# Probabilistic Caching in Wireless D2D Networks: Cache Hit Optimal Versus Throughput Optimal

Zheng Chen, Nikolaos Pappas and Marios Kountouris

**Journal Article**

N.B.: When citing this work, cite the original article.

LINKÖPINGS UNIVERSITET

# Probabilistic Caching in Wireless D2D Networks: Cache Hit Optimal vs. Throughput Optimal

Zheng Chen, Nikolaos Pappas and Marios Kountouris

*Abstract*—Departing from the conventional cache hit optimization in cache-enabled wireless networks, we consider an alternative optimization approach for the probabilistic caching placement in stochastic wireless D2D caching networks taking into account the reliability of D2D transmissions. Using tools from stochastic geometry, we provide a closed-form approximation of cache-aided throughput, which measures the density of successfully served requests by local device caches, and we obtain the optimal caching probabilities with numerical optimization. Compared to the cache-hit-optimal case, the optimal caching probabilities obtained by cache-aided throughput optimization show notable gain in terms of the density of successfully served user requests, particularly in dense user environments.

*Index Terms*—Wireless D2D caching, probabilistic content placement, stochastic geometry.

## I. INTRODUCTION

Proactive caching in wireless networks has attracted enormous attention as a means to reduce cellular data traffic load by bringing content closer to end users before being requested. Besides the conventional "cache the most popular content everywhere" strategy, in random wireless networks with spatially distributed network nodes, one widely used caching strategy is the probabilistic content placement, referred to as *geographic caching* [1] and *independent random caching* [2] in the literature. Recently, caching placement in random wireless networks with Poisson distributed caching helpers has been studied in different scenarios [3], [4]. In contrast to the conventional optimization targeting in cache hit probability maximization [5], in [6] the authors studied the optimization of the average success probability of content delivery in stochastic wireless caching helper networks in the noise-limited regime. The optimal content placement by maximizing the density of successful receptions was recently studied in [7]. It is worth noticing that in D2D assisted wireless networks, caching at user devices leads to a special case of *cache hit* when the requested content is cached within the device itself, which is often overlooked in the literature of wireless D2D caching networks. To the best of our knowledge, this special case was only addressed in [8], referred to as *self offloading*, where authors studied optimal probabilistic caching in a two-tier caching network by maximizing the offloading

Z. Chen is with the Laboratoire de Signaux et Systèmes (L2S, UMR8506), CentraleSupélec - CNRS - Université Paris-Sud, Gif-sur-Yvette, France. Email: zheng.chen@centralesupelec.fr.

N. Pappas is with Department of Science and Technology, Linköping University, Norrköping, Sweden. Email: nikolaos.pappas@liu.se.

M. Kountouris is with the Mathematical and Algorithmic Sciences Lab, France Research Center, Huawei Technologies Co., Ltd. Email: marios.kountouris@huawei.com.

probability, which is equivalent to the cache hit probability in a general sense.

In this letter we investigate probabilistic caching placement in a stochastic wireless D2D caching network with two different objectives: 1) to maximize the cache hit probability, which is the probability that a random user request can be served locally, either by its own cache or by its neighbor devices within a certain distance, and 2) to maximize the cache-aided throughput, which is the average density of successfully served requests by local caches, including when a user request is served by its own device or by a nearby device through D2D transmission. The optimal solutions and their performances are evaluated, which shows that the reliability of D2D transmission plays a critical role in the optimal caching decisions in order to achieve the best throughput-related performance.

## II. NETWORK MODEL

We consider a wireless D2D caching network where mobile user devices are modeled by a homogeneous Poisson Point Process (PPP) $\Phi_\mathrm{u}$ with intensity $\lambda_\mathrm{u}$. Each mobile user has probability $\rho \in [0,1]$ to be active, i.e., making an active request for a file, and the "inactive" devices will serve as potential D2D transmitters.[1] Therefore, the distributions of receivers and potential transmitters follow homogeneous PPPs $\Phi_\mathrm{u}^\mathrm{r}$ and $\Phi_\mathrm{u}^\mathrm{t}$ with intensity $\rho\lambda_\mathrm{u}$ and $(1-\rho)\lambda_\mathrm{u}$, respectively. Each device has cache memory size $M_\mathrm{d}$ and the files are assumed to have equal unit size. Assuming a finite content category $\mathcal{F} = \{f_1, \cdots, f_N\}$, where $f_i$ is the $i$-th most popular file and $N > M_\mathrm{d}$ is the library size. The content popularity follows the Zipf distribution, which is widely used in the literature [1]–[8], i.e., the request probability of file $f_i$ is

$$p_i = \frac{1}{i^\gamma \sum\limits_{j=1}^{N} j^{-\gamma}}, \tag{1}$$

where $\gamma$ is the shape parameter, which defines the skewness of the popularity distribution. We apply the *geographic caching* strategy proposed in [1] on user devices, i.e., each user device independently caches file $f_i$ with a certain probability $q_i$. Denote $\mathbf{q} = [q_1, \ldots, q_N]$ the caching probabilities of file $i \in [1, N]$, we have $\sum\limits_{i=1}^{N} q_i \leq M_\mathrm{d}$ due to the cache storage limit. As a result of independent thinning, the distribution of user devices who has the file $f_i$ follows a homogeneous PPP with intensity $q_i(1-\rho)\lambda_\mathrm{u}$.

When an active user requests for a file in $\mathcal{F}$, a *cache hit* event may happen in two cases:

---

[1]This setting is based on the assumption that a user cannot serve nearby devices when it is waiting to be served.

- **Case 1: self-request**, when the requested file is cached in its own device.
- **Case 2: D2D cache hit**, when the requested file is not cached in its own device, but in one of its nearby devices within a certain distance $R_d$. If there is more than one D2D transmitters which has the requested file, the file is transmitted from the nearest one.

In the case where the requested file can not be found in local device caches, the file is downloaded from the core network to the nearest base station through the backhaul, and then transmitted to the user. Here we assume D2D *overlaid* cellular networks, i.e., D2D and cellular communications thus, cross-tier interference does not exist.

## III. PERFORMANCE METRICS AND ANALYSIS

In this section we define the cache hit probability and the cache-aided throughput as the main performance metrics. Note that the case of finding the requested file of a device in its own cache storage is often overlooked in helper-based D2D caching networks in the literature. This is the major difference between the cache-related performance study in this letter and prior work in [4]–[6].

### A. Cache Hit Probability

The cache hit probability is the probability of a random active user to find its requested file in local caches, including the *self-request* and the *D2D cache hit* cases.

*1) Self-Request:* Denoting by $p_{\text{self}}^d$ the self-request probability of a random user is given by

$$p_{\text{self}} = \sum_{i=1}^{N} p_i q_i. \tag{2}$$

In this case, the request can be handled without the assistance of other devices.

*2) D2D Cache Hit:* As a result of the probabilistic caching, the probability to find a file cached inside a certain area strongly depends on the popularity order of the file and the area size. When a user requests for file $f_i$, the probability to find it cached in the devices within distance $R_d$ is [8], [9]

$$p_{\text{hit},i}^d = 1 - e^{-\pi(1-\rho)\lambda_u q_i R_d^2}. \tag{3}$$

Averaging over all the files in the content library $\mathcal{F}$, we have the D2D cache hit probability $p_{\text{hit}}^d = \sum\limits_{i=1}^{N} p_i (1 - q_i) p_{\text{hit},i}^d$, thus,

$$p_{\text{hit}}^d = \sum_{i=1}^{N} p_i (1 - q_i) \left( 1 - e^{-\pi(1-\rho)\lambda_u q_i R_d^2} \right). \tag{4}$$

The total cache hit probability is given by $p_{\text{hit}} = p_{\text{self}} + p_{\text{hit}}^d$, after replacing (2) and (4), we have

$$p_{\text{hit}} = 1 - \sum_{i=1}^{N} p_i (1 - q_i) e^{-\pi(1-\rho)\lambda_u q_i R_d^2}. \tag{5}$$

### B. Cache-Aided Throughput

Unlike the conventional definition of network throughput, which measures the average number of information successfully transmitted over the network region, here we are interested in studying the average number of requests that can be successfully and simultaneously handled by the local caches per unit area, namely the cache-aided throughput (per area).

Assume that the transmission of each file with equal size takes the same amount of time, one slot for instance. In the self-request case, the request is automatically served with probability one, while in the D2D cache hit case, the success probability of content delivery depends on the received signal-to-interference-plus-noise ratio (SINR). Thus, we have the cache-aided throughput given by

$$\mathcal{T} = \rho\lambda_u \left[ \sum_{i=1}^{N} p_i q_i \cdot 1 + \sum_{i=1}^{N} p_i (1 - q_i) p_{\text{hit},i}^d \cdot p_{\text{suc},i}^d \right], \tag{6}$$

where $p_{\text{suc},i}^d$ is the success probability of D2D transmission for file $f_i$, $\rho\lambda_u$ is the density of user requests in a given time slot. Without loss of generality, conditioning on having a typical active user $k$ to be served at the origin, the received SINR at the typical receiver is given by

$$\text{SINR}_k = \frac{P_d|h_{k,k}|^2 d_{k,k}^{-\alpha}}{\sigma^2 + \sum_{j \in \Phi_t^d \setminus \{k\}} P_d|h_{j,k}|^2 d_{j,k}^{-\alpha}},$$

where $\Phi_t^d$ denotes the set of active D2D transmitters; $P_d$ denotes the device transmission power; $h_{j,k}$ denotes the small-scale channel fading from the transmitter $j$ to the receiver $k$, which follows $\mathcal{CN}(0,1)$ (Rayleigh fading); $d_{j,k}$ denotes the distance from the transmitter $j$ to the receiver $k$; $\sigma^2$ denotes the background thermal noise power.

A file requested by a random user in $\Phi_r^u$ will be found within its nearby devices, but not in its own device with probability $p_{\text{hit}}^d$, as given in (4). Thus, the density of cache-assisted D2D transmissions is $\rho\lambda_u p_{\text{hit}}^d$.[2] Although the distribution of the active D2D transmitters $\Phi_t^d$ is not a homogeneous PPP, the average density of $\Phi_t^d$ can still be approximated by

$$\lambda_t^d \approx \rho\lambda_u p_{\text{hit}}^d. \tag{7}$$

When file $f_i$ is requested by the typical user, we denote $d_i$ the distance to the nearest device who has $f_i$ cached, and approximately consider $\Phi_t^d$ as a homogeneous PPP with intensity $\rho\lambda_u p_{\text{hit}}^d$. For a given SINR target $\theta$ of successful D2D transmission, the D2D success probability is given as

$$p_{\text{suc,i}}^d = \mathbb{P}\left[ \frac{P_d|h_{k,k}|^2 d_i^{-\alpha}}{\sigma^2 + \sum_{k \in \Phi_t^d \setminus \{k\}} P_d|h_{j,k}|^2 d_{j,k}^{-\alpha}} > \theta \right]$$

$$= \mathbb{P}\left[ |h_{k,k}|^2 > \frac{\theta d_i^\alpha}{P_d} \left( \sigma^2 + \sum_{k \in \Phi_t^d \setminus \{k\}} P_d|h_{j,k}|^2 d_{j,k}^{-\alpha} \right) \right]$$

$$\stackrel{(a)}{=} \mathbb{E}_{d_i} \left[ \mathcal{L}_{I_d}(\theta d_i^\alpha) \cdot \exp\left(-\theta\sigma^2 d_i^\alpha/P_d\right) \right]$$

$$\stackrel{(b)}{=} \mathbb{E}_{d_i} \left[ \exp\left(-\frac{\pi\rho\lambda_u p_{\text{hit}}^d d_i^2 \theta^{\frac{2}{\alpha}}}{\text{sinc}(2/\alpha)}\right) \exp\left(-\theta\sigma^2 d_i^\alpha/P_d\right) \right]$$

$$= \int_0^\infty f_{d_i}(r) \exp\left(-\frac{\pi\rho\lambda_u p_{\text{hit}}^d r^2 \theta^{\frac{2}{\alpha}}}{\text{sinc}(2/\alpha)}\right) e^{-\frac{\theta\sigma^2 r^\alpha}{P_d}} \, dr, \tag{8}$$

where $\mathcal{L}_{I_d}(s) = \mathbb{E}\left[\exp\left(-sI_d\right)\right]$ is the Laplace transform of interference $I_d = \sum_{k \in \Phi_t^d \setminus \{k\}} |h_{j,k}|^2 d_{j,k}^{-\alpha}$, and $f_{d_i}(r)$ is the probability density function (PDF) of D2D distance $d_i$, when $f_i$ is requested by the typical user. Here, (a) follows from the complementary cumulative distribution function (CCDF)

---

[2]Note that multiple users might find the same nearest D2D transmitter. In this case the transmitter will multicast the file to the receivers.

of $|h_{k,k}|^2$, which is exponentially distributed with unit mean value; (b) follows from the probability generating functional (PGFL) of PPP [10].

Conditioning on $d_i \leq R_d$ as a result of the maximum D2D distance, the PDF of $d_i$ is given by

$$f_{d_i}(r) = \begin{cases} \frac{2\pi(1-\rho)\lambda_u q_i r}{1-e^{-\pi(1-\rho)\lambda_u q_i R_d^2}} e^{-\pi(1-\rho)\lambda_u q_i r^2} & 0 \leq r \leq R_d \\ 0 & r > R_d. \end{cases}$$ (9)

Substituting (8) and (9) in (6), we obtain the cache-aided throughput averaged over all the files in the content library.

## IV. OPTIMIZATION OF PROBABILISTIC CACHING PLACEMENT

In this section we study the optimal caching probabilities $\mathbf{q} = [q_1, \ldots, q_N]$ by cache hit maximization and by cache-aided throughput optimization, respectively.

### A. Cache Hit Maximization

Based on (5), the optimization problem for maximizing the cache hit probability is defined as

$$\max_{\mathbf{q}} \quad p_{\text{hit}} = 1 - \sum_{i=1}^{N} p_i (1-q_i) e^{-\pi(1-\rho)\lambda_u q_i R_d^2} \quad (10)$$

$$\text{s.t.} \quad 0 \leq q_i \leq 1 \text{ for } i = 1, \ldots, N$$

$$\sum_{i=1}^{N} q_i \leq M_d.$$

The second order derivative of the objective function is strictly negative, thus $p_{\text{hit}}$ is a concave function of $q_i$ for $i = 1, \ldots, N$. Consider the following Lagrangian function

$$\mathcal{L}(\mathbf{q}, \mu) = -1 + \sum_{i=1}^{N} p_i (1-q_i) e^{-\pi(1-\rho)\lambda_u q_i R_d^2} \\ + \mu \left( \sum_{i=1}^{N} q_i - M_d \right),$$ (11)

where $\mu$ is the non-negative Lagrangian multiplier. We solve this optimization problem by applying the Karush-Kuhn-Tucker (KKT) conditions. From $\frac{\partial \mathcal{L}}{\partial q_i} = 0$, we have

$$q_i(\mu) = -\frac{\mathcal{W}\left\{\frac{\mu}{p_i} \exp\left[1 + \pi(1-\rho)\lambda_u R_d^2\right]\right\}}{\pi(1-\rho)\lambda_u R_d^2} + \frac{1}{\pi(1-\rho)\lambda_u R_d^2} + 1,$$ (12)

where $\mathcal{W}$ denotes the Lambert W function [11]. Combined with the condition $0 \leq q_i \leq 1$, let $[x]^+ = \max\{x, 0\}$, we have

$$q_i^\star = \min\left\{[q_i(\mu^\star)]^+, 1\right\},$$ (13)

where $\mu^\star$ can be obtained by the bisection search method under the other KKT condition $\sum_{i=1}^{N} q_i^\star = M_d$.

### B. Cache-aided Throughput Maximization

Due to the complicated expression of $\mathcal{T}$, the optimal caching probabilities that maximize the cache-aided throughput are difficult to obtain, even with numerical methods. Consider the following approximation

$$\mathbb{E}_{d_i}[\exp(-\eta d_i^\delta)] \approx \exp\left(-\eta \mathbb{E}[d_i^2]^{\delta/2}\right),$$ (14)

the success probability $p_{\text{suc},i}^d$ in (8) can be approximated by

$$\hat{p}_{\text{suc},i}^d \approx \exp\left[-\frac{\pi\rho\lambda_u p_{\text{hit}}^d \mathbb{E}[d_i^2]\theta^{2/\alpha}}{\text{sinc}(2/\alpha)}\right] \exp\left[-\frac{\theta\sigma^2 \mathbb{E}[d_i^2]^{\alpha/2}}{P_d}\right].$$ (15)

From the PDF of $d_i$ in (9), we can obtain $\mathbb{E}[d_i^2]$ as follows.

$$\mathbb{E}[d_i^2] = \int_0^{R_d} r^2 \frac{2\pi(1-\rho)\lambda_u q_i r}{1-e^{-\pi(1-\rho)\lambda_u q_i R_d^2}} e^{-\pi(1-\rho)\lambda_u q_i r^2} dr$$

$$= \frac{1}{\pi(1-\rho)\lambda_u q_i} - \frac{R_d^2}{e^{\pi(1-\rho)\lambda_u q_i R_d^2} - 1}.$$ (16)

When $q_i \to 0$, we obtain $\lim_{q_i \to 0} \mathbb{E}[d_i^2] = R_d^2/2$ by applying L'Hôpital's rule.

Then we have the approximated cache-aided throughput as

$$\hat{\mathcal{T}} = \rho\lambda_u \left[\sum_{i=1}^{N} p_i q_i + \sum_{i=1}^{N} p_i(1-q_i)p_{\text{hit},i}^d \cdot \hat{p}_{\text{suc},i}^d\right],$$ (17)

where $\hat{p}_{\text{suc},i}^d$ is given in (15). Our objective is to find $\mathbf{q}^\star = \max_{\mathbf{q}} \hat{\mathcal{T}}$, subject to $0 \leq q_i \leq 1$ and $\sum_{i=1}^{N} q_i \leq M_d$.

This problem is non-convex as it can be seen numerically. Providing an analytical solution to this problem is difficult. Therefore for the the cache-aided throughput maximization we solve it numerically with Simulated Annealing.

## V. NUMERICAL AND SIMULATION RESULTS

For numerical evaluation, we consider the user density between $\lambda_u = [10^{-4}, 10^{-3}]/\text{m}^2$. $\rho = 50\%$ of the users will request for a random file in $\mathcal{F}$ according to the request probabilities $\mathbf{p} = [p_1, \ldots, p_N]$, which follows the Zipf distribution with parameter $\gamma = \{0.5, 1.2\}$. The rest $50\%$ of users act as potential D2D transmitters helping to serve the user requests locally. The device cache capacity is $M_d = 2$ files. The content library has size $N = 20$ files.[3] The D2D searching distance is $R_d = 75$ m. The device transmission power and the background noise power are $P_d = 0.1$ mW and $\sigma^2 = -110$ dB, respectively. The target SINR of successful D2D transmissions is $\theta = 0$ dB.

In Fig. 1 and Fig. 2 we compare the cache-hit-optimal (Section IV-A) and throughput-optimal (Section IV-B) caching probabilities $\mathbf{q}^\star$ in sparse and dense user environments, respectively. The optimal caching probabilities of file $f_i$ for $i = 1, \ldots, N$ are plotted as a function of the popularity order $i$. Interestingly, we observe that with sparse users, throughput-optimal and cache-hit-optimal caching probabilities are very close, while with dense users, each device tends to cache the most popular files with higher probability in order to increase the cache-aided throughput. For instance, in Fig. 2, $q_1^\star$, $q_2^\star$ and $q_3^\star$ in the throughput-optimal case are much higher than in the cache-hit-optimal case. The intuition behind this is that *in dense user regime, due to the excessive D2D interference that leads to very low D2D success probability, users' caching strategy tends to be more "selfish" in the sense that self-request matters more than cache-assisted D2D transmission.*

---

[3]Note that in reality the content library size is very large. Here we take $N = 20$ files to avoid high complexity of the optimization problem. Similar choices can also be found in [6] and [8].
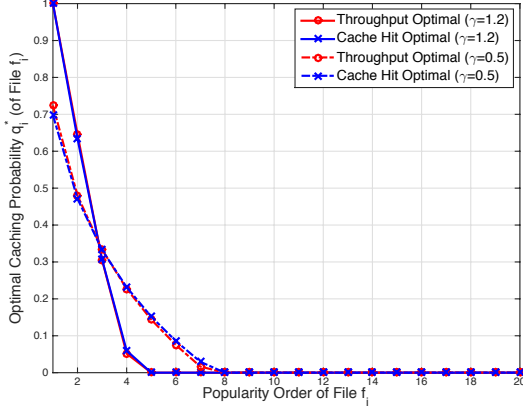
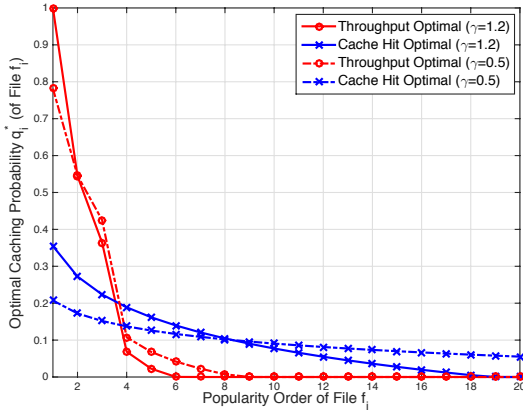Fig. 1. Optimal caching probabilities with sparse devices, $\lambda_u = 10^{-4}/m^2$.



Fig. 2. Optimal caching probabilities with dense devices, $\lambda_u = 10^{-3}/m^2$.



Fig. 3. Simulated cache-aided throughput vs. user device density $\lambda_u$. $\gamma = \{0.5, 1.2\}$

maximizing cache hit probability and maximizing the cache-aided throughput that is defined by the density of successfully served user requests. The main takeaway message is that, in additional to the cache hit probability, the success probability of content delivery is also a critical factor that needs to be taken into account in the optimal caching placement in order to improve the throughput performance in wireless D2D caching networks.

In Fig. 3 we plot the simulated cache-aided throughput obtained with the throughput-optimal caching probabilities. In order to validate the accuracy of the approximation we used in (15), we plot the theoretical values of the approximated cache-aided throughput $\hat{\mathcal{T}}$, which turn out to have negligible error. For the comparison of different caching strategies, we also plot the simulated cache-aided throughput when applying $\mathbf{q}^\star$ obtained with cache hit probability optimization and with the conventional "cache the most popular content" (MPC) strategy. It is obvious that with the throughput-optimal strategy that is aware of the D2D success probability, the achieved cache-aided throughput can be significantly improved compared to the cache hit optimization and the MPC strategies. The gain is more profound in the dense user regime. Another interesting remark is that with dense users and highly concentrated content popularity ($\gamma = 1.2$), the cache-aided throughput with MPC gives better performance than the cache-hit-optimal case, meaning that it is more beneficial to increase the chance of "self-request" than increasing the total cache hit ratio. As summary, our results validate the necessity of taking into account the transmission reliability of cache-aided D2D communication while searching for the optimal content placement.

## VI. Conclusions

In this letter, we studied probabilistic caching placement in stochastic wireless D2D caching networks with two objectives:
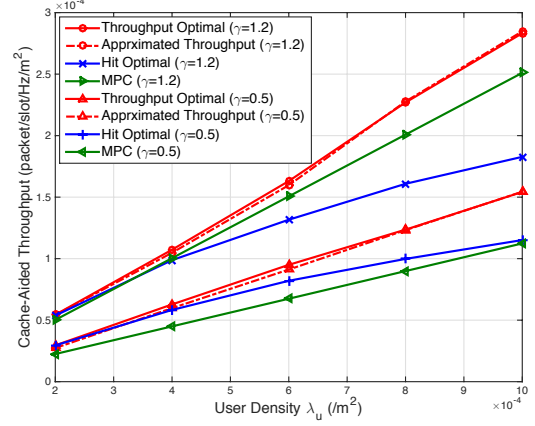
## References

[1] B. Blaszczyszyn and A. Giovanidis, "Optimal geographic caching in cellular networks," in *Proc., IEEE Intl. Conf. on Communications (ICC)*, London, UK, June 2015, pp. 3358–3363.

[2] M. Ji, G. Caire, and A. Molisch, "Wireless device-to-device caching networks: basic principles and system performance," *IEEE Journal on Sel. Areas in Commun.*, vol. 34, no. 1, pp. 176–189, Jan. 2016.

[3] S. H. Chae, J. Y. Ryu, T. Q. S. Quek, and W. Choi, "Cooperative transmission via caching helpers," in *Proc., IEEE Globecom*, San Diego, CA, Dec. 2015, pp. 1–6.

[4] J. Song, H. Song, and W. Choi, "Optimal caching placement of caching system with helpers," in *Proc., IEEE Intl. Conf. on Communications (ICC)*, London, UK, June 2015, pp. 1825–1830.

[5] H. Kang, K. Park, K. Cho, and C. Kang, "Mobile caching policies for device-to-device (D2D) content delivery networking," in *Proc., IEEE Conf. on Computer Commun. Workshops (INFOCOM WKSHPS)*, Toronto, Canada, Apr. 2014, pp. 299–304.

[6] S. H. Chae and W. Choi, "Caching placement in stochastic wireless caching helper networks: Channel selection diversity via caching," *IEEE Trans. on Wireless Commun.*, vol. 15, no. 10, pp. 6626–6637, Oct. 2016.

[7] D. Malak, M. Al-Shalash, and J. G. Andrews, "Optimizing content caching to maximize the density of successful receptions in device-to-device networking," *IEEE Trans. on Communications*, vol. 64, no. 10, pp. 4365–4380, Oct. 2016.

[8] J. Rao, H. Feng, C. Yang, Z. Chen, and B. Xia, "Optimal caching placement for D2D assisted wireless caching networks," in *Proc., IEEE Intl. Conf. on Communications (ICC)*, Kuala Lumpur, Malaysia, May 2016, pp. 1–6.

[9] Z. Chen and M. Kountouris, "D2D caching vs. small cell caching: Where to cache content in a wireless network?" in *IEEE Intl. Workshop on Signal Processing Advances in Wireless Commun.(SPAWC)*, Edinburgh, UK, Jul. 2016, pp. 1–6.

[10] M. Haenggi and R. K. Ganti, "Interference in large wireless networks," *Found. Trends Netw.*, vol. 3, no. 2, pp. 127–248, Feb. 2009.

[11] R. M. Corless, G. H. Gonnet, D. E. G. Hare, D. J. Jeffrey, and D. E. Knuth, "On the Lambert W function," *Advances in Computational Mathematics*, vol. 5, no. 1, pp. 329–359, 1996.