

Coded Load Balancing in Cache Networks

Mahdi Jafari Siavoshani, Farzad Parvaresh, Ali Pourmiri, Seyed Pooya Shariatpanahi

Abstract—We consider load balancing problem in a cache network consisting of storage-enabled servers forming a distributed content delivery scenario. Previously proposed load balancing solutions cannot perfectly balance out requests among servers, which is a critical issue in practical networks. Therefore, in this paper, we investigate a coded cache content placement where coded chunks of original files are stored in servers based on the files popularity distribution. In our scheme, upon each request arrival at the delivery phase, by dispatching enough coded chunks to the request origin from the nearest servers, the requested file can be decoded.

Here, we show that if n requests arrive randomly at n servers, the proposed scheme results in the maximum load of $O(1)$ in the network. This result is shown to be valid under various assumptions for the underlying network topology. Our results should be compared to the maximum load of two baseline schemes, namely, *nearest replica* and *power of two choices* strategies, which are $\Theta(\log n)$ and $\Theta(\log \log n)$, respectively. This finding shows that using coding, results in a considerable load balancing performance improvement, without compromising communications cost performance. This is confirmed by performing extensive simulation results, in non-asymptotic regimes as well.

Keywords

Distributed Caching Servers, Content Delivery Networks, Coded Caching, Request Routing, Load Balancing, Communication Cost.

I. INTRODUCTION

Here, in the first subsection we present the motivation and the main problem we consider in this paper. Then, in the second subsection, we review related works which consider using coding for delivering contents to the users, and discuss how our work differs from them. Finally, we present the paper structure.

A. Motivation and Paper Contributions

The main objective of a Content Delivery Network (CDN) is fulfilling end-users' content requests by forwarding these requests to distributed caching servers. Such forwarding procedure (aka. request routing) from the request origin to CDN servers is done by a *mapping scheme* that decides which request should be answered by which server, considering the

current network state [1], [2]. However, due to the random nature of these requests, addressing the load balancing issue among the servers is of critical importance [3], [4], i.e., overloading a single server with many requests should be avoided. Moreover, to address scalability issues, the main research in this field has been focused on designing effective *distributed* load balancing schemes [5], [6].

Every request-to-server mapping scheme should manage two main metrics when assigning requests to caching servers. The first metric is *communication cost* which is a measure of distance between the assigned server and the request origin, while the second metric is the *maximum load* incurred to servers, which is the number of requests assigned to the most loaded server. Ideally, one aims to design a mapping scheme that results in the minimum communication cost, while evenly distributing the requests among servers. However, due to the random nature of request arrivals this is not always possible. Interestingly, there is an intrinsic trade-off between these two metrics, and managing this fundamental trade-off is a core issue in load balancing schemes [7], [8], [9], [10]. While assigning each request to the nearest eligible server¹ (i.e., the nearest replica strategy), results in the minimum communication cost, it may incur high loads to servers in proximity of request flash crowds. This observation has resulted in proposing a *proximity-aware power of two choices* strategy in [9] and [10], which queries the load of two nearby servers and assigns the request to the server with lesser load. As proved in [10], this will reduce the maximum load of $\Theta(\log n)$ in the nearest replica strategy to $\Theta(\log \log n)$, at a certain increase in communication cost.

In contrast to all previous load balancing schemes in CDNs, in this paper we follow a fundamentally different approach rather than allocating each request to a *single* responding server. We propose to cache coded chunks of each file in distributed caching servers, which either consist of random linear combinations of the original file chunks, or are constructed via more sophisticated Fountain-like codes [11], [12], [13]. In the delivery phase, each request is distributed among a number of nearby servers which have cached the corresponding coded chunks. Then, the requester can decode the whole file if it receives enough coded chunks from such servers.

In particular, we model the cache network topology by an underlying graph of n nodes which represent cache-enabled servers, where the edges are communication links. At the cache content placement phase, each file is divided into ℓ equal-sized chunks, which are then linearly combined to form coded chunks corresponding to this particular file. Then, each server caches coded chunks corresponding to different files,

¹Here, by eligible servers we mean those servers that have cached the requested file.

The authors' names appear in alphabetical order.

M. Jafari Siavoshani is with the Department of Computer Engineering, Sharif University of Technology, Tehran, Iran (email: mjafari@sharif.edu).

F. Parvaresh is with the Department of Electrical Engineering, University of Isfahan, Isfahan, Iran (email: f.parvaresh@eng.ui.ac.ir) and also School of Mathematics, Institute for Research in Fundamental Sciences (IPM), P.O. Box: 1395-5746, Tehran, Iran. F. Parvaresh was supported by a grant from IPM (No. 95680425).

A. Pourmiri is with the Department of Software Engineering, University of Isfahan, Isfahan, Iran, and Department of Computing, Macquarie University, Sydney, Australia (ali.pourmiri@mq.edu.au).

S. P. Shariatpanahi is with the School of Computer Science, Institute for Research in Fundamental Sciences (IPM), Tehran, Iran (email: pooya@ipm.ir).

based on the file popularity profile. At the delivery phase, n file requests arrive uniformly at random at different servers, where files are requested based on the popularity profile. Then, in order to satisfy each request, ℓ nearest coded chunks corresponding to that file are routed to the request origin via the shortest path on the graph. Assuming independent linear combinations in different coded chunks, it is easy to see that each file can be decoded at the request origin server.

In this setup, we consider Zipf file popularity distribution, and network topologies such as grid, random geometric graphs, random d -regular graphs, and hypercube. We prove that our proposed coded scheme will result in the maximum load of $O(1)$, if each file is divided into $\ell = \Theta(\log n)$ chunks. This result should be compared to the maximum load of nearest replica and power of two choices strategies which are $\Theta(\log n)$ and $\Theta(\log \log n)$, respectively [10]. Interestingly, for grid networks, we show that our proposed coded scheme will result in asymptotically the same communication cost as the baseline uncoded nearest replica strategy. Furthermore, we investigate our findings in finite size networks by performing extensive simulations, through which we also confirm the superiority of our proposal compared to the previous proposed schemes.

B. Related Works

Using coding in content delivery scenarios has been proposed in previous works. The authors in [14] and [15] consider the benefit of using network coding in a P2P based scenario and VANETs, respectively. Also, the papers [16] and [17] investigate the role of network coding in Information Centric Networks, and Critical Infrastructure Networks, respectively. Moreover, [18] considers the effect of coding on the multiple multicast problem. Interested reader can also see [19] which is a good review on the role of network coding for multimedia delivery.

On another line of research, the in-network caching idea has been proposed to relieve network congestion as explained in the following. The authors in [20] investigate optimal coded/uncoded cache content placement for minimizing delivery delay in wireless content delivery scenarios. Following the results in [20] many researchers have considered the role of caching in wireless settings such as [21], [22], and [23].

Also, in [24], the authors look at the problem of optimal MDS coded cache placement in wireless edge networks, which is extended to heterogeneous file and cache sizes in [25]. Furthermore the authors in [26] consider MDS coded cache content placement from an energy consumption perspective. Finally, as it is shown in [27], LT codes can be used at the wireless edge caches to minimize backhaul rate.

In all above works, metrics such as the communication rate/cost, delay, and energy consumption in content delivery networks are investigated, and the benefits of using codes are discussed from different perspectives. However, an important aspect which is ignored in previous works is the load balancing issue. In contrast to all the aforementioned research works in this context, our paper is the first work which investigates the role of coding in load balancing performance through an

analytical approach. Here, we derive closed-form asymptotic results for the impact of coding on the communication cost and maximum load of servers in cache networks. Through the framework developed here, we investigate communication cost and maximum load in grid networks, which is then generalized to other network topologies. Moreover, we confirm our analytical findings through extensive simulations, via which we also investigate the effect of network parameters such as number of servers, cache size, popularity profile, and network topology.

C. Paper Structure

The rest of paper is organized as follows. Our system model and performance metrics are introduced in Section II. Then, the proposed coded scheme is presented in Section III and its maximum load and communication cost are asymptotically analyzed for grid networks. Next, our results are extended to more general networks in Section IV. In Section V, simulation results for finite-sized networks are presented. Finally, the paper is concluded in Section VI.

II. SYSTEM MODEL

A. Notation

Throughout the paper, *with high probability* (w.h.p.) refers to an event that happens with probability $1 - 1/n^c$, as $n \rightarrow \infty$, for some constant $c > 0$. Let $G = (V, E)$ be a graph with vertex set V and edge set E . For $u \in V$ let $\deg(u)$ denote the degree of u in G . For every pair of nodes $u, v \in V$, $\text{dist}_G(u, v)$ denotes the length (number of edges along the path) of a shortest path from node u to v in graph G . The neighborhood of u at distance r in G is defined as

$$B_r(u) \triangleq \{v : \text{dist}_G(u, v) \leq r \text{ and } v \in V(G)\}.$$

For a set A we use \bar{A} to denote its complement. To show the complement of an event \mathcal{E} we use $\neg\mathcal{E}$. We use $\text{Po}(\lambda)$ to denote the Poisson distribution with parameter λ . The expected value of a random variable X is denoted by $\mathbf{E}[X]$. The operator “ \circ ” represents the concatenation operator, i.e., the file $W = W_1 \circ W_2$ is generated by appending the file W_2 at the end of the file W_1 . Throughout the paper, $H(\cdot)$ represents the binary entropy (i.e., information content), measured in bits.

For asymptotic notation, we use $g(n) = O(f(n))$ if there exist c and n_0 such that for $n > n_0$ we have $g(n) < cf(n)$. In this case we also write $f(n) = \Omega(g(n))$. Moreover, we write $g(n) = o(f(n))$ if $\lim_{n \rightarrow \infty} g(n)/f(n) = 0$. In this case one can also write $f(n) = \omega(g(n))$. Finally, if $f(n) = O(g(n))$ and $f(n) = \Omega(g(n))$ then we write $f(n) = \Theta(g(n))$.

B. Problem Setting

We consider a cache network consisting of n caching servers (also called cache-enabled servers) and undirected edges connecting neighboring servers forming a graph G . Direct communication is possible only between adjacent servers, and other communications should be carried out in a multi-hop fashion through the network.

Suppose that the cache network is responsible for handling a library of K files $\mathcal{W} = \{W_1, \dots, W_K\}$, each of size F bits, whereas the popularity profile follows a known distribution $\mathcal{P} = \{p_1, \dots, p_K\}$.

We assume that the network operates in two phases, namely, *cache content placement* and *content delivery*. In the cache content placement phase, each server i caches $Z_i \triangleq \Psi_i(W_1, \dots, W_K)$ such that $H(Z_i) \leq MF$ where M is the cache size of each server. Here $\Psi_i(\cdot)$ is a function of files in \mathcal{W} that generates data Z_i to be placed in the cache of server i .

Consider a time block during which n file requests have arrived uniformly at random among the servers (i.e., graph vertices). Let R_i denote the number of requests (demands) arrived at server i . Then, for large n , we have $R_i \sim \text{Po}(1)$ for all $1 \leq i \leq n$.

For the library popularity profile \mathcal{P} , we consider the Zipf distribution with parameter $\gamma \geq 0$, where the request probability of the k -th popular file is inversely proportional to its rank as follows

$$p_k = \frac{1/k^\gamma}{\sum_{l=1}^K 1/l^\gamma}, \quad k = 1, \dots, K, \quad (1)$$

which has been confirmed to be the case in many practical applications [28], [29].

In the content delivery phase, suppose server i has received a set of requests $\mathcal{W}_i \triangleq \{W_{f_{i,1}}, \dots, W_{f_{i,R_i}}\}$, where $f_{i,j}$ is the file index of j th received request of server i . In order to satisfy the demand $W_{f_{i,j}}$, a set of messages $M_{s \rightarrow i}^{(j)}$ for $s \in \{1, \dots, n\}$ will be sent from other servers to server i where $M_{s \rightarrow i}^{(j)}$ (which is a function of Z_s) is the message sent from server s to server i to satisfy j th request of server i . Then, we say that server i can *successfully decode* the request $W_{f_{i,j}}$ if there exists a decoding function $\Phi_{i,j}(\cdot)$ such that

$$\Phi_{i,j} \left(M_{1 \rightarrow i}^{(j)}, \dots, M_{n \rightarrow i}^{(j)}, Z_i \right) = W_{f_{i,j}}.$$

For a given cache content placement, a *delivery strategy* is defined as follows.

Definition 1 (Delivery Strategy). *By assuming full knowledge of cache contents of all the servers, for each file request $f_{i,j}$, the Delivery Strategy determines the message set $\left\{ M_{s \rightarrow i}^{(j)} \right\}_{s=1}^n$, for all $j \in [1 : R_i]$ and $i \in [1 : n]$.*

Now, for each strategy², we define the following metrics.

Definition 2 (Communication Cost and Maximum Load).

- *The (normalized) communication cost (per request) of a strategy is defined as follows*

$$C \triangleq \frac{1}{nF} \sum_{s \in [1:n]} \sum_{i \in [1:n]} \sum_{j \in [1:R_i]} \text{dist}_G(s, i) H \left(M_{s \rightarrow i}^{(j)} \right). \quad (2)$$

²We use the terms “strategy” and “delivery strategy” alternatively.

TABLE I
NOTATIONS SUMMARY.

Notation	Description
n	Number of servers
K	Number of files in the library
$\mathcal{W} = \{W_1, \dots, W_K\}$	Files in the library
$\mathcal{P} = \{p_1, \dots, p_K\}$	Popularity profile
F	File size in bits
M	Server's cache size
C	Communication cost
L	Maximum load

- *The (normalized) maximum load of a strategy is defined as follows*

$$L \triangleq \frac{1}{F} \max_{s \in [1:n]} \sum_{i \in [1:n]} \sum_{j \in [1:R_i]} H \left(M_{s \rightarrow i}^{(j)} \right).$$

It should be noted that since file requests are random, C and L are random variables in the above definitions.

For convenience of readers, Table I summarizes important notations used throughout the paper.

III. CODED LOAD BALANCING IN GRID NETWORKS

In this section, we focus on grid networks where the graph G is a $\sqrt{n} \times \sqrt{n}$ square wireline grid interconnect where \sqrt{n} is assumed to be an integer.

Remark 1. *In this section, for the sake of presentation clarity, we may consider a torus with n servers, which helps to avoid boundary effects of the grid, and all the asymptotic results hold for the grid as well.*

Let us first revisit the baseline schemes with which we compare our proposed scheme's performance. The simplest uncoded scheme which assigns requests to servers is the so-called *nearest replica* strategy, in which files are cached at the servers proportional to their popularity. Then, each request is assigned to the nearest server having a replica of the requested file. More precisely, let us focus on server i at which R_i requests have arrived. For request $j \in \{1, \dots, R_i\}$, i.e., $W_{f_{i,j}}$, we denote the nearest server to i which has cached this request as s_j . Then the delivery strategy will be

$$\begin{aligned} M_{s_j \rightarrow i}^{(j)} &= W_{f_{i,j}}, \\ M_{s \rightarrow i}^{(j)} &= \emptyset \quad \text{for all } s \neq s_j. \end{aligned} \quad (3)$$

Remark 2. *It should be noted that the nearest replica strategy achieves the minimum communication cost among all uncoded schemes. This is true since this scheme minimizes all terms in the summations of (2) (i.e., the definition of communication cost) in Definition 2.*

While this scheme performs well in terms of communication cost, its maximum load will be of order $\log(n)$ (for more details refer to [10]). The second baseline scheme considered in this paper is the one proposed in [10], which is based on the concept of power of two choices [30]. In this scheme, instead of allocating the request to the nearest server having cached that request, the current load of two servers in a limited distance of the request origin is queried, and then the request

is assigned to the server with the lower load. In a certain regime of problem parameters, this will result in the maximum load of $\log \log(n)$, which is achieved at the cost of queries' complexity, and higher communication cost compared to the nearest replica strategy [10].

In contrast, in this section we propose a coded scheme which achieves the maximum load of order $O(1)$, while maintaining almost the same communication cost as the nearest replica strategy, and also does not require querying current load of any servers. To this end, we consider a coded cache placement in the servers using random linear coding, e.g., Fountain-like codes [11], [12], [13].

Let us first take each file W_k of F bits, and partition it into ℓ equal-sized chunks, i.e., $W_k^{(r)}$ where $r \in [1 : \ell]$, each of F/ℓ bits. Thus, each file W_k becomes the concatenation of the chunks $W_k^{(r)}$, $r \in [1 : \ell]$, i.e., $W_k = W_k^{(1)} \circ W_k^{(2)} \circ \dots \circ W_k^{(\ell)}$. Then, we define a random linear operator \mathcal{L} as

$$\mathcal{L}(W_k) \triangleq \sum_{r=1}^{\ell} \alpha_r W_k^{(r)}, \quad (4)$$

where α_r 's are constants chosen uniformly at random over a finite field \mathbb{F}_q with $q = \Theta(2^{F/\ell})$, and the summation is over \mathbb{F}_q (assuming there is an injective mapping from the file chunks $W_k^{(r)}$ to elements of the finite field \mathbb{F}_q). Moreover, we assume that each use of the operator \mathcal{L} is independent of other instances.

In the cache content placement phase, each server stores ℓM coded chunks of the files according to the popularity distribution \mathcal{P} as explained in Algorithm 1.

Algorithm 1 Coded cache content placement for server i

Require: ℓ , M , \mathcal{P} , and \mathcal{W}

- 1: **repeat**
 - 2: Sample k according to distribution \mathcal{P}
 - 3: Store a coded chunk $\mathcal{L}(W_k)$ at server i 's cache
 - 4: **until** server i 's cache is full
-

In the content delivery phase, we assume n requests arrive uniformly at random at network servers. Let us consider server i which has $\text{Po}(1)$ requests, i.e., $R_i \sim \text{Po}(1)$, for large n . For satisfying each request, ℓ nearest (in terms of shortest path in the grid network) coded chunks of the file corresponding to that request should be routed to this server according to Algorithm 2.

Algorithm 2 Coded delivery phase for server i

Require: $\{f_{i,j}\}_{j=1}^{R_i}$, ℓ

- 1: **for** $j = 1 : R_i$ **do**
 - 2: $\mathcal{I} :=$ indices of the ℓ nearest servers caching coded chunks corresponding to file $W_{f_{i,j}}$
 - 3: **for** $s \in \mathcal{I}$ **do**
 - 4: $M_{s \rightarrow i}^{(j)} :=$ a coded chunk of file $W_{f_{i,j}}$ at server s
 - 5: Forward $M_{s \rightarrow i}^{(j)}$ from server s to server i via the shortest path in the grid network
 - 6: **end for**
 - 7: **end for**
-

It is clear that if the field size q is large enough, then Algorithm 2 will successfully send the necessary file chunks to server i , so that server i can reconstruct the requested files by its users with the probability of order $1 - O(1/q)$, e.g., see [31, Lemma 1].

In the next theorem, we characterize the load balancing performance of the proposed scheme.

Theorem 1 (Maximum Load of Grid). *For a $\sqrt{n} \times \sqrt{n}$ grid network, suppose that $\mathcal{P} = \{p_1, p_2, \dots, p_K\}$ be the file popularity distribution and let the number of file chunks be $\ell = \Omega(\log n)$. Then, the maximum load of Algorithm 2 is $O(1)$, w.h.p.*

Proof. In Algorithm 2, we have to find ℓ coded chunks corresponding to each arriving request. Let us focus on server u . Define \mathcal{E}_u to be the event that for all $k \in [1 : K]$, one can find ℓ coded chunks corresponding to W_k in $B_{r_k}(u)$ where $r_k \triangleq \sqrt{\alpha \ell / \tilde{p}_k}$ for some positive constant α . In above $\tilde{p}_k \triangleq 1 - (1 - p_k)^{M \ell}$ is the probability that an arbitrary server has cached at least a coded chunk of W_k . Then, we can state Lemma 1, which appears after the Theorem proof.

Now, for every fixed server $u \in V(G)$, let us define the indicator random variable $Y_{u,j}$, taking one if the j -th request ($j \in [1 : n]$) has asked a coded chunk from server u and zero otherwise. Notice that a server might be asked to respond a request if the following events occur:

- (E1) a request for file W_k is born at a neighborhood of radius r_k of the server, and
- (E2) a coded chunk of the requested file is cached in the server.

Thus, we can write

$$\begin{aligned} \Pr[Y_{u,j} = 1] &= \Pr[Y_{u,j} = 1 | \mathcal{E}] \Pr[\mathcal{E}] \\ &\quad + \Pr[Y_{u,j} = 1 | \neg \mathcal{E}] \Pr[\neg \mathcal{E}] \\ &\stackrel{(a)}{=} \Pr[Y_{u,j} = 1 | \mathcal{E}] \Pr[\mathcal{E}] \\ &\quad + \Pr[Y_{u,j} = 1 | \neg \mathcal{E}] \times o(1/n) \\ &\leq \Pr[Y_{u,j} = 1 | \mathcal{E}] + o(1/n) \\ &= \sum_{k=1}^K \Pr[Y_{u,j} = 1 | \mathcal{E}, W_k \text{ requested}] p_k + o(1/n) \\ &\stackrel{(b)}{\leq} \sum_{k=1}^K \frac{|B_{r_k}(u)|}{n} \tilde{p}_k \cdot p_k + o(1/n) \\ &\stackrel{(c)}{\leq} \frac{2\alpha \ell (1 + o(1))}{n} + o(1/n) \\ &= \frac{2\alpha \ell (1 + o(1))}{n}, \end{aligned}$$

where $\mathcal{E} = \bigcap_{u \in V(G)} \mathcal{E}_u$, (a) follows from Lemma 1, (b) follows from considering the probabilities of events (E1) and (E2) above, and (c) follows from $B_{r_k}(u) = 2r_k^2(1 + o(1))$ since G is a grid.

Now, let $S_u = \sum_{j=1}^n Y_{u,j}$ count the number of requests that are responded by server u , during allocating the total n requests. Hence, we have

$$\mathbf{E}[S_u] = \sum_{j=1}^n \mathbf{E}[Y_{u,j}] \leq 2\alpha \ell (1 + o(1)).$$

Applying a Chernoff bound for S_u implies that

$$\Pr[S_u \geq (1 + \delta)2\alpha\ell] \leq \exp(-\delta^2\alpha\ell) = o(1/n^2),$$

for appropriate choices of constants δ and α , and since $\ell = \Omega(\log n)$. Taking union bound over all servers shows that each server is requested at most $O(\ell(1 + o(1)))$ times where each request involves sending F/ℓ bits. Hence, it has to handle at most $O(1)$ bits. This concludes the proof. \square

Lemma 1. *For the event \mathcal{E}_u defined above, we have $\Pr[\mathcal{E}_u] = 1 - o(n^{-2})$. Then, the event $\mathcal{E} = \bigcap_{u \in V(G)} \mathcal{E}_u$ happens with probability $\Pr[\mathcal{E}] = 1 - o(n^{-1})$.*

Proof. For a given $k \in [1 : K]$, let $X_{u,k}$ denote an indicator random variable taking 1 if node u has cached a coded chunk of W_k and zero otherwise. Thus,

$$\Pr[X_{u,k} = 1] = \tilde{p}_k = 1 - (1 - p_k)^{M\ell},$$

where according to Algorithm 1, p_k is the probability that a coded chunk of W_k has been cached on a server. It is easy to see that $Z_{u,k} = \sum_{v \in B_{r_k}(u)} X_{v,k}$ counts the number of servers that have cached W_k in a neighborhood of u . For every node u , we have

$$|B_{r_k}(u)| = \Theta(r_k^2) = \frac{\alpha\ell}{\tilde{p}_k}.$$

Since $X_{u,k}$'s are i.i.d random variables and by linearity of expectation, we get that

$$\mathbf{E}[Z_{u,k}] = r_k^2 \tilde{p}_k = \alpha\ell.$$

Provided $\alpha > 2$ and $\alpha\ell/8 > \log n$ is an appropriate constant, applying a Chernoff bound results that

$$\begin{aligned} \Pr[Z_{u,k} < \ell] &\leq \Pr[Z_{u,k} \leq \mathbf{E}[Z_{u,k}]/2] \leq e^{-\mathbf{E}[Z_{u,k}]/8} \\ &= e^{-\alpha\ell/8} = o(1/n^3), \end{aligned}$$

where the last equality follows from $\ell = \Omega(\log n)$. By applying union bound over all $K \leq n$ files in the library, we have

$$\sum_{k \in [1, K]} \Pr[Z_{u,k} \leq \ell] = o(1/n^2).$$

This implies that $\Pr[\mathcal{E}_u] = 1 - o(1/n^2)$ and hence by another application of the union bound we get

$$\begin{aligned} \Pr[\mathcal{E}] &= 1 - \Pr[\bigcup_{u \in V(G)} \neg \mathcal{E}_u] \\ &\geq 1 - \sum_{u \in V(G)} \Pr[\neg \mathcal{E}_u] \\ &= 1 - o(1/n). \end{aligned}$$

\square

Next, we analyze the communication cost of the proposed scheme in Algorithm 2. To this end, we first prove a general expression for communication cost of requests with arbitrary popularities in Theorem 2. Then, in Corollary 1 we specialize the result of Theorem 2 to the Zipf popularity profile.

Theorem 2 (Communication Cost). *For a grid network of size $\sqrt{n} \times \sqrt{n}$, suppose that $\mathcal{P} = \{p_1, p_2, \dots, p_K\}$ be the file popularity distribution, and assume $K = O(n)$ and $\ell = \Omega(\log n)$. Define $\tilde{p}_k \triangleq 1 - (1 - p_k)^{M\ell}$ where we assume*

$\sqrt{\ell/\tilde{p}_k} = o(\sqrt{n})$ for every $k \in [1 : K]$. Then, w.h.p., the communication cost for every requested file W_k , $1 \leq k \leq K$, is $\Theta(\sqrt{\ell/\tilde{p}_k})$. Moreover, we have

$$\mathbf{E}[C] = \sum_{k=1}^K \Theta(\sqrt{\ell/\tilde{p}_k}) p_k. \quad (5)$$

Proof. Upper Bound: Suppose that $u \in G$ is an arbitrary server in the grid and a request for file W_k arrives at server u . For every server $v \in G$, let $X_{v,k}$ denote the indicator random variable taking 1 if server v has cached at least one coded chunk of W_k , and zero otherwise. Then, for every positive number r , $Y_u(k, r) = \sum_{v \in B_r(u)} X_{v,k}$ denotes the number of servers that have cached a coded chunk of W_k in $B_r(u)$. Notice that

$$\tilde{p}_k = \Pr[X_{v,k} = 1] = 1 - (1 - p_k)^{M\ell}.$$

Thus,

$$\mathbf{E}[Y_u(k, r)] = |B_r(u)| \cdot \tilde{p}_k.$$

We know that, for a grid network, $|B_r(u)| = 2r(r+1) + 1 = 2r^2(1 + o(1))$. Hence by choosing $r_k = \sqrt{z/\tilde{p}_k}$, where $z = \max[6 \log n, 5\ell]$, we get

$$\mathbf{E}[Y_u(k, r_k)] = 2z(1 + o(1)) \geq 12 \log n(1 + o(1)).$$

Since $X_{v,k}$'s are independent and identical indicator random variables, applying a Chernoff bound for $Y_u(k, r_k)$ yields that

$$\Pr[Y_u(k, r_k) \leq 0.1 \mathbf{E}[Y_u(k, r_k)]] = o(1/n^3).$$

Thus, with probability $1 - o(1/n^3)$

$$Y_u(k, r_k) \geq 0.1 \mathbf{E}[Y_u(k, r_k)] = 0.2z(1 + o(1)) \geq \ell,$$

which means $B_{r_k}(u)$ contains at least ℓ servers that have cached a coded chunk of W_k .

Let $\mathcal{E}_{u,k}$ denote the event that server u requests for file W_k and $Y_u(k, r_k) < \ell$. Now by the union bound over all n servers and $K = O(n)$ files we have,

$$\Pr[\bigcup_{u,k} \mathcal{E}_{u,k}] \leq \sum_{u,k} \Pr[\mathcal{E}_{u,k}] = \sum_{u,k} o(1/n^3) = o(1/n).$$

Since the number of bits in each coded chunk is F/ℓ and with probability $1 - o(1/n)$ we can find all the required coded chunks in $B_{r_k}(u)$, the communication cost defined in Definition 2 is $O(\sqrt{\ell/\tilde{p}_k})$, since $z = O(\ell)$.

Lower Bound: If we set $r_k = o(\sqrt{z/\tilde{p}_k})$, then $\mathbf{E}[Y(k, r_k)] = o(z) = o(\ell)$. Thus, by Markov inequality for any constant $\alpha > 0$

$$\Pr[Y_u(k, r_k) > \alpha \mathbf{E}[Y_u(k, r_k)]] \leq 1/\alpha.$$

So for every u , with probability at least $1 - 1/\alpha$,

$$Y_u(k, r_k) < \alpha \mathbf{E}[Y_u(k, r_k)] < \ell.$$

Let $T_{u,j}$ denote the indicator random variable taking one if the j th request, received by server u , fails to find ℓ coded chunks of the requested file in the set $B_{r_k}(u)$. Thus, we have

$\Pr [T_{u,j} = 1] > 1 - 1/\alpha$. Also let $S = \sum_{u=1}^n \sum_{j=1}^{R_u} T_{u,j}$ denote the total number of failures. Then,

$$\mathbf{E}[S] = n \cdot \Pr [T_{u,j} = 1] > (1 - 1/\alpha)n.$$

Since the requests are independent, another application of the Chernoff bound for random variable S results that

$$\Pr [S < \mathbf{E}[S] / 2] < \exp(-\Omega(n)).$$

Hence, w.h.p. at least $(1 - 1/\alpha)n/2$ requests cannot be responded in the r -neighborhood of the requesting server, and the communication cost for W_k is $\Omega(\sqrt{\ell/\tilde{p}_k})$.

Since the upper and lower bounds meet, the communication cost of requesting file W_k is $\Theta(\sqrt{\ell/\tilde{p}_k})$. Consequently, by averaging over the library with probability distribution \mathcal{P} , the average communication cost is

$$\mathbf{E}[C] = \sum_{k=1}^K \Theta\left(\sqrt{\ell/\tilde{p}_k}\right) p_k.$$

□

The following corollary of Theorem 2 characterizes the communication cost of our problem under the Zipf popularity profile.

Corollary 1. *For a grid network of size $\sqrt{n} \times \sqrt{n}$, suppose that the cache size M of each server is a constant, the number of files is $K = n^\delta$, for any $\delta \in (0, 1]$, and the number of file chunks is $\ell = \Theta(\log n)$. Then, for Zipf popularity distribution with parameter γ , the average communication cost is*

$$\mathbf{E}[C] = \begin{cases} \Theta\left(\sqrt{K/M}\right) & : 0 \leq \gamma < 1, \\ \Theta\left(\sqrt{K/M \log K}\right) & : \gamma = 1, \\ \Theta\left(K^{1-\gamma/2}/\sqrt{M}\right) & : 1 < \gamma < 2, \\ \Theta\left(\log K/\sqrt{M}\right) & : \gamma = 2, \\ \Theta(\sqrt{\ell}) & : \gamma > 2. \end{cases}$$

Proof. To show this, for some small constant $\epsilon > 0$, let us define two sets, say $A \triangleq \{k : 1 \leq k \leq K, p_k M \ell > \epsilon\}$ and its complement \bar{A} . So for every $k \in A$, we have $\tilde{p}_k = \Theta(1)$ and for every $k \in \bar{A}$,

$$\tilde{p}_k = 1 - (1 - p_k)^{M\ell} \approx p_k M \ell.$$

Thus, the average communication cost reduces to the following summation

$$\mathbf{E}[C] = \sum_{k \in A} \Theta(\sqrt{\ell}) p_k + \sum_{k \in \bar{A}} \Theta(\sqrt{p_k/M}).$$

For Zipf distribution, p_k 's are decreasing in k so let us define $k^* = \max A$ if A is not empty and $k^* = 0$, otherwise. Also let $p_0 = 0$. Then,

$$\mathbf{E}[C] = \underbrace{\Theta(\sqrt{\ell}) \sum_{k=0}^{k^*} p_k}_{S_1(\gamma)} + \underbrace{\sum_{k=k^*+1}^K \Theta(\sqrt{p_k/M})}_{S_2(\gamma)}. \quad (6)$$

Now let us estimate $S_1(\gamma)$ in (6). It is clear that $0 \leq \sum_{k=0}^{k^*} p_k < 1$, thus

$$S_1(\gamma) = \begin{cases} O(\sqrt{\ell}) & : 0 \leq \gamma \leq 1, \\ \Theta(\sqrt{\ell}) & : \gamma > 1, \end{cases} \quad (7)$$

where $S_1(\gamma) = \Theta(\sqrt{\ell})$ as $p_1 = \Theta(1)$ when $\gamma > 1$. In what follows we provide an estimation for $S_2(\gamma) = \sum_{k=k^*+1}^K \Theta(\sqrt{p_k/M})$. Let us define $\Lambda(\gamma, s) \triangleq \sum_{k=s+1}^K k^{-\gamma}$. So

$$S_2(\gamma) = \sum_{k=k^*+1}^K \Theta\left(\frac{k^{-\gamma/2}}{\sqrt{M\Lambda(\gamma, 0)}}\right) = \Theta\left(\frac{\Lambda(\gamma/2, k^*)}{\sqrt{M\Lambda(\gamma, 0)}}\right).$$

Next, in order to proceed, we need Lemma 2 stated after the Theorem's proof. From Lemma 2, we know that

$$\Lambda(\gamma, k^*) = \begin{cases} \Theta(K^{1-\gamma}) - \Theta(k^{*1-\gamma}) & : 0 \leq \gamma < 1, \\ \Theta(\log K) - \log k^* & : \gamma = 1, \\ O(1) & : \gamma > 1. \end{cases}$$

Note that by the definition of k^* , we have

$$\frac{\epsilon}{M\ell} \leq p_{k^*} < k^{*-\gamma},$$

and hence $k^* \leq (\frac{M\ell}{\epsilon})^{1/\gamma} = O((\log n)^{1/\gamma})$. This implies that we can ignore terms $k^{*1-\gamma}$ and $\log k^*$ in comparison with $\Theta(K^{1-\gamma})$ and $\Theta(\log K)$, respectively, since $K = n^\delta$, for some constant $\delta \in (0, 1]$. In other words we can write

$$S_2(\gamma) = \Theta\left(\frac{\Lambda(\gamma/2, 0)}{\sqrt{M\Lambda(\gamma, 0)}}\right).$$

Then, we will have

$$S_2(\gamma) = \begin{cases} \Theta\left(\sqrt{K/M}\right) & : 0 \leq \gamma < 1, \\ \Theta\left(\sqrt{K/M \log K}\right) & : \gamma = 1, \\ \Theta\left(K^{1-\gamma/2}/\sqrt{M}\right) & : 1 < \gamma < 2, \\ \Theta\left(\log K/\sqrt{M}\right) & : \gamma = 2, \\ O(1) & : \gamma > 2. \end{cases} \quad (8)$$

Now, considering equalities (7) and (8) completes the proof. □

Lemma 2. *For $\Lambda(\gamma, k) = \sum_{l=k}^K l^{-\gamma}$, where $\gamma \geq 0$, we have*

$$\Lambda(\gamma, k) = \begin{cases} \Theta(K^{1-\gamma}) - \Theta(k^{1-\gamma}) & : 0 \leq \gamma < 1, \\ \Theta(\log K) - \log k & : \gamma = 1, \\ O(1) & : \gamma > 1. \end{cases}$$

Proof. First notice that by the definition of Riemann integration, we can upper and lower bound $\Lambda(\gamma, k)$ as follows

$$\int_k^{K+1} t^{-\gamma} dt \leq \sum_{l=k}^K l^{-\gamma} \leq k^{-\gamma} + \int_k^K t^{-\gamma} dt.$$

This can be simplified to

$$\frac{1}{1-\gamma} [(K+1)^{1-\gamma} - k^{1-\gamma}] \leq \Lambda(\gamma, k)$$

and

$$\Lambda(\gamma, k) \leq k^{-\gamma} + \frac{1}{1-\gamma} [K^{1-\gamma} - k^{1-\gamma}].$$

Now, considering the following three cases $0 \leq \gamma < 1$, $\gamma = 1$, and $\gamma > 1$ will conclude the proof. \square

It is interesting to note that the nearest replica strategy also arrives at almost the same communication cost derived in Corollary 1 (see [10, Theorem 3]).

IV. CODED LOAD BALANCING IN GENERAL NETWORKS

The load balancing behavior of our proposed algorithm in last section, which was proved for grid networks, can be generalized for other networks as well. In this section, we will generalize our results for other networks with some symmetry properties. In order to do this, in this paper, we define *Locally Symmetric Networks* as follows.

Definition 3 (Locally Symmetric). We define a class of graphs \mathcal{G} locally symmetric, if for every $G \in \mathcal{G}$, for an arbitrary integer $r \geq 1$, and for every $u, v \in V(G)$, we have

$$\max \left[\frac{|B_r(u)|}{|B_r(v)|}, \frac{|B_r(v)|}{|B_r(u)|} \right] = O(1).$$

Then, a graph is called locally symmetric if it belongs to such a class.

As will be proved later, *random geometric graphs* and *random regular graphs* are two examples of locally symmetric networks, w.h.p.

To generalize our results, here we need a more general notion of r_k which was introduced in the previous section as the search distance for finding ℓ coded chunks of W_k . Here, for a given network G and $k \in [1 : K]$, we define r_k as the smallest integer such that for every node u in G , we have

$$|B_{r_k}(u)| \geq \frac{32 \log n}{\tilde{p}_k} = \frac{\alpha_k \ell}{\tilde{p}_k},$$

for some constants $\alpha_k, k \in [1 : K]$.

We also need to slightly modify Algorithm 2 for locally symmetric class of graphs, as stated in Algorithm 3.

Algorithm 3 Coded delivery phase for server i in locally symmetric graphs

Require: $\{f_{i,j}\}_{j=1}^{R_i}, \ell, \{r_1, r_2, \dots, r_k\}$

- 1: **for** $j = 1 : R_i$ **do**
 - 2: $\mathcal{J} :=$ set of servers that have cached a coded chunk of $W_{f_{i,j}}$ at distance at most r_k from server i
 - 3: $\mathcal{I} :=$ set of ℓ servers that are selected uniformly at random from \mathcal{J}
 - 4: **for** $s \in \mathcal{I}$ **do**
 - 5: $M_{s \rightarrow i}^{(j)} :=$ a coded chunk of file $W_{f_{i,j}}$ at server s
 - 6: Forward $M_{s \rightarrow i}^{(j)}$ from server s to server i via the shortest path in the network
 - 7: **end for**
 - 8: **end for**
-

Remark 3. Notice that for grid networks, Algorithm 3 is reduced to Algorithm 2. This happens because the sets \mathcal{I} and \mathcal{J} become the same w.h.p.

The load balancing behavior of Algorithm 3 for locally symmetric networks is characterized in the following theorem.

Theorem 3 (Maximum Load of Locally Symmetric Graphs). Suppose that $\mathcal{P} = \{p_1, p_2, \dots, p_K\}$ be the file popularity distribution and let the number of file chunks be $\ell = \Omega(\log n)$. Also, assume that $G = (V, E)$ be a locally symmetric network. Then, the maximum load achieved by Algorithm 3 is $O(1)$, w.h.p.

Proof. For every node $u \in V(G)$, let \mathcal{E}_u denote the event that for every $k \in [1 : K]$, there exist at least $|B_{r_k}(u)|\tilde{p}_k/2$ coded chunks of W_k in $B_{r_k}(u)$. Then, we have Lemma 3 stated after the Theorem's proof, whose proof is similar to that of Lemma 1.

Similar to the proof of Theorem 1, by defining $Y_{u,j}$ (taking one if the j -th request, $j \in [1 : n]$, has asked a coded chunk from server u and zero otherwise), we can write

$$\begin{aligned} \Pr[Y_{u,j} = 1] &= \Pr[Y_{u,j} = 1 | \mathcal{E}] \Pr[\mathcal{E}] \\ &\quad + \Pr[Y_{u,j} = 1 | \neg \mathcal{E}] \Pr[\neg \mathcal{E}] \\ &\stackrel{(a)}{\leq} \Pr[Y_{u,j} = 1 | \mathcal{E}] + o(1/n) \\ &= \sum_{k=1}^K \Pr[Y_{u,j} = 1 | \mathcal{E}, W_k \text{ requested}] p_k + o(1/n) \\ &\stackrel{(b)}{\leq} \sum_{k=1}^K \sum_{w \in B_{r_k}(u)} \frac{\tilde{p}_k \cdot p_k}{n} \cdot \frac{\ell}{X_{w,k}} + o(1/n) \\ &\stackrel{(c)}{\leq} \sum_{k=1}^n \sum_{w \in B_{r_k}(u)} \frac{2p_k \ell}{n|B_{r_k}(w)|} + o(1/n), \end{aligned}$$

where $\mathcal{E} = \bigcap_{u \in V(G)} \mathcal{E}_u$, (a) follows from Lemma 3, in (b) $X_{w,k}$ is defined to be the number of servers in $B_{r_k}(w)$ that have cached a coded chunk of W_k , and (c) follows from the fact that conditioned on \mathcal{E} , we have $X_{k,u} \geq |B_{r_k}(u)| \cdot \tilde{p}_k/2$. Thus, we have

$$\begin{aligned} \Pr[Y_{u,j} = 1] &\leq \sum_{k=1}^n \sum_{w \in B_{r_k}(u)} \frac{2p_k \ell}{nm} + o(1/n) \\ &= \sum_{k=1}^n |B_{r_k}(u)| \frac{2p_k \ell}{nm} + o(1/n) \\ &\stackrel{(a)}{\leq} \sum_{k=1}^n \frac{2c\ell p_k}{n} + o(1/n) \\ &\leq \frac{2c\ell + 1}{n} \end{aligned}$$

where

$$m \triangleq \min_{w \in B_{r_k}(u)} |B_{r_k}(w)|,$$

and in (a) we have used the fact that if G is locally symmetric, then for some constant c we have $|B_{r_k}(u)| \leq cm$.

Now, let $S_u = \sum_{j=1}^n Y_{u,j}$ count the number of requests that are responded by server u , during allocating n requests. Hence, we have

$$\mathbf{E}[S_u] = \sum_{j=1}^n \mathbf{E}[Y_{u,j}] \leq 2c\ell + 1.$$

Applying a Chernoff bound for S_u implies that

$$\Pr[S_u \geq (1 + \delta)2c\ell] \leq \exp(-\delta^2 c\ell) = o(1/n^2),$$

for appropriate choice of constant δ , and since $\ell = \Omega(\log n)$. Thus, taking union bound over all servers shows that each server is requested at most $O(\ell)$ times where each request involves sending F/ℓ bits. Hence, it has to handle at most $O(1)$ bits. This concludes the proof. \square

Lemma 3. *For the event \mathcal{E}_u defined above, we have $\Pr[\mathcal{E}_u] = 1 - o(n^{-2})$. Then, the event $\mathcal{E} = \bigcap_{u \in V(G)} \mathcal{E}_u$ happens with probability $\Pr[\mathcal{E}] = 1 - o(n^{-1})$.*

To apply the above result for other network topologies, it is sufficient to show that those networks are locally symmetric. Here, we consider three well-studied families of networks, namely, hypercube, random geometric graphs and random regular graphs.

Hypercube: A hypercube, also called n -cube, is a network where the nodes can be presented by n -bit binary words. Two nodes are connected if and only if the nodes only differ in one bit.

Proposition 1. *For every given $n > 0$, n -cube is locally symmetric.*

Proof. By the definition of the n -cube, two nodes are connected if and only if they differ in only one bit. It is not hard to see that if u is adjacent to v , then for every n -bit x , $u \oplus x$ is also adjacent to $v \oplus x$. Note that we use “ \oplus ” to denote the XOR of two words. Thus, for each node u and given an integer $r \geq 1$, we have $|B_r(u)| = \sum_{i=0}^r \binom{n}{i}$, which satisfies the locally symmetric graph constraint defined in Definition 3. \square

Random Geometric Graphs: A random geometric graph, denoted by $\text{RG}(n, \lambda)$, is a network with n nodes chosen uniformly at random in the square $[0, \sqrt{n}] \times [0, \sqrt{n}]$. Two nodes are connected if and only if their Euclidean distance is at most λ . It is well-known (e.g., see [32]) that if we set $\lambda > \lambda_c = \sqrt{\frac{\log n + O(1)}{\pi}}$, then, the obtained graph is connected, w.h.p.

Proposition 2. *For a given n , a typical n -node random geometric network is locally symmetric, w.h.p.*

Proof. To show the proposition, we apply a result from [32] showing that if G is a realization of $\text{RG}(n, \lambda)$ with $\lambda = \omega(\sqrt{\log n})$, then for every $u, v \in V(G)$ we have

$$\frac{\text{dist}_{\mathbf{E}}(u, v)}{\lambda} \leq \text{dist}_G(u, v) \leq \frac{\text{dist}_{\mathbf{E}}(u, v)}{\lambda} (1 + o(1))$$

w.h.p., where $\text{dist}_{\mathbf{E}}(\cdot, \cdot)$ denotes the Euclidean distance between u and v . By applying the above inequalities, we have $v \in B_r(u)$ if and only if

$$r\lambda \leq \text{dist}_{\mathbf{E}}(u, v) \leq r\lambda(1 + o(1)).$$

This implies that for each $u \in V(G)$,

$$D_{r\lambda}(u) \subseteq B_r(u) \subseteq D_{r\lambda(1+o(1))}(u),$$

where $D_l(u)$ is a subset of $V(G)$ whose Euclidean distances from u is at most l . Clearly, one may see that the probability that one node falls in $D_l(u)$ is $\frac{\pi l^2}{n}$, where $l^2\pi$ is the area of the disk $D_l(u)$. Since nodes are drawn uniformly and independently at random from $[0, \sqrt{n}] \times [0, \sqrt{n}]$, the expected number of nodes falls in $D_l(u)$ is $l^2\pi$. Applying a Chernoff bound, we conclude that the number of nodes in $D_l(u)$ is $\Theta(l^2\pi)$. By setting l to be either $r\lambda$ or $r\lambda(1 + o(1))$, we get that $|B_r(u)| = \Theta(r^2\lambda^2)$. Therefore, for a typical random geometric graph G , G is locally symmetric. \square

Random d -Regular Graphs: Suppose that for every n and $d \geq 3$, $\mathcal{G}_{n,d}$ is the family of all d -regular graphs with n nodes. Assume that G is a randomly chosen graph from $\mathcal{G}_{n,d}$. Then, it is well-known that w.h.p., for every $r = o(\log_d n)$, and $u \in V(G)$, the subgraph induced by set $B_r(u)$ looks like a d -ary tree rooted at r . Hence, we have $|B_r(u)| = \Theta(d(d-1)^{r-1})$, which holds for every $u \in V(G)$ and $r = o(\log_d n)$ (for more details refer to [33]). Thus, it is clear that G is locally symmetric for $r = o(\log_d n)$ which is sufficient for Theorem 3 to be valid for random d -regular graphs.

V. NUMERICAL ANALYSIS

In this section, we use Monte Carlo simulations to investigate the communication cost and maximum load performance of a content delivery network, under the model described in Section II. Our simulator has been written in Python and can run multiple instances of simulations in parallel [34]. Here, our main goal is twofold. First, we verify our key understandings from asymptotic theoretical results about the benefits of coding. Second, we investigate other aspects of the proposed scheme not revealed in the asymptotic analyses.

In order to clarify the effect of different parameters on the network performance, we start numerical analysis with a uniform popularity distribution, i.e., $\gamma = 0$ for the Zipf distribution, and then investigate the effect of varying γ on the results. In all of the following simulations, each data point is obtained by taking average over 5000 simulation runs.

In Fig. 1 the maximum load L is plotted versus the network size n for the *nearest replica*, *power of two choices*, and *coded load balancing* strategies. Here, the network topology is a grid, the library size is $K = 100$ and we assume a uniform popularity profile. Note that the maximum load of power of two choices scheme is better than the coded scheme when chunk size is $\ell = 4$. However, if chunk size is increased to $\ell = 10$, our proposed scheme surpasses all baseline schemes. Also, Fig. 2 shows similar trend for the cache size $M = 10$.

In order to investigate the role of chunk size ℓ in the performance of proposed scheme, in Fig. 3, we have considered

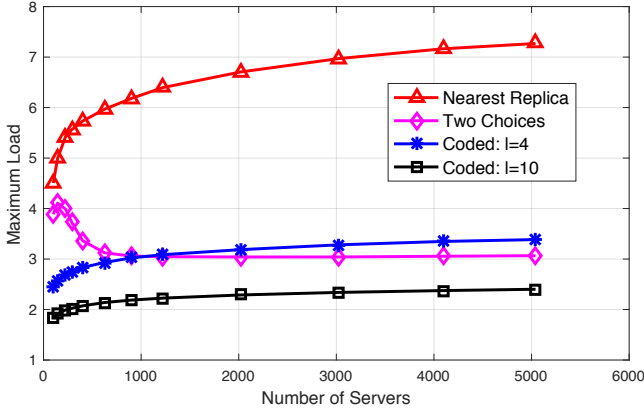


Fig. 1. Comparing the maximum load L of the nearest replica, power of two choices, and coded load balancing strategies versus the number of servers n . The library size is $K = 100$ and the cache size is $M = 2$.

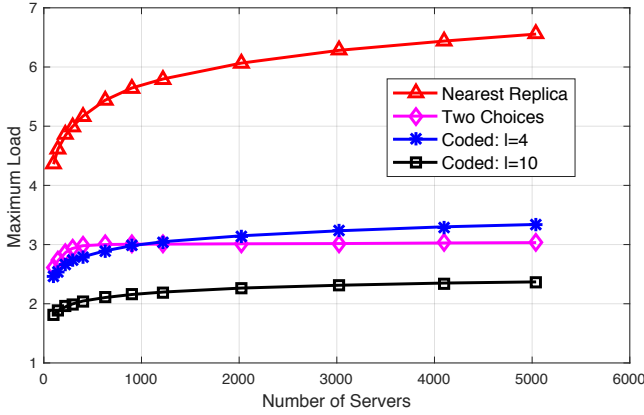


Fig. 2. Comparing the maximum load L of the nearest replica, power of two choices, and coded load balancing strategies versus the number of servers n . The library size is $K = 100$ and the cache size is $M = 10$.

a grid network of size $n = 1024$, a library size of $K = 100$, and cache sizes $M \in \{1, 10\}$. For comparison, the results of nearest replica and power of two choices strategies are also highlighted in this figure. As it is observed in Fig. 3, in order to obtain the coding benefit (i.e., surpassing the power of two choices performance), the chunk size should be above a threshold.

In Fig. 4 the maximum load is depicted versus the cache size M of each server. Here, the network topology is grid, $n = 1024$, and $K = 100$. It is interesting to note that for large enough cache size, the maximum load does not decrease as cache size increases.

Thus far, for the purpose of clarity, we have investigated different aspects of the proposed coded strategy in the case of $\gamma = 0$. Here, we move forward to investigate the role of Zipf parameter γ on the network performance as presented in Fig. 5, in which the maximum load is plotted versus γ . We observe that the maximum load does not depend on the popularity profile (here characterized by γ), which is consistent with our finding in Theorem 1. This is a direct consequence of our *proportional* cache content placement

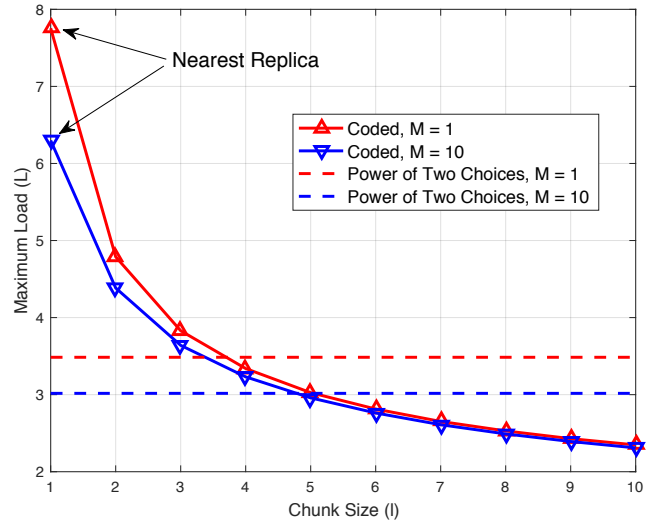


Fig. 3. The maximum load L of the proposed coded load balancing strategy versus the chunk size l . The library size is $K = 100$ and the number of servers is $n = 1024$.

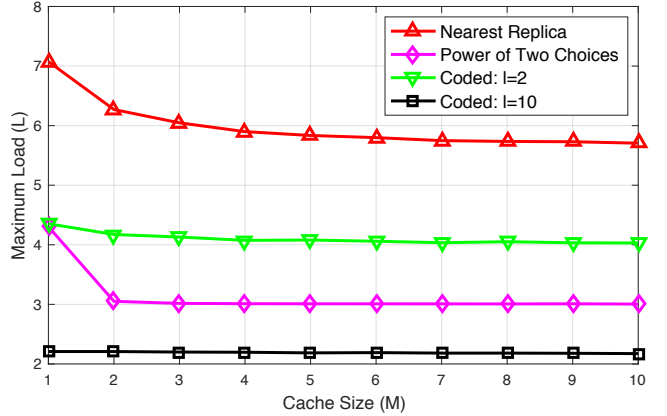


Fig. 4. The maximum load L of proposed coded load balancing strategy versus the cache size M . The library size is $K = 100$ and the number of servers is $n = 1024$.

described in Algorithm 1.

In order to further investigate the validity of our results for other network topologies, in Fig. 6 we have plotted the maximum load versus number of servers for random geometric graph (RGG)³. Comparing this figure with Figure 1 shows that our proposed method performs similarly for grid and random geometric networks.

In all above simulation scenarios, we have verified the superiority of our result compared to previous schemes in terms of maximum load. Here, we compare the communication cost of the proposed scheme with the nearest replica (uncoded) scheme which has the minimum communication cost among previous schemes⁴. To this end, Fig. 7 shows the communication cost versus cache size M for a grid topology of $n = 1024$ nodes, library size $K = 100$, and uniform popularity

³Here, we consider a disk model RGG.

⁴The power of two choices scheme has always higher communication cost than the nearest replica scheme and hence is not plotted in Fig. 7 and Fig. 8.

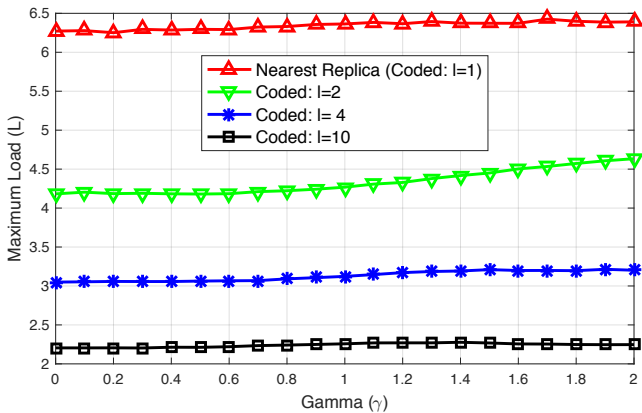


Fig. 5. Maximum load of the coded load balancing scheme versus the Zipf parameter γ for $n = 1024$ and $M = 2$. The network topology is a grid and the file popularity is uniform.

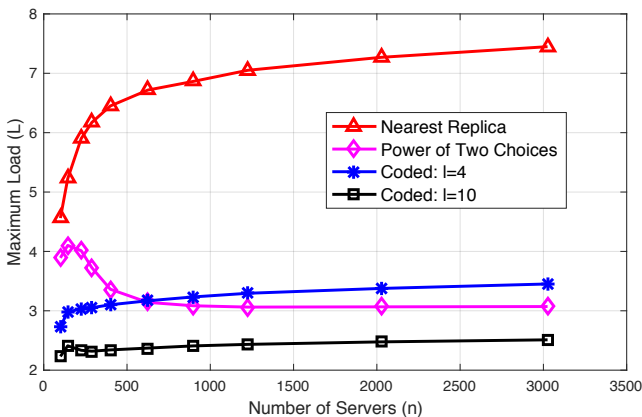


Fig. 6. Maximum load of the coded load balancing scheme versus the number of servers n for the random geometric graph (RGG). The library size is $K = 100$, cache size is $M = 2$, and the popularity profile is uniform.

profile. As it shows, the proposed coded scheme results in a slightly lower communication cost compared to the nearest replica strategy.

Finally, in Fig. 8 communication cost is plotted versus γ . As it is observed from this figure, communication cost of all the schemes show a decreasing trend as γ increases. It is interesting to note that the coded schemes' communication cost is better in the lower γ regime compared to the nearest replica scheme.

VI. DISCUSSIONS AND CONCLUDING REMARKS

We have proposed and investigated a coded cache content placement and content delivery scheme which is shown to surpass the *nearest replica strategy* and *power of two choices* baseline schemes in terms of load balancing performance. By deriving closed-form expressions for a grid network, we have shown that the proposed scheme will result in an almost perfect load balancing performance without sacrificing communication cost (Table II summarizes our results versus results of the two baseline schemes proposed in [10]). Furthermore, we have generalized the above result for a more general class of networks, including Hypercube, Random Geometric

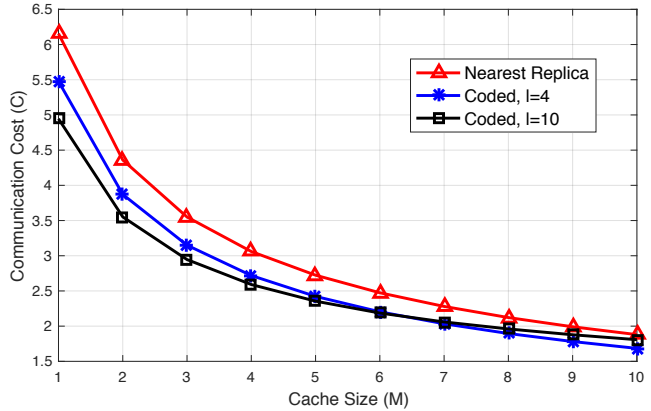


Fig. 7. Communication cost versus cache size M . The network topology is a grid of size $n = 1024$, $K = 100$, and popularity profile is uniform.

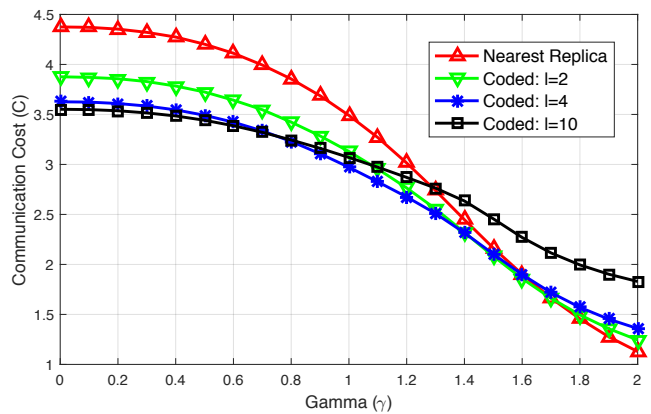


Fig. 8. Communication cost of the coded load balancing scheme versus the Zipf parameter γ for the same set of parameters stated in Fig. 5.

Graphs, and Random d -Regular Graphs. By performing extensive simulations, we have verified our theoretical findings as well as investigated the non-asymptotic performance of the proposed scheme.

Finally, here we comment on the complexity of coding/decoding of the proposed scheme. In the proposed coded caching scheme, at the end of the content delivery phase, each server has received ℓ coded chunks for each of its requests. Then, in order to recover the requested file, it has to calculate the inverse of an $\ell \times \ell = \log n \times \log n$ matrix over a finite field \mathbb{F}_q . Basically, this introduces a complexity of order $O(\ell^{2.37})$ finite field operations [35]. This may not be computationally feasible in certain practical scenarios. Thus, one may ask what are other approaches which achieve the same performance with less computational complexity.

An alternative approach is to use *Fountain-like codes*, originally proposed for packet erasure channels [12]. This coding technique benefits from a non-uniform coding operator (similar to (4) but with non-uniform distribution over the coefficients α_r). The main idea behind these codes is that by optimizing over the coding coefficients distribution, one can design the parity check matrix of these codes such that the complexity of the encoding and decoding algorithms will be significantly reduced. More specifically, in the cache content

TABLE II

THE SCALING RESULTS FOR THE PROPOSED CODING SCHEME VERSUS THE BASELINE SCHEMES OF [10] (FOR GRID TOPOLOGY WITH UNIFORM FILE POPULARITY).

	L	Regime (L)	C	Regime (C)
Nearest replica [10]	$\Theta(\log n)$	$K = n^\delta$ for $0 < \delta < 1$, $M = \Theta(1)$	$\Theta\left(\sqrt{\frac{K}{M}}\right)$	$M \ll K$
Power of two choices [10]	$\Theta(\log \log n)$	$K = n$, $M = n^\alpha$ and $r = n^\beta$ where $\alpha + 2\beta \geq 1 + 2\frac{\log \log n}{\log n}$ for $0 < \alpha, \beta < 1/2$	$\Theta(r)$	$K = n$, $M = n^\alpha$ and $r = n^\beta$ where $\alpha + 2\beta \geq 1 + 2\frac{\log \log n}{\log n}$ for $0 < \alpha, \beta < 1/2$
Coded	$O(1)$	$\ell = \Omega(\log n)$	$\Theta\left(\sqrt{\frac{K}{M}}\right)$	$M = \Theta(1)$, $K = n^\delta$, for $0 < \delta < 1$, and $\ell = \Theta(\log n)$

placement phase in Algorithm 1, one can use the encoder of a Raptor code [13], instead of a uniform operator \mathcal{L} , defined in (4). It is shown that the Raptor codes have linear encoding and decoding complexity in the codewords length [13]. Translating to our problem setting, this leads to the encoding and decoding complexity of order $\ell = \log(n)$ finite field operations, compared to the aforementioned matrix inversion.

REFERENCES

- [1] A. Su, D. R. Choffnes, A. Kuzmanovic, and F. E. Bustamante, "Drafting behind akamai: inferring network conditions based on CDN redirections," *IEEE/ACM Trans. Netw.*, vol. 17, no. 6, pp. 1752–1765, 2009.
- [2] F. Chen, R. K. Sitaraman, and M. Torres, "End-user mapping: Next generation request routing for content delivery," in *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication, SIGCOMM 2015, London, United Kingdom, August 17-21, 2015*, 2015, pp. 167–181.
- [3] S. Manfredi, F. Oliviero, and S. P. Romano, "A distributed control law for load balancing in content delivery networks," *IEEE/ACM Trans. Netw.*, vol. 21, no. 1, pp. 55–68, 2013.
- [4] V. Cardellini, E. Casalicchio, M. Colajanni, and P. S. Yu, "The state of the art in locally distributed web-server systems," *ACM Comput. Surv.*, vol. 34, no. 2, pp. 263–311, Jun. 2002.
- [5] J. Dilley, B. Maggs, J. Parikh, H. Prokop, R. Sitaraman, and B. Weihl, "Globally distributed content delivery," *IEEE Internet Computing*, vol. 6, no. 5, pp. 50–58, 2002.
- [6] E. Nygren, R. K. Sitaraman, and J. Sun, "The akamai network: A platform for high-performance internet applications," *SIGOPS Oper. Syst. Rev.*, vol. 44, no. 3, pp. 2–19, Aug. 2010.
- [7] A.-M. K. Pathan, C. Vecchiola, and R. Buyya, "Load and proximity aware request-redirection for dynamic load distribution in peering cdns," in *OTM 2008 Confederated International Conferences.*, 2008, pp. 62–81.
- [8] J. Tang, W.-P. Tay, and Y. Wen, "Dynamic request redirection and elastic service scaling in cloud-centric media networks," *IEEE Trans. Multimedia*, vol. 16, no. 5, pp. 1434–1445, 2014.
- [9] A. Pourmiri, M. J. Siavoshani, and S. P. Shariatpanahi, "Proximity-aware balanced allocations in cache networks," in *2017 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, May 2017, pp. 1068–1077.
- [10] M. J. Siavoshani, A. Pourmiri, and S. P. Shariatpanahi, "Storage, communication, and load balancing trade-off in distributed cache networks," *IEEE Transactions on Parallel and Distributed Systems*, vol. 29, no. 4, pp. 943–957, April 2018.
- [11] M. Luby, "LT codes," in *Proceedings of the 43rd Symposium on Foundations of Computer Science*, ser. FOCS '02, 2002.
- [12] J. W. Byers, M. Luby, M. Mitzenmacher, and A. Rege, "A digital fountain approach to reliable distribution of bulk data," *SIGCOMM Comput. Commun. Rev.*, vol. 28, no. 4, pp. 56–67, Oct. 1998.
- [13] A. Shokrollahi, "Raptor codes," *IEEE Trans. Information Theory*, vol. 52, no. 6, pp. 2551–2567, 2006.
- [14] C. Gkantsidis and P. R. Rodriguez, "Network coding for large scale content distribution," in *Proceedings IEEE 24th Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM 2005).*, vol. 4, Mar 2005, pp. 2235–2245.
- [15] U. Lee, J.-S. Park, J. Yeh, G. Pau, and M. Gerla, "Code torrent: Content distribution using network coding in VANET," in *Proceedings of the 1st International Workshop on Decentralized Resource Sharing in Mobile Computing and Networking*, ser. MobiShare '06, 2006, pp. 1–5.
- [16] M. Bilal and S. Kang, "Network-coding approach for information-centric networking," *IEEE Systems Journal*, vol. 13, no. 2, pp. 1376–1385, June 2019.
- [17] R. Kumar, V. Babu, and D. Nicol, "Network coding for critical infrastructure networks," in *2018 IEEE 26th International Conference on Network Protocols (ICNP)*, Sep. 2018, pp. 436–437.
- [18] T. Guven, R. J. La, M. A. Shayman, and B. Bhattacharjee, "A unified framework for multipath routing for unicast and multicast traffic," *IEEE/ACM Transactions on Networking*, vol. 16, no. 5, pp. 1038–1051, 2008.
- [19] E. Magli, M. Wang, P. Frossard, and A. Markopoulou, "Network coding meets multimedia: A review," *IEEE Transactions on Multimedia*, vol. 15, no. 5, pp. 1195–1212, Aug 2013.
- [20] K. Shanmugam, N. Golrezaei, A. G. Dimakis, A. F. Molisch, and G. Caire, "Femtocaching: Wireless content delivery through distributed caching helpers," *IEEE Transactions on Information Theory*, vol. 59, no. 12, pp. 8402–8413, Dec 2013.
- [21] F. Yin, A. Wang, D. Liu, and Z. Zhang, "Energy-aware joint user association and resource allocation for coded cache-enabled hetnets," *IEEE Access*, vol. 7, pp. 94 128–94 142, 2019.
- [22] W. Teng, M. Sheng, K. Guo, and Z. Qiu, "Distributed content replacement in small cell networks using continuous-time markov chain," in *ICC 2019 - 2019 IEEE International Conference on Communications (ICC)*, May 2019, pp. 1–6.
- [23] M. Lee, H. Feng, and A. F. Molisch, "Design of caching content replacement in base station assisted wireless d2d caching networks," in *ICC 2019 - 2019 IEEE International Conference on Communications (ICC)*, May 2019, pp. 1–7.
- [24] V. Bioglio, F. Gabry, and I. Land, "Optimizing MDS codes for caching at the edge," in *IEEE Global Communications Conference (GLOBECOM)*, 2015, pp. 1–6.
- [25] J. Liao, K. Wong, M. R. A. Khandaker, and Z. Zheng, "Optimizing cache placement for heterogeneous small cell networks," *IEEE Communications Letters*, vol. 21, no. 1, pp. 120–123, 2017.
- [26] F. Zhang, Y. Sun, Z. Hao, P. Si, R. Yang, Y. Zhang, M. Yu, and M. Yu, "A mds-based caching strategy for heterogeneous network," in *2018 10th International Conference on Communication Software and Networks (ICCSN)*, July 2018, pp. 184–188.
- [27] E. Recayte, F. Lázaro, and G. Liva, "Caching at the edge with LT codes," in *10th IEEE International Symposium on Turbo Codes & Iterative Information Processing, ISTC, Hong Kong, China, December 3-7, 2018*, 2018, pp. 1–5.
- [28] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker, "Web caching and Zipf-like distributions: evidence and implications," in *INFOCOM '99. Eighteenth Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings. IEEE*, vol. 1, Mar 1999, pp. 126–134 vol.1.
- [29] M. Cha, H. Kwak, P. Rodriguez, Y.-Y. Ahn, and S. Moon, "I tube, you tube, everybody tubes: Analyzing the world's largest user generated content video system," in *Proceedings of the 7th ACM SIGCOMM Conference on Internet Measurement*, ser. IMC '07. New York, NY, USA: ACM, 2007, pp. 1–14.
- [30] Y. Azar, A. Z. Broder, A. R. Karlin, and E. Upfal, "Balanced allocations," *SIAM J. Comput.*, vol. 29, no. 1, pp. 180–200, 1999.
- [31] M. J. Siavoshani, C. Fragouli, and S. N. Diggavi, "Subspace proper-

- ties of network coding and their applications,” *IEEE Transactions on Information Theory*, vol. 58, no. 5, pp. 2599–2619, May 2012.
- [32] J. Diaz, D. Mitsche, G. Perarnau, and X. Perez-Gimenez, “On the relation between graph distance and euclidean distance in random geometric graphs,” *Adv. in Appl. Probab.*, vol. 48, no. 3, pp. 848–864, 09 2016. [Online]. Available: <https://projecteuclid.org:443/euclid.aap/1474296318>
- [33] E. Lubetzky and A. Sly, “Cutoff phenomena for random walks on random regular graphs,” *Duke Math. J.*, vol. 153, no. 3, pp. 475–510, 06 2010. [Online]. Available: <https://doi.org/10.1215/00127094-2010-029>
- [34] M. Jafari Siavoshani, “Coded load balancing simulator,” <https://github.com/INL-Laboratory/Coded-Load-Balancing>, 2018.
- [35] D. Coppersmith and S. Winograd, “Matrix multiplication via arithmetic progressions,” *J. Symb. Comput.*, vol. 9, no. 3, pp. 251–280, Mar. 1990.