

Exploiting Mobility in Cache-Assisted D2D Networks: Performance Analysis and Optimization

Rui Wang, Jun Zhang, *Senior Member, IEEE*, S.H. Song, *Member, IEEE*, and Khaled B. Letaief, *Fellow, IEEE*

Abstract

Caching popular content at mobile devices, accompanied by device-to-device (D2D) communications, is one promising technology for effective mobile content delivery. User mobility is an important factor when investigating such networks, which unfortunately was largely ignored in most previous works. Preliminary studies have been carried out, but the effect of mobility on the caching performance has not been fully understood. In this paper, by explicitly considering users' contact and inter-contact durations via an alternating renewal process, we first investigate the effect of mobility with a given cache placement. A tractable expression of the data offloading ratio, i.e., the proportion of requested data that can be delivered via D2D links, is derived, which is proved to be increasing with the user moving speed. The analytical results are then used to develop an effective mobility-aware caching strategy to maximize the data offloading ratio. Simulation results are provided to confirm the accuracy of the analytical results and also validate the effect of user mobility. Performance gains of the proposed mobility-aware caching strategy are demonstrated with both stochastic models and real-life data sets. It is observed that the information of the contact durations is critical to design cache placement, especially when they are relatively short or comparable to the inter-contact durations.

Index Terms

Wireless caching, device-to-device communications, human mobility, renewal process.

This work was supported by the Hong Kong Research Grants Council Grant No. 610113. Part of this work has been presented at IEEE ICC, Paris, France, May 2017 [1].

R. Wang is with Microsoft, Beijing, P. R. China (email: ruiwa@microsoft.com). This work was done when she was with HKUST. J. Zhang, S.H. Song, and K. B. Letaief are with the Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Hong Kong (email: {eejzhang, eeshsong, eekhaled}@ust.hk).

I. INTRODUCTION

The mobile data traffic is growing at an exponential rate, among which mobile video is dominant and will account for 78% of the global data traffic by 2021 [2]. To accommodate such heavy traffic, network densification is a common solution which can significantly improve the spectral efficiency and network coverage [3], [4]. However, it incurs a heavy burden on the backhaul links. Caching popular contents at helper nodes or user devices is a promising approach to reduce the backhaul data traffic, as well as, improving the user experience for video streaming applications [5]–[10]. Exploiting the predictability and reusability of popular content, caching is an effective technology for mobile content delivery, and is the key enabler for content centric wireless networks. In comparison with the commonly considered femto-caching systems [9], [11]–[15], caching at devices enjoys unique advantages. First, the aggregate cache capacity grows with the number of devices, which will subsequently increase the caching performance [8]. Second, device caching can promote device-to-device (D2D) communications, where nearby mobile devices may communicate directly rather than being forced to communicate through the base station (BS), and thus the BS load can be significantly reduced [16].

While there have been lots of studies on D2D caching networks [17]–[20], an important characteristic of mobile users, i.e., the user mobility, has been largely ignored. Specifically, a fixed topology is normally adopted by assuming users to be at fixed locations, which is not realistic. Recently, a few initial studies on mobility-aware caching have appeared [21]–[28]. It has been demonstrated in [27] that mobility-aware D2D caching can help to improve the BS offloading ratio, which, however, was only shown via numerical results. A full understanding of the effect of mobility will require a thorough theoretical analysis, which is not available yet. Moreover, there is some limitation in the mobility model adopted in previous works. For example, it was assumed that a fixed amount of data can be delivered once two users are in contact, while the variation in contact durations was not considered [24], [26]–[28]. As the user mobility will affect both the contact rate and contact duration, it is important to consider their variations when investigating the impact of user mobility, as well as, designing mobility-aware caching strategies. This forms the main objective of this paper.

A. Related Works

Caching in D2D networks has attracted lots of recent attentions. In [17], the scaling behavior of the number of D2D collaborating links was identified. Three concentration regimes, classified

by the concentration of the file popularity, were investigated. The outage-throughput trade-off and optimal scaling laws of both the throughput and outage probability were studied in [18], [19]. One main result was that, with a small file library, the throughput is proportional to the ratio of the cache capacity and file library size, while it is independent of the number of users. Two coded caching schemes, i.e., centralized and decentralized, were proposed in [20], where the contents are delivered via broadcasting. However, a fixed network topology was assumed in most previous works, which is not realistic.

There are some works considering the effect of user mobility, with different mobility models. It was shown through simulations in [27] that a higher user moving speed results in a higher data offloading ratio, while theoretical analysis is missing. In [24], a library of two files was considered while each device randomly caches one file. The coverage probability was derived when a user requests the non-cached file and moves from one location to another. Then, it was showed via numerical results that user mobility has a positive effect on D2D caching. Based on a discrete-time Markov process, the impact of user mobility was investigated in [26]. In this work, several popular locations (e.g., schools and malls) were considered, and it was assumed that users located in the same location are in contact. The throughput-delay scaling law was derived by characterizing the contact rate of the random walk model in [28]. However, in these works, it was assumed that the whole caching content or a complete encoded segment can be delivered once two users are in contact, which failed to take the variation in contact durations into account.

Some preliminary studies also evaluated the effect of user mobility by considering contact durations. In [22], Golrezaei *et al.* validated the performance of their proposed random caching scheme in the mobility scenario via simulations using a random walk model. In [29], the effect of mobility was evaluated via numerical results. It was assumed that a file can be successfully delivered if a user is in contact with another user caching the requested file, and the contact duration is enough to deliver the whole file. Since only the impact on the contact duration was considered while ignoring the contact rate, it showed that mobility has a negative impact on the hit performance. In [23], the effect of mobility was evaluated in D2D networks with coded caching, with the conclusion that mobility can improve the scaling law of throughput. In this work, the timeline was divided into discrete time slots, and it was assumed that one coded segment can be delivered in each time slot while two users may keep contact during multiple time slots. However, this result was based on the assumption that the user locations are random

and independent in each time slot, which failed to take into account the temporal correlation of user mobility.

B. Contributions

In this paper, we investigate a D2D caching network with mobile users, by adopting an alternating renewal process to model the mobility pattern so that both the contact and inter-contact durations are accounted for. Specifically, the timeline for an arbitrary pair of mobile users is divided into *contact durations*, which denote the time intervals when the mobile users are located within the transmission range, and *inter-contact durations*, which denote the time intervals between contact durations [30]. Meanwhile, both the contact and inter-contact durations are assumed to follow exponential distributions. The *data offloading ratio*, which is defined as the proportion of data that can be obtained via D2D links, is adopted as the performance metric. By theoretically analyzing the data offloading ratio, we first evaluate the effect of user mobility, and then, propose a mobility-aware caching strategy. The main contributions are summarized as follows:

- We derive an accurate expression for the data offloading ratio, for which the main difficulty is to deal with multiple alternating renewal processes. We tackle it by using a beta random variable to approximate the *communication duration* of a given user through moment matching.
- We investigate the effect of mobility for a given cache placement. In the low-to-medium mobility scenario, by assuming that the transmission rate does not change with the user speed, it is proved that the data offloading ratio increases with the user speed for any caching strategy that does not cache the same contents at all the users with contacts.
- A cache placement problem is formulated in order to maximize the data offloading ratio. By reformulating the original problem into a submodular maximization problem over a matroid constraint, a greedy algorithm is proposed, which can achieve a $\frac{1}{2}$ -approximation.
- Simulation results validate the accuracy of the derived expression, as well as the effect of user mobility. Moreover, both stochastic models and real-life data sets are used to evaluate the performance of the proposed mobility-aware caching strategy, which is shown to outperform both random and popular caching strategies. Furthermore, it is shown that the variation of contact durations is important while designing caching strategies, especially

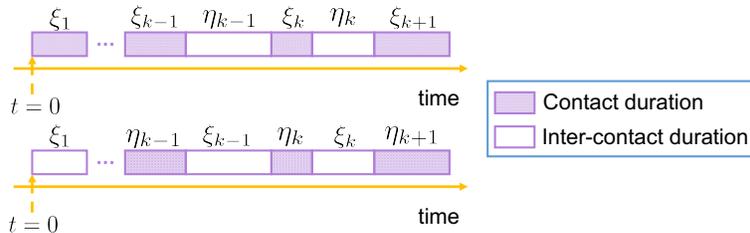


Fig. 1. The timeline for an arbitrary pair of mobile users.

when the average contact duration is relatively short or comparable to the inter-contact duration.

In comparison, our previous work [27] assumed constant contact durations, and did not provide any analytical performance evaluation. With a more realistic mobility model, the cache placement problem formulated in the current paper cannot be solved directly by the algorithm in [27]. Therefore, we improve on both performance analysis and cache placement optimization.

C. Organization

The remainder of this paper is organized as follows. In Section II, we introduce the mobility and caching models, as well as the performance metric. An approximate expression of the data offloading ratio is derived in Section III. The effect of user mobility is investigated in Section IV, and a mobility-aware caching strategy is proposed in Section V. The simulation results are shown in Section VI. Finally, Section VII concludes the paper.

II. SYSTEM MODEL AND PROBLEM FORMULATION

In this section, we first introduce the alternating renewal process to model the user mobility pattern, and discuss the caching and file delivery models. Then, the performance metric, i.e., the data offloading ratio, is defined.

A. Mobility Model

The inter-contact model, which captures the temporal correlation of the user mobility [31], is adopted to model the user mobility pattern. Specifically, the timeline of each pair of users is divided into *contact durations*, i.e., the time intervals when the users are in the transmission range, and *inter-contact durations*, i.e., the times intervals between consecutive contact durations.

Considering that contact durations and inter-contact durations appear alternatively in the timeline of a user pair, similar to [32], an alternating renewal process [33] is applied to model the pairwise contact pattern, as defined below.

Definition 1. Consider a stochastic process with state space $\{U, V\}$. The successive durations for the system to be in states U and V are denoted as $\xi_k, k = 1, 2, \dots$ and $\eta_k, k = 1, 2, \dots$, respectively, which are independent and identically distributed. Specifically, the system starts at state U and remains for ξ_1 , then switches to state V and stays for η_1 , then backs to state U and stays for ξ_2 , and so forth. Let $\psi_k = \xi_k + \eta_k$. The counting process of ψ_k is called an *alternating renewal process*.

As shown in Fig. 1, if a pair of users is in contact at $t = 0$, ξ_k and η_k represent the contact durations and inter-contact durations, respectively. Otherwise, ξ_k and η_k represent the inter-contact durations and contact durations, respectively. It was shown in [34] that exponential curves well fit the distribution of inter-contact durations, while in [35], it was identified that an exponential distribution is a good approximation for the distribution of the contact durations. Thus, the same as [32], we assume that the contact and inter-contact durations follow independent exponential distributions. For simplicity, the timelines of different user pairs are assumed to be independent. Specifically, we consider a network with N_u users, with the user index set denoted as $\mathcal{S} = \{1, 2, \dots, N_u\}$. The contact and inter-contact durations of users $i \in \mathcal{S}$ and $j \in \mathcal{S} \setminus \{i\}$ follow independent exponential distributions with parameters $\lambda_{i,j}^C$ and $\lambda_{i,j}^I$, respectively. If users $i \in \mathcal{S}$ and $j \in \mathcal{S} \setminus \{i\}$ have no contact, the parameters are $\lambda_{i,j}^C = \infty$ and $\lambda_{i,j}^I = 0$.

B. Caching and File Transmission Model

There is a library with N_f files, whose index set is denoted as $\mathcal{F} = \{1, 2, \dots, N_f\}$, each with size F . Each user device has a limited storage capacity with size C , and each file is completely cached or not cached at all at each user device. Specifically, the cache placement is denoted as

$$x_{j,f} = \begin{cases} 1, & \text{if user } j \text{ caches file } f, \\ 0, & \text{if user } j \text{ does not cache file } f, \end{cases} \quad (1)$$

where $j \in \mathcal{S}$ and $f \in \mathcal{F}$. User $i \in \mathcal{S}$ is assumed to request a file $f \in \mathcal{F}$ with probability $p_{i,f}^r$, where $\sum_{f \in \mathcal{F}} p_{i,f}^r = 1$. When a user requests a file f , it will first check its own cache, and then download the file from the users that are in contact and store file f , with a fixed transmission

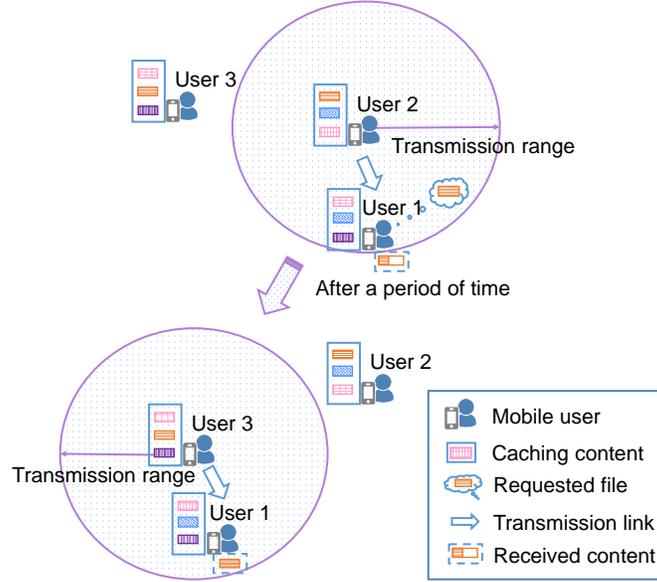


Fig. 2. A sample network with three mobile users.

rate, denoted as r_0 . At every time instant, it can only download from one other user. When a user is in contact with multiple users caching the requested file at the same time, it will randomly choose one to download. Meanwhile, D2D pairs are under the control of the base station, and orthogonal resource allocation is assumed for different simultaneous D2D communication pairs. Thus, there is no inter-user interference. We also assume that each user only requests one file at each time, and a new request is generated after it finishes downloading the previous file. If the user cannot get the whole file within a certain delay threshold, denoted as τ_0 , it will download the remaining part from the BS. We assume that the delay threshold is larger than the time required to download each content (i.e., $\tau_0 > \frac{F}{r_0}$). Fig. 2 shows a sample network, where user 1 gets part of the requested file during the contact with user 2, then gets the whole file after encountering user 3.

C. Performance Metric

The *data offloading ratio*, which is defined as the expected percentage of requested content that can be obtained via D2D links rather than downloading from the BS, is used as the performance metric. Specifically, the data offloading ratio for user $i \in \mathcal{S}$ requesting file $f \in \mathcal{F}$, with $x_{i,f} = 0$, is defined as

$$\mathcal{R}_{i,f} = \mathbb{E}_{D_{i,f}} [\min(D_{i,f}/F, 1)], \quad (2)$$

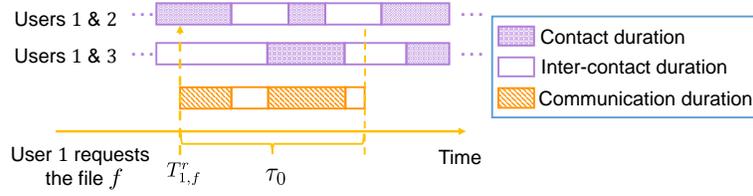


Fig. 3. An illustration of the communication duration.

where $\mathbb{E}[\cdot]$ denotes the expectation and $D_{i,f}$ denotes the amount of requested data that can be delivered via D2D links when user i requests file f . Since a fixed transmission rate is assumed, $D_{i,f}$ can be written as $D_{i,f} = r_0 \tau_{i,f}^c$, where $\tau_{i,f}^c$ is the *communication duration* for user i to download file f from other users caching file f within time τ_0 . We assume that user i can download file f while at least one user caching file f is in contact, where the handover time is ignored. Fig. 3 shows the communication duration of user 1 for the situation in Fig. 2, where $T_{1,f}^r$ is the time instant when user 1 requests a file $f \in \mathcal{F}$. Thus, we get

$$\mathcal{R}_{i,f} = \mathbb{E}_{\tau_{i,f}^c} [\min(r_0 \tau_{i,f}^c / F, 1)]. \quad (3)$$

Then, the data offloading ratio when user i requests file f is $[x_{i,f} + (1 - x_{i,f})\mathcal{R}_{i,f}]$. Considering the file request probabilities, the overall data offloading ratio is

$$\mathcal{R} = \frac{1}{N_u} \sum_{i \in \mathcal{S}} \sum_{f \in \mathcal{F}} p_{i,f}^r [x_{i,f} + (1 - x_{i,f})\mathcal{R}_{i,f}]. \quad (4)$$

While this performance metric has been used in [27] to design mobility-aware caching strategies, the variation of contact durations has not been considered. Furthermore, there is no theoretical analysis available, and it is unclear how user mobility will affect the caching performance. We shall first provide theoretical analysis in Section III for the data offloading ratio with a given cache placement. The analytical result will then be used in Section IV to reveal the effect of user moving speeds, and in Section V to optimize the caching strategy.

III. THEORETICAL ANALYSIS OF DATA OFFLOADING RATIO

The main difficulty in theoretically evaluating the data offloading ratio is to get the distribution of the communication duration. As this distribution is highly complicated, instead of deriving it directly, we develop an accurate approximation. In this section, we first approximate the distribution of the communication duration by a beta distribution, and then an approximation of the data offloading ratio is obtained.

A. Communication Duration Analysis

To assist the analysis of the communication duration, we first define some stochastic processes.

Definition 2. We model the timeline of a pair of users $i, j \in \mathcal{S}, i \neq j$ as an alternating renewal process, and denote $H_{i,j}(t), t \in [0, \infty)$ as the state at time t . Specifically, $H_{i,j}(t) = 1$ means that users i and j are in contact at time instant t . Otherwise, $H_{i,j}(t) = 0$. The durations of staying in states 1 and 0 follow independent exponential distributions with parameters $\lambda_{i,j}^C$ and $\lambda_{i,j}^I$, respectively.

Definition 3. When requesting file $f \in \mathcal{F}$, the timeline of user $i \in \mathcal{S}$ can be divided into intervals when it is in contact with at least one user caching file f , and intervals that user i is out of the transmission range of all the users caching file f . We use a random process to model the timeline of user i requesting file f , and denote $H_i^f(t), t \in [0, \infty)$ as the state at time t . Specifically, $H_i^f(t) = 1$ means that user i can download file f from a user caching file f at time instant t . Otherwise, $H_i^f(t) = 0$.

Based on Definitions 2 and 3, we can get the relationship between $H_i^f(t)$ and $H_{i,j}(t)$ as $H_i^f(t) = 1 - \prod_{j \in \mathcal{S} \setminus \{i\}, x_{j,f}=1} [1 - H_{i,j}(t)]$. Similar to [32], it is reasonable to assume that when a user requests a file, the alternating process between each pair of users has been running for a long time. Thus, denote $T_{i,f}^r, i \in \mathcal{S}$ and $f \in \mathcal{F}$, as the time instant when user i requests file f , and the corresponding communication duration $\tau_{i,f}^c$ can be derived as

$$\tau_{i,f}^c = \lim_{T_{i,f}^r \rightarrow \infty} \int_{T_{i,f}^r}^{T_{i,f}^r + \tau_0} H_i^f(t) dt. \quad (5)$$

In the following lemma, we derive the expectation and variance of the communication duration.

Lemma 1. When user $i \in \mathcal{S}$ requests file $f \in \mathcal{F}$, with $x_{i,f} = 0$, the expectation and variance of its communication duration are

$$\mathbb{E} [\tau_{i,f}^c] = \tau_0 \left(1 - \prod_{j \in \mathcal{S}, x_{j,f}=1} p_{i,j}^I \right), \quad (6)$$

and

$$\begin{aligned} \text{Var} [\tau_{i,f}^c] = & 2 \int_0^{\tau_0} (\tau_0 - u) \prod_{j \in \mathcal{S}, x_{j,f}=1} p_{i,j}^I \left[p_{i,j}^I + (1 - p_{i,j}^I) e^{-u(\lambda_{i,j}^C + \lambda_{i,j}^I)} \right] du \\ & - \tau_0^2 \prod_{j \in \mathcal{S}, x_{j,f}=1} (p_{i,j}^I)^2, \end{aligned} \quad (7)$$

where $p_{i,j}^I = \frac{\lambda_{i,j}^C}{\lambda_{i,j}^C + \lambda_{i,j}^I}$ denotes the probability that users i and j are not in contact at a time instant $t \rightarrow \infty$.

Proof. See Appendix A. □

While the mean and variance have been derived, the exact distribution of the communication duration is still highly intractable. Since it is a bounded random variable, we propose to approximate the distribution of $\tau_{i,f}^c$ by a beta distribution, which has been widely used to model random variables limited to finite ranges. For example, it has been used to model the prior distribution for a Bernoulli trial and the timestamp over a time span [36]. The pdf of a beta distribution is

$$f(x) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}, \quad (8)$$

where α and β are shape parameters, and $B(\cdot, \cdot)$ is the beta function. Specifically, if $\sum_{j \in \mathcal{S} \setminus \{i\}} x_{j,f} > 0$, which means that user i may download file f from at least one other user, we approximate $\tau_{i,f}^c/\tau_0$ by $Z_{i,f}$, where $Z_{i,f} \sim \text{Beta}(\alpha_{i,f}, \beta_{i,f})$, $i \in \mathcal{S}$ and $f \in \mathcal{F}$. Otherwise, $\tau_{i,f}^c = 0$. By matching the first two moments (i.e., $\mathbb{E}[Z_{i,f}] = \mathbb{E}[\tau_{i,f}^c/\tau_0]$ and $\text{Var}[Z_{i,f}] = \text{Var}[\tau_{i,f}^c/\tau_0]$) the parameters of the beta distribution can be derived as¹

$$\begin{cases} \alpha_{i,f} = \frac{\mathbb{E}[\tau_{i,f}^c]^2(\tau_0 - \mathbb{E}[\tau_{i,f}^c])}{\text{Var}[\tau_{i,f}^c]\tau_0} - \frac{\mathbb{E}[\tau_{i,f}^c]}{\tau_0}, \\ \beta_{i,f} = \frac{\tau_0 - \mathbb{E}[\tau_{i,f}^c]}{\mathbb{E}[\tau_{i,f}^c]} \alpha_{i,f}. \end{cases} \quad (9)$$

B. Data Offloading Ratio Analysis

Based on the above approximation, we get an approximate expression of the data offloading ratio in Proposition 1.

Proposition 1. *The data offloading ratio is*

$$\mathcal{R} = \frac{1}{N_u} \sum_{i \in \mathcal{S}} \sum_{f \in \mathcal{F}} p_{i,f}^r [x_{i,f} + (1 - x_{i,f})\mathcal{R}_{i,f}], \quad (10)$$

where $\mathcal{R}_{i,f}$ is approximated by

$$\mathcal{R}_{i,f}^a = 1 - I_{\frac{F}{\tau_0 \tau_0}} \left(\alpha_{i,f}, \frac{\tau_0 - \mathbb{E}[\tau_{i,f}^c]}{\mathbb{E}[\tau_{i,f}^c]} \alpha_{i,f} \right) + \frac{\mathbb{E}[\tau_{i,f}^c] r_0}{F} I_{\frac{F}{\tau_0 r_0}} \left(\alpha_{i,f} + 1, \frac{\tau_0 - \mathbb{E}[\tau_{i,f}^c]}{\mathbb{E}[\tau_{i,f}^c]} \alpha_{i,f} \right), \quad (11)$$

¹The parameters of the beta distribution should be positive, and it can be proved that $\alpha_{i,f} > 0$ and $\beta_{i,f} > 0$, by $e^{-u(\lambda_{i,j}^I + \lambda_{i,j}^C)} \leq 1$.

if $\sum_{j \in \mathcal{S} \setminus \{i\}} x_{j,f} > 0$, and $\mathcal{R}_{i,f}^a = 0$ elsewhere, where $I_q(\cdot, \cdot)$ is the incomplete beta function, $\mathbb{E}[\tau_{i,f}^c]$ is given in (6), and $\alpha_{i,f}$ is given in (9).

Proof. See Appendix B. □

Remark. To evaluate the data offloading ratio in Proposition 1, the main complexity is in calculating the integral in (7) when calculating the variance of the communication duration. Simulations will show that the approximation is quite accurate. This analytical result will serve as the basis for evaluating the effect of user mobility in Section IV and designing the mobility-aware caching strategy in Section V.

IV. EFFECT OF USER MOVING SPEED

In this section, we investigate how changing the user moving speed will affect the data offloading ratio for a given caching strategy. If all the user pairs $i, j \in \mathcal{S}, i \neq j$ with contacts (i.e., $0 < \lambda_{i,j}^I, \lambda_{i,j}^C < \infty$) cache the same contents, D2D communications will not help the content delivery. Thus, in the following, we assume that the contents cached at the users with contacts are not all the same. In other words, there exists a pair of users $i, j \in \mathcal{S}, i \neq j$ with $0 < \lambda_{i,j}^I, \lambda_{i,j}^C < \infty$, and a file $f \in \mathcal{F}$, such that $x_{i,f} = 1$ and $x_{j,f} = 0$. This investigation is based on the approximate expression in (10), and simulations will be provided later to verify the results.

A. Effect of Moving Speed on Communication Duration

We first characterize the relationship between the user moving speed and the parameters $\lambda_{i,j}^C$ and $\lambda_{i,j}^I$ in Lemma 2.

Lemma 2. *When all the user speeds change by μ times, the contact and inter-contact parameters will also change by μ times (i.e., from $\lambda_{i,j}^C$ and $\lambda_{i,j}^I$ to $\mu\lambda_{i,j}^C$ and $\mu\lambda_{i,j}^I$, $i, j \in \mathcal{S}$, respectively).*

Proof. The time for user i to move along a certain path L_i can be given as a curve integral $\int_{L_i} \frac{dz}{v_i(z)}$, where $v_i(z)$ is the speed of user i when passing by a point z on the path L_i . When the speed of user i changes by μ times, the time for moving along the path L_i changes to $\int_{L_i} \frac{dz}{\mu v_i(z)} = \frac{1}{\mu} \int_{L_i} \frac{dz}{v_i(z)}$, which is scaled by $\frac{1}{\mu}$ times. In each contact or inter-contact duration, users i and j move along certain paths. When all the user speeds change by μ times, each contact or inter-contact duration changes by $\frac{1}{\mu}$ times, and thus, the average values also change by $\frac{1}{\mu}$ times.

Since the contact and inter-contact durations are assumed to be exponentially distributed with means $\frac{1}{\lambda_{i,j}^C}$ and $\frac{1}{\lambda_{i,j}^I}$, respectively, the parameters $\lambda_{i,j}^C$ and $\lambda_{i,j}^I$ are scaled by μ times. \square

Considering that a larger μ means that users are moving faster, in the following, we will investigate how changing μ will affect the data offloading ratio. For simplicity, we assume that the transmission rate is a constant, and will not change with the user speed. This is reasonable in the low-to-medium mobility regime. Firstly, when the user speed changes by μ times, we rewrite the expectation and variance of communication duration as in Lemma 3.

Lemma 3. *When the user speed changes by μ times, the expectation of the communication duration when user $i \in \mathcal{S}$ requests file $f \in \mathcal{F}$ is the same as (6), and the corresponding variance is given as*

$$\text{Var}[\tau_{i,f}^c(\mu)] = 2a_0 \sum_{\mathbf{Z} \in \{0, x_{l,f}\}^{N_f} \setminus \mathbf{0}} \frac{a_{\mathbf{Z}}}{\mu \kappa_{\mathbf{Z}}} \left(\tau_0 - \frac{1}{\mu \kappa_{\mathbf{Z}}} + \frac{1}{\mu \kappa_{\mathbf{Z}}} \exp(-\mu \kappa_{\mathbf{Z}} \tau_0) \right), \quad (12)$$

where \mathbf{Z} is an N_f -ary vector, with z_l as the l -th element of \mathbf{Z} , and

$$\begin{cases} a_0 = \prod_{j \in \mathcal{S}, x_{j,f}=1} (p_{i,j}^I)^2, \\ a_{\mathbf{Z}} = \prod_{l=1}^{N_f} \left(\frac{\lambda_{i,l}^I}{\lambda_{i,l}^C} \right)^{z_l}, \\ \kappa_{\mathbf{Z}} = \sum_{l=1}^{N_f} z_l (\lambda_{i,l}^C + \lambda_{i,l}^I). \end{cases} \quad (13)$$

Proof. See Appendix C. \square

Then, the effect of user speed on the communication duration is shown in Lemma 4 .

Lemma 4. *When μ increases, the expectation of the communication duration for user $i \in \mathcal{S}$ requesting file $f \in \mathcal{F}$ (i.e., $\mathbb{E}[\tau_{i,f}^c(\mu)]$) does not change, and the corresponding variance (i.e., $\text{Var}[\tau_{i,f}^c(\mu)]$) decreases, if there is at least one user caching file f has contacts with user i . Accordingly, the parameter $\alpha_{i,f}$ of the beta distribution increases.*

Proof. See Appendix D. \square

Remark. *Intuitively, when the users move faster, they have more contacts with each other within a certain time period τ_0 , while the duration of each contact decreases. That is why the variance of communication durations decreases and the expected value does not change.*

B. Effect of Moving Speed on Data Offloading Ratio

In the following, we evaluate the relationship between $\alpha_{i,f}$ and the data offloading ratio when user i requests file f that is not in its own cache (i.e., $\mathcal{R}_{i,f}$ in (11)) in Lemma 5.

Lemma 5. *When user $i \in \mathcal{S}$ requests file $f \in \mathcal{F}$ and $x_{i,f} = 0$, the data offloading ratio (i.e., $\mathcal{R}_{i,f}$) increases with parameter $\alpha_{i,f}$.*

Proof. See Appendix E. □

Based on Lemmas 4 and 5, we specify the effect of the user speed in Proposition 2.

Proposition 2. *If the transmission rate keeps unchanged, the data offloading ratio increases with the user moving speed.*

Proof. See Appendix F. □

Remark. *The result in Proposition 2 is valid for any caching strategy, only excluding the case that all the user pairs with contacts have the same cache contents. Although a higher mobility shortens the contact durations, it provides more opportunities for the users to share data, and improves the overall data offloading ratio. This can be regarded as a type of diversity gain. Thus, it is important to take advantage of user mobility in the cache-assisted D2D network.*

V. MOBILITY-AWARE CACHING STRATEGY

We have derived the expression for the data offloading ratio and theoretically evaluated the effect of user mobility. These analytical results are also useful for designing mobility-aware caching strategies. Recently, some works started to develop cache placement strategies for D2D networks by taking advantage of user mobility information [25], [27]. However, it was assumed that a complete file or segment can be delivered when two users are in contact, and the variation in the contact duration was ignored. In this section, we formulate a cache placement problem and propose a mobility-aware caching strategy using both the contact and inter-contact information.

A. Problem Formulation

We will develop a caching strategy, which takes advantage of both the contact and inter-contact information to improve the data offloading ratio. With a higher data offloading ratio, more D2D links can be established, and thus, the spectrum efficiency can be improved. Specifically, we

consider that the contact and inter-contact parameters can be estimated via historical data. The mobility-aware cache placement problem is formulated as

$$\max_{x_{i,f}} \mathcal{R}, \quad (14)$$

$$\text{s.t. } F \sum_{f \in \mathcal{F}} x_{i,f} \leq C, i \in \mathcal{S}, \quad (14a)$$

$$x_{i,f} \in \{0, 1\}, i \in \mathcal{S}, \quad (14b)$$

where constraint (14a) implies a limited cache capacity and constraint (14b) means that each file is fully cached or not cached at all.

B. Submodular Functions and Matroid Constraints

Our proposed algorithm is based on submodular optimization. As a typical kind of combinatorial optimization problems, submodular maximization has been extensively investigated [37]–[39]. It has a wide range of applications, including the max- k -cover problem and the max-cut problem [40]. Moreover, in wireless networks, submodular maximization over a matroid constraint has been utilized to design network coding [41] and femto-caching [12]. For completeness, we will present some necessary definitions and properties related to the submodular set function and the matroid constraint. Please refer to [38] for more details.

Definition 4. Let S be a finite ground set, and a function $h : 2^S \rightarrow \mathbb{R}$ is a *submodular set function* if $h(A) + h(B) \geq h(A \cup B) + h(A \cap B)$, $\forall A \subseteq S$ and $\forall B \subseteq S$.

Proposition 3. Let $A \subset S$ and $j, k \in S - A, j \neq k$, and a function $h : 2^S \rightarrow \mathbb{R}$ is a *monotone submodular function* if $h(A \cup \{k\}) - h(A) \geq h(A \cup \{j, k\}) - h(A \cup \{j\}) \geq 0$.

Proposition 3 provides an intuitive explanation of the submodular property, i.e., the marginal gain of a monotone submodular set function decreases with a larger set. It is also useful when proving submodularity.

Definition 5. A pair $\mathcal{M} = \{S, \mathcal{I}\}$, where S is a finite ground set and \mathcal{I} is a collection of subsets of S , is called a *matroid* if

- 1) $\emptyset \in \mathcal{I}$,
- 2) If $A \subseteq B \subseteq S$ and $B \in \mathcal{I}$, then $A \in \mathcal{I}$,
- 3) If $A, B \in \mathcal{I}$ and $|B| > |A|$, then $\exists j \in B - A$ such that $A \cup \{j\} \in \mathcal{I}$.

C. Problem Reformulation

In the following, problem (14) will be reformulated as a submodular maximization problem over a matroid constraint. We first define the ground set as $S = \{y_{j,f} | j \in \mathcal{S} \text{ and } f \in \mathcal{F}\}$, and a cache placement can be represented as a subset of S . Specifically, for a cache placement $Y \subseteq S$, $y_{j,f} \in Y$ means user j caches file f (i.e., $x_{j,f} = 1$), and $y_{j,f} \notin Y$ means user j does not cache file f (i.e., $x_{j,f} = 0$). Accordingly, the data offloading ratio can be rewritten as

$$\mathcal{R}(Y) = \frac{1}{N_u} \sum_{i \in \mathcal{S}} \sum_{f \in \mathcal{F}} p_{i,f}^r [\mathbb{1}(y_{i,f} \in Y) + \mathbb{1}(y_{i,f} \notin Y) \mathcal{R}_{i,f}(Y)], \quad (15)$$

where

$$\mathcal{R}_{i,f}(Y) = \mathbb{E}_{\tau_{i,f}^c} [\min(r_0 \tau_{i,f}^c(Y)) / F, 1]. \quad (16)$$

Lemma 6 shows that $\mathcal{R}(Y)$ is a monotone submodular function by verifying Proposition 3.

Lemma 6. $\mathcal{R}(Y)$ in (15) is a monotone submodular set function on the ground set S .

Proof. See Appendix E. □

Then, Lemma 7 rewrites the constraint in problem (14) as a matroid constraint.

Lemma 7. Let $S_i = \{y_{i,f} | f \in \mathcal{F}\}$, include all files that may be cached at user i , the constraint in problem (14) is equivalent to a matroid constraint $Y \in \mathcal{I}$, where

$$\mathcal{I} = \left\{ Y \subseteq S \mid |Y \cap S_i| \leq C/F, \forall i \in \mathcal{S} \right\}. \quad (17)$$

Proof. Constraint (17) is a partition matroid, which is a typical matroid [38]. □

Thus, problem (14) can be reformulated as the following monotone submodular maximization problem over a matroid constraint.

$$\max_{Y \in \mathcal{I}} \mathcal{R}(Y) = \frac{1}{N_u} \sum_{i \in \mathcal{S}} \sum_{f \in \mathcal{F}} p_{i,f}^r [\mathbb{1}(y_{i,f} \in Y) + \mathbb{1}(y_{i,f} \notin Y) \mathcal{R}_{i,f}(Y)]. \quad (18)$$

D. Greedy Cache Placement Algorithm

To solve a monotone submodular maximization problem, a greedy algorithm provides a $\frac{1}{2}$ -approximation [38], which means that the solution is at least 50% of the optimal one. Although a randomized algorithm can get a higher approximation ratio as $(1 - 1/e)$ [40], its computational complexity is too high as the size of the ground set is $N_u \times N_f$. Thus, it is inapplicable in

practice. Moreover, it will be shown in the simulation that the greedy algorithm provides a near optimal performance.

Denote S^r as the set including all the elements that can be added into the cache placement set Y . We define the priority value of element $y_{j,f} \in S^r$ as the gain of adding $y_{j,f}$ into Y , which is given as

$$g_{j,f} = \mathcal{R}(Y \cup \{y_{j,f}\}) - \mathcal{R}(Y) \\ = \frac{1}{N_u} \left\{ p_{j,f}^r [1 - \mathcal{R}_{j,f}(Y)] + \sum_{i \in \mathcal{S} \setminus \{j\}} p_{i,f}^r \mathbb{1}(y_{i,f} \notin Y) [\mathcal{R}_{i,f}(Y \cup \{y_{j,f}\}) - \mathcal{R}_{i,f}(Y)] \right\}. \quad (19)$$

The main difficulty when developing the greedy algorithm is to efficiently evaluate the priority values, where the key is to evaluate the data offloading ratio when user i requests file f (i.e., $\mathcal{R}_{i,f}(Y)$). In Section III, we have provided an approximation for $\mathcal{R}_{i,f}(Y)$ as (11), which can be used to evaluate the priority values. The procedure of the greedy algorithm is listed in Algorithm 1. It starts from an empty set $Y = \emptyset$, which means that no file is cached. At each iteration, an element $y_{j^*,f^*} \in S^r$ with the maximum priority value is selected and added into Y . Meanwhile, y_{j^*,f^*} is excluded from S^r . If user j^* cannot cache more files, all the elements in S_{j^*} are excluded from S^r . The process continues until no more file can be cached.

Algorithm 1 The Greedy Cache Placement Algorithm

- 1: Set $Y = \emptyset \Leftrightarrow \text{set } x_{j,f} = 0, \forall j \in \mathcal{S} \text{ and } f \in \mathcal{F}$.
 - 2: $S^r = S$.
 - 3: Initialize the priority values $\{g_{j,f} = \mathcal{R}(Y \cup \{y_{j,f}\}) - \mathcal{R}(Y) | j \in \mathcal{S} \text{ and } f \in \mathcal{F}\}$.
 - 4: **while** $|Y| < \frac{N_u C}{F}$ **do**
 - 5: $y_{j^*,f^*} = \arg \max_{y_{j,f} \in S^r} g_{j,f}$.
 - 6: Set $Y = Y \cup \{y_{j^*,f^*}\} \Leftrightarrow \text{set } x_{j^*,f^*} = 1$.
 - 7: $S^r = S^r - \{y_{j^*,f^*}\}$.
 - 8: **if** $|Y \cap S_{j^*}| + 1 > C/F$ **then**
 - 9: $S^r = S^r - S_{j^*}$
 - 10: **end if**
 - 11: Update the priority values $\{g_{j,f^*} | j \in \mathcal{S}, y_{j,f^*} \in S^r\}$.
 - 12: **end while**
-

E. Computational Complexity

There are in total $N_u C/F$ iterations in Algorithm 1. At each iteration, the complexity of updating the priority value is $\mathcal{O}(N_u T_v)$, where T_v is the complexity of calculating the variance in (7), and the complexity to find the maximum priority value is $\mathcal{O}(N_f)$. Thus, the overall computational complexity of Algorithm 1 is $\mathcal{O}(N_u C/F (N_u T_v + N_f))$.

VI. SIMULATION RESULTS

In the simulation, we first validate the analytical results in Sections III and IV, and then, evaluate the performance of the proposed mobility-aware caching strategy. The content request probability is assumed to follow a Zipf distribution with parameter γ_r (i.e., $p_{i,f}^r = \frac{f^{-\gamma_r}}{\sum_{l \in \mathcal{F}} l^{-\gamma_r}}$, $i \in \mathcal{S}$ and $f \in \mathcal{F}$) [8].

A. Data Offloading Ratio

In this part, we evaluate the approximate expression of the data offloading ratio in (10) via simulations. A random caching strategy [42] is applied, where the probabilities of the contents cached at each user are proportional to the file request probabilities. The inter-contact parameters $\lambda_{i,j}^I$, $i \in \mathcal{S}, j \in \mathcal{S} \setminus \{i\}$ are generated according to a gamma distribution as $\Gamma(4.43, 1/1088)$ [43]. Similar as [32], we assume the average inter-contact durations to be 5 times as large as the average contact durations. Thus, the contact parameters are generated according to $\Gamma(4.43 \times 25, 1/1088/5)$. The theoretical results are obtained by (10), and the simulation results are obtained by randomly generating the contact and inter-contact durations according to exponential distributions.

Fig. 4 validates the accuracy of the approximation in (10) by varying the number of users. It is shown that the theoretical results are very close to the simulation results, which means that the approximate expression (10) is quite accurate. Furthermore, the data offloading ratio increases with the number of users, which is brought by the increasing aggregate cache capacity and the content sharing via D2D links. Fig. 5 validates the accuracy of the approximation in (10) by varying the file request parameter, where a larger γ_r means that the requests are more concentrated on the popular files. It is shown that the approximate expression (10) is quite accurate with different values of γ_r , and a larger value of γ_r brings a higher performance gain of caching.

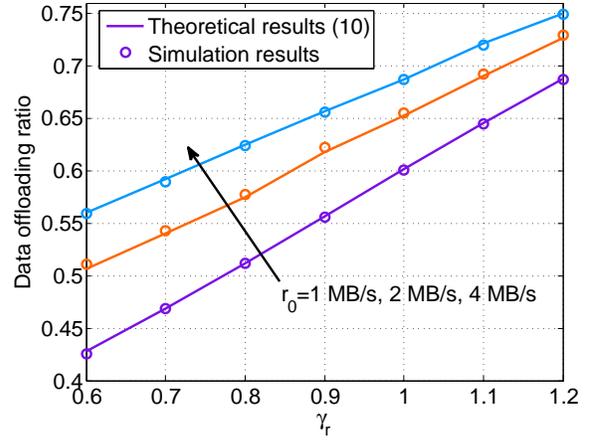
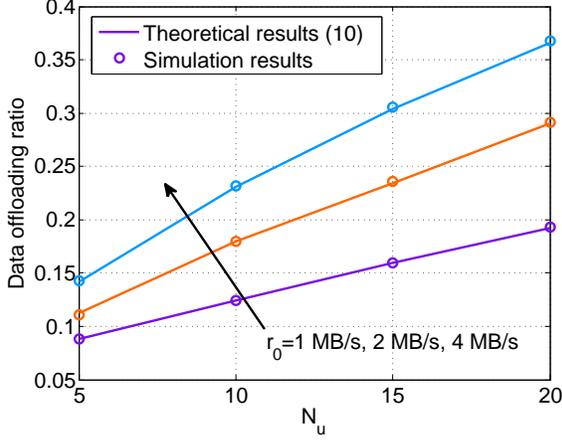


Fig. 4. Data offloading ratio with $N_f = 100$, $\tau_0 = 300$ s, $\gamma_r = 0.6$, $C = 1$ GB, $F = 300$ MB.

Fig. 5. Data offloading ratio with $N_u = 15$, $N_f = 100$, $\tau_0 = 120$ s, $C = 1$ GB, $F = 100$ MB.

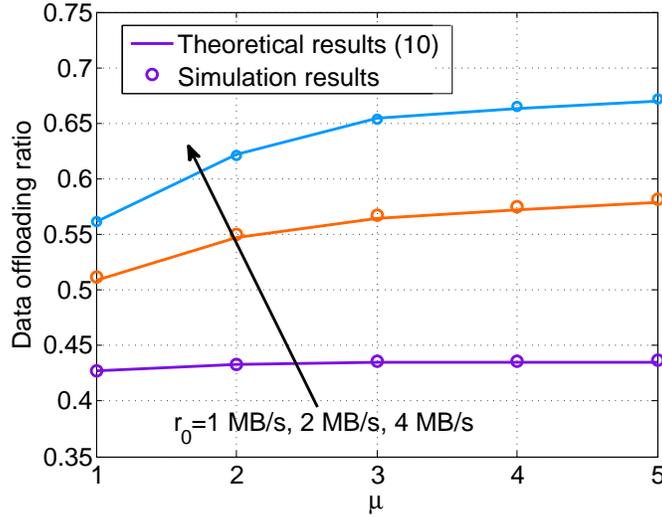


Fig. 6. Data offloading ratio with $N_u = 15$, $N_f = 100$, $\tau_0 = 120$ s, $\gamma_r = 0.6$, $C = 1$ GB, $F = 100$ MB.

B. Effect of User Mobility

In this part, we show the effect of user mobility, assuming the random caching strategy [42]. As in Section IV, we change the user moving speed by μ times. It has been shown in [43] that Gamma distribution can well fit the reciprocal of the aggregate average inter-contact time. Thus, the inter-contact parameters $\lambda_{i,j}^I, i \in \mathcal{S}, j \in \mathcal{S} \setminus \{i\}$ are generated as $\mu \widehat{\lambda}_{i,j}^I$, where $\widehat{\lambda}_{i,j}^I$ follows a gamma distribution as $\Gamma(4.43, 1/1088)$ [43]. Similarly, the contact parameters are generated as $\mu \widehat{\lambda}_{i,j}^C$, where $\widehat{\lambda}_{i,j}^C$ follows $\Gamma(4.43 \times 25, 1/1088/5)$.

In Fig. 6, the effect of μ is demonstrated. Firstly, the small gap between the theoretical and simulation results again verifies the accuracy of the approximate expression in (10). It is also shown that the data offloading ratio increases with μ , which confirms the conclusion in Proposition 2. We also observe that the increasing rate of the data offloading ratio decreases with the user moving speed, and increases with the data transmission rate.

C. Mobility-Aware Caching Strategy

In the following, we evaluate the mobility-aware caching strategy proposed in Section V. The following five caching strategies are compared.

- Optimal mobility-aware caching: The optimal solution of problem (14), which is obtained by a DP algorithm similar to Algorithm 2 in [27].²
- Greedy mobility-aware caching: The suboptimal solution of problem (14) proposed in Section V.
- Greedy mobility-aware caching ignoring contact duration: This caching strategy only utilizes the information of inter-contact durations, while assuming that the whole file can be transmitted within one contact. This strategy is obtained by the suboptimal algorithm proposed in [27].
- Random caching: Each file is randomly cached by each user, where the caching probability is proportional to the file request probability [42].
- Popular caching: All the users cache the most popular files [44].

In Fig. 7, we vary the average contact duration, denoted as \bar{t}^C , where the contact parameters $\lambda_{i,j}^C$ follow a gamma distribution, with expectation $1/\bar{t}^C$ and the same variance as $\Gamma(4.43, 1/1088)$. We observe that the performance of the greedy algorithm is quite close to the optimal one. Meanwhile, the mobility-aware caching strategies outperform the random and popular caching strategies, which demonstrates the advantage of exploiting the mobility information. Moreover, with the information of both the contact and inter-contact durations, the data offloading ratio can be further improved compared to the case ignoring the variation in contact durations. The conclusions do not change with different file request parameters. Furthermore, the gap

²After making two changes, the DP algorithm in [27] can be applied directly. The first one is to make the number of encoded segments of each file (i.e., K_f , $f \in \mathcal{N}_f$, in [27]) equals 1. The other one is to change the utility function in [27] as $U_f(x_{1,f}, \dots, x_{N_u,f}) = \frac{1}{N_u} \sum_{i \in \mathcal{S}} p_{i,f}^r [x_{i,f} + (1 - x_{i,f})\mathcal{R}_{i,f}]$.

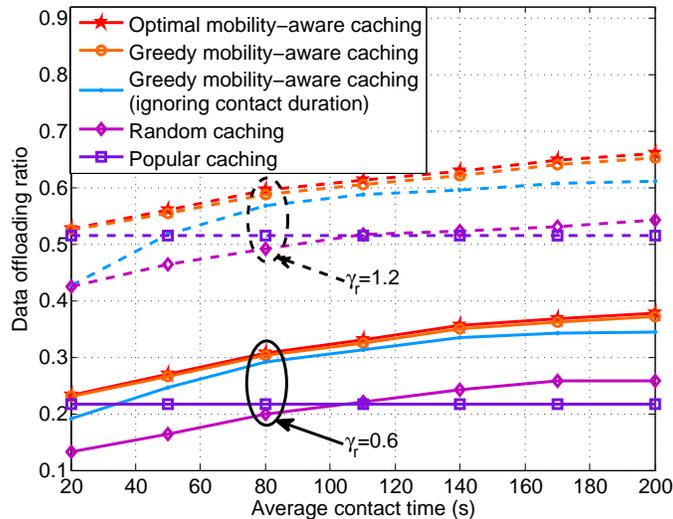


Fig. 7. Data offloading ratio with $N_u = 5$, $N_f = 50$, $\tau_0 = 300$ s, $C = 1$ GB, $F = 300$ MB, $r_0 = 1.5$ MB/s.

between the mobility-aware caching strategies considering and ignoring the contact durations first decreases and then increases with the average contact duration. When the contact duration is short, only part of a file can be delivered within one contact, and thus, it is critical to take the contact duration into account. On the other hand, when ignoring the contact duration, it is implicitly assumed that, a user is in the inter-contact duration with all the others when it requests a file. However, when the contact duration becomes larger, it is more likely that a user can start to download the file once the request is generated. Therefore, it is of critical importance to consider the contact duration, when it is comparable to the inter-contact duration. We also observe that, when the contact duration is very short, the popular caching strategy outperforms the mobility-aware caching strategy that ignores the contact duration, and approaches the proposed mobility-aware caching strategy. In such scenarios, the data shared via D2D links is quite limited, and the cached content is mainly used for a user's own need, and thus, popular caching is a good choice.

We then evaluate the performance of the proposed mobility-aware caching strategy on a real-life data set collected during the INFOCOM 2016 conference [45]. There are 78 selected conference participants, and each is distributed with one iMote, which is a Bluetooth radio device with a transmission range approximately as 30 meters. It is observed that the mobility pattern is quite different during the daytime and nighttime, i.e., users are much more frequently in contact with each other during the daytime. In the simulation, using the daytime data in the first day, we

estimate the contact and inter-contact parameters (i.e., $\lambda_{i,j}^C$ and $\lambda_{i,j}^I$) by the inverse of the average contact and inter-contact durations, respectively. Then, different caching strategies are designed using the estimated contact and inter-contact parameters, and the performance on the daytime in the second day is shown in Fig. 8. We see that the proposed mobility-aware caching strategy outperforms the popular caching and random caching strategies by 5% ~ 35% and 16% ~ 116%, respectively. The performance of the mobility-aware caching strategy that ignores the contact duration is also provided to show that the caching performance can be further improved with the information of contact durations. Moreover, the theoretical analysis of the data offloading ratio of the greedy mobility-aware caching strategy, which is calculated by (10), is also shown in Fig. 8. The small gap between the simulation results and the theoretical value calculated by (10) demonstrates that the alternating renewal process model and the approximated expression of the data offloading ratio in (10) can provide a good approximation of practical performance.

In this work, we simply the communication model for content delivery. To consider the performance in more realistic scenarios, we next investigate the effect of limited radio resources, which will limit the number of simultaneous communicating D2D pairs. We assume that each user can only serve one user's request at each time instant and there are in total 15 resource blocks allocated for D2D communications. For the case that one user receives multiple requests, it will randomly choose one to serve. When the number of D2D links is larger than the number of resource blocks, we will randomly choose 15 D2D links to transmit the requested files. Each user is assumed to make a new request right after the delay threshold of the previous request. The result is shown in Fig. 8, from which we see that the limited radio resources have very little effect on the performance. This is because the number of simultaneous D2D communications links is small, due to the randomness in user mobility.

We also validate the performance of the proposed mobility-aware caching strategy in a campus scenario, based on a data-set collected in the Cambridge campus [46]. The contact and inter-contact parameters (i.e., $\lambda_{i,j}^C$ and $\lambda_{i,j}^I$) are estimated based on the daytime data during the first three days. The performance on the fourth day is shown in Fig. 9. It again verifies the accuracy of the alternating renewal process model and the approximated expression of the data offloading ratio in (10), and shows the inconspicuous effect of limited radio resources. We also observe that the proposed mobility-aware caching strategy outperforms the popular caching and random caching strategies by 6% ~ 25% and 20% ~ 100%, respectively.

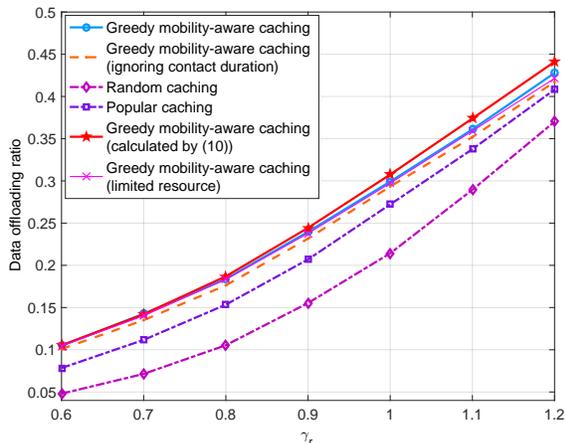


Fig. 8. Data offloading ratio with $N_u = 78$, $N_f = 500$, $\tau_0 = 300$ s, $C = 1$ GB, $F = 300$ MB, $r_0 = 2$ MB/s.

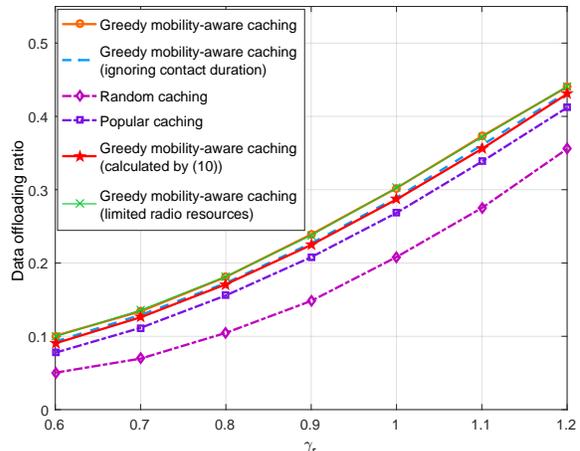


Fig. 9. Data offloading ratio with $N_u = 36$, $N_f = 500$, $\tau_0 = 600$ s, $C = 1$ GB, $F = 300$ MB, $r_0 = 1$ MB/s.

VII. CONCLUSIONS

In this paper, we investigated a D2D caching network with mobile users. A tractable expression of the data offloading ratio was firstly derived and then used to prove that the data offloading ratio increases with the user speed. This result is valid in the scenario with low-to-medium speeds, where the transmission rate does not change with the user moving speed. The extension to the case with varying transmission rates is interesting for future investigation. The analytical results were also applied to develop a mobility-aware caching strategy by utilizing the statistical contact and inter-contact information. Simulation results validated the accuracy of the approximate expression of the data offloading ratio, and demonstrated that the data offloading ratio increases with the user speed, while the increasing rate decreases with the user speed. Moreover, we also observed that it is more critical to take the contact durations into account when they are relatively short or comparable to the inter-contact durations. One future direction is to develop online caching strategies in such systems. It is also interesting to consider more sophisticated resource allocation schemes and rate adaptation, as well as smarter user selection, during the content delivery phase.

APPENDIX

A. Proof of Lemma 1

As the timelines of different user pairs are independent, the expectation of the communication duration when user i requests file f , which is not in its own cache, can be expressed as

$$\mathbb{E}[\tau_{i,f}^c] = \lim_{T_{i,f}^r \rightarrow \infty} \int_{T_{i,f}^r}^{T_{i,f}^r + \tau_0} \left[1 - \prod_{j \in \mathcal{S}, x_{j,f}=1} (1 - \mathbb{E}H_{i,j}(t)) \right] dt. \quad (20)$$

Since the timeline between each pair of users is modeled as an alternating renewal process, according to Chapter 7 in [33], we have $\lim_{t \rightarrow \infty} \Pr[H_{i,j}(t) = 1] = \frac{\lambda_{i,j}^I}{\lambda_{i,j}^C + \lambda_{i,j}^I}$, denoted as $p_{i,f}^I$. Thus, $\lim_{t \rightarrow \infty} \mathbb{E}[H_{i,j}(t)] = p_{i,f}^I$, and then, the expectation in (6) can be obtained. The variance of the communication duration is

$$\text{Var}[\tau_{i,f}^c] = 2 \lim_{T_{i,f}^r \rightarrow \infty} \int_{T_{i,f}^r}^{T_{i,f}^r + \tau_0} \int_{T_{i,f}^r}^{\theta} \Pr[H_i^f(\theta) = 1, H_i^f(t) = 1] dt d\theta - (\mathbb{E}[\tau_{i,f}^c])^2. \quad (21)$$

According to Chapter 7 in [33], $\Pr[H_{i,j}(\theta) = 0 | H_{i,j}(t) = 0] = p_{i,f}^I + p_{i,f}^I e^{-(\lambda_{i,j}^C + \lambda_{i,j}^I)(\theta-t)}$. Then, when $T_{i,f}^r \rightarrow \infty$, we get

$$\begin{aligned} & \Pr[H_i^f(\theta) = 1, H_i^f(t) = 1] \\ &= \Pr[H_i^f(\theta) = 1] - \Pr[H_i^f(\theta) = 1, H_i^f(t) = 0] \\ &= 1 - 2 \prod_{j \in \mathcal{S}, x_{j,f}=1} p_{i,f}^I + \prod_{j \in \mathcal{S}, x_{j,f}=1} p_{i,f}^I \left[p_{i,f}^I + (1 - p_{i,f}^I) e^{-(\lambda_{i,j}^C + \lambda_{i,j}^I)(\theta-t)} \right] \end{aligned} \quad (22)$$

Let $u = \theta - t$ and substitute (22) into (21), then we get (7).

B. Proof of Proposition 1

Based on the beta approximation of the communication duration, we approximate the data offloading ratio of user $i \in \mathcal{S}$ requesting file $f \in \mathcal{F}$ by

$$\mathcal{R}_{i,f}^a = \mathbb{E}_{Z_{i,f}} [\min(r_0 \tau_0 Z_{i,f} / F, 1)]. \quad (23)$$

Since $Z_{i,f} \sim \text{Beta}(\alpha_{i,f}, \beta_{i,f})$, we have

$$\begin{aligned}
\mathcal{R}_{i,f}^a &= \frac{r_0 \tau_0}{F} \int_0^{\frac{F}{r_0 \tau_0}} \frac{(z)^{\alpha_{i,f}} (1-z)^{(\beta_{i,f}-1)}}{B(\alpha_{i,f}, \beta_{i,f})} dz + \int_{\frac{F}{r_0 \tau_0}}^1 \frac{(z)^{(\alpha_{i,f}-1)} (1-z)^{(\beta_{i,f}-1)}}{B(\alpha_{i,f}, \beta_{i,f})} dz, \\
&= \frac{r_0 \tau_0}{F} \cdot \frac{\alpha_{i,f}}{\alpha_{i,f} + \beta_{i,f}} \int_0^{\frac{F}{r_0 \tau_0}} \frac{(z)^{\alpha_{i,f}} (1-z)^{(\beta_{i,f}-1)}}{B(\alpha_{i,f} + 1, \beta_{i,f})} dz \\
&\quad + \int_{\frac{F}{r_0 \tau_0}}^1 \frac{(z)^{(\alpha_{i,f}-1)} (1-z)^{(\beta_{i,f}-1)}}{B(\alpha_{i,f}, \beta_{i,f})} dz, \\
&= \frac{r_0 \tau_0}{F} \cdot \frac{\mathbb{E}[\tau_{i,f}^c]}{\tau_0} I_{\frac{F}{r_0 \tau_0}}(\alpha_{i,f} + 1, \beta_{i,f}) + 1 - I_{\frac{F}{r_0 \tau_0}}(\alpha_{i,f}, \beta_{i,f}), \\
&= 1 - I_{\frac{F}{r_0 \tau_0}} \left(\alpha_{i,f}, \frac{\tau_0 - \mathbb{E}[\tau_{i,f}^c]}{\mathbb{E}[\tau_{i,f}^c]} \alpha_{i,f} \right) + \frac{\mathbb{E}[\tau_{i,f}^c] r_0}{F} I_{\frac{F}{r_0 \tau_0}} \left(\alpha_{i,f} + 1, \frac{\tau_0 - \mathbb{E}[\tau_{i,f}^c]}{\mathbb{E}[\tau_{i,f}^c]} \alpha_{i,f} \right). \tag{24}
\end{aligned}$$

Thus, we get the expression of the data offloading ratio as shown in Proposition 1.

C. Proof of Lemma 3

When the contact and inter-contact parameters are scaled by μ , $p_{i,j}^I(\mu) = \frac{\mu \lambda_{i,j}^C}{\mu(\lambda_{i,j}^C + \lambda_{i,j}^I)} = \frac{\lambda_{i,j}^C}{\lambda_{i,j}^C + \lambda_{i,j}^I}$ remains the same. Thus the expectation in (6) does not change. Denote a_0 , $a_{\mathbf{z}}$, and $\kappa_{\mathbf{z}}$ as in (13), and we get

$$\begin{aligned}
\text{Var}[\tau_{i,f}^c] &= 2 \int_0^{\tau_0} (\tau_0 - u) \left(a_0 + \sum_{\mathbf{z} \in \{0, x_{l,f}\}^{n_f} \setminus \mathbf{0}} a_{\mathbf{z}} e^{-\kappa_{\mathbf{z}} u} \right) du - (\tau_0)^2 a_0, \\
&= 2a_0 \sum_{\mathbf{z} \in \{0, x_{l,f}\}^{n_f} \setminus \mathbf{0}} \frac{a_{\mathbf{z}}}{\kappa_{\mathbf{z}}} \left(\tau_0 - \frac{1}{\kappa_{\mathbf{z}}} + \frac{1}{\kappa_{\mathbf{z}}} \exp(-\kappa_{\mathbf{z}} \tau_0) \right). \tag{25}
\end{aligned}$$

Then, by scaling the contact and inter-contact parameters, we get the result in (12).

D. Proof of Lemma 4

As shown in Lemma 3, when the user speed changes by μ times, the expectation of the communication duration in (6) does not change, while the variance changes. To prove that $\text{Var}[\tau_{i,f}^c(\mu)]$ decreases with μ , we will prove that $\frac{\partial \text{Var}[\tau_{i,f}^c(\mu)]}{\partial \mu} < 0$. The partial derivation of $\text{Var}[\tau_{i,f}^c(\mu)]$ is

$$\frac{\partial \text{Var}[\tau_{i,f}^c(\mu)]}{\partial \mu} = 2a_0 \sum_{\mathbf{z} \in \{0, x_{l,f}\}^{n_f} \setminus \mathbf{0}} \frac{a_{\mathbf{z}}}{\mu^3 \kappa_{\mathbf{z}}} \mathcal{A}_1(x_{\mathbf{z}}), \tag{26}$$

where $\mathcal{A}_1(x_{\mathbf{z}}) = -x_{\mathbf{z}} - x_{\mathbf{z}}e^{-x_{\mathbf{z}}} - 2(e^{-x_{\mathbf{z}}} - 1)$ and $x_{\mathbf{z}} \triangleq \mu\kappa_{\mathbf{Z}}\tau_0 > 0$. Since $\mathcal{A}'_1(x_{\mathbf{z}}) = -1 + (1 + x_{\mathbf{z}})e^{-x_{\mathbf{z}}} < -1 + (1 + x_{\mathbf{z}})\frac{1}{1+x_{\mathbf{z}}} = 0$, $\mathcal{A}_1(x_{\mathbf{z}})$ is a decreasing function of $x_{\mathbf{z}}$. Thus, $\mathcal{A}_1(x_{\mathbf{z}}) < \mathcal{A}_1(0) = 0$. According to (26), when $\exists 0 < \lambda_{i,l}^C, \lambda_{i,l}^I < \infty$ and $x_{l,f} = 1$, where $l \in \mathcal{S}$, there exists $0 < \kappa_{\mathbf{Z}} < \infty$, where $\mathbf{Z} \in \{0, x_{l,f}\}^{n_f} \setminus \mathbf{0}$. Thus, we have $\frac{\partial \text{Var}[\tau_{i,f}^c(\mu)]}{\partial \mu} < 0$. The parameter $\alpha_{i,f}$ given in (9) is a decreasing function of $\text{Var}[\tau_{i,f}^c(\mu)]$, and thus, it increases with μ .

E. Proof of Lemma 5

To simplify the expression in (11), denote $q \triangleq \frac{F}{\tau_0 r_0} \in (0, 1)$, $y \triangleq \frac{\tau_0 - \mathbb{E}[\tau_{i,f}^c]}{\mathbb{E}[\tau_{i,f}^c]} \geq 0$, and $\alpha \triangleq \alpha_{i,f}$. The expression in (11) can be rewritten as a function of α as

$$\mathcal{R}_{i,f} = 1 - \frac{\int_0^q (1 - \frac{u}{q})u^{\alpha-1}(1-u)^{y\alpha-1} du}{B(\alpha, y\alpha)}, \quad (27)$$

where $B(\cdot, \cdot)$ is the beta function. Let $g(\alpha) = 1 - \mathcal{R}_{i,f}$ with the derivative of $g(\alpha)$ given by

$$g'(\alpha) = \frac{1}{B(\alpha, y\alpha)} \left\{ \int_0^q (1 - \frac{u}{q})u^{\alpha-1}(1-u)^{y\alpha-1} [\ln u + y \ln(1-u)] du - \int_0^q (1 - \frac{u}{q})u^{\alpha-1}(1-u)^{y\alpha-1} du D(y, \alpha) \right\}, \quad (28)$$

where $D(y, \alpha) = \psi(\alpha) + y\psi(y\alpha) - (1+y)\psi[(1+y)\alpha]$ and $\psi(\cdot)$ is the digamma function [47]. If $q = 1$, $g'(\alpha) = \frac{\partial [y/(1+y)]}{\partial \alpha} = 0$. Denote $\mathcal{A}_2(q) = \frac{B(\alpha, y\alpha)}{q} g'(\alpha)$, and then we have $\mathcal{A}_2(1) = 0$ and

$$\lim_{q \rightarrow 0^+} \mathcal{A}_2(q) = \lim_{q \rightarrow 0^+} \int_0^q (q-u)u^{\alpha-1}(1-u)^{y\alpha-1} [\ln u + y \ln(1-u)] du \quad (29)$$

Since $q \geq u \geq 0$ and $y \geq 0$, $(q-u)u^{\alpha-1}(1-u)^{y\alpha-1} \geq 0$ and $\ln u + y \ln(1-u) \leq 0$, thus,

$\lim_{q \rightarrow 0^+} \mathcal{A}_2(q) \leq 0$. The derivative of $\mathcal{A}_2(q)$ is

$$\mathcal{A}'_2(q) = \int_0^q u^{\alpha-1}(1-u)^{y\alpha-1} [\ln u + y \ln(1-u)] du - \int_0^q u^{\alpha-1}(1-u)^{y\alpha-1} du D(y, \alpha). \quad (30)$$

Thus, $\mathcal{A}'_2(1) = \frac{\partial B(\alpha, y\alpha)}{\partial \alpha} - \frac{\partial B(\alpha, y\alpha)}{\partial \alpha} = 0$ and $\lim_{q \rightarrow 0^+} \mathcal{A}'_2(q) \leq 0$. Then, we get $\mathcal{A}''_2(q) = q^{\alpha-1}(1-q)^{y\alpha-1} [\ln q + y \ln(1-q) - D(y, \alpha)]$. Let $\mathcal{A}_3(q) = q^{1-\alpha}(1-q)^{1-y\alpha} \mathcal{A}''_2(q)$, then, there is one zero point of $\mathcal{A}'_3(q) = \frac{1-(1+y)q}{q(1-q)}$ in $(0, 1]$. Thus, there is one inflection point of $\mathcal{A}_3(q)$. Considering that $\lim_{q \rightarrow 0^+} \mathcal{A}_3(q) = \lim_{q \rightarrow 1^-} \mathcal{A}_3(q) = -\infty$, the sign of $\mathcal{A}_3(q)$ may be negative, or first negative, then positive, and then negative, when q increases in $(0, 1)$. If $\mathcal{A}_3(q) < 0$, then $\mathcal{A}''_2(q) < 0$, when $q \in (0, 1)$. However, we have $\lim_{q \rightarrow 0^+} \mathcal{A}'_2(q) \leq \mathcal{A}'_2(1)$, which means that $\mathcal{A}'_2(q)$ cannot be a decreasing

function in $(0, 1)$. Thus, the sign of $\mathcal{A}_3(q)$ is first negative, then positive, and then negative, when q increases in $(0, 1)$. Since $\mathcal{A}_2''(q)$ has the same sign with $\mathcal{A}_3(q)$ in $(0, 1)$, $\mathcal{A}_2'(q)$ first decreases, then increases, and then decreases when q increases in $(0, 1)$. Considering that $\lim_{q \rightarrow 0^+} \mathcal{A}_2'(q) \leq 0$ and $\mathcal{A}_2'(1) = 0$, the sign of $\mathcal{A}_2'(q)$ must be first negative, and then positive in $(0, 1)$. Therefore, when q increases in $(0, 1)$, $\mathcal{A}_2(q)$ first decreases, and then increases. Considering that $\lim_{q \rightarrow 0^+} \mathcal{A}_2(q) \leq 0$ and $\mathcal{A}_2(1) = 0$, we have $\mathcal{A}_2(q) < 0$ in $(0, 1)$. Since $g'(\alpha) = \frac{q}{B(\alpha, y\alpha)} \mathcal{A}_2(q)$, we get $g'(\alpha) < 0$ in $(0, 1)$. Thus, $g(\alpha)$ decreases with α , and $\mathcal{R}_{i,f} = 1 - g(\alpha)$ increases with α .

F. Proof of Proposition 2

The data offloading ratio in (10) increases with $\mathcal{R}_{i,f}$ if $x_{i,f} = 0$, $i \in \mathcal{S}$, $f \in \mathcal{F}$. Thus, based on Lemmas 4 and 5, we get that the data offloading ratio when user i requests file f from other users (i.e., $\mathcal{R}_{i,f}$) increases with the user speed when $\exists j \in \mathcal{S}$ such that $0 < \lambda_{i,j}^C, \lambda_{i,j}^I < \infty$ and $x_{j,f} = 1$. Otherwise, $\mathcal{R}_{i,f} = 0$. Accordingly, the data offloading ratio when user i requests file f (i.e., $x_{i,f} + (1 - x_{i,f})\mathcal{R}_{i,f}$) increases with the user speed when $x_{i,f} = 0$, and $\exists j \in \mathcal{S}$ such that $0 < \lambda_{i,j}^C, \lambda_{i,j}^I < \infty$ and $x_{j,f} = 1$. Otherwise, it remains the same. Since we consider that there exists a pair of users $i, j \in \mathcal{S}$, $i \neq j$ with $0 < \lambda_{i,j}^I, \lambda_{i,j}^C < \infty$, and a file $f \in \mathcal{F}$, such that $x_{i,f} = 1$ and $x_{j,f} = 0$, the data offloading ratio increases with the user speed.

G. Proof of Lemma 6

To prove $\mathcal{R}(Y)$ is a monotone submodular set function, we will prove that it satisfies Proposition 3. In the following, we will first prove that $\mathcal{R}_{i,f}(Y)$ in (16) satisfies Proposition 3. Let $A \subseteq \mathcal{S}$, $y_{j_1,f}, y_{j_2,f} \in \mathcal{S} - A$, and $j_1 \neq j_2$, and then we have

$$\begin{aligned} & \mathcal{R}_{i,f}(A \cup \{y_{j_2,f}\}) - \mathcal{R}_{i,f}(A) \\ &= \mathbb{E} \left[\min \left(r_0 \tau_{i,f}^c(A \cup \{y_{j_2,f}\}) / F, 1 \right) \right] - \mathbb{E} \left[\min \left(r_0 \tau_{i,f}^c(A) / F, 1 \right) \right] \\ &= \mathbb{E} \left[\min \left(r_0 \tau_{i,f}^c(A) / F + r_0 D_2 / F, 1 \right) - \min \left(r_0 \tau_{i,f}^c(A) / F, 1 \right) \right], \end{aligned} \quad (31)$$

where D_2 is the duration that user i is in contact with user j_2 , and does not in contact with the users in A , given by

$$D_2 = \lim_{T^r \rightarrow \infty} \int_{T^r}^{T^r + \tau_0} \mathbb{1} [H_{i,j}(t) = 0, \forall y_{j,f} \in A, \text{ and } H_{i,j_2}(t) = 1] dt. \quad (32)$$

Then, we get

$$\mathcal{R}_{i,f}(A \cup \{y_{j_2,f}\}) - \mathcal{R}_{i,f}(A) = \begin{cases} \mathbb{E}[r_0 D_2 / F] & \text{if } \tau_{i,f}^c(A) < F/r_0 - D_2, \\ \mathbb{E}[1 - r_0 \tau_{i,f}^c(A) / F] & \text{if } F/r_0 - D_2 \leq \tau_{i,f}^c(A) \leq F/r_0, \\ 0 & \text{if } \tau_{i,f}^c(A) > F/r_0. \end{cases} \quad (33)$$

Thus, $\mathcal{R}_{i,f}(A \cup \{y_{j_2,f}\}) - \mathcal{R}_{i,f}(A) \geq 0$ and $\mathcal{R}_{i,f}(A \cup \{y_{j_1,f}\}) - \mathcal{R}_{i,f}(A) \geq 0$, i.e., $\mathcal{R}_{i,f}(Y)$ is monotone. Similarly, we have

$$\begin{aligned} & \mathcal{R}_{i,f}(A \cup \{y_{j_1,f}, y_{j_2,f}\}) - \mathcal{R}_{i,f}(A \cup \{y_{j_1,f}\}) \\ &= \begin{cases} \mathbb{E}[r_0 D'_2 / F] & \text{if } \tau_{i,f}^c(A) < F/r_0 - D'_2 - D_1, \\ \mathbb{E}[1 - r_0 \tau_{i,f}^c(A) / F - r_0 D_1 / F] & \text{if } F/r_0 - D'_2 - D_1 \leq \tau_{i,f}^c(A) \leq F/r_0 - D_1, \\ 0 & \text{if } \tau_{i,f}^c(A) > F/r_0 - D_1, \end{cases} \quad (34) \end{aligned}$$

where D_1 is the duration that user i is in contact with user j_1 , and does not in contact with the users in A , given by

$$D_1 = \lim_{T^r \rightarrow \infty} \int_{T^r}^{T^r + \tau_0} \mathbb{1}[H_{i,j}(t) = 0, \forall y_{j,f} \in A, \text{ and } H_{i,j_1}(t) = 1] dt, \quad (35)$$

and D'_2 is the duration that user i is in contact with user j_2 , and does not in contact with the users in $A \cup \{y_{j_1,f}\}$, given by

$$D'_2 = \lim_{T^r \rightarrow \infty} \int_{T^r}^{T^r + \tau_0} \mathbb{1}[H_{i,j}(t) = 0, \forall y_{j,f} \in A \cup \{y_{j_1,f}\}, \text{ and } H_{i,j_2}(t) = 1] dt. \quad (36)$$

The following inequalities reveal the relationship among D_1 , D_2 , and D'_2 . Frist, we have

$$D_2 - D'_2 = \lim_{T^r \rightarrow \infty} \int_{T^r}^{T^r + \tau_0} \mathbb{1}[H_{i,j}(t) = 0, \forall y_{j,f} \in A, \text{ and } H_{i,j_2}(t) = H_{i,j_1}(t) = 1] dt \geq 0, \quad (37)$$

which is the duration that user i is in contact with user j_1 and j_2 , and does not in contact with the users in A . We also have

$$D_1 + D'_2 - D_2 = \lim_{T^r \rightarrow \infty} \int_{T^r}^{T^r + \tau_0} \mathbb{1}[H_{i,j}(t) = 0, \forall y_{j,f} \in A \text{ and } H_{i,j_1}(t) = 1 \text{ and } H_{i,j_2}(t) = 0] dt \geq 0, \quad (38)$$

which is the duration that user i is in contact with user j_1 , and does not in contact with the users in $A \cup \{y_{j_2, f}\}$. Let $\Delta = [\mathcal{R}_{i, f}(A \cup \{y_{j_2, f}\}) - \mathcal{R}_{i, f}(A)] - [\mathcal{R}_{i, f}(A \cup \{y_{j_1, f}, y_{j_2, f}\}) - \mathcal{R}_{i, f}(A \cup \{y_{j_1, f}\})]$, and we get

$$\Delta = \begin{cases} \mathbb{E}[r_0(D_2 - D'_2)/F] & \text{if } \tau_{i, f}^c(A) < F/r_0 - D'_2 - D_1, \\ \mathbb{E}[r_0(\tau_{i, f}^c(A) + D_1 + D_2)/F - 1] & \text{if } F/r_0 - D'_2 - D_1 \leq \tau_{i, f}^c(A) < F/r_0 - D_2, \\ \mathbb{E}[r_0 D_1/F] & \text{if } F/r_0 - D_2 \leq \tau_{i, f}^c(A) < F/r_0 - D_1, \\ \mathbb{E}[1 - r_0 \tau_{i, f}^c(A)/F] & \text{if } F/r_0 - \min(D_1, D_2) \leq \tau_{i, f}^c(A) < F/r_0, \\ 0 & \text{if } \tau_{i, f}^c(A) \geq F/r_0. \end{cases} \quad (39)$$

Thus, when $\tau_{i, f}^c(A) < F/r_0 - D'_2 - D_1$, we have $\Delta = \mathbb{E}[r_0(D_2 - D'_2)/F] \geq 0$. When $F/r_0 - D'_2 - D_1 \leq \tau_{i, f}^c(A) < F/r_0 - D_2$, we have

$$\begin{aligned} \Delta &= \mathbb{E}[r_0(\tau_{i, f}^c(A) + D_1 + D_2)/F - 1] \\ &\geq \mathbb{E}[r_0(\tau_{i, f}^c(A) + D_1 + D'_2)/F - 1] \\ &\geq \mathbb{E}[r_0(F/r_0 - D'_2 - D_1 + D_1 + D'_2)/F - 1] \\ &= 0. \end{aligned} \quad (40)$$

When $F/r_0 - \min(D_1, D_2) \leq \tau_{i, f}^c(A) < F/r_0$, we have $\Delta = \mathbb{E}[1 - r_0 \tau_{i, f}^c(A)/F] \geq 0$. Thus, we get $\Delta \geq 0$, which means that

$$\mathcal{R}_{i, f}(A \cup \{y_{j_2, f}\}) - \mathcal{R}_{i, f}(A) \geq \mathcal{R}_{i, f}(A \cup \{y_{j_1, f}, y_{j_2, f}\}) - \mathcal{R}_{i, f}(A \cup \{y_{j_1, f}\}). \quad (41)$$

Then, we will prove that $\mathcal{R}(Y)$ satisfies Proposition 3. Let $A \subseteq S$, $y_{j_1, f_1}, y_{j_2, f_2} \in S - A$, and $y_{j_1, f_1} \neq y_{j_2, f_2}$, and we first consider the case that $f_1 = f_2 = f$, which gives

$$\begin{aligned} &\mathcal{R}(A \cup \{y_{j_2, f}\}) - \mathcal{R}(A) \\ &= \frac{1}{N_u} \left\{ p_{j_2, f}^r [1 - \mathcal{R}_{j_2, f}(A)] + \sum_{i \in S \setminus \{j_2\}} p_{i, f}^r \mathbb{1}(y_{i, f} \notin A) [\mathcal{R}_{i, f}(A \cup \{y_{j_2, f}\}) - \mathcal{R}_{i, f}(A)] \right\}, \\ &= \frac{1}{N_u} \left\{ p_{j_2, f}^r [1 - \mathcal{R}_{j_2, f}(A)] + \sum_{i \in S \setminus \{j_1, j_2\}} p_{i, f}^r \mathbb{1}(y_{i, f} \notin A) [\mathcal{R}_{i, f}(A \cup \{y_{j_2, f}\}) - \mathcal{R}_{i, f}(A)] \right. \\ &\quad \left. + p_{j_1, f}^r [\mathcal{R}_{j_1, f}(A \cup \{y_{j_2, f}\}) - \mathcal{R}_{j_1, f}(A)] \right\}, \end{aligned}$$

$$\begin{aligned}
&\geq \frac{1}{N_u} \left\{ p_{j_2, f}^r [1 - \mathcal{R}_{j_2, f}(A \cup \{y_{j_1, f}\})] \right. \\
&\quad \left. + \sum_{i \in \mathcal{S} \setminus \{j_1, j_2\}} p_{i, f}^r \mathbb{1}(y_{i, f} \notin A) [\mathcal{R}_{i, f}(A \cup \{y_{j_1, f}, y_{j_2, f}\}) - \mathcal{R}_{i, f}(A \cup \{y_{j_1, f}\})] \right\}, \\
&= \mathcal{R}(A \cup \{y_{j_1, f}, y_{j_2, f}\}) - \mathcal{R}(A \cup \{y_{j_1, f}\}) \geq 0.
\end{aligned} \tag{42}$$

For the case $f_1 \neq f_2$, we have

$$\begin{aligned}
&\mathcal{R}(A \cup \{y_{j_2, f_2}\}) - \mathcal{R}(A) \\
&= \frac{1}{N_u} \left\{ p_{j_2, f_2}^r [1 - \mathcal{R}_{j_2, f_2}(A)] + \sum_{i \in \mathcal{S} \setminus \{j_2\}} p_{i, f_2}^r \mathbb{1}(y_{i, f_2} \notin A) [\mathcal{R}_{i, f_2}(A \cup \{y_{j_2, f_2}\}) - \mathcal{R}_{i, f_2}(A)] \right\}, \\
&= \mathcal{R}(A \cup \{y_{j_1, f_1}, y_{j_2, f_2}\}) - \mathcal{R}(A \cup \{y_{j_1, f_1}\}) \geq 0.
\end{aligned} \tag{43}$$

Accordingly, $\mathcal{R}(Y)$ satisfies Proposition 3, and thus, it is a monotone submodular set function.

REFERENCES

- [1] R. Wang, J. Zhang, S. Song, and K. B. Letaief, "Mobility increases the data offloading ratio in D2D caching networks," in *Proc of IEEE Int. Conf. Commun. (ICC)*, Paris, France, May 2017.
- [2] Cisco Systems Inc., "Cisco visual networking index: Global mobile data traffic forecast update," *White Paper*, Mar. 2017.
- [3] J. G. Andrews, S. Buzzi, W. Choi, S. V. Hanly, A. Lozano, A. C. Soong, and J. C. Zhang, "What will 5G be?" *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1065–1082, Jun. 2014.
- [4] N. Bhushan, J. Li, D. Malladi, R. Gilmore, D. Brenner, A. Damnjanovic, R. Sukhavasi, C. Patel, and S. Geirhofer, "Network densification: the dominant theme for wireless evolution into 5G," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 82–89, Feb. 2014.
- [5] N. Golrezaei, A. F. Molisch, A. G. Dimakis, and G. Caire, "Femtocaching and device-to-device collaboration: A new architecture for wireless video distribution," *IEEE Commun. Mag.*, vol. 51, no. 4, pp. 142–149, Apr. 2013.
- [6] X. Wang, M. Chen, T. Taleb, A. Ksentini, and V. Leung, "Cache in the air: Exploiting content caching and delivery techniques for 5G systems," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 131–139, Feb. 2014.
- [7] X. Peng, J.-C. Shen, J. Zhang, and K. B. Letaief, "Backhaul-aware caching placement for wireless networks," in *Proc. IEEE Global Commun. Conf. (Globecom)*, San Diego, CA, Dec. 2015.
- [8] M. Ji, G. Caire, and A. F. Molisch, "Wireless device-to-device caching networks: Basic principles and system performance," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 99, pp. 176–189, Jul. 2015.
- [9] J. Liu, B. Bai, J. Zhang, and K. B. Letaief, "Cache placement in Fog-RANs: from centralized to distributed algorithms," vol. 16, no. 11, pp. 7039–7051, Aug. 2017.
- [10] M. A. Maddah-Ali and U. Niesen, "Decentralized coded caching attains order-optimal memory-rate tradeoff," *IEEE/ACM Trans. Netw.*, vol. 23, no. 4, pp. 1029–1040, Aug. 2015.
- [11] K. Poularakis and L. Tassiulas, "Exploiting user mobility for wireless content delivery," in *Proc. IEEE Int. Symp. Information Theory (ISIT)*, Istanbul, Turkey, Jul. 2013.

- [12] K. Shanmugam, N. Golrezaei, A. G. Dimakis, A. F. Molisch, and G. Caire, “Femtocaching: Wireless content delivery through distributed caching helpers,” *IEEE Trans. Inf. Theory*, vol. 59, no. 12, pp. 8402–8413, Sep. 2013.
- [13] N. Golrezaei, K. Shanmugam, A. G. Dimakis, A. F. Molisch, and G. Caire, “Femtocaching: Wireless video content delivery through distributed caching helpers,” in *Proc. of IEEE Int. Conf. on Computer Commun. (INFOCOM)*, Orlando, FL, Mar. 2012.
- [14] E. Bastug, M. Bennis, M. Kountouris, and M. Debbah, “Cache-enabled small cell networks: Modeling and tradeoffs,” *EURASIP J. on Wireless Commun. and Netw.*, vol. 2015, no. 1, pp. 1–11, Feb. 2015.
- [15] J. Li, Y. Chen, Z. Lin, W. Chen, B. Vucetic, and L. Hanzo, “Distributed caching for data dissemination in the downlink of heterogeneous networks,” *IEEE Trans. Commun.*, vol. 63, no. 10, pp. 3553–3568, Oct. 2015.
- [16] G. Fodor, E. Dahlman, G. Mildh, S. Parkvall, N. Reider, G. Miklós, and Z. Turányi, “Design aspects of network assisted device-to-device communications,” *IEEE Commun. Mag.*, vol. 50, no. 3, pp. 170–177, Mar. 2012.
- [17] N. Golrezaei, A. G. Dimakis, and A. F. Molisch, “Scaling behavior for device-to-device communications with distributed caching,” *IEEE Trans. Inf. Theory*, vol. 60, no. 7, pp. 4286–4298, Apr. 2014.
- [18] M. Ji, G. Caire, and A. F. Molisch, “Optimal throughput-outage trade-off in wireless one-hop caching networks,” in *Proc. IEEE Int. Symp. Information Theory (ISIT)*, Istanbul, Jul. 2013.
- [19] —, “The throughput-outage tradeoff of wireless one-hop caching networks,” *IEEE Trans. Inf. Theory*, vol. 61, no. 12, pp. 6833–6859, Oct. 2015.
- [20] M. Ji, G. Caire, and A. Molisch, “Fundamental limits of caching in wireless D2D networks,” *IEEE Trans. Inf. Theory*, vol. 62, no. 2, pp. 849 – 869, Feb. 2016.
- [21] R. Wang, X. Peng, J. Zhang, and K. B. Letaief, “Mobility-aware caching for content-centric wireless networks: Modeling and methodology,” *IEEE Commun. Mag.*, vol. 54, no. 8, pp. 77 – 83, Aug. 2016.
- [22] N. Golrezaei, P. Mansourifard, A. F. Molisch, and A. G. Dimakis, “Base-station assisted device-to-device communications for high-throughput wireless video networks,” *IEEE Trans. Wireless Commun.*, vol. 13, no. 7, pp. 3665–3676, Apr. 2014.
- [23] A. Shabani, S. P. Shariatpanahi, V. Shah-Mansouri, and A. Khonsari, “Mobility increases throughput of wireless device-to-device networks with coded caching,” in *Proc. IEEE Int. Conf. on Commun. (ICC)*, Kuala Lumpur, Malaysia, May 2016.
- [24] S. Krishnan and H. S. Dhillon, “Effect of user mobility on the performance of device-to-device networks with distributed caching,” *IEEE Commun. Lett.*, vol. 6, no. 2, pp. 194–197, Apr. 2017.
- [25] S. Hosny, A. Eryilmaz, A. A. Abouzeid, and H. El Gamal, “Mobility-aware centralized D2D caching networks,” in *Proc. Annual Allerton Conference on Commun., Control, and Comput.*, Monticello, IL, USA, Sep. 2016.
- [26] S. Hosny, A. Eryilmaz, and H. El Gamal, “Impact of user mobility on D2D caching networks,” in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, Washington, DC, USA, Dec. 2016.
- [27] R. Wang, J. Zhang, S. Song, and K. B. Letaief, “Mobility-aware caching in D2D networks,” *IEEE Trans. Wireless Commun.*, vol. 16, no. 8, pp. 5001 – 5015, Aug. 2017.
- [28] G. Alfano, M. Garetto, and E. Leonardi, “Content-centric wireless networks with limited buffers: When mobility hurts,” *IEEE/ACM Trans. Netw.*, vol. 24, no. 1, pp. 299–311, Feb. 2016.
- [29] C. Jarray and A. Giovanidis, “The effects of mobility on the hit performance of cached D2D networks,” in *Proc. IEEE WiOpt*, Tempe, AZ, USA, May 2016.
- [30] P. Hui, A. Chaintreau, J. Scott, R. Gass, J. Crowcroft, and C. Diot, “Pocket switched networks and human mobility in conference environments,” in *Proc. ACM Special Interest Group on Data Commun. (SIGCOMM) Workshop*, Philadelphia, PA, Aug. 2005.

- [31] V. Conan, J. Leguay, and T. Friedman, "Fixed point opportunistic routing in delay tolerant networks," *IEEE J. Sel. Areas Commun.*, vol. 26, no. 5, pp. 773–782, Jun. 2008.
- [32] Y. Li and W. Wang, "Can mobile cloudlets support mobile applications?" in *Proc. IEEE Int. Conf. on Computer Commun. (INFOCOM)*, Toronto, Canada, Apr. 2014.
- [33] M. Rausand and A. Høyland, *System reliability theory: Models, statistical methods, and applications*. John Wiley & Sons, 2004.
- [34] V. Conan, J. Leguay, and T. Friedman, "Characterizing pairwise inter-contact patterns in delay tolerant networks," in *Proc. Int. Conf. on Autonomic Computing and Commun. Syst.*, Rome, Italy, 2007.
- [35] M. Zhao, Y. Li, and W. Wang, "Modeling and analytical study of link properties in multihop wireless networks," *IEEE Trans. Commun.*, vol. 60, no. 2, pp. 445 – 455, Feb. 2012.
- [36] X. Wang and A. McCallum, "Topics over time: a non-markov continuous-time model of topical trends," in *Proc. ACM SIGKDD*, PA, USA, Aug. 2006.
- [37] S. Fujishige, *Submodular functions and optimization*. Elsevier, 2005.
- [38] M. L. Fisher, G. L. Nemhauser, and L. A. Wolsey, "An analysis of approximations for maximizing submodular set functions-I," *Math. Prog.*, vol. 14, no. 1, pp. 265–294, Dec. 1978.
- [39] A. Schrijver, *Combinatorial optimization: Polyhedra and efficiency*. Springer Science & Business Media, 2003.
- [40] G. Calinescu, C. Chekuri, M. Pál, and J. Vondrák, "Maximizing a monotone submodular function subject to a matroid constraint," *SIAM J. Comput.*, vol. 40, no. 6, pp. 1740–1766, Dec. 2011.
- [41] S. El Rouayheb, A. Sprintson, and C. Georghiades, "On the index coding problem and its relation to network coding and matroid theory," *IEEE Trans. Inf. Theory*, vol. 56, no. 7, pp. 3187–3195, Jul. 2010.
- [42] B. Blaszczyszyn and A. Giovanidis, "Optimal geographic caching in cellular networks," in *Proc. IEEE Int. Conf. on Commun. (ICC)*, London, UK, Jun. 2015.
- [43] A. Passarella and M. Conti, "Analysis of individual pair and aggregate intercontact times in heterogeneous opportunistic networks," *IEEE Trans. Mobile Comput.*, vol. 12, no. 12, pp. 2483–2495, Oct. 2013.
- [44] H. Ahlehagh and S. Dey, "Video-aware scheduling and caching in the radio access network," *IEEE/ACM Trans. Netw.*, vol. 22, no. 5, pp. 1444–1462, Oct. 2014.
- [45] A. Chaintreau, P. Hui, J. Scott, R. Gass, J. Crowcroft, and C. Diot, "Impact of human mobility on opportunistic forwarding algorithms," *IEEE Trans. Mobile Comput.*, vol. 6, no. 6, pp. 606–620, Jun. 2007.
- [46] J. Leguay, A. Lindgren, J. Scott, T. Friedman, and J. Crowcroft, "Opportunistic content distribution in an urban setting," in *Proc. ACM Special Interest Group on Data Commun. (SIGCOMM) Workshop*, Pisa, Italy, Sep. 2006.
- [47] M. Abramowitz and I. A. Stegun, *Handbook of mathematical functions: with formulas, graphs, and mathematical tables*. Courier Corporation, 1964.