

# Device-to-Device Multicast Content Delivery in cellular networks

Yanli Xu

Department of Information Engineering  
Shanghai Maritime University  
1550 Haigang Avenue  
Shanghai, China  
ylxu@shmtu.edu.cn

Ping Wu

Department of Engineering Sciences  
Uppsala University  
Box 256, 75105  
Uppsala, Sweden  
ping.wu@angstrom.uu.se.

## ABSTRACT

With the explosive increase of video traffic and content sharing applications among mobile devices such as smart phones and mobile tablets, device-to-device (D2D) technique has become attractive in the application to content delivery since it can effectively improve the resource efficiency of cellular network through proximal transmission. In this paper, we investigate D2D multicast content delivery in cellular networks and propose a multicast architecture for D2D content delivery. Based on the architecture, we analyze the performance of content delivery through D2D communications which provides a tool for evaluating the performance of D2D delivery. With these analyses, we can obtain the number of content requests in the network which is served by D2D content delivery and the average serving time through D2D multicast content delivery. Simulation results under different transmission parameters validate the proposed theory. And it demonstrates that the proposed architecture and methods are useful for the application of D2D communication technique to content delivery to offload cellular traffic and improve the resource efficiency of cellular networks.

## CCS Concepts

•**Networks** → *Network management*;

## Keywords

D2D communication; content delivery; multicast; distributed caching.

## 1. INTRODUCTION

With the explosive growth of mobile data, offloading traffic from central nodes, e.g. base stations (BSs), becomes a very interesting topic in future cellular networks. An effective way of achieving this is to enable data to be transmitted directly among user equipments (UEs). This has become possible because of the emergence of device-to-device (D2D)

communication. The D2D communication makes data to be transmitted directly among UEs under the control of BSs [7, 3], and is an effective technology for offloading the increasing demands on local area services. Therefore, it is suited very well for context-aware applications in which content sharing occurs very often. Social network is the second largest traffic contributor in mobile networks and content sharing accounts for 15 percent [1]. The D2D technique can carry out proximal content delivery so that it can largely reduce resource consumption used for the content retrieval from BS [15, 19, 8].

For D2D content delivery, UE may cache contents with a proactive mindset in case that proximal UE or itself needs them. A common proximal content delivery model is based on a cluster. For example, when a UE (UE 1) in the cluster requests some content, another UE (UE 2) that caches the content is selected to send the requested content to this UE 1 [2]. Thus, caching strategy of UE affects the delivery performance which needs to be well designed. It is obvious that more useful content should be cached with larger probability and some works investigate the caching strategy to improve the delivery performance [13, 12, 14]. Besides, interference mitigation needs to be considered for the efficient delivery of contents and transmission power considered for the delivery range. By properly adjusting the size of a cluster and the number of simultaneous requests from UEs [2, 9], the interference can be reduced and successful delivery probability can be improved.

As above mentioned, the introduction of D2D content delivery to cellular networks may provide performance gains. However, the realization D2D content delivery may face many challenges, and up to now, not much related work has been seen to be done yet. In the most existing literature such as [3, 9], the architecture of D2D content delivery was based on D2D unicast with request-and-response pattern. The authors in [4] suggested multicast when discussing role of proactive caching in 5G networks. In this paper, we propose an architecture for D2D content delivery via multicast in cellular networks. With this architecture, the key performance metrics having impact on delivery performance are studied, e.g. the number of caching UEs, the number of requesting UEs and the average serving rate. The theoretical analysis results provide comprehensive understanding of D2D multicast content delivery in cellular networks, and can also be useful tools for determining when D2D multicast should be used for content delivery. The rest parts of the paper are organized as follows. In Section 2 the system model

ACM ISBN .

DOI:

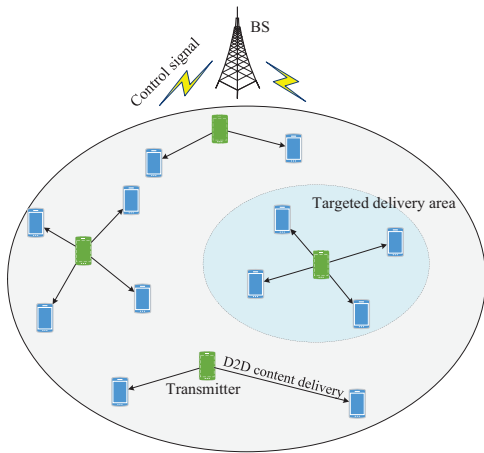


Figure 1: An example for system model

is presented and the D2D delivery problem is formulated. In Section 3, the D2D multicast content delivery in cellular networks is studied. For a random content, the variation of UEs requiring and obtaining it from D2D communication is analyzed, and related results are provided. In addition, the average serving time for a random request of content is studied. Simulation results and related discussions are presented in Section 4. Finally, Section 5 concludes this paper.

## 2. SYSTEM MODEL AND PROBLEM FORMULATION

In this section, the system model is described and the problem on content delivery using D2D communication is formulated. We consider an infinite cellular network where UEs are randomly located. Each BS in the network serves a number of UEs. Without loss of generality, we select a cell in the cellular network as the interest area having the BS located at its center and a radius of  $R_0$ , and study its content delivery performance. Besides, the influences from other cells such as interferences are also considered in the analysis. The service scenario considered is illustrated in Fig.1.

### 2.1 Network Model and Channel Model

UEs in the network are randomly located that may have two alternative transmission modes: (1) cellular mode in which transmission is through BS and (2) D2D mode in which proximal UEs transmit data directly to each other without the relay of BS. Communication links under different modes use orthogonal resources with the consideration of interference mitigation. For a number of D2D links which simultaneously transmit, the same resource is allocated for these links to save resource. Interferences among D2D links can be controlled due to the proximal character of D2D communication. In this situation, Poisson point process (PPP) model is suited and used here for characterizing the spatial distribution of the D2D transmitters that are scheduled simultaneously. Particularly, the distribution of the D2D transmitters follows the PPP distribution with density  $\lambda_s$ . It should be pointed out that the PPP model is commonly used in the characterization of nodes' locations in wireless environment and widely used for the analysis of D2D communication [6, 18, 17].

Considering a communication link  $i \rightarrow j$ , the transmission over the link is regarded to be successful when the signal-to-interference ratio (SIR) at the receiver is larger than a decoding threshold  $\gamma_{th}$ . The thermal noise is assumed to be so small as to be negligible. This assumption may be relaxed (e.g., see [16, 11]) but at the cost of complicating the derived expressions without providing additional insight. For the power control of a transmitter, we assume that the transmitter chooses such transmission power that the signal power at the intended receiver is a pre-specified constant  $E_0$ . In this case, if the communication link distance is  $d$ , then the transmission power is  $E_0 d^\alpha$ , where  $\alpha$  is the path loss exponent. For the link  $i \rightarrow j$ , the SIR at the receiver  $j$  can be expressed by

$$\begin{aligned} \gamma_j &= \frac{E_i d_{ij}^{-\alpha} H_{ij}}{I_j} \\ &= \frac{E_0 H_{ij}}{\sum_{k \in \Omega} E_k d_{kj}^{-\alpha} H_{kj}} \end{aligned} \quad (1)$$

where  $E_i$  denotes the transmission power of node  $i$ ,  $d_{ij}$  is the distance from the transmitter  $i$  to the receiver  $j$ ,  $I_j$  is the interference at the receiver  $j$ ,  $\Omega$  is the interfering transmitter set and  $H_{ij}$  characterizes the fast fading power from  $i$  to  $j$ .

### 2.2 Content Delivery Model

Each UE has a storage capacity called cache storing a number of contents. There are multiple contents coexisting in the network. Without loss of generality, we focus on the analysis of the caching and delivery of a random content. During the lifetime of this content, request for this content is modeled by Poisson arrival process with arriving rate  $\lambda$ , which can easily be generalized to a more practical case based on queuing theory. The mobility of a UE is also modeled by Poisson process, i.e., UE arriving at the network with rate  $\theta$  (called UE arriving rate) and leaving the network with rate  $\mu$  (called UE leaving rate).

Multicast content delivery is used here to serve multiple users which have common interest and require the same content simultaneously. Considering that the number of users requesting the same content at  $t$  may be very small due to the proximal transmission character of D2D communication, then for D2D multicast content delivery, we propose a self-sufficiency cache strategy. In this strategy, transmitting UEs are scheduled to deliver contents during a D2D multicast period as shown in Fig. 1. Proximal UEs are scheduled to receive these multicast contents and cache them proactively in case that itself or proximal UEs need them later. Thereby, many users can receive the contents from one multicast.

For a set of cached contents, a UE transmits them according to a pre-defined distribution. For example, the UE transmits content  $n$  with probability  $B(n)$  over a certain area of interest called targeted delivery area which is a circle with a radius of  $R$  (Fig. 1). UEs located in the targeted delivery areas of the transmitter are the potential receivers which may be scheduled to receive multicast contents. A UE may be located in the targeted delivery areas of several transmitters simultaneously and it only receives the multicast content from a transmitter once a time.

## 3. D2D MULTICAST CONTENT DELIVERY

In this section we propose an architecture for D2D multicast content delivery and analyze the delivery performance based on this architecture. With these analyses, strategy for when should use D2D multicast-based content delivery is proposed.

Contents delivered by a multicast belong to a content set in the network and have a certain lifetime. During the lifetime of the contents belonging to the content set  $F$ , requests for contents from UEs arrive at the system successively. Without loss of generality, we focus on the performance of content delivery for a content  $n \in F$ . For D2D multicast content delivery, the delivery performance depends on the number of UEs which cache the content. When more UEs cache the content, the probability that proximal UEs with content needed by others is larger. In this case, the successful delivery probability is improved with less resource consumption. The more UEs receive the content, the more UEs cache this content in our delivery architecture, and then the more requests will be served. Hence, the delivery performance, i.e., the number of successfully served requests, depends on numbers of UEs that cache and request the content. To evaluate the delivery performance, we investigate the number of UEs caching this content (called caching UE), denoted by  $C(t)$ , and the number of UEs requesting this content (called requesting UE), denoted by  $R(t)$ , at time  $t$ , respectively.

To study  $C(t)$  and  $R(t)$ , the key factor is the derivation of successfully served request rate at  $t$ . When a UE requires content  $n$  at  $t$ , the request is regarded to be served if it successfully receives the content currently or in previous multicast period, that is to say, the request is served if the content is proactively cached by this UE, which is different from the unicast case. Thus, we first study the cache status for a random content at a UE. For a UE  $j$ , the conditional probability that it successfully receives the content  $n$  from its transmitter  $i$  can be calculated by

$$\begin{aligned}\varepsilon_i &= \Pr \{ \gamma_{ij} \geq \gamma_{\text{th}} \} B(n) \\ &= \exp(-\kappa \lambda_s d_{ij}^2) B(n)\end{aligned}\quad (2)$$

$\kappa = \pi m \Gamma(m) \Gamma(1-m)$ ,  $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$  [10] and  $m = 2/\alpha$ .

Since a UE  $j$  may be simultaneously located in the targeted delivery areas of several transmitters which construct its potential transmitter set  $S$ , and it randomly selects one of them, the probability that it successfully receives the content  $n$  from its potential transmitters can be written as

$$\begin{aligned}\varepsilon &= \sum_{j \in S} \Pr \{ j \} \varepsilon_j \\ &= \frac{\int_0^{2\pi} \int_0^R \lambda_s \exp(-\kappa \lambda_s r^2) r dr}{\pi \lambda_s R^2} \\ &= \frac{1 - \exp(-\kappa \lambda_s R^2)}{\kappa \pi \lambda_s R^2}\end{aligned}\quad (3)$$

Assuming that the number of potential receiving UEs in the network is  $N$ , the number of caching UEs without considering UE leaving rate can be expressed in the following manner,

$$D_m(t) = \varepsilon N B(n) \quad (4)$$

which indicates that  $\varepsilon$  affects the number of caching UEs so that it affects the delivery performance. Thus, we define  $\varepsilon$

as the delivery factor and will discuss the effect of  $\varepsilon$  on delivery performance in the later section. Assuming that initial number of UEs in the network as  $N_0$ ,  $N$  can be calculated by

$$N = N_0 + (\theta - \mu)t \quad (5)$$

Inserting (5) into (4),  $D_m(t)$  can be expressed as

$$D_m(t) = \varepsilon [N_0 + (\theta - \mu)t] B(n) \quad (6)$$

When the UE leaving rate  $\mu$  is taken into consideration, the change rate of the number of caching UEs can be written as

$$\begin{aligned}\frac{dC(t)}{dt} &= D_m(t) - \mu C(t) \\ &= \varepsilon [N_0 + (\theta - \mu)t] B(n) - \mu C(t)\end{aligned}\quad (7)$$

Solving (7) for  $C(t)$ , we have

$$C(t) = k e^{-\mu t} + \frac{\varepsilon(\theta - \mu) B(n) t + N_0 \varepsilon B(n)}{\mu} - \frac{\varepsilon(\theta - \mu) B(n)}{\mu^2} \quad (8)$$

Using the initial value  $C(0) = 0$ , we can determine the value of  $k$  as

$$k = \frac{\varepsilon(\theta - \mu) B(n)}{\mu^2} - \frac{N_0 \varepsilon B(n)}{\mu} \quad (9)$$

Substituting (9) into (8),  $C(t)$  with the assumption of an infinite number of UEs in the network can be written as follows.

$$\begin{aligned}C^0(t) &= \frac{\varepsilon(\theta - \mu) B(n)}{\mu} t \\ &+ \frac{N_0 \varepsilon B(n)}{\mu} (1 - e^{-\mu t}) + \frac{\varepsilon(\theta - \mu) B(n)}{\mu^2} (e^{-\mu t} - 1)\end{aligned}\quad (10)$$

From (10), we can find that, when  $\mu$  is small so that the  $2^{\text{nd}}$  and  $3^{\text{rd}}$  terms on the right hand side can be neglected,  $C^0(t)$  increases with  $t$ . As a result, the number of UEs having the content  $n$  in the network will grow with time, that is to say, the number of caching UEs approaches to the total number of UEs. Thus, we consider the constraint of the total number in the network, and express  $C(t)$  in the following manner.

$$C(t) = \min \{ N_0 + \theta t - \mu t, C^0(t) \} \quad (11)$$

We are now studying the number of requesting UEs,  $R(t)$ , which is determined by request rate  $\lambda$  and the served rate  $S_m(t)$  of requests.  $S_m(t)$  depends on the number of requests and the successful transmission probability, i.e., the number of UEs which require the content and successfully receive it. Thus,  $S_m(t)$  can be expressed by

$$\begin{aligned}S_m(t) &= D_m(t) q \\ &= \varepsilon B(n) N q\end{aligned}\quad (12)$$

where  $q$  is the request probability for each UE. Thus, the number of requesting UEs at an instantaneous time follows Binomial distribution. According to probability theory, Binomial distribution can be approximated as Poisson distribution with mean value  $\lambda \approx Nq$ . Thereby,  $S_m(t)$  can be written as

$$S_m(t) = \varepsilon B(n) \lambda \quad (13)$$

From (13), we can see that  $S_m(t)$  is independent to time  $t$ , that is to say, the number of served requests is uniformly

distributed for any arbitrary time period like request arrival process. Since  $S_m(t)$  requests are served at  $t$ , the practical request arrival rate can be seen as the difference between  $\lambda$  and  $S_m(t)$ , which can be expressed as

$$\begin{aligned}\lambda' &= \lambda - S_m(t) \\ &= \lambda \left[ 1 - \frac{1 - \exp(-\kappa\lambda_s R^2)}{\kappa\pi\lambda_s R^2} B(n) \right]\end{aligned}\quad (14)$$

The change rate of the number of requesting UEs,  $R(t)$ , can be expressed as

$$\frac{dR(t)}{dt} = \lambda' - \mu R(t) \quad (15)$$

which is the 1<sup>st</sup> order linear non-homogeneous differential equation and its general solution can be easily obtained using calculus as follows.

$$R(t) = ke^{-\mu t} + \frac{\lambda'}{\mu} \quad (16)$$

Using initial value  $R(0) = 0$ , we have

$$k = -\frac{\lambda'}{\mu} \quad (17)$$

Inserting (17) into (16),  $R(t)$  can be written as

$$\begin{aligned}R(t) &= \frac{\lambda'}{\mu} (1 - e^{-\mu t}) \\ &= \frac{\lambda \left[ 1 - \frac{1 - \exp(-\kappa\lambda_s R^2)}{\kappa\pi\lambda_s R^2} B(n) \right]}{\mu} (1 - e^{-\mu t})\end{aligned}\quad (18)$$

Based on the above analysis, we can obtain the number of served requests during any time period through D2D multicast content delivery. For a time period of  $[0, T]$ , the total number of served requests can be determined as follows.

$$\begin{aligned}S_T &= \int_0^T S_m(t) dt \\ &= \int_0^T \varepsilon B(n) \lambda dt \\ &= T\varepsilon B(n) \lambda\end{aligned}\quad (19)$$

For the steady state of the system, we may let

$$\frac{dC(t)}{dt} = \frac{dR(t)}{dt} = 0 \quad (20)$$

Assuming that the equilibrium values of  $C(t)$  and  $R(t)$  are  $\bar{C}$  and  $\bar{R}$ , then from (7) and (15) we can obtain

$$\bar{C} = \frac{\varepsilon N_0 B(n)}{\mu} \quad (21a)$$

$$\bar{R} = \frac{\lambda'}{\mu} \quad (21b)$$

According to Little's law [5] which indicates that the average number of requests in a stable system is equal to the multiplication of effective arrival rate and serving time per request, the serving time  $T_s$  per content request under D2D multicast can be expressed as

$$\frac{\lambda - \bar{R}}{\lambda} \bar{R} = (\lambda - \bar{R}) T_s \quad (22)$$

**Table 1: An example for mode selection**

1. Calculate $T_s$ using (23)
2. Compare $T_s$ with $T_0$
3. <b>If</b> $T_s < T_0$
4. Choose D2D multicast for content delivery
5. <b>else</b>
6. Choose BS for content delivery
7. <b>end</b>

**Table 2: Simulation parameters**

Simulation parameters	Value
Transmitter density $\lambda_s$	$10^{-7} \sim 10^{-5}$
Path loss factor $\alpha$	3
SINR decoding threshold $\gamma_{th}$	10dB
UE arriving rate $\theta$	$10^{-2}/s$
UE leaving rate $\mu$	$10^{-2}/s$
Request rate $\lambda$	$1/s$
Multicast probability $B(n)$	0.01

Submitting (21b) into (22),  $T_s$  can be written as

$$\begin{aligned}T_s &= \frac{\bar{R}}{\lambda} \\ &= \frac{\lambda'}{\mu\lambda} \\ &= \frac{1 - \frac{1 - \exp(-\kappa\lambda_s R^2)}{\kappa\pi\lambda_s R^2} B(n)}{\mu}\end{aligned}\quad (23)$$

From (23) it follows that the serving time depends on several factors such as UE mobility, transmitter density, targeted delivery area and channel quality. The shorter the serving time the more efficient the D2D delivery. With (23), we can evaluate the serving time under different conditions and determine whether D2D multicast mode should be used for content delivery or not. For example, for a service with QoE (Quality of Experience)  $T_0$  which is the upper constraint for the serving time, the selection of transmission mode can be implemented using the following procedure in Table 1. When the transmission mode is selected, transmission scheduling can be determined accordingly.

#### 4. PERFORMANCE EVALUATION ON D2D MULTICAST CONTENT DELIVERY

To evaluate the theoretical study on D2D communication- s for content delivery in Section 3, simulation is performed with the parameters given in Table 2. Firstly, we investigate how delivery factor  $\varepsilon$  changes against transmitter density  $\lambda_s$  for different settings of targeted delivery area radius  $R$ , path loss factor  $\alpha$ , and decoding threshold  $\gamma_{th}$  to take different values. The results are shown in Fig.2. Here we let  $B(n) = 1$  in order to focus on the influence of transmission parameters on  $\varepsilon$ . From the analysis we know that  $\varepsilon$  has impacts on the number of caching UEs and a larger  $\varepsilon$  leads to a larger number of UEs caching the transmitted content. Therefore, the content request has a larger probability to be satisfied. Hence,  $\varepsilon$  is expected to be improved. From Fig.2, we can observe that  $\varepsilon$  increases with  $\lambda_s$  while this increase becomes slow at higher transmitter density. That's because higher  $\lambda_s$  brings about severer interference and thus reduces

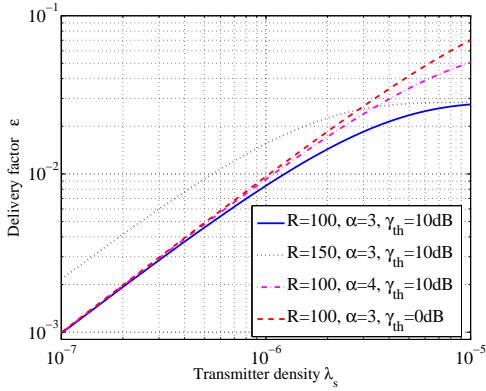


Figure 2: Comparison of delivery factor  $\varepsilon$  against the transmitter density  $\lambda_s$

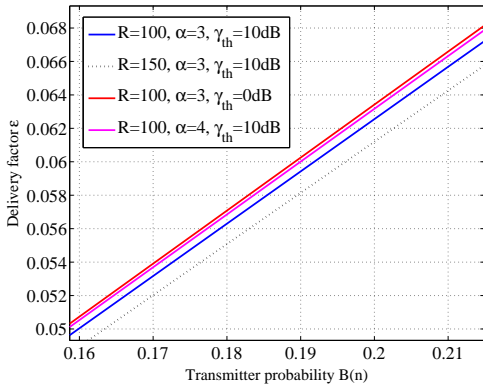


Figure 3: Comparison of delivery factor  $\varepsilon$  against the transmitter probability  $B(n)$

successful delivery probability although higher  $\lambda_s$  makes UE have a larger probability to select a transmitter with a better communication channel. In addition, it is shown that larger  $R$  reduces performance. The reason for this can be that larger  $R$  makes UE have a smaller probability to select a transmitter with better communication channel for a certain transmitter density due to the increase of average transmission distance. Finally, by comparing the dash-dotted curve with the dashed curve, it can be concluded that  $\varepsilon$  is better for lower  $\alpha$  and lower  $\gamma_{th}$  when  $R$  is fixed.

To further evaluate the transmission strategy (self-sufficiency cache strategy), i.e., the impact of factor  $B(n)$  on  $\varepsilon$ ,  $\varepsilon$  is compared against  $B(n)$  for different values of  $\alpha$  and  $\gamma_{th}$  (Fig. 3). From this figure, we can observe that  $\varepsilon$  increases with the increase of  $B(n)$ . Thus, for a popular content or the content which is more commonly interesting for proximal UEs, it is better to enlarge  $B(n)$  to improve the D2D delivery efficiency. This may be an interesting research issue for optimal distributed caching and delivery for D2D communication.

To see how the number of caching UEs  $C(t)$  is affected by  $\mu$ ,  $\theta$  and  $\varepsilon$ ,  $C(t)$  is evaluated when  $\mu$ ,  $\theta$  and  $\varepsilon$  take different values. The results (Fig.4) show that  $C(t)$  is dependent on the mobility of UEs and the content delivery performance. In this figure,  $\varepsilon$  characterizes the effects of transmission environment parameters (such as channel quality, interferences from simultaneous transmitters and used delivery strategy)

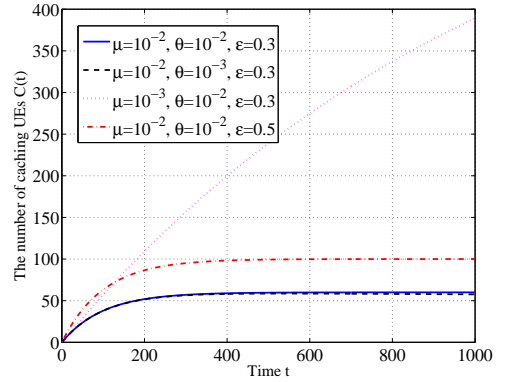


Figure 4: Comparison of the number of caching UEs  $C(t)$  against time  $t$

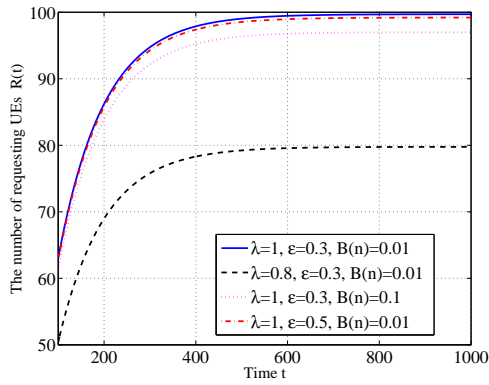
on delivery performance. Comparing the solid curve and the dash-dotted curve, we can observe that  $C(t)$  increases when  $\varepsilon$  increases. Thus, it is effective to improve delivery efficiency through increasing  $\varepsilon$ . Although  $C(t)$  increases with  $\theta$ , it decreases with  $\mu$ . A possible explanation of this is that more UEs in the network make more UEs obtain the network with D2D multicast for the same delivery efficiency.

To understand the influence of  $\lambda$ ,  $\varepsilon$  and  $B(n)$  on the number of requesting UEs,  $R(t)$ , it is evaluated when  $\lambda$ ,  $\varepsilon$  and  $B(n)$  take different values. The results is shown in Fig. 5. Comparing the solid and dotted curves leads to the conclusion that increasing multicast probability  $B(n)$  for a given content leads to increasing the number of served content and reducing the number of request for D2D multicast content delivery. This indicates that we should allocate more resources to more popular contents for D2D multicast content delivery. Comparing the solid and dash-dotted curves, we can observe that  $R(t)$  decreases slightly with the increase of  $\varepsilon$  due to more requests served. Obviously,  $R(t)$  decreases with  $\lambda$  by comparing the solid and dashed curves, which reveals that the requests from UEs need to be scheduled adaptively according to the serving time.

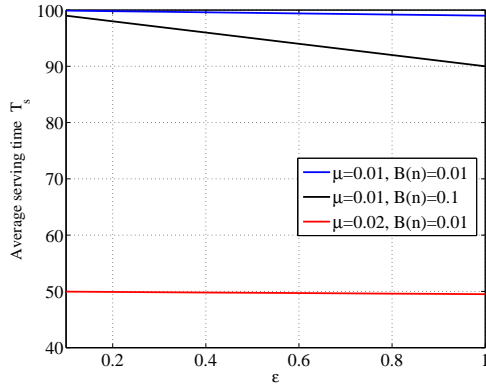
The average serving time  $T_s$  for a request is shown in Fig. 6. It is an important metric for the evaluation of the delivery efficiency of D2D multicast. As shown in this figure,  $T_s$  decreases with the increase of  $\varepsilon$ . It validates that increasing  $\varepsilon$  can improve delivery efficiency. And the methods of increasing  $\varepsilon$  have been discussed in the earlier sections.

## 5. CONCLUSION

In this paper we have investigated device-to-device (D2D) multicast content delivery in cellular networks. Particularly, we have proposed multicast content delivery architecture which takes advantages of proximal transmission and wireless multicast, and then studied the some key metrics, such as the number of caching UEs and the number of requesting UEs, which affect the delivery efficiency of D2D communication. Through this study, the delivery efficiency of D2D multicast has been evaluated, and a criterion has been established to decide when and where to use the D2D multicast for content delivery, which is one of the primary problems that need to be handled when applying D2D communication to content delivery in cellular networks. Furthermore, we have provided the number of served requests and the av-



**Figure 5: Comparison of the number of requesting UEs  $R(t)$  against time  $t$**



**Figure 6: Comparison of the average serving time  $T_s$  against deliver factor  $\epsilon$**

erage serving time for D2D content delivery. Simulations have been conducted and the results show that the delivery performance of D2D multicast is dependent on transmission environment parameters such as path loss factor, transmitter density and decoding threshold. In addition, the delivery performance also depends on the delivery strategy related to the popularity of content, an important factor that needs to be considered in scheduling. The work presented in the paper demonstrates that D2D communication technique can be effective to improve the resource efficiency of cellular networks in the application to content sharing, and the proposed architecture and methods are the concrete implementation to achieve this.

## 6. REFERENCES

- [1] 5G radio access—research and vision. white paper, Ericsson, 2012.
- [2] A. Altieri, P. Piantanida, L. Vega, and C. Galarza. A stochastic geometry approach to distributed caching in large wireless networks. In *Proc. 2014 ISWCS*, pages 863–867, 2014.
- [3] A. Asadi, Q. Wang, and V. Mancuso. A survey on device-to-device communication in cellular networks. *IEEE Communications Surveys and Tutorials*, 16(4):1801–1819, 2014.
- [4] E. Bastug, M. Bennis, and M. Debbah. Living on the edge: The role of proactive caching in 5g wireless networks. *IEEE Transactions on Mobile Computing*, 52(8):82–89, 2014.
- [5] D. Chhajed and T. J. Lowe, editors. *Building Intuition: Insights From Basic Operations Management Models and Principles*. Springer US, 1 edition, 2008.
- [6] Erturk, S. Mukherjee, H. ISHII, and H. Arslan. Distributions of transmit power and SINR in device-to-device networks. *IEEE Communications Letters*, 17(2):273–276, 2013.
- [7] G. Fodor, E. Dahlman, G. Mildh, and S. Parkvall. Design aspects of network assisted device-to-device communications. *IEEE Communications Magazine*, 50(3):170–177, 2012.
- [8] N. Golrezaei, A. Dimakis, and A. Molisch. Scaling behavior for device-to-device communications with distributed caching. *IEEE Transactions on Information Theory*, 60(7):4286–4298, 2014.
- [9] N. Golrezaei, P. Mansourifard, A. Molisch, and A. Dimakis. Base-station assisted device-to-device communications for high-throughput wireless video networks. *IEEE Transactions on Wireless Communications*, 13(7):3665–3676, 2014.
- [10] A. Jeffrey and D. Zwillinger. *Table of integrals, series, and products*. Academic Press, 7th edition, 2007.
- [11] N. Jindal, S. Weber, and J. G. Andrews. Fractional power control for decentralized wireless networks. *IEEE Transactions on Wireless Communications*, 7(12):5482–5492, 2008.
- [12] H. Kang, K. Park, K. Cho, and C. Kang. Mobile caching policies for device-to-device (d2d) content delivery networking. In *Proc. 2014 ICC*, pages 299–304, 2014.
- [13] D. Malak and M. Al-Shalash. Optimal caching for device-to-device content distribution in 5g networks. In *Proc. 2014 Globecom Workshops*, pages 863–868, 2014.
- [14] J. Paakkonen, C. Hollanti, and O. Tirkkonen. Device-to-device data storage for mobile cellular systems. In *Proc. 2013 IEEE Globecom Workshops*, pages 671–676, 2013.
- [15] Y. Sagduyu and Y. Shi. Navigating a mobile social network. *IEEE Wireless Communications*, 22(5):122–128, 2015.
- [16] S. Weber, J. G. Andrews, and N. Jindal. The effect of fading, channel inversion, and threshold scheduling on ad hoc networks. *IEEE Transactions on Information Theory*, 53(11):4127–4149, 2007.
- [17] Y. Xu. On the performance of device-to-device communications with delay constraint. *IEEE Transactions on Vehicular Communications*, to be published. Early Access.
- [18] Y. Xu, Y. Liu, and D. Li. Resource management for interference mitigation in device-to-device communication. *IET Communications*, 9(9):1199–1207, 2015.
- [19] Y. Zhao, Y. Li, Y. Cao, T. Jiang, and N. Ge. Social-aware resource allocation for device-to-device communications underlying cellular networks. *IEEE Transactions on Wireless Communications*, 14(12):6621–6634, 2015.