

Artificial Intelligence Assisted Collaborative Edge Caching in Small Cell Networks

Md Ferdous Pervej*, Le Thanh Tan[†], and Rose Qingyang Hu[‡]

*Department of Electrical and Computer Engineering, North Carolina State University, Raleigh, NC 27695, USA

[†]Commonwealth Cyber Initiative, Old Dominion University, Norfolk, VA 23529, USA

[‡]Department of Electrical and Computer Engineering, Utah State University, Logan, UT 84322, USA

Email: mpervej@ncsu.edu, tle@odu.edu, rose.hu@usu.edu

Abstract—Edge caching is a new paradigm that has been exploited over the past several years to reduce the load for the core network and to enhance the content delivery performance. Many existing caching solutions only consider homogeneous caching placement due to the immense complexity associated with the heterogeneous caching models. Unlike these legacy modeling paradigms, this paper considers heterogeneous content preference of the users with heterogeneous caching models at the edge nodes. Besides, aiming to maximize the cache hit ratio (CHR) in a two-tier heterogeneous network, we let the edge nodes collaborate. However, due to complex combinatorial decision variables, the formulated problem is hard to solve in the polynomial time. Moreover, there does not even exist a ready-to-use tool or software to solve the problem. We propose a modified particle swarm optimization (M-PSO) algorithm that efficiently solves the complex constraint problem in a reasonable time. Using numerical analysis and simulation, we validate that the proposed algorithm significantly enhances the CHR performance when comparing to that of the existing baseline caching schemes.

Index Terms—Cache hit ratio, content delivery network, edge caching, particle swarm optimization, small cell network.

I. INTRODUCTION

Owing to the ever growing requirements of enhanced data rates, quality of service, and latency, wireless communication has evolved from generation to generation. With the exponential increase of the connected devices, existing wireless networks have already been experiencing performance bottleneck. While the general trends are shifting resources towards the edge of the network [2]–[4], study shows that mobile video traffic is one of the dominant applications that prompt this bottleneck [5]–[7]. Caching has become a promising technology to address this performance issue by storing popular contents close to the end users [8], [9]. Therefore, during the peak traffic hours, the requested contents can be delivered from these local nodes ensuring a deflated pressure to the backhaul and the centralized core network yielding reduced latency for content delivery. Thus, an edge cache-enabled network utilizes the much-needed wireless spectrum and wireline bandwidth efficiently. In the ultra-dense network platform, caching at the edge nodes is a powerful mechanism for delivering video traffic.

While the caching solution can significantly benefit next-generation wireless communication, it still comes with various

challenges [10]–[13]. First of all, the content selection has an enormous impact on the cache-enabled platform [4], [14]. Then, choosing the appropriate nodes to store the contents needs to be answered. Due to the broad combinatorial decision parameters, this is an immense challenge for any cache-enabled network platform. Furthermore, owing to the necessity of the system performance metrics, the solution to this combinatorial decision problem may change. Therefore, based on the performance metric, an efficient solution is demanded to handle the issues in a reasonable time. As such, under practical modeling with proper communication protocols, a heterogeneous network platform needs to be adopted for evaluating the caching performance.

There exist several caching solutions in the literature [14]–[17]. Caching policy and cooperative distance were designed in [15], by Lee *et al.*, considering clustered device-to-device (D2D) networks. While the authors showed some brilliant concepts for the caching policy design aiming to maximize (a) energy efficiency and (b) throughput, they only considered the collaboration among the D2D users. Lee *et al.* also proposed a base station (BS) assisted D2D caching network in [14] that maximizes the time-average service rate. However, the authors only considered a single BS underlay D2D communication with homogeneous request probability modeling. Tan *et al.* [16] adopted the collaboration based caching model in the heterogeneous network model. A mobility aware probabilistic edge caching approach was explored in [17]. The authors' proposed model considered the novel idea of collaboration by considering the spatial node distribution and user-mobility. While [16], [17] introduces some splendid concepts of relaying and collaborations, the authors only incorporated homogeneous caching placement strategies.

Unlike these existing works, we investigate heterogeneous content preference model leveraging heterogeneous cache placement strategy in this paper. Particularly, in a small cell network (SCN), we incorporate collaborations among spatially distributed full-duplex (FD) enabled BSs and half-duplex (HD) operated D2D users to maximize the average cache hit ratio (CHR). However, the formulated problem contains intricate combinatorial decision variables that are hard to determine in polynomial time. Therefore, we implement a modified particle swarm optimization (M-PSO) algorithm that effectively solves the grand probabilistic cache placement problem within a reasonable time. To the best of our knowledge, this is the first work to consider heterogeneous user preference with a heterogeneous caching model in a practical SCN that uses

The work of M. F. Pervej, L. T. Tan and R. Q. Hu were supported in part by National Science Foundation under grants NeTS 1423348 and EARS 1547312 as well as in part by the Intel Corporation.

This is the **technical report** of [1].

collaborative content sharing among heterogeneous edge nodes to maximize the CHR.

The outline of this paper is as follows. The system model and the proposed content access protocols are presented in Section II, followed by the CHR analysis in Section III. The optimization problem and the proposed M-PSO algorithm are described in Section V. Section VI gives the performance results, followed by the concluding remarks in Section VII.

II. SYSTEM MODEL AND CONTENT ACCESS PROTOCOLS

This section presents the node distributions and describes the caching properties, followed by the proposed content access protocols.

A. Node Distributions

We consider a practical two-tier heterogeneous network, which consists of macro base stations (MBS) and low-power sBSs (or relays) with underlaid D2D users. The nodes are distributed following an independent homogeneous Poisson point processes (HPPP) model. Let us denote the densities of the D2D user, sBS and MBS by λ_u , λ_b and λ_m , respectively. The sBSs and MBSs operate in the FD mode whereas the D2D users operate in the HD mode. Let us denote the set of D2D users, sBSs and MBSs by \mathcal{U} , \mathcal{B} and \mathcal{M} , respectively. Without any loss of generality, user, sBS and MBS are denoted by $u \in \mathcal{U}$, $b \in \mathcal{B}$, and $m \in \mathcal{M}$, respectively. Besides, the communication ranges of these nodes are denoted by R_u , R_b and R_m , respectively.

The requesting user node is named as the tagged user node. While a user is always associated with the serving MBS, it can also associate with a low powered sBS if the association rules are satisfied. The main benefits of being connected to sBS over MBS are higher data rate, less latency, less power consumption, more effective uses of radio resources, etc. We denote the associated sBS as the tagged sBS for that user. Furthermore, if such a tagged sBS exists for the user, the user maintains its communication with the serving MBS via the tagged sBS. In that case, the sBS can also use its FD mode to deliver requested content from the other sBSs or the cloud via the MBS. If such a tagged sBS does not exist for the user, the user will have to rely on the neighbor sBS nodes and the serving MBS for extracting the requested contents. As all the users may not place a content request at the same time, we assume that only α portions of the users act as tagged users. Without any loss of generality, the requesting user, the associated sBS, and the serving MBS are denoted as u_0 , b_0 and m_0 , respectively.

B. Cache Storage, Caching Policy and Content Popularity

The cache storage of the users, sBSs and MBSs are denoted by \mathcal{C}_d , \mathcal{C}_b and \mathcal{C}_m , respectively. Considering equal-sized contents we investigate a probabilistic caching placement [18] where the users can make a content request from a content directory of $\mathcal{F} = \{f_k\}$, where $k \in \{1, 2, \dots, F\}$. For the caching model, a probabilistic method is considered assuming a heterogeneous caching placement strategy. Let $\eta_{f_k}^{u_i}$, $\eta_{f_k}^{b_j}$ and $\eta_{f_k}^{m_l}$ be the probabilities of storing a content $f_k \in \mathcal{F}$ at the cache store of the user node u_i , the sBS b_j and the MBS m_l ,

respectively. Note that probabilistic caching is highly practical and adopted in many existing works [5], [6], [14]–[18].

The content popularity is modeled by following the Zipf distribution with the probability mass function $P_{f_k} = \frac{k^{-\gamma}}{\sum_{k=1}^F k^{-\gamma}}$. Note that the skewness γ governs this distribution. It is assumed that each user has a different content preference. Therefore, a random content preference order and a random skewness are chosen for each user. While the content order is chosen using random permutation, the parameter, γ , is chosen following Uniform random distribution within a range of maximum γ^{max} and minimum γ^{min} values. Without any loss of generality, the probability that user u_0 requests for content f_k is denoted by $\rho_{f_k}^{u_0}$. This is modeled based on the Zipf distribution.

C. Proposed Content Access Protocol

For accessing the contents, the following practical cases are considered.

Case 1 - Local/self cache hit: If a tagged user requests the content that is previously cached, the user can directly access the content from its own storage.

Case 2 - D2D cache hit: If the required content is not stored in its own storage, the tagged user sends the content request to the neighboring D2D nodes. If any of the neighbors has the content, the user can extract the content from that neighboring user.

Case 3 - sBS cache hit: If the tagged user is under the communication range of any sBS, it maintains its communication via the tagged sBS. In this particular case, we have the following sub-cases:

Case 3.1: If the requested content is in the tagged sBS cache, it can access the content directly from there. We denote this case as a direct cache hit from the tagged sBS.

Case 3.2: If the content is not stored in the tagged sBS cache but is available in one of the neighboring sBSs, the tagged sBS extracts the content from the neighboring sBS via its FD capability and delivers it to the tagged user. We denote this term as soft-sBS (SsBS) cache hit.

Case 3.3: If the requested content is not available in any of the sBSs, the tagged sBS forwards the request to the serving MBS. If the content is in the serving MBS, it is delivered to the tagged sBS and then to the user. This case is denoted as the sBS-MBS cache hit.

Case 3.4: If all of the above sub-cases fail, then the MBS extracts the content from the cloud using its FD capability. The sBS extracts the content from the MBS using its own FD capability and delivers it to the tagged user. This case is denoted as the sBS cache miss.

Case 4 - MBS cache hit: If the tagged user is not in the communication range of any of the sBSs, it has to rely on the serving MBS for its communication. In this case, we consider the following sub-cases:

Case 4.1: If the requested content is available in the MBS cache, the content is directly delivered to the tagged user. This case is denoted as an MBS cache hit.

Case 4.2: If the content is not available in the MBS storage and the above case fails, the MBS extracts the content from the cloud using its FD capability. Then, the content is directly

delivered to the user. This case is referred as an MBS cache miss.

Without loss of generality, *Case 3* and (*Case 4*) are denoted by the indicator function \mathbb{I}_s and \mathbb{I}_m , respectively. Note that, in *Case 3*, if the tagged user is in the communication ranges of multiple sBSs, it gets connected to the one that provides the best received power.

III. EDGE CACHING: CACHE HIT RATIO ANALYSIS

In this section, we analyze and calculate the local cache hit probabilities.

A. Caching Probabilities

We now analyze the cache hit probability at different nodes for the cases mentioned in Section II-C. Note that a cache hit occurs at a node, if a requested content is available in that node.

1) *Case 1 - Local/self cache hit*: The local cache hit probability is denoted as $P_o^u = \eta_{f_k}^{u_0}$, i.e. the probability of storing the content f at the self cache storage of the tagged user.

2) *Case 2 - D2D cache hit*: The cache hit probability for the D2D nodes can be calculated as follows:

$$P_d^u = \left(1 - \eta_{f_k}^{u_0}\right) \left[1 - \prod_{u_i \in \Phi_u \setminus u_0} \left(1 - \eta_{f_k}^{u_i}\right)\right], \quad (1)$$

where $\prod_{u_i \in \Phi_u} \left(1 - \eta_{f_k}^{u_i}\right)$ means that none of the Φ_u active neighbors (D2D nodes) in its communication range have the content. Thus, the complement of that is the probability that at least one of the users stores the content.

3) *Case 3 - sBS cache hit*: In this case, we calculate the cache hit probabilities achieved via the tagged sBS for the respective sub-cases.

Case 3.1: At first, the probability of getting a requested content from the tagged sBS is calculated as follows:

$$P_{b_o}^u = \left(1 - \eta_{f_k}^{u_0}\right) \prod_{u_i \in \Phi_u \setminus u_0} \left(1 - \eta_{f_k}^{u_i}\right) \eta_{f_k}^{b_0}. \quad (2)$$

Case 3.2: The probability of getting a requested content from one of the neighbor sBSs is considered in this sub-case. Essentially, this case states that a cache miss has occurred at the tagged sBS. Mathematically, we express this as follows:

$$P_B^u = \left(1 - \eta_{f_k}^{u_0}\right) \prod_{u_i \in \Phi_u \setminus u_0} \left(1 - \eta_{f_k}^{u_i}\right) \left(1 - \eta_{f_k}^{b_0}\right) \left(1 - \prod_{b_j \in \Phi_b \setminus b_0} \left(1 - \eta_{f_k}^{b_j}\right)\right), \quad (3)$$

where Φ_b is the set of active neighboring sBSs that are in the communication range of the tagged sBS.

Case 3.3: If sub-case 3.1 and 3.2 fail, the content request is forwarded to the serving MBS via the tagged sBS. The cache hit probability, for this case, is calculated as follows:

$$P_{M_{\mathbb{I}_s}}^u = \left(1 - \eta_{f_k}^{u_0}\right) \prod_{u_i \in \Phi_u \setminus u_0} \left(1 - \eta_{f_k}^{u_i}\right) \left(1 - \eta_{f_k}^{b_0}\right) \prod_{b_j \in \Phi_b \setminus b_0} \left(1 - \eta_{f_k}^{b_j}\right) \eta_{f_k}^{m_0}. \quad (4)$$

When $\mathbb{I}_s = 1$ - the tagged user is in the communication range of at least one of the sBS, from the above cases and sub-cases, we calculate the total cache hit probability as follows:

$$P_{\mathbb{I}_s}^u = 1 - \left[\left(1 - \eta_{f_k}^{u_0}\right) \prod_{u_i \in \Phi_u \setminus u_0} \left(1 - \eta_{f_k}^{u_i}\right) \left(1 - \eta_{f_k}^{b_0}\right) \prod_{b_j \in \Phi_b \setminus b_0} \left(1 - \eta_{f_k}^{b_j}\right) \right] \left(1 - \eta_{f_k}^{m_0}\right). \quad (5)$$

Case 3.4: Now, if the content is not even stored in the MBS cache store, it has to be downloaded from the cloud. This case is termed as a cache miss via both sBS and MBS. In this case, the MBS initiates its FD mode and download the content from the cloud. Therefore, the cache miss probability is calculated from (5) as follows:

$$P_{C_{\mathbb{I}_s}}^u = \left[\left(1 - \eta_{f_k}^{u_0}\right) \prod_{u_i \in \Phi_u \setminus u_0} \left(1 - \eta_{f_k}^{u_i}\right) \left(1 - \eta_{f_k}^{b_0}\right) \prod_{b_j \in \Phi_b \setminus b_0} \left(1 - \eta_{f_k}^{b_j}\right) \right] \left(1 - \eta_{f_k}^{m_0}\right). \quad (6)$$

4) *Case 4 - MBS cache hit*: Recall that *Case 4* is only considered when the tagged user is not under the coverage region of any of the sBSs. First, we consider *Case 4.1* - the requested content is available in the MBS cache (i.e. $\mathbb{I}_m = 1$ and $\mathbb{I}_s = 0$). In this sub-case, we calculate the cache hit probability as follows:

$$P_{M_{\mathbb{I}_m}}^u = \left(1 - \eta_{f_k}^{u_0}\right) \prod_{u_i \in \Phi_u \setminus u_i} \left(1 - \eta_{f_k}^{u_i}\right) \eta_{f_k}^{m_0}. \quad (7)$$

Furthermore, we calculate the total local cache hit probability in this case as follows:

$$P_{\mathbb{I}_m}^u = 1 - \left(1 - \eta_{f_k}^{u_0}\right) \prod_{u_i \in \Phi_u \setminus u_0} \left(1 - \eta_{f_k}^{u_i}\right) \left(1 - \eta_{f_k}^{m_0}\right). \quad (8)$$

Note that we derive the cache miss probability of *Case 4.2* as follows:

$$P_{C_{\mathbb{I}_m}}^u = \left(1 - \eta_{f_k}^{u_0}\right) \prod_{u_i \in \Phi_u \setminus u_0} \left(1 - \eta_{f_k}^{u_i}\right) \left(1 - \eta_{f_k}^{m_0}\right). \quad (9)$$

IV. EDGE CACHING: CACHE HIT RATIO ANALYSIS

We determine CHR, followed by successful transmission probabilities in this section.

A. Cache Hit Ratio

We define CHR as the percentage of the served requests of a requester node from the local nodes. In other words, CHR defines the fraction of the requests that are served locally without reaching the cloud. Let us denote the α portion of the users by the set of \mathcal{U}_0 . Recall that $\rho_{f_k}^{u_0}$ denotes the probability that the tagged user u_0 request content f_k . As such, in a heterogeneous caching placement, we determine the fraction of

requests of u_0 that are served from the local nodes as follows:

$$\text{CHR} = \sum_{k=1}^F \rho_{f_k}^{u_0} \left[\eta_{f_k}^{u_0} + \underbrace{P_d^u P_{s,f}^u + \left(P_{M_{\mathbb{I}_M}}^u P_{s,f,\mathbb{I}_m=1}^{m_0} \right)}_{\text{cache hit in case 4}} \right] \mathbb{I}_m + \underbrace{\left(P_{b_0}^u P_{s,f}^{b_0} + P_B^u P_{s,f}^b + P_{M_{\mathbb{I}_s}}^u P_{s,f}^{m_0} \right)}_{\text{cache hit in case 3}} \mathbb{I}_s, \quad (10)$$

where the first term represents the self cache hit, while the second term represents the successfully achieved cache hit from D2D neighbors. The contents inside (\cdot) in the third term and in the fourth term are the successfully achieved cache hit from *Case 3* and *Case 4*, respectively. Moreover, $P_{s,f}^*$ represents the successful transmission probability for the respective * cases. Note that the transmission success probability between two nodes does not depend on the content index. Therefore, we mention the success probability as $P_{s,f}^*$ instead of P_{s,f_k}^* .

B. Probability of Successful Transmission

Now, we calculate the transmission success probabilities among different nodes. When a tagged user request a content, interference comes from - other active D2D users, active sBSs and MBS. The wireless channel between two nodes follows a Rayleigh fading distribution with $\mathcal{CN}(0,1)$. Let us denote the channel between node i and node j by h_{ij} . Let us also denote the threshold SINR for successful communication by ϕ dB. The transmission power of the user, sBS and MBS are denoted by p_u , p_b and p_m , respectively. Moreover, the path loss exponent is denoted by β .

Now, let γ_i^j , d_i^j and I_{ij} denote the SINR at node i served from node j , distance between the nodes and total interference at node i , respectively. We then derive the SINR values for different cases and sub-cases in equation (12). Owing to the space constraint, the detail derivations of these probabilities are omitted. However, the final tight closed form approximations are provided in equation (14-15). Also, note that we do not consider the case of obtaining the content from the cloud, when we calculate CHR. This is due to the fact that we are interested in calculating the percentage of served request from the local nodes only.

V. CACHE HIT RATIO MAXIMIZATION USING PARTICLE SWARM OPTIMIZATION

We present our objective function, followed by the proposed M-PSO algorithm in this section.

A. CHR Maximization Objective Function

To this end, we calculate the average cache hit ratio for the requesting nodes, which is denoted by Σ . The detail derivation of the Σ is shown in (16). Our objective is to maximize the Σ given that the storage constraints are not violated. Thus, we express the objective function in heterogeneous caching model case as follows:

$$\mathbf{P}_1: \quad \text{maximize } \Sigma \quad (11a)$$

$$\eta_{f_k}^{u_i}, \eta_{f_k}^{b_j}, \eta_{f_k}^{m_l}$$

$$\text{s. t. } \sum_{k=1}^F \eta_{f_k}^{u_i} \leq \mathcal{C}_u, \quad \forall u_i \in \{\mathcal{U}\}, f_k \in \{\mathcal{F}\} \quad (11b)$$

$$\sum_{k=1}^F \eta_{f_k}^{b_j} \leq \mathcal{C}_b, \quad \forall b_j \in \{\mathcal{B}\}, f_k \in \{\mathcal{F}\} \quad (11c)$$

$$\sum_{k=1}^F \eta_{f_k}^{m_l} \leq \mathcal{C}_m, \quad \forall m_l \in \{\mathcal{M}\}, f_k \in \{\mathcal{F}\} \quad (11d)$$

$$0 \leq \eta_{f_k}^{u_i} \leq 1, \quad 0 \leq \eta_{f_k}^{b_j} \leq 1, \quad 0 \leq \eta_{f_k}^{m_l} \leq 1, \quad (11e)$$

where the constraints in (11b-11d) ensure the physical storage size limitations of the user, the sBS and the MBS, respectively, while the constraints in (11e) are due to the probability range in $[0,1]$.

We intend to find optimal caching placements variables that deliver us the optimal solutions. The motivation of \mathbf{P}_1 is to ensure that the requested contents are delivered locally instead of overwhelming the core network during busy traffic hours. However, in general, problem \mathbf{P}_1 is non-convex [17] by nature and may not be solved efficiently in a polynomial-time due to the nonlinear and combinatorial content placement variables. Had we have binary decision parameters, it is not hard to see that the \mathbf{P}_1 would have been reduced to a Knapsack problem, which is widely recognized as an NP-complete problem. Nevertheless, each of our decision variables is a probability that is in $[0,1]$. There is an infinite number of possible values, to determine the optimal solution from, in this range. Therefore, the use of typical metaheuristic solutions such as genetic algorithms may not be a suitable choice. Thanks to particle swarm optimization (PSO) technique, we can leverage its fundamental concept to get to a modified version of it that is suitable for a complex combinatorial problem such as \mathbf{P}_1 . We discuss our proposed modified PSO (M-PSO) framework in what follows.

B. Modified-Particle Swarm Optimization Algorithm

PSO is a swarm intelligence approach that guarantees to converge [19]. In this meta-heuristic algorithm, all possible sets of candidate solutions are named as the particles - denoted by i . Each particle has a position - denoted by x_i . Furthermore, it maintains a personal best position of each particle and the global best positions of the entire swarm. These two terms are denoted by p_i^{best} and g^{best} , respectively. The algorithm evolves, with an exploration and exploitation manner, by adding a velocity term - v_i^t at each particle's previous position aiming to converge at the global optima. The following two simple equations, thus, govern the PSO algorithm.

$$v_i^{t+1} = av_i^t + \psi_1 \epsilon_1 (p_i^{best} - x_i) + \psi_2 \epsilon_2 (g^{best} - x_i), \quad (17)$$

$$x_i^{t+1} = x_i^t + v_i^t, \quad (18)$$

where a , ψ_1 and ψ_2 are the parameters that needs to be selected properly. Moreover, ϵ_1 and ϵ_2 are two Unifrom random variables. Note that ψ_1 and ψ_2 are positive acceleration coefficients, which are also known as the cognitive and social learning factors [17], respectively. While this is a general framework for the PSO algorithm, it may not be used directly in constraint optimization [20]. In our objective function, each particle must have a position matrix - each dimension of which must not violate the restrictions. Therefore, in the following,

$$\gamma_{u_0}^u = \frac{p_u h_{u_0 u_i} d_{u_0 u_i}^{-\beta}}{\sigma^2 + I_{u_0 u}}, \quad \gamma_{u_0}^{b_0} = \frac{p_b h_{u_0 b_0} d_{u_0 b_0}^{-\beta}}{\sigma^2 + I_{u_0 b_0}}, \quad \gamma_{b_0}^b = \frac{p_b h_{b_0 b_i} d_{b_0 b_i}^{-\beta}}{\sigma^2 + I_{b_0 b}}, \quad \gamma_{b_0}^{m_0} = \frac{p_m h_{b_0 m_0} d_{b_0 m_0}^{-\beta}}{\sigma^2 + I_{b_0 m_0}}, \quad \gamma_{u_0}^{m_0} = \frac{p_m h_{u_0, m_0} d_{u_0, m_0}^{-\beta}}{\sigma^2 + I_{u_0, m_0}}, \quad (12a)$$

where the interference, I_{ij} s, are calculated as follows:

$$I_{u_0 u} = \sum_{u \in \Phi_u \setminus \{u_0, u_i\}} p_u h_{u_0 u} d_{u_0 u}^{-\beta} + \sum_{b_l \in \{\mathcal{B}\}_{b_0}^B} p_{b_l} h_{u_0 b_l} d_{u_0 b_l}^{-\beta} \mathbb{I}_{b_l} + \sum_{b_l \in \{\mathcal{B}\}_{b_0}^B, u_i \in \Phi_u \setminus \{u_0, u_i\}} p_m h_{u_0 m} d_{u_0 m}^{-\beta} \mathbb{I}_{m_0} \quad (13a)$$

$$I_{u_0 b_0} = \sum_{u \in \Phi_u \setminus u_0} p_u h_{u_0 u} d_{u_0 u}^{-\beta} + \sum_{b_l \in \{\mathcal{B}\}_{b_1}^B \setminus b_0} p_{b_l} h_{u_0 b_l} d_{u_0 b_l}^{-\beta} \mathbb{I}_{b_l} + \sum_{b_l \in \{\mathcal{B}\}_{b_1}^B \setminus b_0, u_i \in \Phi_u \setminus u_0} p_m h_{u_0 m} d_{u_0 m}^{-\beta} \mathbb{I}_{m_0} \quad (13b)$$

$$I_{b_0 b} = \sum_{u \in \Phi_u \setminus u_0} p_u h_{b_0 u} d_{b_0 u}^{-\beta} + \sum_{b \in \{\mathcal{B}\} \setminus \{b_0, b_i\}} p_b h_{b_0 b} d_{b_0 b}^{-\beta} + h_{b_0 b_0} \zeta p_b + \sum_{b_l \in \{\mathcal{B}\}_{b_1}^B \setminus \{b_0, b_i\}, u \in \Phi_u \setminus u_0} p_m h_{u_0 m} d_{u_0 m}^{-\beta} \mathbb{I}_{m_0} \quad (13c)$$

$$I_{b_0 m_0} = \sum_{u \in \Phi_u \setminus u_0} p_u h_{b_0 u} d_{b_0 u}^{-\beta} + \sum_{b \in \{\mathcal{B}\} \setminus b_0} p_b h_{b_0 b} d_{b_0 b}^{-\beta} + h_{b_0 b_0} \zeta p_b + \sum_{b_l \in \{\mathcal{B}\}_{b_1}^B \setminus b_0, u \in \Phi_u \setminus u_0} p_m h_{u_0 m} d_{u_0 m}^{-\beta} \mathbb{I}_{m_0} \quad (13d)$$

$$I_{u_0, m_0} = \sum_{u \in \Phi_u \setminus u_0} p_u h_{u_0 u} d_{u_0 u}^{-\beta} + \sum_{b_l \in \{\mathcal{B}\}, u_i \in \Phi_u \setminus u_0} p_m h_{u_0 m} d_{u_0 m}^{-\beta} \mathbb{I}_{m_0}. \quad (13e)$$

$$\mathbf{P}_{s,f}^u = \frac{A}{B} [1 - \exp(-\pi R_u^2 B)], \quad \mathbf{P}_{s,f}^{b_0} = \frac{A_1}{B_1} [1 - \exp(-\pi R_b^2 B_1)], \quad \mathbf{P}_{s,f, \mathbb{I}_m=1}^{m_0} = \frac{A_2}{B_2} [1 - \exp(-\pi R_m^2 B_2)], \quad (14a)$$

$$\mathbf{P}_{s,f}^b = \left[\int_{r>0} \left\{ f_{d_1}(r) \exp\left(\frac{-\pi \alpha \lambda_u \left(\phi \frac{p_u}{p_b}\right)^{\frac{2}{\beta}} r^2}{\text{sinc}\left(\frac{2}{\beta}\right)}\right) \exp\left(\frac{-\pi \lambda_b \left(\phi \frac{p_b}{p_b}\right)^{\frac{2}{\beta}} r^2}{\text{sinc}\left(\frac{2}{\beta}\right)}\right) \exp(-\phi r^\beta \bar{\zeta}) \exp\left(\frac{-\pi \lambda_m \left(\phi \frac{p_m}{p_b}\right)^{\frac{2}{\beta}} r^2}{\text{sinc}\left(\frac{2}{\beta}\right)}\right) \right\} dr \right] \left\{ \frac{A_1}{B_1} [1 - \exp(-\pi R_b^2 B_1)] \right\}, \quad (14b)$$

$$\mathbf{P}_{s,f}^{m_0} = \left[\int_{r>0} \left\{ f_{d_2}(r) \exp\left(\frac{-\pi \alpha \lambda_u \left(\phi \frac{p_u}{p_m}\right)^{\frac{2}{\beta}} r^2}{\text{sinc}\left(\frac{2}{\beta}\right)}\right) \exp\left(\frac{-\pi \lambda_b \left(\phi \frac{p_b}{p_m}\right)^{\frac{2}{\beta}} r^2}{\text{sinc}\left(\frac{2}{\beta}\right)}\right) \exp(-\phi r^\beta \bar{\zeta}) \exp\left(\frac{-\pi \lambda_m r^2 \phi^{\frac{2}{\beta}}}{\text{sinc}\left(\frac{2}{\beta}\right)}\right) \right\} dr \right] \left\{ \frac{A_1}{B_1} [1 - \exp(-\pi R_b^2 B_1)] \right\}, \quad (14c)$$

where $\bar{\zeta}$ is the self-interference [16] due to FD communication. Moreover, A , B , A_1 , B_1 , A_2 and B_2 are calculated as follows:

$$A = \frac{(1-\alpha)\lambda_u}{1 - \exp[-\pi(1-\alpha)\lambda_u R_u^2]}, \quad B = \lambda_u \left((1-\alpha) + \frac{\alpha \phi^{\frac{2}{\beta}}}{\text{sinc}(2/\beta)} \right) + \frac{\lambda_b \left(\phi \frac{p_b}{p_u}\right)^{\frac{2}{\beta}}}{\text{sinc}\left(\frac{2}{\beta}\right)} + \frac{\lambda_m \left(\phi \frac{p_m}{p_u}\right)^{\frac{2}{\beta}}}{\text{sinc}\left(\frac{2}{\beta}\right)}, \quad (15a)$$

$$A_1 = \frac{\lambda_b}{1 - \exp(\pi \lambda_b R_b^2)}, \quad B_1 = \lambda_b \left[1 + 2\phi^{\frac{2}{\beta}} \int_{\phi^{-\frac{2}{\beta}}}^{\infty} \left(\frac{1}{1+u^{\frac{\beta}{2}}} \right) du \right] + \frac{\alpha \lambda_u \left(\phi \frac{p_u}{p_b}\right)^{\frac{2}{\beta}}}{\text{sinc}\left(\frac{2}{\beta}\right)} + \frac{\lambda_m \left(\phi \frac{p_m}{p_b}\right)^{\frac{2}{\beta}}}{\text{sinc}\left(\frac{2}{\beta}\right)}, \quad (15b)$$

$$A_2 = \frac{\lambda_m}{1 - \exp(\pi \lambda_m R_m^2)}, \quad B_2 = \lambda_m \left[1 + \frac{\phi^{\frac{2}{\beta}}}{\text{sinc}\left(\frac{2}{\beta}\right)} \right] + \frac{\alpha \lambda_u \left(\phi \frac{p_u}{p_m}\right)^{\frac{2}{\beta}}}{\text{sinc}\left(\frac{2}{\beta}\right)}. \quad (15c)$$

$$\Sigma = \frac{1}{|\mathcal{Z}_0|} \sum_{u_0 \in \mathcal{Z}_0} \sum_{k=1}^F \rho_{f_k}^{u_0} \left\{ \eta_{f_k}^{u_0} + (1-\eta_{f_k}^{u_0}) \left[1 - \prod_{u_i \in \Phi_u \setminus u_0} (1-\eta_{f_k}^{u_i}) \right] \mathbf{P}_{s,f}^u + \left((1-\eta_{f_k}^{u_0}) \prod_{u_i \in \Phi_u \setminus u_0} (1-\eta_{f_k}^{u_i}) \eta_{f_k}^{b_0} \mathbf{P}_{s,f}^{b_0} + (1-\eta_{f_k}^{u_0}) \prod_{u_i \in \Phi_u \setminus u_0} (1-\eta_{f_k}^{u_i}) (1-\eta_{f_k}^{b_0}) \left(1 - \prod_{b_j \in \Phi_b \setminus b_0} (1-\eta_{f_k}^{b_j}) \right) \mathbf{P}_{s,f}^b + (1-\eta_{f_k}^{u_0}) \prod_{u_i \in \Phi_u \setminus u_0} (1-\eta_{f_k}^{u_i}) (1-\eta_{f_k}^{b_0}) \prod_{b_j \in \Phi_b \setminus b_0} (1-\eta_{f_k}^{b_j}) \eta_{f_k}^{m_0} \mathbf{P}_{s,f}^{m_0} \right) \mathbb{I}_s + \left((1-\eta_{f_k}^{u_0}) \prod_{u_i \in \Phi_u \setminus u_i} (1-\eta_{f_k}^{u_i}) \eta_{f_k}^{m_0} \mathbf{P}_{s,f, \mathbb{I}_m=1}^{m_0} \right) \mathbb{I}_m \right\}. \quad (16)$$

we modify the PSO algorithm to solve our optimization problem efficiently.

Let P be numbers of particles. Let $\boldsymbol{\eta}_f^{u_i}$ denote the caching probabilities of user u_i for all contents $f_k \in \{\mathcal{F}\}$. Then, this

parameter has a size of $F \times 1$. Similarly, for all sBS and MBS, let $\boldsymbol{\eta}_f^{bj}$ and $\boldsymbol{\eta}_f^{mj}$ denote their caching placement probabilities for all contents. Then, all of these parameters can be stacked into a matrix with dimension of $(|\mathcal{U}| + |\mathcal{B}| + |\mathcal{M}|) \times |\mathcal{F}|$, which is the exact shape of each particle. Let the current position of each of these particles be denoted by \mathbf{X}_i^t . Note that in this case, each particle's position \mathbf{X}_i^t has a shape of $(|\mathcal{U}| + |\mathcal{B}| + |\mathcal{M}|) \times |\mathcal{F}|$. Let $\mathbf{V}_i^t \in \mathbb{R}^{(|\mathcal{U}| + |\mathcal{B}| + |\mathcal{M}|) \times |\mathcal{F}|}$ denote the velocity. Furthermore, the personal best position of particle i is denoted by $\mathbf{P}_i^{\text{best}}$, while the global best for the entire swarm is denoted by \mathbf{G}^{best} . Therefore, each particle updates its velocity with social and individual cognition. We use the following equation to govern these updates.

$$\mathbf{V}_i^{t+1} = a\mathbf{V}_i^t + \psi_1 \left[\mathcal{E}_1 \odot (\mathbf{P}_i^{\text{best}} - \mathbf{X}_i^t) \right] + \psi_2 \left[\mathcal{E}_2 \odot (\mathbf{G}^{\text{best}} - \mathbf{X}_i^t) \right], \quad (19)$$

where a , ψ_1 and ψ_2 are the parameters as described in (17). Moreover, \mathcal{E}_1 and \mathcal{E}_2 are two matrices with sizes of $\mathbb{R}^{(|\mathcal{U}| + |\mathcal{B}| + |\mathcal{M}|) \times |\mathcal{F}|}$, where their element is drawn from Uniform random distribution. Finally, \odot represents Hadamard product.

The position of each particle is then updated by the velocity similar to (17). However, as we have the constraints as in (11b)-(11e), we need to modify this equation accordingly. Let $\mathbf{X}_{i_{\text{int}}}^{t+1}$ denote an intermediate updated position of particle i as shown in the following expression.

$$\mathbf{X}_{i_{\text{int}}}^{t+1} = \mathbf{X}_i^t + \mathbf{V}_i^t. \quad (20)$$

We consider this intermediate position to keep each particle's position in the feasible search space. Besides, we also perform the necessary normalization and scaling. Note that from this intermediate particle position leads to a normalized particle position. This parameter is then used as the current particle position \mathbf{X}_i^t . Moreover, the ultimate goal for each particle is to converge at an optimal position \mathbf{X}_i^* (i.e. the global best \mathbf{G}^{best}).

Algorithm 1 summarizes the steps of the proposed algorithm. Note that our proposed algorithm can be implemented to solve any similar hard combinatorial problems.

We model the algorithm such a way that we deal with the normalized particle position and velocity. The constraints guide us to restrict the particle position in a probability range, while the summation cannot exceed the cache storage capacity of the respective node. Therefore, we consider to limit the initial values in the range of $[0, 1/\mathcal{C}^j]$. By doing so, when we perform the necessary scaling, the obtained number does not violate the probability range. Then, we correspondingly initialize the particle position and the velocity in steps 4 and 5 following this notion. Furthermore, the caching probabilities of the nodes in dimension j are limited to \mathcal{C}^j , in steps 22 and 24, hence, we choose the random number of contents, $\text{randint}(\mathcal{C}^j)$, to be stored with higher probability values. We stress the fact that, although our proposed M-PSO is a modified version of PSO, it inherits all properties of the original PSO algorithm. As such, it is not hard to analyze the convergence and complexities of our proposed algorithm following the analysis of the original PSO algorithm [21].

VI. RESULTS AND DISCUSSION

For the simulation, user are considered to be distributed in a 2D plane following a HPPP of intensity, $\lambda_u \in [10^{-4}, 10^{-3}]$

Algorithm 1 CHR Maximization using M-PSO

```

1: for each particle,  $i = 1, 2, \dots, P$  do
2:    $\mathbf{X}_i = [ ], \mathbf{V}_i = [ ]$ 
3:   for each dimension  $j = 1, 2, \dots, D$  do ▷
4:      $D = |\mathcal{U}| + |\mathcal{B}| + |\mathcal{M}|$ 
5:     initialize the particles positions,  $\mathbf{x}_{ji}$  with uniform
6:     random vector of size  $\mathbb{R}^{|\mathcal{F}|}$  by making sure  $\sum_{k=1}^F x_{ji}[k] = 1$ 
7:     and  $0 \leq x_{ji}[k] \leq \frac{1}{\mathcal{C}^j}, \forall k \in \mathcal{F}$ ; then set  $\mathbf{X}_i[j, :] \leftarrow \mathbf{x}_{ji}$  ▷
8:      $\mathcal{C}^j$  is the cache storage of the node in  $j^{\text{th}}$  dimension
9:     initialize particles velocity,  $\mathbf{v}_{ji}$  with uniform random
10:    vector of size  $\mathbb{R}^{|\mathcal{F}|}$  by making sure  $\sum_{k=1}^F v_{ji}[k] = 1$ 
11:    and  $0 \leq v_{ji}[k] \leq \frac{1}{\mathcal{C}^j}, \forall k \in \mathcal{F}$ ; then set  $\mathbf{V}_i[j, :] \leftarrow \mathbf{v}_{ji}$ 
12:    end for
13:    set particle best position,  $\mathbf{P}_i^{\text{best}}$  as the initial position
14:    if  $\Sigma(\mathbf{P}_i^{\text{best}}) > \Sigma(\mathbf{G}^{\text{best}})$  then
15:       $\mathbf{G}^{\text{best}} \leftarrow \mathbf{P}_i^{\text{best}}$ 
16:    end if
17:  end for
18:  while termination criteria has not met do
19:    for each particle,  $i$  do
20:      for each dimension,  $j = 1, 2, \dots, D$  do
21:        draw uniform random vectors,  $\boldsymbol{\epsilon}_1$  and  $\boldsymbol{\epsilon}_2$  of size
22:         $\mathbb{R}^{|\mathcal{F}|}$ 
23:        set  $\mathbf{v}_{ji} \leftarrow a\mathbf{v}_{ji} + \psi_1 \left[ \boldsymbol{\epsilon}_1 \odot (\mathbf{P}_i^{\text{best}} - \mathbf{x}_{ji}) \right] +$ 
24:         $\psi_2 \left[ \boldsymbol{\epsilon}_2 \odot (\mathbf{G}^{\text{best}} - \mathbf{x}_{ji}) \right]$ 
25:        set  $\mathbf{V}_i[j, :] \leftarrow \mathbf{v}_{ji}$ 
26:      end for
27:      update particles intermediate position,  $\mathbf{X}_{i_{\text{int}}}$ 
28:       $\mathbf{X}_i^{\text{scl}} = [ ], \mathbf{P}_i^{\text{scl\_best}} = [ ], \mathbf{G}^{\text{scl\_best}} = [ ]$ 
29:      for each dimension  $j = 1, 2, \dots, D$  do
30:        random_hike  $\leftarrow \text{randint}(\mathcal{C}^j)$ 
31:        for  $i$  in  $\text{len}(\text{random\_hike})$  do
32:           $\mathbf{X}_{i_{\text{int}}}[j, \text{randint}(F)] \leftarrow \frac{\sum_{k=1}^F \mathbf{X}_{i_{\text{int}}}[j, :]}{\mathcal{C}^j}$ 
33:        end for
34:         $\mathbf{X}_i[j, :] \leftarrow \frac{\mathbf{X}_{i_{\text{int}}}[j, :]}{\sum_{k=1}^F \mathbf{X}_{i_{\text{int}}}[j, :]}; \mathbf{X}_i^{\text{scl}}[j, :] \leftarrow \mathcal{C}^j \mathbf{X}_i[j, :]$  ▷
35:        Normalized particle position
36:         $\mathbf{P}_i^{\text{scl\_best}}[j, :] \leftarrow \mathcal{C}^j \mathbf{P}_i^{\text{best}}[j, :]$ 
37:         $\mathbf{G}_i^{\text{scl\_best}}[j, :] \leftarrow \mathcal{C}^j \mathbf{G}_i^{\text{best}}[j, :]$ 
38:      end for
39:      if  $\Sigma(\mathbf{X}_i^{\text{scl}}) > \Sigma(\mathbf{P}_i^{\text{scl\_best}})$  then
40:         $\mathbf{P}_i^{\text{best}} \leftarrow \mathbf{X}_i$ 
41:        do necessary scaling following step 27
42:        if  $\Sigma(\mathbf{P}_i^{\text{scl\_best}}) > \Sigma(\mathbf{G}^{\text{scl\_best}})$  then
43:           $\mathbf{G}^{\text{best}} \leftarrow \mathbf{P}_i^{\text{best}}$ 
44:        end if
45:      end if
46:    end for
47:  end while
48:  return  $\mathbf{G}^{\text{best}}$  and do necessary scaling following step 28
49:  and return  $\mathbf{G}^{\text{scl\_best}}$ 

```

(per m^2). The low powered sBS are drawn following another HPPP of intensity, $\lambda_b \in [10^{-6}, 10^{-5}]$, (per m^2). For the MBS, $\lambda_m = 1.5^{-7}$ (per m^2) is considered. The coverage radii of the user, sBS and MBS are taken as $R_u = 15 m$, $R_b = 150 m$, $R_m = 500 m$, respectively. Total contents in the catalog $|\mathcal{F}| = [10, 50]$, $\alpha \in [0.2, 0.5]$ and the skewness, γ of

the Zipf distribution is considered to be selected uniformly in between $\{0.1, 2.5\}$. For the M-PSO algorithm, we set $a = 0.9$ and $\psi_1 = \psi_2 = 0.4$. Moreover, $p_u = 23$ dBm, $p_b = 26$ dBm, $p_m = 43$ dBm, $\phi = 10^{-8}$ dB, $\beta = 4$, $\zeta = 0.01$ and $\sigma^2 = -174$ dBm/Hz are considered. We apply Monte Carlo simulation methods while performing our evaluation. In the following, we use the proposed M-PSO algorithm to attain the optimal caching placement solution. After that, we study its performances for our hard-combinatorial maximization problem.

A. Cache Placement

To show the effectiveness of the proposed algorithm, we first validate that the obtained results do not violate any of the constraints. The obtained global best $\mathbf{G}^{\text{scl_best}}$, using Algorithm 1, is therefore scrutinized as follows. Note that it must not violate any of the caching storage constraints of the edge nodes. Besides, each of the caching probabilities must be in the range of $[0, 1]$. Furthermore, each node must store different copies of the content. Notice that we have applied all of these constraints in our proposed algorithm. Therefore, it is expected that the obtained results will satisfy these constraints. The caching probabilities of the 1st and 2nd respective D2D users, sBSs and MBSs are illustrated in Fig. 1. Notice that each node stores different copies. Moreover, caching probabilities and storage constraints are also satisfied.

B. Performance Analysis

We study the performance of our proposed M-PSO algorithm and make a fair comparison to the following benchmark caching schemes in this sub-section.

Random Caching Scheme: In the random caching scheme, contents are stored randomly, while satisfying the constraints.

Equal Caching Scheme: In the equal caching scheme, each content is placed with the same probability.

To show the effectiveness of our proposed algorithm, we only consider 100 iterations. In Fig. 2, we demonstrate the results obtained from using our proposed algorithm, random caching scheme, and equal caching scheme. With only 100 iterations, we achieve $\approx 24\%$ better performance than the benchmark caching schemes. Therefore, we claim that our proposed algorithm achieves better system performance than the other baseline caching schemes within a minimal number of iterations. In the following, we use our algorithm to evaluate the system performance in terms of different parameter setting.

1) *Impact of the Catalog Size:* Recall that if the requested content is delivered from one of the cache-enabled edge nodes, a cache hit occurs. Therefore, we aim to store as many to-be-requested contents as possible into the local edge nodes. We consider the catalog size in the set of $[10, 20, 30, 40, 50]$. Furthermore, the intensities are set as $\lambda_u = 10^{-4}$, $\lambda_b = 10^{-5}$ and $\lambda_m = 10^{-7}$. Also, the total number of iterations is chosen in the set of $100 \times [1, 10, 20, 40, 80]$ for the catalog size in $[10, 20, 30, 40, 50]$, respectively. Note that if the catalog size increases, the number of possible combinations also increases. Therefore, whenever the content catalog increases, we slightly increase the total number of iterations. Also, if the total number of contents increases and we have only a limited number of cache-enabled nodes, the chance of storing the contents locally decreases, meaning that more content requests

need to be served from the cloud. Therefore, the Σ should decrease if the content catalog increases. Moreover, if the percentage of the requester nodes increases, the performance should degrade as we consider the heterogeneous preference of the users. Fig. 3 also shows that if we increase the catalog size, $|\mathcal{F}|$ or the number of requesters (α), then the Σ decreases.

2) *Impact of the Storage Size:* We now investigate the impact of the cache sizes of the edge nodes on the system performance. Remember that if cache size increases, more content can be stored at the cache-enabled nodes. Therefore, increasing the cache size of the users means that users store more contents at their local storage. As these storage sizes increases, the job of the proposed M-PSO algorithm is to determine the optimal caching placements. The simulation results, presented in Fig. 4, validate that as the storage size increases more content is locally stored leading to an improved CHR. Notice that increasing MBS cache size provides lesser CHR gain than increasing the cache size of the D2D users (or, the sBS). This is because the total number of MBS are typically very lower than the available D2D (or, sBS) nodes.

VII. CONCLUSION

Caching solution helps to achieve better system performances. However, the hard combinatorial decision-making problem of placing the contents at the local nodes is challenging. The grand problem is effectively solved with good accuracy by using the artificial intelligence based technique. Considering heterogeneous content preferences in a real-world network platform, the proposed algorithm converges fast and achieves a much better performance than the existing benchmark caching schemes.

ACKNOWLEDGMENT

The authors sincerely thank Shaju Shah for the critical and helpful discussions during this work.

REFERENCES

- [1] M. F. Pervej, L. T. Tan, and R. Q. Hu, "Artificial intelligence assisted collaborative edge caching in modern small cell networks," in *Proc. IEEE Globecom*, Dec. 2020.
- [2] M. F. Pervej and S.-C. Lin, "Dynamic power allocation and virtual cell formation for Throughput-Optimal vehicular edge networks in highway transportation," in *Proc. IEEE ICC Workshops*, June 2020.
- [3] M. F. Pervej and S.-C. Lin, "Eco-Vehicular edge networks for connected transportation: A distributed multi-agent reinforcement learning approach," in *Proc. IEEE VTC2020-Fall*, Oct. 2020.
- [4] M. F. Pervej, L. T. Tan, and R. Q. Hu, "User preference learning aided collaborative edge caching for small cell networks," in *Proc. IEEE Globecom*, Dec. 2020.
- [5] 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, July 2014.
- [6] A. F. Molisch, G. Caire, D. Ott, J. R. Foerster, D. Bethanabhotla, and M. Ji, "Caching eliminates the wireless bottleneck in video aware wireless networks," *Advances in Electrical Engineering*, vol. 2014, 2014.
- [7] V. A. Siris and D. Dimopoulos, "Multi-source mobile video streaming with proactive caching and d2d communication," in *Proc. WoWMoM*. IEEE, 2015.
- [8] Y. Hao, L. Hu, Y. Qian, and M. Chen, "Profit maximization for video caching and processing in edge cloud," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 7, pp. 1632–1641, May 2019.
- [9] J. Du, C. Jiang, E. Gelenbe, H. Zhang, Y. Ren, and T. Q. S. Quek, "Double auction mechanism design for video caching in heterogeneous ultra-dense networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 3, pp. 1669–1683, Feb. 2019.

Caching Probabilities at Different Nodes

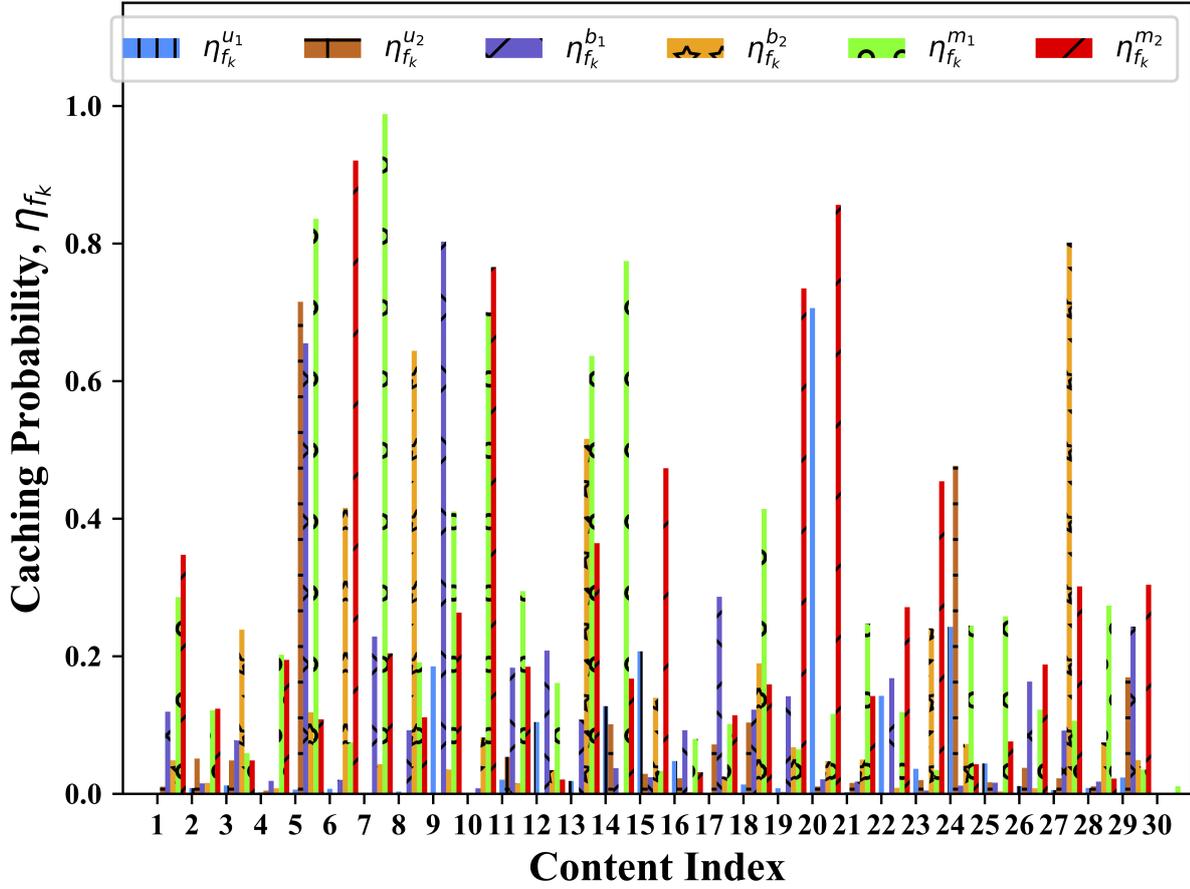


Fig. 1. Obtained caching probabilities at the local nodes when $\mathcal{C}_d = 2$, $\mathcal{C}_b = 4$ and $\mathcal{C}_m = 8$

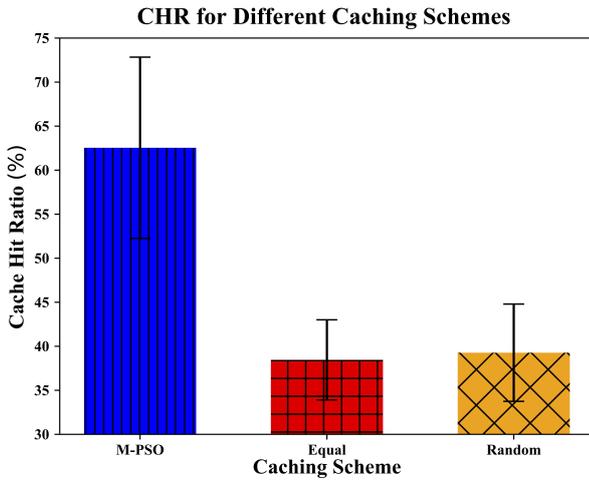


Fig. 2. CHR using the proposed M-PSO algorithms for 100 iteration, $|\mathcal{F}| = 30$, $\mathcal{C}_d = 2$, $\mathcal{C}_b = 4$ and $\mathcal{C}_m = 8$

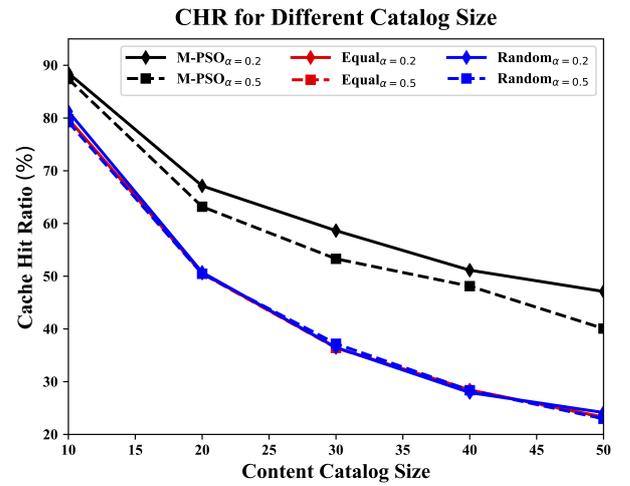


Fig. 3. Impact of catalog size: CHR with $\mathcal{C}_d = 2$, $\mathcal{C}_b = 4$ and $\mathcal{C}_m = 8$

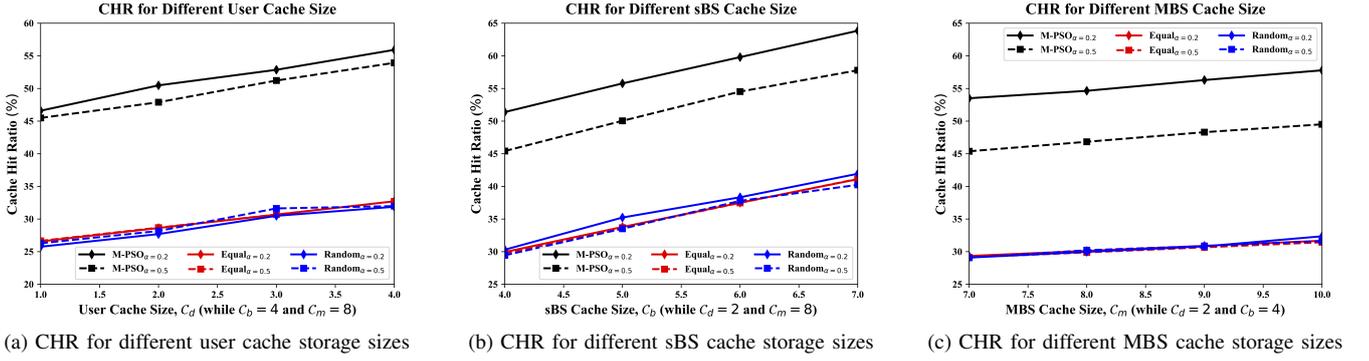


Fig. 4. Impact of cache size on CHR

- [10] J. Xu, J. Liu, B. Li, and X. Jia, "Caching and prefetching for web content distribution," *Computing in science & engineering*, vol. 6, no. 4, pp. 54–59, 2004.
- [11] M. Sheng, C. Xu, J. Liu, J. Song, X. Ma, and J. Li, "Enhancement for content delivery with proximity communications in caching enabled wireless networks: Architecture and challenges," *IEEE Commun. Mag.*, vol. 54, no. 8, pp. 70–76, 2016.
- [12] X. Wang, M. Chen, T. Taleb, A. Ksentini, and V. C. Leung, "Cache in the air: Exploiting content caching and delivery techniques for 5g systems," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 131–139, 2014.
- [13] D. Liu, B. Chen, C. Yang, and A. F. Molisch, "Caching at the wireless edge: design aspects, challenges, and future directions," *IEEE Commun. Mag.*, vol. 54, no. 9, pp. 22–28, 2016.
- [14] M. Lee, H. Feng, and A. F. Molisch, "Dynamic caching content replacement in base station assisted wireless d2d caching networks," *IEEE Access*, vol. 8, pp. 33 909–33 925, Feb 2020.
- [15] M. Lee and A. F. Molisch, "Caching policy and cooperation distance design for base station-assisted wireless d2d caching networks: Throughput and energy efficiency optimization and tradeoff," *IEEE Trans. Wireless Commun.*, vol. 17, no. 11, pp. 7500–7514, Nov 2018.
- [16] L. T. Tan, R. Q. Hu, and Y. Qian, "D2d communications in heterogeneous networks with full-duplex relays and edge caching," *IEEE Trans. Ind. Informat.*, vol. 14, no. 10, pp. 4557–4567, 2018.
- [17] L. T. Tan, R. Q. Hu, and L. Hanzo, "Heterogeneous networks relying on full-duplex relays and mobility-aware probabilistic caching," *IEEE Trans. Commun.*, pp. 1–1, 2019.
- [18] B. Blaszczyszyn and A. Giovanidis, "Optimal geographic caching in cellular networks," in *Proc. ICC*, June 2015.
- [19] J. Kennedy, "Particle swarm optimization," *Encyclopedia of machine learning*, pp. 760–766, 2010.
- [20] X. Hu, R. C. Eberhart, and Y. Shi, "Engineering optimization with particle swarm," in *Proc. IEEE Swarm Intell. Symp.* IEEE, 2003.
- [21] M. Clerc and J. Kennedy, "The particle swarm - explosion, stability, and convergence in a multidimensional complex space," *IEEE Trans. Evolution. Comput.*, vol. 6, no. 1, pp. 58–73, 2002.