

# Adaptive Delivery in Caching Networks

Seyed Ali Saberali, Hamidreza Ebrahimzadeh Saffar,  
Lutz Lampe, *Senior Member, IEEE*, and Ian Blake, *Fellow, IEEE*

**Abstract**—The problem of content delivery in caching networks is investigated for scenarios where multiple users request identical files. An adaptive method is proposed for the delivery of redundant demands in caching networks. Based on the redundancy pattern in the current demand vector, the proposed method decides between the transmission of uncoded messages or the coded messages of [1] for delivery. Moreover, a lower bound on the delivery rate of redundant requests is derived. The performance of the adaptive method is investigated through numerical examples and Monte Carlo simulations. It is shown that the adaptive method considerably reduces the performance gap to the lower bound for specific ranges of network parameters.

**Index Terms**—Adaptive delivery algorithm, average delivery rate, coded caching, correlated requests, redundant demands.

## I. INTRODUCTION

LOCAL caching is a promising technique to meet the unprecedented traffic demands in the next generation communication networks [1]–[6]. Caching networks operate in two phases, commonly referred to as placement and delivery phases. In the placement phase, caches fill their memories with parts of the popular files during the off-peak hours. The delivery phase, however, is performed when the network is congested. In this phase, each cache provides its users with the parts of the requested files that it has available. The remaining parts are conventionally delivered to the users through unicast transmissions performed by a central server on a broadcast channel. In a more recent approach, known as coded caching [2], the central server uses coded-multicasting to deliver the requested content, to further reduce the network congestion.

In [2], the authors derived an information-theoretic formulation for the caching problem and proposed a centralized scheme for coded caching. In a later work [1], a decentralized coded caching method was proposed which became the building block of several caching schemes developed for more complicated scenarios [6]–[9]. [1], [2] used the peak delivery rate as the figure of merit of the network. The peak rate occurs when all the users request distinct files, given that the number of files is greater than the number of caches.

Average delivery rate is another significant network performance metric, and depends on the statistics of the user requests. [5]–[8] proposed caching schemes to decrease the average delivery rate in scenarios with non-uniform file popularities. The statistics of user requests also affect the design of caching networks by increasing the chance of multiple users requesting identical files. In such a scenario, the delivery method can be modified to benefit from the redundancies in the user demands, and further reduce the delivery rate.

Redundant demands are likely when either the files have significantly different popularity levels or the user requests

are positively correlated. For the case of non-uniform file popularities, the schemes in [6]–[8] do not take the effect of redundant requests into account. This is because the delivery in all these schemes is based on the delivery of [1], which is designed for the demand vectors with distinct requests. In addition to non-uniform popularity levels, correlated requests are also likely in many practical scenarios. A considerable amount of multimedia requests are made through social networks like Facebook and Instagram, where users with common friends and interests are likely to request the same content.

In this paper, we investigate the delivery of redundant demands in caching networks. We use the placement schemes of [1], [2] to ensure that the peak delivery rate does not exceed the delivery rates of [1], [2], and the link capacity is satisfied. Further, these placement schemes are natural candidates when the file popularities are uniform or little prior knowledge about the popularities is available during the placement phase.

For the delivery phase, we propose an adaptive scheme based on *message selection*. Upon receiving a demand vector from the users and based on the redundancy pattern of the requests made, the server decides whether to use uncoded messages or the coded messages of [1] to deliver each part of the requested files. This distinguishes our work from [1], [2], as our proposed delivery takes the specifics of the current demand vector into account to decide on the form of the server messages. However, [1], [2] use a fixed structure to compose the server messages for all demand vectors. We show the superiority of our adaptive method through numerical examples and Monte Carlo simulations. We also derive a lower bound on the delivery rate of redundant requests. In some cases, the adaptive method shrinks the gap between the average rate of the non-adaptive scheme and the lower bound by 50%.

The remainder of this paper is organized as follows. In Section II, we present the network model. The adaptive delivery scheme is derived in Section III. Section IV presents numerical examples and simulation results.

## II. NETWORK MODEL

Assume a network with a central server and  $K$  caches. The server is able to communicate with the caches through a broadcast link. We denote the set of all caches in the network by  $\mathcal{K}$ . A library of  $N \geq K$  files is given, where each file is  $F$  bits long. All files are available at the central server. Each cache has a memory capacity of  $M \times F$  bits. We define  $q \triangleq \frac{M}{N}$ .

*Placement Phase:* Placement takes place only once and remains unchanged during the delivery phase. After the placement, the distribution of bits in the caches can be described as follows. For a given file  $n$  and a given subset of caches  $\mathcal{S} \subset \mathcal{K}$ , denote by  $V_{\mathcal{S}}^n$  the subset of bits of file  $n$  that are exclusively stored at the caches in  $\mathcal{S}$ . The resulting subsets of bits *partition* the set of all the bits of every file into  $2^K$

The authors are with the Department of Electrical and Computer Engineering, University of British Columbia, Vancouver, Canada.



## B. Simplified Adaptive Delivery

A simplified version of the message selection step can be formulated by only taking the number of distinct requests  $L$  into account, and ignoring the redundancy pattern of the demand vector. Then, because of the symmetry, we set  $y_S^n = y_s$  for all  $n$  and all  $S : |S| = s$ . This leads to

$$\begin{aligned} & \underset{y_s}{\text{minimize}} && Ly_0 + \sum_{s=1}^{K-1} \binom{K}{s+1} y_s \\ & \text{subject to} && \sum_{s=0}^K \binom{K}{s} y_s = 1 \\ & && 0 \leq y_s \leq x_s, \quad s = 1, \dots, K \\ & && 0 \leq y_0 \leq 1 \end{aligned} \quad (6)$$

as the simplified message selection problem.

*Proposition 1:* Let  $\hat{s} = \lfloor \frac{K-L}{L+1} \rfloor$ . Optimal parameters for the simplified message selection problem of (6) are given by

$$y_s^* = \begin{cases} \sum_{i=1, \dots, \hat{s}} \binom{K}{i} x_i, & s = 0 \\ 0, & s = 1, \dots, \hat{s} \\ x_s, & s = \hat{s} + 1, \dots, K \end{cases}. \quad (7)$$

*Proof:* If we transfer bits from the subsets  $V_S^n : |S| = s$  to  $V_\emptyset^n$ , the resulting change in the rate will be  $L \binom{K}{s} x_s - \binom{K}{s+1} x_s$ . We transfer the bits only if this difference is negative. This is the case when  $s \leq \hat{s}$ . This results in the parameters of (7). ■

Algorithm 2 shows the simplified adaptive delivery scheme.

### Algorithm 2 Simplified Adaptive Delivery Algorithm

**Require:**  $\{V_S^n\}_{n,S}$  // From the placement phase

- 1: **Procedure** SimplifiedAdaptiveDelivery( $d_1, \dots, d_K$ )
- 2:  $L = \text{size}(\text{unique}(d_1, \dots, d_K))$  // Number of distinct requests
- 3:  $\hat{s} \leftarrow \lfloor \frac{K-L}{L+1} \rfloor$
- 4: **for**  $d_k \in \mathcal{D}$  **do**
- 5:  $\hat{V}_\emptyset^{d_k} \leftarrow \cup_{S: s \leq \hat{s}} V_S^{d_k}$  // Corresponds to the first rule of (7)
- 6: **for**  $\mathcal{S} \subset \mathcal{K} : |\mathcal{S}| > 0$  **do**
- 7: **if**  $|\mathcal{S}| \leq \hat{s}$  **then**
- 8:  $\hat{V}_\emptyset^{d_k} \leftarrow \emptyset$  // Corresponds to the second rule of (7)
- 9: **else**
- 10:  $\hat{V}_\mathcal{S}^{d_k} \leftarrow V_\mathcal{S}^{d_k}$  // Corresponds to the third rule of (7)
- 11: **end if**
- 12: **end for**
- 13: **end for**
- 14: Use [1, Algorithm 1] with  $\{\hat{V}_\mathcal{S}^n\}_{n,S}$  instead of  $\{V_S^n\}_{n,S}$

To recap, the proposed adaptive delivery algorithms adjust their use of coded and uncoded messages based on the redundancies in the demand vector. This is in contrast to the delivery of [1], [2], where the server always uses the same structure for construction of the coded messages, regardless of the redundancy pattern in the demand vector.

## C. Lower Bound

Let  $R_L^*(M)$  denote the minimum rate that is achievable for every possible demand vector with  $L$  distinct requests. Proposition 2 gives a lower bound on  $R_L^*(M)$ .

*Proposition 2 (Cutset Bound):* Assume that  $K$  caches request  $L \leq K$  distinct files. Then,  $R_L^*(M)$  must satisfy

$$R_L^*(M) \geq \max_{s \in \{1, \dots, L\}} \left( s - \frac{s}{\lfloor N/s \rfloor} M \right). \quad (8)$$

*Proof:* We modify the cutset bound argument of [2, Sec. VI] to bound  $R_L^*(M)$ . Let  $\mathcal{S}$  be a subset of caches with  $|\mathcal{S}| = s$ , such that there are no two caches in  $\mathcal{S}$  with identical user

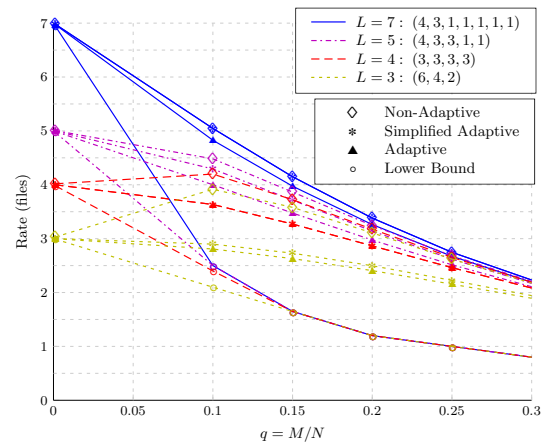


Fig. 1. Rates of different delivery schemes for  $K = 12$ .

requests. Assume that these caches request files  $1, \dots, s$  from the library. Let  $X_1$  denote the server's input to the shared link which determines files  $1, \dots, s$ . Similarly, assume that the same users request files  $(i-1)s+1, \dots, is$ , and the server input  $X_i$  determines the files requested. Let  $i = 1, \dots, \lfloor N/s \rfloor$ . Consider the cut separating  $X_1, \dots, X_{\lfloor N/s \rfloor}$  and the caches in  $\mathcal{S}$  from the corresponding users. The total information available to the users in the cut should be more than or equal to the total information requested. Thus,  $\lfloor N/s \rfloor R_L^*(M) + sM \geq s \lfloor N/s \rfloor$ . Since  $s$  accepts any value between 1 and  $L$ , (8) results. ■

## IV. NUMERICAL EXAMPLES AND SIMULATION RESULTS

We now compare the performance of the proposed adaptive methods and the non-adaptive method of [1, Algorithm 1] through numerical examples. Notice that by the rate of non-adaptive method, we refer to the rate of [1] or [2] depending on whether the decentralized or centralized placement is used, respectively. This rate is calculated by (3).

Fig. 1 shows the delivery rates of the non-adaptive and adaptive schemes, as well as the lower bound (8) for a network of  $K = 12$  caches. The same decentralized placement is used for all cases with the parameters in (2). We consider several redundancy patterns for the demand vector with different values of  $L$ . In Fig. 1, we observe a considerable improvement in the delivery rate for  $M/N \leq 0.25$ , when the adaptive methods are used. This improvement in the rate is more considerable when  $L$  is smaller. For instance, the performance gap to the lower bound decreases by almost 50% when  $L = 3$ . Notice that for a symmetric redundancy patterns like  $(3, 3, 3, 3)$ , both adaptive methods lead to the same delivery rate. As the pattern gets more asymmetric, the gap between the rates of the original and simplified adaptive methods increases. Also, observe that for some cases, the rate of the non-adaptive method increases with the storage capacity for small  $M/N$ . This shows the inefficiency of [1, Algorithm 1] to deliver redundant requests.

For the second numerical example, we use the centralized placement and plot the delivery rates resulted from the different methods versus  $L$ . Fig. 2 shows the results. Notice that the rate of original adaptive method depends not only on  $L$ , but also on the redundancy pattern. Hence, in this example, for every value of  $L$ , the delivery rate of the original adaptive method is averaged over all the redundancy patterns with  $L$  distinct requests, assuming that the requests are independent

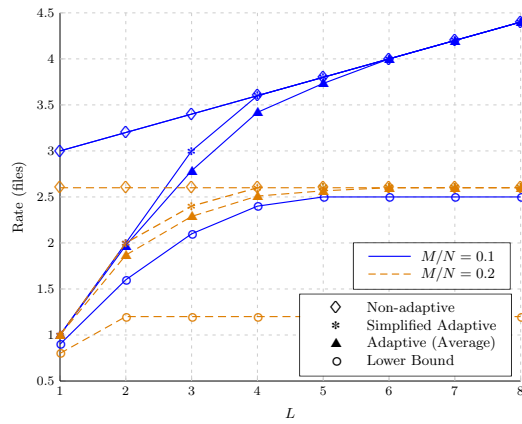


Fig. 2. Rates of different delivery schemes for  $K = 8$ .

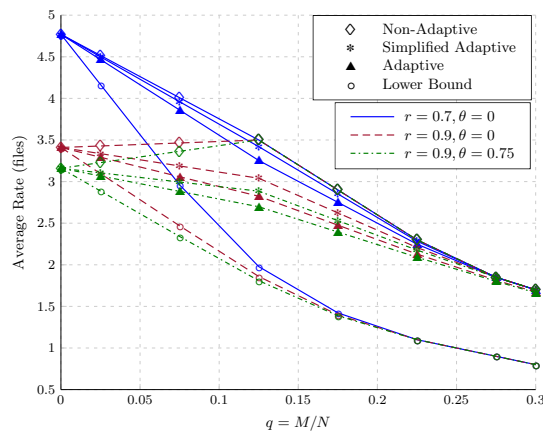


Fig. 3. Average rates of the different delivery schemes for  $K = 8$ .

and file popularities are uniform. Observe that the superiority of the adaptive method over the non-adaptive method of [2] is more significant for small  $L$ . In particular, we observe a sharp decrease in the adaptive delivery rate when  $L$  gets smaller than  $K/2$ . This suggests that the adaptive methods are considerably more effective for highly redundant requests.

We now investigate the average rates of the different delivery methods through a stochastic modelling of the dynamics of a caching network with correlated user requests. Consider a graph representation of the network where vertices represent the caches. An (undirected) edge between two vertices shows that the requests of the corresponding caches are correlated. To model the correlations, let  $\mathcal{N}(k)$  denote the set of the last files requested by the neighbour caches of cache  $k$ . We assume that cache  $k$  requests a file, either based on its neighbours' previous requests with probability  $r$ , or independently with probability  $1-r$ . In the former case,  $k$  chooses a file from  $\mathcal{N}(k)$  uniformly at random. However, when choosing independently,  $k$  picks a file  $n$  from the library based on the popularity distribution of files  $p_n$ . Hence, the chance of requesting file  $n$  by cache  $k$  is

$$\hat{p}_{n,k} = \begin{cases} r/|\mathcal{N}(k)| + (1-r)p_n, & n \in \mathcal{N}(k) \\ (1-r)p_n, & \text{otherwise} \end{cases} \quad (9)$$

For our simulations, we use Gibbs sampling [11, Sec. 24.2] to generate  $10^4$  sample vectors from the induced joint distribution of user demands. We set  $K = 8$  and  $N = 10^3$ ,

$(r, \theta)$	(0.7, 0)	(0.9, 0)	(0.9, 0.75)
Maximum $\hat{p}_{ij}$	0.19	0.34	0.34
Average $\hat{p}_{ij}$	0.16	0.32	0.31
Average $L$	4.80	3.41	3.18

TABLE I. User requests' statistics in simulations of Fig. 3.

and assume a complete graph for the network. Further, we mainly use uniform distribution for the popularity of files. We also consider a scenario where the placement phase is performed based on a uniform popularity distribution, while the actual file popularities in the delivery phase follow a non-uniform Zipf distribution with parameter  $\theta$ . Note that a Zipf distribution with  $\theta = 0$  is identical to the uniform distribution, and increasing  $\theta$  makes the distribution more non-uniform. We use  $\theta$  and  $r$  to control the popularity distribution and the dependency level of the users' requests, respectively. To characterize the resulting correlation levels among the caches' requests in our simulations, we empirically calculate the correlation coefficients  $-1 \leq \hat{p}_{ij} \leq 1$  [12, Section 4.1] between the requests of the different caches  $i$  and  $j$ . A larger  $\hat{p}_{ij}$  implies a higher chance that caches  $i$  and  $j$  request the same content, which leads to more redundancy in the demand vector. Table I presents the average and the maximum  $\hat{p}_{ij}$  over all the different  $i$  and  $j$  pairs ( $i \neq j$ ), in our simulations.

Fig. 3 shows the resulting average delivery rates. It also shows a lower bound on the average rate calculated by averaging the lower bounds of (8) for the sample demand vectors. We observe that as requests become more correlated (larger  $r$ ) and the file popularities get more non-uniform (larger  $\theta$ ), the adaptive method makes a larger improvement in the rate. Also, observe that the adaptive schemes are effective in decreasing the average delivery rate for  $M/N < 0.25$ . The improvement in the performance gap to the lower bound can be as large as 50% for specific choices of parameters.

## REFERENCES

- [1] M. A. Maddah-Ali and U. Niesen, "Decentralized coded caching attains order-optimal memory-rate tradeoff," *IEEE/ACM Trans. Networking*, vol. 23, pp. 1029–1040, Aug. 2015.
- [2] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Trans. Inf. Theory*, vol. 60, pp. 2856–2867, May 2014.
- [3] E. Bastug, M. Bennis, and M. Debbah, "Living on the edge: The role of proactive caching in 5G wireless networks," *IEEE Commun. Mag.*, vol. 52, pp. 82–89, Aug. 2014.
- [4] K. Shanmugam, N. Golrezaei, A. G. Dimakis, A. F. Molisch, and G. Caire, "Femtocaching: Wireless content delivery through distributed caching helpers," *IEEE Trans. Inf. Theory*, vol. 59, pp. 8402–8413, Dec. 2013.
- [5] M. Ji, A. M. Tulino, J. Llorca, and G. Caire, "On the average performance of caching and coded multicasting with random demands," in *Proc. 11th International Symposium on Wireless Communications Systems (ISWCS)*, pp. 922–926, Aug. 2014.
- [6] J. Hachem, N. Karamchandani, and S. Diggavi, "Content caching and delivery over heterogeneous wireless networks," in *IEEE Conference on Computer Communications (INFOCOM)*, pp. 756–764, Apr. 2015.
- [7] U. Niesen and M. A. Maddah-Ali, "Coded caching with nonuniform demands," in *IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPs)*, pp. 221–226, Apr. 2014.
- [8] J. Zhang, X. Lin, and X. Wang, "Coded caching under arbitrary popularity distributions," in *Proc. Information Theory and Applications Workshop (ITA)*, Feb. 2015.
- [9] N. Karamchandani, U. Niesen, M. A. Maddah-Ali, and S. Diggavi, "Hierarchical coded caching," in *Proc. IEEE Int. Symp. Information Theory*, pp. 2142–2146, June 2014.
- [10] S. Boyd and L. Vandenberghe, *Convex Optimization*. New York, NY, USA: Cambridge University Press, 2004.
- [11] K. P. Murphy, *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012.
- [12] R. Peck, *Statistics: Learning from Data*. Cengage Learning, 2014.