

# A Fundamental Storage-Communication Tradeoff in Distributed Computing with Straggling Nodes

Qifa Yan, Michèle Wigger, Sheng Yang, Xiaohu Tang

► To cite this version:

Qifa Yan, Michèle Wigger, Sheng Yang, Xiaohu Tang. A Fundamental Storage-Communication Tradeoff in Distributed Computing with Straggling Nodes. 2019. hal-02288594

HAL Id: hal-02288594

<https://hal.telecom-paris.fr/hal-02288594>

Preprint submitted on 15 Sep 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Fundamental Storage-Communication Tradeoff in Distributed Computing with Straggling Nodes

Qifa Yan, Michèle Wigger

LTCI, Télécom ParisTech

75013 Paris, France

Email: {qifa.yan, michele.wigger}  
@telecom-paristech.fr

Sheng Yang

L2S, CentraleSupélec

91190 Gif-sur-Yvette, France

Email: sheng.yang@centralesupelec.fr

Xiaohu Tang

Information Security and National

Computing Grid Laboratory,

Southwest Jiaotong University,

611756, Chengdu, Sichuan, China

Email: xhutang@swjtu.edu.cn

**Abstract**—The optimal storage-computation tradeoff is characterized for a MapReduce-like distributed computing system with straggling nodes, where only a part of the nodes can be utilized to compute the desired output functions. The result holds for arbitrary output functions and thus generalizes previous results that restricted to linear functions. Specifically, in this work, we propose a new information-theoretical converse and a new matching coded computing scheme, that we call *coded computing for straggling systems* (CCS).

## I. INTRODUCTION

Distributed computing has emerged as one of the most important paradigms to speed up large-scale data analysis tasks such as machine learning. A well-known computing framework, MapReduce, can deal with tasks with large data size. In such systems, the tasks are typically decomposed into computing map and reduce functions, where the map functions can be computed by different nodes across the network, and the final outputs are computed by combining the outputs of the map functions by means of reduce functions. Such architectures can be used, for example, to perform learning tasks in neural networks [1].

Recently, Li *et al.* proposed a scheme named coded distributed computing (CDC), which through clever coding reduces the communication load for data shuffling between the map and reduce phases [2]. The CDC scheme achieves the optimal *storage-communication tradeoff*, i.e., the smallest communication load under a total memory constraint. More recently, this result was extended in [3], [4] to account also for the computation load during the map phase.

In this paper, we consider a variant of this problem where each node takes a random time to compute its desired map functions. In this case, waiting that all nodes have finished their computation can be too time consuming. Instead, the data shuffling and reduce computations are started as soon as any set of  $Q$  nodes, for  $Q$  a constant in  $\{1, \dots, K\}$ , and  $K$  the number of total nodes, has terminated the map computations. One of the difficulties in such systems is that when assigning the map computations to the nodes, it is yet unclear which set of  $Q$  nodes will perform the data shuffling and compute the reduce functions. In this sense, the assignment of map computations needs to be convenient irrespective of the set of  $Q$  nodes that computes the reduce functions.

This setup with straggling nodes, hereafter referred to as the *straggling system*, has been introduced in [5]. That work used maximum separable codes (MDS) to set up coded computing. Similar codes were also used in [6], [7], and it was shown that they achieve the order optimal storage-communication tradeoff in such straggling systems when the global task is to compute *linear functions*.

Many important tasks in practice, such as computations in neural networks, are however non-linear. This motivates us to investigate the MapReduce framework with straggling nodes for general map and reduce functions. In this work, we completely characterize the optimal storage-communication tradeoff for systems with straggling nodes. This optimal tradeoff is achieved by our new scheme that we name *coded computing for straggling systems* (CCS). We also present an information-theoretical converse matching our new CCS scheme.

Related to the present work are also distributed computing schemes with straggling nodes for master-worker network models, where the worker nodes perform the map computations and the master node the reduce computations. Specifically, [8]–[10] proposed coding schemes for systems with linear map and reduce functions and [11] proposed coding schemes for calculating the derivative of likelihood functions.

*Notations:* For positive integers  $n, k$  where  $k \leq n$ , we will use the notations  $C_n^k \triangleq \frac{n!}{k!(n-k)!}$ ,  $[n] \triangleq \{1, 2, \dots, n\}$ , and  $[k : n] \triangleq \{k, k+1, \dots, n\}$ . We use sans-serif font, e.g.,  $K$ , to denote system parameters. The binary field is denoted by  $\mathbb{F}_2$  and the  $n$  dimensional vector space over  $\mathbb{F}_2$  is denoted by  $\mathbb{F}_2^n$ . The operator  $|\cdot|$  is used in the following way: for a set  $\mathcal{A}$ , we use  $|\mathcal{A}|$  to denote its cardinality, while for a signal  $X$ , we use  $|X|$  to denote its length measured in bits. The notation  $\Pr\{\cdot\}$  is used to denote probability of an event.

## II. SYSTEM MODEL

In this section, we define our model, referred to as the  $(K, Q)$  straggling system. This model is parameterized by positive integers  $K, Q, N, D, U, V, W$ . Consider a system that aims to compute  $D$  output functions through  $K$  distributed computing nodes, denoted by  $\mathcal{K} \triangleq [K]$ . Each output function  $\phi_d$  ( $d \in [D]$ ) takes all files in the library  $\mathcal{W} = \{w_n : n \in [N]\}$

as inputs, where  $w_n$  is a file of size  $W$  bits, and outputs a bit stream of length  $U$ , i.e.,

$$u_d = \phi_d(w_1, \dots, w_N) \in \mathbb{F}_2^U,$$

where  $\phi_d : \mathbb{F}_2^{NW} \rightarrow \mathbb{F}_2^U$ . Assume that the computation of each output function  $\phi_d$  can be decomposed as:

$$\phi_d(w_1, \dots, w_N) = h_d(f_{d,1}(w_1), \dots, f_{d,N}(w_N)),$$

where

- The *map function*  $f_{d,n} : \mathbb{F}_2^W \rightarrow \mathbb{F}_2^V$  maps the file  $w_n$  into a binary stream of length  $V$ , called intermediate values (IVA), i.e.,

$$v_{d,n} \triangleq f_{d,n}(w_n) \in \mathbb{F}_2^V, \quad \forall n \in [N].$$

- The *reduce function*  $h_d : \mathbb{F}_2^{NV} \rightarrow \mathbb{F}_2^U$  maps the IVAs  $\{v_{d,n}\}_{n=1}^N$  into the output stream

$$u_d = h_d(v_{d,1}, \dots, v_{d,N}).$$

Hence, the computation is carried out through three phases.

1) **Map Phase:** Each node  $k \in \mathcal{K}$  stores a subset of files  $\mathcal{M}_k \subseteq \mathcal{W}$ , and tries to compute all the IVAs from the files in  $\mathcal{M}_k$ , denoted by  $\mathcal{C}_k$ :

$$\mathcal{C}_k \triangleq \{v_{d,n} : d \in [D], w_n \in \mathcal{M}_k\}. \quad (1)$$

Each node has a random time to compute its corresponding IVAs. To limit latency of the system, the distributed computing scheme proceeds with the shuffle and reduce phases as soon as a fixed number of  $Q$  nodes has terminated the map computations. This subset of nodes will be called *active set*.

Notice that  $Q$  is a fixed system parameter. Moreover, for simplicity we assume that each subset  $\mathcal{Q} \subseteq \mathcal{K}$  of size  $|\mathcal{Q}| = Q$  is active with same probability. Define

$$\mathfrak{T}_t \triangleq \{\mathcal{T} \subseteq \mathcal{K} : |\mathcal{T}| = t\}, \quad \forall t \in [K],$$

and let  $\mathbf{Q}$  denote the active set. Then,

$$\Pr\{\mathbf{Q} = \mathcal{Q}\} = \frac{1}{C_K^Q}, \quad \forall \mathcal{Q} \in \mathfrak{T}_Q.$$

We also assume that all the files can be recovered<sup>1</sup> from any active set of size  $Q$ . Hence, for any file  $w_n \in \mathcal{W}$ , the number of nodes storing this file must satisfy  $t_n \geq K - Q + 1$ ,  $\forall n \in [N]$ , and hence  $r \geq K - Q + 1$ .

The output functions  $\phi_1, \dots, \phi_D$  are then uniformly assigned to the nodes in  $\mathbf{Q}$ . Denote the indices of the output functions assigned to a given node  $k \in \mathbf{Q}$  by  $\mathcal{D}_{\mathbf{Q},k}$ .<sup>2</sup> Thus,  $\mathbf{D}_{\mathbf{Q}} \triangleq \{\mathcal{D}_{\mathbf{Q},k}\}_{k \in \mathbf{Q}}$  forms a partition of  $[D]$ , and each  $\mathcal{D}_{\mathbf{Q},k}$  is of cardinality  $\frac{D}{Q}$ . Denote the set of all uniform partitions of  $[D]$  by  $\mathfrak{B}$ .

2) **Shuffle Phase:** The nodes in  $\mathbf{Q}$  proceed to exchange their computed IVAs. Each node  $k$  multicasts a signal

$$X_k^{\mathbf{Q}} = \varphi_{\mathbf{Q}}^{(k)}(\mathcal{C}_k, \mathbf{D}_{\mathbf{Q}})$$

<sup>1</sup>In this paper, we exclude the ‘‘outage’’ event in which some active set cannot compute the given function due to missing files.

<sup>2</sup>Here we assume for simplicity that  $Q$  divides  $D$ . Note that otherwise we can always add empty functions for the assumption to hold.

to all other nodes in  $\mathbf{Q}$ , where  $\varphi_{\mathbf{Q}}^{(k)} : \mathbb{F}_2^{|\mathcal{C}_k|V} \times \mathfrak{B} \rightarrow \mathbb{F}_2^{|\mathcal{X}_k^{\mathbf{Q}}|}$  is the encoding function of node  $k$ . Thus each active node  $k \in \mathbf{Q}$  receives the signals  $X^{\mathbf{Q}} \triangleq \{X_k^{\mathbf{Q}} : k \in \mathbf{Q}\}$  error-free.

3) **Reduce Phase:** With the received signals  $X^{\mathbf{Q}}$  exchanged during the shuffle phase and the IVAs  $\mathcal{C}_k$  computed locally during the map phase, node  $k$  restores all the IVAs

$$\{(v_{d,1}, \dots, v_{d,N})\}_{d \in \mathcal{D}_{\mathbf{Q},k}} = \psi_{\mathbf{Q}}^{(k)}(X^{\mathbf{Q}}, \mathcal{C}_k, \mathbf{D}_{\mathbf{Q}}),$$

where  $\psi_{\mathbf{Q}}^{(k)} : \mathbb{F}_2^{\sum_{k \in \mathbf{Q}} |\mathcal{X}_k^{\mathbf{Q}}|} \times \mathbb{F}_2^{|\mathcal{C}_k|V} \times \mathfrak{B} \rightarrow \mathbb{F}_2^{\frac{NDV}{Q}}$ . Subsequently, it proceeds to compute

$$u_d = h_d(v_{d,1}, \dots, v_{d,N}), \quad \forall d \in \mathcal{D}_{\mathbf{Q},k}. \quad (2)$$

*Remark 1:* Notice that a decomposition into map and reduce functions is always possible. In fact, trivially, one can set the map and reduce functions to be the identity and output functions respectively, i.e.,  $f_{d,n}(w_n) = w_n$ , and  $h_d = \phi_d$ ,  $\forall n \in [N]$ ,  $d \in [D]$ , in which case  $V = W$ . However, to mitigate the communication cost in shuffle phase, one would prefer a decomposition such that the length of the IVAs is small but suffices to compute the final outputs.

To measure the storage and computation costs, we introduce the following definitions:

*Definition 1 (Storage Space):* Storage space  $\bar{r}$  is defined as the total number of files stored across the  $K$  nodes normalized by the total number of files  $N$ , i.e.,

$$\bar{r} \triangleq \frac{\sum_{k=1}^K |\mathcal{M}_k|}{N}.$$

*Definition 2 (Communication Load):* Communication load  $\bar{L}$  is defined as the average total number of bits sent in the shuffle phase, normalized by the total number of bits of all intermediate values, i.e.,

$$\bar{L} = \mathbf{E} \left[ \frac{\sum_{k \in \mathbf{Q}} |\mathcal{X}_k^{\mathbf{Q}}|}{NDV} \right],$$

where the expectation is taken with respect to the active set  $\mathbf{Q}$ .

*Definition 3 (Optimal SC Tradeoff):* A pair of real numbers  $(r, L)$  is achievable if for any  $\epsilon > 0$ , there exist positive integers  $N, D, U, V, W$ , a storage design  $\{\mathcal{M}_k\}_{k=1}^K$  of storage space less than  $r + \epsilon$ , a set of uniform assignments of output functions  $\{\mathbf{D}_{\mathbf{Q}}\}_{\mathbf{Q} \in \mathfrak{T}_Q}$ , and a collection of encoding functions  $\{\{\varphi_{\mathbf{Q}}^{(k)}\}_{k \in \mathbf{Q}}\}_{\mathbf{Q} \in \mathfrak{T}_Q}$  with communication load less than  $L + \epsilon$ , such that all the output functions  $\phi_1, \dots, \phi_D$  can be computed successfully. For a fixed  $Q \in [K]$ , we define the optimal storage-communication (SC) tradeoff as

$$L_{K,Q}^*(r) \triangleq \inf \{L : (r, L) \text{ is achievable}\}.$$

The goal of this paper is to characterize  $L_{K,Q}^*(r)$  for  $r \in [K - Q + 1, K]$  for each  $Q \in [K]$ .

### III. MAIN RESULT

We summarize our main result in the following theorem.

*Theorem 1:* For a  $(K, Q)$  straggling system, with a given integer storage space  $r \in [K - Q + 1 : K]$ , the optimal storage-communication tradeoff is

$$L_{K,Q}^*(r) = \left(1 - \frac{r}{K}\right) \cdot \sum_{l=r+Q-K}^{\min\{r, Q-1\}} \frac{1}{l} \cdot \frac{C_r^l \cdot C_{K-r-1}^{Q-l-1}}{C_{K-1}^{Q-1}}.$$

In general, for  $r \in [K - Q + 1, K]$ ,  $L_{K,Q}^*(r)$  is given by the lower convex envelope formed by the above points.

Theorem 1 characterizes the storage-communication tradeoff for all positive integers  $Q \leq K$ . Fig. 1 shows the curves  $\{L_{K,Q}^*(r) : Q \in [K]\}$  for  $K = 10$ . Notice that, when  $Q = 1$ , the curve reduces to a single point  $(K, 0)$ , while when  $Q = K$ , the curve corresponds to the optimal tradeoff without straggling node obtained in [2]. In fact, in the case  $Q = K$ , our proposed scheme becomes the CDC scheme in [2].

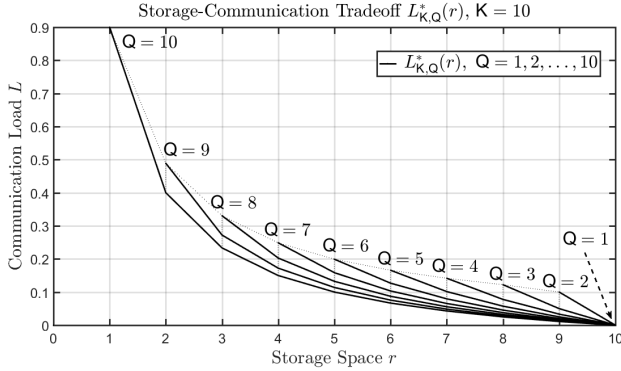


Fig. 1. Optimal Storage-Communication Tradeoff  $L_{K,Q}^*(r)$  for  $Q \in [K]$  when  $K = 10$ .

It is worth pointing out that the CCS scheme achieves the tradeoff without any assumption on the structure of map and reduce functions. Therefore, the files have to be stored in an uncoded fashion. Our converse relies on this assumption. If map and the reduce functions have certain properties, for example, linearity, a better SC tradeoff can be obtained with coded storage [6], [7]. A more detailed comparison of CCS scheme with linear map and reduce functions can be found in the long version of this paper [12].

We will see that, the storage design  $\{\mathcal{M}_k\}_{k=1}^K$  does not depend on the parameter  $Q$  but only on the available storage space  $r$  (See (3)). The proposed storage design is thus universally optimal irrespective of the size of the active set. In practice, the map phase can thus be carried out even without knowing how many nodes will be participating in the reduce phase. It turns out that, the CCS scheme can be generalized to a more general class of coded computing schemes with this property by means of a combination concept called placement delivery array [13], see [12] for details.

### IV. ACHIEVABILITY

In this section, we first present our proposed scheme, and then analyze its performance to prove the achievability part of

Theorem 1.

#### A. Coded Computing for Straggling Systems (CCS)

Let  $r$  be an integer value in  $[K - Q + 1 : K]$ . Consider  $N$  a multiple of  $C_K^r$ . Partition the files into  $C_K^r$  batches, each containing  $\eta_r \triangleq \frac{N}{C_K^r}$  files. Each batch is then associated with a subset  $\mathcal{T}$  of  $\mathcal{K}$  of cardinality  $r$ . Let  $\mathcal{W}_{\mathcal{T}}$  denote the batch of  $\eta_r$  files associated with set  $\mathcal{T}$ . Then,

$$\mathcal{W} = \{w_1, \dots, w_N\} = \bigcup_{\mathcal{T} \in \mathfrak{T}_r} \mathcal{W}_{\mathcal{T}}.$$

We now describe the map, shuffle, and reduce procedures of the CCS scheme.

1) *Map Phase:* Each node  $k$  stores batches  $\mathcal{W}_{\mathcal{T}}$  such that  $k \in \mathcal{T}$ , i.e.,

$$\mathcal{M}_k = \bigcup_{\mathcal{T}: \mathcal{T} \in \mathfrak{T}_r, k \in \mathcal{T}} \mathcal{W}_{\mathcal{T}}, \quad (3)$$

and tries to compute all IVAs in (1).

2) *Shuffle Phase:* Assume the random active set  $\mathcal{Q}$  be  $Q$ . Pick any uniform assignment of output functions  $\mathbf{D}_{\mathcal{Q}} = \{\mathcal{D}_{\mathcal{Q},k}\}_{k \in \mathcal{Q}}$ . For any subset  $\mathcal{T} \subseteq \mathcal{K}$  of size  $r$  and any  $k \in \mathcal{K}$ , let

$$\mathcal{V}_{\mathcal{T},k} \triangleq \{v_{d,n} : d \in \mathcal{D}_k, w_n \in \mathcal{W}_{\mathcal{T}}\}.$$

The shuffle phase is accomplished in  $\min\{K - Q + 1, K - r\}$  rounds, where each round is indexed by an integer  $l \in [r + Q - K : \min\{r, Q - 1\}]$ . Consider round  $l$ , where we define all subsets of  $\mathcal{K}$  of size  $r$  with exactly  $l$  active nodes:

$$\mathfrak{T}_{\mathcal{Q},r}^l \triangleq \{\mathcal{T} \subseteq \mathcal{K} : |\mathcal{T}| = r, |\mathcal{T} \cap \mathcal{Q}| = l\}.$$

For any  $j \in \mathcal{Q}$  and any  $\mathcal{T} \in \mathfrak{T}_{\mathcal{Q},r}^l$  with  $j \notin \mathcal{T}$ , evenly split the IVAs  $\mathcal{V}_{\mathcal{T},j}$  into  $l$  disjoint packets that we label as  $\{V_{\mathcal{T},j}^i : i \in \mathcal{T} \cap \mathcal{Q}\}$ . So,

$$\mathcal{V}_{\mathcal{T},j} = \{V_{\mathcal{T},j}^i : i \in \mathcal{T} \cap \mathcal{Q}\}. \quad (4)$$

Each node  $k \in \mathcal{Q}$  then multicasts the bit-wise XOR

$$X_{\mathcal{T}}^k \triangleq \bigoplus_{j \in \mathcal{T} \cap \mathcal{Q}} V_{(\mathcal{T} \setminus \{j\}) \cup \{k\}, j}^k, \quad \forall \mathcal{T} \in \mathfrak{T}_{\mathcal{Q},r}^l \text{ s.t. } k \notin \mathcal{T}. \quad (5)$$

Node  $k$  can compute all these signals because it can compute all IVAs  $\{\mathcal{V}_{\mathcal{T}',j} : \mathcal{T}' \in \mathfrak{T}_{\mathcal{Q},r}^l \text{ and } k \in \mathcal{T}'\}$ , see (3).

3) *Reduce Phase:* Any given node  $k \in \mathcal{Q}$  can recover all IVAs

$$\{\mathcal{V}_{\mathcal{T},k} : \mathcal{T} \in \mathfrak{T}_r, k \in \mathcal{T}\}$$

locally, see again (3). To obtain the missing IVAs

$$\{\mathcal{V}_{\mathcal{T},k} : \mathcal{T} \in \mathfrak{T}_r, k \notin \mathcal{T}\}, \quad (6)$$

it again proceeds in rounds  $l = r + Q - K, \dots, \min\{r, Q - 1\}$ . In round  $l$ , it forms for each  $\mathcal{T} \in \mathfrak{T}_{\mathcal{Q},r}^l$  s.t.  $k \notin \mathcal{T}$  and  $j \in \mathcal{T} \cap \mathcal{Q}$  the IVA packet

$$V_{\mathcal{T},k}^j = \bigoplus_{i \in (\mathcal{T} \setminus \{j\}) \cap \mathcal{Q}} V_{(\mathcal{T} \setminus \{i\}) \cup \{k\}, i}^j \oplus X_{(\mathcal{T} \setminus \{j\}) \cup \{k\}}^j. \quad (7)$$

Notice that,

- 1) Node  $k$  can form the signal (7) since it obtained the signal  $X_{(\mathcal{T} \setminus \{j\}) \cup \{k\}}^j$  from node  $j$ , and it can compute the IVAs  $\{V_{(\mathcal{T} \setminus \{i\}) \cup \{k\}, i}^j : i \in (\mathcal{T} \setminus \{j\}) \cap \mathcal{Q}\}$  locally.
- 2) Node  $k$  can decode all the IVAs in (6), since for any  $\mathcal{T} \in \mathfrak{T}_r$  such that  $k \notin \mathcal{T}$ ,  $r + Q - K \leq |\mathcal{T} \cap \mathcal{Q}| \leq \min\{r, Q - 1\}$ , and thus,  $\mathcal{T}$  must be an element of  $\mathfrak{T}_{\mathcal{Q}, r}^l$ , for some  $l \in [r + Q - K : \min\{r, Q - 1\}]$ .

Finally, it computes the output functions through (2).

### B. Performance Analysis

By (3), each node  $k$  stores  $C_{K-1}^{r-1}$  batches, each of which containing  $\eta_r$  files. Thus, the storage space is

$$\bar{r} = \frac{\sum_{k=1}^K |\mathcal{M}_k|}{N} = \frac{K \cdot \binom{K-1}{r-1} \cdot \eta_r}{N} = r.$$

For each  $l \in [r + Q - K : \min\{r, Q - 1\}]$ , by (4) and (5), there are in total  $(Q - l) \cdot C_{\mathcal{Q}}^l \cdot C_{K-Q}^{r-l}$  signals sent in the round  $l$ , each of size  $\frac{\eta_r}{l} \cdot \frac{D}{Q} \cdot V$  bits. Notice that, the communication load does not depend on the realization of the active set  $\mathcal{Q}$  but only on its size  $Q$ . Hence, the average communication load is

$$\begin{aligned} \bar{L} &= \sum_{l=r+Q-K}^{\min\{r, Q-1\}} \frac{(Q-l) \cdot C_{\mathcal{Q}}^l \cdot C_{K-Q}^{r-l}}{NDV} \cdot \frac{\eta_r}{l} \cdot \frac{D}{Q} \cdot V \\ &= \left(1 - \frac{r}{K}\right) \cdot \sum_{l=r+Q-K}^{\min\{r, Q-1\}} \frac{1}{l} \cdot \frac{C_r^l \cdot C_{K-r-1}^{Q-l-1}}{C_{K-1}^{Q-1}}. \end{aligned}$$

We have proved the achievability when  $r \in [K - Q + 1 : K]$ . For general  $r \in [K - Q + 1, K]$ , the lower convex envelope can be achieved with memory- and time- sharing.

### V. CONVERSE

Let  $\{\mathcal{M}_k\}_{k=1}^K$  be a storage design and  $(r, L)$  be a storage-communication pair achieved based on  $\{\mathcal{M}_k\}_{k=1}^K$ . For each  $s \in [K - Q + 1 : K]$ , define

$$a_{\mathcal{M}, s} \triangleq \sum_{\mathcal{I} \subseteq \mathcal{K} : |\mathcal{I}|=s} \left| \left( \bigcap_{k \in \mathcal{I}} \mathcal{M}_k \right) \setminus \left( \bigcup_{\bar{k} \in \mathcal{K} \setminus \mathcal{I}} \mathcal{M}_{\bar{k}} \right) \right|,$$

i.e.,  $a_{\mathcal{M}, s}$  is the number of files stored  $s$  times across all the nodes. Then by definition,  $a_{\mathcal{M}, s}$  satisfies

$$a_{\mathcal{M}, s} \geq 0, \quad (8)$$

$$\sum_{s=K-Q+1}^K a_{\mathcal{M}, s} = N, \quad (9)$$

$$\sum_{s=K-Q+1}^K s a_{\mathcal{M}, s} = \bar{r}N. \quad (10)$$

For any  $\mathcal{Q} \in \mathfrak{T}_{\mathcal{Q}}$  and any  $l \in [Q]$ , define

$$b_{\mathcal{M}, l}^{\mathcal{Q}} \triangleq \sum_{\mathcal{I} \subseteq \mathcal{Q} : |\mathcal{I}|=l} \left| \left( \bigcap_{k \in \mathcal{I}} \mathcal{M}_k \right) \setminus \left( \bigcup_{\bar{k} \in \mathcal{Q} \setminus \mathcal{I}} \mathcal{M}_{\bar{k}} \right) \right|,$$

i.e.,  $b_{\mathcal{M}, l}^{\mathcal{Q}}$  is the number of files stored exactly  $l$  times in the nodes of  $\mathcal{Q}$ . Since any file that is stored  $s$  times across the all

nodes, has  $l$  occurrences in exactly  $C_s^l \cdot C_{K-s}^{Q-l}$  subsets of  $\mathfrak{T}_{\mathcal{Q}}$ , we have

$$\sum_{\mathcal{Q} \in \mathfrak{T}_{\mathcal{Q}}} b_{\mathcal{M}, l}^{\mathcal{Q}} = \sum_{s=\max\{l, K-Q+1\}}^{K-Q+1} a_{\mathcal{M}, s} \cdot C_s^l \cdot C_{K-s}^{Q-l}. \quad (11)$$

Consider the case that the active set is  $\mathcal{Q}$ , all the output functions are computed through the nodes in  $\mathcal{Q}$ . For the sub-system with computing nodes in  $\mathcal{Q}$ , under the storage  $\{\mathcal{M}_k\}_{k \in \mathcal{Q}}$ , the optimal communication load  $L_{\mathcal{M}, \mathcal{Q}}^*$  has the following lower bound by [2, Lemma 1]:

$$L_{\mathcal{M}, \mathcal{Q}}^* \geq \sum_{l=1}^Q \frac{b_{\mathcal{M}, l}^{\mathcal{Q}}}{N} \cdot \frac{Q-l}{Ql}.$$

Then the communication load

$$\begin{aligned} \bar{L} &= \mathbf{E} \left[ \frac{\sum_{k \in \mathcal{Q}} |X_k^{\mathcal{Q}}|}{NDV} \right] \\ &= \sum_{\mathcal{Q} \in \mathfrak{T}_{\mathcal{Q}}} \mathbf{E} \left[ \frac{\sum_{k \in \mathcal{Q}} |X_k^{\mathcal{Q}}|}{NDV} \middle| \mathbf{Q} = \mathcal{Q} \right] \cdot \Pr\{\mathbf{Q} = \mathcal{Q}\} \\ &\geq \sum_{\mathcal{Q} \in \mathfrak{T}_{\mathcal{Q}}} L_{\mathcal{M}, \mathcal{Q}}^* \cdot \Pr\{\mathbf{Q} = \mathcal{Q}\} \\ &\geq \frac{1}{C_K^Q} \cdot \sum_{\mathcal{Q} \in \mathfrak{T}_{\mathcal{Q}}} \sum_{l=1}^Q \frac{b_{\mathcal{M}, l}^{\mathcal{Q}}}{N} \cdot \frac{Q-l}{Ql} \\ &= \frac{1}{C_K^Q} \cdot \sum_{l=1}^Q \left( \sum_{\mathcal{Q} \in \mathfrak{T}_{\mathcal{Q}}} \frac{b_{\mathcal{M}, l}^{\mathcal{Q}}}{N} \right) \cdot \frac{Q-l}{Ql} \\ &\stackrel{(a)}{=} \sum_{l=1}^Q \left( \sum_{s=\max\{l, K-Q+1\}}^{K-Q+1} \frac{a_{\mathcal{M}, s}}{N} \cdot \frac{C_s^l C_{K-s}^{Q-l}}{C_K^Q} \right) \cdot \frac{Q-l}{Ql} \\ &= \sum_{s=K-Q+1}^K \frac{a_{\mathcal{M}, s}}{N} \cdot \sum_{l=s+Q-K}^{\min\{s, Q\}} \frac{Q-l}{Ql} \cdot \frac{C_s^l C_{K-s}^{Q-l}}{C_K^Q}, \quad (12) \end{aligned}$$

where (a) follows from (11).

Let  $Z_K^{\mathcal{Q}}(x)$  be the function defined on the interval  $[K - Q + 1, K]$  by connecting the points  $(s, Z_K^{\mathcal{Q}}(s))$  sequentially, where

$$Z_K^{\mathcal{Q}}(s) \triangleq \sum_{l=s+Q-K}^{\min\{s, Q\}} \frac{Q-l}{Ql} \cdot C_s^l \cdot C_{K-s}^{Q-l} \quad (13)$$

is defined for any  $s \in [K - Q + 1 : K]$ . In the Appendix, we will prove the following claim.

*Claim 1:* The function  $Z_K^{\mathcal{Q}}(x)$  is a convex decreasing function over  $[K - Q + 1, K]$ .

By (8) and (9),  $\frac{a_{\mathcal{M}, s}}{N}$  are coefficients between 0 and 1 whose summation is 1. Hence by (12) and Claim 1, we have

$$\begin{aligned} \bar{L} &\geq \frac{1}{C_K^Q} \cdot \sum_{s=K-Q+1}^K \frac{a_{\mathcal{M}, s}}{N} \cdot Z_K^{\mathcal{Q}}(s) \\ &\geq \frac{1}{C_K^Q} \cdot Z_K^{\mathcal{Q}} \left( \sum_{s=K-Q+1}^K \frac{s a_{\mathcal{M}, s}}{N} \right) \\ &\stackrel{(a)}{=} \frac{Z_K^{\mathcal{Q}}(\bar{r})}{C_K^Q}, \end{aligned}$$

where (a) follows from (10). Combining the facts  $L + \epsilon \geq \bar{L}$ ,  $r + \epsilon \geq \bar{r}$  for any  $\epsilon > 0$  and  $Z_K^Q(x)$  is continuous over the interval  $[K - Q + 1, K]$ , we obtain

$$L \geq \frac{Z_K^Q(r)}{C_K^Q}, \quad \forall r \in [K - Q + 1, K].$$

In particular, when  $r \in [K - Q + 1 : K]$ , we have

$$L \geq \left(1 - \frac{r}{K}\right) \cdot \sum_{l=r+Q-K}^{\min\{r, Q-1\}} \frac{1}{l} \cdot \frac{C_r^l \cdot C_{K-r-1}^{Q-l-1}}{C_{K-1}^{Q-1}}.$$

#### ACKNOWLEDGMENT

The work of Q. Yan and M. Wigger has been supported by the ERC under grant agreement 715111.

#### APPENDIX PROOF OF CLAIM 1

We prove Claim 1 by verifying that, the sequence  $\{Z_K^Q(s)\}_{s=K-Q+1}^K$  is a discrete monotonically decreasing convex sequence, i.e.,

$$Z_K^Q(s+1) - Z_K^Q(s) < 0, \quad (14)$$

$$Z_K^Q(s+1) - Z_K^Q(s) > Z_K^Q(s) - Z_K^Q(s-1). \quad (15)$$

Notice that by (13),

$$\begin{aligned} Z_K^Q(s) &= \sum_{l=s+Q-K}^{\min\{s, Q\}} \frac{C_s^l \cdot C_{K-s}^{Q-l}}{l} - \sum_{l=s+Q-K}^{\min\{s, Q\}} \frac{C_s^l \cdot C_{K-s}^{Q-l}}{Q} \\ &\stackrel{(a)}{=} \sum_{l=s+Q-K}^{\min\{s, Q\}} \frac{C_s^l \cdot C_{K-s}^{Q-l}}{l} - \frac{C_K^Q}{Q}, \end{aligned}$$

where in (a), we used the equality  $\sum_{l=s+Q-K}^{\min\{s, Q\}} C_s^l \cdot C_{K-s}^{Q-l} = C_K^Q$ . Thus the constant  $\frac{C_K^Q}{Q}$  shows up in both  $Z_K^Q(s+1)$  and  $Z_K^Q(s)$  for all  $s \in [K - Q + 1 : K - 1]$ , then

$$\begin{aligned} &Z_K^Q(s+1) - Z_K^Q(s) \\ &= \sum_{l=s+1+Q-K}^{\min\{s+1, Q\}} \frac{C_{s+1}^l \cdot C_{K-s-1}^{Q-l}}{l} - \sum_{l=s+Q-K}^{\min\{s, Q\}} \frac{C_s^l \cdot C_{K-s}^{Q-l}}{l} \\ &\stackrel{(a)}{=} \sum_{l=s+1+Q-K}^{\min\{s+1, Q\}} \frac{(C_s^l + C_s^{l-1}) \cdot C_{K-s-1}^{Q-l}}{l} \\ &\quad - \sum_{l=s+Q-K}^{\min\{s, Q\}} \frac{C_s^l \cdot (C_{K-s-1}^{Q-l} + C_{K-s-1}^{Q-l-1})}{l} \\ &= \sum_{l=s+1+Q-K}^{\min\{s+1, Q\}} \frac{C_s^{l-1} \cdot C_{K-s-1}^{Q-l}}{l} - \sum_{l=s+Q-K}^{\min\{s, Q-1\}} \frac{C_s^l \cdot C_{K-s-1}^{Q-l-1}}{l} \\ &\stackrel{(b)}{=} \sum_{l=s+Q-K}^{\min\{s, Q-1\}} \frac{C_s^l \cdot C_{K-s-1}^{Q-l-1}}{l+1} - \sum_{l=s+Q-K}^{\min\{s, Q-1\}} \frac{C_s^l \cdot C_{K-s-1}^{Q-l-1}}{l} \\ &= - \sum_{l=s+Q-K}^{\min\{s, Q-1\}} \frac{C_s^l \cdot C_{K-s-1}^{Q-l-1}}{l(l+1)} \\ &< 0, \end{aligned} \quad (16)$$

where in (a), we applied the equality  $C_{n+1}^{m+1} = C_n^{m+1} + C_n^m$  and in (b), we used the variable change  $l' = l - 1$  in the first summation.

Moreover, following similar steps as above, by (16), for  $s \in [K - Q + 2 : K - 1]$ , we can prove

$$\begin{aligned} &\left(Z_K^Q(s+1) - Z_K^Q(s)\right) - \left(Z_K^Q(s) - Z_K^Q(s-1)\right) \\ &= \sum_{l=Q+s-1-K}^{\min\{s, Q\}-1} \frac{C_{s-1}^l \cdot C_{K-s}^{Q-l-1}}{l(l+1)} - \sum_{l=Q+s-K}^{\min\{s, Q-1\}} \frac{C_s^l \cdot C_{K-s-1}^{Q-l-1}}{l(l+1)} \\ &= \sum_{l=s+Q-K}^{\min\{s, Q-1\}} \frac{2C_{s-1}^{l-1} \cdot C_{K-s-1}^{Q-l-1}}{(l-1)l(l+1)} \\ &> 0. \end{aligned}$$

We have proved (14) and (15). As a result, the piecewise linear function  $Z_K^Q(x)$  connecting points  $\{(s, Z_K^Q(s))\}_{s=K-Q+1}^K$  sequentially is a convex function.

#### REFERENCES

- [1] Y. Liu, J. Yang, Y. Huang, L. Xu, S. Li, and M. Qi, "MapReduce based parallel neural networks in enabling large scale machine learning," *Comput. Intell. Neurosci.*, 2015.
- [2] S. Li, M. A. Maddah-Ali, Q. Yu, and A. S. Avestimehr, "A fundamental tradeoff between computation and communication in distributed computing," *IEEE Trans. Inf. Theory*, vol. 64, no. 1, pp. 109–128, Jan. 2018.
- [3] Q. Yan, S. Yang, and M. Wigger, "Storage, computation, and communication: A fundamental tradeoff in distributed computing," arXiv:1806.07565
- [4] Y. H. Ezzeldin, M. Karmoose, and C. Fragouli, "Communication vs distributed computation: An alternative trade-off curve," in *Proc. IEEE Inf. Theory Workshop (ITW)*, Kaohsiung, Taiwan, pp. 279–283, Nov. 2017.
- [5] K. Lee, M. Lam, R. Pedarsani, D. Papailiopoulos, and K. Ramchandran, "Speeding up distributed machine learning using codes," *IEEE Trans. Inf. Theory*, vol. 64, no. 3, pp. 1514–1529, Mar. 2018.
- [6] S. Li, M. A. Maddah-Ali, and A. S. Avestimehr, "A unified coding framework for distributed computing with straggling servers," in *Proc. IEEE Globecom Works (GC Wkshps)*, Washington, DC, USA, Dec. 2016.
- [7] J. Zhang and O. Simeone, "Improved latency-communication trade-off for map-shuffle-reduce systems with stragglers," arXiv:1808.06583
- [8] Q. Yu, M. Maddah-Ali, and S. Avestimehr, "Polynomial codes: An optimal design for high-dimensional coded matrix multiplication," in *Proc. The 31st Annual Conf. Neural Inf. Processing System (NIPS)*, Long Beach, CA, USA, May 2017.
- [9] Q. Yu, M. Maddah-Ali, and S. Avestimehr, "Straggler mitigation in distributed matrix multiplication: Fundamental limits and optimal coding," in *Proc. IEEE Int. Symp. Inf. Theory*, Vail, CO, USA, pp. 2022–2026, Jun. 2018.
- [10] T. Baharav, K. Lee, O. Ocal, and K. Ramchandran, "Straggler-proofing massive-scale distributed matrix multiplication with  $d$ -dimensional product codes," in *Proc. IEEE Int. Symp. Inf. Theory*, Vail, CO, USA, pp. 1993–1997, Jun. 2018.
- [11] N. Raviv, I. Tamo, R. Tandon, and A. G. Dimakis, "Gradient coding from cyclic MDS codes and expander graphs," arXiv:1707.03858.
- [12] Q. Yan, M. Wigger, S. Yang, and X. Tang, "A fundamental storage-communication tradeoff in distributed computing with straggling nodes," arXiv:1901.07793.
- [13] Q. Yan, M. Cheng, X. Tang, and Q. Chen, "On the placement delivery array design for centralized coded caching scheme," *IEEE Trans. Inf. Theory*, vol. 63, no. 9, pp. 5821–5833, Sep. 2017.