

The evolution of overlapping communities in dynamic mobile networks

Nam P. Nguyen, Yilin Shen, Thang N. Dinh, *Student member, IEEE* and My T. Thai, *Member, IEEE*

Abstract—Many practical problems on Mobile networks, such as routing strategies in MANETs, sensor reprogramming in WSNs and worm containment in online social networks (OSNs) share an ubiquitous, yet interesting feature in their organizations: community structure. Knowledge of this structure provides us not only crucial information about the network principles, but also key insights into designing more effective algorithms for practical problems enabled by Mobile networking. However, understanding this interesting feature is extremely challenging on dynamic networks where changes to their topologies are frequently introduced, and especially when network communities in reality usually overlap with each other.

We focus on the following questions (1) Can we effectively detect the overlapping community structure in a *dynamic network*? (2) Can we quickly and adaptively update the network structure only based on its history without recomputing from scratch? (3) How does the detection of network communities help mobile applications? We propose *AFOCS*, a two-phase framework for not only detecting quickly but also tracing effectively the evolution of overlapped network communities in dynamic mobile networks. With the great advantages of the overlapping community structure, *AFOCS* significantly helps in reducing up to 7 times the infection rates in worm containment on OSNs, and up to 11 times overhead while maintaining good delivery time and ratio in forwarding strategies in MANETs.

Index Terms—Overlapping community structure, Dynamic network, Adaptive algorithm, Worm containment, Routing in MANETs.



1 INTRODUCTION

The rapid and exceptional growing of mobile network has called for a deeper understanding of its organization principles, in order to develop better techniques for a wide range of problems enabled by mobile networking. Many practical problems, such as forwarding and routing strategies in MANETs [1], sensor reprogramming in WSNs [2] and worm containment in cellular networks [3], [4] appear to share an ubiquitous and interesting property: the property of containing community structure, i.e., there are groups of devices or people that frequently communicate more with each other than with the others in the underlying organizations.

In a general concept, a community is a group of tight-knit nodes having more internal than external connections [5], [6]. For instance, a community in MANETs often comprises of sensors or mobile devices that are frequently transmitting data to each other than to other devices. Similarly, since people have a natural tendency to form groups of communication, a community in a cellular network usually consists of mobile devices that often call or text each other. The detection of network community structure, as a result, provides us a better knowledge about its characteristics as well as its organization principles, thereby providing more efficient solutions for mobile networking problems such as forwarding in MANETs or worm containment in OSNs.

Particularly, how does community structure help in Mobile applications? Indeed, detecting network communities is of considerable advantage in mobile networks. Let us consider the worm containment in cellular networks [3], or in OSNs [4], [7]. Nowadays, many social applications such as Facebook, Twitter and FourSquare, are able to run on open-API enabled mobile devices like PDAs and Iphones. However, if such an application is infected with malicious software, such as worms or viruses, this openness will also make it easier for their propagation. A possible solution to prevent worms from spreading out wider is to send patches to critical users and let them redistribute to the others. Intuitively, the smaller the set of important users for sending patches, the better. But how can we effectively choose that set of minimal size? This is where community structure comes into the picture and helps. In particular, we show that selecting users in the boundaries of the overlapped nodes gives a tighter and more efficient set of influential users, thus significantly lowers the number of sent patches as well as overhead information, which are essential in cellular networks and OSNs.

Another great advantage of community structure can be found in the forwarding problem in communication networks. Due to the high mobility of devices, an apparent challenge for a forwarding method is to quickly forward the message from the source to the destination, without introducing too many duplicate messages or overhead information. Since people tend to form groups of communication, there are also communities of tightly connected devices in the underlying network as a refec-tion. A good forwarding strategy, as soon as it discovers the network structure, can actively forward messages to

• Nam P. Nguyen, Yilin Shen, Thang N. Dinh and My T. Thai are with the Department of Computer and Information Science and Engineering, University of Florida, Gainesville, FL.
E-mail: {nanguyen, yshen, tdinh, mythai}@cise.ufl.edu

devices sharing more common community labels with the destination, rather than simply sending messages to those in the same community as the destination. With the helpful knowledge of network communities, this strategy will considerably reduce the number of duplicate messages while maintaining good delivery ratios, as we shall see in Section 8. This example, again, amplifies the importance of an efficient community detection method in mobile networks.

Mobile networks in reality are highly dynamic and thus, their communities are not always disjoint from each other. Indeed, their communities often overlap with each other since some active devices can participate in multiple groups at the same time, thereby reassemble the concept of overlapping community structure. Furthermore, most practical models for mobile network problems evolve frequently over time due to the high mobility of participating devices. Although any slight change does not seem to have a significant effect on the network structure, the evolution of the mobile network over a long duration might lead to an unpredictable transformation of its communities, particularly when they can overlap. This drives a crucial need of reidentification. However, the rapid changing network topology makes this an extremely challenging problem, especially on dynamic mobile networks.

A naive solution to the above problem would try to repeatedly execute one of the available static methods [8], [9], [10] to find new communities whenever the network changes; doing so, nonetheless, suffers from some major disadvantages (1) the huge consumption of time and computing resources on large networks and (2) the almost same reactions to some local parts of the network. Intuitively, a much better approach should adaptively update the current community structure based on its history and the network changes only, thus can eventually avoid the hassle of redetection.

Motivated by this intuition and the applicability of overlapping community structure, we propose *AFOCS* (Adaptive *FOCS*), an adaptive framework for detecting, updating and tracing the evolution of overlapping communities in dynamic mobile networks. Our two-phase framework first identifies all possible basic network communities with *FOCS* (short for Finding Overlapping Community Structure), and then employs *AFOCS* to adaptively update these structures as the network evolves. Since only *AFOCS* will stay up and handle all changes introduced to the network, this adaptive phase is the main focus of the paper, and hence composes the name of our framework.

In order to effectively handle network changes, *AFOCS* decomposes them into simpler events in such a way that each event can be quickly handled. Thanks to this feature, *AFOCS* can eventually obviate the need of reidentifying the network community structure every time. Both *FOCS* and *AFOCS* require β , the *overlapping threshold*, as the only input for their entire operations. This requirement is essential since network communities

can overlap at different scales, and as a result, we do need a control parameter in order to certify how much the overlap means to them.

The contributions of our work are threefold: First, we propose *AFOCS*, a two-phase adaptive framework for not only detecting and updating the overlapping network communities but also tracing their evolution over time. Theoretical analyses show *AFOCS* partially achieves more than 0.74% internal density of the optimal solution. Second, We evaluate *AFOCS* on both synthesized and real traces in comparison to both the state-of-the-art and the most popular static detection methods *COPRA* [10] and *CFinder* [8], as well as to recent adaptive methods *FacetNet* [11], *iLCD* [12] and *OSLOM* [13]. Empirical results show that *AFOCS* achieves both competitively results and high quality community structures in a timely manner. Finally, with *AFOCS*, we suggest a community based forwarding strategy for communication networks that reduces up to 11x overhead information while maintaining competitively delivery time and ratio. We also propose a new social-aware patching scheme for containing worms in OSNs, which helps reducing up to 7x the infection rates on Facebook dataset.

2 RELATED WORK

Community detection in complex networks has attracted huge attention since its introduction. In general, one can classify detection methods in two main categories including non-overlapping versus overlapping communities, and on static networks versus on dynamic networks. Many efficient methods have been proposed for detecting both non-overlapping and overlapping communities on *static networks*, among which *CFinder* [8] and *COPRA* [10] have remarked themselves as the most popular and most effective methods once fed with correct parameters [14]. A recent work [15] detailed a survey and benchmark on those algorithms.

Detecting communities on *dynamic networks*, both on overlapping and disjoint structures, has so far been an untrodden area. [4] proposed *QCA*, an adaptive method that can update and trace the network structure through a series of changes. This method is quick and effective, however, is not able to detect overlapped communities. [11] proposed *FacetNet*, a framework for analyzing communities in dynamic networks based on the optimization of snapshot costs. *FaceNet* is guaranteed to converge to a local optimal solution; however, its convergence speed is slow and its input asks for the number of network communities which are usually unknown in practice. [16] proposed *Stream – Group*, an incremental method to solve the community mining and detect the change points in weighted dynamic graphs. This method is modularity-based thus may inherit the resolution limit while discovering network communities. In another attempt, [17] suggested a particle-and-density based clustering method for dynamic networks,

based on the extended modularity and the concepts of nano-community and l -quasi-clique-by-clique. Apart from that, [12] proposed *iLCD* to find the overlapping network communities by adding edges and then merging similar ones. However, this model might not be sufficient in consideration with the dynamic behaviors of the network when new nodes are introduced or removed, or when existing edges are removed from the network. In [13], the author presented OSLOM, a framework for testing the statistical significance of a cluster with respect to a global null model (e.g., a random graph). To expand a community, OSLOM locally computes the value r for each neighbor node and tries to include that node into the current community. Comparison between *AFOCS* and these aforementioned dynamic methods is conducted in section 6.

3 PROBLEM FORMULATION

3.1 Basic notations

Let $G = (V, E)$ be an undirected unweighted graph representing a network where V is the set of N nodes and E is the set of M connections. Denote by $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ the network community structure, i.e., a collection of subsets of V where each $C_i \in \mathcal{C}$ and its induced subgraph form a community of G . In contrast with the disjoint community structure, we allow $C_i \cap C_j \neq \emptyset$ so that network communities can overlap with each other. For a node $u \in V$, let d_u , $N(u)$ and $Com(u)$ denote its degree, its neighbors and its set of community labels, respectively. For any $C \subseteq V$, let C^{in} and C^{out} denote the set of links having both endpoints in C and the set of links having exactly one endpoint in C , respectively. Finally, the terms *node-vertex* as well as *edge-link-connection* are used interchangeably.

3.2 Dynamic network model

Let $G_0 = (V_0, E_0)$ be the original input network and $G_t = (V_t, E_t)$ be a time dependent network snapshot recorded at time t . Denote by ΔV_t and ΔE_t the sets of nodes and edges to be added to or removed from the network at time t . Furthermore, let $\Delta G_t = (\Delta V_t, \Delta E_t)$ describe this change in terms of the whole network. The network snapshot at next time step $t + 1$ is expressed as a combination of the previous one together with the change, i.e., $G_{t+1} = G_t \cup \Delta G_t$. Finally, a *dynamic network* \mathcal{G} is defined as a sequence of network snapshots changing over time: $\mathcal{G} = (G_0, G_1, G_2, \dots)$.

3.3 Density function

In order to quantify the goodness of an identified community, we use the popular density function Ψ [18] defined as: $\Psi(C) = \frac{|C^{in}|}{\binom{|C|}{2}}$ where $C \subseteq V$. The more C approaches a clique of its size, the higher its density value $\Psi(C)$. In order to set up a threshold on the internal

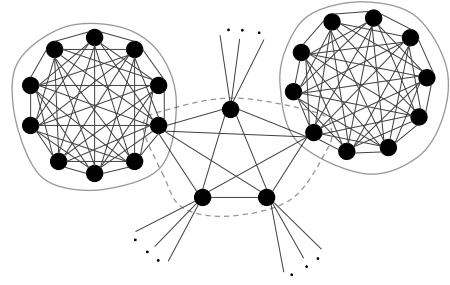


Fig. 1. Overlapped v.s. non-overlapped community structure. The central clique violates the general concept of community in both strong and weak senses

density that suffices for C to be a local community, we propose a function $\tau(C)$ defined as follows:

$$\tau(C) = \frac{\sigma(C)}{\binom{|C|}{2}} \text{ where } \sigma(C) = \binom{|C|}{2}^{1 - \frac{1}{\binom{|C|}{2}}}$$

Here $\sigma(C)$ is the threshold on the number of inner connections that suffices for C to be a local community. Particularly, a subgraph induced by C is a local community iff $\Psi(C) \geq \tau(C)$ or equivalently $|C^{in}| \geq \sigma(C)$.

Several functions with the same purpose have been introduced in the literature, for instance, in the work of [9], [19], and it is worth noting down the main differences between them and ours. First and foremost, our functions $\tau(C)$ and $\sigma(C)$ locally process on the candidate community C only and neither require any predefined thresholds or user-input parameters. Secondly, by Proposition 1, $\sigma(C)$ and $\tau(C)$ are increasing functions and closely approach C 's full connectivity as well as its maximal density. That makes $\sigma(C)$ and $\tau(C)$ relaxation versions of the traditional density function, yet useful ones as we shall see in the experiments.

Proposition 1: The function $f(n) = n^{1 - \frac{1}{n}}$ is strictly increasing for $n \geq 4$ and $\lim_{n \rightarrow \infty} f(n) = n$.

3.4 Objective function

Our objective is to find a community assignment for the set of nodes V which maximizes the overall internal density function $\Psi(\mathcal{C}) = \sum_{C \in \mathcal{C}} \Psi(C)$ since the higher the internal density of a community is, the clearer its structure would be. Unlike the case of disjoint community structure, in which the number of connections crossing communities should be less than those inside them, our objective does not take into account the number of outgoing links from each community.

To understand the reason, let us consider a simple example pictured in Figure 1. In the overlapping community structure point of view, it is clear that every clique should form a community on its own, and each community shares with the central clique exactly one node. However, in the disjoint community structure point of view, any vertex at the central clique has n internal and $2n$ external connections, which violates the

concept of a community in the strong sense. Furthermore, the internal connectivity of the central clique is also dominated by its external density, which implies the concept of a community in weak sense is also violated. (A community C is in *weak sense* if $|C^{in}| > |C^{out}|$, and in *strong sense* if any node in C has more links inward than outward C [20]).

3.5 Problem Definition

Given a dynamic network $\mathcal{G} = (G_0, G_1, G_2, \dots)$ where G_0 is the input network and G_1, G_2, \dots are network snapshots obtained through a collection of network topology changes $\Delta G_1, \Delta G_2, \dots$ over time. The problem asks for an adaptive framework to efficiently detect and update the network overlapping community structure \mathcal{C}_t at any time point t by only utilizing the information from the previous snapshot \mathcal{C}_{t-1} , as well as tracing the evolution of the network communities.

In the next section, we present our main contribution: an adaptive framework for (1) identifying basic overlapped community structure in a network snapshot and (2) updating as well as tracing the evolution of the network communities in a dynamic network model. First, we describe *FOCS*, a procedure to identify the basic communities in a static network, and then discuss in great detail how *AFOCS* adaptively updates these basic communities to cater with the evolution of the dynamic network.

4 BASIC COMMUNITY STRUCTURE

We describe *FOCS*, the first phase of our framework that quickly discovers the basic overlapping network community structure. In general, *FOCS* works toward the classification of network nodes into different groups by first locating all possible densely connected parts of the network (4.1), and then combining those who highly overlap with each other, i.e., those share a significant substructure (4.2). Finally, a final refinement to group unassigned nodes into different communities is conducted in (4.3).

In *FOCS*, β (the input overlapping threshold) defines how much substructure two communities can share. Note that *FOCS* fundamentally differs from [21] in the way it allows $|C_i \cap C_j| \geq 2$ for any subsets C_i, C_j of V , and consequently allows network communities to overlap not only at a single vertex but also at a part of the whole community.

4.1 Locating local communities

Local communities are connected parts of the network whose internal densities are greater than a certain level. In *FOCS*, this level is automatically determined based on the function $\tau(\cdot)$ and the size of each corresponding part. Particularly, a local community is defined based on a connection (u, v) when the number of internal connections in the subgraph induced by $C \equiv \{u, v\} \cup$

$(N(u) \cap N(v))$ exceeds $\sigma(C)$, or in other words, when $\Psi(C) \geq \tau(C)$ as illustrated in Figure 2(a).

However, there is a problem that might eventually arise: the containment of sub communities in an actual bigger one. Intuitively, one would like to detect a bigger community unified by smaller ones if the bigger community is itself densely connected. In order to filter this undesired case, we impose $\Psi(\bigcup_{i=1}^s C_i) < \tau(\bigcup_{i=1}^s C_i) \quad \forall s = 1 \dots |C|$ (note that some of these unifications do not contain all the nodes). In addition, we allow this locating procedure to skip over tiny communities of size less than 4. This condition is carried out from Proposition 1. This makes sense in terms of mobile or social networks where a group of mobile devices or a social community usually has size larger than 3, and intuitively agrees with the finding of [22], [23]. Thus, the condition $|C| \geq 4$ is imposed for any community C we discuss hereafter. The tiny communities will then be identified later. Alg. 1 describes this procedure.

Algorithm 1 Locating local communities

Input: $G = (V, E)$
Output: A collection of raw communities \mathcal{C}_r .

```

1: for  $(u, v) \in E$  do
2:   if  $Com(u) \cap Com(v) = \emptyset$  then
3:     Let  $C = \{u, v\} \cup N(u) \cap N(v)$ ;
4:     if  $|C^{in}| \geq \sigma(C)$  and  $|C| \geq 4$  then
5:       Check  $C$ 's connectivity if  $|C| = 5$ ;
6:       Define  $C$  a local community;
7:       /*Include  $C$  into the raw community structure*/
8:        $\mathcal{C}_r = \mathcal{C}_r \cup \{C\}$ ;
9:     end if
10:  end if
11: end for

```

Lemma 1: All local communities C 's detected by Alg. 1 satisfy $\Psi(C) \geq \tau(4) \approx 0.74$. Furthermore, other communities satisfying these conditions will also be detected by Alg. 1.

Theorem 1: The local community structure \mathcal{C}_r detected by Alg. 1 satisfies $\Psi(\mathcal{C}_r) \geq \tau(4) \times \Psi(OPT)$ where *OPT* is the optimal dense community assignment satisfying $\Psi(S) \geq \tau(4)$ for any $S \in OPT$.

Lemma 2: The time complexity of Alg. 1 is $O(dM)$ where $d = \max_{v \in V} d_v$ (Note: All proofs are included in the appendix).

4.2 Combining overlapping communities

After Alg. 1 finishes, the raw network community structure is pictured as a collection of (possibly overlapped) dense parts of the network together with outliers. As some of those dense parts can possibly share significant substructures, we need to merge them if they are highly overlapped. To this end, we introduce the overlapping score of two communities defined as follow

$$OS(C_i, C_j) = \frac{|I_{ij}|}{\min\{|C_i|, |C_j|\}} + \frac{|I_{ij}^{in}|}{\min\{|C_i^{in}|, |C_j^{in}|\}}$$

where $I_{ij} = C_i \cap C_j$. Basically, $OS(C_i, C_j)$ values how important the common nodes and links shared between C_i and C_j mean to the smaller community. In

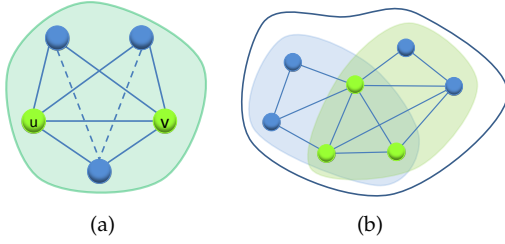


Fig. 2. (a) A local community C defined by a link (u, v) . Here $\Psi(C) = 0.9 > \tau(C) = 0.794$ (b) Merging two local communities sharing a significant substructure (OS score = $1.027 > \beta = 0.8$)

comparison with the distance metric suggested in [24], our overlapping score not only takes into account the fraction of common nodes but also values the fraction of common connections, which is crucial in order to combine network communities. Furthermore, $OS(\cdot, \cdot)$ is symmetric and scales well with the size of any community, and the higher the overlapping score, the more those communities in consideration should be merged. In this merging process, we combine communities C_i and C_j if $OS(C_i, C_j) \geq \beta$ (Figure 2(b)).

Algorithm 2 Combining local communities

Input: Raw community structure C_r
Output: A refined community structure C_f .

- 1: $C_f \leftarrow C_r$;
- 2: $Done \leftarrow false$;
- 3: **for** $C_i, C_j \in C_r$ and $!Done$ **do**
- 4: $Done \leftarrow true$;
- 5: **if** $OS(C_i, C_j) > \beta$ **then**
- 6: $C \leftarrow$ Combine C_i and C_j ;
- 7: /*Update the current structure*/
- 8: $C_f = (C_f \setminus \{C_i \cup C_j\}) \cup C$;
- 9: $Done \leftarrow False$;
- 10: **end if**
- 11: **end for**

The time complexity of Alg. 2 is $O(N_0^2)$ where N_0 is the number of local communities. Clearly, $N_0 \leq M$ and thus, it can be $O(M^2)$. However, when the intersection of two communities is upper bounded, by Lemma 3 we know that the number of local communities is also upper bounded by $O(N)$, and thus, the time complexity of Alg. 2 is $O(N^2)$. In our experiments, we observe that the running time of this procedure is, indeed, much less than $O(N^2)$.

Lemma 3: The number of raw communities detected in Alg. 1 is $O(N)$ when the number of nodes in the intersection of any two communities is upper bounded by a constant α .

4.3 Revisiting unassigned nodes

Even when the above two procedures are executed, there would still exist leftover nodes or edges due to their less attraction to the rest of the network. Because of its size constraint, the first procedure skips over tiny communities of sizes less than four and thus, may leave out some nodes unlabeled. These nodes will not be touched

in the second phase since they do not belong to any local communities, and consequently, will remain unassigned afterwards. Therefore, we need to revisit those nodes to either group them into appropriate communities or classify them as outliers based on their connectivity structures.

Alternatively, this process can be thought of as a community trying to hire adjacent unassigned nodes which are similar to the host community. To this end, we need a community fitness function in order to quantify the *similarity* between a node u and a neighbor community C . We find the fitness function $F_S = \frac{|S^{in}|}{2|S^{in}| + |S^{out}|}$ (where $S \subseteq V$) commonly used in [25], [9], [24] performs competitively in both synthesized and real-world datasets. Taking into account this fitness function, a community C will keep hiring any unassigned adjacent vertex of maximum similarity in a greedy manner, provided the newly joined vertex does not shrink down the community's current fitness value. If there is no such node, C is defined as a final network community. Nodes remained unlabeled through this last procedure are identified as outliers. The detailed algorithm is presented in Alg. 3.

Algorithm 3 Revisit Unassigned Nodes

Input: The refined community structure $C_f = \{C'_1, C'_2, \dots, C'_t\}$
Output: The basic community structure $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$

- 1: $\mathcal{C} = C_f$;
- 2: **for** $u \in V$ and $Com(u) = \emptyset$ **do**
- 3: Let $NC(u) = \{C_j \in \mathcal{C} | u \text{ is adjacent to } C_j\}$;
- 4: **for** $C_j \in NC(u)$ **do**
- 5: **if** $F_{C_j \cup \{u\}} \geq F_{C_j}$ **then**
- 6: $C_j \leftarrow C_j \cup \{u\}$;
- 7: $Com(u) \leftarrow Com(u) \cup \{j\}$;
- 8: **end if**
- 9: **end for**
- 10: **if** $Com(u) = \emptyset$ **then**
- 11: Classify u as an outlier;
- 12: **end if**
- 13: **end for**

5 DETECTING EVOLVING NETWORK COMMUNITIES

We describe *AFOCS*, the second phase and also the main focus of our detection framework. In particular, we use *AFOCS* to adaptively update and trace the network communities, which were previously initialized by *FOCS*, as the dynamic network evolves over time. Note that *FOCS* is executed only once on G_0 , after that *AFOCS* will take over and handle all changes introduced to the network.

Let us first discuss the various behaviors of the community structure when the network topology evolves over time. Suppose $G = (V, E)$ and $\mathcal{C} = \{C_1, C_2, \dots, C_n\}$ is the current network and its corresponding overlapping community structure, respectively. We use the term *intra links* to denote edges whose two endpoints belong to the same community, *inter links* to denote those with endpoints connecting different disjoint communities and the term *hybrid links* to stand for the others. For each community C of G , the number of connections joining C

with the others are lesser than the number of connections within C itself by definition

Intuitively, the addition of intra links or removal of inter links between communities of G will strengthen them and consequently, will make the structure of G more clear. Similarly, removing intra links from or introducing inter links to a community of G will decrease its internal density and as a result, loosen its internal structure. However, when two communities have less distraction to each other, adding or removing links makes them more attractive to each other and therefore, leaves a possibility that they can overlap with each other or can be combined to form a new community. The updating process, as a result, is very complicated and challenging since any insignificant change in the network topology could possibly lead to an unpredictable transformation of the network community structure.

In order to reflect these changes to a complex network, its underlying graph model is frequently updated by either inserting or removing a node or a set of nodes, or an edge or a set of edges. A scrutiny look into these events reveals that the introduction or removal of a set of nodes (or edges) can furthermore be decomposed as a collection of node (or edge) insertions (or removals), in which only a node (or only an edge) is inserted (or removed) at a time. Therefore, changes to the network at each time step can be viewed as a collection of simpler events whose details are as follow:

- *newNode* ($V + u$): A new node u and its adjacent edge(s) are introduced
- *removeNode* ($V - u$): A node u and its adjacent edge(s) are removed from the network.
- *newEdge* ($E + e$): A new edge e connecting two existing nodes is introduced.
- *removeEdge* ($E - e$): An edge e in the network is removed.

As we mentioned earlier, our adaptive framework initially requires a basic community structure C_0 . To obtain this basic structure, we apply *FOCS* algorithm at the first network snapshot, i.e., we execute *FOCS* on the network G_0 and then let *AFOCS* adaptively handle this structure as the network evolves.

5.1 Handling a new node

Let us discuss the first case when a new node u and its associated links are introduced to the network. Possibilities are (1) u may come with no adjacent edge or (2) with many of them connecting one or more possibly overlapped communities. If u has no adjacent edge, we simply join u in the set of outliers and preserve the current community structure.

The interesting case happens, and it usually does, when u comes with multiple links connecting one or more existing communities. Since network communities can overlap each other, we need to determine which ones u should join in in order to maximize the gained internal density. But how can we quickly and effectively do so?

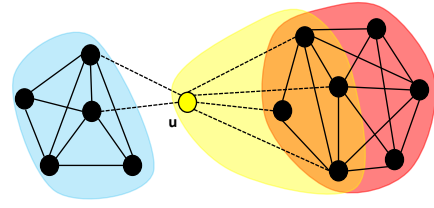


Fig. 3. When a new node u is introduced, u could gather some nodes from an existing community (red) to form a new community (yellow)

By Lemma 4, we give a necessary condition for a new node in order to join in an existing community, i.e., our algorithm will join node u in C if the number of connections u has to C suffices: $d_{ui} > \max\{\frac{2|C_i^{in}|}{|C_i|-1}, f(|C_i| + 1) - |C_i^{in}|\}$. However, failing to satisfy this condition does not necessarily imply that u should not belong to C , since it can potentially gather some substructure of C to form a new community (Figure 3). Thus, we also need to handle this possibility. Alg. 4 presents the algorithm.

Algorithm 4 Handling a new node u

Input: The current community structure C_{t-1}
Output: An updated structure C_t .

- 1: $C_1, C_2, \dots, C_k \leftarrow$ Adjacent communities of u ;
- 2: **for** $i = 1$ **do to** k
- 3: **if** $d_{ui} > \max\{\frac{2|C_i^{in}|}{|C_i|-1}, f(|C_i| + 1) - |C_i^{in}|\}$ **then**
- 4: $C_i \leftarrow C_i \cup \{u\}$;
- 5: **else**
- 6: $C \leftarrow N(u) \cap C_i$;
- 7: **if** $\Psi(C) \geq \tau(C)$ and $|C| \geq 4$ **then**
- 8: $C_i \leftarrow C_i \cup \{u\}$;
- 9: **end if**
- 10: **end if**
- 11: **end for**
- 12: /*Checking new community from outliers*/
- 13: **for** $v \in N(u)$ and $Com(v) \cap Com(u) = \emptyset$ **do**
- 14: $C \equiv N(u) \cap N(v)$;
- 15: **if** $\Psi(C) \geq \tau(C)$ and $|C| \geq 4$ **then**
- 16: Define C a new community;
- 17: **end if**
- 18: **end for**
- 19: Merging overlapping communities on C_1, C_2, \dots, C_k ;
- 20: Update C_t ;

Lemma 4: Suppose u is a newly introduced node with d_{ui} connections to each adjacent community C_i . u will join in C_i if $d_{ui} > \max\{\frac{2|C_i^{in}|}{|C_i|-1}, f(|C_i| + 1) - |C_i^{in}|\}$.

The analysis of Alg. 4 is shown by Lemma 5. In particular, we show that this procedure achieves at least 0.74% internal density of the optimal assignment for u , given the prior community structure.

Lemma 5: Alg. 4 produces a community assignment that, prior to the community combination process, achieves $\Psi(C_t) \geq \tau(4) \times \Psi(OPT(u)_t)$ where $OPT(u)_t$ is the optimal community assignment for u at time t , given the prior community structure C_{t-1} .

5.2 Handling a new edge

In case where a new edge $e = (u, v)$ connecting two existing vertices u and v is introduced, we divide it further into two four smaller cases: (1) e is solely inside

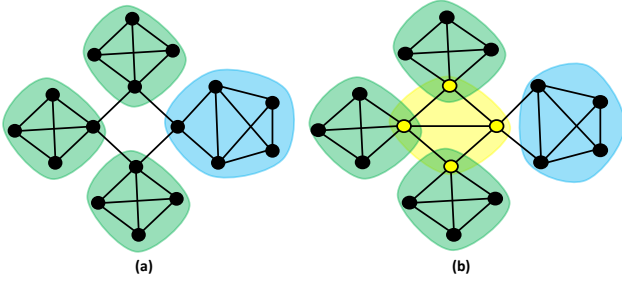


Fig. 4. (a) The network with 4 disjoint communities (b) When the central edge is added, the central nodes form a new community (yellow)

a single community C (2) e is within the intersection of two (or more) communities (3) e is joining two separated communities and (4) e is crossing overlapped communities. If e is totally inside a community C , its presence will strengthen C 's internal density and by Lemma 6, we know that adding e should not split the current community C into smaller substructures. The same reaction applies in the second subcase when e is within the intersection of two communities since their inner densities are both increased. Thus, in these first two cases, we leave the current network structure intact.

Handling the last two subcases is complicated since any of them can either have no effect on the current network structure or unpredictably form a new network community, and furthermore can overlap or merge with the others (Figure 4). However, there is still a possibility that the introduction of this new link, together with some substructure of C_u or C_v , suffices to form a new community that can overlap with not only C_u and C_v but also with some of the others. The other subcases can be handled similarly. Alg. 5 describe this procedure.

Algorithm 5 Handling a new edge (u, v)

Input: The current community structure \mathcal{C}_{t-1} .
Output: An updated community structure \mathcal{C}_t .

- 1: if $((u, v) \in \text{a single community OR } (u, v) \in C_u \cap C_v)$ then
- 2: $\mathcal{C}_t \leftarrow \mathcal{C}_{t-1}$;
- 3: else if $\text{Com}(u) \cap \text{Com}(v) = \emptyset$ then
- 4: $C \leftarrow N(u) \cap N(v)$;
- 5: if $\Psi(C) \geq \tau(C)$ then
- 6: Define C a new community;
- 7: Check for combining on $\text{Com}(u)$, $\text{Com}(v)$ and C ;
- 8: else
- 9: for $D \in \text{Com}(u)$ (or $D' \in \text{Com}(v)$) do
- 10: if $\Psi(D \cup \{v\}) \geq \tau(D)$ (or $\Psi(D' \cup \{u\}) \geq \tau(D')$) then
- 11: $D \leftarrow D \cup \{v\}$ (or $D' \leftarrow D' \cup \{u\}$)
- 12: end if
- 13: end for
- 14: Merging overlapping communities for D 's (or D');
- 15: end if
- 16: Update \mathcal{C}_t ;
- 17: end if

Lemma 6: If an new edge (u, v) is introduced solely inside a community C , it should not split C into smaller substructures.

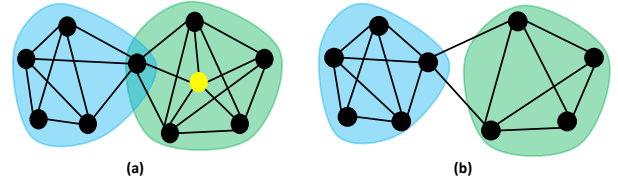


Fig. 5. (a) Two overlapped communities (b) When the central node is removed, the new structure consists of two disjoint communities

5.3 Removing an existing node

When an existing node u is about to be removed from the network, all of its adjacent edges will also be removed as a consequence. If u is an outlier, we can simply exclude u and its corresponding links from the current structure and safely keep the network communities unchanged.

In unfortunate situations where u is not an outlier, the problem becomes very challenging in the sense that the resulting community is complicated: it can either be unchanged, or broken into smaller communities, or could probably be merged with the other communities. To give a sense of this effect, let's consider two examples illustrated in Figure 5. In the first example, when C is almost a full clique, the removal of any node will not break it apart. However, if we remove a node that tends to connect the others within a community, the leftover module is broken into a smaller one together with a node that will later be merged to one of its nearby communities. Therefore, identifying the leftover structure of C is a crucial task once a vertex u in C is removed.

To quickly handle this task, we first examine the internal density of C excluding the removed node u . If the number of internal connections still suffices, e.g. $\Psi(C \setminus \{u\}) \geq \tau(C \setminus \{u\})$, we can safely keep the current network communities intact. Otherwise, we apply Alg. 1 on the subgraph induced by $C \setminus \{u\}$ to quickly identify the leftover modules in C , and then let these modules hire a set of unassigned nodes $\Psi(C)$ that help them increasing their inner densities. Finally, we locally check for community combination, if any, by using an algorithm similar to Alg. 2.

Algorithm 6 Removing a node u

Input: The current community structure \mathcal{C}_{t-1} .
Output: An updated structure \mathcal{C}_t .

- 1: for $C \in \text{Com}(u)$ and $\Psi(C \setminus \{u\}) < \tau(C \setminus \{u\})$ do
- 2: $LC \leftarrow$ Local communities by Alg 1 on $C \setminus \{u\}$;
- 3: for $C_i \in LC$ and $|C_i| \geq 4$ do
- 4: $S_i \leftarrow$ Nodes such that $\Psi(C_i \cup S_i) \geq \tau(C_i \cup S_i)$;
- 5: $C_i \leftarrow C_i \cup S_i$;
- 6: end for
- 7: Merging overlapping communities on LC ;
- 8: end for
- 9: Update \mathcal{C}_t ;

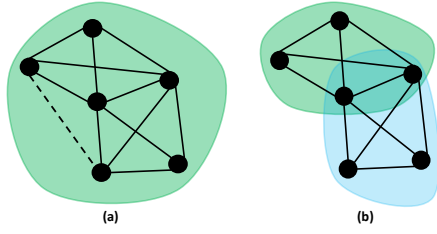


Fig. 6. (a) The original community (b) When the dotted edge is removed, the community is broken into two overlapped communities

5.4 Removing an edge

In the last situation when an edge $e = (u, v)$ is about to be removed, we divide it further into four subcases similar to those of a new edge (1) e is between two disjoint communities (2) e is inside a sole community (3) e is within the intersection of two (or more) communities and finally (4) e is crossing overlapping communities.

In the first subcase, when e is crossing two disjoint communities, its removal will make the network structure more clear (since we now have less connections between groups), and thus, the current communities should be keep unchanged. When e is totally within a sole community C , handling its removal is complicated since this can lead to an unpredictable transformation of the host module: C could either be unchanged or broken into smaller modules if it contains substructures which are less attractive to each other, as depicted in Figure 6. Therefore, the problem of identify the structure of the remaining module becomes the central part for not only this case but also for the others.

To quickly handle these tasks, we first verify the inner density of the remaining module and, again utilize the local community location method (Alg. 1) to locally identify the leftover substructures. Next, we check for community combination since these structures can possibly overlap with existing network communities. The detailed procedure is described in Alg. 7.

Algorithm 7 Removing an edge (u, v)

Input: The current structure \mathcal{C}_{t-1} .

Output: An updated community structure \mathcal{C}_t .

- 1: if (u, v) is an isolated edge then
 - 2: $\mathcal{C}_t = (\mathcal{C}_{t-1} \setminus \{u, v\}) \cup \{u\} \cup \{v\}$;
 - 3: else if $d_u = 1$ (or $d_v = 1$) then
 - 4: $\mathcal{C}_t = (\mathcal{C}_{t-1} \setminus \mathcal{C}(u)) \cup \{u\} \cup \mathcal{C}(v)$;
 - 5: else if $C \equiv \mathcal{C}(u) \cap \mathcal{C}(v) = \emptyset$ then
 - 6: $\mathcal{C}_t = \mathcal{C}_{t-1}$;
 - 7: else if $\Psi(\mathcal{C} \setminus (u, v)) < \tau(\mathcal{C} \setminus (u, v))$ then /*Here $C \neq \emptyset$ */
 - 8: $LC \leftarrow$ Local communities by Alg 1 on $\mathcal{C} \setminus (u, v)$;
 - 9: Define each $L \in LC$ a local community of \mathcal{C}_{t-1} ;
 - 10: Merging overlapping community on L 's;
 - 11: end if
 - 12: Update \mathcal{C}_t ;
-

5.5 Remarks

Note that the ultimate goal of our framework is to adaptively detect and update the community structure

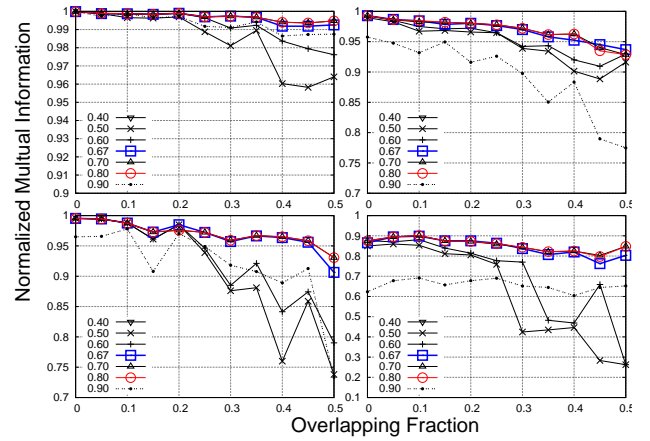


Fig. 7. NMI scores for different values of β . $N = 5000$ (top), $N = 1000$ (bottom), $\mu = 0.1$ (left), $\mu = 0.3$ (right).

as the network evolves, i.e., to mainly deal with the dynamics of a mobile network. As a result, we mainly put our focus on *AFOCS*. Although *FOCS*, the first detection phase, appears to be a centralized algorithm, it is executed only once at the very first network snapshot whereas *AFOCS* stays up and locally handles all changes as the network evolves over time. That said, we do not execute *FOCS* again. Furthermore, *AFOCS* can be run independently with *FOCS*, i.e., one can use any localized detection algorithm to identify a basic community structure at the first phase. Thus, *AFOCS* can be easily apply to mobile network problems, as presented in sections 7 and 8.

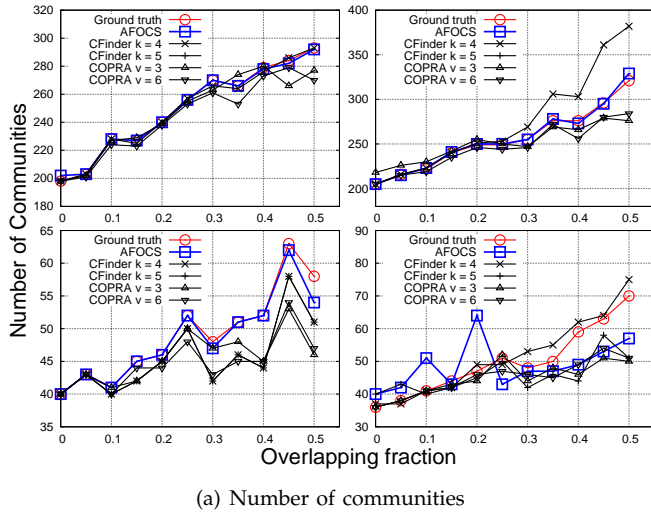
6 EXPERIMENTAL RESULTS

In this section, we first present the empirical results of *AFOCS* in comparison with two *static* detection methods: *CFinder* - the most popular method [8], and *COPRA* - the most effective method [10]. We next compare the performance of *AFOCS* with other dynamic methods including *OSLOM* [13], *FacetNet* [11] and *iLCD* [12].

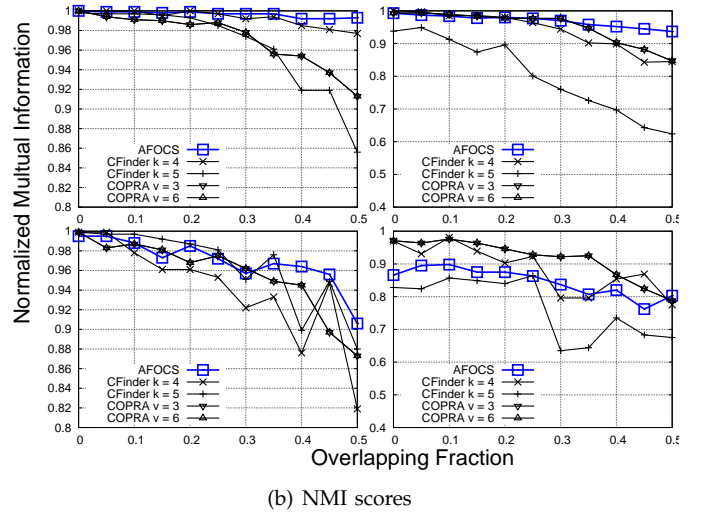
Data Sets: We use networks generated by the well-known LFR overlapping benchmark [15], the ‘de facto’ standard for evaluating overlapping community detection algorithms. Generated networks follow power-law degree distributions and contain embedded overlapping communities (the ground truth) of varying sizes that capture the internal characteristics of real-world networks.

Set up: To fairly compare with *COPRA* and to avoid being biased, we keep the parameters close to [10]: the minimum and maximum community sizes are $c_{min} = 10$ and $c_{max} = 50$, each vertex belongs to at most two communities, $o_m = 2$. $N = 1000$ and $N = 5000$ The mixing rate $\mu = 0.1$ and $\mu = 0.3$. The *overlapping fraction* γ , which determines the fraction of overlapped nodes, is from 0 to 0.5. Since *COPRA* is nondeterministic, we run it 10 times on each instance and select the best result.

Metrics: We evaluate following metrics.



(a) Number of communities



(b) NMI scores

Fig. 8. Comparison among *AFOCS*, *COPRA* and *CFinder* methods. $N = 5000$ (top), $N = 1000$ (bottom), $\mu = 0.1$ (left), $\mu = 0.3$ (right).

(1) The generalized *Normalized Mutual Information* (*NMI*) [9] specially built for overlapping communities. *NMI* scores the similarity between the detected network communities and the ground truth. This is a standardized measure since $NMI(U, V) = 1$ if structures U and V are identical and 0 if they are totally separated.

(2) The *number of communities*, ignoring singleton communities and unassigned nodes. A good community detection method should produce roughly the same number of communities with the known ground truth.

6.1 Choosing the overlapping threshold β

The overlapping threshold β is the only input parameter required by our framework, and thus, determining its appropriate value plays an important role in assessing *AFOCS*'s performance. To best determine this threshold, we run *AFOCS* on generated networks with different values of β , and record the similarities between the detected communities and the ground-truth via *NMI* scores (Figure 7). Of course, the higher *NMI* scores imply the better β values.

As a threshold parameter, β controls how much sub-structure communities can have in common. The smaller values of β imply the more we allow network communities to overlap with each other, and vice versa. Similarly, β can be thought of as the zooming scale of the network structure where lower β 's reveal the coarser and higher β 's reveal the finer structure. As depicted in Figure 7, the best values for β are ranging from 0.67 to 0.80, among which $\beta = 0.70$ yields the best community similarity (*NMI* scores are ranging from 0.8 to 1) in all of the generated networks. Therefore, we fix the overlapping threshold in *AFOCS* to be 0.70 hereafter.

6.2 Reference to static methods

We show our results in groups of four. For each case we vary the overlapping fraction γ from 0 to 0.5 and analyze

the results found by *AFOCS*, *CFinder*, and *COPRA*. We only present results when corresponding parameters give top performance for *CFinder* (clique size $k = 4, 5$) and *COPRA* (max. communities per vertex $v = 3, 6$).

Figure 8(a) shows the number of communities found by *AFOCS*, *COPRA* and *CFinder* and the ground truth. It reveals from this figure that the numbers of communities found by *AFOCS*, marked with squares, are the closest and almost identical to the ground truth as the overlapping fraction gets higher. There is an exception when $N = 1000$ and $\mu = 0.3$ which we will discuss later. As one can infer from Figure 8(b), *AFOCS* achieves the highest performance among all methods with much more stable. A common trend in this test is the performances of all methods degrade (1) when the mixing rate μ increases, i.e., when the community structure becomes more ambiguous or (2) when the size of network decreases while the mixing rate μ stays the same. Even though *AFOCS* is not very competitive only when both negative factors happen in the bottom-right chart as $N = 1000$ and $\mu = 0.3$, it is in general the best performer.

The significant gap is observed when the mixing rate gets higher ($\mu = 0.3$) and the network size gets smaller ($N = 1000$). *AFOCS* provides less numbers of communities than those of the ground truth but with much higher overlapping rates. The reason is with a larger mixing rate μ , a node will have more edges connecting vertices in other communities, thus increases the chance that *AFOCS* will merge highly overlapped communities. Hence, *AFOCS* creates less but with larger size communities. We note that this 'weakness' of *AFOCS* is controversial as when the mixing rate increases, the ground truth does not necessarily coincide with the structure implied by the network's topology. Extensive experiments show the ability of *AFOCS* in identifying high quality overlapping communities. In addition,

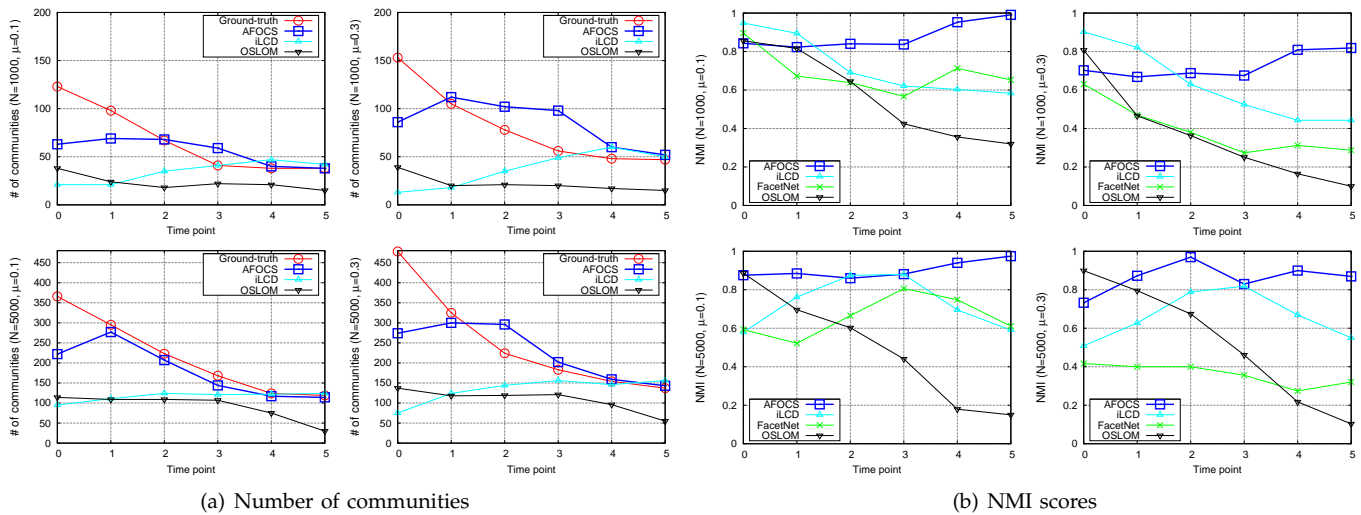


Fig. 9. Comparison among *AFOCS*, *iLCD*, *FacetNet* and *OSLOM* dynamic methods.

we found *AFOCS* runs substantially faster than the other competitors: on the network containing 63K nodes, *AFOCS* is 150x faster than *COPRA* while *CFinder* is unable to finish its tasks.

6.3 Reference to other dynamic methods

We next observe the performance of *AFOCS* in reference to two dynamic methods *FacetNet* [11], *iLCD* [12] and *OSLOM* [13]. Since the ground-truth communities are known on synthesized datasets, fair comparisons among three methods can be obtained via their NMI scores and running times. Of course, the higher its NMI scores with less time consuming, the better the method seems to be.

Each synthesized dynamic network is simulated via 5 snapshots, in which the basic communities are formed by using 50% of the network data with approximately 770 edges added to each growing snapshot at a time. Since *FacetNet* requires the number of communities a priori, we input this method the actual number as mined from the ground-truth. For *iLCD* and *OSLOM* methods, we keep the default setting as provided in their deliverable.

The NMI scores of four methods are presented in Figure 9(b) and 9(a). It reveals from these subfigures that the NMI scores of *AFOCS* are higher than those of *FacetNet*, *iLCD* and *OSLOM*. In particular, the NMI scores of *AFOCS* are about just 5-7% lag behind that of *OSLOM* and *iLCD* in the first 2 network snapshots, while are much better than the others at the end of the evolution. The *OSLOM*' NMI values are very high at the very beginning, however, they tend to decrease quickly as more connections and nodes are introduced. The NMI scores of *iLCD* and *FacetNet* tend to fluctuate and also decrease significantly at the last snapshot. *AFOCS*, in the other trend, keeps its NMI scores high and wealthy, especially at the end of the network evolution. This implies communities discovered by *AFOCS* are of higher similarity to the ground-truth than the other dynamic methods, especially in the long run.

The number of communities found by all methods are reported in Figure 9(a). Of course, the closer these detected numbers of communities to the ground-truth, the better the method are believed to be. As revealed in the subfigures of Figure 9(a), these quantities discovered by *AFOCS* tend to closely approach the actually numbers, even when the mixing rates are high (right figures). The highest similarity between these numbers of communities is possibly the best explanation for the high NMI scores of *AFOCS* over the other competitors.

We next take a look at the running time of all methods in these synthesized networks. *AFOCS* requires at most 5 seconds to finish updating each network snapshot whereas *FacetNet* asks for more than 25 seconds (5x more time consuming) in the networks with just 5000 nodes. *iLCD* and *OSLOM* also perform fast in these generated datasets; however, the similarity of the detected communities and the ground-truth is surprisingly poor, as revealed from the results. Therefore, in terms of dynamic approaches, we strongly believe that *AFOCS* achieves competitive community detection results in a timely manner. These results also provide us the confidence when applying *AFOCS* to analyze real-world networks.

7 COMMUNITY-BASED FORWARDING IN COMMUNICATION NETWORKS

We present a practical application where the detection of overlapping network communities plays a vital role in forwarding strategies in communication networks. With the helpful knowledge of the network community structure, we propose a new community-based forwarding algorithm that significantly reduces the number of duplicate messages while maintaining competitive delivery times and ratios, which are essential factors of a forwarding strategy.

Many routing methods based on the discovery of network community structure have been proposed in the

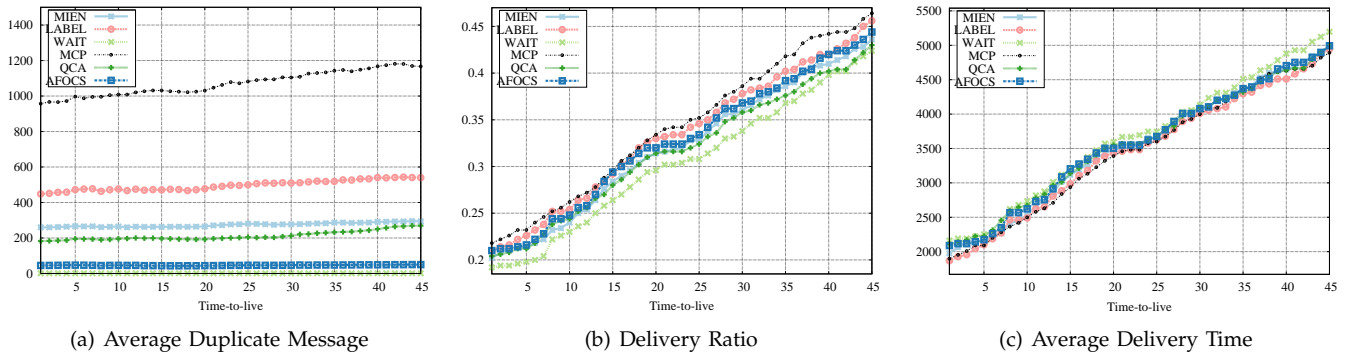


Fig. 10. Experimental results on the Reality Mining data set

literature [26], [27], [28]. However, the community detection cores in those strategies encounter (1) the lack of knowledge about overlapping communities and (2) the repeated identification of communities as the network evolves. The second issue is computationally costly and time consuming, thus may reduce the performance of those forwarding strategies.

7.1 Message forwarding strategy

Let us first discuss how our new forwarding algorithm works in practice and then how *AFOCS* helps it to overcome the above limitations. We use *AFOCS* to detect overlapping communities and keep it up-to-date as the network changes. Each node in a community is assigned the same label and each overlapped node u has a set of corresponding labels $Com(u)$. During the network operation, if a device u carrying the message meets another device v who indeed shares more common community labels with the destination than u , i.e., $|Com(v) \cap Com(dest)| > |Com(u) \cap Com(dest)|$, then u will forward the message to v . The same actions then apply to v as well as to devices that v meets.

The intuition behinds this strategy is that if v shares more communities with the destination nodes, it is likely that v will have more chances to deliver the message to the destination. By doing in this way, we not only have higher chances to correctly forward the messages but also generate much less duplicate messages. Due to its adaptive nature and the ability of identifying overlapping communities, *AFOCS* helps our algorithm to overcome the above shortcomings naturally. This explains why our forwarding algorithm can significantly reduce the number of duplicate messages while maintaining very competitive delivery times and ratios.

7.2 Experiment setup

We compare six forwarding strategies (1) *MIEN*: A recently proposed social-aware routing strategy on MANETs [1] (2) *LABEL*: A node will forward the messages to another node if it is in the same community as the destination [29] (3) *WAIT*: The source node waits and keeps forwarding the message until it meets the

destination (4) *MCP*: A node keeps forwarding the messages until they reach the maximum number of hops (5) *QCA*: A *LABEL* version utilizing *QCA* [4] as the adaptive disjoint community detection method and lastly (6) *AFOCS*: Our newly proposed forwarding algorithm equipped with *AFOCS* as a community detection and update core.

Results of *WAIT* and *MCP* algorithms provide us the lower and upper bounds of important factors: message delivery ratio, time redundancy and message redundancy. Our experiments are performed on the Reality Mining dataset provided by the MIT Media Lab [30]. This dataset contains communication, proximity, location, and activity information from 100 students at MIT over the course of the 2004-2005 academic year. In particular, we take into account the Bluetooth information to construct the underlying communication network and evaluate the performance of the above six routing strategies.

In each experiment, 500 message sending requests are randomly generated and distributed in different time points. To control the forwarding process, we use *hop-limit*, *time-to-live*, and *max-copies* parameters. A message cannot be forwarded more than *hop-limit* hops in the network or exist in the process longer than *time-to-live*, otherwise it will be automatically discarded. Moreover, the maximum number of same messages a device can forward to the others is restricted by *max-copies*. Experiments results are repeated and results are averaged for consistency.

7.3 Results

Our results are presented in Figures 10(a), 10(b), 10(c). The first observation reveals that our proposed forwarding algorithm achieves the lowest number of duplicate messages as depicted in Figure 10(a), and even far better than the second best method *QCA*. On average, only 46.5 duplicate messages are generated by *AFOCS* during evaluation process in contrast with 212.2 of *QCA*, 274.2 of *MIEN*, 496.4 of *LABEL* and the huge 1071.0 overhead messages of *MCP*. Thus, on the number of duplicate messages, *AFOCS* strikingly achieves improvement factors of 4.5x, 5x, 11x and 23x over these mentioned

strategies, respectively. These extremely low overhead strongly imply the efficiency of *AFOCS* in communication networks.

Figures 10(b) and 10(c) present our results on the other two important factors, the message delivery ratios and delivery times. These figures supportively indicate that *AFOCS* achieves competitive results on both of these vital factors. In general, *AFOCS* is the second best strategy with almost no noticeable different between itself and the leader method *LABEL*. On average, *AFOCS* gets 33% of the total messages delivered in 3569.2s and only a little bit lags over *MCP* (34% in 3465.3s) and *LABEL* (slightly over 33% in 3462.7s), and is far better than *MIEN* (32% in 3537.6s) and *QCA* (32% in 3572.2s). This can be explained by the advantages of knowing the overlapping community structure: the disjoint network communities in *QCA* and *MIEN* can possibly have messages forwarded to the wrong communities when the destination changes its membership. With the ability of quickly updating the network structure, *AFOCS* can efficiently cope with this change and thus, can still provide the most updated forwarding information.

In summary, *AFOCS* helps our forwarding strategy to reduce up to 11x the number of duplicate messages while keeping good average delivery ratio and time. These experimental results are highly competitive and supportively confirm the effectiveness of *AFOCS* and our new routing algorithm on communication networks.

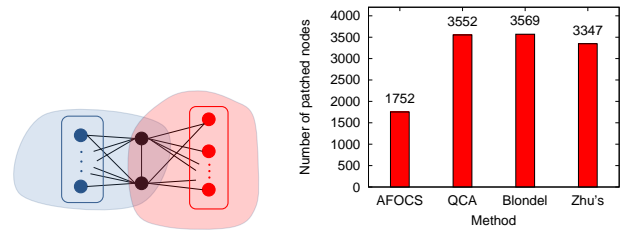
8 CONTAINING WORMS USING OVERLAPPING COMMUNITIES

We show another application of *AFOCS* in worm containment problem on OSNs. OSNs are good places for people to socialize online or to stay in touch with friends and colleagues. However, when some of the users are infected with malicious software, such as viruses or worms, OSNs are also fertile grounds for their rapid propagations. Since mobile devices are able to access online social applications nowadays, worms and viruses now can target computers [4] and mobile devices [3].

Recently, community structure-based methods have been proven to be effective solutions to prevent worms from spreading out wider on not only social networks [4], [7] but also cellular networks [3]. Due to the high and low frequencies of interactions inside and between communities, worms spread out quicker within a community than between communities. Therefore, an appropriate reaction should first contain worms into only infected communities, and then prevent them from getting outside. This strategy can be accomplished by patching the most *influential members* who are well-connected not only to members of their community but also to people in other communities.

8.1 Setup

In our experiments, we use Facebook network dataset collected in [31]. This data set contains friendship information and wall posts among New Orleans regional



(a) Influential users selection (b) Number of patched nodes

Fig. 11. *OverCom* patching scheme

network, spanning from Sep 2006 to Jan 2009. The data set contains more than 63.7K nodes (users) connected by more than 1.5 million friendship links. We keep other parameters as well as the “Koobface” worm propagation model the same as [7] for comparison convenience. With the advantages of knowledge overlapping communities, we are able to develop a better and more efficient patching scheme. In particular, we enhance the patching scheme presented in in [7] to take the advantage of the overlap regions: nodes in the boundary of overlapped regions are selected for patching (Figure 11(a)). Alg 8 details the adjusted scheme.

Algorithm 8 *OverCom* Patching Scheme

Input: $G = (V, E)$ and $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ detected by *AFOCS*
Output: A set of patched nodes IS .

```

1:  $IS = \emptyset$ ;
2: for  $C_i, C_j \in \mathcal{C}$  do
3:   if  $C_i \cap C_j \neq \emptyset$  then
4:     %Choose the neighbors of overlapped nodes as influential ones%
5:      $IS = IS \cup N(u) \quad \forall u \in C_i \cap C_j$ ;
6:   end if
7: end for
8: %Patch distribution procedure%
9: for  $u \in IS$  do
10:  Send patches to  $u$ ;
11:  Let  $u$  redistribute patches to  $w \in IS \setminus N(u)$ ;
12: end for

```

8.2 Results

We compare the *OverCom* patching scheme and overlapping communities found by *AFOCS* to those using disjoint communities proposed by Blondel *et al* [32], *QCA* by Nguyen *et al* [4] and Clustering based method suggested by Zhu *et al* [3]. The number of patched nodes is shown in Figure 11(b). Both the number of patched nodes and the infection rates decline remarkably. In particular, the number of nodes to send patch in *AFOCS* is substantially smaller by half of those required by *Blondel*, *QCA* as well as *Zhu's* methods: only 1725 nodes over 63K nodes in the networks are needed to be patched by *OverCom* patching scheme, while the other schemes require nearly twice ($\geq 3,300$ nodes). The reason behind this improvement is due to the nature of our *AFOCS* framework, the neighbors of the overlapped nodes should not be to far away from the center of each community, thus they can easily redistribute the patches once received.

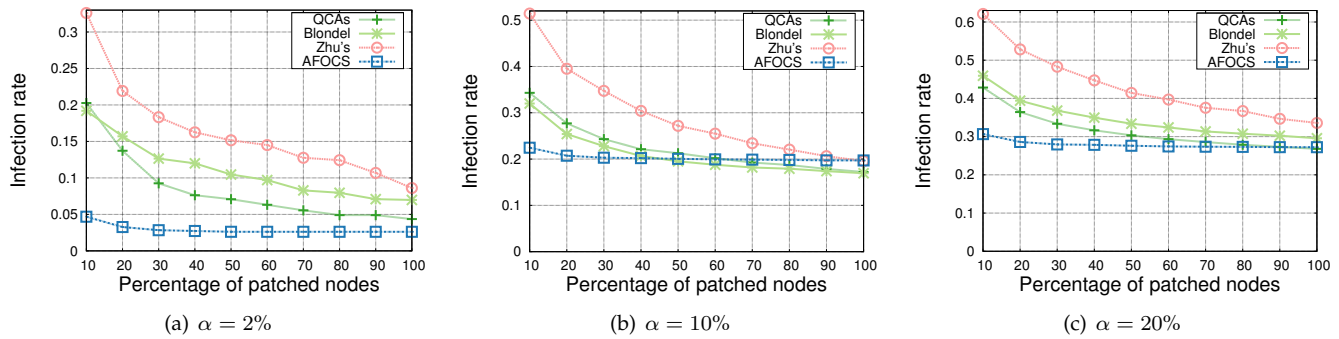


Fig. 12. Infection rates between four methods

We next present the achieved infection rates with alarming thresholds (the fraction of infected nodes over all nodes) $\alpha = 2\%$, 10% and 20% , respectively. This threshold alarms the distribution process as soon as the infected rate goes beyond α . The results are reported in Figures 12(a), 12(b), 12(c), respectively. In general, the higher α (i.e., the longer we wait), the more nodes we have to send patches and the higher infection rate. OverCom with AFOCS achieves the lowest infection rates in almost all the experiments and just a little bit lag behind when $\alpha = 10\%$. In particular, when $\alpha = 2\%$, AFOCS helps OverCom to remarkably reduce from 1.6x up to 4.3x the infection rates of QCA, from 2.6x up to 4x the infection rates of Blondel and 3.2x to 7x those of Zhu's method. When $\alpha = 10\%$, AFOCS + OverCom achieves average improved rates of 9% over QCA, 5% over Blondel and 43% over Zhu's methods. As $\alpha = 20\%$, the average improvements are 12%, 23% and 53%, respectively. Due to the nature of the event handling processes, the neighbors of overlapped nodes are not located far away from the rest of their communities. As a result, they can help to distribute patches to more users in the communities, hence help to lower the infection rates of AFOCS. These improvement factors, again, confirm the effectiveness of our proposed method.

9 CONCLUSION

In this paper, we presented AFOCS, a two-phase framework for detecting network overlapping communities as well as tracing their evolution in dynamic mobile networks. Analyses show that AFOCS partially achieves no less than 83% internal density of the optimal community assignment. Experiments on synthesis and real-world data traces show good results. We show two mobile applications, namely forwarding and routing in MANETs and worm containment on OSNs, in which AFOCS significantly helps to increase the performances up to 11x and 7x, respectively. These results confirm the effectiveness of AFOCS as well as its applicability in mobile applications. In our future work, we plan to improve the performance of AFOCS by taking into account the stability of the network community structure. In particular, we aim to discover not only overlapping communities

but also stable ones whose internal interactions remain significant over a long period of time, or over a random perturbation. We hope that this new feature will further improve AFOCS on mobile networks.

ACKNOWLEDGMENTS

This work is supported by the NSF under CAREER Award #0953284, by the DTRA YIP grant HDTRA1-09-1-0061, and by the DTRA grant HDTRA1-08-10.

REFERENCES

- [1] T. N. Dinh, Y. Xuan, and M. T. Thai, "Towards social-aware routing in dynamic communication networks," *IPCCC*, '09.
- [2] B. Pasztor, L. Mottola, C. Mascolo, G. Picco, S. Ellwood, and D. Macdonald, "Selective reprogramming of mobile sensor networks through social community detection," *Wireless Sensor Networks*, '10.
- [3] Z. Zou, G. Cao, S. Zhu, S. Ranjan, and A. Nucci, "A social network based patching scheme for worm containment in cellular networks," in *INFOCOM*, '09.
- [4] N. P. Nguyen, T. N. Dinh, Y. Xuan, and M. T. Thai, "Adaptive algorithms for detecting community structure in dynamic social networks," *INFOCOM*, '11.
- [5] M. Girvan and M. E. J. Newman, "Community structure in social and biological networks," *PNAS*, '02.
- [6] M. A. Porter, J.-P. Onnela, and P. J. Mucha, "Communities in networks," *Notices of the AMS*, '09.
- [7] N. P. Nguyen, Y. Xuan, and M. T. Thai, "A novel method for worm containment on dynamic social networks," *MILCOM*, '10.
- [8] G. Palla, I. Derenyi, I. Farkas, and T. Vicsek, "Uncovering the overlapping community structure of complex networks in nature and society," *Nature*, '05.
- [9] A. Lancichinetti, S. Fortunato, and K. Jnos, "Detecting the overlapping and hierarchical community structure in complex networks," *New J. of Phys.*, '09.
- [10] S. Gregory, "Finding overlapping communities in networks by label propagation," *New J. of Physics*, '10.
- [11] Y.-R. Lin, Y. Chi, S. Zhu, H. Sundaram, and B. L. Tseng, "Analyzing communities and their evolutions in dynamic social networks," *ACM TKDD*, '09.
- [12] R. Cazabet, F. Amblard, and C. Hanachi, "Detection of overlapping communities in dynamical social networks," in *SocialCom*, '10.
- [13] A. Lancichinetti, F. Radicchi, J. Ramasco, and S. Fortunato, "Finding statistically significant communities in networks," *PLoS ONE* 6: e18961, '11.
- [14] L. Peel, "Estimating network parameters for selecting community detection algorithms," *FUSION*, '10.
- [15] A. Lancichinetti and S. Fortunato, "Community detection algorithms: A comparative analysis," *Phys. rev. E*, 80, '09.
- [16] D. Duan, Y. Li, Y. Jin, and Z. Lu, "Community mining on dynamic weighted directed graphs," in *CNIKM*, '09.
- [17] M.-S. Kim and J. Han, "A particle-and-density based evolutionary clustering method for dynamic networks," *VLDB Endow.*, 2009.

- [18] S. Fortunato and C. Castellano, "Community structure in graphs," *arXiv*, '07.
- [19] A. Lzr, D. bel, and T. Vicsek, "Modularity measure of networks with overlapping communities," *Euro. Let.*, '10.
- [20] F. Radicchi, C. Castellano, F. Cecconi, V. Loreto, and D. Parisi, "Defining and identifying communities in networks," *Proc. Natl. Acad. Sci. USA*, '04.
- [21] J. P. B. Y-Y Ahn and S. Lehmann, "Link communities reveal multiscale complexity in networks," *Nature*, '10.
- [22] S. Fortunato, "Community detection in graphs," *Phys. Reports*, '10.
- [23] J. Leskovec, K. J. Lang, A. Dasgupta, and M. W. Mahoney, "Statistical properties of community structure in large social and information networks," in *WWW*, '08.
- [24] C. Lee, F. Reid, A. McDaid, and N. Hurley, "Detecting highly overlapping community structure by greedy clique expansion," in *KDD*, '10.
- [25] M. Goldberg, S. Kelley, M. Magdon-Ismail, K. Mertsalov, and A. Wallace, "Finding overlapping communities in social networks," in *SOCIALCOM*, '10.
- [26] P. Hui, E. Yoneki, S. Y. Chan, and J. Crowcroft, "Distributed community detection in delay tolerant networks," in *MobiArch*, '07.
- [27] E. M. Daly and M. Haahr, "Social network analysis for routing in disconnected delay-tolerant manets," in *MobiHoc*, '07.
- [28] P. Hui, J. Crowcroft, and E. Yoneki, "Bubble rap: social-based forwarding in delay tolerant networks," in *MobiHoc*, '08.
- [29] P. Hui and J. Crowcroft, "How small labels create big improvements," *PERCOMW*, '07.
- [30] E. Nathan and A. Pentland, "Reality mining: sensing complex social systems," *Personal Ubiquitous Comput.*, '06.
- [31] B. Viswanath, A. Mislove, M. Cha, and K. P. Gummadi, "On the evolution of user interaction in facebook," in *2nd ACM SIGCOMM Workshop on Social Networks*, '09.
- [32] V. D. Blondel, J. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *J. Stat. Mech.: Theory and Experiment*, '08.

PLACE
PHOTO
HERE

is a recipient of DoD Young Investigator Awards and NSF CAREER awards. She is a member of the IEEE.

My T. Thai (M'06) received her Ph.D. degree in computer science from the University of Minnesota, Twin Cities, in 2006. She is an associate professor in the CISE Department, University of Florida. Her current research interests include algorithms and optimization on network science and engineering. She also serves as an associate editor for the Journal of Combinatorial Optimization (JOCO) and Optimization Letters and a conference chair of COCOON 2010 and several workshops in an area of network science. She

PLACE
PHOTO
HERE

Nam P. Nguyen received his B.S. and M.S. degrees in Applied Mathematics from Vietnam National University, HCMC (in 2007) and University of Ohio, USA (in 2009). He is currently a Ph.D. student at the CISE Department, University of Florida, under the supervision of Dr. My T. Thai. His research interests include Dynamic complex networks, Social networks; Cascading failures; Combinatorial optimization and Approximation Algorithms.

PLACE
PHOTO
HERE

Yilin Shen is currently a PhD student at the Department of Computer and Information Science and Engineering, University of Florida, under the supervision of Dr. My T. Thai. His research focuses on vulnerability assessment and security of complex networks, including communication networks, wireless sensor networks and social networks, and designing approximation algorithms for network optimization problems. He is a student member of the IEEE.

PLACE
PHOTO
HERE

Thang N. Dinh (S'11) received the B.A. degree in Information Technology from Vietnam National University, Hanoi, Vietnam (2007). He is currently a Ph.D. student at the CISE Department, University of Florida, under the supervision of Dr. My T. Thai. His research focuses on designing combinatorial optimization methods for dynamic complex networks and mobile ad hoc network including network vulnerability, dynamic community structure, and fast information propagation. He is a IEEE student member.

APPENDIX

PROOF OF LEMMA 1

Alg. 1 will examine every edge $(u, v) \in E$ (except those whose endpoints are already in the same community), and by this greedy nature, any local community it detects has $|C| > 4$ and $\Psi(C) \geq \tau(C) \geq \tau(4) \approx 0.74$.

We now show that any community C satisfying $|C| \geq 4$ and $\Psi(C) \geq \tau(C) \geq \tau(4)$ will also be detected by Alg. 1. Suppose otherwise, that is there exists a community C satisfying these two conditions and is not detected by Alg. 1. To prove that this is not the case, we do the following (1) Construct a community D which is not detected by Alg. 1 with $|D| = n \equiv |C|$ and $\Psi(D)$ is maximized, and (2) show that $\Psi(D) < \tau(D)$. Because $|D| = |C|$, it implies $\tau(D) = \tau(C)$. However, since $\Psi(D)$ is maximized, $\Psi(D) \geq \Psi(C)$ which in turn implies $\Psi(C) \leq \Psi(D) < \tau(D) = \tau(C)$. This raises a contradiction to our original assumption, and thus concludes the proof.

To construct D , we do as follow (i) make D a clique of size n , and (ii) remove edges from D one by one until D cannot be detected by Alg. 1. By doing in this way, $\Psi(D)$ is maximized iff the number of removed edges is minimized.

It is easy to find the least number of edges we have to remove from D is $n/2$ if n is even and $n/2 - 1$ if n is odd. Therefore, $m_D = n(n-1)/2 - n/2$ if n is even, and $m_D = n(n-1)/2 - (n-1)/2$ if n is odd. Now, $\Psi(D) < \tau(D)$ iff $m_D < \binom{n(n-1)}{2}^{1 - \frac{2}{n(n-1)}}$. Let $f(n)$ be the difference between the left and the right hand sides, we show that $f(n) < 0$ as n increases. Taking the derivative of $f(n)$ gives $f'(4) < 0$ and $f(n) < f(4) < 0$ for all even $n > 4$, and $f'(7) < 0$ and $f(n) < f(7) < 0$ for all odd $n > 7$. When $n = 5$, $f(5) > 0$ but this is the only exception and thus, can be handled easily in line 5 of Alg. 1. Therefore, we have $\Psi(D) < \tau(D)$, and hence, the conclusion follows.

PROOF OF THEOREM 1

Let C_r be the local community structure returned by Alg. 1, and OPT be the optimal solution of the dense community assignment satisfying $\Psi(S) \geq \tau(4)$ for any $S \in OPT$. Let $k = |OPT|$. Clearly $\Psi(OPT) \leq k$. By Lemma 1, we know that Alg. 1 can detect as many communities as OPT but probably with less internal density. Moreover, since Alg. 1 only skips over edges in a community, it ensures that no real community is a substructure of a bigger one. Hence, we have $\Psi(C_r) \geq \tau(4) \times k \approx 0.74 \times \Psi(OPT)$. This also implies that Alg. 1 is an 0.74-approximation algorithm for finding local densely connected communities.

PROOF OF LEMMA 2

Time to examine an edge (u, v) is $|N(u)| + |N(v)| = d_u + d_v$. However, when u and v are in the same community, (u, v) will be skipped. Therefore, the total time complexity is upper bounded by $d \sum_{u \in V} d_u = O(dM)$.

PROOF OF LEMMA 3

For each $C_i \in \mathcal{C}$, decompose it into overlapped and non-overlapped parts, denoted by C_i^{ov} and C_i^{nov} . We have $C_i = C_i^{ov} \cup C_i^{nov}$ and $C_i^{ov} \cap C_i^{nov} = \emptyset$. Therefore, $|C_i| = |C_i^{ov}| + |C_i^{nov}|$.

Now, $\sum_{C_i \in \mathcal{C}} |C_i| = \sum_{C_i \in \mathcal{C}} (|C_i^{ov}| + |C_i^{nov}|) \leq N + \sum_{i < j} |C_i^{ov} \cap C_j^{nov}|$, where $N = \sum_{C_i \in \mathcal{C}} |C_i^{nov}| + |\bigcup_{C_i \in \mathcal{C}} |C_i^{ov}||$. For an upper bound of the second term, rewrite $\sum_{i < j} |C_i^{ov} \cap C_j^{nov}| \leq N + \sum_{|C_i \cap C_j| \geq 2} |C_i \cap C_j| \leq N(1 + \alpha)$ where $\alpha = \max\{|C_i \cap C_j| : |C_i \cap C_j| \geq 2\}$

Hence, $\sum_{C_i \in \mathcal{C}} |C_i| \leq N(2 + \alpha)$. Let N_0 be the number of raw communities, it follows that $N_0 \min\{|C_i|\} \leq \sum_{C_i \in \mathcal{C}} |C_i| \leq (2 + \alpha)N$. Since $\min\{|C_i|\} \geq 4$, we have $N_0 \leq \frac{(2+\alpha)N}{4} = O(N)$.

PROOF OF LEMMA 4

Prior to u joining to C_i , the internal density is $\Psi(C_i) = \frac{2|C_i^{in}|}{|C_i|(|C_i|-1)}$. Similarly, after u joining in C_i , the density function is $\Psi(C_i \cup \{u\}) = \frac{2|C_i^{in}| + 2d_{u,i}}{|C_i|(|C_i|+1)}$. Taking the difference between these two quantities gives $\Psi(C_i \cup \{u\}) > \Psi(C_i) \iff d_{u,i} > \frac{2|C_i^{in}|}{|C_i|-1}$. Moreover, u should also satisfy $\Psi(C_i \cup \{u\}) \geq \tau(C_i \cup \{u\})$, which in turn implies $d_{u,i} \geq f(|C_i| + 1) - |C_i^{in}|$. Therefore, $d_{u,i} > \max\{\frac{2|C_i^{in}|}{|C_i|-1}, f(|C_i| + 1) - |C_i^{in}|\}$.

PROOF OF LEMMA 5

Let C_1, C_2, \dots, C_k be the communities (including the newly formed ones) in \mathcal{C}_t that Alg. 4 assigns the new node u to. Note that in the optimal solution $OPT(u)_t$, the number of communities u belongs to should not exceed k since each C_i is also a candidate for $OPT(u)_t$ (of course, $OPT(u)_t$ could possibly rearrange nodes differently). Therefore, the optimal internal density gained is upper bounded by k . On the other hand, Alg. 4 makes sure that each community C_i that u joins in should have $\Psi(C_i) \geq \tau(C_i) \geq \tau(4)$ since $|C_i| \geq 4$. Thus, Alg. 4 will achieve at least $\tau(4) \times k \approx 0.74 \times \Psi(OPT(u)_t)$.

PROOF OF LEMMA 6

Suppose otherwise, that is C is divided into smaller parts C_1 and C_2 . Prior to the introduction of (u, v) , we have $\Psi(C) = \Psi(C_1 \cup C_2) \geq \tau(C) = \tau(C_1 \cup C_2)$. Now, when C_1 and C_2 are formed, they imply that $\Psi(C_1 \cup C_2 + (u, v)) < \tau(C_1 \cup C_2 + (u, v))$. Putting all together, we have $\tau(C_1 \cup C_2 + (u, v)) = \tau(C_1 \cup C_2) > \Psi(C_1 \cup C_2 + (u, v)) > \Psi(C) > \tau(C_1 \cup C_2)$, which raises a contradiction. Thus, the conclusion follows.