

Proactive Seeding for Information Cascades in Cellular Networks

Francesco Malandrino
Dipartimento di Elettronica,
Politecnico di Torino, Italy
malandrino@tlc.polito.it

Maciej Kurant, Athina Markopoulou
EECS Dept. and CalIT2
University of California, Irvine
{mkurant, athina}@uci.edu

Cedric Westphal, Ulas C. Kozat
DOCOMO USA Labs
Palo Alto, CA
{cwestphal,kozat}@docomolabs-usa.com

Abstract—In today’s Internet, online social networks (OSNs) play an important role in informing users about content. At the same time, mobile devices provide ubiquitous access to this content through the cellular infrastructure. In this paper, we propose Proactive Seeding— a technique for minimizing the peak load of cellular networks, by proactively pushing (“seeding”) some content to the users before they actually request it. We exploit the fact that the interest in the content is spread over OSNs, which makes it, to certain extent, predictable. We develop a family of algorithms that take as input information cascades, and possibly the background traffic load or the local connectivity among mobiles, and select which nodes to seed and when. We prove that Proactive Seeding is optimal when the prediction of information cascades is perfect. In our simulations, driven by traces from Twitter and cellular networks, Proactive Seeding leads to 20%-50% reduction in the peak load.

I. INTRODUCTION

Cellular traffic is growing exponentially, tripling every year, with a share of video traffic increasing from 50% now to an expected 66% by 2015 [1]. Credit Suisse reported in [2] that 23% of base stations globally have utilization rates of more than 80 to 85% in busy hours, up from 20% last year. This dramatic increase in demand is generating serious problems for 3G networks and these problems are likely to remain in 4G networks as well. Another aggravating fact for the operators is that the cellular network traffic greatly fluctuates throughout the day, following strong daily and weekly patterns, as we show in Fig. 4(c). Since the cellular network is provisioned for *peak traffic*, any capability that can distribute the network load more evenly over time would significantly address the current as well as future capacity shortcomings for the operator.

At the same time, in today’s Internet, online social networks (OSNs)¹ are becoming an increasingly important way in which users are informed about content. This is not surprising: people tend to value highly the content recommended by friends or people with similar interests (*e.g.*, members of the same groups), and are also likely to recommend it further to others.

The growth of cellular traffic and of OSN’s importance are inherently related. Indeed, mobile devices are quickly becoming the primary mean to access OSNs. For example, one third of all Facebook users regularly access the service from their mobile devices and they generate twice as much activity than

non-mobile users [3]. Consequently, the interest diffusion over OSNs translates directly into increased cellular traffic.

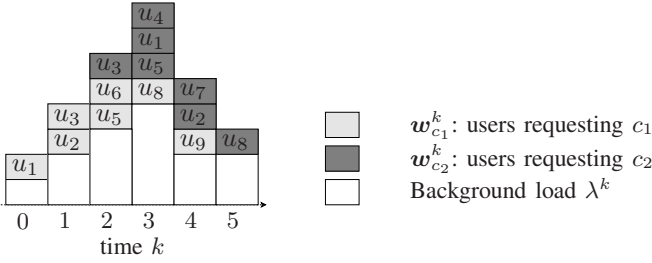
Cellular operators may try to exploit the knowledge of such interest diffusion to alleviate the peak demand in cellular traffic. One approach is to *delay* some of the traffic, *e.g.*, by limiting the diffusion of interest [4] or by using techniques that trade-off user delay for traffic load [5,6].

We take a different approach and aim at serving *impatient* users, *i.e.*, users that expect the content right after they demand it and do not tolerate large delays and jitters. Our key observation is that given the vast information often available to the cellular operator (*e.g.*, address books, session logs, location history, partnerships with OSNs, etc.), the network can detect information cascades and predict the future demand. Consider, for example, the case of Youtube videos: Google reported that up to 200 million Youtube videos per day were delivered to mobile devices in 2010 [1]. Many views of these videos are due to the spread of their URLs over various OSNs. The evolution of such cascades of forwarded URLs depends on the structure of the OSN, similarity of users and other relevant features. With this information, it is possible to model and predict the diffusion of interest [7,8]. For example, in [9], the authors apply machine learning techniques to Twitter traces, and predict more than half of URL-based cascades of tweets with only a 15% false positive rate.

In this paper, we propose Proactive Seeding, a technique for reducing the peak load in cellular networks, while providing users with low (or zero) access latency. Proactive Seeding exploits the predictability of future demand by proactively pushing (“seeding”) the content to users before, and no later than, they request it. This allows us to move some cellular traffic from the busiest hours to times with lower load and thus to reduce its peaks, as illustrated in Fig. 1. Proactive Seeding is optimal in the offline setting (*i.e.*, assuming perfect knowledge of all information cascades), in the sense that it minimizes the peak load while delivering the content to a user no later than she requests it. In our simulation driven by traces from Twitter and cellular networks, Proactive Seeding leads to 20%-50% reduction in the cellular peak load. In the case of imperfect prediction, where the gains are naturally reduced, we show that the conservative approach of underestimating the future demand still guarantees positive gains. Finally, we show how Proactive Seeding can be combined with techniques [10,11] that exploit the local device-to-device (D2D) connectivity (over WiFi or Bluetooth), and that such a combination performs better than

¹By OSNs here we refer to online social networks such as Facebook and Twitter, websites with social networking features such as Digg.com, blogs, email communication, and other online networks that exploit social ties for interest diffusion.

(a) without Proactive Seeding



(b) with Proactive Seeding

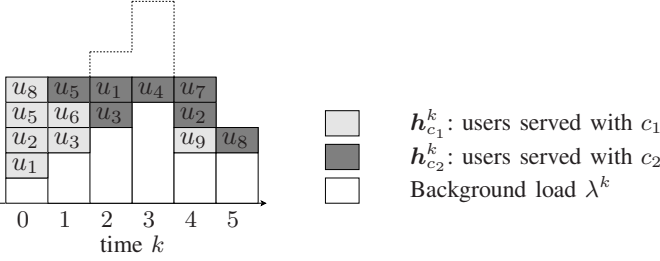


Fig. 1. Illustration of Proactive Seeding in a system with two types of contents $\mathcal{C} = \{c_1, c_2\}$ disseminated among 9 users $\mathcal{U} = \{u_1 \dots u_9\}$, in presence of the background load λ^k . (a) The diffusion of interest between the users in content c_1 (bright gray) and c_2 (dark gray). For example, $u_3 \in \mathcal{w}_{c_2}^2$ means that user u_3 becomes interested in content c_2 at time $k=2$. Without Proactive Seeding, users request and pull the content through cellular right when they get interested in it ($\mathcal{h}_c^k \equiv \mathcal{w}_c^k$), which results in an uneven total cellular load (the total height of bars). (b) Proactive Seeding serves some users before they actually become interested in the content ($\mathcal{W}_c^k \subseteq \mathcal{H}_c^k$). The total load becomes more even in time and its peaks decrease (here by 3 units).

each technique separately.

The structure of the rest of the paper is as follows. In Section II, we provide the formal problem statement. In Section III, we present the Proactive Seeding solutions under the assumption that demand can be perfectly predicted. In Section IV, we modify our framework to allow for imperfect, probabilistic estimation of the prediction. In Section V, we present our evaluations results. We overview the related literature in Section VI, then conclude the paper with Section VII.

II. PROBLEM STATEMENT

We distinguish between two components of cellular traffic: (i) background load and (ii) predictable traffic.

A. Background cellular load

We refer as background (cellular) load to all traffic which is out of our control: its content cannot be predicted (at least not with a reasonable accuracy) and/or served before the actual request occurs. For example, phone conversations and other types of real-time traffic contribute to background load. We denote by λ^k the total amount of background load at time frame k , $0 \leq k \leq K$.

We illustrate λ^k by white bars in Fig. 1; note that because the content composing it cannot be predicted or served earlier, λ^k remains unchanged in Fig. 1(b).

B. Predictable cellular traffic

In contrast, the predictable cellular traffic is all the traffic that can somehow be predicted and thus proactively served. Denote

by \mathcal{U} the set of all users, and by \mathcal{C} the set of all existing pieces of predictable content. We assume that transmitting a single piece $c \in \mathcal{C}$ of content to a single user $u \in \mathcal{U}$ takes exactly a single unit of cellular traffic.² Now, denote by $\mathcal{w}_c^k \subseteq \mathcal{U}$ the set of users that demand (“want”) the content $c \in \mathcal{C}$ exactly at time frame k . In other words, \mathcal{w}_c^k describes the diffusion of interest in content c (typically over OSNs). Let

$$\mathcal{W}_c^k = \bigcup_{m=0}^k \mathcal{w}_c^m \quad (\mathcal{W}_c^k \subseteq \mathcal{U}) \quad (1)$$

be the cumulative version of \mathcal{w}_c^k , i.e., the set of all users that have requested c until frame k . Finally, we denote by $k(u, c)$ the time when user u demands content c , i.e., such that $u \in \mathcal{w}_c^{k(u, c)}$.

In the example in Fig. 1(a), $\mathcal{w}_{c_1}^2 = \{u_5, u_6\}$ and, consequently, $k(u_5, c_1) = k(u_6, c_1) = 2$.

C. Transmission schedule

In this paper, we decouple the diffusion of interest in the content (i.e., demand) from the actual delivery process. To this end, we denote by $\mathcal{h}_c^k \subseteq \mathcal{U}$ the set of users that get (“have”) content c over cellular network exactly at frame k . Its cumulative version

$$\mathcal{H}_c^k = \bigcup_{m=0}^k \mathcal{h}_c^m \quad (\mathcal{H}_c^k \subseteq \mathcal{U})$$

is the set of all users that have c at frame k . In the other words, \mathcal{h}_c^k is a *schedule* that determines when the cellular operator sends content c to which users.

For example, in Fig. 1(b), $\mathcal{h}_{c_1}^1 = \{u_3, u_6\}$ and $\mathcal{h}_{c_2}^1 = \{u_5\}$.

D. User Impatience

In this work, we consider the case where all users are *impatient*: a user $u \in \mathcal{U}$ wants to enjoy content $c \in \mathcal{C}$ right after she becomes interested in it. This means that u should receive c at time l not larger than $k(u, c)$, i.e., $u \in \mathcal{h}_c^l$ such that $l \leq k(u, c)$. This is achieved by guaranteeing that

$$\mathcal{W}_c^k \subseteq \mathcal{H}_c^k \quad \text{for every } k \text{ and } c. \quad (2)$$

For example, in Fig. 1(b), we push content c_1 to user u_5 at time $k = 0 < k(u_5, c_1) = 2$, which is allowed by Eq.(2). In contrast, sending it at time $k > 2 = k(u_5, c_1)$ would violate the constraint in Eq.(2).

E. Objective

Using the notation above, the *total cellular traffic/load* at time k can be decomposed as the sum of background cellular load and total predictable traffic, i.e.,

$$\text{total cellular load} = \lambda^k + \sum_{c \in \mathcal{C}} |\mathcal{h}_c^k|. \quad (3)$$

Our objective is to minimize the peak of total cellular load, i.e.,

$$\text{minimize} \quad \max_{0 \leq k \leq K} \left(\lambda^k + \sum_{c \in \mathcal{C}} |\mathcal{h}_c^k| \right) \quad (4)$$

²In practice, the content spread over OSNs may greatly vary in size: a ten-minutes-long Youtube movie is orders of magnitude bigger than a photograph. All the equations can be easily modified to reflect heterogeneous content size, at the cost of notation clarity.

subject to the user impatience constraint in Eq.(2).

Note that because we have no control over the diffusion of interest w_c^k , we can affect Eq.(4) only by choosing the schedule h_c^k . We give an example of such an optimized schedule in Fig. 1(b). In particular, we (*i.e.*, the cellular operator) predict which users will be interested in content c , and proactively seed some of them with c when the cellular load is relatively small, *e.g.*, during the previous night. This allows us to reshape the cellular traffic and reduce its peaks, but not the total traffic.

III. PROACTIVE SEEDING ALGORITHMS

In this section, we focus on the *offline* case, where we have perfect knowledge of the future diffusion of interest, *i.e.*, we know w_c^k for all time frames k and pieces of content c . The offline case serves as a baseline for understanding the maximum achievable gains. It also serves as a building block for the more realistic, *online* scenario, where prediction of the future is imperfect, described in Sec. IV.

A. Special Case: single content, no background load

Let us first consider the simplest, yet intuitive case: there is only a single content ($C = \{c\}$) and no background load ($\lambda^k = 0$). An example of the demand curve corresponding to such a cascade (*e.g.*, a single content flash-crowd) is shown in Fig. 2: the total number of users interested in the content increases until reaches a peak and then decreases.

In this special case, objective Eq.(4) is equivalent to minimizing $\max_k(|h_c^k|)$ subject to the user impatience constraint Eq.(2). Intuitively, this entails delivering the content more evenly over time. Ideally, we would like to send the content with a constant seeding rate $|h_c^k|$ and thus at linear $|H_c^k|$. This rate should be the lowest possible, while still satisfying Eq.(2). Because $C = \{c\}$, Eq.(2) is satisfied if $|W_c^k| \leq |H_c^k|$ for every k . Consequently, $|H_c^k|$ should be linear and never smaller than $|W_c^k|$. This leads to an intuitive geometric solution: Draw a straight line that crosses point $(-1,0)$ and is tangential to $|W_c^k|$. The optimal service rate $|h_c^k|$ is determined by the point where the line crosses the y-axis. We show an example in Fig. 2.

It is also easy to see that this optimal rate $|h_c^k|$ is also provided by the following formula

$$|h_c^k| = \left[\max_{l=k}^K \frac{|W_c^0| - |H_c^l|}{l+1} \right]. \quad (5)$$

B. General Case: multiple contents, background traffic

The simple geometric solution from Sec. III-A does not directly extend to the general case, *i.e.*, in presence of arbitrary background cellular load $\lambda^k > 0$ and multiple contents $|C| > 1$. For example, Eq.(5) would not necessarily satisfy the user impatience constraint Eq.(2) for each of the $|C| > 1$ contents separately.

To address these problems, we propose the Proactive Seeding algorithm, shown in Alg. 1. We construct the seeding schedule h_c^k iteratively, starting from an empty set (line 1). In lines 2-6, we create a list L of existing user-content pairs (u, c) , sorted according to the growing want times $k(u, c)$. Note that user u may appear in L multiple times, *i.e.*, exactly once for each

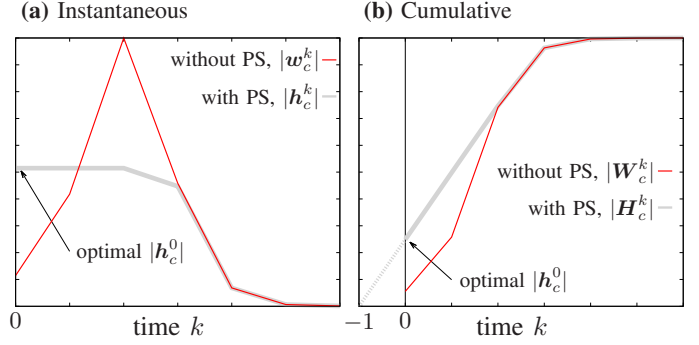


Fig. 2. Geometric interpretation of optimal Proactive Seeding (PS) under a single content cascade ($C = \{c\}$), with no background cellular load ($\lambda^k = 0$), as described in Sec. III-A. The curve represents a typical cascade on the Facebook social graph (see Sec. V-C). We minimize the peak instantaneous cellular load in (a) while satisfying the impatience constraint Eq.(2), by proactively seeding the users at a constant rate, until the cascade passes. The optimal seeding rate $|h_c^0|$ can be found by studying the cumulative version (b) of the time evolution, where a line anchored at point $(-1,0)$ and tangential to $|W_c^k|$, crosses the y-axis at point $(0, |h_c^0|)$.

Algorithm 1 Proactive Seeding

Require: $w_c^k \forall c, k, \lambda^k \forall k$ *future demand and load*
1: $h_c^k \leftarrow \emptyset \forall c, k$
2: $L \leftarrow \emptyset$
3: **for all** (u, c) such that $u \in W_c^K$ **do**
4: $L \leftarrow L \cup \{(u, c)\}$
5: **end for**
6: **sort** L by increasing $k(u, c)$
7: **for all** (u, c) in L **do** *water-filling*
8: $k^* \leftarrow \arg \min_{0 \leq l \leq k(u, c)} (\lambda^l + \sum_c |h_c^l|)$
9: $h_c^{k^*} \leftarrow h_c^{k^*} \cup \{u\}$
10: **end for**
11: **return** $h_c^k \forall c, k$ *optimal*

content c she is interested in. Lines 7-9 implement a water-filling type of algorithm, where for each pair (u, c) we find the time frame $k^* \leq k(u, c)$ with the smallest total cellular load $\lambda^{k^*} + \sum_c |h_c^{k^*}|$. We then schedule this pair (u, c) at time k^* by adding u to $h_c^{k^*}$ (line 9). Finally, once all existing pairs (u, c) are scheduled, Proactive Seeding returns the seeding schedule h_c^k for all contents c and time frames k .

We illustrate the output of Proactive Seeding in the example of Fig. 1(b). The sorted list L resulting after line 6 is $L = [(u_1, c_1), (u_2, c_1), (u_3, c_1), (u_5, c_1), (u_6, c_1), (u_3, c_2), (u_8, c_1), (u_5, c_2), (u_1, c_2), (u_4, c_2), (u_9, c_1), (u_2, c_2), (u_7, c_2), (u_8, c_2)]$. For pair (u_1, c_1) , we have $k(u_1, c_1) = 0$, and therefore lines 8-9 result in $k^* = 0$ and $h_{c_1}^0 = \{u_1\}$, respectively. When processing the second element in L , (u_2, c_1) , we have $\lambda^l + \sum_c |h_c^l| = 2$ for both $l = 0$ and $l = 1$. We arbitrarily break this tie by setting $k^* = 0$, which results in $h_{c_1}^0 = \{u_1, u_2\}$. The third pair (u_3, c_1) has now a unique $k^* = 1$, and is scheduled therein. The process continues until L is exhausted.

This schedule h_c^k returned by Proactive Seeding is optimal:

Theorem 1 (Optimality of Proactive Seeding). *The seeding schedule $h_c^k, \forall c, k$, created by Proactive Seeding minimizes the peak load (objective in Eq.(4)), while satisfying the user impatience constraint Eq.(2) for each content c separately.*

Proof: First, note that the frame k^* chosen for user u in line 8 is not greater than the time $k(u, c)$ when u actually wants the content. Therefore, by construction, the schedule created by Proactive Seeding always satisfies the user impatience constraint Eq.(2) for every content c separately.

We now have to prove that the objective Eq.(4) is met by Proactive Seeding. Denote by $L(j)$ the set of all pairs (u, c) such that $k(u, c) = j$ and by $L(i, j) = \bigcup_{m=i}^j L(m)$. Denote by $\mathbf{h}(j)$ the transmission schedule constructed by Proactive Seeding just after processing the pairs $L(j)$ in lines 7-9. In other words, $\mathbf{h}(j)$ schedules all contents for all users that want it not later than at time j . Consequently, $\mathbf{h}(K)$ denotes the entire schedule, $\mathbf{h}(K) \equiv \bigcup_{c,k} \mathbf{h}_c^k$. We prove the optimality of Proactive Seeding by induction on j , as follows.

Initialization ($j = 0$): For every pair $(u, c) \in L(0)$, line 8 automatically sets $k^* = 0$. Consequently, $\mathbf{h}(0)$ schedules all pairs $L(0)$ at time slot 0. This is the only feasible solution, thus the optimal one.

Induction step: Assume that $\mathbf{h}(j)$ is optimal for all pairs $L(0, j)$. We now must prove that $\mathbf{h}(j+1)$ is optimal for all pairs $L(0, j+1)$.

Denote by $\max(\mathbf{h}(j))$ the peak total cellular load resulting from $\mathbf{h}(j)$. Either an optimal allocation will increase the peak rate at $j+1$, or keep it constant. Thus we can distinguish two cases, as follows:

Case 1: It is possible to schedule the pairs $L(j+1)$ such that $\max(\mathbf{h}(j+1)) = \max(\mathbf{h}(j))$. In this case, lines 7-9 guarantee that this equality holds under Proactive Seeding, by iteratively choosing the least loaded time slots. Now, because $\max(\mathbf{h}(j))$ is optimal, it is the smallest value that does not violate the impatience constraint Eq.(2). So $\mathbf{h}(j+1)$ cannot be lower than $\max(\mathbf{h}(j))$ without violating Eq.(2). Consequently, $\max(\mathbf{h}(j+1)) = \max(\mathbf{h}(j))$ implies the optimality of $\mathbf{h}(j+1)$.

Case 2: It is *not* possible to schedule the pairs $L(j+1)$ such that $\max(\mathbf{h}(j+1)) = \max(\mathbf{h}(j))$. We can now distinguish two sub-cases, depending of the background load at time $j+1$:

Case 2.1: If $\max(\mathbf{h}(j+1)) = \lambda^{j+1}$ is achievable, then lines 7-9 of Proactive Seeding will achieve that by iteratively choosing the least loaded time slots. In this case, the peak load is equal to the background load λ^{j+1} . Such a peak load is optimal, because, by definition, background load cannot be changed.

Case 2.2: If $\max(\mathbf{h}(j+1)) = \lambda^{j+1}$ is *not* achievable, then lines 7-9 guarantee that $\max(\mathbf{h}(j+1)) - \min(\mathbf{h}(j+1)) \leq 1$, where $\min(\mathbf{h}())$ denotes the minimal total cellular load resulting from $\mathbf{h}()$. Consequently, $\max(\mathbf{h}(j+1))$ cannot be decreased and $\mathbf{h}(j+1)$ is thus optimal. ■

Note: Although optimal in the sense of objective Eq.(4), Proactive Seeding does not guarantee that the users will be served in the order they request the content; it may schedule user u before user w , even if $k(u, c) > k(w, c)$. For example, in Fig. 1 user u_3 wants content c_1 before user u_5 , but is scheduled to receive it after u_5 , as we show in Fig. 1(b). However, it is easy to see that an additional step that reshuffles the users to enforce the “first-want-first-serve” (*i.e.*, chronological) order, preserves the optimality and feasibility of the resulting schedule \mathbf{h}_c^k .

C. Extension: D2D-aware Proactive Seeding

In addition to their cellular connections, it is often the case that some users are within physical proximity of each other and can establish direct device-to-device (or D2D [12]) connections between them, *e.g.*, via ad-hoc 802.11 or Bluetooth. If these users are interested in the same content, they can exploit their D2D connectivity, and thus offload the cellular network. Several variants of this idea have been studied in the past, *e.g.*, in [10,11,13,14]. What makes this particularly promising, in our context, is the fact that there is a correlation between geographical proximity and proximity on the social graph [15]. We show below (and later, in simulations) that these techniques can be combined with Proactive Seeding, and address two complementary aspects: using the D2D connections helps to offload the total aggregated cellular load, while Proactive Seeding helps to smooth the load over time.

The D2D connectivity graph changes over time. We denote by $\mathcal{N}^k(u)$ all D2D neighbors of user u at time k . Consider time $k(u, c)$ when user u becomes interested in content c . We will assume that each mobile user behaves as follows:

- 1) If u has been seeded with c before, no action is needed.
- 2) Otherwise, u attempts to pull c from its current local neighbors $\mathcal{N}^{k(u,c)}(u)$. This is possible only if at least one of these neighbors has c , *i.e.*, if $\mathcal{N}^{k(u,c)}(u) \cap \mathbf{H}_c^{k(u,c)} \neq \emptyset$.
- 3) Otherwise, u fetches c through the cellular network.

Depending on the extent to which the operator is aware of D2D connectivity, different optimizations are possible:

1) D2D-unaware Proactive Seeding: In this simplest scenario, the operator does not have information about the location of users and thus performs Proactive Seeding without taking proximity into account. Consequently, user u can benefit from D2D, in an opportunistic way, *i.e.*, only if u has not been seeded earlier (*i.e.*, if $u \in \mathbf{h}_c^{k(u,c)} \cap \mathbf{w}_c^{k(u,c)}$), which results in

$$\mathbf{h}_c^k \leftarrow \mathbf{h}_c^k \setminus \left\{ u \in \mathbf{h}_c^k \cap \mathbf{w}_c^k : \mathcal{N}^{k(u,c)}(u) \cap \mathbf{H}_c^{k(u,c)} \neq \emptyset \right\}.$$

In the example of Fig. 1, user u_4 will pull content c_2 from its D2D neighbors $\mathcal{N}^3(u_4)$ at time $k = 3$ if at least one of them is in $\{u_1, u_3, u_5\} = \mathbf{H}_{c_2}^2$ (*i.e.*, already has c_2).

2) D2D-aware Proactive Seeding: In this scenario, the operator has information about location and thus proximity of users³ and takes it into account while seeding. In particular, it applies Proactive Seeding but avoids seeding user u if u will be able to get the content from its neighbors. This can be achieved by the following refinement of schedule \mathbf{h}_c^k :

$$\mathbf{h}_c^k \leftarrow \mathbf{h}_c^k \setminus \left\{ u \in \mathbf{h}_c^k : \mathcal{N}^{k(u,c)}(u) \cap \mathbf{H}_c^{k(u,c)} \neq \emptyset \right\}.$$

In the example of Fig. 1, we will seed user u_5 with content c_2 at time $k = 1$. If we know that $u_5 \in \mathcal{N}^3(u_1)$, *i.e.*, that u_1 and u_5 will form a D2D connection at time $k = 3$ (*i.e.*, when u_1 wants c_2) then then we can exclude u_1 from $\mathbf{h}_{c_2}^2$.

³This information can be obtained either directly from the cellular network or can be contributed by the user *e.g.*, via applications on OSNs (such as FourSquare, Facebook Places) or on a smartphone, in exchange for the service.

IV. DEALING WITH UNCERTAINTY

In Sec. III, we developed an optimal seeding strategy given the full and precise knowledge of the future (i) cellular background load, and (ii) predictable traffic pattern. Clearly, the performance of Proactive Seeding will strongly depend on the quality of our estimation of the predictable traffic w_c^k . Many prediction techniques have been proposed in the literature and developing new ones is out of the scope of this paper. Instead, in this section, we review some existing techniques, and we show how they can be incorporated in Proactive Seeding.

A. Interest diffusion on OSNs

In this paper, we are interested in the content that becomes popular through social ties.⁴ One can exploit the structure of the social network and information about interest diffusion, in order to predict information cascades. Such a prediction can then serve as input (instead of the offline knowledge) to our predictive seeding algorithms.

There is a rich literature on predicting the diffusion of interest in social networks, see *e.g.*, [7,8]. In our context, predicting the future progress of a cascade related to content c , can be modeled as finding the probability

$$\mathbb{P}(w_c^{k+1}, w_c^{k+2}, \dots \mid w_c^k, w_c^{k-1}, \dots, w_c^0, I_{other}), \quad (6)$$

where $w_c^k, w_c^{k-1}, \dots, w_c^0$ is the observed history at the current time k , and I_{other} represents any other available piece of information. Below, we comment on how some of the existing approaches translate into the Eq.(6) probabilities.

1) *The threshold model:* In the threshold model [7], each user u is associated with a threshold $0 \leq \theta_u \leq 1$. u becomes interested in the content at time $k+1$ if at least a (weighted) fraction of θ_u of her neighbors are interested in it at time k . This model is deterministic, *i.e.*, the probabilities in Eq.(6) are either 0 or 1.

2) *The cascade model:* In the cascade model [7,8], each edge (u, w) of the social graph is associated with an activation probability $q_{u,w}$. If user u gets interested in the content at time k , then the edge (u, w) is used exactly once to determine whether user w will become interested in the content at frame $k+1$, which happens with probability $q_{u,w}$. In other words, given the activation probabilities $q_{u,w}$ (*i.e.*, I_{other}) and the history $w_c^k, w_c^{k-1}, \dots, w_c^0$, the cascade model gives us the following probabilities, concerning the next time frame:

$$\mathbb{P}(w_c^{k+1} \mid w_c^k, w_c^{k-1}, \dots, w_c^0, I_{other}), \quad (7)$$

which is a special case of Eq.(6).

3) *Machine learning:* Another line of research focuses on machine learning techniques that make use of all the available information. For example, in [9], the authors, based on the observed history, manage to accurately predict more than half of future re-tweets (of URL links) with 15% false positives.

⁴An alternative approach to learn w_c^k could be by studying the download patterns of individual users. For example, assume that user u regularly visits www.bbc.co.uk (or checks out her Facebook updates) everyday in the morning. We can then seed u with some heavier content (graphics, videos) over night. A machine learning approach could help us choose whom to seed and with what content.

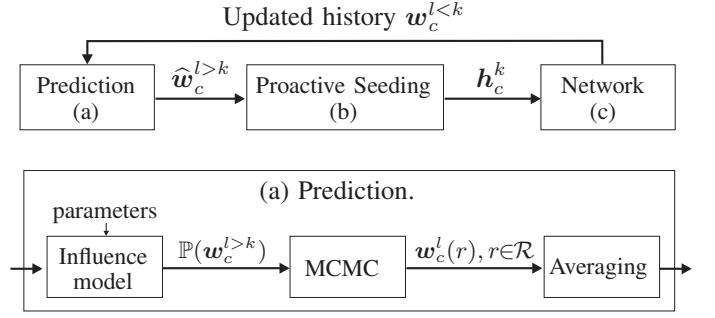


Fig. 3. Adaptive Proactive Seeding. (top) High-level overview. (bottom) The “Prediction” block.

B. From probabilities to Proactive Seeding

Given the knowledge of probabilities in Eq.(6), we follow the procedure presented in Fig. 3. First, at the current time k , we use Eq.(6) to calculate the most likely future $\hat{w}_c^{l>k}$ (Fig. 3(a)). Next, we plug $\hat{w}_c^{l>k}$ into Proactive Seeding (Fig. 3(b)), which returns us the schedule h_c^k for the current time frame. Finally, we implement h_c^k and collect the actual evolution of demand w_c^k that is used to refine our calculations in the next time frame (Fig. 3(c)). This means that our scheme is *adaptive* – at every iteration it updates the history by the current state of the network and recalculates h_c^k .

Our prediction includes all times l between the current time k and time K . K is the latest time for which at least one realization of the interest diffusion process has at least one user interested in content c , *i.e.*, $|w_c^K| \geq 1$. For instance, for the cascade influence model, K is trivially upper-bounded by the total number of users, *i.e.*, $K \leq |U|$.

In Fig. 3(bottom), we show in more detail the “Prediction” block from Fig. 3. Given the knowledge of Eq.(6), we are, in principle, able to calculate exactly the expected future demand $\mathbb{E}[w_c^{l>k}]$. In practice, however, the solution space is too big (especially if the number $|U|$ of users or the final time K are large) to do it precisely. Instead, we run an MCMC (Monte Carlo Markov Chain) simulation, *i.e.*, we use Eq.(6) to generate a number of realizations $w_c^{l>k}(r)$, $r \in \mathcal{R}$. This step is illustrated by the middle block in Fig. 3(bottom). Next, we average over all $|\mathcal{R}|$ realizations (right-most block in Fig. 3, bottom), as follows.

First, we estimate the *number of users* $|\widehat{W}_c^K|$ that eventually become interested the content, by the average over all the realizations:

$$|\widehat{W}_c^K| = \frac{1}{|\mathcal{R}|} \cdot \sum_{r \in \mathcal{R}} |W_c^K(r)|.$$

Next, we decide *which users* will become interested in the content, by taking $|\widehat{W}_c^K|$ users with the highest observed probabilities $\widehat{\mathbb{P}}(u \in W_c^K) = \frac{1}{|\mathcal{R}|} \cdot |\{r \in \mathcal{R} : u \in W_c^K(r)\}|$ to request it. Finally, we interpret as $k(u, c)$ the time that is the most frequent across the realizations in R :

$$\widehat{k}(u, c) = \arg \max_{0 \leq k \leq K} |\{r \in \mathcal{R} : u \in w_c^k(r)\}|.$$

The above process provides an estimate \hat{w}_c^k of the future demand, which we use as input to Proactive Seeding, as

in Fig. 3(b).

V. EVALUATION

In this section, we evaluate the performance of Proactive Seeding through simulation.

A. Performance Metric

Without Proactive Seeding, user u fetches the content c over cellular when she wants it, which yields $\mathbf{h}_c^k \equiv \mathbf{w}_c^k$ and the peak cellular load equal to $\max_k (\lambda^k + \sum_c |\mathbf{w}_c^k|)$. In contrast, with Proactive Seeding, the peak cellular load drops to $\max_k (\lambda^k + \sum_c |\mathbf{h}_c^k|)$. Our main performance metric is the relative *gain* in peak cellular load, defined as

$$\gamma = \frac{\max_k (\lambda^k + \sum_c |\mathbf{w}_c^k|) - \max_k (\lambda^k + \sum_c |\mathbf{h}_c^k|)}{\max_k (\lambda^k + \sum_c |\mathbf{w}_c^k|)}.$$

Clearly, the larger the amount of the predictable traffic, the bigger gain γ we can expect. We therefore denote by ρ the ratio of the unpredictable traffic (aggregate over all contents) over the aggregate predictable traffic, *i.e.*,

$$\rho = \frac{\text{aggregated unpredictable traffic}}{\text{aggregated predictable traffic}} = \frac{\sum_k \lambda^k}{\sum_k \sum_c |\mathbf{w}_c^k|}. \quad (8)$$

B. Offline Scenario (using Twitter, Cellular and D2D traces)

First, we consider the offline case, with large-scale simulations fed by real traces of (a) interest diffusion process in Twitter [9], (b) background traffic from a US cellular operator [16], and (c) mobility [17]. This allows us to evaluate Proactive Seeding in presence of cellular background load and techniques that exploit D2D connectivity. We assume a priori knowledge of (a), (b), (c), and we evaluate how much gain γ is achieved by Proactive Seeding.

1) *Description of Datasets:* (a) *Predictable traffic* π^k : We use the Twitter trace from [9], where the authors collected the tweets that carry a URL (which defines our content), over a period of 300 hours (12.5 days). For our simulations, we kept only the “re-tweets” (indicated by an RT tag), which allows us to directly follow the cascades of interests in valuable (non-spam) content on Twitter (see also RT-cascades in [9]). Furthermore, in order to be able to observe the full evolution of such cascades, we exclude the URLs that appear in the first three or the last three hours of the trace. This leaves us with around 2.5M of tweets from 554K different users, sharing about 9000 contents (URLs). In Fig. 4(a), we show the evolution of two typical cascades from that trace. The “cascade” behavior is easy to see: the URL’s popularity quickly increases over time, reaches a peak, and then declines. However, when we aggregate all the 9000 cascades together in Fig. 4(b), the individual cascade shapes are not visible anymore; instead, the aggregated predictable traffic π^k clearly follows the daily pattern.⁵

(b) *Background cellular load* λ^k : As background load λ^k , we take a cellular traffic trace coming from a major operator in

⁵Recall, however, that our constraint Eq.(2) is defined for each content, not for the aggregated traffic.

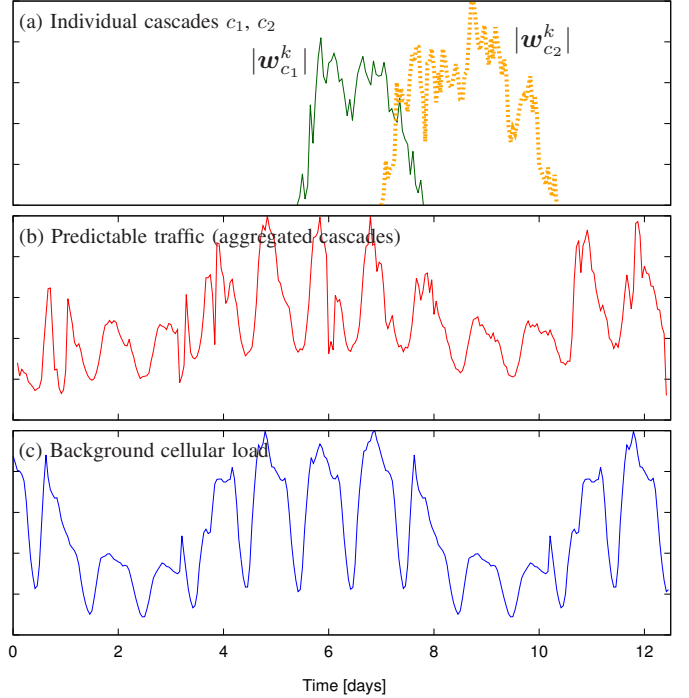


Fig. 4. Traces used in offline simulations. (a) Example of two individual Twitter cascades; (b) All 9000 Twitter cascades together [9]; (c) Background cellular load from a US operator [16]. For the sake of readability, all figures are normalized with respect to the peak value of the data they represent (*i.e.*, they do not have the same scale).

one US state [16].⁶ Because this trace covers one full week (at a resolution of 1 hour), we replicate it, concatenate, and shift to match the 12.5 days of the Twitter trace. The result is presented in Fig. 4(c). Similarly to Twitter, the cellular background load follows weekly and daily patterns.

(c) *D2D connectivity:* We use the Infocom06 contact trace [17] to simulate the device-to-device (D2D) connectivity. The trace logs the D2D contacts between 78 devices (iMotes) distributed to the attendees, over a period of three days.

For each content c , we randomly map the users \mathbf{H}_c^K (*i.e.*, eventually requesting c) to the users in the trace. Because of the limited size and duration of the trace, we replicate these users when $|\mathbf{H}_c^K| > 78$, and we repeat the connectivity pattern when the diffusion of interest in content c lasts for more than 3 days. Finally, users u and w are defined neighbors in our connectivity graph at hour k , *i.e.*, $w \in \mathbf{N}^k(u)$ and $u \in \mathbf{N}^k(w)$, if u and w encounter each other within this hour (according to the Infocom06 trace).

The above mapping matches users U with nodes in the mobility trace in a purely random way. We also experimented with D2D connectivity graphs that reflect various levels of correlations between physical proximity and friendship. The results were similar and are omitted for lack of space.

⁶Strictly speaking, the trace [16] represents the total cellular traffic. For simplicity of presentation (*e.g.*, independence of ρ), we interpret this trace as the background cellular load λ^k . We have also considered in simulations this trace as the total load, subtracting π^k to get the background load. The results in both cases are very similar.

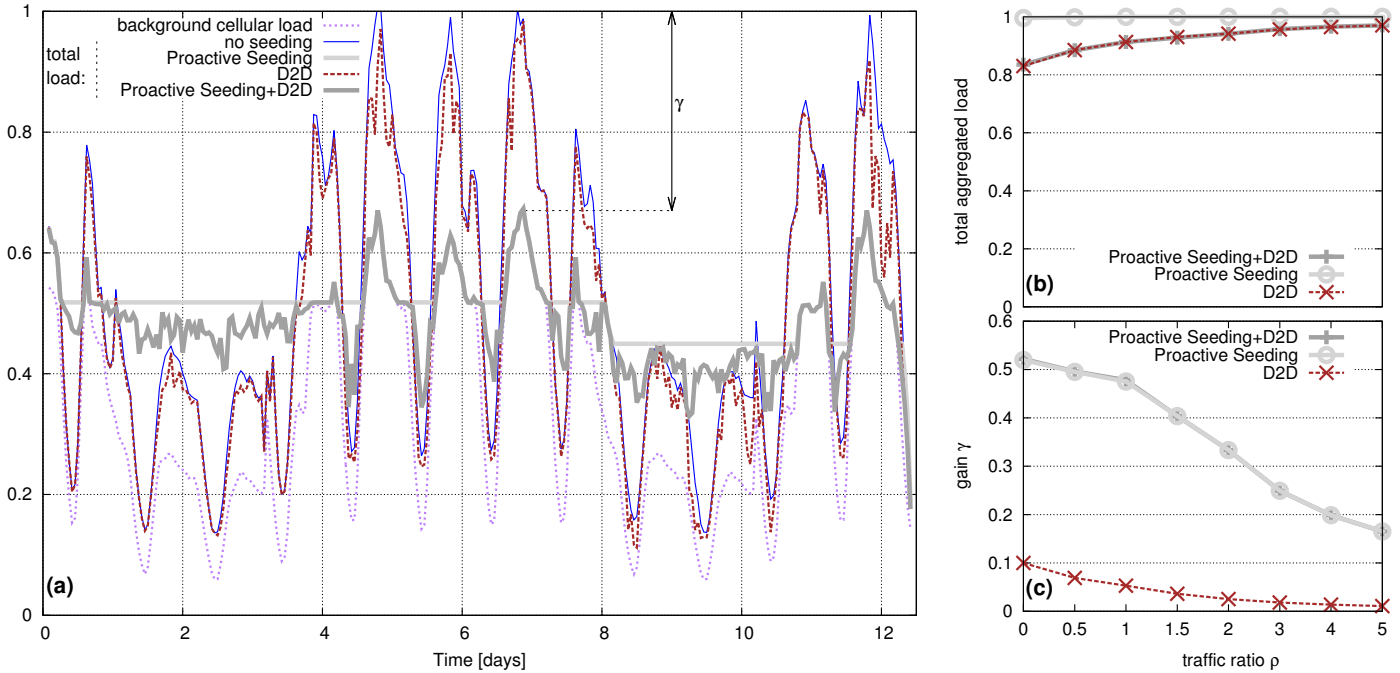


Fig. 5. Offline simulations driven by traces of (i) Twitter cascades (predictable traffic), (ii) background cellular load, and (iii) mobility. (a) Per hour time evolution of the total cellular load $\lambda^k + \pi^k$ under various scenarios, for traffic ratio of $\rho = 2$. (b) Aggregated cellular traffic as a function of ρ . (c) Gain γ as a function of ρ .

2) *Results:* In Fig. 5(a) we focus on a case when $\rho = 2$, *i.e.*, the background load is twice the predictable traffic, and depict the time evolution of the total load on the 3G network in the following cases:

- no seeding: All users get the content they are interested in through the cellular network (*i.e.*, $\mathbf{h}_c^k = \mathbf{w}_c^k, \forall c, k$).
- Proactive Seeding: Proactive Seeding algorithm is used to schedule predictable traffic. D2D is disabled.
- D2D: Users exploit the D2D connectivity as explained in Sec. III-C, but Proactive Seeding is disabled.
- Proactive Seeding + D2D: predictable traffic is scheduled using Proactive Seeding *and* users exploit D2D links if available.

The no-seeding scenario results in a cellular load that is very uneven over time, with high peaks and periods of very low usage. Under D2D, we observe a slight reduction in the network load, with the peaks almost unchanged. In contrast, Proactive Seeding effectively reshapes the total cellular traffic, reducing the peaks by exploiting the less busy periods. Note that the peak load (around day 9) corresponds to a peak in the *background* load, which confirms that Proactive Seeding is optimal with respect to objective Eq.(4) (as we proved in Theorem 1). Finally, when we combine Proactive Seeding and D2D, we observe a further reduction in the network load.

Fig. 5(b) and Fig. 5(c) show how the aggregated (*i.e.*, over the whole trace duration) load and the gain γ depend on the ratio ρ between predictable and background load. Unsurprisingly, the higher ρ , the less beneficial Proactive Seeding becomes. Proactive Seeding effectively reduces the peak load (Fig. 5(c)), but has no impact on the aggregated load (Fig. 5(b)). The effect of D2D is quite the opposite. Applying both Proactive Seeding

and D2D, we get the best of both worlds: *i.e.*, a significant reduction in both the peak and the aggregated load.

C. The Online Case (using Diffusion Models on OSNs)

Sec. V-B assumed full knowledge of the entire traces. In this section, we consider the case where the future can be predicted only with some amount of uncertainty, as described in Sec. IV. For ease of explanation, we assume no background load and a single content c and we focus on evaluating the effect of uncertainty on the results.

1) *Social Graphs (Datasets):* We use datasets from two different graphs, each capturing a different type of social tie.

- *Facebook:* The New Orleans network of the Facebook social graph [18], consisting of 63K vertices and 816K edges. The rationale for using this data set is that friends in Facebook share links and thus participate in spreading information about content.
- *Email:* a trace of e-mail contacts, consisting of 1133 nodes and 5452 edges. The rationale behind using this datasets is that emails often contain links that propagate in a viral way, leading to information cascades.

2) *Social Influence (Models):* Using each of the previous graphs, we simulate interest diffusion through the cascade model [7,8] described in Sec. IV-A2. We assume that 5% of users are interested in the content at time $k = 0$. The activation probability for each edge (u, w) is set to $q_{u,w} = 0.1$. (We have also tried a range of parameters, omitted for lack of space, and results were qualitatively similar.)

3) *Uncertainty about the model and its parameters:* Although the cascade model provides us with a probabilistic output, there are several other major sources of uncertainty

about the future, which naturally lead to errors in the prediction. In particular, in practice, (i) we can never know exactly the model driving the spread of information and (ii) we can never know precisely the parameters of such a model. We capture these two effects in our simulations by introducing a multiplicative noise ν to the probabilities Eq.(6), *i.e.*, we set $\mathbb{P}() \leftarrow \min(1, \nu \mathbb{P}())$. For example, $\nu = 1.2$ results in a systematic overestimation of the future demand by 20%, and $\nu = 0.8$ underestimates it by 20%.

4) *Results:* In Fig. 6, we present results for the Facebook (left) and Email (right) graphs. Although the two networks are very different in size and structure, they exhibit the same qualitative behavior, with a clear cascade evolution. The way Proactive Seeding works is easy to observe: the users known (or assumed) to request the content during the peak time are served during earlier frames, thus reducing the peak load.

For both networks, we compare the ideal (*i.e.*, offline) performance with the adaptive (*i.e.*, online) case, in which the demand is not known a priori. In the latter, we consider three values of the noise ν . If our prediction is not systematically biased ($\nu = 1$), the online performance of Proactive Seeding is close to the optimal (offline). In contrast, systematically overestimating ($\nu > 1$) or underestimating ($\nu < 1$) the future demand leads to less gain γ , but with qualitatively different effects. *Overestimating* the demand means serving users that will never need the content, thus wasting network and user resources. In the extreme case, it may even lead to a negative gain, *i.e.*, a peak load $\max_k |h_c^k|$ greater than the peak demand $\max_k |w_c^k|$. On the other hand, *underestimating* the demand is conservative, as moves towards the no-seeding case. The gain γ can decrease, but is still above zero. Therefore, as a practical take-away from our online evaluation, we can recommend to tune the prediction parameters so as to underestimate rather than overestimation the demand.

Fig. 6 also allows us to see how the adaptiveness, *i.e.*, the fact that at each time frame k we feed the actual set \mathbf{W}_c^k of users interested in the content back to the prediction algorithm, allows us to recover from prediction errors. If $\nu > 1$, we tend to overestimate the number of users interested in the content at the begin of the cascade. However, as we observe the actual number of interested users, we are able to correct the error, and schedule fewer users in the subsequent frames. Conversely, if $\nu < 1$, we start seeding fewer users than we should, and we make it up for this error later. Notice however that both such cases imply a peak load that is higher than the ideal (*i.e.*, offline) one.

VI. RELATED WORK

Proactive Seeding touches upon several research areas. We now review the closest ones and how they relate to our work.

Opportunistic communication. When several users are interested in the same content and they are in proximity of each other, some of them may be able to use device-to-device connections, *e.g.*, through WiFi or Bluetooth, to get the content, instead of their cellular connection. This opportunistic communication results in offloading the cellular network. In [5], device-to-device and cellular connections are used to disseminate dynamic content, so as to maximize the “freshness” of the content. The connectivity of nodes are taken into account

in order to select the right users to act as relays. As an example, a node with many neighbors is more likely to be selected as a relay. The work in [13] considers a similar scenario and assumes that social ties among the users are strongly correlated with their physical proximity and similar interests. [6] offloads the cellular network through proximity connections, while still meeting strict deadlines. With respect to these works, we have a different goal – decreasing the *peak* load on the cellular network – and a stronger constraint, *i.e.*, the user impatience.

Socially-aware forwarding. Another body of work [19]–[24] exploits the principle that social ties affect the mobility, and eventually the proximity, of users. Evidence has been provided, for example in [15], which shows that there is a significant correlation between similar interests and geographical proximity, for four different OSNs (BrightKite, FourSquare, LiveJournal and Twitter). Therefore, knowledge about social ties, can be taken into account to optimize routing for content delivery.

[19] presents Bubble rap – a routing protocol for DTNs. Devices detect the centrality of the community the user belongs to, based on the frequency of contact. This is then used for routing decisions. [21,22] use social information to optimize content discovery in a publish/subscribe setting: the more social users are given a special role in the delivery process. [20,23,24] exploit social information to route queries and to decide which items should be cached or duplicated.

In our work, we exploit social ties for a different purpose, namely predicting the content requests in order to proactively serve them. Furthermore, we limit the amount of information that users disclose to their peers (*e.g.*, users do not broadcast their whole list of topics of interest, as in [20]).

Interest diffusion in social networks. There is a large body of literature on diffusion in networks, including but not limited to technological networks. The classic work in [25] reviews several influence models and proposes an algorithm for selecting which nodes to seed so as to maximize the diffusion, given the social structure. This is different from our objective in this paper (to minimize the peak of the cascade) as well as in the fact that seeding is done only once in the beginning, while we adaptively seed at every time slot.

Such influence models are motivated by the many studies of information diffusion on actual social networks. For example, [26] identifies and studies several cascades on the Flickr social network. [27] analyzes 1.5 million YouTube videos, showing that not all popular videos are “social” and that highly social videos rise to, and fall from, their peak popularity more quickly than less social videos. Somewhat related to our work, [4] considers information cascades caused by social influence and shows which links to select and limit this influence, so as to delay the peak of the load caused on the cellular network.

Predicting content popularity. Forecasting the popularity of content, with or without taking into account network effects, is another active research area. [28] presents methods for predicting the popularity of items given historical access data, but without taking into account the network effect, for the YouTube and Digg social networks. [9] collects a dataset of 22M tweets, containing 15M URLs and presents a methodology (based on influence models) which predicts more than half of the tweets in the dataset with only 15% false positives.

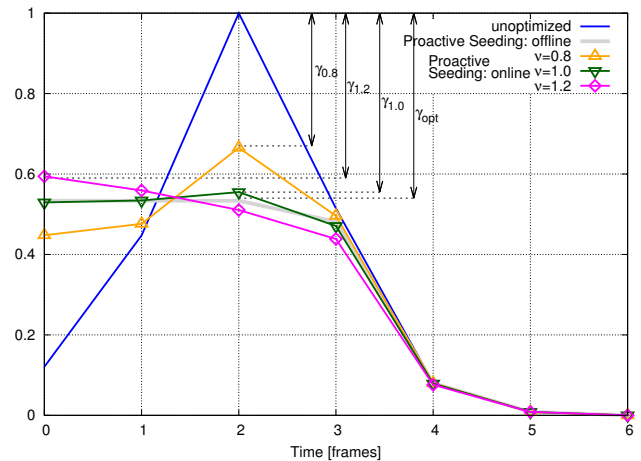
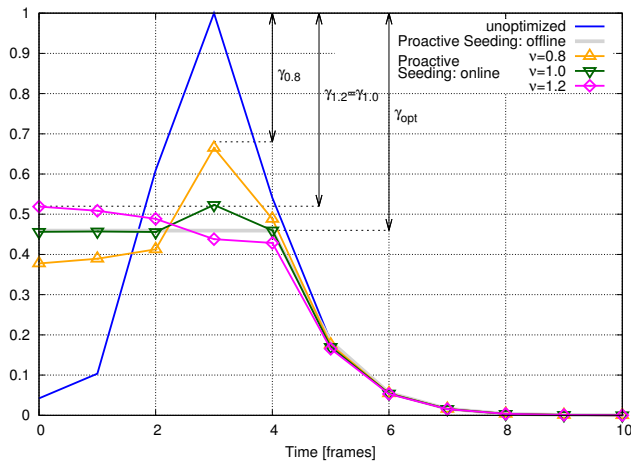


Fig. 6. Online simulations on the Facebook (left) and Email (right) graphs.

In this paper, we use the dataset collected in [9] for simulations of the offline scenario. More generally, we rely on prediction models as a part of our machinery, but we do not develop one ourselves.

VII. CONCLUSION

We presented proactive seeding for information cascades in social media - as a new technique to reduce the peak demand in cellular networks. In the special case of single content with no background load, the optimal solution that minimizes the peak load turns out to have an intuitive interpretation. In the general case of multiple contents with known background traffic, we provide a greedy algorithm and prove its optimality, in the offline case. In the online case, we investigated the performance of the proposed solutions by replacing the actual future demand by the predicted demand. Our evaluation showed robustness, especially when underestimating the total demand. We also extended our algorithm to take into account D2D communication, when this is available, thus offloading the total cellular traffic, in addition to reducing the peak load. Our evaluation over real traces indicate that proactive seeding via predicting social cascades significantly reduces the peak load as much as 50%.

ACKNOWLEDGEMENTS

We would like to thank the authors of [9] and [16], for providing the Twitter and 3G traffic traces, respectively.

REFERENCES

- [1] "Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2010-2015," *Cisco white paper*, 2011.
- [2] Credit Suisse, "U.S. wireless networks running at 80% of capacity," <http://benton.org/node/81874>, July 2011.
- [3] Facebook, "Facebook statistics," <http://www.facebook.com/press/info.php?statistics>, July 2011.
- [4] H. Sharara, C. Westphal, S. Radosavac, and U. C. Kozat, "Utilizing Social Influence in Content Distribution Networks," *IEEE ICC*, 2011.
- [5] S. Ioannidis, A. Chaintreau, and L. Massoulie, "Optimal and Scalable Distribution of Content Updates over a Mobile Social Network," *IEEE INFOCOM*, Apr. 2009.
- [6] J. Whitbeck, Y. Lopez, J. Leguay, V. Conan, and M. De Amorim, "Relieving the wireless infrastructure: When opportunistic networks meet guaranteed delays," *Arxiv preprint arXiv:1007.5459*, 2010.

- [7] D. Kempe, J. Kleinberg, and E. Tardos, "Maximizing the spread of influence through a social network," in *ACM SIGKDD*, 2003.
- [8] J. Goldenberg, B. Libai, and E. Muller, "Talk of the network: A complex systems look at the underlying process of word-of-mouth," *Marketing Letters*, vol. 12, no. 3, pp. 211–223, 2001.
- [9] W. Galuba, K. Aberer, D. Chakraborty, Z. Despotovic, and W. Kellerer, "Outtweeting the twitterers—predicting information cascades in microblogs," in *WOSN*, 2010.
- [10] L. Yin and G. Cao, "Supporting cooperative caching in ad hoc networks," *IEEE Trans. on Mobile Computing*, vol. 5, no. 1, pp. 77–89, Jan. 2006.
- [11] H. Gupta and S. Das, "Benefit-Based Data Caching in Ad Hoc Networks," *IEEE Trans. on Mobile Computing*, vol. 7, no. 3, pp. 289–304, Mar. 2008.
- [12] 3GPP Work Item Description TSG-RAN, "Study on LTE device to device discovery and communication," RP-110707 LTE-D2D RAN, May 2011.
- [13] B. Han, P. Hui, V. Kumar, M. Marathe, G. Pei, and A. Srinivasan, "Cellular traffic offloading through opportunistic communications: a case study," in *ACM CHANTS*. ACM, 2010.
- [14] Qualcomm, "Flashling," <http://www.qualcomm.com/stories/2011/02/08/phone-conversations-minus-people>, February 2011.
- [15] S. Scellato, C. Mascolo, M. Musolesi, and V. Latora, "Distance matters: Geo-social metrics for online social networks," in *WOSN*, 2010.
- [16] M. Shafiq, L. Ji, and A. Liu, "Characterizing and modeling internet traffic dynamics of cellular devices," *ACM SIGMETRICS*, 2011.
- [17] J. Scott, R. Gass, J. Crowcroft, P. Hui, C. Diot, and A. Chaintreau, "CRAWDAD data set cambridge/haggle (v. 2009-05-29)," Downloaded from <http://crawdad.cs.dartmouth.edu/cambridge/haggle>, May 2009.
- [18] B. Viswanath, A. Mislove, M. Cha, and K. Gummadi, "On the evolution of user interaction in facebook," in *WOSN*, Barcelona, Spain, 2009.
- [19] P. Hui, J. Crowcroft, and E. Yoneki, "Bubble rap: social-based forwarding in delay tolerant networks," *IEEE Trans. on Mobile Computing*, 2010.
- [20] C. Boldrini, M. Conti, and A. Passarella, "ContentPlace: social-aware data dissemination in opportunistic networks," in *ACM MSWiM*, 2008.
- [21] E. Yoneki, P. Hui, S. Chan, and J. Crowcroft, "A socio-aware overlay for publish/subscribe communication in delay tolerant networks," in *ACM MSWiM*, 2007.
- [22] P. Costa, C. Mascolo, M. Musolesi, and G. Picco, "Socially-aware routing for publish-subscribe in delay-tolerant mobile ad hoc networks," *IEEE JSAC*, vol. 26, no. 5, pp. 748–760, Jun. 2008.
- [23] E. Jaho and I. Stavrakakis, "Joint interest- and locality-aware content dissemination in social networks," in *IEEE/IFIP WOSN*, Feb. 2009.
- [24] A. J. Mashhadi, S. B. Mokhtar, and L. Capra, "Habit: Leveraging human mobility and social network for efficient content dissemination in Delay Tolerant Networks," in *IEEE WoWMoM*, Jun. 2009.
- [25] D. Easley and J. Kleinberg, *Networks, Crowds, and Markets: Reasoning About a Highly Connected World*. Cambridge University Press, 2010.
- [26] M. Cha, A. Mislove, and K. P. Gummadi, "A measurement-driven analysis of information propagation in the flickr social network," *WWW*, 2009.
- [27] T. Broxton, Y. Interian, J. Vaver, and M. Wattenhofer, "Catching a Viral Video," in *IEEE ICDM Workshops*, 2010.
- [28] G. Szabó and B. Huberman, "Predicting the popularity of online content," *Communications of the ACM*, vol. 53, no. 8, pp. 80–88, Aug. 2008.