

Joint Optimization for Social Content Delivery in Wireless Networks

Xiangnan Weng, John S. Baras

Institute for Systems Research and Department of Electrical and Computer Engineering

University of Maryland, College Park, Maryland 20742

Email: {wengxn, baras}@umd.edu

Abstract—Over the last decade, success of social networks has significantly reshaped how people consume information. Recommendation of contents based on user profiles is well-received. However, as users become increasingly mobile, little is done to consider the constraints of the wireless environment. In this paper, we investigate a centralized wireless content delivery system, aiming to optimize overall user experience given the capacity constraints of the wireless networks, by deciding what contents to deliver and how. We propose a scheduling framework incorporating content deliverability. Results indicate this novel joint optimization approach outperforms existing systems, which separates recommendation and delivery, especially when the wireless network is operating at maximum capacity.

Index Terms—Social network, wireless network, optimization.

I. INTRODUCTION

Social networks, which enable sharing information among users in various forms, have long existed on the Internet. However, it is not until the past decade that we witnessed their huge commercial and social success. Facebook, Twitter, YouTube, along with innumerable social networks, greatly facilitate information exchange for users all across the world. Their successes rely heavily, if not solely, on content recommendation based on user profiles, including but not limited to social connections, engagement history, social groups, etc. On one hand, social network providers are collecting more data about users than ever, to deliver high quality contents for users to consume; on the other hand, users are more willing to share data with their trusted social network providers, in return for better experience. Therefore, researchers in this area are working diligently to maximize user experience by recommending individualized contents that are most likely to incur user engagements, with all the available data at hand.

As more and more users are moving towards their mobile devices and become mostly mobile, wireless networks become the new frontier for social network applications. The trend of an ever growing mobile population is inevitable and irreversible, due to the flexibility provided by wireless technologies: with seamless and ubiquitous connection, users could consume information whenever and wherever they want or need.

Unfortunately, most of the social network applications were designed for wired connections, assuming unlimited capacity and reliable transmission. Unlike wired connections, wireless networks are limited by insufficient radio frequency resources

and changing channel. As multimedia contents de facto dominate contents consumed in social network applications, the disparity, between the capacity constraints of wireless networks and the assumptions of guaranteed delivery of contents, results in poor mobile experience and loss of user engagements, or ultimately users themselves. Even with new generations of access technology (LTE-Advanced and beyond), users are still unable to fetch their contents when the number of active users increases, because the wireless network becomes congested as it approaches its capacity limits and quality of service degrades drastically. Hence, we are motivated to provide services in these scenarios so that users could still consume contents, even though the contents delivered might not be the best contents recommended by the social network applications.

Existing solutions fail to utilize the broadcast nature of wireless networks. In light of the preliminary research and industrial observations of social network applications, we note that users exhibit patterns of temporal, spatial and social correlations: they are likely to consume similar contents, though not at the exact same time. This fact leads us towards a novel solution that multicasts and precaches contents to groups of users to reduce redundant transmissions. It also suggests that under congested circumstances, the system shall deliver contents that maximize overall user experience across the system, but not necessarily optimal for each individual user. In other words, we are more tempted to deliver contents in a collectively optimal way when the wireless networks are congested. This approach requires information from both social and wireless networks to jointly optimize system decisions.

Certain preliminary work has been done concerning the joint user experience optimization incorporating social and wireless networks. In [1], an opportunistic protocol based on a wireless ad-hoc system is proposed to forward contents with respect to possible future transmissions that could further disseminate the contents. This protocol considers the physical feasibility of wireless transmission, but does not discuss how to allocate wireless resources in different wireless conditions. The work of [2] employs hyperbolic embedding to stabilize the network. It examined the performance of a greedy back pressure algorithm in different topologies. It provided an algorithm that could achieve optimal throughput within the capacity region. However, as stated in [3], the capacity region is hard to obtain and system utility is generally not optimized under stabilizing policies.

In this paper, we investigate the content delivery problem with capacity constraints in a system with centralized wireless infrastructure. We need to decide what contents to transmit to which users and how to transmit. Traditional systems include two separated stages: (i) the social network application chooses the best recommended contents for users; (ii) the wireless network allocates resources for transmission. As stated before, this unicast approach provides optimal results only when the wireless resource is sufficient and/or when users have absolutely nothing in common. However, if the wireless network is congested, this approach fails to work.

We introduce our system model and evaluation framework in Section II. The performance with and without look-ahead at the wireless layer is analyzed in III-B. In Section IV, we propose a scalable solution by utilizing limited multicast modes and eliminating content limits for users. Then, we analyze its performance and sensitivity against real-world data in Section V. Extension to hybrid systems handling different types of delivery is discussed in Section VI. Finally, we present conclusions in Section VII.

II. PROBLEM FORMULATION

A. System Model

We consider a centralized system that both selects contents for users according to rewards given wireless capacity constraints and delivers the contents to users via a wireless network, as illustrated in Fig.1. All the users are served using the same base station and we assume the bandwidth of the wired connections between the base station and content server(s) is sufficient enough that the base station could access contents as if they are stored locally. This assumption is valid in practice because base stations are generally connected to the Internet via fiber optic cables. Trivially, with modern chips, the storage on user devices is sufficiently large to precache all the contents scheduled for delivery.

(a) The normalized reward (or reward for simplicity) earned from delivering content j to user i is denoted by f_{ij} , $0 \leq f_{ij} \leq 1$. This number is obtained from social network applications via various big data techniques. The reward is earned only after the first successful delivery of the whole content; partial or repeated delivery does not earn (additional) reward. It is important to note that the reward matrix might not be fully filled due to either lack of interest or sufficient data, in which case we simply denote them as no rewards ($f_{ij} = 0$).

(b) The system is time-slotted with slot length T and the scheduling horizon is T_H time slots. This is consistent with modern cellular systems (GSM, UMTS, LTE, LTE-Advanced), which are all time-slot based to simplify overhead of control plane. Generally, the channel state is considered to remain unchanged within the slot. We consider decisions for all the users served by the base station at the beginning of each time slot t , with the bandwidth of the wireless network B .

(c) The wireless channel is slow fading (i.e. it remains unchanged during each scheduling time slot, as stated in (b)) and is described by signal-to-interference-noise ratio $\text{SINR}_i^{(t)}$ for

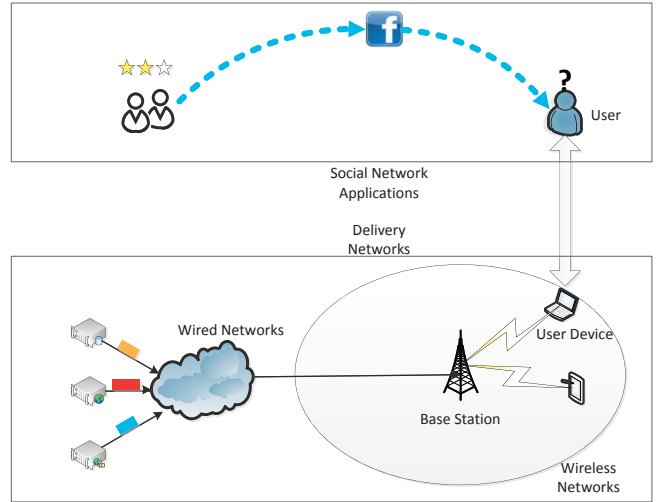


Fig. 1. System Model: The social network applications are responsible to select contents for users according to their profiles while the delivery networks (including wired backbone and wireless last-hop networks) are responsible to deliver to user devices.

user i at time slot t , where $P_i^{(t)}$, $I_i^{(t)}$, $N_i^{(t)}$ are received power strength, interference power and noise power respectively.

$$\text{SINR}_i^{(t)} = \frac{P_i^{(t)}}{I_i^{(t)} + N_i^{(t)}} \quad (1)$$

To reduce system complexity, contents shall be transmitted within one time slot; otherwise, management of multicast groups is too complicated to implement, given the changing channel state between slots.

(d) For simplicity, we assume contents could be reliably delivered with given transmission rate not exceeding the Shannon limit. Therefore, given bandwidth B and slot length T , the maximum bits deliverable to a user with SINR in the time slot is denoted by W and it follows Shannon's law:

$$\mathcal{R}(\text{SINR}) = \log_2(1 + \text{SINR}) \quad (2)$$

$$W = B \cdot T \cdot \mathcal{R}(\text{SINR}) \quad (3)$$

Table I summarizes definitions of parameters.

TABLE I
SUMMARY OF VARIABLES

Notations	Definition
M	Number of users.
N	Number of contents.
f_{ij}	Reward for delivering content j to user i .
$\alpha_{ij}^{(t)}$	Binary decision variable whether to transmit content j to user i at time slot t .
$B^{(t)}$	Total available bandwidth at time slot t .
$s_j^{(t)}$	Wireless resource allocated for content j at time slot t .
$\text{SINR}_i^{(t)}$	Signal-to-Interference-Noise ratio of user i at time slot t
W_j	Size of content j in bits

B. Problem Formulation

We formulate the content delivery problem as a mixed integer programming (MIP) problem.

$$\begin{aligned}
& \underset{\alpha_{ij}^{(t)}, s_j^{(t)}}{\text{maximize}} && \sum_{t=1}^{T_H} \sum_{i,j} \alpha_{ij}^{(t)} \cdot f_{ij} \\
& \text{subject to} && \alpha_{ij}^{(t)} \in \{0, 1\} && \forall i, j, t \\
& && \sum_{t=1}^{T_H} \alpha_{ij}^{(t)} \leq 1 && \forall i, j \\
& && \sum_j \alpha_{ij}^{(t)} \leq N_0 && \forall i, t \\
& && \alpha_{ij}^{(t)} \cdot W_j \leq s_j^{(t)} \cdot \mathcal{R}(\text{SINR}_i^{(t)}) && \forall i, j \\
& && \sum_j s_j^{(t)} \leq B \cdot T && \forall t \\
& && s_j^{(t)} \geq 0 && \forall j, t
\end{aligned} \tag{4}$$

The objective is to maximize the overall system reward, with the constraints of: (i) decisions are binary; (ii) content j is delivered to user i at most once; (iii) at most N_0 content is delivered for each user in one time slot; (iv) quality of service is satisfied; and (v) resource allocated does not exceed system capacity.

III. SCHEDULING FRAMEWORK

A. Decision at Each Time Slot

For general cases, the MIP problem (4) is NP-hard and solving it takes exponential time. Moreover, we care more about the online version of the optimization, i.e. at time slot t_0 , we only have access to current and historic channel information $\{\text{SINR}_i^{(t)}, \forall i; t = 1, \dots, t_0\}$. Solving the offline version is of little practical use for scheduling.

Therefore, at each time slot t_0 we introduce the T_L -step lookahead (no lookahead with $T_L = 0$) version of the MIP problem given the wireless channel profile of each user. Scheduling decisions $\{\alpha_{ij}^{(t_0)}\}, \{s_j^{(t_0)}\}$ at each time slot are obtained by solving this lookahead version (5). All the auxiliary symbols with tilde are the lookahead version of the corresponding parameters in the MIP problem (4).

$$\begin{aligned}
& \underset{\tilde{\alpha}_{ij}^{(t)}, \tilde{s}_j^{(t)}}{\text{maximize}} && \sum_{t=t_0}^{t_0+T_L} \sum_{i,j} \tilde{\alpha}_{ij}^{(t)} \cdot \tilde{f}_{ij}^{(t_0)} \\
& \text{subject to} && \tilde{\alpha}_{ij}^{(t)} \in \{0, 1\} && \forall i, j, t \\
& && \sum_{t=t_0}^{t_0+T_L} \tilde{\alpha}_{ij}^{(t)} \leq 1 && \forall i, j \\
& && \sum_j \tilde{\alpha}_{ij}^{(t)} \leq N_0 && \forall i, t \\
& && \tilde{\alpha}_{ij}^{(t)} \cdot W_j \leq \tilde{s}_j^{(t)} \cdot \mathcal{R}(\widetilde{\text{SINR}}_i^{(t)}) && \forall i, j \\
& && \sum_j \tilde{s}_j^{(t)} \leq B \cdot T && \forall t \\
& && \tilde{s}_j^{(t)} \geq 0 && \forall j, t
\end{aligned} \tag{5}$$

We rewrite the single transmission constraints in (4) by changing reward values at each time slot. Initially,

$$\tilde{f}_{ij}^{(1)} = f_{ij}, \forall i, j \tag{6}$$

If content j was successfully transmitted to user i at slot t_0 , the reward drops to zero to avoid future transmission(s).

$$\tilde{f}_{ij}^{(t_0+1)} = \tilde{f}_{ij}^{(t_0)} \cdot (1 - \alpha_{ij}^{(t_0)}), \forall i, j \tag{7}$$

We predict SINR in future scheduling slots for each user using historic channel information:

$$\widetilde{\text{SINR}}_i^{(t)} = \begin{cases} \text{SINR}_i^{(t_0)} & t = t_0 \\ \phi(\text{SINR}_i^{(1)}, \dots, \text{SINR}_i^{(t_0)}) & t = t_0 + 1, \dots, t_0 + T_L \end{cases} \tag{8}$$

The actual decisions (what contents to deliver $\{\alpha_{ij}^{(t_0)}\}$ and how to allocate wireless resources $\{s_j^{(t_0)}\}$) taken at time slot t_0 are the decisions at the first slot obtained from the solution of the lookahead version of optimization (5). It is trivial to prove that the choice of decisions satisfies all the constraints in MIP (4).

$$\alpha_{ij}^{(t_0)} = \tilde{\alpha}_{ij}^{(t_0)}, \forall i, j \tag{9}$$

$$s_j^{(t_0)} = \tilde{s}_j^{(t_0)}, \forall j \tag{10}$$

B. Results and Analysis

In this part, we compare our proposed scheduling system with the traditional layered design, where the social network applications attempt to provide users with the most rewarding contents regardless of the status of the wireless networks, while the wireless networks attempt to deliver the contents with best effort regardless of the rewards of contents.

1) *Simulation Setup*: There are $M = 30$ users and $N = 20$ contents in the system. Reward values f_{ij} 's are independent and uniformly distributed in $[0, 1]$. Note that this is already the largest scale within which we could obtain the optimal solution to compare the results. Content size is independent and uniformly distributed in $[10, 20]$ Mbits. The scheduling time slot has length of $T = 1s$ and the scheduling horizon is $T_H = 10$. Each user could receive at most $N_0 = 1$ content in every scheduling time slot. System level parameters are shown in Table II [4].

TABLE II
SYSTEM LEVEL SIMULATION PARAMETERS

Simulation Parameter	Value
UE distribution	UEs dropped with uniform density within the macro coverage area.
Carrier frequency	2.0 GHz
Channel model	Typical Urban (TU)
Inter-site distance	1500 m
Noise power spectral density	-174 dBm/Hz
Macro BS transmit power	40 W (46 dBm)
Macrocell path loss model	$128.1 + 37.6 \log_{10} R$ (R in km)
Macrocell shadowing model	Log normal fading with std. 10 dB
Macro BS antenna gain	15 dBi

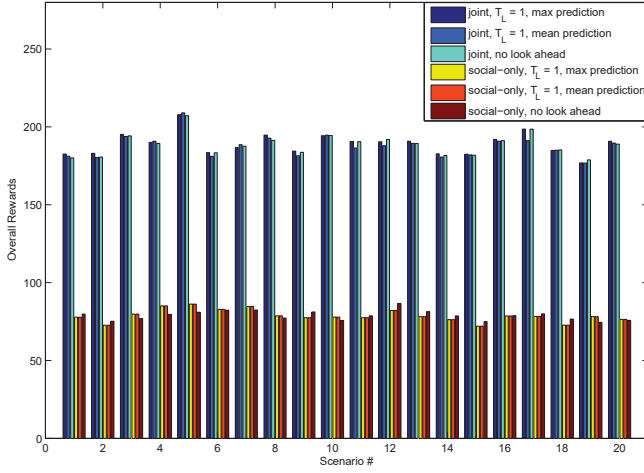


Fig. 2. Comparison of Overall System Rewards ($B = 25$ MHz)

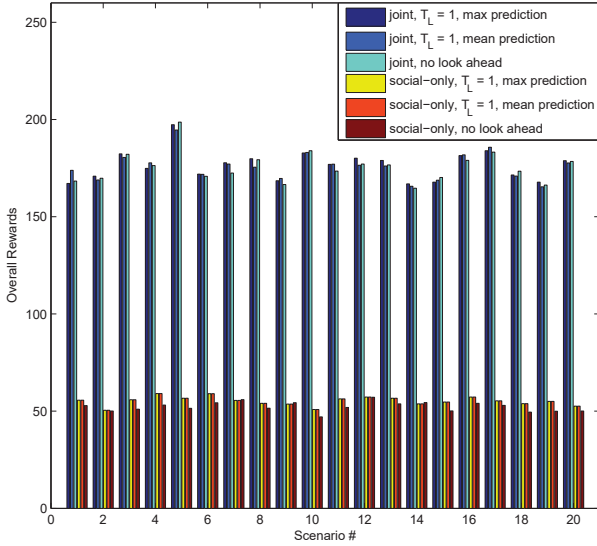


Fig. 3. Comparison of Overall System Rewards ($B = 15$ MHz)

2) *Results*: We simulate various instances against different bandwidth, as shown in Fig. 2, 3. Clearly, joint optimization outperforms the traditional system by a wide margin. The joint optimization gain increases as the available wireless resources decrease.

Surprisingly, one-step lookahead scheduling (whether we use mean or max function as SINR prediction function) does not outperform no-lookahead scheduling. Due to the extra computation cost (it takes at least 10 times of computation time to obtain optimization results), it is sufficient to schedule without looking ahead.

IV. SCALING UP

For real-world systems, the available multicast transmission modes are limited and pre-determined. Assume the system has K available transmission modes, and the associated data transmission rates and minimum channel quality requirements

are R_k and SINR_k^{th} , respectively. More formally, we substitute Shannon's Law in the data rate function (2) with a step function (slightly abusing notation $R_0 = 0$), denoting K available modes with the order convention $R_{k-1} < R_k$, $\text{SINR}_{k-1}^{th} < \text{SINR}_k^{th}$, $\forall k = 1, \dots, K$:

$$\mathcal{R}(\text{SINR}) = \sum_{k=1}^K (R_k - R_{k-1}) \cdot u(\text{SINR} - \text{SINR}_k^{th}) \quad (11)$$

Therefore, if we give up the constraints of the maximum number of contents that are allowed to be transmitted to a user, we could further reduce the scheduling problem to deciding on which wireless transmission mode we use to transmit what contents. In this way, the complexity of the problem is significantly reduced. It would only rely on the number of transmission modes K and contents N , dropping the number of users M .

For this complexity-reduction version, at slot t_0 , the reward $\hat{f}_{ij}^{(t_0)}$ for transmitting content j in wireless mode k is induced from summing up all the rewards of the users that meet the quality of service requirement of this mode.

$$\hat{f}_{jk}^{(t_0)} = \sum_{i: \text{SINR}_i^{(t_0)} \geq \text{SINR}_k^{th}} \tilde{f}_{ij}^{(t_0)} \quad (12)$$

The new scheduling formulation is thus:

$$\begin{aligned} & \underset{\hat{\alpha}_{jk}^{(t_0)}, s_j^{(t_0)}}{\text{maximize}} && \sum_{j,k} \hat{\alpha}_{jk}^{(t_0)} \cdot \hat{f}_{jk}^{(t_0)} \\ & \text{subject to} && \hat{\alpha}_{jk}^{(t_0)} \in \{0, 1\} \quad \forall j, k \\ & && \sum_{k=1}^K \hat{\alpha}_{jk}^{(t_0)} \leq 1 \quad \forall j \\ & && \hat{\alpha}_{jk}^{(t_0)} \cdot W_j \leq s_j^{(t_0)} \cdot R_k \quad \forall j, k \\ & && \sum_j s_j^{(t_0)} \leq B \cdot T \\ & && s_j^{(t_0)} \geq 0 \quad \forall j \end{aligned} \quad (13)$$

The reverse mapping between system decision and decisions for each individual user is:

$$\alpha_{ij}^{(t_0)} = \sum_{\substack{i,j: \tilde{f}_{ij}^{(t_0)} > 0 \\ k: \text{SINR}_k^{th} \leq \text{SINR}_i^{(t_0)}}} \hat{\alpha}_{jk}^{(t_0)} \quad (14)$$

Apparently, this user-aggregated version of scheduling framework only relies on N (the number of contents to be scheduled) and K (the number of multicast transmission modes in the system), thus it is robust against increase in number of users in the system. However, when we relax the constraints on the number of contents delivered to each individual user, it is important to consider fairness among users. The system is now evaluated on both average overall reward delivered to each user and fairness among users.

For each user, the overall reward delivered during the scheduling horizon $[0, T_H]$ is calculated as:

$$u_i = \sum_{t=1}^{T_H} \sum_j \alpha_{ij}^{(t)} f_{ij} \quad (15)$$

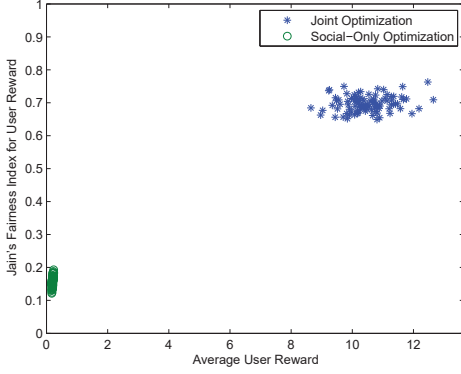


Fig. 4. Comparison of User Reward ($B = 20$ MHz) with Full Random Reward Matrix

Denote the user average reward and its variation as \bar{u}, σ_u^2 respectively.

$$\bar{u} = \frac{1}{M} \sum_{i=1}^M u_i \quad (16)$$

$$\sigma_u^2 = \frac{1}{M} \sum_{i=1}^M (u_i - \bar{u})^2 \quad (17)$$

We use Jain's fairness index [5] as fairness metric.

$$J(\vec{u}) = \frac{\bar{u}^2}{\bar{u}^2 + \sigma_u^2} \quad (18)$$

where $\vec{u} = (u_1, \dots, u_M)$. The larger the fairness index, the more fare for a scheduling result.

Fig. 4 shows the performance comparison of our proposed scheduling framework and traditional layered design, on a full random reward matrix. Each point in the graph represents one simulation instance. Clearly, the joint optimization framework is better in terms of average overall rewards delivered per user and Jain's fairness metric.

V. REAL-WORLD USER REWARDS

In real-world applications, the reward matrix $F = [f_{ij}]$ is usually sparse due to multiple factors: (1) naturally, users exhibit diverse interests towards different contents; (2) technically, it is extremely difficult, if not impossible, to gather enough information about users, in order to obtain accurate and comprehensive prediction. In our system, we do not distinguish between unknown reward or lack of interests, by assigning zero reward value, i.e. $f_{ij} = 0$, to avoid such transmissions (proof is trivial due to the formulation of the optimization problem).

We further define the sparseness of the reward matrix as follows:

$$\eta = \frac{|\{f_{ij} : f_{ij} > 0\}|}{M \cdot N} \quad (19)$$

Obviously, for full random reward matrix, $\mathbb{P}[\eta = 1] = 1$.

In light of this observation, we need to examine the performance of our scheduling framework on sparse real-world data sets.

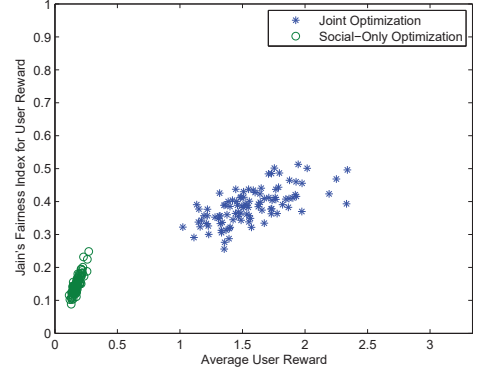


Fig. 5. Comparison of User Reward ($B = 20$ MHz) for ML-1M ($\eta = 5\%$)

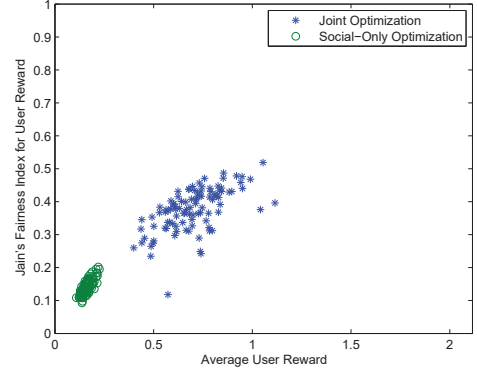


Fig. 6. Comparison of User Reward ($B = 20$ MHz) for Yahoo! Music ($\eta = 2\%$)

Our empirical analysis is based on two data sets:

(1) MovieLens [6]: users' ratings of different movies. This data set has two sub data sets. (i) ML-1M, with sparseness of $\eta = 5\%$, number of total users 6000 and number of total contents 4000. (ii) ML-10M, with sparseness of $\eta = 1\%$, number of total users 72,000 and number of total contents 10,000.

(2) The Yahoo! Music ratings for User-Selected and Randomly Selected Songs, version 1.0 data set, which is available through the Yahoo! Webscope data sharing program. This data set has sparseness of $\eta = 2\%$, number of total users 15,400 and number of total contents 1000.

A. Performance

From Fig.5 and Fig.6, it is clear that our scheduling framework still works well for real-world data sets. However, as the sparsity of reward matrix increases ($\eta \downarrow$), the performance gain compared to existing system diminishes. This result is intuitive, as multicast works best only when the number of users that are interested in the same contents is large enough.

B. Sensitivity and Saturation on Wireless Resources

In this part, we demonstrate the performance sensitivity with respect to available wireless resource (in our system, band-

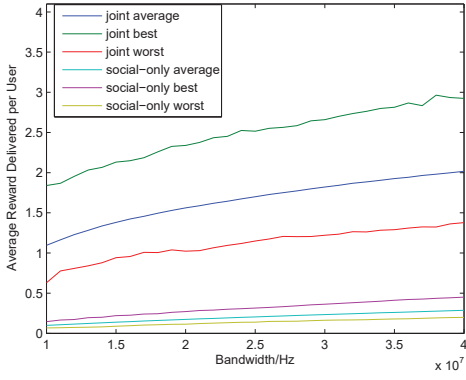


Fig. 7. Sensitivity of Average User Reward for ML-1M ($\eta = 5\%$)

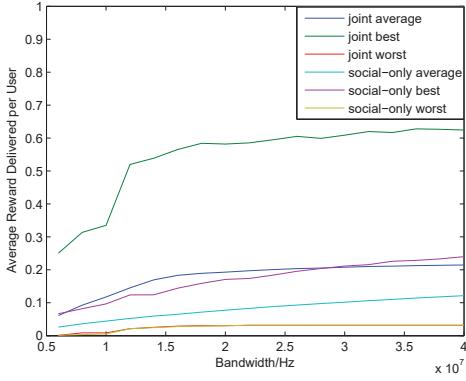


Fig. 8. Sensitivity of Average User Reward for ML-10M ($\eta = 1\%$)

width B). Intuitively, when the wireless resource is sufficient, the performance solely relies on the numerical reward values obtained from the social network, as we could deliver all the contents. However, if the wireless resource is insufficient, the performance of the system shall be reduced due to incapability to deliver the contents.

We plot the user average reward against bandwidth in Fig.7 to analyze the benefits of adding more wireless resources to the system. Fig.8 shows the performance saturation for our joint optimization approach when the delivery demand is low compared to the bandwidth.

VI. HYBRID SYSTEMS

Until now, we have been discussing systems that involve no active user requests for delivery. This assumption is generally valid for pure recommender system, in which users passively consume contents delivered to their devices. We define this delivery method as ‘push’ operation. However, practical systems are also required to handle active user requests within a designated timeframe, which we define as ‘pull’ operation. The hybrid systems shall handle both ‘push’ and ‘pull’ functions smoothly.

Now we formally define ‘pull’ operation. The user request q could be described using a quadruple $q = (t_0, i, j, t_{\max}) \in$

\mathbb{N}^4 , denoting the request originates at slot t_0 and demands the system to deliver content j to user i no later than t_{\max} (deadline). Obviously, in real systems, $t_0 \leq t_{\max}$. We further assert that all requests are valid, assuming all requests fulfilled in the past shall not occur again. This assumption could easily be implemented by serving the transmitted content(s) from the cache of user’s device. Denote the set of new and active requests at slot t as $Q_t^{\text{new}}, Q_t^{\text{active}}$, respectively, $Q_t^{\text{served}} \subseteq Q_t^{\text{active}}$ as the set of requests served and $Q_t^{\text{expired}} \subseteq Q_t^{\text{active}}$ as the set of requests expired at the end of the slot. Then,

$$Q_t^{\text{served}} = \left\{ q \in Q_t^{\text{active}} : \alpha_{q_2 q_3}^{(t)} = 1 \right\} \quad (20)$$

$$Q_t^{\text{expired}} = \left\{ q \in Q_t^{\text{active}} : q_4 = t \right\} \quad (21)$$

$$Q_{t+1}^{\text{active}} = (Q_{t+1}^{\text{new}} \cup Q_t^{\text{active}}) \setminus (Q_t^{\text{served}} \cup Q_t^{\text{expired}}) \quad (22)$$

with the trivial convention $Q_0^{\text{active}} = Q_0^{\text{served}} = Q_0^{\text{expired}} = \phi$.

The validity of requests is described as: $\forall q \in Q_{t_0}^{\text{new}}$, we have $q_1 = t_0$ and $\nexists q' \in \bigcup_{t=1}^{t_0-1} Q_t^{\text{served}}$, s.t. $q_2 = q'_2, q_3 = q'_3$.

The objective for ‘pull’ operation is to serve as many requests as possible. We could write the optimization problem for ‘pull’ system as:

$$\begin{aligned} & \underset{Q_t^{\text{served}}}{\text{maximize}} && \sum_{t=1}^{T_H} |Q_t^{\text{served}}| \\ & \text{subject to} && W_{q_3} \leq s_{q_3}^{(t)} \cdot \mathcal{R}(\text{SINR}_{q_2}^{(t)}) \quad \forall q \in Q_t^{\text{served}} \\ & && \sum_j s_j^{(t)} \leq B \cdot T \quad \forall t \\ & && s_j^{(t)} \geq 0 \quad \forall j, t \end{aligned} \quad (23)$$

To integrate this ‘pull’ system with the existing ‘push’ system in Section IV, we add additional reward for requests with respect to their expiration time. Reward transition in (7) is rewritten in part for the ‘pull’ request as:

$$\tilde{f}_{ij}^{(t_0+1)} = \begin{cases} f_{ij} \cdot (1 - \alpha_{ij}^{(t_0)}) & \exists (*, i, j, t_0) \in Q_{t_0}^{\text{expired}} \\ (\tilde{f}_{ij}^{(t_0)} + \lambda \Gamma(i, j, t_0)) \cdot (1 - \alpha_{ij}^{(t_0)}) & \text{otherwise} \end{cases} \quad (24)$$

where the incentive function Γ is exclusively awarded to active user requests:

$$\Gamma(i, j, t_0) = \begin{cases} \gamma(t_d - t_0) & \exists (*, i, j, t_d) \in Q_{t_0}^{\text{active}}, t_0 < t_d \\ 0 & \text{otherwise} \end{cases} \quad (25)$$

To balance between the two types of operations, we introduce a ‘push’ weight coefficient $\lambda \in \mathbb{R}_+$ to configure the bias of the hybrid system with respect to user requests, which reduces to pure ‘push’ system (no active user requests) when $\lambda = 0$, or pure ‘pull’ system (best-effort to serve all active user requests) when $\lambda = \infty$.

In this way, we seamlessly integrate the recommender system and active user requests together to form a hybrid delivery system that could handle both types of content delivery. Essentially, if we did not yet schedule transmission for the

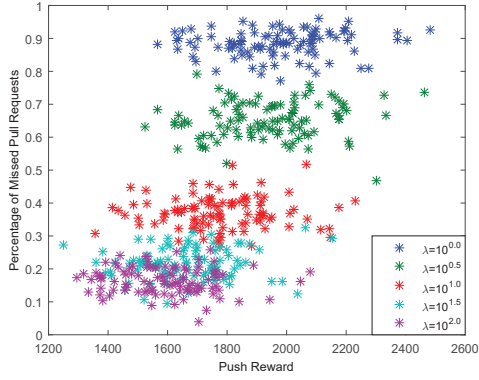


Fig. 9. Performance for Hybrid System (ML-1M, $T_d = 5$)

request(s), we shall add additional reward as time goes by, to steer the system in the direction to accomplish such requests before deadlines.

To capture the urgency when approaching the request deadline, we could wisely choose the γ function to reflect it. One possible function is:

$$\gamma(t) = \frac{1}{t} \quad (26)$$

In hybrid systems, we are evaluating the performance of both the overall user rewards delivered by the ‘push’ operation and total served requests coming from the ‘pull’ operation.

In the simulation, new user requests are generated at the beginning of each scheduling time slot with same expiration time T_d .

$$q_4 = q_1 + T_d, \forall q \quad (27)$$

The number of requests per slot follows independent and identical uniform distribution $U[0, 3]$. The requests at slot t_0 are uniformly selected in random from unfulfilled set of user-content pairs $\{(i, j) : \tilde{f}_{ij}^{(t_0)} > 0\}$ so that the transmission has not been scheduled before slot t_0 .

We plot the performance (percentage of missed ‘pull’ requests and overall user rewards delivered) of our hybrid system with scheduling horizon $T_H = 120s$ and different deadlines T_d in Fig.9,10 for comparison. Obviously, the larger the ‘push’ weight coefficient λ , the fewer missed user requests, but the less overall system reward.

VII. CONCLUSIONS

In this paper, we investigate a social content delivery system given the constraints from the wireless networks. We propose a framework that evaluates the performance of the system in terms of overall delivered rewards. To optimize system performance, we need to schedule contents and wireless resources according to the solution of a MIP problem, which requires exponential time to obtain the optimal solution. We further reduce the complexity of the joint optimization approach to rely only on the number of candidate contents and system transmission modes, regardless of the number of users, by aggregating user rewards at each supported transmission modes.

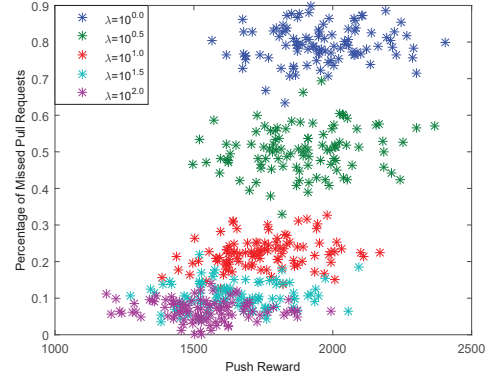


Fig. 10. Performance for Hybrid System (ML-1M, $T_d = 10$)

The simulation results indicate that our joint optimization approach provides better system performance than the traditional layered systems that schedule contents without considering constraints of wireless networks. The joint optimization gain is significant when the resources of the wireless network are comparatively insufficient, either due to (i) number of users is large, and/or (ii) conditions of wireless channel are bad for large number of users. The scheduling is also fair among users in terms of overall rewards received by each user during the scheduling horizon.

We also investigate the performance of a hybrid system, by introducing additional rewards for user generated requests (‘pull’ operation) with deadlines. This hybrid system provides a natural way to balance the resource allocation between suggested contents generated by recommender system and actual user requests. It further proves that our joint optimization framework is a suitable scheduling solution for social content delivery.

Future work for this paper includes extending the discussion to the cases with multiple cells and time-variant reward matrix.

ACKNOWLEDGMENT

Research partially supported by grants US AFOSR MURI FA9550-09-1-0538, AFOSR MURI FA-9550-10-1-0573, NSF CNS-1035655, NIST 70NANB11H148, and DARPA contract FA8750-14-C-0019 Phase 2.

REFERENCES

- [1] K. Lin, C. Chen, and C. Chou, “Preference-aware content dissemination in opportunistic mobile social networks,” in *INFOCOM, 2012 Proceedings IEEE*, pp. 1960–1968, 2012.
- [2] E. Stai, J. Baras, and S. Papavassiliou, “Social networks over wireless networks,” in *Conference on Decision and Control (CDC), 2012 Proceedings IEEE*, 2012.
- [3] L. Tassiulas and A. Ephremides, “Stability properties of constrained queueing systems and scheduling policies for maximum throughput in multihop radio networks,” *Automatic Control, IEEE Transactions on*, vol. 37, no. 12, pp. 1936–1948, 1992.
- [4] 3GPP, “Evolved Universal Terrestrial Radio Access (E-UTRA); Radio Frequency (RF) requirements for LTE Pico Node B,” TR 36.931, 3rd Generation Partnership Project (3GPP), 9 2014.
- [5] R. Jain, A. Durreesi, and G. Babic, “Throughput fairness index: An explanation,” tech. rep., The Ohio State University, 1999.
- [6] *MovieLens dataset*, <http://www.grouplens.org/data/>, 2003.