# Don't Be Greedy: Leveraging Community Structure to Find High Quality Seed Sets for Influence Maximization[*]

Rico Angell[1] and Grant Schoenebeck[1]

[1]University of Michigan

September 22, 2016

## Abstract

We consider the problem of maximizing the spread of influence in a social network by choosing a fixed number of initial seeds — a central problem in the study of network cascades. The majority of existing work on this problem, formally referred to as the *influence maximization problem*, is designed for submodular cascades. Despite the empirical evidence that many cascades are non-submodular, little work has been done focusing on non-submodular influence maximization.

We propose a new heuristic for solving the influence maximization problem and show via simulations on real-world and synthetic networks that our algorithm outputs more influential seed sets than the state-of-the-art greedy algorithm in many natural cases, with average improvements of 7% for submodular cascades, and 55% for non-submodular cascades. Our heuristic uses a dynamic programming approach on a hierarchical decomposition of the social network to leverage the relation between the spread of cascades and the community structure of social networks. We verify the importance of network structure by showing the quality of the hierarchical decomposition impacts the quality of seed set output by our algorithm. We also present "worst-case" theoretical results proving that in certain settings our algorithm outputs seed sets that are a factor of $\Theta(\sqrt{n})$ more influential than those of the greedy algorithm, where $n$ is the number of nodes in the network. Finally, we generalize our algorithm to a message passing version that can be used to find seed sets that have at least as much influence as the dynamic programming algorithms.

## 1   Introduction

A *cascade* is a fundamental social network process in which a number of nodes, or agents, start with some property that they then may spread to neighbors. Network structure has been shown relevant for a wide array of real world cascade processes including the adoption of products [7], farming technology [16], medical practices [15], participation in microfinancing [4], and the spread of information over social networks [28].

How to place a limited number of initial seeds, in order to maximize the spread of the resulting cascade, is a natural question known as INFLUENCEMAXIMIZATION [18, 39, 24, 25, 36]. This problem requires as input a network, a cascade process, and the number of initial seeds. For example, which students can most effectively be enrolled in an intervention to decrease student conflict at a school [38]?

To study INFLUENCEMAXIMIZATION, we first need to understand how cascades spread. While many cascade models have been proposed [2, 35, 44], they can be roughly divided into two categories: *submodular* and *non-submodular*.

1

In submodular cascade models, such as the Independent Cascade model defined in Section 2 [24, 25, 36], a node's marginal probability of becoming infected after a new neighbor is infected decreases when the number of previously infected neighbors increases [24]. In non-submodular cascade models the marginal probability of being infected may increase as more neighbors are infected. For example, in the Threshold model [22], each node has a threshold for the number of infected neighbors after which it too will become infected. If a node has a threshold of 2, then the first infected neighbor has zero marginal impact, but the second infected neighbor causes this node to become infected with probability 1. Unlike submodular cascades, non-submodular cascades require well-connected regions to spread [8].

For INFLUENCEMAXIMIZATION in submodular cascades, a straightforward greedy algorithm efficiently finds a seed set with influence at least a $(1 - 1/e)$ fraction of the optimal; but for general non-submodular casacdes, it is NP-hard even to approximate INFLUENCEMAXIMIZATION to within a $n^{1-\epsilon}$ factor of optimal [24].

Unfortunately, empirical research shows that most cascades are non-submodular [40, 3, 29], and in this case little is known about INFLUENCEMAXIMIZATION other than worst-case hardness. INFLUENCEMAXIMIZATION becomes qualitatively different in the non-submodular setting. In the submodular case, one should put as much distance between the $k$ initial adopters as possible, lest they erode each other's effectiveness. However, in the non-submodular case, it may be advantageous to place the initial adopters close together to create synergy and yield more adoptions. Thus, the intuition that it is better to saturate one market first, and then expand implicitly assumes non-submodular influence. However, this synergy renders the problem intractable.

As we will illustrate, greedy approaches can perform poorly in these settings. However, much of the work following Kempe et al. [24], which proposed the greedy algorithm, has attempted to make *greedy approaches* efficient and scalable [11, 12, 34, 14, 42]. New ideas seem necessary to design effective heuristics for non-submodular INFLUENCEMAXIMIZATION.

We observe that structural problems for networks—such as community detection—are also, in general NP-complete, but many efficient heuristics already exist [23, 13]. There are reasons to believe that such problems are not intractable in cases likely to arise in practice [1]. This work asks whether we can design heuristics for INFLUENCEMAXIMIZATION that work well for both submodular and non-submodular cascades, and what new algorithmic techniques might efficiently find hidden synergies necessary to maximize influence.

## 1.1 Contributions

We provide a new heuristic for solving INFLUENCEMAXIMIZATION designed to work for both submodular and non-submodular cascades. Our algorithm takes as input not only a network, but a hierarchical decomposition of the network. It then uses a dynamic programming technique to search for an influential seed set of nodes. We provide the following results concerning our algorithm:

1. We show theoretically that in certain cases, our algorithm outputs seed sets that are a factor of $\Theta(\sqrt{n})$ more influential than those of the state-of-the-art greedy algorithm, where $n$ is the number of nodes in the network. This illustrates the intuition behind our algorithm, as well as the poor performance of greedy.

2. We empirically compare our algorithm to the greedy algorithm via simulations on real-world and synthetic networks for a variety of cascade models. Our algorithm appears to do at least as well as greedy and substantially better for non-submodular cascades. Our algorithm achieves average improvements of 7% for submodular cascades and 55% for non-submodular cascades, performing 266% better in one exceptional case.

3. We verify the importance of network structure by showing that the quality of the hierarchical decomposition impacts the quality of our algorithm's output.

Finally, we define a generalization of our algorithm to a "message-passing" algorithm. Because this algorithm is generalization, it finds seeds sets of strictly higher quality than the dynamic programming algorithm, but empirically has a longer running time. We find that the results it returns

are only marginally better than the dynamic program. We believe that it may be more amenable to speeding up with heuristics, but leave such studies to future work.

## 1.2 Related Work

Following the work of Kempe et al. [24], which proposed the greedy algorithm, extensive work has constructed *efficient and scalable* algorithms and heuristics INFLUENCEMAXIMIZATION [11, 12, 37, 34, 14, 42].

The heuristic algorithms presented in [11, 12] rely on input parameters from the user that sacrifice accuracy for speed. The authors state that fine tuning the input parameters can make solving INFLUENCEMAXIMIZATION fast and accurate. Borgs et al. provably show fast running times when the influence function is the independent cascade model [6]. Tang et al. extend this work to provide an algorithm that maintains the same theoretical guarantees as the greedy algorithm presented in [24] and is efficient in practice [42]. Lucier et al. show how to parallelize (in a model based on Map Reduce) the subproblem of determining the influence of a particular seed [34]. Additional work has been done to speed up algorithms for solving INFLUENCEMAXIMIZATION by providing techniques to efficiently compute the total influence of a seed set [26, 14].

Leskovec et al. [30] consider the analogous problem of effectively placing sensors in a network in order to effectively detect an outbreak in the network. They present the algorithm CELF that uses a greedy approach, but leverages the submodularity of the cascade to reduce the amount of time it takes to evaluate the spread of the cascade. Moreover, CELF is built upon by the work in [19, 20], which present modifications to CELF to make an even more cost effective solution to INFLUENCEMAXIMIZATION. Nguyen and Zheng present an algorithm based on belief propagation for INFLUENCEMAXIMIZATION [37]. The algorithm works by systematically removing edges until the resulting graph is a tree, and then running a belief propagation algorithm on the scaled-down network. The authors of [37] show that the performance of their algorithm is not substantially worse than that of the greedy algorithm.

In contrast to the aforementioned work, our goal is not to deliver an algorithm that is more efficient and scalable, but rather to present an algorithm that finds higher quality seed sets. We are unaware of any other work that claims to substantially out-perform the greedy algorithm with respect to the quality of solution. Additionally, with the exception of [37], the prior work is based on a greedy-like approach. Our algorithm uses a dynamic programming framework, and is fundamentally different.

Other variations of INFLUENCEMAXIMIZATION have also been considered. The works [10, 32, 33] consider the problem where the time the cascade takes to spread is constrained. Seeman and Singer study the special case where only a subset of the nodes in the network are available to be infected [41]. INFLUENCEMAXIMIZATION has also been studied as a game between two different infectors [5, 21].

## 2 Preliminaries

A real function $f$ on sets is ***submodular*** if the marginal gain of from adding an element to a set $A$ is at least as large as the marginal gain from adding the same element to a superset $B$ of $A$. Formally, $f$ is submodular if for all $A$, $B$, $u$ where $A \subseteq B$ we have $f(A \cup \{u\}) - f(A) \geq f(B \cup \{u\}) - f(B)$.

**Cascade Model:** A cascade model is a triple $(G, F, S)$ where $G = (V, E)$ is an unweighted graph; $F = \{f_v : \{0,1\}^{|\Gamma(v)|} \to [0,1]\}_{v \in V}$ is a collection of ***local influence functions***, where $f_v$ takes in the set of infected neighbors of a node $v$, and produces a real value which encodes the "influence" of this set on $v$; and $S$ is the subset of the vertices that are initially infected. The cascade will proceed in rounds. In round 0, the set $S$ is infected and each of the remaining vertices is assigned a threshold value $\theta_v \in [0,1]$ drawn uniformly at random. At each subsequent round, a vertex $v$ becomes infected if and only if $f_v(T) \geq \theta_v$, where $T$ is the set of $v$'s infected neighbors. We will require $f_v$ to be monotone for each $v$.

We denote the **global influence function** as $\sigma(S)$ which is the *expected* total number of infected vertices due to the influence of the initial seed set $S$.

It can be shown that if $f_v$ is submodular for each $v$, then the global influence function $\sigma$ is submodular too [36]. Thus, we say that a model of cascade is **submodular** if $f_v$ is submodular for each $v$, and is **non-submodular** otherwise.

The same research that shows the $f$ usually fail to be submodular [40, 3, 29] shows that this submodularity fails in one particular way: the second adopting neighbor is, on average, more influential than the first; and that after this point, each subsequent adopting neighbor's marginal influence decreases. We call such functions **2-submodular**. Formally, $f$ is **2-submodular** if for all $A$, $B$, $A \subseteq B$, $|A|, |B| \geq 1$ and $u \notin A, B$, we have $f(A \cup \{u\}) - f(A) \geq f(B \cup \{u\}) - f(B)$; and for $v \neq u$, we have $f(\{u\}) - f(\emptyset) \leq f(\{u, v\}) - f(\{v\})$.

Any nonzero influence function $f_v$ can be turned into a 2-submodular function by sufficiently decreasing the value of $f_v(\cdot)$ on singleton sets.

For any local influence function $f_v$, we define the $q$-**deflated** version $f_v^{q-defl}$ of $f_v$ as follows:

$$f_v^{q-defl}(S) = \left\{ \begin{array}{ll} q \cdot f_v(S) & |S| = 1 \\ f_v(S) & \text{o.w.} \end{array} \right.$$

**Specific Cascade Models:** The two popular cascade models studied in the INFLUENCEMAXI-MIZATION literature are the Independent Cascade model (ICM) and the Linear Threshold model (LTM). In the **Independent Cascade model**, each newly infected node infects each currently uninfected neighbor in the subsequent round with some fixed probability $p$. Thus, for all $v$,

$$f_v^{ICM}(S) = 1 - (1 - p)^{|S|}.$$

In the **Linear Threshold model**, each node has a threshold $\theta_v \in [0, 1]$, each of $v$'s neighbors $u$ has influence $b_{u,v}$ on $v$ such that $\sum_{u \in \Gamma(v)} b_{u,v} \leq 1$, and $v$ becomes infected when the sum of the influences of the infected neighbors meets or surpasses $v$'s threshold.

We define the **Deflated Independent Cascade model** (DICM) which takes two parameters: $p, q \in [0, 1]$ to be the $q$-deflated version of the Independent Cascade model.

In the **S-Cascade model** (SCM) we have that

$$f_d^{SCM}(S) = \frac{\left(\frac{|S|}{2d}\right)^2}{\left(\frac{|S|}{2d}\right)^2 + \left(1 - \frac{|S|}{d}\right)^2}.$$

This is a modified version of the Tullock Cost function [43] with power 2.

We note that the Independent Cascade model and the Linear Threshold model are submodular, while the $q$-Deflated Independent Cascade model (for $q < 1 - p/2$) and S-Cascade are not. Figure 1 illustrates the local influence functions of the various cascade models.

**Synthetic Network Model:** Most existing synthetic models fail to have meaningful asymmetry between nodes, or significant community structure, or both. Therefore, we design our own synthetic network model. The **directed $(d, \ell, t)$-hierarchical network model** creates a random network on $2^d$ nodes as follows: We create an edge-weighted complete binary tree of depth $d$, each leaf representing a vertex of the graph. The weights are drawn i.i.d from a Binomial$(\ell, 1/2)$ distribution. Each node $v$ issues $t$ random edges, each generated via a random walk — illustrated below. Each random walk starts at $v$. At each step in the walk, an outgoing edge is chosen proportional to its weight (we disallow exiting the node along the same edge that the walk arrived at the node). The walk terminates when it arrives at a leaf node. If a terminating node is duplicated, we draw again, which keeps the graph simple.
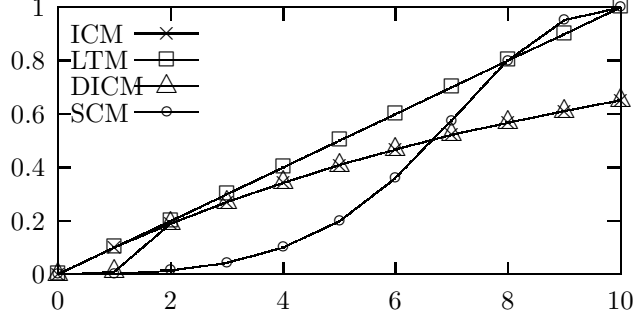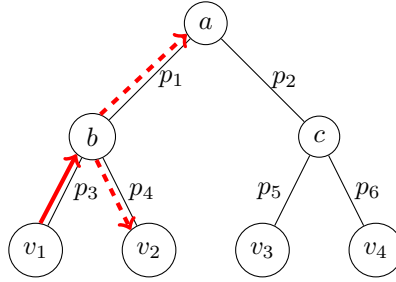
Figure 1: The local influence functions where the parameters of the ICM and DICM and the influence weights of LTM are all .1; and the vertex's degree in LTM and SCM is 10.



As $\ell$ grows larger, this approaches the hierarchical Kleinberg model [27]. But for moderately sized $\ell$, there is a non-trivial amount of asymmetry introduced into the graph — some subcommunities are more influential than others.

**Hierarchical Decomposition:** We define a **hierarchical decomposition** of a graph $G$ to be a rooted full binary tree $T = (V_T, E_T)$ where the leaves of $T$ correspond to the vertices of $G$. Let $r \in V_T$ be the root of $T$. For a tree node $v \in V_T$, define $\boldsymbol{T(v)}$ to be the subset of vertices in $G$ corresponding to the leaves of the subtree rooted at $v$. Let the **height** of $v \in V_T$ be defined as the length of the path to $v$'s deepest descendent.

We use the recently proposed cost function of Dasgupta [17] to evaluate the quality of a hierarchical decomposition. Let $\text{lca}_T(u, v)$ be the least common ancestor of $u, v \in V$ in the tree $T$. Then we define

$$Cost(T) = \sum_{\{u,v\} \in E} |T(\text{lca}(u, v))|$$

which sums the number of leaves in the smallest subtree containing each edge.

**Influence Maximization:** An **InfluenceMaximization Instance** consists of a graph $G = (V, E)$, an influence function $\sigma$, and an integer $k$. Given an InfluenceMaximization Instance, the goal is to find a set $S$ of $k$ nodes as to maximize $\sigma(S)$.

The **_greedy algorithm_** [24] for an INFLUENCEMAXIMIZATION Instance start with a tentative seed set $T = \emptyset$ and for $k$ rounds, simply adds $\arg\max_{v \in V} \sigma(T \cup \{v\})$ to $T$.

# 3  DPIM: Dynamic Programming Influence Maximization Algorithm

The Dynamic Programming Influence Maximization Algorithm (DPIM), formally specified in Algorithm 1, takes as input a graph $G$, and corresponding hierarchical decomposition $T$, an integer $k$,

5

and a global influence function $\sigma(\cdot)$ and outputs a subset of vertices $S \subseteq V$ such that $|S| = k$ and $S$ is a highly influential set of seeds. DPIM seeks to maximize the total influence of a fixed-sized seed set $S$ by performing dynamic programming upon $T$.

For each node $v \in T$, and each $i \in \{0, 1, \ldots, \min(|T(v)|, k)\}$, the algorithm stores $A[v, i]$, a choice of $i$ seeds in $T(v)$ which seeks to maximize the total influence in $G$. Starting at the leaves of the tree, and moving up level by level until reaching the root, DPIM processes each tree node. For each leaf node $v \in T$, we store $A[v, 0] = \emptyset$ and $A[v, 1] = \{v\}$. For each internal node $v \in T$, which has children $v_L$ and $v_R$, we set $A[v, i] = A[v_L, j] \cup A[v_R, i - j]$ where $j \in \{0, 1, \ldots, i\}$ is selected as to maximize $\sigma(A[v_L, j] \cup A[v_R, i - j])$.

---

**ALGORITHM 1:** DPIM: Dynamic Programming Influence Maximization Algorithm

---

**Input**: $G = (V, E), T = (V_T, E_T), \sigma(\cdot), k$
**Output**: $S \subset V$ such that $|S| = k$
Let $A[\cdot, \cdot] = V_T \times [k] \to 2^V$, such that $A[v_T, j]$ stores a choice of $j$ seeds in $T(v_T)$.
Let $h$ be the height $T$.
**for** *each height* $i = 0, 1, \ldots, h$ **do**
    **for** *each node* $v \in V_T$ *with height* $i$ **do**
        **if** $i = 0$ **then**
            $A[v, 0] = \emptyset$
            $A[v, 1] = \{v\}$
        **else**
            Let $v_L, v_R$ be the left and right children of $v$, respectively.
            **for** *each* $i = 0, 1, \ldots, \min\{|T(v)|, k\}$ **do**
                $j = \underset{j \in \{0, 1, \ldots, i\}}{\arg\max} \ \sigma(A[v_L, j] \cup A[v_R, i - j])$
                $A[v, i] = A[v_L, j] \cup A[v_R, i - j]$
            **end**
        **end**
    **end**
**end**
**return** $A[r, k]$

---

The analysis of the running time for DPIM is straightforward.

**Theorem 1.** *Given a graph $G = (V, E)$ with $|V| = n, |E| = m$, fixed positive integers $k, r$, and a hierarchical decomposition $T$,* DPIM *calls the $\sigma(\cdot)$ oracle $O(nk^2)$ times.*

*Proof.* Observe that, for each node in $T$, DPIM makes $O(k^2)$ queries to $\sigma(\cdot)$. The number of nodes in $T$ is exactly $2n - 1$. Hence, the number of oracle calls in DPIM is $O(nk^2)$. $\square$

Note that this is a factor of $k$ more than the greedy algorithm, which requires only $O(nk)$ calls to the oracle. The execution of a single query to $\sigma(\cdot)$ can be approximated by repeatedly, $r$ times, simulating the cascade process and returning the average number of infected vertices. This can be done in time $O(mr)$ because simulating the cascade requires at most simulating the cascade on each edge in $G$. However, there are often techniques to speed up the oracle access beyond simply running the cascade [6], but they are beyond the scope of this work.

## 3.1   Motivation:

When DPIM analyzes a node $v_T \in V_T$, and the two subtrees of $v_T$ are disjoint. If the inputs from the subtrees are optimal, then the output at of the tree will be optimal as well. That is, DPIM correctly decides how many nodes to allocate to each disconnected component. Although this extreme case of disconnected subtrees may be rare, we expect the performance to degrade gently as the number of inter-community edges increases. Thus, DPIM finds the globally influential vertices by combining the knowledge it gained about each of the subtrees. Because we expect that the community structure

is important to how the cascade spreads [9], if the subtrees are dense "community-like" structures, we expect our algorithm to do well. Moreover, Community structure is hierarchical in real social networks [13], which dynamic programming can exploit. However, in the same setting where you have a collection of disjoint graphs, greedy may be suboptimal (especially if the cascade is non-submodular). It may spread the seeds among different communities, where the optimal thing to do is saturate one community.

We illustrate this is two ways: first empirically on a random graph, and then with a theorem.

### 3.1.1   Empirical Motivating Illustration



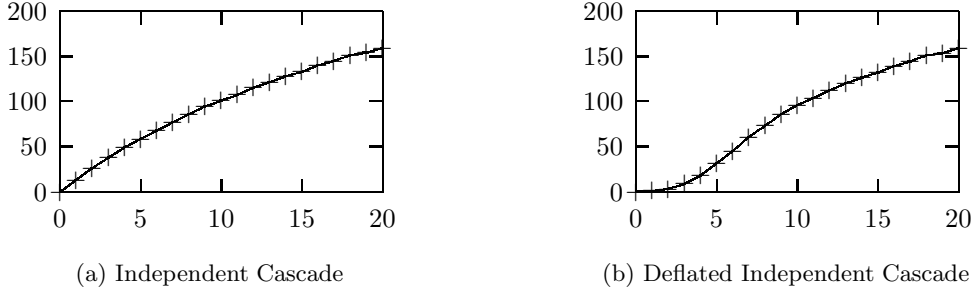| (a) Independent Cascade | (b) Deflated Independent Cascade |

Figure 2:  Global influence of a random seed set on both the Independent Cascade and 0.1-Deflated Independent Cascade on a random graph with 10,000 vertices and 50,000 edges. The x-axis is the size of the initial seed set and the y-axis is the total number of infections at the end to the simulation. Notice the Global non-submodularity of Deflated Independent Cascade.

Consider Figure 2 which illustrates both an Independent Cascade and 0.1-Deflated Independent Cascade on a random graph. Notice that for this very natural graph, a small change in the local influence function creates a large change the global influence function. This is especially true when only infecting a few seeds. After infecting 10 seeds, there is barely any difference (101.4 versus 96.1). However, in the Independent Cascade, the marginal impact of adding one seed, 13.3, is always greater than the average marginal impact of the first ten seeds, 10.1. But in the Deflated Independent Cascade, the marginal impact of adding one seed, 1.3, is not a good predictor of the average marginal impact of the first ten seeds, 9.6, which is over seven times that of adding just one node! Adding 10 seeds in this graph can create a synergy that is not available when just considering one seed.

The greedy algorithm, which only ever considers one node at a time, may have trouble finding these synergies. In particular, if this random graph is only part of a larger network where there are several other regions where adding one seed yields 2 infections in expectation, the greedy algorithm will never explore adding nodes to this section of the network where synergies are possible.

In contrast our DPIM specifically considers what happens if it adds a large number of nodes to a particular "community-like" region. These dense regions are exactly where we expect synergies to occur.

The greedy algorithm can make bad initial choices that are hard to mitigate later. If the greedy algorithm adds too many seeds to graph locations that seem initially promising, but fail to provide synergies, it may not be able to catch up to the DPIM even if it sometimes happens upon these synergies. In contract, DPIM can later reallocate nodes away from a subtree that is not performing well.

### 3.1.2   Theoretical Comparison in "Worst-case" scenario

In this section, we illustrate a particular setting where Algorithm 1 provably has considerably better performance than the greedy algorithm proposed in [24]. This makes the aforementioned intuition rigorous.

**Theorem 2.** *There exist* INFLUENCEMAXIMIZATION *instances in which the influence of the seed set output by Algorithm 1 is a factor of* $\Theta(\sqrt{N})$ *larger than the influence of the seed set output by the greedy algorithm, where $N$ is the number of verices in the graph.*

*Proof.* Consider the graph $G = (V, E)$ with $N = 2n^2 + n + 2$ vertices consisting of the following three connected components:

- two stars of size $n^2 + 1$, and

- one clique of size $n$.

Consider the complex contagion model with influence function $F = \{f_v\}$ such that

- $f_v(S_v) = \frac{1}{n^2}$ if $|S_v| = 1$, and

- $f_v(S_v) = 1$ if $|S_v| \geq 2$,

In the INFLUENCEMAXIMIZATION instance with parameter $k = 2$, we evaluate the performance of both the greedy algorithm and Algorithm 1.

In the greedy algorithm, the vertex with maximum marginal influence is selected. The two obvious potential targets are the either one of the centers of the two stars, or any vertex in the clique. If the center vertex of a star is infected, each of the remaining $n^2$ vertices will be infected with probability $\frac{1}{n^2}$, so the expected number of the infected vertices is $1 + n^2 \times \frac{1}{n^2} = 2$. On the other hand, if we infect any vertex in the clique, the infected vertex will not infect any other vertices in the clique with probability $\left(1 - \frac{1}{n^2}\right)^{n-1}$, and it will infect at least one other vertex with probability $1 - \left(1 - \frac{1}{n^2}\right)^{n-1}$ in which case all the remaining vertices in the clique will be infected. Therefore, in the case we choose a vertex in the clique, the expected number of vertices infected is

$$\mathbb{E}[\sigma(S)] = 1 \times \left(1 - \frac{1}{n^2}\right)^{n-1} + n \times \left(1 - \left(1 - \frac{1}{n^2}\right)^{n-1}\right)$$

$$= 1 + (n-1)\left(1 - \left(1 - \frac{1}{n^2}\right)^{n-1}\right)$$

$$= 1 + (n-1)\frac{1}{n^2}\left(\sum_{i=0}^{n-2}\left(1 - \frac{1}{n^2}\right)^i\right)$$

$$< 1 + (n-1)\frac{1}{n^2} \cdot (n-1)$$

$$< 2.$$

Thus, infecting the center of a star has higher marginal influence, so the greedy algorithm will choose the two centers of the two stars to infect, and an expected total number of 4 vertices will be infected.

On the other hand, assuming that the hierarchical decomposition $T$ contains a subtree $T(v)$ consisting only of the clique (which it should as it is disconnected from the rest of the graph): Algorithm 1 will choose 2 vertices from the clique causing a total number of $n$ infected vertices, which is also the optimal strategy. To see this, note that Algorithm 1 will process $v$, and at this point will compute $\sigma(S)$ where $S$ contains two nodes from the clique. By the nature of Algorithm 1, this solution, or a better one, will be considered (inductively) for each vertex that has $v$ as a descendent. Thus the output seed set is at least as influential as placing two seeds in the clique, which has influence $n = \Theta(\sqrt{N})$. Consequently, the performance of Algorithm 1 is better than the performance of the greedy algorithm by a factor of

$$\frac{n}{4} = \Theta(n) = \Theta(\sqrt{N}).$$

$\square$

Notice that the structure of the graph in the example above is related to some of the common structures found in social networks, which contain many different communities with different internal structures.

# 4 Experimental Results

## 4.1 Experimental Setup

We execute DPIM and the greedy algorithm from [24] on a variety of networks and cascades to test the relative quality of solutions.

### 4.1.1 Cascade Models

We adopt the two common submodular cascade models from the literature: the linear threshold model and the independent cascade model, defined in Section 2, and two non-submodular cascades:

I) *Independent Cascade* (IC): We uniformly assign the probability $p = 1\%$, thus $v$ with $\ell$ infected neighbors is infected with probability $1 - (0.99)^\ell$.

II) *Linear Threshold* (LT): For each node $v$, we assign each of $u \in \Gamma(v)$ to have $1/|\Gamma(v)|$ influence on $v$.

III) *Deflated Independent Cascade* (DIC): We uniformly assign the probability $p = 1\%$, thus $v$ with $\ell = 1$ infected neighbors is infected with probabilty 0.001 and with $\ell \geq 2$ infected neighbors is infected with probability $1 - (0.99)^\ell$.

IV) *S-Cascade model* (SCM): The influence on any given node $v$ is

$$\frac{(x/2)^2}{(x/2)^2 + (1-x)^2},$$

where $x$ is the fraction of $v$'s neighbors that are infected.

Both of these algorithms require access to an oracle for $\sigma(\cdot)$, which is also required to evaluate the effectiveness of the algorithms. To implement this oracle, we simulate the cascade 100 times, resampling the randomness for the cascade each time (using pseudorandomness from the standard C++ library) and return the average number of infections.

### 4.1.2 Networks

We use two real-world networks (from [31]) and two synthetic networks, summarized in Table 1. In the arXiv collaboration network (ca-GrQc), the vertices are authors of e-print scientific articles and edges represent coauthorship relations. The ego-Facebook network is largest such network provided by [31]. This network denotes the facebook friendship ties from a single person's (ego's) set of friends. The ego vertex has been removed. Furthermore, we generate two synthetic networks by first sampling from directed $(d, \ell, t)$-hierarchical network model using parameters $(10, 50, 50)$ and $(11, 50, 50)$ (synthetic-1 & synthetic-2, resp.), and then making the graph simple and undirected in the natural way.

| Name | Nodes | Edges |
|---|---|---|
| synthetic-1 | 1,024 | 51,200 |
| synthetic-2 | 2,048 | 102,400 |
| ca-GrQc | 5,276 | 28,827 |
| ego-Facebook | 1,034 | 53,498 |

Table 1: Networks used to evaluate the effectiveness of our algorithm.

### 4.1.3 Algorithms for Hierarchical Decomposition

Lastly, in order to evaluate our algorithm, we present 4 algorithms for generating a hierarchical decomposition of any network. The algorithms we used in our simulations are implemented as follows:

I) *Random Pair*: Each node starts in its own partition, and partitions are joined randomly until all of the nodes are contained in one partition.

II) *Random Edge*: Each node starts in its own partition, and partitions are joined by contracting a random edge between partitions. If no edges remain between the partitions, partitions are merged randomly until all of the nodes are contained in one partition.

III) *Jaccard Similarity*: Each node starts in its own partition, and pairs of partitions $(A, B)$, for $A, B \subset V$ are joined based on which pair maximizes

$$\frac{|\Gamma(A) \cap \Gamma(B)|}{|\Gamma(A) \cup \Gamma(B)|},$$

where $\Gamma(X \subset V) = \bigcup_{v \in X} \Gamma(v)$.

IV) *METIS-based*: The whole network starts as one partition; using METIS [23], partitions are recursively divided into two partitions until each partition contains only a single node.

## 4.2 Algorithm Evaluation

### 4.2.1 Performance of DPIM

The results of the simulations we ran are shown in Figure 3. For each execution of DPIM, we used the METIS-based hierarchical decomposition algorithm to construct a hierarchical decomposition of the network.

Considering cascades across all four networks with seed set size 20, DPIM increases influence on average by 8% for ICM, 6% for LTM, 22% for DICM, and 88% for SCM. Surprisingly, DPIM performs marginally better than the greedy algorithm even for submodular cascades. As predicted, when the cascade is non-submodular, DPIM outperforms the greedy algorithm by a significant amount. However, gains were more impressive for synthetic-2, ca-GrQc, and ego-Facebook — including an 266% increase in influence for the SCM cascade on the ego-Facebook network — than for synthetic-1, where we see only marginal improvement even when the cascade is non-submodular. Table 3 contains the approximated expected total influence values for each simulations rounded to the nearest integer.

In addition, we present the running times of each simulation (both DPIM and Greedy) in Table 2. Theoretically, the greedy algorithm queries $\sigma(\cdot)$ $O(nk)$ times, and DPIM queries $\sigma(\cdot)$ $O(nk^2)$ times. Despite this, the empirical results show that DPIM is roughly a small constant factor slower than the greedy algorithm, and occasionally much faster.

|  | synthetic-1 | | synthetic-2 | | ca-GrQc | | ego-Facebook | |
|---|---|---|---|---|---|---|---|---|
|  | Greedy | DPIM | Greedy | DPIM | Greedy | DPIM | Greedy | DPIM |
| **ICM** | 23,720 | 90,945 | 49,885 | 162,892 | 39,641 | 350,006 | 275,516 | 758,843 |
| **LTM** | 24,121 | 92,786 | 52,018 | 165,369 | 25,626 | 56,156 | 308,637 | 83,584 |
| **DICM** | 18,536 | 69,993 | 42,176 | 134,744 | 6,108 | 26,672 | 266,859 | 641,962 |
| **SCM** | 13,676 | 36,765 | 23,749 | 69,737 | 10,274 | 46,196 | 38,279 | 49,892 |

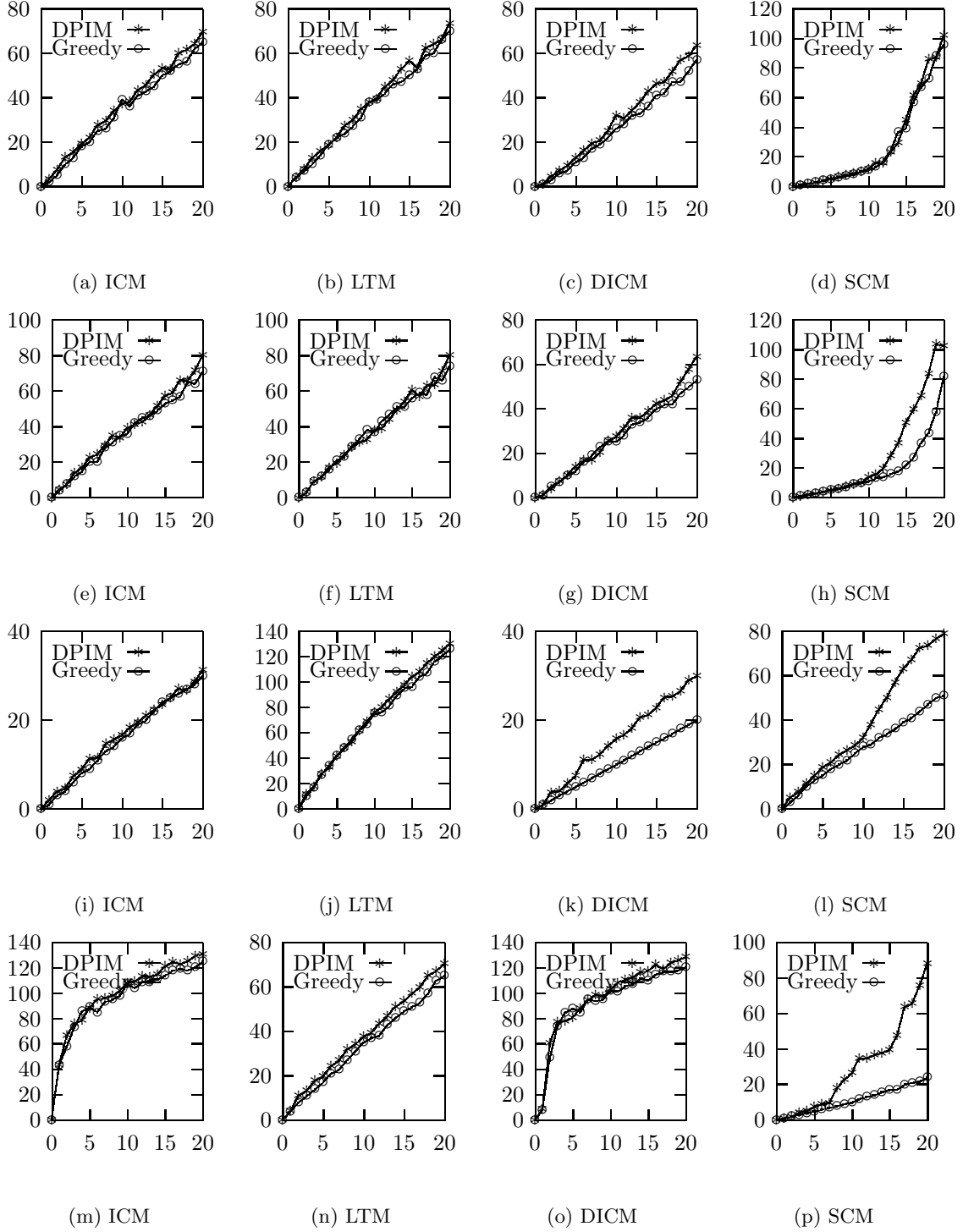Table 2: Running times in seconds of each of the simulations.

Figure 3: **Comparison of performance:** DPIM **vs. Greedy**. The rows from top to bottom correspond to synthetic-1, synthetic-2, ca-GrQc, and ego-Facebook, respectively. For each plot, the x-axis is $k$ and the y-axis is the number of total infections at the end of the cascade.

|  | synthetic-1 | | synthetic-2 | | ca-GrQc | | ego-Facebook | |
|---|---|---|---|---|---|---|---|---|
|  | Greedy | DPIM | Greedy | DPIM | Greedy | DPIM | Greedy | DPIM |
| **ICM** | 65 | 70 | 71 | 80 | 30 | 31 | 125 | 131 |
| **LTM** | 70 | 74 | 74 | 81 | 126 | 130 | 65 | 70 |
| **DICM** | 57 | 64 | 53 | 64 | 20 | 30 | 121 | 129 |
| **SCM** | 96 | 102 | 82 | 103 | 51 | 79 | 24 | 88 |

Table 3: Expected total influence of the final seed sets of size 20 chosen by both algorithms for each of the simulations (rounded to the nearest integer).

### 4.2.2  Comparison of Hierarchical Decomposition Algorithms

For each network, we tested how DPIM performed over the various hierarchical decomposition algorithms with SCM as the cascade model. In addition, we recorded the cost of each hierarchical decomposition (see Section 2 for details). The details of the performance of our algorithm are shown in Figure 4 and the cost of each of the hierarchical decompositions is noted in Table 4. Observe that DPIM performs significantly better with the hierarchical decompositions that have lower costs.

Despite performing better when the hierarchical decomposition better represents the community structure of the network, there seem to be diminishing returns beyond a certain cost for each network. The METIS-based algorithm and the Jaccard Similarity algorithm produce hierarchical decompositions that differ in cost by a relatively large amount, but the difference in performance of our algorithm between the two different hierarchical decompositions is small.

Since a lower cost valuation implies that the hierarchical decomposition better represents the structure of the network, this result provides evidence that our algorithm is leveraging the community structure provided by the hierarchical decomposition.

|  | synthetic-1 | synthetic-2 | ca-GrQc | ego-Facebook |
|---|---|---|---|---|
| **Random Pair** | 35,063,945 | 139,590,514, | 99,792,540 | 36,387,204 |
| **Random Edge** | 34,084,502 | 133,759,197 | 57,800,479 | 31,092,978 |
| **Jaccard Similarity** | 22,508,941 | 88,531,982 | 27,695,152 | 13,455,920 |
| **METIS-based** | 5,267,974 | 12,185,692 | 18,425,481 | 8,968,620 |

Table 4: Cost of the hierarchical decompositions produced by each algorithm for each network.



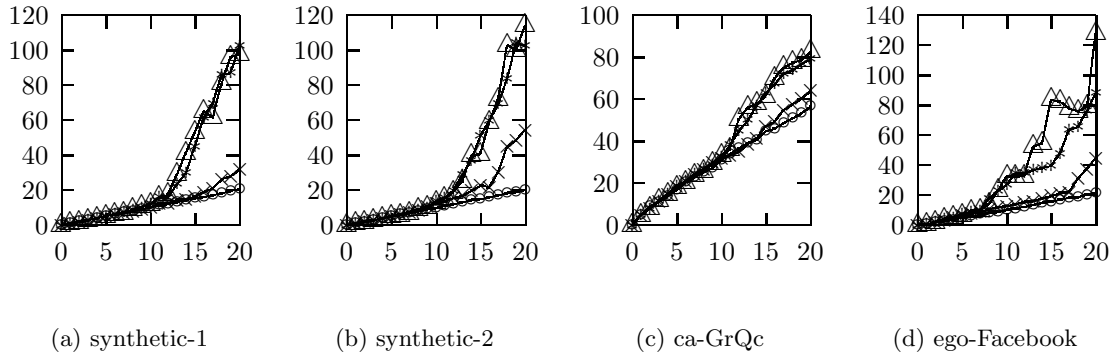(a) synthetic-1          (b) synthetic-2          (c) ca-GrQc          (d) ego-Facebook

Figure 4: **Comparison of Hierarchical Decomposition Algorithms**. For each plot, the x-axis is $k$ and the y-axis is the number of total infections at the end of the cascade. In each plot, the circle is Random Pair, the cross is Random Edge, the triangle is Jaccard Similarity, and the asterisk is METIS-based. Each seed set was choosen by DPIM with each of the respective hierarchical decomposition algorithm as input, and the cascade used was SCM.

# 5  Message Passing Algorithm

In this section, we describe a Message Passing Algorithm (MPA), formally presented in Algorithm 3, which is a generalization of DPIM. This algorithm provides a new perspective on DPIM, and can be made to perform better at the cost of additional run time.

**Directional Subtrees and Their Recursive Decomposition:**  First, we must present notation and ideas that we will use to formally describe MPA. For a node $v_T \in V_T$, we define $L(v_T), R(v_T)$, and $U(v_T)$ to be the left child, right child, and parent (up), respectively. We use $D$ as a placeholder for an element in $\{L, R, U\}$ (left, right, up directions), so that $D(v_T)$ refers to either $L(v_T), R(v_T)$, or $U(v_T)$. Let $\mathbb{D} = \{(L, R), (L, U), (R, U)\}$.

Furthermore, for node, $v_T \in V_T$ and a direction $D \in \{L, R, U\}$ we define the backward direction $B_D^{(v_T)} \in \{L, R, U\}$ to be the direction one takes from $D(v_T)$ in order to return to $v_T$ so that $B_D^{(v_T)}(D(v_T)) = v_T$. We naturally define $\neg B_D^{(v_T)} \in \mathbb{D}$ such that neither of the coordinates of $\neg B_D^{(v_T)}$ are $B_D^{(v_T)}$.

Lastly, for each node $v_T \in V_T$ we define three *directional subtrees*: *left* — $T_L(v_T) = T(L(v_T))$, *right* — $T_R(v_T) = T(R(v_T))$, and *up* — $T_U(v_T) = V \setminus (T(v_T))$; which partition the vertices of $G$ because $T(v_T) = T_L(v_T) \cup T(R(v_T))$. A key observation is that each of these three subtrees, can be decomposed into smaller directions subtrees of $L(v_T), R(v_T)$, and $U(v_T)$ respectively. For example, for $D \in \{L, R\}$ we have that $D(v_T) = T_L(D(v_T)) \cup T_R(D(v_T))$. And while $v_T$ is not the root of $V_T$, we have $T_U(v_T) = \bigcup_{D \in \neg B_U^{v_T}} T_D(U(v_T))$.

**Data Structure:**  MPA will sequentially, locally update a data-structure, which we describe now. Let $A : V_T \times \mathbb{D} \times [k] \to [k] \times [k]$ where both output coordinates of $A(v_T, (D_1, D_2), \ell)$ are non-negative integers that sum to at most $\ell$. Intuitively, $A(v_T, (D_1, D_2), \ell)$ answers: how should $\ell$ seeds should be split between $T_{D_1}(v_T)$ and $T_{D_2}(v_T)$? The answer is that we should put $A(v_T, (D_1, D_2), \ell)_i$ seeds in $T_{D_i}(v_T)$ for each $i \in \{1, 2\}$. For the root vertex $r \in V_T$, which has no parent, we insist that $A(r, (D, U), \ell)_2 = 0$ for $D \in \{L, R\}$.

**Obtaining a Seed Set:**  Once we have the data-structure $A$, we define a function that answers the more natural question: Let $C : V_T \times \mathbb{D} \times [k] \to 2^V$ where $C(v_T, (D_1, D_2), \ell)$ is a subset of $T_{D_1}(v_T) \cup T_{D_2}(v_T)$ of at most $\ell$ vertices. Intuitively, the question that $C$ answers is: where should we allocate $\ell$ seed vertices in $T_{D_1}(v_T) \cup T_{D_2}(v_T)$? Given an instance of $A$, $C(v_T, (D_1, D_2), \ell)$ recursively follows the advice of $A$ in placing seeds in $T_{D_1}(v_T)$ and $T_{D_2}(v_T)$. Intuitively, we can do this using $A$ and the recursive tree decomposition. Formally, we divide the definition of $C(v_T, (D_1, D_2), \ell)$ into two cases:

I) $(D_1, D_2) = (L, R)$: If $v_T$ is a leaf node, then

$$C(v_T, (L, R), \ell) = \left\{ \begin{array}{ll} T(v_T) & \ell \geq 1 \\ \emptyset & \ell = 0 \end{array} \right. .$$

Otherwise,

$$\begin{aligned} C(v_T, (L, R), \ell) = &C(L(v_T), (L, R), A[v_T, (L, R), \ell]_1) \cup \\ &C(R(v_T), (L, R), A[v_T, (L, R), \ell]_2). \end{aligned}$$

II) $(D_1, D_2) = (D, U)$ for $D \in \{L, R\}$: If $v_T = r$, the root of $T$, then

$$C(v_T, (D, U), \ell) = C(D(v_T), (L, R), \ell).$$

If $v_T$ is a leaf node, then it is undefined (or just $\emptyset$). Otherwise,

$$\begin{aligned} C(v_T, (D, U), l) = &C(D(v_T), (L, R), A(v_T, (D, U), l)_1) \cup \\ &C(U(v_T), \neg B_D^{(v_T)}, A(v_T, (D, U), l)_2). \end{aligned}$$

Note that $C(r, (L, R), k)$ denotes where to place $k$ initial seeds over all of $V$.

**Local Update Subroutine:** Now that we have defined $A$ and shown how it defines initial seeds sets, we will define a way to locally update $A$ to improve the allocation. Intuitively, we update $A(v_T, (D_1, D_2), \ell)$ by trying all possible outputs whose coordinates sum to $\ell$ and returning the best. Let $D_3$ be the sole element of $\{L, R, U\} \setminus \{D_1, D_2\}$. We will first allocate $k - \ell$ vertices to $T_{D_3}(v_T)$ according to $C$. Note that this tree can be decomposed into two directional subtrees of $D_3(v_T)$ upon which $C$ is defined. Next, we try the different divisions between $T_{D_1}(v_T)$ and $T_{D_2}(v_T)$ for the remaining $\ell$ vertices. This procedure is formally defined in Algorithm 2.

---

**ALGORITHM 2:** update: Local Update Subroutine

**Input**: $I = (G = (V, E), T = (V_T, E_T), \sigma(\cdot), k), A, v_T, (D_1, D_2) \in \mathbb{D}$
**Output**: An updated instance of $A$
Let $D_3$ be the sole element of $\{L, R, U\} \setminus \{D_1, D_2\}$.
**for** *each* $i = 0, 1, \ldots, k$ **do**
 $\quad j = $
 $\quad \underset{j \in \{0, 1, \ldots, i\}}{\arg\max} \ \sigma\left( C\left(D_1(v_T), \neg R_{D_1}^{(v_T)}, j\right) \cup C\left(D_2(v_T), \neg R_{D_2}^{(v_T)}, i - j\right) \cup C\left(D_3(v_T), \neg R_{D_3}^{(v_T)}, k - i\right) \right)$
 $\quad A(v_T, (D_1, D_2), i) := (j, i - j)$
**end**
**return** $A$

---

**Generalizing DPIM:** Given that we have a local update, we need only define: (a) an initial configuration of $A$, (b) a sequence of updates, and (c) a terminating condition (if the sequence is infinite). For example, our DPIM is a special case where every entry in $A$ is $(0, 0)$ initially, and then we update $A(v_T, (L, R), \cdot)$ by going up the tree.

---

**ALGORITHM 3:** MPA: Message Passing Algorithm

**Input**: $I = (G = (V, E), T = (V_T, E_T), \sigma(\cdot), k)$
**Output**: $S$, a set of $k$ nodes of $V$
Let $A(\cdot) = (0, 0)$
**for** *each height* $i = 1, 2, \ldots h$ **do**
 $\quad$**for** *each node* $v_T \in V_T$ *with height* $i$ **do**
 $\quad\quad$update$(I, A, v_T, (L, R))$;
 $\quad$**end**
**end**
Let $S, S' := C(r_T, (L, R), k)$;
**do**
 $\quad$S := S';
 $\quad$**for** *each height* $i = h - 1, \ldots, 1$ **do**
 $\quad\quad$**for** *each node* $v_T \in V_T$ *with height* $i$ **do**
 $\quad\quad\quad$update$(I, A, v_T, (L, U))$;
 $\quad\quad\quad$update$(I, A, v_T, (R, U))$;
 $\quad\quad$**end**
 $\quad$**end**
 $\quad$**for** *each height* $i = 1, 2, \ldots h$ **do**
 $\quad\quad$**for** *each node* $v_T \in V_T$ *with height* $i$ **do**
 $\quad\quad\quad$update$(I, A, v_T, (L, R))$;
 $\quad\quad$**end**
 $\quad$**end**
 $\quad$$S' = C(r_T, (L, R), k)$;
**while** $\sigma(S') > \sigma(S)$;
**return** $C(r_T, (L, R), k)$

---

In Algorithm 3, we give a possible schedule where we go up and down the tree until no additional improvements are found. Given a graph, decomposition, influence function, $k$, and current allocation $A$, we run DPIM to initialize $A$. Then, we evaluate each node on each level of $T$ down and up the tree, finishing at the root node, checking whether there has been any improvements in the influence of the seed sets, and repeating until there are no more improvements. Note that this algorithm must terminate, since requires improvement in each round and there are only a finite number of possible configurations.

**Intuition:** When DPIM sets $A(v_T, (L, R), \ell)$ for $\ell < k$, it does not know where the other seeds (outside $T(v_T)$ will be placed. Thus, a better decision may be available with this added information. How should an algorithm decide how to place elements in $V \setminus T(v_T)$? If we only know $A(\cdot, (L, R), \cdot)$, this does not appear to be enough information. At first blush, we must know how many seeds to allocate to the other subtree of $U(v_T)$, and how many to send further up the tree. This is exactly what $A(U(v_T), \neg B_U^{v_T}, \cdot)$ tells us! Of course, all the $A(v_T, \cdot, \cdot)$'s seem interdependent. Thus, we create a local update algorithm to gradually refine each of them.

## 5.1 Message Passing Algorithm Evaluation

We evaluate MPA by comparing it's effectiveness against DPIM for all four cascades. Due to the *significant time complexity* of MPA, we were only able to simulate MPA on two small networks — we sampled from the directed $(d, l, t)$-hierarchical network model using parameters $(8, 15, 15)$ and have a smaller ego-network from Facebook that has 150 nodes and 3,386 edges. The results from these simulations can be seen in Figures 5 and 6. Notice that we have plots for $k = 5, 10, 15, 20$, since MPA optimizes for the $k$ exactly and does not guarentee that is is finding the best seeds for seeds sets of size $1, \ldots, k - 1$. This way we get to see how effective MPA is at multiple sizes of seed sets.

As can be seen in Figures 5 and 6, the improvement that MPA has over DPIM is marginal for the ICM, LTM, and DICM cascades (averageing improvements of 3%, 1%, 3%, respectively). On the other hand, MPA finds a seed set that has significantly more influence when the cascade is SCM (with an average improvement of 16%). It seems like for an algorithm to find a high-quality seed set for the SCM cascade, the algorithm must be able to find intercommunity synergies which are the only way that SCM will propogate throughout networks.

# 6 Conclusions

We have given a heuristic which exploits the hierarchical community structure of networks to find influential seed sets. We have shown, using both real-world and synthetic networks, that our algorithm outperforms the state of the art, with large gains for non-submodular influence maximization. We have also exhibited "worst-case" theoretical instances where our algorithm produces sets that are $\Theta(\sqrt{n})$ more influential. Lastly, we have generalized our heuristic to a message passing algorithm.

One possible direction of future exploration is to try additional hierarchical decomposition techniques. Interestingly, DPIM can be seen as a way to *test* hierarchical decomposition techniques. Decompositions that perform better are intuitively finding a better decomposition.

DPIM is typically not as fast as naive greedy, and to be useful in practice, it would greatly help if it were more scalable. We believe that this will prove to be the case. For example, we could stop the recursion before exploring the entire hierarchical decomposition. We might stop dividing if a subtree does not appear to have any additional community structure, and then run a heuristic (such a degree or greedy) to process the rest of the subtree. The intuition here is that dynamic programming works best where the network has strong community structure, so where no structure exists, it may not provide much added benefit. Additionally, the same techniques that have made the greedy algorithm more scalable might be adopted to our dynamic programming and message passing frameworks.

(a) ICM      (b) LTM      (c) DICM      (d) SCM

(e) ICM      (f) LTM      (g) DICM      (h) SCM

(i) ICM      (j) LTM      (k) DICM      (l) SCM

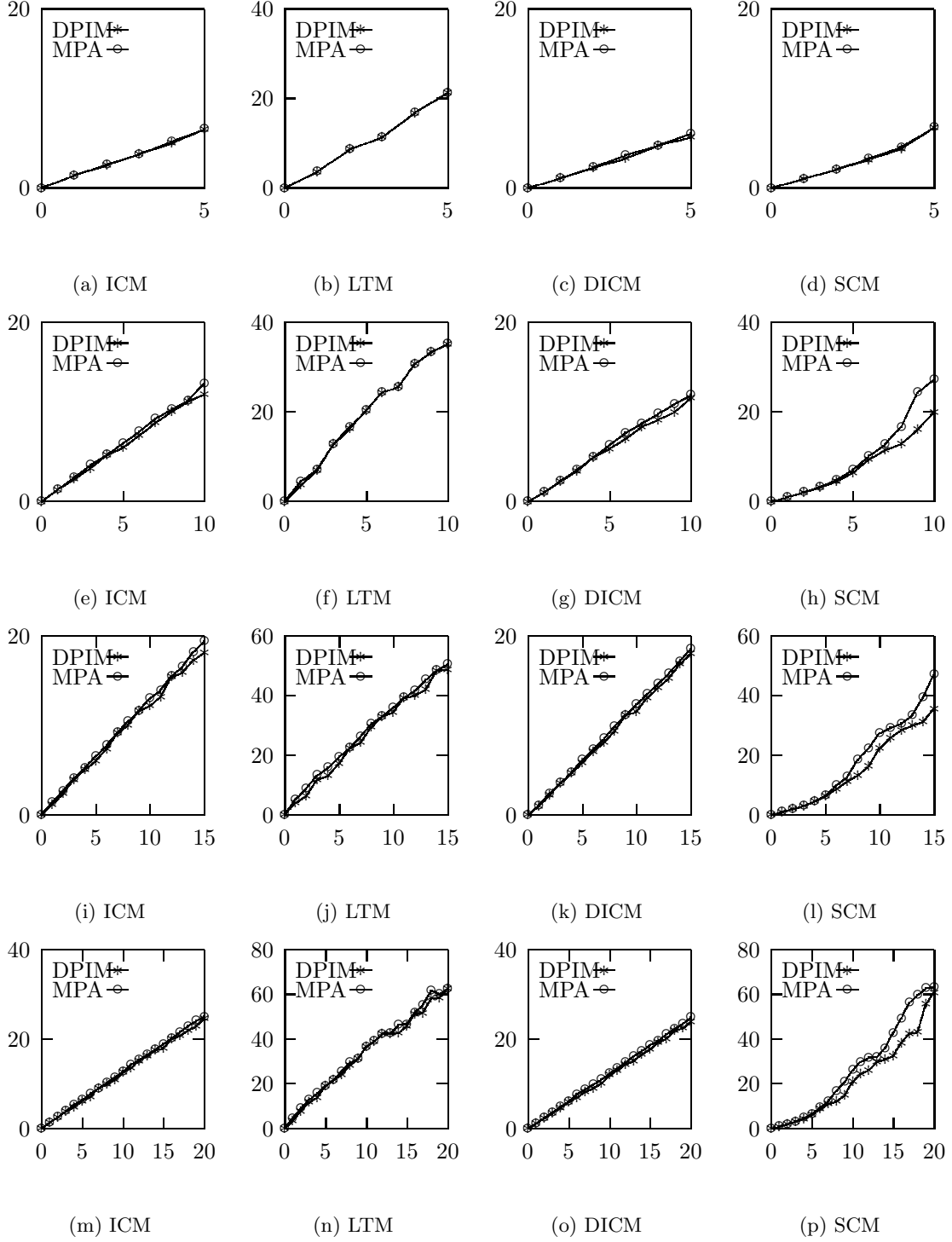(m) ICM      (n) LTM      (o) DICM      (p) SCM

Figure 5: **Evaluation of MPA on a synthetic network.** We evaluate the effectiveness of MPA by comparing the results to DPIM. The network was a sample from the directed $(d, l, t)$- hierarchical network model using parameters (8,15,15). For each plot, the x-axis is $k$ and the y-axis is the expected total influence from the seed set chosen by each respective algorithm.
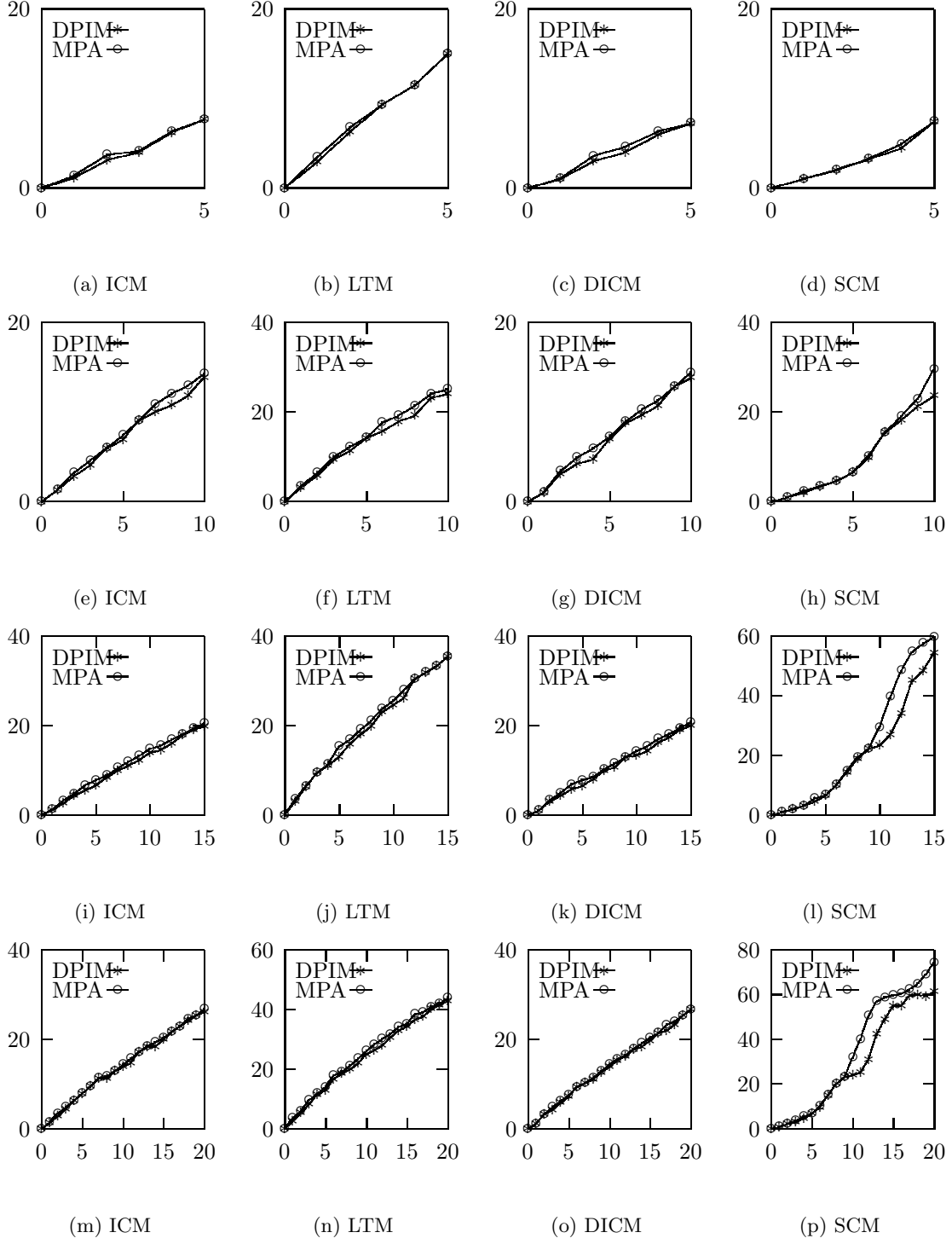
(a) ICM    (b) LTM    (c) DICM    (d) SCM

(e) ICM    (f) LTM    (g) DICM    (h) SCM

(i) ICM    (j) LTM    (k) DICM    (l) SCM

(m) ICM    (n) LTM    (o) DICM    (p) SCM

Figure 6: **Evaluation of** MPA **on an ego-network.** We evaluate the effectiveness of MPA by comparing the results to DPIM. The network is another ego-network from Facebook that has 150 nodes and 3,386 edges. For each plot, the x-axis is $k$ and the y-axis is the expected total influence from the seed set chosen by each respective algorithm.

17

# References

[1] Sanjeev Arora, Rong Ge, Sushant Sachdeva, and Grant Schoenebeck. Finding overlapping communities in social networks: Toward a rigorous approach. In *ACM EC '12*, 2012.

[2] W. B. Arthur. Competing technologies, increasing returns, and lock-in by historical events. *Economic Journal*, 99(394):pp. 116–131, 1989.

[3] Lars Backstrom, Daniel P. Huttenlocher, Jon M. Kleinberg, and Xiangyang Lan. Group formation in large social networks: membership, growth, and evolution. In *KDD '06*, pages 44–54, 2006.

[4] Abhijit Banerjee, Arun G Chandrasekhar, Esther Duflo, and Matthew O Jackson. The diffusion of microfinance. *Science*, 341(6144), 2013.

[5] Shishir Bharathi, David Kempe, and Mahyar Salek. Competitive influence maximization in social networks. In *WINE '07*, pages 306–311, Berlin, Heidelberg, 2007. Springer-Verlag.

[6] Christian Borgs, Michael Brautbar, Jennifer T Chayes, and Brendan Lucier. Maximizing social influence in nearly optimal time. *SODA '14*, 2014.

[7] J. J. Brown and P. H. Reingen. Social ties and word-of-mouth referral behavior. *Journal of Consumer Research*, 14:350–362, 1987.

[8] D. Centola. The spread of behavior in an online social network experiment. *Science*, 329(5996):1194–1197, 2010.

[9] Damon Centola and Michael Macy. Complex contagions and the weakness of long ties. *American Journal of Sociology*, 2007.

[10] Wei Chen, Wei Lu, and Ning Zhang. Time-critical influence maximization in social networks with time-delayed diffusion process. *AAAI '12*, 2012.

[11] Wei Chen, Yajun Wang, and Siyu Yang. Efficient influence maximization in social networks. In *KDD '09*. ACM, 2009.

[12] Wei Chen, Yifei Yuan, and Li Zhang. Scalable influence maximization in social networks under the linear threshold model. In *ICDM '10*, pages 88–97. IEEE, 2010.

[13] Aaron Clauset, Mark EJ Newman, and Cristopher Moore. Finding community structure in very large networks. *Physical review E*, 70(6):066111, 2004.

[14] Edith Cohen, Daniel Delling, Thomas Pajor, and Renato F Werneck. Sketch-based influence maximization and computation: Scaling up with guarantees. In *CIKM '14*, pages 629–638. ACM, 2014.

[15] J. Coleman, E. Katz, and H. Menzel. The diffusion of an innovation among physicians. *Sociometry*, 20:253–270, 1957.

[16] T. G. Conley and C. R. Udry. Learning about a new technology: Pineapple in Ghana. *American Economic Review*, 100(1):35–69, 2010.

[17] Sanjoy Dasgupta. A cost function for similarity-based hierarchical clustering. In *STOC '16*, pages 118–127, New York, NY, USA, 2016. ACM.

[18] P. Domingos and M. Richardson. Mining the network value of customers. In *7th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 57–66, 2001.

[19] Amit Goyal, Wei Lu, and Laks VS Lakshmanan. Celf++: optimizing the greedy algorithm for influence maximization in social networks. In *WWW '11*, pages 47–48. ACM, 2011.

[20] Amit Goyal, Wei Lu, and Laks VS Lakshmanan. Simpath: An efficient algorithm for influence maximization under the linear threshold model. In *ICDM '11*, pages 211–220. IEEE, 2011.

[21] Sanjeev Goyal and Michael Kearns. Competitive contagion in networks. In *STOC '12*, pages 759–774, 2012.

[22] Mark Granovetter. Threshold models of collective behavior. *American Journal of Sociology*.

[23] George Karypis and Vipin Kumar. MeTis: Unstructured Graph Partitioning and Sparse Matrix Ordering System, Version 4.0. http://www.cs.umn.edu/~metis, 2009.

[24] David Kempe, Jon M. Kleinberg, and Éva Tardos. Maximizing the spread of influence through a social network. In *KDD '03*, pages 137–146, 2003.

[25] David Kempe, Jon M. Kleinberg, and Éva Tardos. Influential nodes in a diffusion model for social networks. In *ICALP '05*, pages 1127–1138, 2005.

[26] Masahiro Kimura and Kazumi Saito. Tractable models for information diffusion in social networks. In *PKDD '06*. 2006.

[27] Jon Kleinberg. Small-world phenomena and the dynamics of information. *NIPS '02*, 1:431–438, 2002.

[28] K. Lerman and R. Ghosh. Information contagion: An empirical study of the spread of news on Digg and Twitter social networks. In *ICWSM '10*, pages 90–97, 2010.

[29] Jure Leskovec, Lada A. Adamic, and Bernardo A. Huberman. The dynamics of viral marketing. In *ACM EC '06*, pages 228–237, 2006.

[30] Jure Leskovec, Andreas Krause, Carlos Guestrin, Christos Faloutsos, Jeanne VanBriesen, and Natalie Glance. Cost-effective outbreak detection in networks. In *KDD '07*, pages 420–429. ACM, 2007.

[31] Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. http://snap.stanford.edu/data, June 2014.

[32] Bo Liu, Gao Cong, Dong Xu, and Yifeng Zeng. Time constrained influence maximization in social networks. In *ICDM '12*, pages 439–448. IEEE, 2012.

[33] Bo Liu, Gao Cong, Yifeng Zeng, Dong Xu, and Yeow Meng Chee. Influence spreading path and its application to the time constrained social influence maximization problem and beyond. *IEEE, Knowledge and Data Engineering*, 2014.

[34] Brendan Lucier, Joel Oren, and Yaron Singer. Influence at scale: Distributed computation of complex contagion in networks. In *KDD '15*, pages 735–744. ACM, 2015.

[35] S. Morris. Contagion. *Review of Economic Studies*, 67(1):57–78, 2000.

[36] Elchanan Mossel and Sébastien Roch. Submodularity of influence in social networks: From local to global. *SIAM J. Comput.*, 39(6):2176–2188, 2010.

[37] Huy Nguyen and Rong Zheng. Influence spread in large-scale social networks–a belief propagation approach. In *Machine Learning and Knowledge Discovery in Databases*, pages 515–530. Springer, 2012.

[38] Elizabeth Levy Paluck, Hana Shepherd, and Peter M. Aronow. Changing climates of conflict: A social network experiment in 56 schools. *Proceedings of the National Academy of Sciences*, 113(3):566–571, 2016.

[39] M. Richardson and P. Domingos. Mining knowledge-sharing sites for viral marketing. In *KDD '12*, pages 61–70, 2002.

[40] Daniel M Romero, Brendan Meeder, and Jon Kleinberg. Differences in the mechanics of information diffusion across topics : Idioms , political hashtags , and complex contagion on twitter. In *WWW '11*, pages 695–704. ACM, 2011.

[41] Lior Seeman and Yaron Singer. Adaptive seeding in social networks. In *FOCS '13*, pages 459–468. IEEE, 2013.

[42] Youze Tang, Xiaokui Xiao, and Yanchen Shi. Influence maximization: Near-optimal time complexity meets practical efficiency. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pages 75–86. ACM, 2014.

[43] Gordon Tullock. *Towards a theory of the rent-seeking society*, chapter Efficient Rent Seeking. Texas A&M University Press, 1980.

[44] D. J. Watts. A simple model of global cascades on random networks. *Proceedings of the National Academy of Sciences*, 99(9):5766–5771, 2002.