

Research Article

A Privacy-Preserving Caching Scheme for Device-to-Device Communications

Yuqing Zhong ¹, Zhaohua Li ², and Liping Liao ³

¹Communication Research Center, Guangzhou Power Supply Bureau, Guangzhou 510600, China

²Guangdong Electric Power Design and Research Institute, China Energy Construction Group, Guangzhou 510600, China

³Guangdong Polytechnic Normal University, Guangzhou, China

Correspondence should be addressed to Liping Liao; liping1110@hotmail.com

Received 26 October 2020; Revised 25 November 2020; Accepted 19 January 2021; Published 8 February 2021

Academic Editor: Weizhi Meng

Copyright © 2021 Yuqing Zhong et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

With device-to-device (D2D) communication, user equipment can share data with each other without the involvement of network infrastructures. In order to maintain the Quality of Service (QoS) and Quality of Experience (QoE) for user applications in D2D communications, most existing schemes use proactive content caching that needs to predict content popularity before making caching decisions which may result in privacy leakage, since the information of users is collected to train a deep learning-based model to predict content popularity. Therefore, it is crucial to guarantee secure data collection in machine learning-based framework. In this paper, we propose a privacy-preserving D2D caching scheme with a passive content caching strategy based on node importance, which can deliver more efficient caching and prevent the potential leakage of user privacy. The scheme is based on software-defined networking (SDN), in which the controller is responsible for calculating node importance of devices according to the information of requests and encounters collected by SDN switches. Base station will decide which device can establish reliable and secure communication with content requester based on historical information. The simulation results show that the proposed strategy can outperform other D2D caching strategies in terms of cache hit rate and data rate.

1. Introduction

During the recent years, with the rapid development of the mobile Internet, mobile data traffic has increased exponentially. A Cisco VNI report predicted that 79% of global mobile data traffic will be mainly generated from access to video content by 2022 [1]. The current wireless network is facing huge challenges. To reduce the backhaul traffic and base station load, device-to-device (D2D) communication technology has emerged. User devices can directly establish D2D communication links with other devices within the communication range in the D2D communication network without the use of base stations or other access points.

Research results have shown that most mobile traffic is generated from repeated access to popular content [2, 3]. By deploying caches in the core network [4], access network [5–8], and user device [9, 10] in the 5G mobile communication network architecture, the popular content can be

cached at the network edges, which can effectively reduce network congestion and improve network performance. The content is cached to the user devices, and a D2D caching network is constructed. Thus, user devices can share content via the D2D communication link without the use of access points, which can effectively reduce backhaul traffic and base station load. The caching mechanism can determine the caching position and caching contents, is the kernel of the D2D caching network, and determines the caching performance. However, due to the heterogeneity of D2D caching network, it is more vulnerable to security and data privacy threats. It is important to find the trade-off between the performance of security and the cost of protection. Most existing D2D caching mechanisms belong to a proactive caching mechanism that requires content popularity in advance. Before a user sends a request, the content with higher popularity is cached into the D2D caching network in advance. Some owners of mobile devices may be curious

about the content cached in their devices and scan the caching content, which may result in privacy leaking [11].

During recent years, short videos distributed by users have attracted massive traffic. The QuestMobile report indicates that the time spent by users in short video applications was 5.5% of the total time spent on mobile applications in 2017. With the quick development of short video applications, such as Bilibili and TikTok, a user can distribute short videos whenever and wherever possible. Short videos can be distributed and requested in a random and burst manner. After popular videos are distributed, they are accessed by massive users in a short period. For example, one user distributed a short video on TikTok, and it was accessed 30,000 times in a couple of minutes. In such cases, a proactive caching strategy cannot predict the popularity of popular video content in time. Thus, the cached contents cannot be updated in real time, which results in massive ineffective caching and wastes network resources.

Motivated by this, we propose a privacy-preserving D2D caching scheme with passive content caching based on node importance to update the cached contents in real time, increase the cache hit rate, and preserve the user privacy. Since the proposed caching scheme adopts a passive caching strategy, the device can only cache the content that it requests, which prevents the leakage of user privacy. By using the network coding technology [12], the diversity of the cached contents can also be improved in limited cache size. It has been proved that using network coding in content caching can improve the security and performance of the caching system [13]. The SDN switch collects the history request of the terminal devices and meeting information among devices, and the SDN controller can compute the node importance of the terminal devices using the history information. Base station decides which content holder can establish reliable and secure D2D communication link with the content requester. The cached contents can be updated in real time based on user requests and node importance.

The contributions of this paper include the following:

- (i) Firstly, we introduce a privacy-preserving D2D caching scheme with passive content caching based on node importance to improve the security and performance of D2D caching network. The device with higher node importance and social trust will be selected to establish reliable and secure D2D communication link with the content requester.
- (ii) Then, we define the node importance as the weighted sum of the physical intimacy and request similarity between devices, which also reflects the social trust of the device.
- (iii) Finally, we evaluate the performance of the proposed D2D caching strategy and other two well-known caching strategies. The simulation results show that the caching strategy based on node importance proposed by this paper could effectively improve network performance compared to the other two caching strategies.

The remainder of the paper is organized as follows. We introduce related works in Section 2. In Section 3, we define the node importance and propose a privacy-preserving D2D caching strategy based on node importance. Simulation results are presented in Section 4. Finally, we conclude the paper in Section 5.

2. Related Work

If the popular content is cached in the user devices, it can effectively reduce the backhaul traffic and the downloading delay for users. Thus, service quality and user experience can be improved. Currently, most D2D caching mechanisms are based on a proactive caching strategy, and the content popularities are assumed to be known. Golrezaei et al. [14, 15] divided the D2D network into multiple D2D clusters, and only the devices in one cluster can establish the D2D communication link. Based on this, the authors proposed two in-cluster D2D caching strategies to improve network performance of cellular networks, including deterministic cache and random cache based on Zipf. In the deterministic caching mechanism [14], k devices can cache nonrepeated k contents with the top popularities in one virtual cluster and each device only caches one content. In the random caching mechanism based on the Zipf distribution [14], within one virtual cluster, each device can independently and randomly cache content and the popularities of the content cached in the cluster obeying Zipf distribution. Wang et al. proposed a novel D2D caching strategy based on mobile perception, which takes the mobility of users into account. The low-speed and high-speed moving user devices cache content with top popularities, and user devices with middle-speed moving cache content with lower popularities. Thus, the offloading rate can be improved [16]. Chen et al. modelled the offloading benefits and energy consumption of content holders and proposed a proactive caching strategy and user-oriented protocol to obtain higher offloading benefits with lower energy consumption [17]. Malak et al. extended the caching mechanism based on the geographical position and proposed a space-based caching strategy to improve the cache hit rate. To reduce cache redundancy and improve the diversity of cached contents, all the devices in the mutual exclusion area cannot cache the same content [18]. Wu et al. proposed a distributed D2D caching strategy that considers the characteristics of different requests and demands of physical links [19]. Besides a proactive caching strategy that assumes that the content popularity is known, partial D2D caching mechanisms predict content popularity by algorithms such as machine learning. After the future content popularities are predicted, the caching mechanism is determined. Jiang et al. modelled the D2D caching as the multiagent and multiarm gaming machine and determined the cache by using reinforcement learning to reduce downloading delay for users [20]. Li et al. optimized content caching and content distribution jointly to reduce transmission delay and power consumption. They deployed two potential recurrent neural network models, echo state network (ESN) and long short-term memory (LSTM), to predict mobility of users and future popularity of

content and then determined the cached contents and cache position. The authors also proposed a content distribution mechanism based on deep reinforcement learning to improve user experience [21]. In another study [22], the authors proposed a proactive caching strategy based on the association between users and content. They predicted the content popularity by using machine learning and collaborative filtering technology. The content with higher popularities is precached to the base stations and user devices in the low peak period to alleviate the backhaul congestion. Chen and Yang proposed a D2D caching strategy based on user preference to improve the offloading rate, which predicts user preferences by using a collaborative filtering algorithm based on the model and then makes caching decisions [23].

In summary, most D2D caching strategies are based on proactive caching mechanisms and the content popularities must be known or be predicted. Then, the caching decision problem is transformed into optimization problems to find a solution. The prediction accuracy of future content popularities decides the caching performance. However, the rise of short video applications makes the proactive caching strategy fail to precisely predict content popularities and update cached content in real time. Thus, caching performance is reduced. Motivated by this, we propose a passive D2D caching strategy based on node importance, and the cached content can be updated in real time according to user requests. Moreover, network coding technology is employed, which can improve the diversity of the cached content and increase the cache hit rate and data rate without increasing the cache size.

3. System Design

3.1. Defining Node Importance. In this paper, we propose a software-defined passive D2D caching strategy (NIC) based on node importance. Different devices have different node importance in the D2D network. In this system, the SDN switch collects meeting information among user devices and information of content requested by user devices to compute the physical intimacy and request similarity among devices. The node importance is defined as the weighted sum of the physical intimacy [24] and request similarity. For the physical network layer, the user devices with higher node importance will have a higher probability to establish stable D2D communication links with other devices in the future. For content requesting, the user devices with higher node importance have a higher probability to request the same content as other devices. In other words, the user devices with higher node importance have a higher probability of providing other user devices with requested content in the future. Thus, the devices with higher importance will cache content with higher popularities, and the original blocks are cached to reduce the network coding and decoding time and computing consumption. The devices with lower importance will cache contents with lower popularities, and the coding blocks are cached to distribute contents in the network more reasonably, improve caching diversity, and increase the caching efficiency without increasing cache size.

The physical intimacy between user devices indicates the probability of establishing reliable D2D communication links between two user devices in the future. When device D_i and device D_j are within the D2D communication range, the two devices may have D2D communication potential, which is recorded as one meeting between them. The duration is recorded as the meeting duration. Zhang et al. [24] proved that the user meeting time obeys the gamma distribution $\Gamma(k, \theta)$, namely, $X \sim \Gamma(k = (M_{ij}^2/I_{ij}), \theta = (I_{ij}/M_{ij}))$, wherein

$$M_{ij} = \frac{f_n X_n}{N_{ij}},$$

$$I_{ij} = \frac{\sum_n (X_n - M_{ij})^2}{N_{ij}},$$
(1)

where X_n indicates the n^{th} meeting duration of device D_i and device D_j and N_{ij} indicates the meeting count between device D_i and device D_j . The physical intimacy $c_{ij} \in [0, 1]$ can be expressed as [24]

$$c_{ij} = 1 - \int_0^{X_{\min}} f(u; k, \theta) du = 1 - \frac{\gamma(k, (X_{\min}/\theta))}{\Gamma(k)},$$

$$\gamma\left(k, \frac{X_{\min}}{\theta}\right) = \int_0^{(X_{\min}/\theta)} t^{k-1} e^{-t} dt,$$
(2)

where X_{\min} is the minimal meeting duration required by two user devices to successfully transfer one file via the D2D communication link. Thus, the average physical intimacy \bar{c}_i between device D_i and other devices in the D2D caching network is given by following:

$$\bar{c}_i = \frac{\sum_{j=1}^n c_{ij}}{n},$$
(3)

where n is the quantity of the devices within the D2D communication range of device D_i . The higher the physical intimacy of the user device is, the higher its probability of establishing a reliable D2D communication link with other user devices in the future will be. As shown in Figure 1, the average physical intimacy \bar{c}_1 of user device 1 is

$$\bar{c}_1 = \frac{\sum_{j=2}^5 c_{1j}}{4}.$$
(4)

Based on the history requests of the device, the SDN controller can compute the request similarity $s_{ij} \in [0, 1]$ between user devices by using cosine similarity, which is expressed as follows:

$$s_{ij} = \cos(\mathbf{w}_i, \mathbf{w}_j) = \frac{\mathbf{w}_i \cdot \mathbf{w}_j}{\|\mathbf{w}_i\|_2 \|\mathbf{w}_j\|_2},$$
(5)

where \mathbf{w}_i and \mathbf{w}_j indicate the interest vector of device D_i and device D_j , respectively. The average request similarity \bar{s}_i of device D_i and the other devices in the D2D caching network is

$$\bar{s}_i = \frac{\sum_{j=1}^n s_{ij}}{n}, \quad (6)$$

where n is the quantity of the devices within the D2D communication range of device D_i . The higher the average request similarity of the user device is, the higher the overlapping degree of requested contents with other devices will be. In this paper, we propose a passive caching mechanism, namely, the user device only caches the ever-requested content to satisfy other users' requests by providing other users with the desired content. The cached content in the user device with higher average request similarity will have a higher probability to be requested by other users. In this paper, we define the node importance as the weighted sum of the physical intimacy and request similarity. The devices with higher node importance will have a higher probability to provide other user devices with requested content and receive higher caching benefits. The node importance I_i of user device D_i is given by the following equation:

$$I_i = \alpha \bar{c}_i + \beta \bar{s}_i, \quad (7)$$

where \bar{c}_i is the normalized average physical intimacy of device D_i , \bar{s}_i is the normalized average request similarity of device D_i , and $\alpha \in [0, 1]$ and $\beta \in [0, 1]$ are the design parameters and indicate the importance of physical intimacy and request similarity. In this paper, α and β are set as 0.5.

3.2. D2D Caching Strategy Based on Node Importance. In this system, to improve the caching efficiency, the base station divides the content into m content blocks with the same size. When device D_i requests content f , it will first send the request packet to the base station. If the D2D caching network includes this content, the base station will locate a group of user devices with hit caches and determine the caching plan according to the node importance. If I_i of device D_i is higher, then device D_i caches the original blocks. If the node importance I_i of device D_i is lower, then device D_i caches the coding blocks. This is because the device with higher node importance has higher request similarity with other devices, namely, the cached content has higher probabilities to be requested by other devices. The devices with higher node importance have a higher probability to successfully establish D2D communication links with other users and can ensure successful content transfer. Therefore, the devices with higher node importance will cache content with higher popularities and cache original content blocks to improve the cache hit rate and reduce decoding/coding time and computing consumption. The devices with lower node importance will cache contents with lower popularities and cache coding blocks. Each coding block should include all the original block information. This can improve the content diversity of the caching system without increasing the cache size.

In the D2D caching network, the base station maintains the D2D caching network information and locates a group of user device with hit caches, namely, the content holder. \mathbf{Info}_k indicates the caching information of cached content

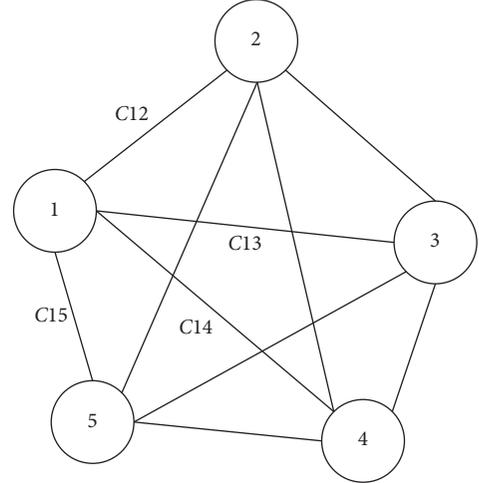


FIGURE 1: Example of physical intimacy.

f_k in the D2D caching network, and $\mathbf{Info}_k = \mathbf{Info}_k^1, \dots, \mathbf{Info}_k^N$, where N is the number of user devices with cached content f_k . The content holders will be ranked by the node importance in descending order. \mathbf{Info}_k^1 indicates the caching information of content f_k cached in the user devices with highest node importance, and the caching information is expressed as follows:

$$\mathbf{Info}_k^1 = \{D_i, I_i, \mathbf{V}_i, n_i\}, \quad (8)$$

where D_i is the device ID, I_i is the node importance of device D_i , n_i indicates the quantity of the cached original blocks or coding blocks, $\mathbf{V}_i = \{v_{i1}, \dots, v_{im}\}$, and $v_{ij} \in \{0, 1\}$ indicates if the cached content block in device D_i includes the information of the original content block j . If it is included, then $v_{ij} = 1$; otherwise, $v_{ij} = 0$. When \mathbf{V}_i is an all-1 vector and $n_i \neq m$, it indicates that the user device caches the coding blocks of content f_k .

When device D_i requests content f_k , the base station will find the D2D cache information table and check if the D2D caching network can satisfy the user requests. If the quantity of the content blocks cached in the D2D caching network is more than or equal to m , then the base station selects a group of content holders with higher node importance, and the selected content holders will establish D2D communication links with the device D_i to transfer corresponding content blocks. If the D2D caching network cannot satisfy the user requests, the base station will send $(m - m')$ content blocks to respond to device D_i , where m' is the quantity of the content blocks cached in the D2D caching network, as described in Algorithm 1. The complexity of Algorithm 1 is $O(n)$, where n is the number of content holders.

In this system, the base station selects a group of content holders with higher node importance and social trust to establish reliable and secure D2D communication links to transmit desired content. Moreover, the base station is also responsible for making caching decisions and instructing how the user devices cache the received content blocks. If the request is from users with higher node importance, then it indicates that caching the content in content requester will

Input: α, β ;
Output: f_q, nc // f_q is the cache identifier, nc is the code identifier, 0 indicates the original block, and 1 is the coding block;

- (1) Initialization: $f_q = 0, nc = 0$;
- (2) The SDN switch collects the request records and interactive information of the devices and periodically sends it to the SDN controller;
- (3) The SDN controller computes the node importance of the device according to the history information collected by SDN switches, namely, I_i ;
- (4) **While** BS receives the request from device D_i for content f_k **do**
- (5) **if** the node importance of device D_i is higher **then**
- (6) Make $f_q = 1, nc = 0$;
- (7) **else**
- (8) Make $f_q = 1, nc = 1$;
- (9) **end if**
- (10) **if** $m' (m' \geq m)$ content blocks are cached in the D2D caching network **then**
- (11) BS locates a group of cached content holder D_j with the top node importance to respond to the user request;
- (12) **for** each cached content holder D_j **do**
- (13) **if** the cached content is the original block **then**
- (14) make $f_q = 1$;
- (15) **else**
- (16) make $f_q = 0$;
- (17) **end if**
- (18) BS sends data packets (f_k, n_j, f_q) to the content holder; // n_j is the number of content blocks to be sent by content holder D_j to the requester
- (19) After content holder D_j receives the data packets from BS, it will establish the D2D communication link with device D_i and transfer the data packet $(f_k, \text{block}(s), f_q)$; // $\text{block}(s)$ is the coding block or content block;
- (20) **end for**
- (21) **else**
- (22) **if** $nc = 0$ **then**
- (23) BS sends $(m - m')$ original blocks to the content requester device D_i , namely, $(f_k, \text{blocks}, f_q)$;
- (24) **else**
- (25) BS sends $(m - m')$ coding blocks to the content requester device D_i , namely, $(f_k, \text{blocks}, f_q)$;
- (26) **end if**
- (27) **end if**
- (28) **end while**

ALGORITHM 1: D2D caching strategy based on node importance.

bring benefits with a higher probability, i.e., higher cache hit rate. When the base station responds to the device request, it will send the original block and instruct the user devices to cache the original block. Otherwise, when the base station responds to the user request, it sends the coding blocks generated with all the original blocks and instructs the user devices to cache the coding blocks. For details, refer to Algorithm 1. To prevent users from receiving the linearly dependent coding blocks, the user device only caches the coding blocks sent by the base station and does not cache coding blocks obtained by D2D communication links.

When content requester D_i receives the content blocks from the base stations or other devices, it will decide whether to cache the received contents based on cache identifier f_q . If $f_q = 1$, then the content will be cached locally; otherwise, the content will not be cached. Since the cache size of the device is limited, when the caches are replaced, the user device codes the replaced content blocks into a coding block. In this way, all the original block

information of this content will be reserved while the cache size is released to improve the diversity of the caching contents.

4. Experimental Results and Analysis

In the simulation test, the radius of the base station was 500 m and a total of 100 devices were provided. The maximum D2D communication distance was 100 m. The D2D communication belongs to the in-band communication, namely, the D2D communication shares the bandwidth with cellular communication [25]. The cache size of the device included $\{1, 2, 5, 8, 10, 15\}$ files, and the number of files was 500. The user request obeyed the Poisson distribution, the popularities of the content obeyed the Zipf distribution, and the Zipf parameter was $\alpha \in \{0.56, 0.8, 1, 1.2, 1.5\}$. In the simulation test, the passive caching strategy NIC in this paper was compared with the passive caching strategy, Leave Copy Everywhere (LCE), and the proactive caching strategy, Most Popular Cache (MPC). For the LCE, all the received

content will be cached by requester. For the MPC, the most popular content will be cached into devices in advance. The strategies were assessed by the cache hit rate and data rate. The cache hit rate was defined as the ratio of the number of requests responded by user devices to the total requests sent by user devices, which is an important parameter to assess the performance of the caching system. When device D_i requests contents from the base station via cellular communication, the data rate $R_{B,i}$ is defined as follows [24]:

$$R_{B,i} = W \log_2 \left(1 + \frac{P_B |h_{Bi}|^2}{\sum_{j' \neq i} \beta_{j'i} P_{j'} |h_{j'i}|^2 + N_0} \right). \quad (9)$$

When device D_i requests contents from device D_j via D2D communication, the data rate is defined as follows [24]:

$$R_{j,i} = W \log_2 \left(1 + \frac{P_j |h_{ji}|^2}{P_B |h_{Bi}|^2 + \sum_{j' \neq j} \beta_{j'i} P_{j'} |h_{j'i}|^2 + N_0} \right), \quad (10)$$

where P_B , P_j , and $P_{j'}$ indicate the transmission power of the base station, device D_j , and device $D_{j'}$; N_0 is the Gauss white noise; $|h_{Bi}|^2$ and $|h_{j'i}|^2$ are the path loss and are related to communication distance; $\beta_{j'i} = 1$ indicates interference; and $\beta_{j'i} = 0$ indicates no interference.

The influences of the cache size on the three caching mechanisms are shown in Figures 2 and 3. As shown in Figure 2, the cache hit rate of the three caching mechanisms increased with the growth of the cache size. This is because with the growth of the cache size of the user device, more files can be cached in the D2D caching network to make more requests responded to by other terminals. In this case, it can reduce the base station load and backhaul traffic to improve the data rate, as shown in Figure 3. As shown in Figures 2 and 3, the cache hit rate and data rate of the NIC caching strategy proposed in this paper were higher than those of the other two caching strategies. The NIC strategy deploys caches and implements differential caching strategies according to the content popularities and node importance. Thus, the content can be distributed more reasonably in the D2D caching network and more content can be obtained via D2D communication. It can reduce the base station load and improve the caching hit rate. Compared to the other two caching strategies, the NIC mechanism can consider the physical intimacy between user devices and make the cache hit node closer to the request devices and obtain a higher data rate.

The influences of the Zipf parameters on the three caching mechanisms are shown in Figures 4 and 5. The bigger the Zipf parameter is, the more similar the user requests will be, and there will be a greater repeated request time for a small part of the files. As time elapses, the cached contents in the D2D network increase, and the cached contents are centralized in a small number of files. Then,

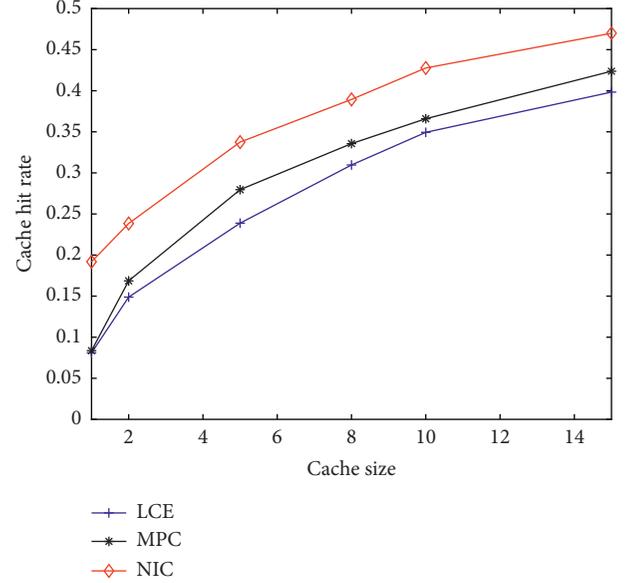


FIGURE 2: Influences of cache size on the cache hit rate.

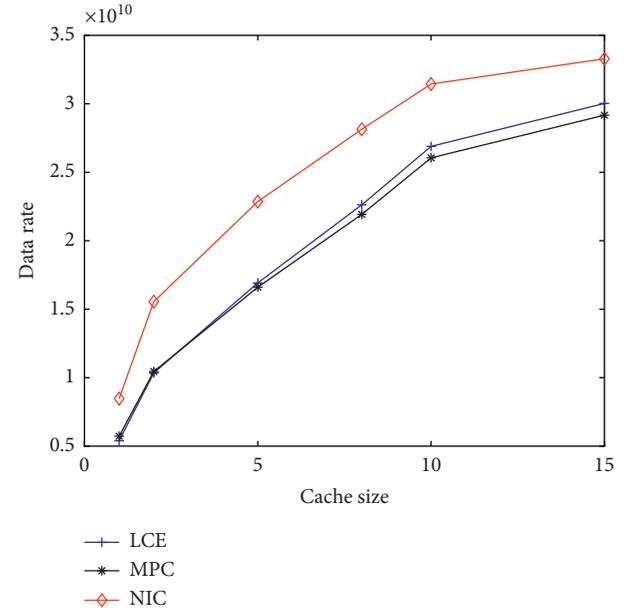


FIGURE 3: Influences of cache size on the data rate.

users have a higher probability to receive files via the D2D communication link. Thus, the cache hit rate and data rate of the three caching mechanisms grow as Zipf parameters enlarge. The caching performance of NIC was always superior to the other two caching strategies, which was more significant when the Zipf parameter was smaller. The reason for this is that NIC can improve the diversity of the cached contents and improve caching performance by using network coding technology without increasing the cache size.

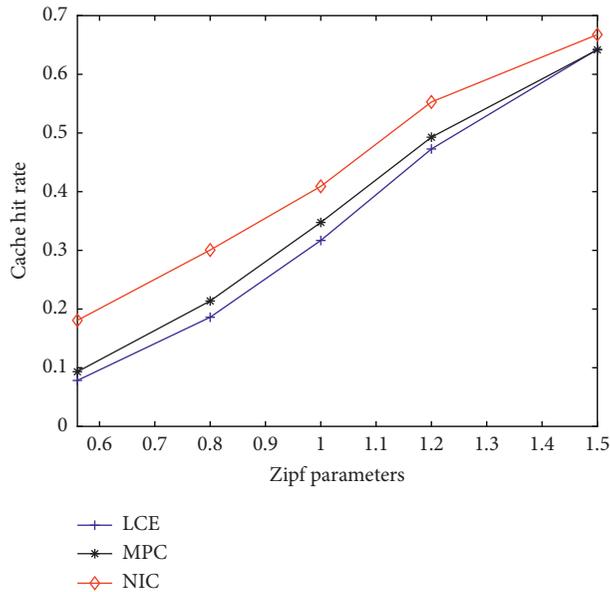


FIGURE 4: Influences of the Zipf parameters on the cache hit rate.

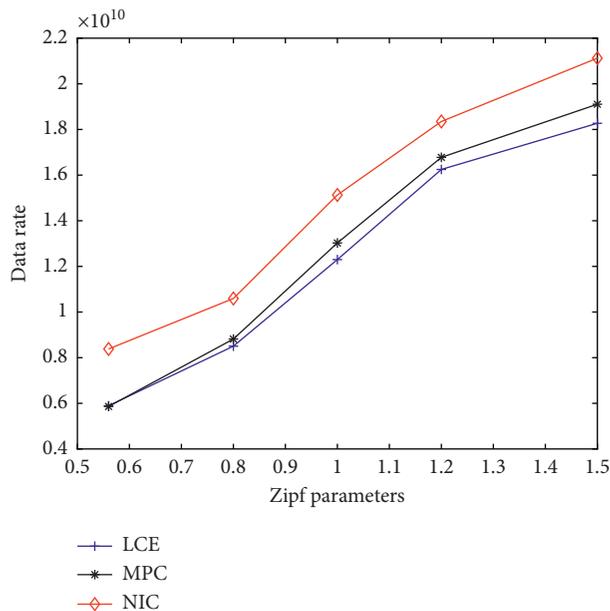


FIGURE 5: Influences of the Zipf parameters on the data rate.

5. Conclusion

In this paper, we propose a privacy-preserving D2D caching strategy based on the node importance, which can improve the diversity of the cached content by using the network coding technology and preserve the privacy of users and data. The SDN switch collects the information of requests and meeting information of the devices, and the SDN controller can compute the physical intimacy and request similarity between user devices and other devices by using the history information collected by SDN switches to obtain the node importance of the device. The node importance is

defined as the weighted sum of the physical intimacy and request similarity. The nodes with higher importance have higher cache benefits and will cache the original blocks; the nodes with lower importance have lower cache benefits and will cache the coding blocks to make the cached contents be distributed more reasonably in the network. Base station will decide which device can establish reliable and secure communication with the requester based on historical information, which reflects the importance and social trust of devices. The simulation results show that the caching strategy based on node importance in this paper could improve the cache hit rate and data rate and effectively improve the performance and security of the caching network compared to the other two proactive and passive caches.

Data Availability

No data were used to support this study.

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

Acknowledgments

This research was funded by the National Natural Science Foundation of China (62002072) and the Guangzhou Science and Technology Project (GZHKJXM20170117).

References

- [1] C. Cisco, *Visual Networking Index: Global Mobile Data Traffic Forecast Update*, 2016 WhitePaper, German, Netherlands, 2017.
- [2] L. Qiu and G. Cao, "Popularity-Aware caching increases the capacity of wireless networks," *IEEE Transactions on Mobile Computing*, vol. 19, no. 1, pp. 173–187, 2020.
- [3] Y. Zhou, L. Chen, C. Yang, and D. M. Chiu, "Video popularity dynamics and its implication for replication," *IEEE Transactions on Multimedia*, vol. 17, no. 8, pp. 1273–1285, 2015.
- [4] I. Parvez, A. Rahmati, I. Guvenc, A. I. Sarwat, and H. Dai, "A survey on low latency towards 5G: RAN, core network and caching solutions," *IEEE Communications Surveys and Tutorials*, vol. 20, no. 4, pp. 3098–3130, 2018.
- [5] K. Guo, C. Yang, and T. Liu, "Caching in base station with recommendation via Q-learning," in *Proceedings of the 2017 IEEE Wireless Communications and Networking Conference (WCNC)*, pp. 1–6, San Francisco, CA, USA, March 2017.
- [6] P. Blasco and D. Gündüz, "Learning-based optimization of cache content in a small cell base station," in *Proceedings of the 2014 IEEE International Conference on Communications (ICC)*, pp. 1897–1903, Sydney, Australia, June 2014.
- [7] 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 Communications Magazine*, vol. 52, no. 2, pp. 131–139, 2014.
- [8] Y. Wang, Y. Chen, H. Dai, Y. Huang, and L. Yang, "A learning-based approach for proactive caching in wireless communication networks," in *Proceedings of the 2017 9th International Conference on Wireless Communications and*

- Signal Processing (WCSP)*, pp. 1–6, Nanjing, China, October 2017.
- [9] B. Bai, L. Wang, Z. Han, W. Chen, and T. Svensson, “Caching based socially-aware D2D communications in wireless content delivery networks: a hypergraph framework,” *IEEE Wireless Communications*, vol. 23, no. 4, pp. 74–81, 2016.
- [10] M. Afshang, H. S. Dhillon, and P. H. J. Chong, “Fundamentals of cluster-centric content placement in cache-enabled device-to-device networks,” *IEEE Transactions on Communications*, vol. 64, no. 6, pp. 2511–2526, 2016.
- [11] L. Xiao, X. Wan, C. Dai, X. Du, X. Chen, and M. Guizani, “Security in mobile edge caching with reinforcement learning,” *IEEE Wireless Communications*, vol. 25, no. 3, pp. 116–122, 2018.
- [12] R. Ahlswede, N. Ning Cai, S.-Y. R. Li, and R. W. Yeung, “Network information flow,” *IEEE Transactions on Information Theory*, vol. 46, no. 4, pp. 1204–1216, 2000.
- [13] J. Wang, J. Ren, K. Lu, J. Wang, S. Liu, and C. Westphal, “A minimum cost cache management framework for information-centric networks with network coding,” *Computer Networks*, vol. 110, pp. 1–17, 2016.
- [14] 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 Transactions on Wireless Communications*, vol. 13, no. 7, pp. 3665–3676, 2014.
- [15] 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 Communications Magazine*, vol. 51, no. 4, pp. 142–149, 2013.
- [16] R. Wang, J. Zhang, S. H. Song, and K. B. Letaief, “Mobility-Aware caching in D2D networks,” *IEEE Transactions on Wireless Communications*, vol. 16, no. 8, pp. 5001–5015, 2017.
- [17] B. Chen, C. Yang, and A. F. Molisch, “Cache-Enabled device-to-device communications: offloading gain and energy cost,” *IEEE Transactions on Wireless Communications*, vol. 16, no. 7, pp. 4519–4536, 2017.
- [18] D. Malak, M. Al-Shalash, and J. G. Andrews, “Spatially correlated content caching for device-to-device communications,” *IEEE Transactions on Wireless Communications*, vol. 17, no. 1, pp. 56–70, 2018.
- [19] K. Wu, M. Jiang, F. She, and X. Chen, “Relay-aided request-aware distributed packet caching for device-to-device communication,” *IEEE Wireless Communications Letters*, vol. 8, no. 1, pp. 217–220, 2019.
- [20] W. Jiang, G. Feng, S. Qin, T. S. P. Yum, and G. Cao, “Multi-agent reinforcement learning for efficient content caching in mobile D2D networks,” *IEEE Transactions on Wireless Communications*, vol. 18, no. 3, pp. 1610–1622, 2019.
- [21] L. Li, Y. Xu, J. Yin et al., “Deep reinforcement learning approaches for content caching in cache-enabled D2D networks,” *IEEE Internet of Things Journal*, vol. 7, no. 1, pp. 544–557, 2020.
- [22] E. Bastug, M. Bennis, and M. Debbah, “Living on the edge: the role of proactive caching in 5G wireless networks,” *IEEE Communications Magazine*, vol. 52, no. 8, pp. 82–89, 2014.
- [23] B. Chen and C. Yang, “Caching policy optimization for D2D communications by learning user preference,” in *Proceedings of the 2017 IEEE 85th Vehicular Technology Conference (VTC Spring)*, pp. 1–6, Sydney, Australia, June 2017.
- [24] Y. Zhang, E. Pan, L. Song, W. Saad, Z. Dawy, and Z. Han, “Social network aware device-to-device communication in wireless networks,” *IEEE Transactions on Wireless Communications*, vol. 14, no. 1, pp. 177–190, 2015.
- [25] A. Asadi, Q. Wang, and V. Mancuso, “A survey on device-to-device communication in cellular networks,” *IEEE Communications Surveys & Tutorials*, vol. 16, no. 4, pp. 1801–1819, 2014.