}$$

## 4.3 Tie Strength Calculation

Measuring tie strength will be an aggregation of a selection of indicators based on those discussed in Section 3.2. An evidence-based strategy is used to evaluate whether each measure supports or contradicts the presence of a strong tie. The evidence is represented as a tuple (s, c), where s is the supporting evidence and c is the contradicting evidence. The trust in a piece of evidence is measured as a ratio of supporting and contradicting evidences [19]. Below elaborates on specific tie strength indicators, representing them as a ratio of supporting and contradicting evidences and bring them into the context of DDTMs.

**Frequency.** The frequency indicator may be based on the frequency with which a node is encountered. The supporting evidence of a strong tie strength is defined as the total number of times node n has encountered node m. The contradicting evidence is defined as the amount of encounters node n has observed, where node m was not the encountered node. The frequency indicator, where f(m) is the number of times node n encountered node m and F(n) is the total number of encounters node n has observed, is given by

$$FI_n(m) = \frac{f(m)}{F(n) - f(m)}.$$
 (10)

**Intimacy/closeness.** The duration indicator can be based on the amount of time the node has spent connected to a given node. The supporting evidence of intimacy/closeness is defined as a measure of how much time node n has been connected to node m. The contradicting evidence is defined as the total amount of time node n has been connected to nodes in the network, where node m was not the encountered node. If d(m) is the total amount of time node n has been connected to node m and D(n) is the total amount of time node n has been connected across all encountered nodes, then the intimacy/closeness indicator can be expressed as

$$ICI_n(m) = \frac{d(m)}{D(n) - d(m)}.$$
(11)

**Recency.** The recency indicator is based on how recently node n has encountered node m. According to the social network analysis theory, a strong tie must be maintained in order to stay strong, which is captured by the recency indicator. The supporting evidence rec(m) is how recently node n last encountered node m. This is defined as the length of time between node n encountered node m and the time node n has been on the network. The contradicting evidence is the amount of time node n has been on the network. The recency indicator, where L(n) is the total amount of time node n has been a part of the network, is given as

$$RecI_n(m) = \frac{rec(m)}{L(n) - rec(m)}.$$
(12)

Tie strength is a combination of a selection of indicators and the above metrics are all combined in order to evaluate an overall single tie strength measure. A strong tie should, ideally, be high in all measures. Consequently, the tie strength of node n for an encountered node m, where T is the set of social indicators, is given by

$$TieStrength_n(m) = \sum_{I \in T} I_n(m).$$
 (13)

## 4.4 Node Utility Calculation

The SimBetTS utility captures the overall improvement a node represents when compared to an encountered node across all measures presented in Sections 4.1, 4.2, and 4.3.

Selecting which node represents the best carrier for the message becomes a multiple attribute decision problem across all measures, where the aim is to select the node that provides the maximum utility for carrying the message. This is achieved using a pairwise comparison matrix on the normalized relative weights of the attributes.

The tie strength utility  $TSUtil_n$ , the betweenness utility  $BetUtil_n$ , and the similarity utility  $SimUtil_n$  of node n for delivering a message to destination node d compared to node m are given by

$$SimUtil_n(d) = \frac{Sim_n(d)}{Sim_n(d) + Sim_m(d)},$$
(14)

$$BetUtil_n = \frac{Bet_n}{Bet_n + Bet_m},\tag{15}$$

$$TSUtil_n(d) = \frac{TieStrength_n(d)}{TieStrength_n(d) + TieStrength_m(d)}.$$
 (16)

The  $SimBetTSUtil_n(d)$  is given by combining the normalized relative weights of the attributes, where  $U = \{TSUtil, SimUtil, BetUtil\}$ :

$$SimBetTSUtil_n(d) = \sum_{u \in U} u_n(d).$$
(17)

All utility values are considered of equal importance, and the SimBetTSUtil is the sum of the contributing utility values.

Replication may be used in order to increase the probability of message delivery. Messages are assigned a replication value R. When two nodes meet, if the replication value R > 1 then a message copy is made. The value of R is divided between the two copies of the message. This division is dependent on the SimBetTS utility value of each node; therefore, the division of the replication number for destination d between node n and node m is given by

$$R_{n} = \left\lfloor R_{cur} \times \frac{SimBetTSUtil_{n}(d)}{SimBetTSUtil_{n}(d) + SimBetTSUtil_{m}(d)} \right\rfloor,$$

$$R_{m} = R_{cur} - R_{n}.$$
(18)

Consequently, the node with the higher utility value receives a higher replication value. If R = 1, then the forwarding becomes a single-copy strategy. For evaluation purposes, a replication value of R = 4 is used. Replication improves the probability of message delivery, however comes at the cost of increased resource consumption. If replication is used, in order to avoid sending messages to nodes that are already carrying the message, the Summary Vector must include a list of message identifiers it is currently carrying for each destination node. However, a benefit of this replication methodology is that the message is only split if a carrying node encounters a node with a greater SimBetTSUtil value, as a result if the message reaches the most appropriate node before the replication has occurred, then the maximum number of replicas will not be used.

## **5** SIMULATION RESULTS

In this section, we describe the simulation setup used to evaluate SimBetTS Routing. The first experiment evaluates the utility of each metric discussed in Section 4 in terms of overall message delivery and the resulting distribution of the delivery load across the nodes on the network. The second experiment illustrates the routing protocol behavior and utility values that become important in the decision to transfer messages between nodes. The third experiment demonstrates the delivery performance of SimBetTS Routing compared to that of Epidemic Routing and PROPHET Routing.

#### 5.1 Simulation Setup

Cambridge University and Intel Research conducted three experiments of Bluetooth encounters using participants carrying iMote devices as part of the Haggle project [8], [24]. The experiments lasted between three and five days and only included a small number of participants. Encounters between Bluetooth devices were logged. Log entries include the encounter start time, the deviceID, and the deviceID of the encountering and encountered nodes. Logs are only available for the participating nodes, but the data set also includes encounters between participants carrying the iMote devices and external Bluetooth contacts. These include anyone who had an active Bluetooth device in the

	Intel	Cambridge	Infocom
Participants	8	12	41
Total Devices	128	223	264
Total Encounters	2,264	6,736	28,250
Average Encounter Duration	885 secs	572 secs	231 secs
Duration	3 days	5 days	3 days

TABLE 1 Haggle Data Set

vicinity of the iMote carriers. The trace file of these sightings was used in order to generate an event-based simulation where a Bluetooth sighting in the trace is assumed to be a contact where nodes can exchange information.

Table 1 summarizes the experiment details for each data set.

In order to evaluate the delivery performance of SimBetTS Routing, the event-based simulations are used to explore routing performance. In order to provide an opportunity to gather information about the nodes within the network, 15 percent of the simulation duration is used as a warm-up phase. After the warm-up phase, each node on the network generates messages to uniformly randomly selected destination nodes at a rate of 20 messages per hour from the start of the message sending phase when the sending node first reappears on the network until the last time the sending node appears on the network during the sending phase. The sending phase lasts for 70 percent of the simulation duration, allowing an additional 15 percent at the end in order to allow messages that have been generated time to be delivered. The simulations were run for all three protocols 10 times with different random number seeds. In this case, it is assumed that each Bluetooth encounter provides an opportunity to exchange all routing data and all messages.

#### 5.2 Evaluating Social Metrics for Information Flow

The aim of the first experiment is to evaluate the utility of each SimBetTS Routing utility component for routing and the benefit of combining these three metrics in order to improve delivery performance. We evaluate the delivery performance based on the following criteria.

**Total number of messages delivered.** This metric captures the total number of messages delivered by routing

based on the utility metrics of betweenness, similarity, tie strength, and finally, SimBetTS Routing, which combines all three utilities.

**Distribution of delivery load.** This metric is the percentage of messages delivered per node, which signifies the distribution of load across the network. If a high proportion of the messages is delivered by a small subset of nodes, then these nodes may become points of congestion. The aim of SimBetTS Routing is to maximize delivery performance while avoiding heavy congestion on central nodes in the network.

## 5.2.1 Intel Data Set

The first experiment included eight researchers and interns working at Intel Research, Cambridge and lasted three days. A total of 128 nodes, including the eight participants, were encountered for the duration of the experiment.

Fig. 2 shows the total number of messages delivered and the load distribution for the Intel Data Set. As can be seen from Fig. 2b when evaluating betweenness, similarity, and tie strength utility, betweenness achieves the highest delivery performance. Fig. 2a, however, shows the disadvantage of using betweenness utility in routing where over 50 percent of the total messages delivered in the network were delivered by a single node. Routing based on similarity and tie strength show a better distribution of load, but this comes at the price of reduced overall delivery performance. SimBetTS Routing achieves the highest delivery performance by combining the three metrics. Additionally, the load distribution shows that routing based on the combined metrics reduces congestion on highly central nodes.

## 5.2.2 Cambridge Data Set

The second data set included 12 doctoral students and faculty from Cambridge University Computer Laboratory [8]. The experiment lasted five days, and a total of 223 devices, including the participants, were encountered during the experiment.

Fig. 3 shows a similar trend in terms of delivery performance and distribution. SimBetTS Routing achieves the highest overall delivery performance. The load distribution of SimBetTS Routing when compared to betweenness utility is improved but to a lesser extent than the Intel data set. Similarity and tie strength both show similar



Fig. 2. Performance of protocol utility components: Intel data. (a) Distribution of message delivery. (b) Total message delivery.



Fig. 3. Performance of protocol utility components: Cambridge data. (a) Distribution of message delivery. (b) Total message delivery.

behavior where a single node delivers a high proportion of the messages. This illustrates that in this particular social network, a single node was highly important in terms of message delivery. A large proportion of the extra messages delivered by combining the utility metrics is delivered by the node that shows a sharp increase in Fig. 3. As a consequence, this result is not directly due to a reduction in load distribution, but rather that these messages were not delivered by using tie strength alone.

## 5.2.3 InfoCom Data Set

The third experiment collected data of encounters using iMotes equipped with Bluetooth distributed among conference attendees at the IEEE 2005 INFOCOM [24]. The participants were chosen in order to represent a range of different groups belonging to different organizations. Participants were asked to carry the devices with them for the duration of the conference. The experiment lasted three days and encounters between 264 devices were recorded.

Fig. 4 shows that routing based on betweenness utility results in a single node delivering as much as 65 percent of the total messages delivered. SimBetTS Routing reduces this load significantly, where less than 30 percent of the messages are delivered by this central node. As with the previous data sets, SimBetTS Routing achieves the highest overall delivery.

This experiment has demonstrated that SimBetTS Routing achieves superior delivery performance by combining the three utility metrics of betweenness, similarity, and tie strength when compared to routing based on the individual metrics alone. Additionally, we have shown that although betweenness centrality achieves high delivery performance, routing based on centrality alone results in significant load on a small subset of the nodes on the network. Combining the utility metrics shows a better load distribution engaging more nodes in message delivery.

## 5.3 Demonstrating Routing Protocol Behavior

The goal of this experiment is to illustrate the benefit of utilizing a multicriteria decision method for combining the values in order to identify the node that represents the best trade-off across these three metrics. The SimBetTS Routing protocol is based on the premise that when forwarding a message for a given destination, the message is forwarded to a node with a higher probability of encountering the destination node. A node with a high tie strength has a higher probability of encountering the destination node. If a node with a high tie strength cannot be found, the utility metrics of social similarity and betweenness centrality are used. To demonstrate that the SimBetTS Routing protocol achieves this behavior, we divide the paths taken by the message from source to destination into three categories:

- Similarity + tie strength. The sending node has a nonzero tie strength value for the destination node and has a nonzero similarity to the destination node.
   Similarity: The sending node has never encountered
- the destination node resulting in a zero tie strength



Fig. 4. Performance of protocol utility components: InfoCom data. (a) Distribution of message delivery. (b) Total message delivery.



Fig. 5. Utility values at first hop and last hop along the message path divided into the categories of similarity + tie strength, similarity, and none. (a) Intel data set. (b) Cambridge data set. (c) InfoCom data set.

value; however, the sending node has a nonzero similarity value.

 None. The sending node has never encountered the destination node and has no common neighbors, hence a zero similarity value.

Fig. 5 shows the utility values of the first hop node on the message delivery path and the utility values of the final delivery node along the message path for each of these categories. It is important to note that these values are the relative utility values of the node selected to forward the message compared to the currently carrying node. As a result a high utility value at the first hop and a lower utility value at the final hop do not mean a lower absolute value, but lower when compared to the currently carrying node.

### 5.3.1 Similarity and Tie Strength

In the category of Similarity and Tie strength, a combination of tie strength and betweenness utilities is the deciding factor in forwarding to the next hop. Similarity has a lower impact, most likely due to the fact that if the sending node has a social similarity and a tie strength value for the destination node, then the nodes encountered by the sending node most likely also have a social similarity. In all data sets, it can be seen that the tie strength utility is the highest contributing utility when forwarding to the last hop delivering node.

## 5.3.2 Similarity

In the category of Similarity where the sending node has a nonzero similarity value to the destination node, the highest contributing utility value is consistently the tie strength utility in both the first hop and the last hop delivering node. It can also be seen that the betweenness utility is also a contributing utility on the first hop; however, for the final delivering node betweenness utility is much reduced. This is the desired behavior of the protocol, because betweenness utility should only have a high impact in finding a node with a high tie strength to the destination node and then should have a lower impact on the forwarding decision.

#### 5.3.3 None

In the category where the sending node has no social similarity to the destination node and has also never encountered the destination node, we can see the benefit of the routing protocol. In all data sets, a combination of betweenness and similarity is the most contributing utility value. It was found that, in this category, tie strength has a zero utility value in all cases at the first hop; however, it can be consistently seen that the final delivering node had a high tie strength utility. As a result, we can conclude that the routing protocol functions as intended. Whenever the destination node is unknown, a combination of betweenness utility and similarity utility will navigate the message to a node in the network that has a higher tie strength for the destination node.

This experiment has demonstrated the navigation behavior of the SimBetTS Routing protocol. The message is forwarded to nodes with a high betweenness and social similarity, until a node with a high tie strength for the destination node is found. In all cases, the tie strength utility for the final hop is the highest contributing utility value, and in all cases, the betweenness utility value is much reduced in its influence of the forwarding decision as the message is routed closer to the destination.

#### 5.4 Delivery Performance of Combined Metrics

The goal of the third experiment is to evaluate the delivery performance of SimBetTS Routing protocol compared to two protocols: Epidemic Routing [51] and PRoPHET Routing [36]. The default parameters for PRoPHET Routing were used as defined in [36].

Two versions of SimBetTS Routing and PROPHET are evaluated: a single-copy version and a multicopy version. In the single-copy strategy, when two nodes meet, messages are exchanged between nodes where messages are forwarded to the node with the highest utility. The node that



Fig. 6. Protocol performance for Intel data. (a) Delivery performance. (b) Average number of hops. (c) Average end-to-end delivery delay. (d) Total number of forwards.

has forwarded the message must then delete the message from the message queue. In the multicopy strategy, replication is used where messages are assigned a replication value R. For evaluation purposes, a replication value of R = 4 is used.

**Total number of messages delivered.** The ultimate goal of the SimBetTS Routing design is to achieve delivery performance as close to Epidemic Routing as possible. This is because Epidemic Routing always finds the best possible path to the destination and therefore represents the baseline for the best possible delivery performance.

Average end-to-end delay. End-to-End delay is an important concern in SimBetTS Routing design. Long endto-end delays mean that the message must occupy valuable buffer space for longer, and consequently, a low end-to-end delay is desirable. Again, Epidemic Routing presents a good baseline for the minimum end-to-end delay possible.

Average number of hops per message. It is desirable to minimize the number of hops a message must take in order to reach the destination. Wireless communication is costly in terms of battery power, and as a result, minimizing the number of hops also minimizes the battery power expended in forwarding the message.

**Total number of forwards.** This value represents the network overhead in terms of how many times a message forward occurs. PRoPHET and SimBetTS are expected to perform similarly in this respect, as both only assume the existence of one copy of the message on the network.

Epidemic Routing, however, assumes the existence of multiple copies and continues forwarding a given message until each node is carrying a copy. This means Epidemic Routing is costly in terms of the number of transmissions required along with the amount of buffer space required on each node.

## 5.4.1 Intel Data Set

Fig. 6 graphs the results for this simulation. As shown in Fig. 6a, it is clear that Epidemic Routing outperforms both SimBetTS Routing and PRoPHET. Single-copy SimBetTS Routing delivers fewer messages than Epidemic but shows improvement when compared to single-copy PRoPHET. Multicopy SimBetTS Routing shows significant improvement with the addition of replication resulting in an improvement of nearly 50 percent than the single-copy strategy. Multicopy PRoPHET shows much less improvement with the addition of replication, achieving results comparable to single-copy SimBetTS Routing. All protocols show a number of plateaus where no messages are delivered before the 150 mark and at the 200 mark. These plateaus represent nighttime when the devices are not within the range of other devices.

Fig. 6b shows the average number of hops taken by the delivered messages. All protocols result in a small number of hops of around 3. Single-copy and multicopy PRoPHET show the largest and smallest number of average hops, respectively. Single-copy and multicopy SimBetTS Routing



Fig. 7. Control data overhead for Intel data.

show very similar average hop values, meaning the path lengths found by single-copy SimBetTS Routing were close to that achieved by multicopy SimBetTS. Epidemic achieves short paths from the source to destination of just below an average of 3. SimBetTS Routing shows a distinct increase in the average number of hops after the 300 mark, which coincides with a large number of messages being delivered. This indicates that these messages required a larger number of hops in order to be delivered, thus raising the average.

Fig. 6c shows the average message end-to-end delay. As expected, Epidemic Routing shows the lowest average end-to-end delay. Single-copy PROPHET shows the highest overall delay. Single-copy SimBetTS Routing shows similar delays. Both multicopy SimBetTS Routing and PROPHET show reduced delays with multicopy SimBetTS Routing resulting in low delays near those achieved by Epidemic Routing. All protocols show a significant increase in delay around the 250 mark, which coincides with a steep increase in the number of messages delivered, meaning these delivered messages required a longer amount of time to be delivered, thus increasing the average end-to-end delay.

The true disadvantage of Epidemic Routing becomes clear when examining the total number of forwards. The overhead associated with Epidemic Routing is so great that it has been omitted from Fig. 6d in order for the differentiation of SimBetTS and PRoPHET to be seen. Epidemic Routing continues to forward messages throughout the network and, as a result, incurs a great deal of overhead. Single-copy SimBetTS and PRoPHET only have a single copy of each message on the network, resulting in a significantly lower number of forwards. Multicopy SimBetTS Routing and PRoPHET as expected show an increase in the number of forwards with the use of replication. However, when compared to that of Epidemic Routing, multicopy SimBetTS Routing results in approximately 98 percent reduction in number of forwards. SimBetTS results in a higher number of forwards when compared to PRoPHET; however, this result should be viewed in the context that PRoPHET delivered significantly less messages overall.

Fig.7a shows the protocol control data overhead. Epidemic Routing overhead is so large it makes it difficult to differentiate between the overhead generated by single-copy SimBetTS Routing and PROPHET. This is because the summary vector must contain a list of all messages the node is currently carrying, and as the number of messages on the network increases, so does the control data. Single-copy SimBetTS Routing and PRoPHET generate significantly lower overhead due to the fact that nodes exchange information about message destination rather than explicit message identifiers. As a result, there is an upper limit on the amount of bytes required. This upper bound depends on the node population. Multicopy SimBetTS Routing and PRoPHET show an increase in control data due to the necessary exchange of message identifiers, thus reducing the benefit of exchanging only routing data, as is the case for the singlecopy strategy. However, the overhead is still significantly lower than that of Epidemic Routing due to the fact that nodes do not carry a copy of every message in the network.

## 5.4.2 Cambridge Data Set

Fig. 8 graphs the results from the Cambridge data set performance tests. Plateaus are again evident where message delivery remains level. There are three such plateaus around the 200, 300, and 400 marks that are again nighttime where devices are inactive. The first night occurred before the message sending phase commenced, and the 200, 300, and 400 marks represent the second, third, and fourth nights of the experiment, respectively. As expected, Epidemic Routing shows superior delivery performance compared to that of PRoPHET and SimBetTS Routing as shown in Fig. 8a. Single-copy SimBetTS Routing outperforms both single-copy and multicopy PRoPHET. Single-copy SimBetTS Routing achieves improved delivery performance of 100 percent when compared to single-copy PRoPHET. More dramatically, multicopy SimBetTS Routing achieved improved delivery performance of 300 percent when compared to that of multicopy PRoPHET.

Fig. 8b shows the average number of hops. As with the previous experiment, all protocols achieve message delivery in a relatively small number of hops of around 3-4. Epidemic Routing results in the highest average. It could be assumed that this is due to the higher number of messages delivered by Epidemic Routing resulting in a number of messages that required a larger number of hops in order to reach the destination. However, upon inspection of the message delivery graph, the average number of hops remains level for Epidemic Routing even after a large proportion of the messages is delivered. This can be explained by the fact that Epidemic Routing finds the optimum path in terms of delivery delay rather than the shortest path. Multicopy PRoPHET results in the least number of hops. Interestingly, multicopy SimBetTS Routing results in a higher average when compared to single-copy SimBetTS Routing, but this increase around the 250 mark coincides with a sharp increase in message delivery. After this time, the average number of hops starts to decrease.

Fig. 8c shows the average message delay. Epidemic Routing shows the highest overall delay, but it can be noted that the delay increases most sharply around the 350 mark. Upon inspection of the delivery graph (Fig. 8a), this increase can be explained by the increase in total messages delivered, representing the fact that the messages delivered during this time phase were unable to be delivered in a shorter time. Approximately 30 percent of the encountered nodes do not appear in the network until after the 300 mark. PROPHET shows the lowest overall delivery delay due to



Fig. 8. Protocol performance for Cambridge data. (a) Delivery performance. (b) Average number of hops. (c) Average end-to-end delivery delay. (d) Total number of forwards.

the fact that it delivered a smaller proportion of the total messages. SimBetTS Routing results in an average end-toend delay between that of PRoPHET and Epidemic Routing. Similar to the average hop count, multicopy SimBetTS Routing shows an increase in average end-to-end delay when compared to single-copy SimBetTS Routing. The increase in delay is clearly related to the increase in message delivery and follows a trend almost identical to that of Epidemic Routing.

Fig. 8d shows the total number of message forwards throughout the simulation. As with the previous data set, Epidemic Routing results in a large number of forwards, which is to be expected with a flooding protocol. Singlecopy SimBetTS Routing and PROPHET result in a relatively low number of forwards. This value increases for multicopy SimBetTS Routing and PROPHET. Unlike the previous data set, multicopy PROPHET results in a higher number of forwards than that of multicopy SimBetTS Routing. As a result, it can be determined that this value is dependent on the dynamics of the underlying network rather than a deterministic side effect of the protocol.

Fig. 9 shows the total overhead associated with the control data for each protocol. Epidemic Routing as expected increases dramatically as the number of messages on the network increases. Both single-copy SimBetTS Routing and PRoPHET protocols increase at approximately the same rate; however, single-copy SimBetTS Routing results in the lowest number of bytes. Multicopy SimBetTS Routing results in a larger amount of control data when

compared to multicopy PRoPHET; however, both protocols follow a similar trend.

## 5.4.3 InfoCom Data Set

The overall delivery performance in Fig. 10a shows that Epidemic Routing achieves significant improvement after the 150 mark. This can be explained by the fact that approximately 34 percent of the node population first appear after this time frame. At this point, SimBetTS Routing and PROPHET start to build up encounter information with regard to these nodes. This explains why SimBetTS Routing starts to outperform PROPHET after the 150 mark, because messages destined for nodes that



Fig. 9. Control data overhead for Cambridge data.



Fig. 10. Protocol performance for InfoCom data. (a) Delivery performance. (b) Average number of hops. (c) Average end-to-end delivery delay. (d) Total number of forwards.

have yet to be encountered are already routed to more central nodes, which are more likely to gather encounter information about the previously unseen nodes. Multicopy SimBetTS Routing achieves message delivery close to Epidemic and outperforms PRoPHET.

Fig. 10b shows the average number of hops of delivered messages. Epidemic Routing results in hop lengths of approximately 4. Single-copy SimBetTS Routing and PRoPHET achieve similar results of between five and six hops. Multicopy PRoPHET results in the lowest number of hops. In contrast, multicopy SimBetTS Routing shows significant increase, resulting in hops of approximately 9. This increase is due to an increased path length of the additional messages delivered by multicopy SimBetTS Routing that remained undelivered for singlecopy SimBetTS Routing. Fig. 10c shows the average endto-end delay. Even with the increased path lengths used by multicopy SimBetTS Routing, the average end-to-end delay is similar to that of Epidemic Routing. The sharpest increase occurs for Epidemic after the 150 mark, which is most likely when messages are delivered to nodes previously unseen on the network. SimBetTS Routing shows a similar increase. PRoPHET shows the lowest average delay. However, the fact that it shows delays lower than that of Epidemic Routing illustrates that the increase in message delay shown by Epidemic is caused by the additional messages delivered, which required a greater amount of time before it was possible to deliver these messages.

Fig. 10c shows the average message delay. Similar to the Intel and Cambridge data sets, Epidemic Routing shows the greatest average delay due the higher proportion of messages delivered. The sharpest increase occurs after for Epidemic after the 150 mark, which is most likely due to message delivery of messages destined for the previously unseen on the network. SimBetTS Routing shows a similar increase just before the 250 mark, which coincides with an increase in message delivery. Consequently, this increase also coincides with the message delivery to nodes that have yet to be encountered.

Fig. 10d shows the total number of messages forwarded in the network. As expected, Epidemic Routing results in the highest number of total forwards and has been omitted. SimBetTS Routing and PROPHET result in a similar number of forwards when comparing the single-copy strategy; however, multicopy PROPHET results in a higher number of forwards than multicopy SimBetTS.

Fig. 11 shows the total control data. The overhead of Epidemic Routing increases dramatically as the number of messages on the network also increases. As seen with the other data sets, single-copy SimBetTS Routing and PRoPHET result in a similar amount of overhead. Multicopy SimBetTS Routing results in a larger amount of control data when compared to multicopy PRoPHET. However, this result should be viewed in the context that SimBetTS results in an improvement of message delivery of nearly 50 percent compared to multicopy PRoPHET.



Fig. 11. Baseline control data overhead for InfoCom data.

From this experiment, we can determine that single-copy and multicopy SimBetTS Routing show higher message delivery when compared to that of PRoPHET. Multicopy SimBetTS Routing achieves delivery performance similar to that of Epidemic Routing with short path lengths and low end-to-end delay. The use of replication comes at the cost of an increased number of forwards and an increase in control data. However, when compared to that of Epidemic Routing, these metrics are still relatively low in terms of overhead.

## 6 CONCLUSIONS

We have presented social network analysis techniques and shown how they may be applied to routing in DDTMs. Three metrics have been defined: betweenness utility, similarity utility, and tie strength utility. Each of these metrics has been evaluated individually, and the results show that betweenness utility results in the best overall delivery performance. However, it results in congestion on highly central nodes. SimBetTS Routing combines these three metrics, which results in improved overall delivery performance with the additional advantage that the load on central nodes is reduced and better distributed across the network.

We have evaluated SimBetTS Routing compared to two other protocols, Epidemic Routing and PRoPHET. Sim-BetTS Routing achieves delivery performance comparable to Epidemic Routing, without the additional overhead. We have also demonstrated that SimBetTS Routing outperforms the PRoPHET routing protocol in terms of overall delivery performance. The evaluation consisted of three separate real-world trace experiments. Although the first two data sets are relatively small networks, the consistency of the results across all three scenarios shows that, given the existence of an underlying social structure, SimBetTS Routing provides message delivery with greatly reduced overhead when compared to Epidemic Routing. Similar results can be seen when analyzing the larger MIT Reality Mining data set covered in our previous work [10]. Additionally, these utility metrics may be applied in other distributed systems where global topology information is unavailable, especially where the underlying networks exhibit small-world characteristics.

## REFERENCES

- L.A. Adamic and E. Adar, "Friends and Neighbors on the Web," Social Networks, vol. 25, no. 3, pp. 211-230, July 2003.
- [2] A. Beaufour, M. Leopold, and P. Bonnet, "Smart-Tag Based Data Dissemination," Proc. First ACM Int'l Workshop Wireless Sensor Networks and Applications (WSNA '02), pp. 68-77, Sept. 2002.
- [3] M. Benassi, A. Greve, and J. Harkola, "Looking for a Network Organization: The Case of GESTO," J. Market-Focused Management, vol. 4, no. 3, pp. 205-229, Oct. 1999.
- [4] P. Blumstein and P. Kollock, "Personal Relationships," Ann. Rev. Sociology, vol. 14, pp. 467-490, 1988.
- [5] S.P. Borgatti, "Centrality and Network Flow," Social Networks, vol. 27, no. 1, pp. 55-71, Jan. 2005.
- [6] J.J. Brown and P.H. Reingen, "Social Ties and Word-of-Mouth Referral Behavior," *The J. Consumer Research*, vol. 14, no. 3, pp. 350-362, Dec. 1987.
- [7] J. Burgess, B. Gallagher, D. Jensen, and B.N. Levine, "Maxprop: Routing for Vehicle-Based Disruption-Tolerant Networking," *Proc. IEEE INFOCOM*, Mar. 2006.
- [8] A. Chaintreau, P. Hui, J. Crowcroft, C. Diot, R. Gass, and J. Scott, "Impact of Human Mobility on the Design of Opportunistic Forwarding Algorithms," *Proc. IEEE INFOCOM*, 2006.
- [9] S. Corson and J. Macker, Mobile Ad Hoc Networking (MANET): Routing Protocol Performance Issues and Evaluation Considerations, IETF RFC 2501, Jan. 1999.
- [10] E.M. Daly and M. Haahr, "Social Network Analysis for Routing in Disconnected Delay-Tolerant MANETs," Proc. ACM MobiHoc, Sept. 2007.
- [11] H. Dubois-Ferriere, M. Grossglauser, and M. Vetterli, "Age Matters: Efficient Route Discovery in Mobile Ad Hoc Networks Using Encounter Ages," *Proc. ACM MobiHoc*, pp. 257-266, 2003.
- [12] M. Everett and S.P. Borgatti, "Ego Network Betweenness," Social Networks, vol. 27, no. 1, pp. 31-38, Jan. 2005.
- [13] L.C. Freeman, "A Set of Measures of Centrality Based on Betweenness," Sociometry, vol. 40, no. 1, pp. 35-41, Mar. 1977.
- [14] L.C. Freeman, "Centrality in Social Networks Conceptual Clarification," Social Networks, vol. 1, no. 3, pp. 215-239, 1978-1979.
- [15] R.H. Frenkiel, B.R. Badrinath, J. Borres, and R.D. Yates, "The Infostations Challenge: Balancing Cost and Ubiquity in Delivering Wireless Data," *IEEE Personal Comm.*, vol. 7, no. 2, pp. 66-71, 2000.
- [16] J. Ghosh, H.Q. Ngo, and C. Qiao, "Mobility Profile Based Routing within Intermittently Connected Mobile Ad Hoc Networks (ICMAN)," Proc. Int'l Conf. Wireless Comm. and Mobile Computing (IWCMC '06), pp. 551-556, 2006.
- [17] N. Glance, D. Snowdon, and J.-L. Meunier, "Pollen: Using People as a Communication Medium," *Computer Networks*, vol. 35, no. 4, pp. 429-442, Mar. 2001.
- [18] M.S. Granovetter, "The Strength of Weak Ties," The Am. J. Sociology, vol. 78, no. 6, pp. 1360-1380, May 1973.
- [19] E.L. Gray, "A Trust-Based Reputation Management System," PhD thesis, Apr. 2006.
- [20] M. Grossglauser and M. Vetterli, "Locating Nodes with Ease: Last Encounter Routing in Ad Hoc Networks through Mobility Diffusion," *Proc. IEEE INFOCOM*, vol. 3, pp. 1954-1964, 2003.
- [21] R. Handorean, C. Gill, and G.-C. Roman, "Accommodating Transient Connectivity in Ad Hoc and Mobile Settings," *Proc. Second Int'l Conf. Pervasive Computing (PERVASIVE '04)*, A. Ferscha and F. Mattern, eds., pp. 305-322, Mar. 2004.
- [22] W. Hsu and A. Helmy, "On Nodal Encounter Patterns in Wireless LAN Traces," Proc. Fourth Int'l Symp. Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOpt '06), pp. 1-10, Apr. 2006.
- [23] W.-J. Hsu and A. Helmy, "Impact: Investigation of Mobile-User Patterns across University Campuses Using WLAN Trace Analysis," Proc. IEEE INFOCOM, Aug. 2005.
- [24] P. Hui, A. Chaintreau, J. Scott, R. Gass, J. Crowcroft, and C. Diot, "Pocket Switched Networks and Human Mobility in Conference Environments," *Proc. ACM SIGCOMM Workshop Delay-Tolerant Networking (WDTN '05)*, pp. 244-251, Aug. 2005.
- [25] P. Hui and J. Crowcroft, "How Small Labels Create Big Improvements," Proc. Fifth Ann. IEEE Int'l Conf. Pervasive Computing and Comm. Workshops (PerCom Workshops '07), pp. 65-70, 2007.
- [26] S. Jain, K. Fall, and R. Patra, "Routing in a Delay Tolerant Network," SIGCOMM Computer Comm. Rev., vol. 34, no. 4, pp. 145-158, Oct. 2004.

- [27] D.B. Johnson and D.A. Maltz, *Dynamic Source Routing in Ad-Hoc Wireless Networks*, vol. 353, pp. 153-181. Kluwer Academic Publishers, 1996.
- [28] A. Khelil, P.J. Marron, and K. Rothermel, "Contact-Based Mobility Metrics for Delay-Tolerant Ad Hoc Networking," *Proc.* 13th IEEE Int'l Symp. Modeling, Analysis, and Simulation of Computer and Telecomm. Systems (MASCOTS '05) pp. 435-444, 2005.
- [29] Y.-B. Ko and N.H. Vaidya, "Location-Aided Routing (LAR) in Mobile Ad Hoc Networks," Wireless Networks, vol. 6, no. 4, pp. 307-321, Sept. 2000.
- [30] J. Lebrun, C.-N. Chuah, D. Ghosal, and M. Zhang, "Knowledge-Based Opportunistic Forwarding in Vehicular Wireless Ad Hoc Networks," *Proc. IEEE 61st Conf. Vehicular Technology (VTC-Spring* '05), vol. 4, pp. 2289-2293, May 2005.
- [31] J. Leguay, T. Friedman, and V. Conan, "Evaluating Mobility Pattern Space Routing for DTNs," *Proc. IEEE INFOCOM*, Apr. 2006.
- [32] Q. Li and D. Rus, "Sending Messages to Mobile Users in Disconnected Ad-Hoc Wireless Networks," Proc. ACM MobiCom, pp. 44-55, Aug. 2000.
- [33] D. Liben-Nowell and J. Kleinberg, "The Link Prediction Problem for Social Networks," *Proc. 12th Int'l Conf. Information and Knowledge Management (CIKM '03)*, pp. 556-559, Nov. 2003.
  [34] N. Lin, P. Dayton, and P. Reenwald, "Analyzing the Instrumental
- [34] N. Lin, P. Dayton, and P. Reenwald, "Analyzing the Instrumental Use of Relations in the Context of Social Structure," *Sociological Methods and Research*, vol. 7, no. 2, pp. 149-166, 1978.
- [35] N. Lin, J.C. Vaughn, and W.M. Ensel, "Social Resources and Occupational Status Attainment," *Social Forces*, vol. 59, no. 4, pp. 1163-1181, June 1981.
- [36] A. Lindgren, A. Doria, and O. Schelén, "Probabilistic Routing in Intermittently Connected Networks," *Proc. First Int'l Workshop Service Assurance with Partial and Intermittent Resources (SAPIR '04)*, pp. 239-254, Aug. 2004.
- [37] P.V. Marsden, "Egocentric and Sociocentric Measures of Network Centrality," Social Networks, vol. 24, no. 4, pp. 407-422, Oct. 2002.
- [38] P.V. Marsden and K.E. Campbell, "Measuring Tie Strength," Social Forces, vol. 63, no. 2, pp. 482-501, Dec. 1984.
- [39] S. Merugu, M. Ammar, and E. Zegura, "Routing in Space and Time in Networks with Predictable Mobility," technical report, July 2004.
- [40] S. Milgram, "The Small World Problem," Psychology Today, vol. 2, pp. 60-67, May 1967.
- [41] M. Musolesi, S. Hailes, and C. Mascolo, "Adaptive Routing for Intermittently Connected Mobile Ad Hoc Networks," *Proc. Sixth IEEE Int'l Symp. World of Wireless Mobile and Multimedia Networks* (WOWMOM '05), pp. 183-189, 2005.
- [42] M.E.J. Newman, "Clustering and Preferential Attachment in Growing Networks," *Physical Rev. E*, vol. 64, no. 2, July 2001.
- [43] M.E.J. Newman, "A Measure of Betweenness Centrality Based on Random Walks," *Social Networks*, vol. 27, no. 1, pp. 39-54, Jan. 2005.
- [44] S.-Y. Ni, Y.-C. Tseng, Y.-S. Chen, and J.-P. Sheu, "The Broadcast Storm Problem in a Mobile Ad Hoc Network," *Proc. ACM MobiCom*, pp. 151-162, Aug. 1999.
- [45] C.E. Perkins and E.M. Royer, "Ad-Hoc On-Demand Distance Vector Routing," Proc. Second IEEE Workshop Mobile Computing Systems and Applications (WMCSA '99), pp. 90-100, Feb. 1999.
- [46] C.E. Perkins and P. Bhagwat, "Highly Dynamic Destination-Sequenced Distance-Vector Routing (DSDV) for Mobile Computers," Proc. ACM SIGCOMM, vol. 24, pp. 234-244, Oct. 1994.
- [47] R.C. Shah, S. Roy, S. Jain, and W. Brunette, "Data Mules: Modeling a Three-Tier Architecture for Sparse Sensor Networks," Proc. First IEEE Int'l Workshop Sensor Network Protocols and Applications (SNPA '03), pp. 30-41, May 2003.
- [48] T. Spyropoulos, K. Psounis, and C.S. Raghavendra, "Single-Copy Routing in Intermittently Connected Mobile Networks," Proc. First Ann. IEEE Comm. Soc. Conf. Sensor and Ad Hoc Comm. and Networks (SECON '04), pp. 235-244, Oct. 2004.
- [49] T. Spyropoulos, K. Psounis, and C.S. Raghavendra, "Spray and Wait: An Efficient Routing Scheme for Intermittently Connected Mobile Networks," *Proc. ACM SIGCOMM Workshop Delay-Tolerant Networking (WDTN '05)*, pp. 252-259, Aug. 2005.
- [50] K. Tan, Q. Zhang, and W. Zhu, "Shortest Path Routing in Partially Connected Ad Hoc Networks," Proc. IEEE Global Telecomm. Conf. (GLOBECOM '03), vol. 2, pp. 1038-1042, Dec. 2003.
- [51] A. Vahdat and D. Becker, "Epidemic Routing for Partially Connected Ad Hoc Networks," technical report, Apr. 2000.

- [52] D.J. Watts and S.H. Strogatz, "Collective Dynamics of "Small-World" Networks," *Nature*, vol. 393, no. 6684, pp. 440-442, June 1998.
- [53] G. Yan, T. Zhou, B. Hu, Z.Q. Fu, and B.H. Wang, "Efficient Routing on Complex Networks," *Physical Rev. E (Statistical, Nonlinear, and Soft Matter Physics)*, vol. 73, no. 4, Apr. 2006.
  - [4] W. Zhao, M. Ammar, and E. Zegura, "A Message Ferrying Approach for Data Delivery in Sparse Mobile Ad Hoc Networks," *Proc. ACM MobiHoc*, pp. 187-198, May 2004.



Elizabeth M. Daly received the BABAI degree in computer engineering and the MSc degree in networks and distributed systems from Trinity College Dublin in 2001 and 2002, respectively, and the PhD degree in 2007. She is currently with the Lotus Connections Team, IBM Dublin Software Laboratory and is also a part-time lecturer at Trinity College Dublin. Her research interests include social networks, peer-to-peer, distributed search, and trust.



**Mads Haahr** received the BSc and MSc degrees from the University of Copenhagen in 1996 and 1999, respectively, and the PhD degree from Trinity College Dublin in 2004. He is currently a lecturer at Trinity College Dublin. He is the editor-in-chief of *Crossings: Electronic Journal of Art and Technology* and also built and operates RANDOM.ORG. His current research interests are in large-scale self-organizing distributed and mobile systems and in sensor-augmented arte-

facts. He is a member of the IEEE.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.