



**HAL**  
open science

## Fundamental Limits of Decentralized Data Shuffling

Kai Wan, Daniela Tuninetti, Mingyue Ji, Giuseppe Caire, Pablo Piantanida

► **To cite this version:**

Kai Wan, Daniela Tuninetti, Mingyue Ji, Giuseppe Caire, Pablo Piantanida. Fundamental Limits of Decentralized Data Shuffling. IEEE Transactions on Information Theory, 2020, 66 (6), pp.3616-3637. <10.1109/tit.2020.2966197>. <hal-02951653v4>

**HAL Id: hal-02951653**

**<https://centralesupelec.hal.science/hal-02951653v4>**

Submitted on 22 Jun 2023

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.



HAL Authorization

# Fundamental Limits of Decentralized Data Shuffling

Kai Wan, *Member, IEEE*, Daniela Tuninetti, *Senior Member, IEEE*, Mingyue Ji, *Member, IEEE*, Giuseppe Caire, *Fellow, IEEE*, and Pablo Piantanida, *Senior Member, IEEE*

**Abstract**—Data shuffling of training data among different computing nodes (workers) has been identified as a core element to improve the statistical performance of modern large-scale machine learning algorithms. Data shuffling is often considered as one of the most significant bottlenecks in such systems due to the heavy communication load. Under a master-worker architecture (where a master has access to the entire dataset and only communication between the master and the workers is allowed) coding has been recently proved to considerably reduce the communication load. This work considers a different communication paradigm referred to as *decentralized data shuffling*, where workers are allowed to communicate with one another via a shared link. The decentralized data shuffling problem has two phases: workers communicate with each other during the *data shuffling* phase, and then workers update their stored content during the *storage* phase. The main challenge is to derive novel converse bounds and achievable schemes for decentralized data shuffling by considering the asymmetry of the workers' storages (i.e., workers are constrained to store different files in their storages based on the problem setting), in order to characterize the fundamental limits of this problem.

For the case of uncoded storage (i.e., each worker directly stores a subset of bits of the dataset), this paper proposes converse and achievable bounds (based on distributed interference alignment and distributed clique-covering strategies) that are within a factor of  $3/2$  of one another. The proposed schemes are also exactly optimal under the constraint of uncoded storage for either large storage size or at most four workers in the system.

**Index Terms**—Decentralized Data shuffling, uncoded storage, distributed clique covering.

## I. INTRODUCTION

RECENT years have witnessed the emergence of big data and machine learning with wide applications in both business and consumer worlds. To cope with such a large size/dimension of data and the complexity of machine learning

A short version of this paper was presented the 56th Annual Allerton Conference (2018) on Communication, Control, and Computing in Monticello, USA.

K. Wan and G. Caire are with the Electrical Engineering and Computer Science Department, Technische Universität Berlin, 10587 Berlin, Germany (e-mail: kai.wan@tu-berlin.de; caire@tu-berlin.de). The work of K. Wan and G. Caire was partially funded by the European Research Council under the ERC Advanced Grant N. 789190, CARENET.

D. Tuninetti is with the Electrical and Computer Engineering Department, University of Illinois at Chicago, Chicago, IL 60607, USA (e-mail: danielat@uic.edu). The work of D. Tuninetti was supported in parts by NSF 1527059 and NSF 1910309.

M. Ji is with the Electrical and Computer Engineering Department, University of Utah, Salt Lake City, UT 84112, USA (e-mail: mingyue.ji@utah.edu). The work of M. Ji was supported by NSF 1817154 and NSF 1824558.

P. Piantanida is with CentraleSupélec–French National Center for Scientific Research (CNRS)–Université Paris-Sud, 91192 Gif-sur-Yvette, France, and with Montreal Institute for Learning Algorithms (MILA) at Université de Montréal, QC H3T 1N8, Canada (e-mail: pablo.piantanida@centralesupelec.fr). This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 797805.

algorithms, it is increasingly popular to use distributed computing platforms such as Amazon Web Services Cloud, Google Cloud, and Microsoft Azure services, where large-scale distributed machine learning algorithms can be implemented. The approach of data shuffling has been identified as one of the core elements to improve the statistical performance of modern large-scale machine learning algorithms [1], [2]. In particular, data shuffling consists of re-shuffling the training data among all computing nodes (workers) once every few iterations, according to some given learning algorithms. However, due to the huge communication cost, data shuffling may become one of the main system bottlenecks.

To tackle this communication bottleneck problem, under a master-worker setup where the master has access to the entire dataset, coded data shuffling has been recently proposed to significantly reduce the communication load between master and workers [3]. However, when the whole dataset is stored across the workers, data shuffling can be implemented in a distributed fashion by allowing direct communication between the workers<sup>1</sup>. In this way, the communication bottleneck between a master and the workers can be considerably alleviated. This can be advantageous if the transmission capacity among workers is much higher than that between the master and workers, and the communication load between this two setups are similar.

In this work, we consider such a *decentralized data shuffling* framework, where workers, connected by the same communication bus (common shared link), are allowed to communicate<sup>2</sup>. Although a master node may be present for the initial data distribution and/or for collecting the results of the training phase in a machine learning application, it is not involved in the data shuffling process which is entirely managed by the worker nodes in a distributed manner. In the following, we will review the literature of coded data shuffling (which we shall refer to as centralized data shuffling) and introduce the decentralized data shuffling framework studied in this paper.

### A. Centralized Data Shuffling

The coded data shuffling problem was originally proposed in [3] in a *master-worker centralized model* as illustrated in Fig. 1a. In this setup, a master, with the access to the whole

<sup>1</sup>In practice, workers communicate with each other as described in [1].

<sup>2</sup>Notice that putting all nodes on the same bus (typical terminology in Compute Science) is very common and practically relevant since this is what happens for example with Ethernet, or with the Peripheral Component Interconnect Express (PCI Express) bus inside a multi-core computer, where all cores share a common bus for intercommunication. The access of such bus is regulated by some collision avoidance protocol such as Carrier Sense Multiple Access (CSMA) [4] or Token ring [5], such that once one node talks at a time, and all other listen. Therefore, this architecture is relevant in practice.

dataset containing  $N$  data units, is connected to  $K = N/q$  workers, where  $q := N/K$  is a positive integer. Each shuffling epoch is divided into *data shuffling* and *storage update* phases. In the data shuffling phase, a subset of the data units is assigned to each worker and each worker must recover these data units from the broadcasted packets of the master and its own stored content from the previous epoch. In the storage update phase, each worker must store the newly assigned data units and, in addition, some information about other data units that can be retrieved from the storage content and master transmission in the current epoch. Such additional information should be strategically designed in order to help the coded delivery of the required data units in the following epochs. Each worker can store up to  $M$  data units in its local storage. If each worker directly copies some bits of the data units in its storage, the storage update phase is said to be *uncoded*. On the other hand, if the workers store functions (e.g., linear combinations) of the data objects' bits, the storage update is said to be *coded*. The goal is, for a given  $(M, N, q)$ , to find the best two-phase strategy that minimizes the communication load during the data shuffling phase regardless of the shuffle.

The scheme proposed in [3] uses a random uncoded storage (to fill users' extra memories independently when  $M > q$ ) and a coded multicast transmission from the master to the workers, and yields a gain of a factor of  $O(K)$  in terms of communication load with respect to the naive scheme for which the master simply broadcasts the missing, but required data bits to the workers.

The centralized coded data shuffling scheme with coordinated (i.e., deterministic) uncoded storage update phase was originally proposed in [6], [7], in order to minimize the worst-case communication load  $R$  among all the possible shuffles, i.e.,  $R$  is smallest possible such that any shuffle can be realized. The proposed schemes in [6], [7] are optimal under the constraint of uncoded storage for the cases where there is no extra storage for each worker (i.e.,  $M = q$ ) or there are less than or equal to three workers in the systems. Inspired by the achievable and converse bounds for the single-bottleneck-link caching problem in [8]–[10], the authors in [11] then proposed a general coded data shuffling scheme, which was shown to be order optimal to within a factor of 2 under the constraint of uncoded storage. Also in [11], the authors improved the performance of the general coded shuffling scheme by introducing an aligned coded delivery, which was shown to be optimal under the constraint of uncoded storage for either  $M = q$  or  $M \geq (K - 2)q$ .

Recently, inspired by the improved data shuffling scheme in [11], the authors in [12] proposed a linear coding scheme based on interference alignment, which achieves the optimal worst-case communication load under the constraint of uncoded storage for all system parameters. In addition, under the constraint of uncoded storage, the proposed coded data shuffling scheme in [12] was shown to be optimal for any shuffles (not just for the worst-case shuffles) when  $q = 1$ .

## B. Decentralized Data Shuffling

An important limitation of the centralized framework is the assumption that workers can only receive packets from the

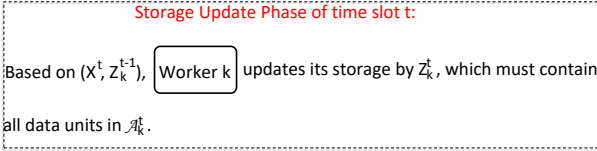
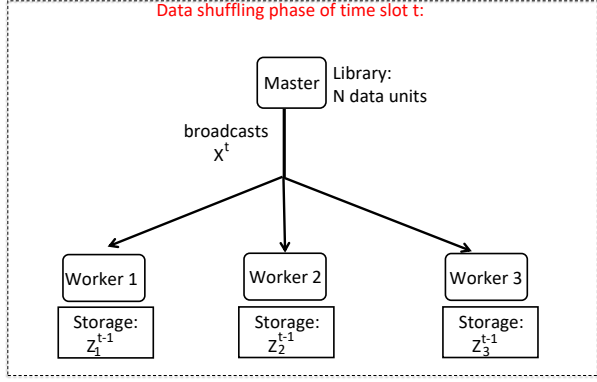
master. Since the entire dataset is stored in a decentralized fashion across the workers at each epoch of the distributed learning algorithm, the master may not be needed in the data shuffling phase if workers can communicate with each other (e.g., [1]). In addition, the communication among workers can be much more efficient compared to the communication from the master node to the workers [1], e.g., the connection between the master node and workers is via a single-ported interface, where only one message can be passed for a given time/frequency slot. In this paper, we propose the *decentralized data shuffling* problem as illustrated in Fig. 1b, where only communications among workers are allowed during the shuffling phase. This means that in the data shuffling phase, each worker broadcasts well designed coded packets (i.e., representations of the data) based on its stored content in the previous epoch. Workers take turn in transmitting, and transmissions are received error-free by all other workers through the common communication bus. The objective is to design the data shuffling and storage update phases in order to minimize the total communication load across all the workers in the worst-case shuffling scenario.

*Importance of decentralized data shuffling in practice:* In order to make the decentralized topology work in practice, we need to firstly guarantee that all the data units are already stored across the nodes so that the communication among computing nodes is sufficient. This condition is automatically satisfied from the definition of the decentralized data shuffling problem. Although the decentralized coded data shuffling incurs a larger load compared to its centralized counterpart, in practice, we may prefer the decentralized coded shuffling framework. This is due to the fact that the transmission delay/latency of the data transmission in real distributed computing system may depend on other system properties besides the total communication load, and the decentralized topology may achieve a better transmission delay/latency. This could be due to that 1) the connection between the master node and the worker clusters is normally via a single-ported interference, where only one message can be transmitted per time/frequency slot [1]; 2) computing nodes are normally connected (e.g., via grid, or ring topologies) and the link bandwidth is generally much faster, in addition, computing nodes can transmit in parallel.

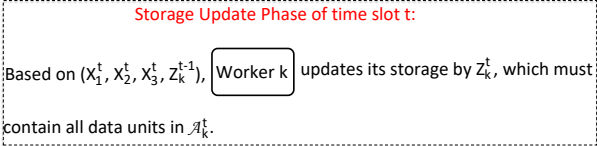
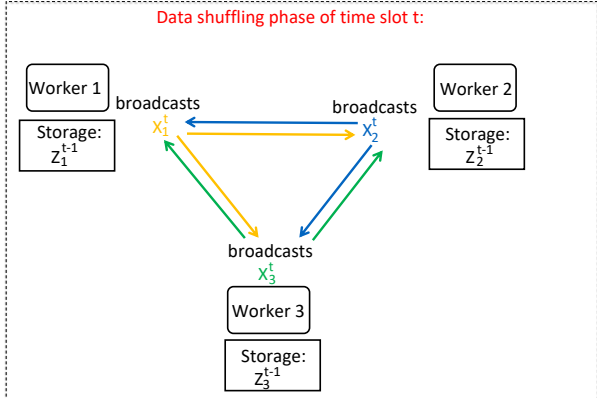
## C. Relation to Device-to-device (D2D) Caching and Distributed Computing

The coded decentralized data shuffling problem considered in this paper is related to the *coded device-to-device (D2D) caching problem* [13] and the *coded distributed computing problem* [14] – see also Remark 1 next.

The coded caching problem was originally proposed in [8] for a shared-link broadcast model. The authors in [13] extended the coded caching model to the case of D2D networks under the so-called protocol model. By choosing the communication radius of the protocol model such that each node can broadcast messages to all other nodes in the network, the delivery phase of D2D coded caching is resemblant (as far as the topology of communication between the nodes is



(a) Centralized data shuffling.



(b) Decentralized data shuffling.

Fig. 1: The system models of the 3-worker centralized and decentralized data shuffling problems in time slot  $t$ . The data units in  $\mathcal{A}_k^t$  are assigned to worker  $k$ , where  $k \in \{1, 2, 3\}$  at time  $t$ .

concerned) to the shuffling phase of our decentralized data shuffling problem.

Recently, the scheme for coded D2D caching in [13] has been extended to the coded distributed computing problem [14], which consists of two stages named *Map* and *Reduce*. In the Map stage, workers compute a fraction of intermediate computation values using local input data according to the designed Map functions. In the Reduce stage, according to the designed Reduce functions, workers exchange among each other a set of well designed (coded) intermediate computation values, in order to compute the final output

results. The coded distributed computing problem can be seen as a coded D2D caching problem under the constraint of uncoded and symmetric cache placement, where the symmetry means that each worker uses the same cache function for each file. A converse bound was proposed in [14] to show that the proposed coded distributed computing scheme is optimal in terms of communication load. This coded distributed computing framework was extended to cases such as computing only necessary intermediate values [15], [16], reducing file partitions and number of output functions [16], [17], and considering random network topologies [18], random connection graphs [19], [20], stragglers [21], storage cost [22], and heterogeneous computing power, function assignment and storage space [23], [24].

Compared to coded D2D caching and coded distributed computing, the decentralized data shuffling problem differs as follows. On the one hand, an *asymmetric* constraint on the stored contents for the workers is present (because each worker must store all bits of each assigned data unit in the previous epoch, which breaks the symmetry of the stored contents across data units of the other settings). On the other hand, each worker also needs to *dynamically update its storage* based on the received packets and its own stored content in the previous epoch. Therefore the decentralized data shuffling problem over multiple data assignment epochs is indeed a dynamic system where the evolution across the epochs of the node stored content plays a key role, while in the other problems reviewed above the cache content is static and determined at a single initial placement setup phase.

#### D. Relation to Centralized, Distributed, and Embedded Index Codings

In a *distributed index coding* problem [25], [26], there are multiple senders connected to several receivers, where each sender or receiver can access to a subset of messages in the library. Each receiver demands one message and according to the users' demands and side informations, the senders cooperatively broadcast packets to all users to satisfy the users' demands. The difference in a *centralized index coding* problem [27] compared to the distributed one is that only one sender exists and this sender can access the whole library. Very recently, the authors in [28] considered a special case of distributed index coding, referred to as *embedded index coding*, where each node acts as both a sender and a receiver in the system. It was shown in [28] that a linear code for this embedded index coding problem can be obtained from a linear index code for the centralized version of the problem by doubling the communication load.

The centralized and decentralized data shuffling phases with uncoded storage are special cases of centralized and embedded index coding problems, respectively. By using the construction in [28] we could thus design a code for the decentralized data shuffling problem by using the optimal (linear) code for the centralized case [12]; this would give a decentralized data shuffling scheme with a load twice that of [12]. It will be clarified later (in Remark 2) that the proposed decentralized data shuffling schemes are strictly better than the

those derived with the construction in [28]. This is so because the construction in [28] is general, while our design is for the specific topology considered.

### E. Contributions

In this paper, we study the decentralized data shuffling problem, for which we propose converse and achievable bounds as follows.

- 1) *Novel converse bound under the constraint of uncoded storage.* Inspired by the induction method in [14, Thm.1] for the distributed computing problem, we derive a converse bound under the constraint of uncoded storage. Different from the converse bound for the distributed computing problem, in our proof we propose a novel approach to account for the additional constraint on the ‘‘asymmetric’’ stored content.
- 2) *Scheme A: General scheme for any M.* By extending the general centralized data shuffling scheme from [11] to our decentralized model, we propose a general decentralized data shuffling scheme, where the analysis holds for any system parameters.
- 3) *Scheme B: Improved scheme for  $M \geq (K - 2)q$ .* It can be seen later that Scheme A does not fully leverage the workers’ stored content. With the storage update phase inspired by the converse bound and also used in the improved centralized data shuffling scheme in [11], we propose a two-step scheme for decentralized data shuffling to improve on Scheme A. In the first step we generate multicast messages as in [8], and in the second step we encode these multicast messages by a linear code based on distributed interference alignment (see Remark 3).  
By comparing our proposed converse bound and Scheme B, we prove that Scheme B is exactly optimal under the constraint of uncoded storage for  $M \geq (K - 2)q$ . Based on this result, we can also characterize the exact optimality under the constraint of uncoded storage when the number of workers satisfies  $K \leq 4$ .
- 4) *Scheme C: Improved scheme for  $M = 2q$ .* The delivery schemes proposed in [8], [11], [13] for coded caching with a shared-link, D2D caching, and centralized data shuffling, all belong to the class of *clique-covering* method from a graph theoretic viewpoint. By a non-trivial extension from a *distributed clique-covering approach* for the two-sender distributed index coding problems [29] to our decentralized data shuffling problem for the case  $M = 2q$ , we propose a novel decentralized data shuffling scheme. The resulting scheme outperforms the previous two schemes for this specific storage size.
- 5) *Order optimality under the constraint of uncoded storage.* By combing the three proposed schemes and comparing with the proposed converse bound, we prove the order optimality of the combined scheme within a factor of  $3/2$  under the constraint of uncoded storage.

### F. Paper Organization

The rest of the paper is organized as follows. The system model and problem formulation for the decentralized data

shuffling problem are given in Section II. Results from decentralized data shuffling related to our work are compiled in Section III. Our main results are summarized in Section IV. The proof of the proposed converse bound can be found in Section V, while the analysis of the proposed achievable schemes is in Section VI. Section VII concludes the paper. The proofs of some auxiliary results can be found in the Appendix.

### G. Notation Convention

We use the following notation convention. Calligraphic symbols denote sets, bold symbols denote vectors, and sans-serif symbols denote system parameters. We use  $|\cdot|$  to represent the cardinality of a set or the length of a vector;  $[a : b] := \{a, a + 1, \dots, b\}$  and  $[n] := \{1, 2, \dots, n\}$ ;  $\oplus$  represents bit-wise XOR;  $\mathbb{N}$  denotes the set of all positive integers.

## II. SYSTEM MODEL

The  $(K, q, M)$  decentralized data shuffling problem illustrated in Fig. 1b is defined as follows. There are  $K \in \mathbb{N}$  workers, each of which is charged to process and store  $q \in \mathbb{N}$  data units from a dataset of  $N := Kq$  data units. Data units are denoted as  $(F_1, F_2, \dots, F_N)$  and each data unit is a binary vector containing  $B$  i.i.d. bits. Each worker has a local storage of  $MB$  bits, where  $q \leq M \leq Kq = N$ . The workers are interconnected through a noiseless multicast network.

The computation process occurs over  $T$  time slots/epochs. At the end of time slot  $t - 1$ ,  $t \in [T]$ , the content of the local storage of worker  $k \in [K]$  is denoted by  $Z_k^{t-1}$ ; the content of all storages is denoted by  $Z^{t-1} := (Z_1^{t-1}, Z_2^{t-1}, \dots, Z_K^{t-1})$ . At the beginning of time slot  $t \in [T]$ , the  $N$  data units are partitioned into  $K$  disjoint batches, each containing  $q$  data units. The data units indexed by  $\mathcal{A}_k^t \subseteq [N]$  are assigned to worker  $k \in [K]$  who must store them in its local storage by the end of time slot  $t \in [T]$ . The dataset partition (i.e., data shuffle) in time slot  $t \in [T]$  is denoted by  $\mathcal{A}^t = (\mathcal{A}_1^t, \mathcal{A}_2^t, \dots, \mathcal{A}_K^t)$  and must satisfy

$$|\mathcal{A}_k^t| = q, \quad \forall k \in [K], \quad (1a)$$

$$\mathcal{A}_{k_1}^t \cap \mathcal{A}_{k_2}^t = \emptyset, \quad \forall (k_1, k_2) \in [K]^2 : k_1 \neq k_2, \quad (1b)$$

$$\bigcup_{k \in [K]} \mathcal{A}_k^t = [N] \quad (\text{dataset partition}). \quad (1c)$$

If  $q = 1$ , we let  $\mathcal{A}_k^t = \{d_k^t\}$  for each  $k \in [K]$ .

We denote the worker who must store data unit  $F_i$  at the end of time slot  $t$  by  $u_i^t$ , where

$$u_i^t = k \text{ if and only if } i \in \mathcal{A}_k^t. \quad (2)$$

The following two-phase scheme allows workers to store the requested data units.

*Initialization:* We first focus on the initial time slot  $t = 0$ , where a master node broadcasts to all the workers. Given partition  $\mathcal{A}^0$ , worker  $k \in [K]$  must store all the data units  $F_i$  where  $i \in \mathcal{A}_k^0$ ; if there is excess storage, that is, if  $M > q$ , worker  $k \in [K]$  can store in its local storage parts of the data units indexed by  $[N] \setminus \mathcal{A}_k^0$ . The storage function for worker  $k \in [K]$  in time slot  $t = 0$  is denoted by  $\psi_k^0$ , where

$$Z_k^0 := \psi_k^0 \left( \mathcal{A}^0, (F_i : i \in N) \right) \quad (\text{initial storage placement}) : \quad (3a)$$

$$H\left(Z_k^0\right) \leq \text{MB}, \forall k \in [\mathsf{K}] \quad (\text{initial storage size constraint}), \quad (3b)$$

$$H\left(\left(F_i : i \in \mathcal{A}_k^0\right) \mid Z_k^0\right) = 0 \quad (\text{initial storage content constraint}). \quad (3c)$$

Notice that the storage initialization and the storage update phase (which will be described later) are without knowledge of later shuffles. In subsequent time slots  $t \in [\mathsf{T}]$ , the master is not needed and the workers communicate with one another.

*Data Shuffling Phase:* Given global knowledge of the stored content  $Z^{t-1}$  at all workers, and of the data shuffle from  $\mathcal{A}^{t-1}$  to  $\mathcal{A}^t$  (indicated as  $\mathcal{A}^{t-1} \rightarrow \mathcal{A}^t$ ) worker  $k \in [\mathsf{K}]$  broadcasts a message  $X_k^t$  to all other workers, where  $X_k^t$  is based only on the its local storage content  $Z_k^{t-1}$ , that is,

$$H\left(X_k^t \mid Z_k^{t-1}\right) = 0 \quad (\text{encoding}). \quad (4)$$

The collection of all sent messages is denoted by  $X^t := (X_1^t, X_2^t, \dots, X_{\mathsf{K}}^t)$ . Each worker  $k \in [\mathsf{K}]$  must recover all data units indexed by  $\mathcal{A}_k^t$  from the sent messages  $X^t$  and its local storage content  $Z_k^{t-1}$ , that is,

$$H\left(\left(F_i : i \in \mathcal{A}_k^t\right) \mid Z_k^{t-1}, X^t\right) = 0 \quad (\text{decoding}). \quad (5)$$

The rate  $\mathsf{K}$ -tuple  $(R_1^{\mathcal{A}^{t-1} \rightarrow \mathcal{A}^t}, \dots, R_{\mathsf{K}}^{\mathcal{A}^{t-1} \rightarrow \mathcal{A}^t})$  is said to be *feasible* if there exist delivery functions  $\phi_k^t : X_k^t = \phi_k^t(Z_k^{t-1})$  for all  $t \in [\mathsf{T}]$  and  $k \in [\mathsf{K}]$  satisfying the constraints (4) and (5), and such that

$$H\left(X_k^t\right) \leq \text{BR}_k^{\mathcal{A}^{t-1} \rightarrow \mathcal{A}^t} \quad (\text{load}). \quad (6)$$

*Storage Update Phase:* After the data shuffling phase in time slot  $t$ , we have the storage update phase in time slot  $t \in [\mathsf{T}]$ . Each worker  $k \in [\mathsf{K}]$  must update its local storage based on the sent messages  $X^t$  and its local stored content  $Z_k^{t-1}$ , that is,

$$H\left(Z_k^t \mid Z_k^{t-1}, X^t\right) = 0 \quad (\text{storage update}), \quad (7)$$

by placing in it all the recovered data units, that is,

$$H\left(\left(F_i : i \in \mathcal{A}_k^t\right) \mid Z_k^t\right) = 0, \quad (\text{stored content}). \quad (8)$$

Moreover, the local storage has limited size bounded by

$$H\left(Z_k^t\right) \leq \text{MB}, \forall k \in [\mathsf{K}], \quad (\text{storage size}). \quad (9)$$

A storage update for worker  $k \in [\mathsf{K}]$  is said to be *feasible* if there exist functions  $\psi_k^t : Z_k^t = \psi_k^t(\mathcal{A}_k^t, Z_k^{t-1}, X^t)$  for all  $t \in [\mathsf{T}]$  and  $k \in [\mathsf{K}]$  satisfying the constraints in (7), (8) and (9).

Note: if for any  $k_1, k_2 \in [\mathsf{K}]$  and  $t_1, t_2 \in [\mathsf{T}]$  we have  $\Psi_{k_1}^{t_1} \equiv \Psi_{k_2}^{t_2}$  (i.e.,  $\Psi_{k_1}^{t_1}$  is equivalent to  $\Psi_{k_2}^{t_2}$ ), the storage phase is called *structural invariant*.

*Objective:* The objective is to minimize the *worst-case total communication load*, or just load for short in the following, among all possible consecutive data shuffles, that is we aim to characterized  $R^*$  defined as

$$R^* := \lim_{\mathsf{T} \rightarrow \infty} \min_{\substack{\psi_k^{t'}, \phi_k^{t'} \\ t' \in [\mathsf{T}], k \in [\mathsf{K}]}} \max_{(\mathcal{A}^0, \dots, \mathcal{A}^{\mathsf{T}})} \left\{ \max_{t \in [\mathsf{T}]} \sum_{k \in [\mathsf{K}]} R_k^{\mathcal{A}^{t-1} \rightarrow \mathcal{A}^t} \right\}$$

$$\left. \begin{array}{l} \text{the rate } \mathsf{K}\text{-tuple and the storage are feasible} \end{array} \right\}. \quad (10)$$

The minimum load under the constraint of uncoded storage is denoted by  $R_{\text{u}}^*$ . In general,  $R_{\text{u}}^* \geq R^*$ , because the set of all general data shuffling schemes is a superset of all data shuffling schemes with uncoded storage.

**Remark 1** (Decentralized Data Shuffling vs D2D Caching). *The D2D caching problem studied in [13] differs from our setting as follows:*

- 1) *in the decentralized data shuffling problem one has the constraint on the stored content in (8) that imposes that each worker stores the whole requested files, which is not present in the D2D caching problem; and*
- 2) *in the D2D caching problem each worker fills its local cache by accessing the whole library of files, while in the decentralized data shuffling problem each worker updates its local storage based on the received packets in the current time slot and its stored content in the previous time slot as in (7).*

*Because of these differences, achievable and converse bounds for the decentralized data shuffling problem can not be obtained by trivial renaming of variables in the D2D caching problem.*  $\square$

### III. RELEVANT RESULTS FOR CENTRALIZED DATA SHUFFLING

Data shuffling was originally proposed in [3] for the centralized scenario, where communications only exists between the master and the workers, that is, the  $\mathsf{K}$  decentralized encoding conditions in (4) are replaced by  $H(X^t \mid F_1, \dots, F_N) = 0$  where  $X^t$  is broadcasted by the master to all the workers. We summarize next some key results from [11], which will be used in the following sections. We shall use the subscripts “u,cen,conv” and “u,cen,ach” for converse (conv) and achievable (ach) bounds, respectively, for the centralized problem (cen) with uncoded storage (u). We have

- 1) *Converse for centralized data shuffling:* For a  $(\mathsf{K}, q, M)$  centralized data shuffling system, the worst-case communication load under the constraint of uncoded storage is lower bounded by the lower convex envelope of the following storage-load pairs [11, Thm.2]

$$\left( \frac{M}{q} = m, \frac{R}{q} = \frac{\mathsf{K} - m}{m} \right)_{\text{u,cen,conv}}, \quad \forall m \in [\mathsf{K}]. \quad (11)$$

- 2) *Achievability for centralized data shuffling:* In [11] it was also shown that the lower convex envelope of the following storage-load pairs is achievable with uncoded storage [11, Thm.1]

$$\left( \frac{M}{q} = 1 + g \frac{\mathsf{K} - 1}{\mathsf{K}}, \frac{R}{q} = \frac{\mathsf{K} - g}{g + 1} \right)_{\text{u,cen,ach}}, \quad \forall g \in [0 : \mathsf{K}]. \quad (12)$$

The achievable bound in (12) was shown to be within a factor  $\frac{\mathsf{K}}{\mathsf{K}-1} \leq 2$  of the converse bound in (11) under the constraint of uncoded storage [11, Thm.3].

3) *Optimality for centralized data shuffling*: It was shown in [12, Thm.4] that the converse bound in (11) can be achieved by a scheme that uses linear network coding and interference alignment/elimination. An optimality result similar to [12, Thm.4] was shown in [11, Thm.4], but only for  $m \in \{1, K-2, K-1\}$ ; note that  $m = K$  is trivial.

Although the scheme that achieves the load in (12) is not optimal in general, we shall next describe its inner workings as we will generalize it to the case of decentralized data shuffling.

*Structural Invariant Data Partitioning and Storage*: Fix  $g \in [0 : K]$  and divide each data unit into  $\binom{K}{g}$  non-overlapping and equal-length sub-blocks of length  $B/\binom{K}{g}$  bits. Let each data unit be  $F_i = (G_{i,\mathcal{W}} : \mathcal{W} \subseteq [K] : |\mathcal{W}| = g)$ ,  $\forall i \in [N]$ . The storage of worker  $k \in [K]$  at the end of time slot  $t$  is as follows,<sup>3</sup>

$$\begin{aligned} Z_k^t &= \underbrace{\left( (G_{i,\mathcal{W}} : \forall \mathcal{W}, \forall i \in \mathcal{A}_k^t) \cup (G_{i,\mathcal{W}} : k \in \mathcal{W}, \forall i \in [N] \setminus \mathcal{A}_k^t) \right)}_{\text{required data units}} \cup \underbrace{\left( (G_{i,\mathcal{W}} : k \notin \mathcal{W}, \forall i \in \mathcal{A}_k^t) \cup (G_{i,\mathcal{W}} : k \in \mathcal{W}, \forall i \in [N]) \right)}_{\text{other data units}} \\ &= \underbrace{\left( (G_{i,\mathcal{W}} : k \notin \mathcal{W}, \forall i \in \mathcal{A}_k^t) \cup (G_{i,\mathcal{W}} : k \in \mathcal{W}, \forall i \in [N]) \right)}_{\text{variable part of the storage}} \cup \underbrace{\left( (G_{i,\mathcal{W}} : k \in \mathcal{W}, \forall i \in [N]) \right)}_{\text{fixed part of the storage}}. \end{aligned} \quad (13)$$

Worker  $k \in [K]$  stores all the  $\binom{K}{g}$  sub-blocks of the required  $q$  data units indexed by  $\mathcal{A}_k^t$ , and also  $\binom{K-1}{g-1}$  sub-blocks of each data unit indexed by  $[N] \setminus \mathcal{A}_k^t$  (see (13)), thus the required storage space is

$$M = q + (N - q) \frac{\binom{K-1}{g-1}}{\binom{K}{g}} = \left( 1 + g \frac{K-1}{K} \right) q. \quad (15)$$

It can be seen (see (14) and also Table I) that the storage of worker  $k \in [K]$  at time  $t \in [T]$  is partitioned in two parts: (i) the ‘‘fixed part’’ contains all the sub-blocks of all data points that have the index  $k$  in the second subscript; this part of the storage will not be changed over time; and (ii) the ‘‘variable part’’ contains all the sub-blocks of all required data points at time  $t$  that do not have the index  $k$  in the second subscript; this part of the storage will be updated over time.

*Initialization (for the achievable bound in (12))*: The master directly transmits all data units. The storage is as in (14) given  $\mathcal{A}^0$ .

*Data Shuffling Phase of time slot  $t \in [T]$  (for the achievable bound in (12))*: After the end of storage update phase at time  $t-1$ , the new assignment  $\mathcal{A}^t$  is revealed. For notation convenience, let

$$G'_{k,\mathcal{W}} = \left( G_{i,\mathcal{W}} : i \in \mathcal{A}_k^t \setminus \mathcal{A}_k^{t-1} \right), \quad (16)$$

for all  $k \in [K]$  and all  $\mathcal{W} \subseteq [K]$ , where  $|\mathcal{W}| = g$  and  $k \notin \mathcal{W}$ . Note that in (16) we have  $|G'_{k,\mathcal{W}}| \leq B \frac{q}{\binom{K}{g}}$ , with equality (i.e.,

<sup>3</sup> Notice that here each sub-block  $G_{i,\mathcal{W}}$  is stored by workers  $\{u_i^t\} \cup \mathcal{W}$ . In addition, later in our proofs of the converse bound and proposed achievable schemes for decentralized data shuffling, the notation  $F_{i,\mathcal{W}}$  denotes the sub-block of  $F_i$ , which is stored by workers in  $\mathcal{W}$ .

worst-case scenario) if and only if  $\mathcal{A}_k^t \cap \mathcal{A}_k^{t-1} = \emptyset$ . To allow the workers to recover their missing sub-blocks, the central server broadcasts  $X^t$  defined as

$$X^t = (W_{\mathcal{J}}^t : \mathcal{J} \subseteq [K] : |\mathcal{J}| = g+1), \quad (17)$$

$$\text{where } W_{\mathcal{J}}^t = \bigoplus_{k \in \mathcal{J}} G'_{k,\mathcal{J} \setminus \{k\}}, \quad (18)$$

where in the multicast message  $W_{\mathcal{J}}^t$  in (18) the sub-blocks  $G'_{k,\mathcal{W}}$  involved in the sum are zero-padded to meet the length of the longest one. Since worker  $k \in \mathcal{J}$  requests  $G'_{k,\mathcal{J} \setminus \{k\}}$  and has stored all the remaining sub-blocks in  $W_{\mathcal{J}}^t$  defined in (18), it can recover  $G'_{k,\mathcal{J} \setminus \{k\}}$  from  $W_{\mathcal{J}}^t$ , and thus all its missing sub-blocks from  $X^t$ .

*Storage Update Phase of time slot  $t \in [T]$  (for the achievable bound in (12))*: Worker  $k \in [K]$  evicts from the (variable part of its) storage the sub-blocks  $(G_{i,\mathcal{W}} : k \notin \mathcal{W}, \forall i \in \mathcal{A}_k^{t-1} \setminus \mathcal{A}_k^t)$  and replaces them with the sub-blocks  $(G_{i,\mathcal{W}} : k \notin \mathcal{W}, \forall i \in \mathcal{A}_k^t \setminus \mathcal{A}_k^{t-1})$ . This procedure maintains the structural invariant storage structure of the storage in (14).

*Performance Analysis (for the achievable bound in (12))*: The total worst-case communication load satisfies

$$R \leq q \frac{\binom{K}{g+1}}{\binom{K}{g}} = q \frac{K-g}{g+1}, \quad (19)$$

with equality (i.e., worst-case scenario) if and only if  $\mathcal{A}_k^t \cap \mathcal{A}_k^{t-1} = \emptyset$  for all  $k \in [K]$ .

## IV. MAIN RESULTS

In this section, we summarize our main results for the decentralized data shuffling problem. We shall use the subscripts ‘‘u,dec,conv’’ and ‘‘u,dec,ach’’ for converse (conv) and achievable (ach) bounds, respectively, for the decentralized problem (dec) with uncoded storage (u). We have:

1) *Converse*: We start with a converse bound for the decentralized data shuffling problem under the constraint of uncoded storage.

**Theorem 1** (Converse). *For a  $(K, q, M)$  decentralized data shuffling system, the worst-case load under the constraint of uncoded storage is lower bounded by the lower convex envelope of the following storage-load pairs*

$$\left( \frac{M}{q} = m, \frac{R}{q} = \frac{K-m}{m} \frac{K}{K-1} \right)_{\text{u,dec,conv}}, \quad \forall m \in [K]. \quad (20)$$

Notice that the proposed converse bound is a piecewise linear curve with the corner points in (20) and these corner points are successively convex.

The proof of Theorem 1 can be found in Section V and is inspired by the induction method proposed in [14, Thm.1] for the distributed computing problem. However, there are two main differences in our proof compared to [14, Thm.1]: (i) we need to account for the additional constraint on the stored content in (8), (ii) our storage update phase is by problem definition in (8) asymmetric across data units, while it is symmetric in the distributed computing problem.

TABLE I: Example of file partitioning and storage in (14) at the end of time slot  $t$  for the decentralized data shuffling problem with  $(K, q, M) = (3, 1, 7/3)$  and  $\mathcal{A}^t = (3, 1, 2)$  where  $g = 2$ .

Workers	Sub-blocks of $F_1$	Sub-blocks of $F_2$	Sub-blocks of $F_3$
Worker 1 stores	$G_{1,\{1,2\}}, G_{1,\{1,3\}}$	$G_{2,\{1,2\}}, G_{2,\{1,3\}}$	$G_{3,\{1,2\}}, G_{3,\{1,3\}}, G_{3,\{2,3\}}$
Worker 2 stores	$G_{1,\{1,2\}}, G_{1,\{1,3\}}, G_{1,\{2,3\}}$	$G_{2,\{1,2\}}, G_{2,\{2,3\}}$	$G_{3,\{1,2\}}, G_{3,\{2,3\}}$
Worker 3 stores	$G_{1,\{1,3\}}, G_{1,\{2,3\}}$	$G_{2,\{1,2\}}, G_{2,\{1,3\}}, G_{2,\{2,3\}}$	$G_{3,\{1,3\}}, G_{3,\{2,3\}}$

2) *Achievability*: We next extend the centralized data shuffling scheme in Section III to our decentralized setting.

**Theorem 2** (Scheme A). *For a  $(K, q, M)$  decentralized data shuffling system, the worst-case load under the constraint of uncoded storage is upper bounded by the lower convex envelope of the following storage-load pairs*

$$\left( \frac{M}{q} = 1 + g \frac{K-1}{K}, \frac{R}{q} = \frac{K-g}{g} \right)_{\text{u,dec,ach}}, \quad \forall g \in [K-1]. \quad (21)$$

$$\text{and (smallest storage)} \quad \left( \frac{M}{q} = 1, \frac{R}{q} = K \right)_{\text{u,dec,ach}}, \quad (22)$$

$$\text{and (largest storage)} \quad \left( \frac{M}{q} = K, \frac{R}{q} = 0 \right)_{\text{u,dec,ach}}. \quad (23)$$

The proof is given in Section VI-A.

A limitation of Scheme A in Theorem 2 is that, in time slot  $t \in [T]$  worker  $k \in [K]$  does not fully leverage all its stored content. We overcome this limitation by developing Scheme B described in Section VI-B.

**Theorem 3** (Scheme B). *For a  $(K, q, M)$  decentralized data shuffling system, the worst-case load under the constraint of uncoded storage for  $M \geq (K-2)q$  is upper bounded by the lower convex envelope of the following storage-load pairs*

$$\left( \frac{M}{q} = m, \frac{R}{q} = \frac{K-m}{m} \frac{K}{K-1} \right)_{\text{u,dec,ach}}, \quad \forall m \in \{K-2, K-1, K\}. \quad (24)$$

We note that Scheme B is neither a direct extension of [11, Thm.4] nor of [12, Thm.4] from the centralized to the decentralized setting. As it will become clear from the details in Section VI-B, our scheme works with a rather simple way to generate the multicast messages transmitted by the workers, and it applies to any shuffle, not just to the worst-case one. In Remark 4, we also extend this scheme for the general storage size regime.

Scheme B in Theorem 3 uses a distributed clique-covering method to generate multicast messages similar to what is done for D2D caching [8], where distributed clique cover is for the side information graph (more details in Section V-A). Each multicast message corresponds to one distributed clique and includes one linear combination of all nodes in this clique. However, due to the asymmetry of the decentralized data shuffling

problem (not present in D2D coded caching), the lengths of most distributed cliques are small and thus the multicast messages based on cliques and sent by a worker in general include only a small number of messages (i.e., small multicast gain). To overcome this limitation, the key idea of Scheme C for  $M/q = 2$  (described in Section VI-C) is to augment some of the cliques and send them in  $M/q = 2$  linear combinations.

**Theorem 4** (Scheme C). *For a  $(K, q, M)$  decentralized data shuffling system, the worst-case load under the constraint of uncoded storage for  $M/q = 2$  is upper bounded by*

$$\left( \frac{M}{q} = 2, \frac{R}{q} = \frac{2K(K-2)}{3(K-1)} \right)_{\text{u,dec,ach}}. \quad (25)$$

It will be seen later that the proposed schemes only use binary codes, and only XOR operations are needed for the decoding procedure.

Finally, we combine the proposed three schemes (by considering the one among Schemes A, B or C that attains the lowest load for each storage size).

**Corollary 1** (Combined Scheme). *For a  $(K, q, M)$  decentralized data shuffling system, the achieved storage-load tradeoff of the combined scheme is the lower convex envelope of the corner points is as follows:*

- $M = q$ . With Scheme A, the worst-case load is  $q \frac{K-m}{m} \frac{K}{K-1}$ .
- $M = 2q$ . With Scheme C, the worst-case load is  $q \frac{K-m}{m} \frac{K}{K-1} \frac{4}{3}$ .
- $M = (1 + g \frac{K-1}{K})q$  where  $g \in [2 : K-3]$ . With Scheme A, the worst-case load is  $q \frac{K-g}{g}$ .
- $M = mq$  where  $m \in [K-2 : K]$ . With Scheme B, the worst-case load is  $q \frac{K-m}{m} \frac{K}{K-1}$ .

3) *Optimality*: By comparing our achievable and converse bounds, we have the following exact optimality results.

**Theorem 5** (Exact Optimality for  $M/q \geq K-2$ ). *For a  $(K, q, M)$  decentralized data shuffling system, the optimal worst-case load under the constraint of uncoded storage for  $M/q \in [K-2, K]$  is given in Theorem 1 and is attained by Scheme B in Theorem 3.*

Note that the converse bound on the load for the case  $M/q = 1$  is trivially achieved by Scheme A in Theorem 2.

From Theorem 5 we can immediately conclude the following.

**Corollary 2** (Exact Optimality for  $K \leq 4$ ). *For a  $(K, q, M)$  decentralized data shuffling system, the optimal worst-case load under the constraint of uncoded storage is given by Theorem 1 for  $K \leq 4$ .*

Finally, by combining the three proposed achievable schemes, we have the following order optimality result proved in Section VI-D.

**Theorem 6** (Order Optimality for  $K > 4$ ). *For a  $(K, q, M)$  decentralized data shuffling system under the constraint of uncoded storage, for the cases not covered by Theorem 5, the combined scheme in Corollary 1 achieves the converse bound in Theorem 1 within a factor of  $3/2$ . More precisely, when  $m q \leq M \leq (m+1)q$ , the multiplicative gap between the achievable load in Corollary 1 and the converse bound in Theorem 1 is upper bounded by*

- $4/3$ , if  $m = 1$ ;
- $1 - \frac{1}{K} + \frac{1}{2}$ , if  $m = 2$ ;
- $1 - \frac{1}{K} + \frac{1}{m-1}$ , if  $m \in [3 : K-3]$ ;
- $1$ , if  $m \in \{K-2, K-1\}$ .

4) Finally, by directly comparing the minimum load for the centralized data shuffling system (the master-worker framework) in (12) with the load achieved by the combined scheme in Corollary 1, we can quantify the communication cost of peer-to-peer operations (i.e., the multiplicative gap on the minimum worst-case load under the constraint of uncoded storage between decentralized and centralized data shufflings), which will be proved in Section VI-D.

**Corollary 3.** *For a  $(K, q, M)$  decentralized data shuffling system under the constraint of uncoded storage, the communication cost of peer-to-peer operations is no more than a factor of 2. More precisely, when  $K \leq 4$ , this cost is  $\frac{K}{K-1}$ ; when  $K \geq 5$  and  $m q \leq M \leq (m+1)q$ , this cost is upper bounded by*

- $\frac{4K}{3(K-1)}$ , if  $m = 1$ ;
- $1 + \frac{K}{2(K-1)}$ , if  $m = 2$ ;
- $1 + \frac{K}{(m-1)(K-1)}$ , if  $m \in [3 : K-3]$ ;
- $\frac{K}{K-1}$ , if  $m \in \{K-2, K-1\}$ .

**Remark 2** (Comparison to the direct extension from [28]). *As mentioned in Section I-D, the result in [28] guarantees that from the optimal (linear) centralized data shuffling scheme in [12] one can derive a linear scheme for the decentralized setting with twice the number of transmissions (by the construction given in [28, Proof of Theorem 4]), that is, the following storage-load corner points can be achieved,*

$$\left( \frac{M}{q} = m, \frac{R}{q} = 2 \frac{K-m}{m} \right), \forall m \in [K]. \quad (26)$$

*The multiplicative gap between the data shuffling scheme in (26) and the proposed converse bound in Theorem 1, is  $\frac{2(K-1)}{K}$ , which is close to 2 when  $K$  is large. Our proposed combined scheme in Corollary 1 does better: for  $K \leq 4$ , it exactly matches the proposed converse bound in Theorem 1, while for  $K > 4$  it is order optimal to within a factor of  $3/2$ .*

*In addition, the multiplicative gap  $\frac{2(K-1)}{K}$  is independent of the storage size  $M$ . It is shown in Theorem 6 that the multiplicative gap between the combined scheme and the converse decreases towards 1 when  $M$  increases.*

*Similar observation can be obtained for the communication cost of peer-to-peer operations. With the data shuffling scheme*

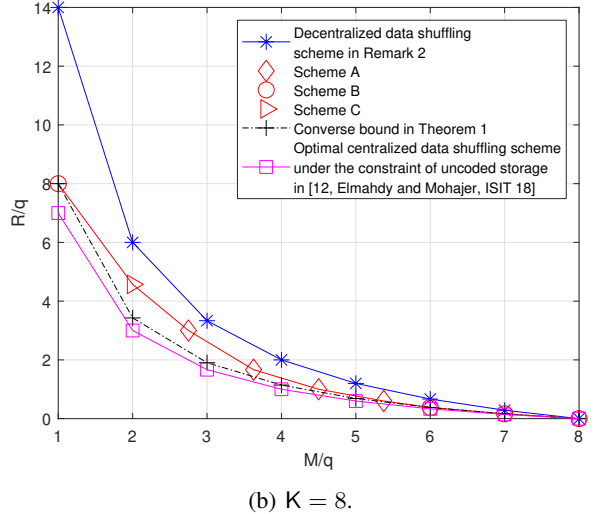
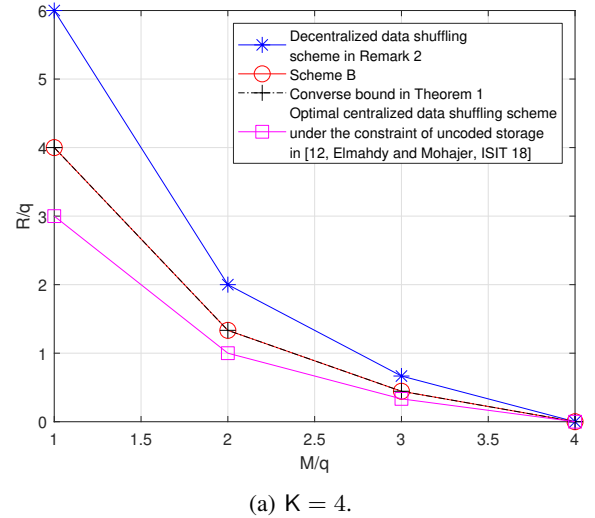


Fig. 2: The storage-load tradeoff for the decentralized data shuffling problem.

in (26), we can only prove this cost upper is bounded by 2, which is independent of  $M$  and  $K$ . With the combined scheme, it is shown in Corollary 3 that this cost decreases towards 1 when  $M$  and  $K$  increase.

We conclude this section by providing some numerical results. Fig. 2 plots our converse bound and the best convex combination of the proposed achievable bounds on the worst-case load under the constraint of uncoded storage for the decentralized data shuffling systems with  $K = 4$  (Fig. 2a) and  $K = 8$  (Fig. 2b) workers. For comparison, we also plot the achieved load by the decentralized data shuffling scheme in Remark 2, and the optimal load for the corresponding centralized system in (11) under the constraint of uncoded storage. For the case of  $K = 4$  workers, Theorem 1 is tight under the constraint of uncoded storage. For the case of  $K = 8$  workers, Scheme B meets our converse bound when  $M/q \in [6, 8]$ , and also trivially when  $M/q = 1$ .

## V. PROOF OF THEOREM 1: CONVERSE BOUND UNDER THE CONSTRAINT OF UNCODED STORAGE

We want to lower bound  $\max_{\mathcal{A}^t} \sum_{k \in [K]} R_k^{\mathcal{A}^{t-1} \rightarrow \mathcal{A}^t}$  for a fixed  $t \in [T]$  and a fixed  $\mathcal{A}^{t-1}$ . It can be also checked that this lower bound is also a lower bound on the worst-case total communication load in (10) among all  $t \in [T]$  and all possible  $(\mathcal{A}^0, \dots, \mathcal{A}^T)$ . Recall that the excess storage is said to be *uncoded* if each worker simply copies bits from the data units in its local storage. When the storage update phase is uncoded, we can divide each data unit into sub-blocks depending on the set of workers who store them.

### A. Sub-block Division of the Data Shuffling Phase under Uncoded Storage

Because of the data shuffling constraint in (1), all the bits of all data units are stored by at least one worker at the end of any time slot. Recall that the worker who must store data unit  $F_i$  at the end of time slot  $t$  is denoted by  $u_i^t$  in (2). In the case of excess storage, some bits of some files may be stored by multiple workers. We denote by  $F_{i,\mathcal{W}}$  the sub-block of bits of data unit  $F_i$  exclusively stored by workers in  $\mathcal{W}$  where  $i \in [N]$  and  $\mathcal{W} \subseteq [K]$ . By definition, at the end of step  $t-1$ , we have that  $u_i^{t-1}$  must be in  $\mathcal{W}$  for all sub-blocks  $F_{i,\mathcal{W}}$  of data unit  $F_i$ ; we also let  $F_{i,\mathcal{W}} = \emptyset$  for all  $\mathcal{W} \subseteq [K]$  if  $u_i^{t-1} \notin \mathcal{W}$ . Hence, at the end of step  $t-1$ , each data unit  $F_i$  can be written as

$$F_i = \{F_{i,\mathcal{W}} : \mathcal{W} \subseteq [K], u_i^{t-1} \in \mathcal{W}\}, \quad (27)$$

and the storage content as

$$\begin{aligned} Z_k^{t-1} &= \{F_{i,\mathcal{W}} : \mathcal{W} \subseteq [K], \{u_i^{t-1}, k\} \subseteq \mathcal{W}, i \in [N]\} \\ &= \underbrace{\{F_i : i \in \mathcal{A}_k^{t-1}\}}_{\text{required data units}} \cup \underbrace{\{F_{i,\mathcal{W}} : i \notin \mathcal{A}_k^{t-1}, \{u_i^{t-1}, k\} \subseteq \mathcal{W}\}}_{\text{other data units}}. \end{aligned} \quad (28)$$

We note that the sub-blocks  $F_{i,\mathcal{W}}$  have different content at different times (as the partition in (27) is a function of  $\mathcal{A}^{t-1}$  through  $(u_1^{t-1}, \dots, u_N^{t-1})$ ); however, in order not to clutter the notation, we will not explicitly denote the dependance of  $F_{i,\mathcal{W}}$  on time. Finally, please note that the definition of sub-block  $F_{i,\mathcal{W}}$ , as defined here for the converse bound, is not the same as  $G_{i,\mathcal{W}}$  defined in Section VI for the achievable scheme (see Footnote 3).

### B. Proof of Theorem 1

We are interested in deriving an information theoretic lower bound on the worst-case communication load. We will first obtain a number of lower bounds on the load for some carefully chosen shuffles. Since the load of any shuffle is at most as large as the worst-case shuffle, the obtained lower bounds are valid lower bounds for the worst-case load as well. We will then average the obtained lower bounds.

In particular, the shuffles are chosen as follows. Consider a permutation of  $[K]$  denoted by  $\mathbf{d} = (d_1, \dots, d_K)$  where  $d_k \neq k$  for each  $k \in [K]$  and consider the shuffle

$$\mathcal{A}_k^t = \mathcal{A}_{d_k}^{t-1}, \quad \forall k \in [K]. \quad (29)$$

We define  $X_S^t$  as the messages sent by the workers in  $S$  during time slot  $t$ , that is,

$$X_S^t := \{X_k^t : k \in S\}. \quad (30)$$

From Lemma 1 in the Appendix with  $S = [K]$ , which is the key novel contribution of our proof and that was inspired by the induction argument in [14], we have

$$\begin{aligned} \frac{R_{\text{pl}}^*}{B} &\geq H(X_{[K]}^t) \\ &\geq \sum_{m=1}^K \sum_{k \in [K]} \sum_{i \in \mathcal{A}_k^t} \sum_{\mathcal{W} \subseteq [K] \setminus \{k\}: u_i^{t-1} \in \mathcal{W}, |\mathcal{W}|=m} \frac{|F_{i,\mathcal{W}}|}{m}. \end{aligned} \quad (31)$$

To briefly illustrate the main ingredients on the derivation of (31), we provide the following example.

**Example 1** ( $K = N = 3$ ). We focus on the decentralized data shuffling problem with  $K = N = 3$ . Without loss of generality, we assume

$$\mathcal{A}_1^{t-1} = \{3\}, \quad \mathcal{A}_2^{t-1} = \{1\}, \quad \mathcal{A}_3^{t-1} = \{2\}. \quad (32)$$

Based on  $(\mathcal{A}_1^{t-1}, \mathcal{A}_2^{t-1}, \mathcal{A}_3^{t-1})$ , we can divide each data unit into sub-blocks as in (27). More precisely, we have

$$\begin{aligned} F_1 &= \{F_{1,\{2\}}, F_{1,\{1,2\}}, F_{1,\{2,3\}}, F_{1,\{1,2,3\}}\}; \\ F_2 &= \{F_{2,\{3\}}, F_{2,\{1,3\}}, F_{2,\{2,3\}}, F_{2,\{1,2,3\}}\}; \\ F_3 &= \{F_{3,\{1\}}, F_{3,\{1,2\}}, F_{3,\{1,3\}}, F_{3,\{1,2,3\}}\}. \end{aligned}$$

At the end of time slot  $t-1$ , each worker  $k \in [3]$  stores  $F_{i,\mathcal{W}}$  if  $k \in \mathcal{W}$ . Hence, we have

$$\begin{aligned} Z_1^{t-1} &= \{F_{1,\{1,2\}}, F_{1,\{1,2,3\}}, F_{2,\{1,3\}}, F_{2,\{1,2,3\}}, \\ &\quad F_{3,\{1\}}, F_{3,\{1,2\}}, F_{3,\{1,3\}}, F_{3,\{1,2,3\}}\}; \\ Z_2^{t-1} &= \{F_{1,\{2\}}, F_{1,\{1,2\}}, F_{1,\{2,3\}}, F_{1,\{1,2,3\}}, F_{2,\{2,3\}}, \\ &\quad F_{2,\{1,2,3\}}, F_{3,\{1,2\}}, F_{3,\{1,2,3\}}\}; \\ Z_3^{t-1} &= \{F_{1,\{1,3\}}, F_{1,\{1,2,3\}}, F_{2,\{3\}}, F_{2,\{1,3\}}, F_{2,\{2,3\}}, \\ &\quad F_{2,\{1,2,3\}}, F_{3,\{1,3\}}, F_{3,\{1,2,3\}}\}. \end{aligned}$$

Now we consider a permutation of  $[3]$  denoted by  $\mathbf{d} = (d_1, d_2, d_3)$  where  $d_k \neq k$  for each  $k \in [3]$  and assume the considered permutation is  $(2, 3, 1)$ . Based on  $\mathbf{d}$ , from (29), we consider the shuffle

$$\mathcal{A}_1^t = \mathcal{A}_2^{t-1} = \{1\}, \quad \mathcal{A}_2^t = \mathcal{A}_3^{t-1} = \{2\}, \quad \mathcal{A}_3^t = \mathcal{A}_1^{t-1} = \{3\}. \quad (33)$$

We first prove

$$H(X_{[2]}^t | Z_3^{t-1}, F_3) \geq |F_{1,\{2\}}|. \quad (34)$$

More precisely, by the decoding constraint in (5), we have

$$H(F_{1,\{2\}} | Z_1^{t-1}, X_{[3]}^t) = 0, \quad (35)$$

which implies

$$H(F_{1,\{2\}} | Z_1^{t-1}, X_{[3]}^t, Z_3^{t-1}, F_3) = 0. \quad (36)$$

Since  $F_{1,\{2\}}$  is not stored by workers 1 and 3, we have

$$|F_{1,\{2\}}| \leq H(F_{1,\{2\}}, X_{[2]}^t | Z_1^{t-1}, X_3^t, Z_3^{t-1}, F_3) \quad (37a)$$

$$= H(X_{[2]}^t | Z_1^{t-1}, X_3^t, Z_3^{t-1}, F_3) + H(F_{1,\{2\}} | Z_1^{t-1}, X_{[3]}^t, Z_3^{t-1}, F_3) \quad (37b)$$

$$= H(X_{[2]}^t | Z_1^{t-1}, X_3^t, Z_3^{t-1}, F_3) \quad (37c)$$

$$\leq H(X_{[2]}^t | X_3^t, Z_3^{t-1}, F_3) \quad (37d)$$

$$= H(X_{[2]}^t | Z_3^{t-1}, F_3), \quad (37e)$$

where (37c) comes from (36) and (37e) comes from the fact that  $X_3^t$  is a function of  $Z_3^{t-1}$ . Hence, we prove (34).

Similarly, we can also prove

$$H(X_{\{1,3\}}^t | Z_2^{t-1}, F_2) \geq |F_{3,\{1\}}|. \quad (38)$$

$$H(X_{\{2,3\}}^t | Z_1^{t-1}, F_1) \geq |F_{2,\{3\}}|. \quad (39)$$

In addition, we have

$$H(X_{[3]}^t) = \frac{1}{3} \sum_{k \in [3]} \left( H(X_k^t) + H(X_{[3] \setminus \{k\}}^t | X_k^t) \right) \quad (40a)$$

$$\geq \frac{1}{3} \left( H(X_{[3]}^t) + \sum_{k \in [3]} H(X_{[3] \setminus \{k\}}^t | X_k^t) \right), \quad (40b)$$

and thus

$$2H(X_{[3]}^t) \geq \sum_{k \in [3]} H(X_{[3] \setminus \{k\}}^t | X_k^t) \quad (40c)$$

$$\geq \sum_{k \in [3]} H(X_{[3] \setminus \{k\}}^t | Z_k^{t-1}) \quad (40d)$$

$$= \sum_{k \in [3]} H(X_{[3] \setminus \{k\}}^t, F_k | Z_k^{t-1}) \quad (40e)$$

$$= \sum_{k \in [3]} H(F_k | Z_k^{t-1}) + H(X_{[3] \setminus \{k\}}^t | F_k, Z_k^{t-1}), \quad (40f)$$

where (40d) comes from the fact that  $X_k^t$  is a function of  $Z_k^{t-1}$  and conditioning cannot increase entropy, and (40e) comes from the decoding constraint for worker  $k$  in (5).

Let us focus on worker 1 and we have

$$H(F_1 | Z_1^{t-1}) = |F_{1,\{2\}}| + |F_{1,\{2,3\}}|. \quad (41)$$

In addition, we have  $H(X_{\{2,3\}}^t | Z_1^{t-1}, F_1) \geq |F_{2,\{3\}}|$  from (39). Hence,

$$\begin{aligned} & H(F_1 | Z_1^{t-1}) + H(X_{\{2,3\}}^t | Z_1^{t-1}, F_1) \\ & \geq |F_{1,\{2\}}| + |F_{1,\{2,3\}}| + |F_{2,\{3\}}|. \end{aligned} \quad (42)$$

Similarly, we have

$$\begin{aligned} & H(F_2 | Z_2^{t-1}) + H(X_{\{1,3\}}^t | Z_2^{t-1}, F_2) \\ & \geq |F_{2,\{3\}}| + |F_{2,\{1,3\}}| + |F_{3,\{1\}}|; \end{aligned} \quad (43)$$

$$\begin{aligned} & H(F_3 | Z_3^{t-1}) + H(X_{\{1,2\}}^t | Z_3^{t-1}, F_3) \\ & \geq |F_{3,\{1\}}| + |F_{3,\{1,2\}}| + |F_{1,\{2\}}|. \end{aligned} \quad (44)$$

By taking (42)-(44) to (40f), we have

$$H(X_{[3]}^t) \geq \frac{1}{2} \sum_{k \in [3]} H(F_k | Z_k^{t-1}) + H(X_{[3] \setminus \{k\}}^t | F_k, Z_k^{t-1}) \quad (45a)$$

$$= |F_{1,\{2\}}| + |F_{2,\{3\}}| + |F_{3,\{1\}}| + \frac{|F_{1,\{2,3\}}|}{2} + \frac{|F_{2,\{1,3\}}|}{2}$$

$$+ \frac{|F_{3,\{1,2\}}|}{2}, \quad (45b)$$

coinciding with (31).  $\square$

We now go back to the general proof of Theorem 1. We next consider all the permutations  $\mathbf{d} = (d_1, \dots, d_K)$  of  $[K]$  where  $d_k \neq k$  for each  $k \in [K]$ , and sum together the inequalities in the form of (31). For an integer  $m \in [K]$ , by the symmetry of the problem, the sub-blocks  $F_{i,\mathcal{W}}$  where  $i \in [N]$ ,  $u_i^{t-1} \in \mathcal{W}$  and  $|\mathcal{W}| = m$  appear the same number of times in the final sum. In addition, the total number of these sub-blocks in general is  $N \binom{K-1}{m-1}$  and the total number of such sub-blocks in each inequality in the form of (31) is  $N \binom{K-2}{m-1}$ . So we obtain

$$R_u^* \geq \sum_{m=1}^K \sum_{i \in [N]} \sum_{\mathcal{W} \subseteq [K]: u_i^{t-1} \in \mathcal{W}, |\mathcal{W}|=m} \frac{\binom{K-2}{m-1}}{m \binom{K-1}{m-1}} |F_{i,\mathcal{W}}| \frac{qK}{NB} \quad (46)$$

$$= \sum_{m=1}^K qK x_m \frac{1 - (m-1)/(K-1)}{m} \quad (47)$$

$$= \sum_{m=1}^K q x_m \frac{K-m}{m} \frac{K}{K-1}, \quad (48)$$

where we defined  $x_m$  as the total number of bits in the sub-blocks stored by  $m$  workers at the end of time slot  $t-1$  normalized by the total number of bits NB, i.e.,

$$0 \leq x_m := \sum_{i \in [N]} \sum_{\mathcal{W} \subseteq [K]: u_i^{t-1} \in \mathcal{W}, |\mathcal{W}|=m} \frac{|F_{i,\mathcal{W}}|}{NB}, \quad (49)$$

which must satisfy

$$\sum_{m \in [K]} x_m = 1 \quad (\text{total size of all data units}), \quad (50)$$

$$\sum_{m \in [K]} m x_m \leq \frac{KM}{N} = \frac{M}{q} \quad (\text{total storage size}). \quad (51)$$

We then use a method based on Fourier-Motzkin elimination [30, Appendix D] for to bound  $R_u^*$  from (47) under the constraints in (50) and (51), as we did in [9] for coded caching with uncoded cache placement. In particular, for each integer  $p \in [K]$ , we multiply (50) by  $\frac{-N(2Kp-p^2+K-p)}{p(p+1)}$  to obtain

$$\frac{-N(2Kp-p^2+K-p)}{p(p+1)} = \sum_{m=1}^K \frac{-N(2Kp-p^2+K-p)}{p(p+1)} x_m, \quad (52)$$

and we multiply (51) by  $\frac{NK}{(K-1)p(p+1)}$  to have

$$\frac{NK}{(K-1)p(p+1)} \frac{KM}{N} \geq \sum_{m=1}^K \frac{NK}{(K-1)p(p+1)} m x_m. \quad (53)$$

We then add (52), (53), and (47) to obtain,

$$\begin{aligned} R_u^* & \geq \sum_{m=1}^K \frac{NK(p-m)(p+1-m)}{m(K-1)p(p+1)} x_m - \frac{NK}{(K-1)p(p+1)} \frac{KM}{N} \\ & \quad + \frac{N(2Kp-p^2+K-p)}{p(p+1)} \end{aligned} \quad (54)$$

$$\geq -\frac{NK}{(K-1)p(p+1)} \frac{KM}{N} + \frac{N(2Kp - p^2 + K - p)}{p(p+1)}. \quad (55)$$

Hence, for each integer  $p \in [K]$ , the bound in (55) becomes a linear function in  $M$ . When  $M = qp$ , from (55) we have  $R_u^* \geq \frac{N(K-p)}{(K-1)p}$ . When  $M = q(p+1)$ , from (55) we have  $R_u^* \geq \frac{N(K-p-1)}{(K-1)(p+1)}$ . In conclusion, we prove that  $R_u^*$  is lower bounded by the lower convex envelope (also referred to as ‘‘memory sharing’’) of the points  $(M = qm, R = \frac{N(K-m)}{(K-1)m})$ , where  $m \in [K]$ .

This concludes the proof of Theorem 1.

### C. Discussion

We conclude this session the following remarks:

- 1) The corner points from the converse bound are of the form  $(M/q = m, R/q = \frac{K \binom{K-2}{m-1}}{m \binom{K-1}{m-1}})$ , which may suggest the following placement.

At the end of time slot  $t-1$ , each data unit is partitioned into  $\binom{K-1}{m-1}$  equal-length sub-blocks of length  $B/\binom{K-1}{m-1}$  bits as  $F_i = (F_{i,\mathcal{W}} : \mathcal{W} \subseteq [K], |\mathcal{W}| = m, u_i^{t-1} \in \mathcal{W})$ ; by definition  $F_{i,\mathcal{W}} = \emptyset$  if either  $u_i^{t-1} \notin \mathcal{W}$  or  $|\mathcal{W}| \neq m$ . Each worker  $k \in [K]$  stores all the sub-blocks  $F_{i,\mathcal{W}}$  if  $k \in \mathcal{W}$ ; in other words, worker  $k \in [K]$  stores all the  $\binom{K-1}{m-1}$  sub-blocks of the desired data units, and  $\binom{K-2}{m-2}$  sub-blocks of the remaining data units.

In the data shuffling phase of time slot  $t$ , worker  $k \in [K]$  must decode the missing  $\binom{K-1}{m-1} - \binom{K-2}{m-2} = \binom{K-2}{m-1}$  sub-blocks of data unit  $F_j$  for all  $j \in \mathcal{A}_k^t \setminus \mathcal{A}_k^{t-1}$ . An interpretation of the converse bound is that, in the worst case, the total number of transmissions is equivalent to at least  $\frac{qK}{m} \binom{K-2}{m-1}$  sub-blocks.

We will use this interpretation to design the storage update phase our proposed Schemes B and C.

- 2) The converse bound is derived for the objective of minimizing the ‘‘sum load’’  $\sum_{k \in [K]} R_k^{A^{t-1} \rightarrow A^t}$ , see (10). The same derivation would give a converse bound for the ‘‘largest individual load’’  $\max_{k \in [K]} R_k^{A^{t-1} \rightarrow A^t}$ . In the latter case, the corner points from converse bound are of the form  $(M/q = m, R/q = \frac{\binom{K-2}{m-1}}{m \binom{K-1}{m-1}})$ . This view point may suggest that, in the worst case, all the individual loads  $R_k^{A^{t-1} \rightarrow A^t}$  are the same, i.e., the burden of communicating missing data units is equally shared by all the workers.

Our proof technique for Theorem 1 could also be directly extended to derive a converse bound on the *average load* (as opposed to the worst-case load) for all the possible shuffles in the decentralized data shuffling problem when  $N = K$ .

## VI. ACHIEVABLE SCHEMES FOR DECENTRALIZED DATA SHUFFLING

In this section, we propose three schemes for the decentralized data shuffling problem, and analyze their performances.

### A. Scheme A in Theorem 2

Scheme A extends the general centralized data shuffling scheme in Section III to the distributed model. Scheme A achieves the load in Theorem 2 for each storage size  $M = (1 + g \frac{K-1}{K})q$ , where  $g \in [K-1]$ ; the whole storage-load tradeoff piecewise curve is achieved by memory-sharing<sup>4</sup> between these points (given in (21)) and the (trivially achievable) points in (22)-(23).

*Structural Invariant Data Partitioning and Storage:* This is the same as the one in Section III for the centralized case.

*Initialization:* The master directly transmits all data units. The storage is as in (14) given  $\mathcal{A}^0$ .

*Data Shuffling Phase of time slot  $t \in [T]$ :* The data shuffling phase is inspired by the delivery in D2D caching [13]. Recall the definition of sub-block  $G'_{k,\mathcal{W}}$  in (16), where each sub-block is known by  $|\mathcal{W}| = g$  workers and needed by worker  $k$ . Partition  $G'_{k,\mathcal{W}}$  into  $g$  non-overlapping and equal-length pieces  $G'_{k,\mathcal{W}} = \{G'_{k,\mathcal{W}}(j) : j \in \mathcal{W}\}$ . Worker  $j \in \mathcal{J}$  broadcasts

$$W_{j,\mathcal{J}}^t = \bigoplus_{k \in \mathcal{J} \setminus \{j\}} G'_{k,\mathcal{J} \setminus \{k\}}(j), \quad \forall \mathcal{J} \subseteq [K] \text{ where } |\mathcal{J}| = g+1, \quad (56)$$

in other words, one linear combination  $W_{\mathcal{J}}^t$  in (17) for the centralized setting becomes  $g+1$  linear combinations  $W_{j,\mathcal{J}}^t$  in (56) for the decentralized setting, but of size reduced by a factor  $g$ . Evidently, each sub-block in  $W_{j,\mathcal{J}}^t$  is stored in the storage of worker  $j$  at the end of time slot  $t-1$ . In addition, each worker  $k \in \mathcal{J} \setminus \{j\}$  knows  $G'_{k_1,\mathcal{J} \setminus \{k_1\}}(j)$  where  $k_1 \in \mathcal{J} \setminus \{k, j\}$  such that it can recover its desired block  $G'_{k,\mathcal{J} \setminus \{k\}}(j)$ .

Since  $|G'_{k,\mathcal{W}}| \leq qB/\binom{K}{g}$ , the worst-case load is

$$R \leq q \frac{(g+1) \binom{K}{g+1}}{g \binom{K}{g}} = q \frac{K-g}{g} =: R_{\text{Ach.A}}, \quad (57)$$

as claimed in Theorem 2, where the subscript ‘‘Ach.A’’ in (57) denotes the worst-case load achieved by Scheme A.

*Storage Update Phase of time slot  $t \in [T]$ :* The storage update phase is the same as the general centralized data shuffling scheme in Section III, and thus is not repeated here.

### B. Scheme B in Theorem 3

During the data shuffling phase in time slot  $t$  of Scheme A, we treat some sub-blocks known by  $g+1$  workers as if they were only known by  $g$  workers (for example, if  $u_i^{t-1} \notin \mathcal{W}$ ,  $G_{i,\mathcal{W}}$  is stored by workers  $\{u_i^{t-1}\} \cup \mathcal{W}$ , but Scheme A treats  $G_{i,\mathcal{W}}$  as if it was only stored by workers in  $\mathcal{W}$ ), which may be suboptimal as more multicasting opportunities may be leveraged. In the following, we propose Scheme B to remedy for this shortcoming for  $M = mq$  for  $m \in \{K-2, K-1\}$ .

<sup>4</sup> Memory-sharing is an achievability technique originally proposed by Maddah-Ali and Niesen in [8] for coded caching systems, which is used to extend achievability in between discrete memory points. More precisely, focus one storage size  $M' = \alpha M_1 + (1-\alpha)M_2$ , where  $\alpha \in [0, 1]$ ,  $M_1 = (1 + g \frac{K-1}{K})q$ ,  $M_2 = (1 + (g+1) \frac{K-1}{K})q$ , and  $g \in [K-1]$ . We can divide each data unit into two non-overlapping parts, with  $\alpha B$  and  $(1-\alpha)B$  bits, respectively. The first and second parts of the  $N$  data units are stored and transmitted based on the proposed data shuffling scheme for  $M_1$  and  $M_2$ , respectively.

*Structural Invariant Data Partitioning and Storage:* Data units are partitions as inspired by the converse bound (see discussion in Section V-C), which is as in the improved centralized data shuffling scheme in [11]. Fix  $m \in [K]$ . Partition each data unit into  $\binom{K-1}{m-1}$  non-overlapping equal-length sub-block of length  $B/\binom{K-1}{m-1}$  bits. Write  $F_i = (F_{i,\mathcal{W}} : \mathcal{W} \subseteq [K], |\mathcal{W}| = m, u_i^t \in \mathcal{W})$ , and set  $F_{i,\mathcal{W}} = \emptyset$  if either  $u_i^t \notin \mathcal{W}$  or  $|\mathcal{W}| \neq m$ . The storage of worker  $k \in [K]$  at the end of time slot  $t$  is as follows,

$$Z_k^t = \left( \underbrace{(F_{i,\mathcal{W}} : i \in \mathcal{A}_k^t, \forall \mathcal{W})}_{\text{required data units}} \cup \underbrace{(F_{i,\mathcal{W}} : i \notin \mathcal{A}_k^t, k \in \mathcal{W})}_{\text{other data units}} \right), \quad (58)$$

that is, worker  $k \in [K]$  stores all the  $\binom{K-1}{m-1}$  sub-blocks of data unit  $F_i$  if  $i \in \mathcal{A}_k^t$ , and  $\binom{K-2}{m-2}$  sub-blocks of data unit  $F_j$  if  $j \notin \mathcal{A}_k^t$  (the sub-blocks stored are such that  $k \in \mathcal{W}$ ), thus the required storage space is

$$M = q + (N - q) \frac{\binom{K-2}{m-2}}{\binom{K-1}{m-1}} = q + (m-1) \frac{N - q}{K-1} = mq. \quad (59)$$

In the following, we shall see that it is possible to maintain the storage structure in (58) after the shuffling phase.

*Initialization:* The master directly transmits all data units. The storage is as in (58) given  $\mathcal{A}^0$ .

*Data Shuffling Phase of time slot  $t \in [T]$  for  $m = K-1$ :* Each data unit has been partitioned into  $\binom{K-1}{m-1} = K-1$  sub-blocks and each sub-block is stored by  $m = K-1$  workers. Similarly to Scheme A, define the set of sub-blocks needed by worker  $k \in [K]$  at time slot  $t$  and not previously stored as

$$F'_{k,[K] \setminus \{k\}} = (F_{i,\mathcal{W}} : i \in \mathcal{A}_k^t \setminus \mathcal{A}_k^{t-1}, \mathcal{W} = [K] \setminus \{k\}), \forall k \in [K]. \quad (60)$$

Since  $F'_{k,[K] \setminus \{k\}}$  (of length  $qB/(K-1)$  bits in the worst case) is desired by worker  $k$  and known by all the remaining  $m = K-1$  workers, we partition  $F'_{k,[K] \setminus \{k\}}$  into  $m = K-1$  pieces (of length  $qB/(K-1)^2$  bits in the worst case), and write  $F'_{k,[K] \setminus \{k\}} = (F'_{k,[K] \setminus \{k\}}(j) : j \in [K] \setminus \{k\})$ . Worker  $j \in [K]$  broadcasts the single linear combination (of length  $qB/(K-1)^2$  bits in the worst case) given by

$$W_j^t = \bigoplus_{k \neq j} F'_{k,[K] \setminus \{k\}}(j). \quad (61)$$

Therefore, the worst-case satisfies

$$R \leq \frac{K}{(K-1)^2} q = \frac{K-m}{m} \frac{K}{K-1} q \Big|_{m=K-1} =: R_{\text{Ach.B}}|_{M=(K-1)q} \quad (62)$$

which coincides with the converse bound.

*Storage Update Phase of time slot  $t \in [T]$  for  $m = K-1$ :* In time slot  $t-1 > 0$ , we assume that the above storage configuration of each worker  $k \in [K]$  can be done with  $Z_k^{t-2}$  and  $X_j^{t-1}$  where  $j \in [K] \setminus \{k\}$ . We will show next that at the end of time slot  $t$ , we can re-create the same configuration of storage, but with permuted data units. Thus by the induction method, we prove the above storage update phase is also structural invariant.

For each worker  $k \in [K]$  and each data unit  $F_i$  where  $i \in \mathcal{A}_k^t \setminus \mathcal{A}_k^{t-1}$ , worker  $k$  stores the whole data unit  $F_i$  in its storage. For each data unit  $F_i$  where  $i \in \mathcal{A}_k^{t-1} \setminus \mathcal{A}_k^t$ , instead of storing the whole data unit  $F_i$ , worker  $k$  only stores the bits of  $F_i$  which was stored at the end of time slot  $t-1$  by worker  $u_i^t$ . For other data units, worker  $k$  does not change the stored bits. Hence, after the storage phase in time slot  $t$ , we can re-create the same configuration of storage as the end of time slot  $t-1$  but with permuted data units.

*Data Shuffling Phase of time slot  $t \in [T]$  for  $m = K-2$ :* We partition the  $N$  data units into  $q$  groups as  $[N] = \bigcup_{i \in [q]} \mathcal{H}_i$ , where each group contains  $K$  data units, and such that for each group  $\mathcal{H}_i, i \in [q]$ , and each worker  $k \in [K]$  we have  $|\mathcal{H}_i \cap \mathcal{A}_k^t| = 1$  and  $|\mathcal{H}_i \cap \mathcal{A}_k^{t-1}| = 1$ . In other words, the partition is such that, during the data shuffling phase of time slot  $t$ , among all the  $K$  data units in each group, each worker requests exactly one data unit and knows exactly one data unit. Such a partition can always be found [12, Lemma 7]. The dependance of  $\mathcal{H}_i$  on  $t$  is not specified so not to clutter the notation. For group  $\mathcal{H}_i, i \in [q]$ , we define

$$\mathcal{U}(\mathcal{H}_i) := \{k \in [K] : \mathcal{H}_i \cap \mathcal{A}_k^t \subseteq \mathcal{A}_k^{t-1}\}, \quad \forall i \in [q], \quad (63)$$

as the set of workers in the group who already have stored the needed data point (i.e., who do not need to shuffle). Since each worker has to recover at most one data unit in each group, the delivery in each group is as if  $q = 1$ . Hence, to simplify the description, we focus on the case  $q = 1$ , in which case there is only one group and thus we simplify the notation  $\mathcal{U}(\mathcal{H}_i)$  to just  $\mathcal{U}$ . We first use the following example to illustrate the main idea.

**Example 2.** Consider the  $(K, q, M) = (5, 1, 3)$  decentralized data shuffling problem, where  $m = M/q = 3$ . Let  $\mathcal{A}^{t-1} = (5, 1, 2, 3, 4)$ . During the storage update phase in time slot  $t-1$ , we partition each data unit into 6 equal-length sub-blocks, each of which has  $B/6$  bits, as

$$F_1 = (F_{1,\{1,2,3\}}, F_{1,\{1,2,4\}}, F_{1,\{1,2,5\}}, F_{1,\{2,3,4\}}, F_{1,\{2,3,5\}}, F_{1,\{2,4,5\}}), \quad (64a)$$

$$F_2 = (F_{2,\{1,2,3\}}, F_{2,\{1,3,4\}}, F_{2,\{1,3,5\}}, F_{2,\{2,3,4\}}, F_{2,\{2,3,5\}}, F_{2,\{3,4,5\}}), \quad (64b)$$

$$F_3 = (F_{3,\{1,2,4\}}, F_{3,\{1,3,4\}}, F_{3,\{1,4,5\}}, F_{3,\{2,3,4\}}, F_{3,\{2,4,5\}}, F_{3,\{3,4,5\}}), \quad (64c)$$

$$F_4 = (F_{4,\{1,2,5\}}, F_{4,\{1,3,5\}}, F_{4,\{1,4,5\}}, F_{4,\{2,3,5\}}, F_{4,\{2,4,5\}}, F_{4,\{3,4,5\}}), \quad (64d)$$

$$F_5 = (F_{5,\{1,2,3\}}, F_{5,\{1,2,4\}}, F_{5,\{1,2,5\}}, F_{5,\{1,3,4\}}, F_{5,\{1,3,5\}}, F_{5,\{1,4,5\}}), \quad (64e)$$

and each worker  $k$  stores  $F_{i,\mathcal{W}}$  if  $k \in \mathcal{W}$ .

In the following, we consider various shuffles in time slot  $t$ . If one sub-block is stored by some worker in  $\mathcal{U}$ , we let this worker transmit it and the transmission is equivalent to centralized data shuffling; otherwise, we will introduce the proposed scheme B to transmit it.

We first consider  $\mathcal{A}^t = (1, 2, 3, 4, 5)$ . For each set  $\mathcal{J} \subseteq [K]$  of size  $|\mathcal{J}| = m + 1 = K - 1 = 4$ , we generate

$$V_{\mathcal{J}}^t = \bigoplus_{k \in \mathcal{J}} F_{d_k^t, \mathcal{J} \setminus \{k\}}, \quad (65)$$

where  $d_k^t$  represents the demanded data unit of worker  $k$  in time slot  $t$  if  $q = 1$ . The details are illustrated in Table II. For example, when  $\mathcal{J} = \{1, 2, 3, 4\}$ , we have  $V_{\{1,2,3,4\}}^t = F_{1,\{2,3,4\}} + F_{2,\{1,3,4\}} + F_{3,\{1,2,4\}} + F_{4,\{1,2,3\}}$  where  $F_{4,\{1,2,3\}}$  is empty and we replace  $F_{4,\{1,2,3\}}$  by B/6 zero bits. Since  $F_{1,\{2,3,4\}}$ ,  $F_{2,\{1,3,4\}}$ , and  $F_{3,\{1,2,4\}}$  are all stored by worker 4,  $V_{\{1,2,3,4\}}^t$  can be transmitted by worker 4. Similarly, we let worker 3 transmit  $V_{\{1,2,3,5\}}^t$ , worker 2 transmit  $V_{\{1,2,4,5\}}^t$ , worker 1 transmit  $V_{\{1,3,4,5\}}^t$ , and worker 5 transmit  $V_{\{2,3,4,5\}}^t$ . It can be checked that each worker can recover its desired sub-blocks and the achieved load is  $5/6$ , which coincides with the proposed converse bound.

Let us then focus on  $\mathcal{A}^t = (5, 2, 3, 4, 1)$ . For this shuffle,  $\mathcal{A}_1^{t-1} = \mathcal{A}_1^t$  such that worker 1 needs not to decode anything from what the other workers transmit. We divide all desired sub-blocks into two sets, stored and not stored by worker 1 as follows

$$\begin{aligned} \mathcal{S}_{\{1\}} &= \{F_{1,\{1,2,3\}}, F_{1,\{1,2,4\}}, F_{2,\{1,3,4\}}, F_{2,\{1,3,5\}}, F_{3,\{1,2,4\}}, \\ &\quad F_{3,\{1,4,5\}}, F_{4,\{1,2,5\}}, F_{4,\{1,3,5\}}\}, \\ \mathcal{S}_{\emptyset} &= \{F_{1,\{2,3,4\}}, F_{2,\{3,4,5\}}, F_{3,\{2,4,5\}}, F_{4,\{2,3,5\}}\}. \end{aligned}$$

Since the sub-blocks in  $\mathcal{S}_{\{1\}}$  are all stored by worker 1, we can treat worker 1 as a central server and the transmission of  $\mathcal{S}_{\{1\}}$  is equivalent to centralized data shuffling with  $K_{\text{eq}} = 4$ ,  $M_{\text{eq}} = 2$  and  $q_{\text{eq}} = 1$ , where  $\mathcal{U}_{\text{eq}} = \emptyset$ . For this centralized problem, the data shuffling schemes in [11], [12] are optimal under the constraint of uncoded storage. Alternatively, we can also use the following simplified scheme. By generating  $V_{\{1,2,3,4\}}^t$  as in (65), and it can be seen that  $V_{\{1,2,3,4\}}^t$  is known by workers 1 and 4. Similarly,  $V_{\{1,2,3,5\}}^t$  is known by workers 1 and 3,  $V_{\{1,2,4,5\}}^t$  is known by workers 1 and 2, and  $V_{\{1,3,4,5\}}^t$  is known by workers 1 and 5. Hence, we can let worker 1 transmit  $V_{\{1,2,3,4\}}^t \oplus V_{\{1,2,3,5\}}^t$ ,  $V_{\{1,2,3,4\}}^t \oplus V_{\{1,2,4,5\}}^t$ , and  $V_{\{1,2,3,4\}}^t \oplus V_{\{1,3,4,5\}}^t$ . Hence, each worker can recover  $V_{\{1,2,3,4\}}^t$ ,  $V_{\{1,2,3,5\}}^t$ ,  $V_{\{1,2,4,5\}}^t$ , and  $V_{\{1,3,4,5\}}^t$ . We then consider the transmission for  $\mathcal{S}_{\emptyset} = \{F_{1,\{2,3,4\}}, F_{2,\{3,4,5\}}, F_{3,\{2,4,5\}}, F_{4,\{2,3,5\}}\}$ , which is equivalent to decentralized data shuffling with  $K_{\text{eq}} = 4$ ,  $M_{\text{eq}} = 3$  and  $q_{\text{eq}} = 1$ , where  $\mathcal{U}_{\text{eq}} = \emptyset$  defined in (63). Hence, we can use the proposed Scheme B for  $m = K - 1$ . More precisely, we split each sub-block in  $V_{\{2,3,4,5\}}^t$  into 3 non-overlapping and equal-length sub-pieces, e.g.,  $F_{2,\{3,4,5\}} = \{F_{2,\{3,4,5\}}(3), F_{2,\{3,4,5\}}(4), F_{2,\{3,4,5\}}(5)\}$ . We then let

worker 2 transmit  $F_{3,\{2,4,5\}}(2) \oplus F_{4,\{2,3,5\}}(2) \oplus F_{1,\{2,3,4\}}(2)$ ;  
 worker 3 transmit  $F_{2,\{3,4,5\}}(3) \oplus F_{4,\{2,3,5\}}(3) \oplus F_{1,\{2,3,4\}}(3)$ ;  
 worker 4 transmit  $F_{2,\{3,4,5\}}(4) \oplus F_{3,\{2,4,5\}}(4) \oplus F_{1,\{2,3,4\}}(4)$ ;  
 worker 5 transmit  $F_{2,\{3,4,5\}}(5) \oplus F_{3,\{2,4,5\}}(5) \oplus F_{4,\{2,3,5\}}(5)$ .

In conclusion, the total load for  $\mathcal{A}^t = (5, 2, 3, 4, 1)$  is  $\frac{3}{6} + \frac{2}{9} = \frac{13}{18}$ .

TABLE II: Multicast messages for Example 2. Empty sub-blocks are colored in magenta.

For $\mathcal{A}^t = (1, 2, 3, 4, 5) = (F_1, F_2, F_3, F_4, F_5)$	
Worker 1 wants	$(F_{1,\{2,3,4\}}, F_{1,\{2,3,5\}}, F_{1,\{2,4,5\}}, F_{1,\{3,4,5\}}) = 0$
Worker 2 wants	$(F_{2,\{1,3,4\}}, F_{2,\{1,3,5\}}, F_{2,\{3,4,5\}}, F_{2,\{1,4,5\}}) = 0$
Worker 3 wants	$(F_{3,\{1,2,4\}}, F_{3,\{1,4,5\}}, F_{3,\{2,4,5\}}, F_{3,\{1,2,5\}}) = 0$
Worker 4 wants	$(F_{4,\{1,2,5\}}, F_{4,\{1,3,5\}}, F_{4,\{2,3,5\}}, F_{4,\{1,2,3\}}) = 0$
Worker 5 wants	$(F_{5,\{1,2,3\}}, F_{5,\{1,2,4\}}, F_{5,\{1,3,4\}}, F_{5,\{2,3,4\}}) = 0$
$V_{\{1,2,3,4\}}^t$	$= F_{1,\{2,3,4\}} + F_{2,\{1,3,4\}} + F_{3,\{1,2,4\}} + F_{4,\{1,2,3\}}$
$V_{\{1,2,3,5\}}^t$	$= F_{1,\{2,3,5\}} + F_{2,\{1,3,5\}} + F_{3,\{1,2,5\}} + F_{5,\{1,2,3\}}$
$V_{\{1,2,4,5\}}^t$	$= F_{1,\{2,4,5\}} + F_{2,\{1,4,5\}} + F_{4,\{1,2,5\}} + F_{5,\{1,2,4\}}$
$V_{\{1,3,4,5\}}^t$	$= F_{1,\{3,4,5\}} + F_{3,\{1,4,5\}} + F_{4,\{1,3,5\}} + F_{5,\{1,3,4\}}$
$V_{\{2,3,4,5\}}^t$	$= F_{2,\{3,4,5\}} + F_{3,\{2,4,5\}} + F_{4,\{2,3,5\}} + F_{5,\{2,3,4\}}$
For $\mathcal{A}^t = (5, 2, 3, 4, 1) = (F_5, F_2, F_3, F_4, F_1)$	
Worker 1 wants	$(F_{5,\{2,3,4\}}, F_{5,\{2,3,5\}}, F_{5,\{2,4,5\}}, F_{5,\{3,4,5\}}) = 0$
Worker 2 wants	$(F_{2,\{1,3,4\}}, F_{2,\{1,3,5\}}, F_{2,\{3,4,5\}}, F_{2,\{1,4,5\}}) = 0$
Worker 3 wants	$(F_{3,\{1,2,4\}}, F_{3,\{1,4,5\}}, F_{3,\{2,4,5\}}, F_{3,\{1,2,5\}}) = 0$
Worker 4 wants	$(F_{4,\{1,2,5\}}, F_{4,\{1,3,5\}}, F_{4,\{2,3,5\}}, F_{4,\{1,2,3\}}) = 0$
Worker 5 wants	$(F_{1,\{1,2,3\}}, F_{1,\{1,2,4\}}, F_{1,\{2,3,4\}}, F_{1,\{1,3,4\}}) = 0$
$V_{\{1,2,3,4\}}^t$	$= F_{5,\{2,3,4\}} + F_{2,\{1,3,4\}} + F_{3,\{1,2,4\}} + F_{4,\{1,2,3\}}$
$V_{\{1,2,3,5\}}^t$	$= F_{5,\{2,3,5\}} + F_{2,\{1,3,5\}} + F_{3,\{1,2,5\}} + F_{1,\{1,2,3\}}$
$V_{\{1,2,4,5\}}^t$	$= F_{5,\{2,4,5\}} + F_{2,\{1,4,5\}} + F_{4,\{1,2,5\}} + F_{1,\{1,2,4\}}$
$V_{\{1,3,4,5\}}^t$	$= F_{5,\{3,4,5\}} + F_{3,\{1,4,5\}} + F_{4,\{1,3,5\}} + F_{1,\{1,3,4\}}$
$V_{\{2,3,4,5\}}^t$	$= F_{2,\{3,4,5\}} + F_{3,\{2,4,5\}} + F_{4,\{2,3,5\}} + F_{1,\{2,3,4\}}$
For $\mathcal{A}^t = (5, 1, 3, 4, 2) = (F_5, F_1, F_3, F_4, F_2)$	
Worker 1 wants	$(F_{5,\{2,3,4\}}, F_{5,\{2,3,5\}}, F_{5,\{2,4,5\}}, F_{5,\{3,4,5\}}) = 0$
Worker 2 wants	$(F_{1,\{1,3,4\}}, F_{1,\{1,3,5\}}, F_{1,\{3,4,5\}}, F_{1,\{1,4,5\}}) = 0$
Worker 3 wants	$(F_{3,\{1,2,4\}}, F_{3,\{1,4,5\}}, F_{3,\{2,4,5\}}, F_{3,\{1,2,5\}}) = 0$
Worker 4 wants	$(F_{4,\{1,2,5\}}, F_{4,\{1,3,5\}}, F_{4,\{2,3,5\}}, F_{4,\{1,2,3\}}) = 0$
Worker 5 wants	$(F_{2,\{1,2,3\}}, F_{2,\{1,3,4\}}, F_{2,\{2,3,4\}}, F_{2,\{1,2,4\}}) = 0$
$V_{\{1,2,3,4\}}^t$	$= F_{5,\{2,3,4\}} + F_{1,\{1,3,4\}} + F_{3,\{1,2,4\}} + F_{4,\{1,2,3\}}$
$V_{\{1,2,3,5\}}^t$	$= F_{5,\{2,3,5\}} + F_{1,\{1,3,5\}} + F_{3,\{1,2,5\}} + F_{2,\{1,2,3\}}$
$V_{\{1,2,4,5\}}^t$	$= F_{5,\{2,4,5\}} + F_{1,\{1,4,5\}} + F_{4,\{1,2,5\}} + F_{2,\{1,2,4\}}$
$V_{\{1,3,4,5\}}^t$	$= F_{5,\{3,4,5\}} + F_{3,\{1,4,5\}} + F_{4,\{1,3,5\}} + F_{2,\{1,3,4\}}$
$V_{\{2,3,4,5\}}^t$	$= F_{1,\{3,4,5\}} + F_{3,\{2,4,5\}} + F_{4,\{2,3,5\}} + F_{2,\{2,3,4\}}$

Finally, we consider  $\mathcal{A}^t = \{5, 1, 3, 4, 2\}$ . For this shuffle,  $\mathcal{A}_1^{t-1} = \mathcal{A}_1^t$  and  $\mathcal{A}_2^{t-1} = \mathcal{A}_2^t$  such that workers 1 and 2 need not to decode anything from other workers. We divide all desired sub-blocks into three sets

$$\mathcal{S}_{\{1,2\}} = \{F_{2,\{1,2,3\}}, F_{3,\{1,2,4\}}, F_{4,\{1,2,5\}}\},$$

stored by workers 1 and 2,

$$\mathcal{S}_{\{1\}} = \{F_{2,\{1,3,4\}}, F_{3,\{1,4,5\}}, F_{4,\{1,3,5\}}\},$$

stored by worker 1 and not by worker 2, and

$$\mathcal{S}_{\{2\}} = \{F_{2,\{2,3,4\}}, F_{3,\{2,4,5\}}, F_{4,\{2,3,5\}}\}$$

stored by worker 2 and not by worker 1.

The transmission for  $\mathcal{S}_{\{1,2\}}$  is equivalent to a centralized data shuffling with  $K_{\text{eq}} = 3$ ,  $M_{\text{eq}} = 1$  and  $q_{\text{eq}} = 1$ . We use the following simplified scheme to let worker 1 transmit  $V_{\{1,2,3,4\}}^t \oplus V_{\{1,2,3,5\}}^t$  and  $V_{\{1,2,3,4\}}^t \oplus V_{\{1,2,4,5\}}^t$  such that each worker can recover  $V_{\{1,2,3,4\}}^t$ ,  $V_{\{1,2,3,5\}}^t$ , and  $V_{\{1,2,4,5\}}^t$  (as illustrated in Table II,  $V_{\{1,2,3,4\}}^t = F_{3,\{1,2,4\}}$ ,  $V_{\{1,2,3,5\}}^t = F_{2,\{1,2,3\}}$ , and  $V_{\{1,2,4,5\}}^t = F_{4,\{1,2,5\}}$ ). Similarly, for  $\mathcal{S}_{\{1\}}$ , we let worker 1 transmit  $V_{\{1,3,4,5\}}^t$ . For  $\mathcal{S}_{\{2\}}$ , we let worker 2 transmit  $V_{\{2,3,4,5\}}^t$ . In conclusion the total load for  $\mathcal{A}^t = \{5, 1, 3, 4, 2\}$  is  $\frac{2}{6} + \frac{1}{6} + \frac{1}{6} = \frac{2}{3}$ .  $\square$

Now we are ready to introduce Scheme B for  $m = K - 2$  as a generalization of Example 2. Recall that, from our earlier discussion, we can consider without loss of generality  $q = 1$ , and that  $\mathcal{U}$  represents the set of workers who need not to

recover anything from others. We divide all desired sub-blocks by all workers into non-overlapping sets

$$\mathcal{S}_{\mathcal{K}} := \{F_{d_k, \mathcal{W}} : k \in [K] \setminus \mathcal{U}, |\mathcal{W}| = m + 1, \mathcal{W} \cap \mathcal{U} = \mathcal{K}, k \notin \mathcal{W}\}, \quad (66)$$

where  $\mathcal{K} \subseteq \mathcal{U}$ . We then encode the sub-blocks in each set in  $\mathcal{S}_{\mathcal{K}}$  in (66) as follows:

- For each  $\mathcal{K} \subseteq \mathcal{U}$  where  $\mathcal{K} \neq \emptyset$ , the transmission for  $\mathcal{S}_{\mathcal{K}}$  is equivalent to a centralized data shuffling problem with  $K_{\text{eq}} = K - |\mathcal{U}|$ ,  $q_{\text{eq}} = 1$  and  $M_{\text{eq}} = m - |\mathcal{K}|$ , where  $\mathcal{U}_{\text{eq}} = \emptyset$ . It can be seen that  $K_{\text{eq}} - M_{\text{eq}} \leq 2$ . Hence, we can use the optimal centralized data shuffling schemes in [11], [12].

Alternatively, we propose the following simplified scheme. For each set  $\mathcal{J} \subseteq [K]$  of size  $|\mathcal{J}| = m + 1 = K - 1$ , where  $\mathcal{J} \cap \mathcal{U} = \mathcal{K}$ , we generate  $V_{\mathcal{J}}^t$  as in (65). Each sub-block in  $\mathcal{S}_{\mathcal{K}}$  appears in one  $V_{\mathcal{J}}^t$ , where  $\mathcal{J} \subseteq [K]$ ,  $|\mathcal{J}| = K - 1$  and  $\mathcal{J} \cap \mathcal{U} = \mathcal{K}$ . It can be seen that for each worker  $j \in [K] \setminus \mathcal{U}$ , among all  $V_{\mathcal{J}}^t$  where  $\mathcal{J} \subseteq [K]$ ,  $|\mathcal{J}| = K - 1$  and  $\mathcal{J} \cap \mathcal{U} = \mathcal{K}$ , worker  $j$  knows one of them (which is  $V_{[K] \setminus \{u_{d_j}^{t-1}\}}^t$ ). We denote all sets  $\mathcal{J} \subseteq [K]$  where  $|\mathcal{J}| = K - 1$  and  $\mathcal{J} \cap \mathcal{U} = \mathcal{K}$ , by  $\mathcal{J}_1(\mathcal{K}), \mathcal{J}_2(\mathcal{K}), \dots, \mathcal{J}_{\binom{K-|\mathcal{U}|}{K-1-|\mathcal{K}|}}(\mathcal{K})$ . For  $\mathcal{S}_{\mathcal{K}}$ , we choose one worker in  $\mathcal{K}$  to transmit  $V_{\mathcal{J}_1(\mathcal{K})}^t \oplus V_{\mathcal{J}_2(\mathcal{K})}^t, \dots, V_{\mathcal{J}_1(\mathcal{K})}^t \oplus V_{\mathcal{J}_{\binom{K-|\mathcal{U}|}{K-1-|\mathcal{K}|}}(\mathcal{K})}^t$ , such that each worker in  $K \setminus \mathcal{U}$  can recover all  $V_{\mathcal{J}}^t$  where  $\mathcal{J} \subseteq [K]$ ,  $|\mathcal{J}| = K - 1$  and  $\mathcal{J} \cap \mathcal{U} = \mathcal{K}$ .

- For  $\mathcal{K} = \emptyset$ , the transmission for  $\mathcal{S}_{\mathcal{K}}$  is equivalent to decentralized data shuffling with  $K_{\text{eq}} = K - |\mathcal{U}|$ ,  $q_{\text{eq}} = 1$  and  $M_{\text{eq}} = m = K - 2$ , where  $\mathcal{U}_{\text{eq}} = \emptyset$ . Hence, in this case  $\mathcal{U} \leq 2$ .

If  $|\mathcal{U}| = 2$ , we have  $M_{\text{eq}} = K_{\text{eq}}$  and thus we do not transmit anything for  $\mathcal{S}_{\emptyset}$ .

If  $|\mathcal{U}| = 1$ , we have  $M_{\text{eq}} = K_{\text{eq}} - 1$  and thus we can use Scheme B for  $m = K - 1$  to transmit  $\mathcal{S}_{\emptyset}$ .

Finally, we consider  $|\mathcal{U}| = \emptyset$ . For each set  $\mathcal{J} \subseteq [K]$  where  $|\mathcal{J}| = m + 1 = K - 1$ , among all the workers in  $\mathcal{J}$ , there is exactly one worker in  $\mathcal{J}$  where  $u_{d_k}^{t-1} \notin \mathcal{J}$  (this worker is assumed to be  $k$  and we have  $u_{d_k}^{t-1} = [K] \setminus \mathcal{J}$  with a slight abuse of notation). We then let worker  $k$  transmit  $V_{\mathcal{J}}^t$ .

In conclusion, by comparing the loads for different cases, the worst-cases are when  $\mathcal{A}_k^{t-1} \cap \mathcal{A}_k^t = \emptyset$  for each  $k \in [K]$  and the worst-case load achieved by Scheme B is

$$qK / \binom{K-1}{K-2} = q \frac{K-m}{m} \frac{K}{K-1} \Big|_{m=K-1} =: R_{\text{Ach.B}} \Big|_{M=(K-2)q}, \quad (67)$$

which coincides with the converse bound.

*Storage Update Phase of time slot  $t \in [T]$  for  $m = K - 2$ :*

The storage update phase for  $m = K - 2$  is the same as the one of scheme B for  $m = K - 1$ .

**Remark 3** (Scheme B realizes distributed interference alignment). *In Example 2, from  $\mathcal{A}^{t-1} = (5, 1, 2, 3, 4)$  to  $\mathcal{A}^t = (1, 2, 3, 4, 5)$ , by the first column of Table II, we see that each worker desires  $K - 2 = 3$  of the sub-blocks that need to*

*be shuffled. Since each worker cannot benefit from its own transmission, we see that the best possible scheme would have each worker recover its  $K - 2 = 3$  desired sub-blocks from the  $K - 1 = 4$  “useful transmissions,” that is, the unwanted sub-blocks should “align” in a single transmission, e.g., for worker 1, all of its unwanted sub-blocks are aligned in  $V_{\{2,3,4,5\}}^t$ . From the above we see that this is exactly what happens for each worker when  $K - 2 = m$ . How to realize distributed interference alignment seems to be a key question in decentralized data shuffling.  $\square$*

**Remark 4** (Extension of Scheme B to other storage sizes). *We can extend Scheme B to any storage size by the following three steps:*

- We partition the  $N$  data units into  $q$  groups, where each group contains  $K$  data units, and such that during the data shuffling phase each worker requests exactly one data unit and knows exactly one data unit among all the  $K$  data units in each group.
- For each group  $\mathcal{H}_i$ , we partition all desired sub-blocks by all workers into sets depending on which workers in  $\mathcal{U}(\mathcal{H}_i)$  know them. Each set is denoted by  $\mathcal{S}_{\mathcal{K}}(\mathcal{H}_i)$ , which is known by workers in  $\mathcal{K} \subseteq \mathcal{U}(\mathcal{H}_i)$ , and is defined similarly to (66).
- For each set  $\mathcal{S}_{\mathcal{K}}(\mathcal{H}_i)$ ,
  - if  $\mathcal{K} \neq \emptyset$ , the transmission is equivalent to centralized data shuffling with  $K_{\text{eq}} = K - |\mathcal{U}(\mathcal{H}_i)|$ ,  $q_{\text{eq}} = 1$  and  $M_{\text{eq}} = M - |\mathcal{K}|$ . We can use the optimal centralized data shuffling scheme in [12];
  - if  $\mathcal{K} = \emptyset$ , for each set  $\mathcal{J} \subseteq ([K] \setminus \mathcal{U}(\mathcal{H}_i))$ , where  $|\mathcal{J}| = m + 1$ , we generate the multicast messages  $V_{\mathcal{J}}^t$  as defined in (65). If there exists some empty sub-block in  $V_{\mathcal{J}}^t$ , we let the worker who demands this sub-block transmit  $V_{\mathcal{J}}^t$ . Otherwise,  $V_{\mathcal{J}}^t$  is transmitted as Scheme B for  $m = K - 1$ .

*Unfortunately, a general closed-form expression for the load in this general case is not available as it heavily depends on the shuffle.*

*Note: in Example 3 next we show the transmission for  $\mathcal{K} = \emptyset$ , can be further improved by random linear combinations.  $\square$*

**Example 3.** Consider the  $(K, q, M) = (5, 1, 2)$  decentralized data shuffling problem, where  $m = M/q = 2$ . From the outer bound we have  $R_{\text{u}}^* \geq 15/8$ ; if each sub-block is of size  $1/\binom{K-1}{m-1} = 1/4$ , the outer bound suggests that we need to transmit at least  $15/2 = 7.5$  sub-blocks in total.

Let  $\mathcal{A}^{t-1} = (5, 1, 2, 3, 4)$ . During the storage update phase in time slot  $t - 1$ , we partition each data unit into 4 equal-length sub-blocks, each of which has  $B/4$  bits, as

$$F_1 = (F_{1,\{1,2\}}, F_{1,\{2,3\}}, F_{1,\{2,4\}}, F_{1,\{2,5\}}), \quad (68a)$$

$$F_2 = (F_{2,\{1,3\}}, F_{2,\{2,3\}}, F_{2,\{3,4\}}, F_{2,\{3,5\}}), \quad (68b)$$

$$F_3 = (F_{3,\{1,4\}}, F_{3,\{2,4\}}, F_{3,\{3,4\}}, F_{3,\{4,5\}}), \quad (68c)$$

$$F_4 = (F_{4,\{1,5\}}, F_{4,\{2,5\}}, F_{4,\{3,5\}}, F_{4,\{4,5\}}), \quad (68d)$$

$$F_5 = (F_{5,\{1,2\}}, F_{5,\{1,3\}}, F_{5,\{1,4\}}, F_{5,\{1,5\}}), \quad (68e)$$

and each worker  $k$  stores  $F_{i,\mathcal{W}}$  if  $k \in \mathcal{W}$ .

Let  $\mathcal{A}^t = (1, 2, 3, 4, 5)$ . During the data shuffling phase in time slot  $t$ , each worker must recover 3 sub-blocks of the desired data unit which it does not store, e.g., worker 1 must recover  $(F_{1,\{2,3\}}, F_{1,\{2,4\}}, F_{1,\{2,5\}})$ , worker 2 must recover  $(F_{2,\{1,3\}}, F_{2,\{3,4\}}, F_{2,\{3,5\}})$ , etc.

For each set  $\mathcal{J} \subseteq [K]$  where  $|\mathcal{J}| = m + 1 = 3$ , we generate  $V_{\mathcal{J}}^t = \bigoplus_{k \in \mathcal{J}} F_{k, \mathcal{J} \setminus \{k\}}$  as in (65). More precisely, we have

$$V_{\{1,2,3\}}^t = F_{1,\{2,3\}} \oplus F_{2,\{1,3\}}, \quad (\text{can be sent by worker 3}), \quad (69a)$$

$$V_{\{1,2,4\}}^t = F_{1,\{2,4\}}, \quad (69b)$$

$$V_{\{1,2,5\}}^t = F_{1,\{2,5\}} \oplus F_{5,\{1,2\}}, \quad (\text{can be sent by worker 2}), \quad (69c)$$

$$V_{\{1,3,4\}}^t = F_{3,\{1,4\}}, \quad (69d)$$

$$V_{\{1,3,5\}}^t = F_{5,\{1,3\}}, \quad (69e)$$

$$V_{\{1,4,5\}}^t = F_{4,\{1,5\}} \oplus F_{5,\{1,4\}}, \quad (\text{can be sent by worker 1}), \quad (69f)$$

$$V_{\{2,3,4\}}^t = F_{2,\{3,4\}} \oplus F_{3,\{2,4\}}, \quad (\text{can be sent by worker 4}), \quad (69g)$$

$$V_{\{2,3,5\}}^t = F_{2,\{3,5\}}, \quad (69h)$$

$$V_{\{2,4,5\}}^t = F_{4,\{2,5\}}, \quad (69i)$$

$$V_{\{3,4,5\}}^t = F_{3,\{4,5\}} \oplus F_{4,\{3,5\}}, \quad (\text{can be sent by worker 5}). \quad (69j)$$

We deliver these multicast messages with a two-phase scheme, as follows.

- Phase 1. It can be seen that the multicast messages like  $V_{\{1,2,3\}}^t$  in (69) (which is known by worker 3 only) can be sent by one specific worker. Similarly, we can let workers 2, 1, 4 and 5 broadcast  $V_{\{1,2,5\}}^t$ ,  $V_{\{1,4,5\}}^t$ ,  $V_{\{2,3,4\}}^t$  and  $V_{\{3,4,5\}}^t$ , respectively.
- Phase 2. After the above Phase 1, the remaining messages are known by two workers. For example,  $V_{\{1,2,4\}}^t = F_{1,\{2,4\}}$  is known by workers 2 and 4; we can let worker 2 transmit  $V_{\{1,2,4\}}^t$ . If we do so as Scheme B, since each multicast message in (69) has  $B/4$  bits and there are 10 multicast messages in (69), the total load is  $10/4$ .

In this example Phase 2 of Scheme B can be improved as follows. The problem with the above strategy (i.e., assign each multicast message to a worker) is that we have not leveraged the fact that, after Phase 1, there are still five sub-blocks to be delivered (one demanded per worker, namely  $F_{1,\{2,4\}}, F_{2,\{3,5\}}, F_{3,\{1,4\}}, F_{4,\{2,5\}}, F_{5,\{1,3\}}$ ), each of which is known by two workers. Therefore, we can form random linear combinations so that each worker can recover all three of the unstored sub-blocks. In other words, if a worker receives from each of the remaining  $K - 1 = 4$  workers  $3/4 \times$  “size of a sub-block” linear equations then it solved for the three missing sub-blocks, that is, each worker broadcasts  $\frac{3B}{16}$  random linear combinations of all bits of the two sub-blocks he stores among  $F_{1,\{2,4\}}, F_{2,\{3,5\}}, F_{3,\{1,4\}}, F_{4,\{2,5\}}, F_{5,\{1,3\}}$ . It can be checked that for worker 1, the received  $3B/4$  random linear combinations from other workers are linearly independent known  $F_{1,\{2,4\}}$  and  $F_{5,\{1,3\}}$  as

$B \rightarrow \infty$ , such that it can recover all these five sub-blocks. By the symmetry, each other worker can also recover these five sub-blocks.

In conclusion, the total load with this two-phase is  $5(1 + 3/4) \times 1/4 = \frac{35}{16} < 10/4$ , which is achieved by Scheme B. By comparison, the load of Scheme A is  $\frac{27}{8}$  and the converse bound under the constraint of uncoded storage in Theorem 1 is  $\frac{15}{8}$ .

As a final remark, note that the five sub-blocks in Phase 2 is symmetric, i.e., the number of sub-blocks stored by each worker is the same and the one known by each worker is also the same. In general case, this symmetry may not hold (thus the conditional linearly independent of the received random linear combinations by each worker may not hold), and the generalization of this scheme is part of ongoing work.  $\square$

**Remark 5** (Limitation of Scheme B for small storage size). *The data shuffling phase with uncoded storage can be represented by a directed graph, where each sub-block demanded by some worker is represented by a node in the graph. A directed edge exists from node  $a$  to node  $b$  in the graph if the worker demanding the data represented by node  $b$  has the data represented by node  $a$  in its storage. By generating a directed graph for the data shuffling phase as described in Section II, we can see that each multicast message in (65) is the sum of the sub-blocks contained in a clique, where a clique is a set of nodes where each two of them are connected in two directions. The sub-blocks in each multicast message in (65) also form a distributed clique, where a distributed clique is a clique whose nodes are all known by some worker.*

Consider the case where  $K = N$  is much larger than  $M = m = 2$  (i.e., small storage size regime). Consider a “cyclic shuffle” of the form  $\mathcal{A}_1^{t-1} = \{K\}$  and  $\mathcal{A}_k^{t-1} = \{k - 1\}$  for  $k \in [2 : K]$ , to  $\mathcal{A}_k^t = \{k\}$  for  $k \in [K]$ . Each data unit is split into  $\binom{K-1}{m-1} = K - 1$  sub-blocks and each worker needs to recover  $\binom{K-2}{m-1} = K - 2$  sub-blocks.

If during the data shuffling phase in time slot  $t$  we generate the multicast messages as above, only 2 of the  $K - 2$  demanded sub-blocks are in a distributed clique of size  $m = 2$ . More precisely, let us focus on worker 1 who needs to recover  $F_1$ . Notice that each sub-block of  $F_1$  is stored by worker 2, and each of its demanded sub-blocks  $F_{1,\{2,j\}}$  where  $j \in [K] \setminus \{1, 2\}$ , is in the multicast message

$$V_{\{1,2,j\}}^t = F_{1,\{2,j\}} \oplus F_{2,\{1,j\}} \oplus F_{j,\{1,2\}}. \quad (70)$$

When  $j = 3$ , it can be seen that  $F_{3,\{1,2\}}$  is empty because all sub-blocks of  $F_3$  are stored by worker 4, and thus  $V_{\{1,2,3\}}^t = F_{1,\{2,3\}} \oplus F_{2,\{1,3\}}$  could be transmitted by worker 3. However, when  $j = 4$ , both sub-blocks  $F_{2,\{1,4\}}$  and  $F_{4,\{1,2\}}$  are empty, and thus  $V_{\{1,2,4\}}^t = F_{1,\{2,4\}}$ . Similarly, among all the  $K - 2$  demanded sub-blocks by worker 1, only  $F_{1,\{2,3\}}$  and  $F_{1,\{2,K\}}$  are in cliques including two sub-blocks, while the remaining  $K - 4$  ones are in cliques including only one sub-block.

If the delivery strategy is to assign each multicast message, or distributed clique, to a worker, we see that most of the distributed cliques have a multicast coding gain of 1 (where the multicast coding gain is the gain on the transmitted load compared to uncoded/direct transmission, e.g., if we transmit

one sub-block in a distributed clique of length 2, the multicast coding gain to transmit this sub-block is 2). Hence, Scheme B is generally inefficient for  $m \in [2 : K - 3]$ . In this paper we mainly use Scheme B for  $m \in \{K - 2, K - 1, K\}$  for which it is optimal under the constraint of uncoded storage.  $\square$

### C. Scheme C in Theorem 4

To overcome the limitation of Scheme B described in Remark 5, in the following we propose Scheme C for  $M/q = m = 2$  based on an *unconventional distributed clique-covering* strategy for the two-sender distributed index coding problems proposed in [29].

The storage update phase of Scheme C is the same as Scheme B, which is structural invariant, and thus we only describe the transition from time slot  $t - 1$  to time slot  $t$ . The main idea is not to use the conventional distributed clique-covering method which transmits distributed cliques (e.g., the multicast messages in Scheme B), because most of the distributed cliques only contain one sub-block, and most sub-blocks are not in any distributed clique including more than one sub-blocks, as explained in Remark 5. Instead, we propose a novel decentralized data shuffling scheme to increase the efficiency (coded multicasting gain) of the transmissions. The main idea is that we search for cliques with length  $m = 2$ ; if this clique is not a distributed clique, we add one node to the clique and transmit the three nodes by two binary sums, such that each node can be recovered by the corresponding user; if this clique is a distributed clique we proceed as in Scheme B.

Before introducing the details of Scheme C, we first recall the unconventional clique-covering method for the two-sender distributed index coding problems proposed in [29, Theorem 8, Case 33].

**Proposition 1** (Unconventional Distributed Clique-covering in [29]). *In a distributed index coding problem with two senders (assumed to be  $k_1$  and  $k_2$ ), as illustrated Fig. 3, there are three messages  $a$  (demanded by receiver  $u(a)$ ),  $b$  (demanded by receiver  $u(b)$ ), and  $c$  (demanded by receiver  $u(c)$ ), where  $u(a)$  stores  $b$ ,  $u(b)$  stores  $a$ , and  $u(c)$  stores  $a$ . Sender  $k_1$  stores  $a$  and  $c$ , and sender  $k_2$  stores  $b$  and  $c$ . We can let worker  $k_1$  transmit  $a \oplus c$  and worker  $k_2$  transmit  $b \oplus c$ , such that workers  $u(a)$ ,  $u(b)$  and  $u(c)$  can recover node  $a$ ,  $b$  and  $c$ , respectively.*

*Indeed, receiver  $u(a)$  knows  $b$  such that it can recover  $c$ , and then recover  $a$ . Similarly receiver  $u(b)$  can recover  $b$ . User  $u(c)$  knows  $a$ , such that it can recover  $c$  from  $a \oplus c$ .*

In the decentralized data shuffling problem, each data unit is divided into sub-blocks depending on which subset of workers stored them before the data shuffling phase; each sub-block desired by a worker is an independent message in the corresponding distributed index coding problem; thus the data shuffling phase is a  $K$ -sender distributed index coding problem that contains a number of messages that in general is doubly exponential in the number of workers in the original decentralized data shuffling problem. Hence, it is non-trivial to use Proposition 1 in the decentralized data shuffling problem.

We then illustrate the main idea of Scheme C by means of an example.

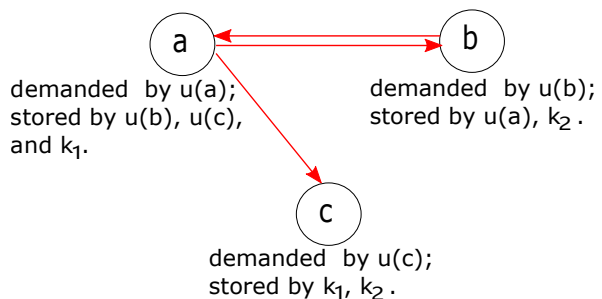


Fig. 3: Directed graph of the two-sender ( $k_1$  and  $k_2$ ) distributed index problem in Proposition 1. A direct edge from node  $a$  to node  $b$ , means that receiver  $u(b)$  demanding message  $b$  stores message  $a$  demanded by receiver  $u(a)$ .

**Example 4.** We consider the same example as Remark 4, where  $K = 5$ ,  $q = 1$  and  $M = 2$ . Let  $\mathcal{A}_1^{t-1} = \{5\}$  and  $\mathcal{A}_k^{t-1} = \{k - 1\}$  for  $k \in [2 : 5]$ , and  $\mathcal{A}_k^t = \{k\}$  for  $k \in [5]$ . The data units are split as in (68).

In the data shuffling phase of time slot  $t$ , the distributed clique-covering method in Scheme B has many sub-blocks which are not in any distributed clique including more than one sub-block (e.g.,  $F_{1,\{2,4\}}$ ), as explained in Remark 5. Moreover, the sub-blocks  $F_{1,\{2,3\}}$  and  $F_{3,\{1,4\}}$  are in a clique in the graph, but none of the workers can transmit  $F_{1,\{2,3\}} \oplus F_{3,\{1,4\}}$ , and thus it is not a distributed clique. However, if we add  $F_{1,\{2,4\}}$  to the group, and transmit the two sums  $F_{1,\{2,3\}} \oplus F_{1,\{2,4\}}$  (by worker 2) and  $F_{3,\{1,4\}} \oplus F_{1,\{2,4\}}$  (by worker 4), we see that worker 1 (who knows  $F_{3,\{1,4\}}$ ) can recover  $F_{1,\{2,4\}}$  from the second sum, and then recover  $F_{1,\{2,3\}}$  from the first sum. Similarly, worker 3 (who knows  $F_{1,\{2,3\}}$ ) can recover  $F_{1,\{2,4\}}$  from the first sum, and then recover  $F_{3,\{1,4\}}$  from the second sum. It can be seen that  $F_{1,\{2,3\},1}$  is message  $a$ ,  $F_{3,\{1,4\},1}$  is message  $b$ , and  $F_{1,\{2,4\},1}$  is message  $c$  in Proposition 1, while  $u(a)$  and  $u(c)$  are both worker 1,  $u(b)$  is worker 3. In addition, workers 2 and 4 serve as senders  $k_1$  and  $k_2$ .

Similarly, the sub-blocks  $F_{1,\{2,4\}}$  and  $F_{4,\{1,5\}}$  are in a clique in the graph, but none of the workers can transmit  $F_{1,\{2,4\}} \oplus F_{4,\{1,5\}}$ . However, if we add  $F_{1,\{2,5\}}$  to the group, and transmit  $F_{1,\{2,4\}} \oplus F_{1,\{2,5\}}$  (by worker 2) and  $F_{1,\{2,5\}} \oplus F_{4,\{1,5\}}$  (by worker 5), then worker 1 (who knows  $F_{4,\{1,5\}}$ ) can recover  $F_{1,\{2,5\}}$  from the second sum, and then recover  $F_{1,\{2,4\}}$  from the first sum; also, worker 4 (who knows  $F_{1,\{2,4\}}$ ) can recover  $F_{1,\{2,5\}}$  from the first sum, and then recover  $F_{4,\{1,5\}}$  from the second sum.

In general, we have the following data shuffling scheme in time slot  $t$ . Recall that  $u_i^{t-1}$  denotes the worker who should recover  $F_i$  at the end of time slot  $t - 1$ . We partition each sub-block into 3 equal-length sub-pieces as  $F_{i,\mathcal{W}} = (F_{i,\mathcal{W},1}, F_{i,\mathcal{W},2}, F_{i,\mathcal{W},3})$  (recall  $|\mathcal{W}| = m = 2$  and  $u_i^{t-1} \in \mathcal{W}$ ). In general, we consider a pair  $(a, b)$ , where  $a \in [5]$  and  $b \in [5] \setminus \{a, u_a^{t-1}\}$ , i.e.,  $a$  is a worker with the request  $F_a$  in time slot  $t$ , while  $b$  is another worker which is not the one who requests  $F_a$  in time slot  $t - 1$ . We have two cases:

- 1) If  $a \neq u_b^{t-1}$ , we find the group of sub-blocks  $F_{a,\{u_a^{t-1}, b\}}$ ,  $F_{a,\{u_a^{t-1}, u_b^{t-1}\}}$ , and  $F_{b,\{a, u_b^{t-1}\}}$ . We pick

one of untransmitted sub-pieces for each of these three sub-blocks. Assume the picked sub-pieces are  $n_1$ ,  $n_2$ , and  $n_3$ , respectively. We let worker  $u_a^{t-1}$  transmit  $F_{a,\{u_a^{t-1},b\},n_1} \oplus F_{a,\{u_a^{t-1},u_b^{t-1}\},n_2}$ , and let worker  $u_b^{t-1}$  transmit  $F_{a,\{u_a^{t-1},u_b^{t-1}\},n_2} \oplus F_{b,\{u_a^{t-1},b\},n_3}$ . It can be seen that worker  $a$  can recover  $F_{a,\{u_a^{t-1},b\},n_1}$  and  $F_{a,\{u_a^{t-1},u_b^{t-1}\},n_2}$  while worker  $b$  can recover  $F_{b,\{u_a^{t-1},b\},n_3}$ .

- 2) If  $a = u_b^{t-1}$ , we assume  $c = u_a^{t-1}$  and  $d = u_c^{t-1}$  (e.g, if  $a = 1$  and  $b = 5$ , we have  $c = u_1^{t-1} = 2$  and  $d = u_c^{t-1} = 3$ ), i.e., worker  $a$  requests  $F_a$  in time slot  $t$ , worker  $c$  requests  $F_a$  in time slot  $t-1$  and requests  $F_c$  in time slot  $t$ , worker  $d$  requests  $F_c$  in time slot  $t-1$  and requests  $F_d$  in time slot  $t$ . We find the group of sub-blocks  $F_{a,\{c,b\}}$ ,  $F_{a,\{c,d\}}$ , and  $F_{c,\{a,d\}}$ . Notice that  $F_{a,\{c,d\}}$  and  $F_{c,\{a,d\}}$  form a distributed clique. We pick one of untransmitted sub-pieces for each of these three sub-blocks (assumed to be  $n_1$ ,  $n_2$ , and  $n_3$ , respectively). We let worker  $c$  transmit  $F_{a,\{c,b\},n_1}$ , and let worker  $d$  transmit  $F_{a,\{c,d\},n_2} \oplus F_{c,\{a,d\},n_3}$ . It can be seen that worker  $a$  can recover  $F_{a,\{c,b\},n_1}$  and  $F_{a,\{c,d\},n_2}$  while worker  $c$  can recover  $F_{c,\{a,d\},n_3}$ .

By this construction, see also Table III, each sub-block appears in three groups such that each of its sub-pieces is transmitted. We use two binary sums to encode each group containing 3 sub-pieces, such that the coding gain of this scheme is  $2/3$ . The achieved worst-case load is  $5/2$ , while the achieved loads by Schemes A and B are  $\frac{27}{8}$  and  $\frac{10}{4}$ , respectively.

The introduced scheme in this example has the same load as Scheme B. However, in general when  $M = 2q$ , the coding gain of the new scheme is  $2/3$ , and this is independent of  $K$ . For the other schemes, the coding gain of Scheme B is close to 1 if  $K$  is large, and the same holds for Scheme A (as it will be shown in Section VI-D). Therefore, Scheme C is preferable to the other schemes when  $M = 2q$  and  $K$  is large.  $\square$

We are now ready to generalize the scheme in Example 4 to the case  $M/q = m = 2$ .

*Structural Invariant Data Partitioning and Storage:* This is the same as in Scheme B (described in Section VI-B), i.e., each sub-block of  $F_i$ ,  $i \in [N]$ , is stored by worker  $u_i^t$  and by another worker in  $[K] \setminus \{u_i^t\}$ . The length of each sub-block is  $\frac{B}{\binom{K-1}{m-1}} = \frac{B}{K-1}$ .

*Data Shuffling Phase of time slot  $t \in [T]$ :* As in Scheme B, we partition the  $N$  data units into  $q$  equal-length groups such that, during the data shuffling phase of time slot  $t$ , among all the  $K$  data units in each group, each worker requests exactly one data unit and knows exactly one data unit. To simplify the description, as in Scheme B, we focus on one group and remove the  $\mathcal{H}_i$  in the notation. In addition, we can restrict attention to  $\mathcal{U} = \emptyset$ , because if  $\mathcal{U} \neq \emptyset$  we can divide all desired sub-blocks by all workers into sets as we did in Scheme B. For each set which is known by some worker in  $\mathcal{U}$ , the transmission for this set is equivalent to centralized data shuffling. Thus we only need to consider the set which is not known by any worker in  $\mathcal{U}$ , and the transmission for this set is equivalent to decentralized data shuffling with  $K_{\text{eq}} = K - |\mathcal{U}|$ ,

$q_{\text{eq}} = 1$  and  $M_{\text{eq}} = m$ . Hence, for the simplicity, we only consider  $\mathcal{U} = \emptyset$  for Scheme C.

We define a set of pair of workers as

$$\mathcal{Y} := \left\{ (a, b) : a \in [K], b \in ([K] \setminus \{u_a^{t-1}, a\}) \right\}. \quad (71)$$

For each vector  $(a, b) \in \mathcal{Y}$ , we divide  $F_{d_a^t, \{b, u_{d_a^t}^{t-1}\}}$  into 3 non-overlapping and equal-length sub-pieces,  $F_{d_a^t, \{b, u_{d_a^t}^{t-1}\}, 1}$ ,  $F_{d_a^t, \{b, u_{d_a^t}^{t-1}\}, 2}$  and  $F_{d_a^t, \{b, u_{d_a^t}^{t-1}\}, 3}$ . For each vector  $(a, b) \in \mathcal{Y}$ , we consider two cases:

- Case  $u_{d_b^t}^{t-1} \neq a$ : For each one of  $F_{d_a^t, \{u_{d_a^t}^{t-1}, b\}}$ ,  $F_{d_b^t, \{u_{d_b^t}^{t-1}, a\}}$ , and  $F_{d_c^t, \{u_{d_c^t}^{t-1}, u_{d_b^t}^{t-1}\}}$ , we select one of its non-transmitted sub-pieces and assume they are  $F_{d_a^t, \{u_{d_a^t}^{t-1}, b\}, n_1}$ ,  $F_{d_b^t, \{u_{d_b^t}^{t-1}, a\}, n_2}$ , and  $F_{d_c^t, \{u_{d_c^t}^{t-1}, u_{d_b^t}^{t-1}\}, n_3}$ . By Proposition 1,

worker  $u_{d_a^t}^{t-1}$  transmits  $F_{d_a^t, \{u_{d_a^t}^{t-1}, b\}, n_1} \oplus F_{d_c^t, \{u_{d_c^t}^{t-1}, u_{d_b^t}^{t-1}\}, n_3}$ ;

worker  $u_{d_b^t}^{t-1}$  transmits  $F_{d_b^t, \{u_{d_b^t}^{t-1}, a\}, n_2} \oplus F_{d_c^t, \{u_{d_c^t}^{t-1}, u_{d_b^t}^{t-1}\}, n_3}$ ,

such that each of the above linear combinations can be decoded by its requesting worker. For example in Table III, for pair  $(1, 3)$ , we let worker 2 transmit  $F_{1, \{2, 3\}, 1} \oplus F_{1, \{2, 4\}, 1}$ , and worker 4 transmit  $F_{1, \{2, 4\}, 1} \oplus F_{3, \{1, 4\}, 1}$ .

- Case  $u_{d_b^t}^{t-1} = a$ : Let  $u_{d_a^t}^{t-1} = c$  and  $u_{d_c^t}^{t-1} = d$ . For each one of  $F_{d_a^t, \{c, b\}}$ ,  $F_{d_a^t, \{c, d\}}$ , and  $F_{d_c^t, \{a, d\}}$ , we select one of its non-transmitted sub-pieces and assume they are  $F_{d_a^t, \{c, b\}, n_1}$ ,  $F_{d_a^t, \{c, d\}, n_2}$ , and  $F_{d_c^t, \{a, d\}, n_3}$ . By Proposition 1,

worker  $c$  transmits  $F_{d_a^t, \{c, b\}, n_1}$ ;

worker  $d$  transmit  $F_{d_a^t, \{c, d\}, n_2} \oplus F_{d_c^t, \{a, d\}, n_3}$ ,

such that each of the above sub-pieces can be decoded by its requesting worker. For example in Table III, for vector  $(1, 5)$ , we let worker 2 transmit  $F_{1, \{2, 5\}, 2}$  and worker 3 transmit  $F_{1, \{2, 3\}, 2} \oplus F_{2, \{1, 3\}, 1}$ . Notice that in the case if  $d = b$ , we have  $F_{d_a^t, \{c, b\}} = F_{d_a^t, \{c, d\}}$ , and thus we transmit two different sub-pieces of  $F_{d_a^t, \{c, b\}}$  for the vector  $(a, b)$ .

Next we prove that after considering all the pairs in  $\mathcal{Y}$ , each sub-piece of sub-block  $F_{d_{a_1}^t, \{u_{d_{a_1}^t}^{t-1}, b_1\}}$  has been transmitted for  $(a_1, b_1) \in \mathcal{Y}$ . For each  $(a_1, b_1) \in \mathcal{Y}$ , if there exists one worker  $c \notin \{a_1, b_1\}$  such that  $c = u_{d_{a_1}^t}^{t-1}$ ,  $b_1 = u_{d_c^t}^{t-1}$ , and  $a_1 = u_{d_{b_1}^t}^{t-1}$  (in Example 4, this pair  $(a_1, b_1)$  does not exist), we transmit two sub-pieces of  $F_{d_{a_1}^t, \{u_{d_{a_1}^t}^{t-1}, b_1\}}$  in the transmission for the pair  $(a_1, b_1)$  and one sub-piece in the transmission for the pair  $(b_1, c)$ . Otherwise, we transmit the three sub-pieces of  $F_{d_{a_1}^t, \{u_{d_{a_1}^t}^{t-1}, b_1\}}$  in the following three transmissions for three different pairs:

- 1) The transmission for the pair  $(a_1, b_1)$ .
- 2) The transmission for the pair  $(a_1, b_2)$  where  $u_{d_{b_2}^t}^{t-1} = b_1$  if the requested data unit by worker  $u_{d_{a_1}^t}^{t-1}$  in time slot  $t$  was not stored by worker  $b_1$  at the end of time slot

TABLE III: Transmission of Scheme C for Example 4

Considered vectors	Groups of sub-pieces	First sum	Second sum
(1, 3)	$F_{1,\{2,3\},1}, F_{1,\{2,4\},1}, F_{3,\{1,4\},1}$	$F_{1,\{2,3\},1} \oplus F_{1,\{2,4\},1}$	$F_{1,\{2,4\},1} \oplus F_{3,\{1,4\},1}$
(1, 4)	$F_{1,\{2,4\},2}, F_{1,\{2,5\},1}, F_{4,\{1,5\},1}$	$F_{1,\{2,4\},2} \oplus F_{1,\{2,5\},1}$	$F_{1,\{2,5\},1} \oplus F_{4,\{1,5\},1}$
(1, 5)	$F_{1,\{2,5\},2}, F_{1,\{2,3\},2}, F_{2,\{1,3\},1}$	$F_{1,\{2,5\},2}$	$F_{1,\{2,3\},2} \oplus F_{2,\{1,3\},1}$
(2, 4)	$F_{2,\{3,4\},1}, F_{2,\{3,5\},1}, F_{4,\{2,5\},1}$	$F_{2,\{3,4\},1} \oplus F_{2,\{3,5\},1}$	$F_{2,\{3,5\},1} \oplus F_{4,\{2,5\},1}$
(2, 5)	$F_{2,\{3,5\},2}, F_{2,\{1,3\},2}, F_{5,\{1,2\},1}$	$F_{2,\{3,5\},2} \oplus F_{2,\{1,3\},2}$	$F_{2,\{1,3\},2} \oplus F_{5,\{1,2\},1}$
(2, 1)	$F_{2,\{1,3\},3}, F_{2,\{3,4\},2}, F_{3,\{2,4\},1}$	$F_{2,\{1,3\},3}$	$F_{2,\{3,4\},2} \oplus F_{3,\{2,4\},1}$
(3, 5)	$F_{3,\{4,5\},1}, F_{3,\{1,4\},2}, F_{5,\{1,3\},1}$	$F_{3,\{4,5\},1} \oplus F_{3,\{1,4\},2}$	$F_{3,\{1,4\},2} \oplus F_{5,\{1,3\},1}$
(3, 1)	$F_{3,\{1,4\},3}, F_{3,\{2,4\},2}, F_{1,\{2,3\},3}$	$F_{3,\{1,4\},3} \oplus F_{3,\{2,4\},2}$	$F_{3,\{2,4\},2} \oplus F_{1,\{2,3\},3}$
(3, 2)	$F_{3,\{2,4\},3}, F_{3,\{4,5\},2}, F_{4,\{3,5\},1}$	$F_{3,\{2,4\},3}$	$F_{3,\{4,5\},2} \oplus F_{4,\{3,5\},1}$
(4, 1)	$F_{4,\{1,5\},2}, F_{4,\{2,5\},2}, F_{1,\{2,4\},3}$	$F_{4,\{1,5\},2} \oplus F_{4,\{2,5\},2}$	$F_{4,\{2,5\},2} \oplus F_{1,\{2,4\},3}$
(4, 2)	$F_{4,\{2,5\},3}, F_{4,\{3,5\},2}, F_{2,\{3,4\},3}$	$F_{4,\{2,5\},3} \oplus F_{4,\{3,5\},2}$	$F_{4,\{3,5\},2} \oplus F_{2,\{3,4\},3}$
(4, 3)	$F_{4,\{3,5\},3}, F_{4,\{1,5\},3}, F_{5,\{1,4\},1}$	$F_{4,\{3,5\},3}$	$F_{4,\{1,5\},3} \oplus F_{5,\{1,4\},1}$
(5, 2)	$F_{5,\{1,2\},2}, F_{5,\{1,3\},2}, F_{2,\{3,5\},3}$	$F_{5,\{1,2\},2} \oplus F_{5,\{1,3\},2}$	$F_{5,\{1,3\},2} \oplus F_{2,\{3,5\},3}$
(5, 3)	$F_{5,\{1,3\},3}, F_{5,\{1,4\},2}, F_{3,\{4,5\},3}$	$F_{5,\{1,3\},3} \oplus F_{5,\{1,4\},2}$	$F_{5,\{1,4\},2} \oplus F_{3,\{4,5\},3}$
(5, 4)	$F_{5,\{1,4\},3}, F_{5,\{1,2\},3}, F_{1,\{2,5\},3}$	$F_{5,\{1,4\},3}$	$F_{5,\{1,2\},3} \oplus F_{1,\{2,5\},3}$

$t - 1$ , (e.g., in Table III, let  $(a_1, b_1) = (1, 4)$ , one sub-piece of  $F_{1,\{2,4\}}$  appears in the transmission for vector  $(a_1, b_2) = (1, 3)$ ).

Otherwise, the transmission for the pair  $(a_1, b_3)$  where  $u_{d_{b_3}^{t-1}} = a_1$  (e.g., in Table III, let  $(a_1, b_1) = (1, 3)$ , one sub-piece of  $F_{1,\{2,3\}}$  appears in the transmission for vector  $(a_1, b_3) = (1, 5)$ ).

- 3) The transmission for the pair  $(b_1, a_1)$  if  $u_{d_{b_1}^{t-1}} \neq a_1$  (e.g., in Table III, let  $(a_1, b_1) = (1, 3)$ , one sub-piece of  $F_{1,\{2,3\}}$  appears in the transmission for vector  $(b_1, a_1) = (3, 1)$ ).

Otherwise, the transmission for the pair  $(b_1, b_4)$  where  $u_{d_{b_4}^{t-1}} = b_1$  (e.g., in Table III, let  $(a_1, b_1) = (1, 5)$ , one sub-piece of  $F_{1,\{2,5\}}$  appears in the transmission for vector  $(b_1, b_4) = (5, 4)$ ).

This shows that  $F_{d_{a_1}^t, \{u_{d_{a_1}^{t-1}}, b_1\}}$  is transmitted. In Scheme C, we transmit each three sub-pieces in two sums, and thus the multicast coding gain is  $2/3$ .

Finally, by comparing the loads for different cases, the worst-cases are when  $\mathcal{A}_k^{t-1} \cap \mathcal{A}_k^t = \emptyset$  for each  $k \in [K]$  and the worst-case load achieved by Scheme C is

$$q \frac{2K(K-2)}{3(K-1)} =: R_{\text{Ach.C}}|_{M=2q}, \quad (72)$$

which is optimal within a factor  $4/3$  compared to the converse bound  $\frac{K(K-2)}{2(K-1)}$  under the constraint of uncoded storage for  $M/q = 2$ .

*Storage Update Phase of time slot  $t \in [T]$ :* The storage update phase of Scheme C is the same as the one of Scheme B.

#### D. Optimality Results of the Combined Achievable Scheme

Since the proposed converse bound is a piecewise linear curve with corner points  $(mq, q \frac{K-m}{m} \frac{K}{K-1})$  for  $m \in [K]$  and these corner points are successively convex, it follows immediately that the combined scheme in Corollary 1 is optimal under the constraint of uncoded storage when  $M/q = 1$  or  $M/q \in [K-2, K]$ , thus proving Theorem 5 (and also Corollary 2).

In order to prove Theorem 6 (i.e., an order optimality result for the cases not covered by Theorem 5 or Corollary 2)

we proceed as follows. Recall that the proposed converse bound is a piecewise linear curve with successively convex corner points, and that the straight line in the storage-load tradeoff between two achievable points is also achievable by memory-sharing. Hence, in the following, we focus on each corner point of the converse bound  $(mq, q \frac{K-m}{m} \frac{K}{K-1})$  where  $m \in [K]$ , and characterize the multiplicative gap between the combined scheme and the converse when  $M = mq$ . Thus the multiplicative gap between the achievability and the converse curves is upper bounded by the maximum of the obtained  $K$  gaps.

We do not consider the corner points where  $m \in \{1, K-2, K-1, K\}$  because the optimality of the combined scheme has been proved. We have:

- $M = 2q$ : It was proved in Section VI-C that the multiplicative gap between the Scheme C and the converse bound is  $4/3$ .
- Interpolate the achievable bound for Scheme A in (21) between  $M_1 = (1 + g \frac{K-1}{K})q$  and  $M_2 = (1 + (g+1) \frac{K-1}{K})q$  to match the converse bound in (20) at  $M_3 = (g+1)q$  where  $g \in [2 : K-4]$ : For each  $g \in [2 : K-4]$ , we have

$$M_1 = \left(1 + g \frac{K-1}{K}\right)q, \quad R_{\text{Ach.A}}(M_1) = q \frac{K-g}{g}, \quad (73)$$

$$M_2 = \left(1 + (g+1) \frac{K-1}{K}\right)q, \quad R_{\text{Ach.A}}(M_2) = q \frac{K-g-1}{g+1}. \quad (74)$$

By memory-sharing between  $(M_1, R_{\text{Ach.A}}(M_1))$  and  $(M_2, R_{\text{Ach.A}}(M_2))$  with coefficient  $\alpha = (K-1-g)/(K-1)$ , we get

$$M_3 = \alpha M_1 + (1-\alpha)M_2|_{\alpha=1-g/(K-1)} = (1+g)q, \quad (75)$$

as in the converse bound for  $m = g+1 \in [3 : K-3]$ , and

$$R_{\text{Ach.A}}(M_3) = \alpha R_{\text{Ach.A}}(M_1) + (1-\alpha)R_{\text{Ach.A}}(M_2) = \frac{qK-g}{g} \frac{K-1-g}{K-1} + \frac{qg}{K-1} \frac{K-g-1}{g+1} = \frac{q(K-g-1)(Kg+K-g)}{(K-1)g(g+1)}. \quad (76)$$

From (the proof of) Theorem 1, we know that

$$R_{\text{Out}}(M_3) \geq N \frac{1 - g/(K-1)}{g+1} = q \frac{K}{K-1} \frac{K-g-1}{g+1}. \quad (77)$$

Hence, from (76) and (77), we have

$$\begin{aligned} \frac{R_{\text{Ach.A}}(M_3)}{R_{\text{Out}}(M_3)} &\leq \frac{Kg + K - g}{gK} = 1 - \frac{1}{K} + \frac{1}{g} \\ &\leq 1 - 0 + \frac{1}{2} = \frac{3}{2} \quad (\text{since } g \geq 2). \end{aligned} \quad (78)$$

We then focus on  $m q \leq M \leq (m+1)q$ , where  $m \in [K-1]$ . The converse bound in Theorem 1 for  $m q \leq M \leq (m+1)q$  is a straight line between  $(M', R') = (m q, \frac{(K-m)K}{m(K-1)})$  and  $(M'', R'') = ((m+1)q, \frac{(K-m-1)K}{(m+1)(K-1)})$ .

- $m = 1$ . The combined scheme can achieve  $(M', R')$  and  $(M'', 4R''/3)$ . Hence, by memory-sharing, the multiplicative gap between the combined scheme and the converse bound is less than  $4/3$ .
- $m = 2$ . The combined scheme can achieve  $(M', 4R'/3)$  and  $(M'', (1 - \frac{1}{K} + \frac{1}{2})R'')$ . Hence, by memory-sharing, the multiplicative gap between the combined scheme and the converse bound is less than  $1 - \frac{1}{K} + \frac{1}{2}$ .
- $m \in [3 : K-4]$ . The combined scheme can achieve  $(M', (1 - \frac{1}{K} + \frac{1}{m-1})R')$  and  $(M'', (1 - \frac{1}{K} + \frac{1}{m})R'')$ . Hence, by memory-sharing, the multiplicative gap between the combined scheme and the converse bound is less than  $1 - \frac{1}{K} + \frac{1}{m-1}$ .
- $m = K-3$ . The combined scheme can achieve  $(M', (1 - \frac{1}{K} + \frac{1}{m-1})R')$  and  $(M'', R'')$ . Hence, by memory-sharing, the multiplicative gap between the combined scheme and the converse bound is less than  $1 - \frac{1}{K} + \frac{1}{m-1}$ .
- $m \in \{K-2, K-1\}$ . The combined scheme can achieve  $(M', R')$  and  $(M'', R'')$ . Hence, the combined scheme coincides with the converse bound.

This concludes the proof of Theorem 6.

*On communication cost of peer-to-peer operations:* By comparing the decentralized data shuffling converse bound and the optimal centralized data shuffling load (denoted by  $R_{\text{Opt.Cen}}(M)$ ), we have  $R_{\text{Out}}(M)/R_{\text{Opt.Cen}}(M) = K/(K-1)$  for any  $q \leq M \leq Kq$ . In addition, the maximum multiplicative gap between the achieved load by the combined scheme and  $R_{\text{Out}}(M)$ , is  $\max\left\{1 - \frac{1}{K} + \frac{1}{m-1}, \frac{4}{3}\right\}$ , where  $m \geq 3$ . Hence, the maximum multiplicative gap between the achieved load by the combined scheme and  $R_{\text{Opt.Cen}}(M)$  is

$$\begin{aligned} &\frac{K}{K-1} \max\left\{1 - \frac{1}{K} + \frac{1}{m-1}, \frac{4}{3}\right\} \\ &= \max\left\{1 + \frac{K}{(m-1)(K-1)}, \frac{4K}{3(K-1)}\right\}, \end{aligned} \quad (79)$$

which is no more than  $5/3$  if  $K \geq 5$ . In addition, when  $K \leq 4$ , by Corollary 2, the combined scheme is optimal such that the communication cost of peer-to-peer operations is  $K/(K-1) \leq 2$ . In general, the communication cost of peer-to-peer operations is no more than a factor of 2 as stated in Corollary 3.

In addition, with a similar proof as above to analyse each storage size regime  $m q \leq M \leq (m+1)q$  where  $m \in [K-1]$ , we prove Corollary 3.

## VII. CONCLUSIONS

In this paper, we introduced the decentralized data shuffling problem and studied its fundamental limits. We proposed a converse bound under the constraint of uncoded storage and three achievable schemes. In general, under the constraint of uncoded storage, our schemes are order optimal to within a factor of  $3/2$ , and exactly optimal for small and large storage sizes, or the systems with no more than four workers.

## APPENDIX

We define

$$\mathcal{V}_{\mathcal{S}} := \{k \in \mathcal{S} : d_k \in \mathcal{S}\}, \quad \forall \mathcal{S} \subseteq [K], \quad (80)$$

where  $\mathcal{V}_{\mathcal{S}}$  represents the subset of workers in  $\mathcal{S}$  whose demanded data units in time slot  $t$ , indexed by  $\mathcal{A}_k^t = \mathcal{A}_{d_k}^{t-1}$  in (29), were stored by some workers in  $\mathcal{S}$  at the end of time slot  $t-1$ <sup>5</sup>.

In addition, we also define  $Y_{\mathcal{S}}^t$  as the sub-blocks that any worker in  $\mathcal{S}$  either needs to store at the end of time slot  $t$  or has stored at the end of time slot  $t-1$ , that is,

$$\begin{aligned} Y_{\mathcal{S}}^t &:= \left\{F_i : i \in \bigcup_{k \in \mathcal{S}} \mathcal{A}_k^t\right\} \cup \left\{Z_k^{t-1} : k \in \mathcal{S}\right\} \\ &= \left\{F_i : i \in \bigcup_{k \in \mathcal{S}} (\mathcal{A}_k^t \cup \mathcal{A}_k^{t-1})\right\} \\ &\cup \left\{F_{i, \mathcal{W}} : i \notin \bigcup_{k \in \mathcal{S}} (\mathcal{A}_k^t \cup \mathcal{A}_k^{t-1}), \mathcal{W} \cap \mathcal{S} \neq \emptyset\right\}. \end{aligned} \quad (81)$$

With the above definitions and recalling  $X_{\mathcal{S}}^t$  defined in (30) represents the messages sent by the workers in  $\mathcal{S}$  during time slot  $t$ , we have the following lemma:

**Lemma 1** (Induction Lemma). *For each non-empty set  $\mathcal{S} \subseteq [K]$ , we have*

$$\begin{aligned} &H\left(X_{\mathcal{S}}^t | Y_{[K] \setminus \mathcal{S}}^t\right) \\ &\geq \sum_{m=1}^{|\mathcal{S}|} \sum_{k \in \mathcal{V}_{\mathcal{S}}} \sum_{i \in \mathcal{A}_k^t} \sum_{\substack{\mathcal{W} \subseteq \mathcal{S} \setminus \{k\}: \\ u_i^{t-1} \in \mathcal{W}, |\mathcal{W}|=m}} \frac{|F_{i, \mathcal{W}}|}{m}. \end{aligned} \quad (82)$$

Lemma 1 is the key novel contribution of our proof. The bound in (82) can be intuitively explained as follows:  $H(X_{\mathcal{S}}^t | Y_{[K] \setminus \mathcal{S}}^t)$  is lower bounded by the size of the requested sub-blocks by the workers in  $\mathcal{V}_{\mathcal{S}}$  (instead of in  $\mathcal{S}$  as in the distributed computing problem [14]) because each requested data unit by the workers in  $\mathcal{S} \setminus \mathcal{V}_{\mathcal{S}}$  was requested in the previous time slot by some workers in  $[K] \setminus \mathcal{S}$  because of the storage constraint in (8) and the definition of  $\mathcal{V}_{\mathcal{S}}$  in (80).

This lemma is proved by induction, inspired by [14].

*Proof:* Case  $|\mathcal{S}| = 1$ : If  $\mathcal{S} = \{k\}$  where  $k \in [K]$ , we have that  $\mathcal{V}_{\{k\}} = \emptyset$  and thus the RHS of (82) is 0; thus (82) holds for  $|\mathcal{S}| = 1$  because entropy is non-negative.

<sup>5</sup>For example, if  $K = 4$  and  $(d_1, \dots, d_4) = (2, 3, 4, 1)$ , we have  $\mathcal{V}_{\{2,3\}} = \{2\}$  because  $d_2 = 3$  and thus the requested data unit by worker 2 in time slot  $t$  was stored by worker 3 at the end of time slot  $t-1$ ; similarly, we have  $\mathcal{V}_{\{2,4\}} = \emptyset$  and  $\mathcal{V}_{\{1,2,4\}} = \{1, 4\}$ .

Case  $|\mathcal{S}| \leq s$ : Assume that (82) holds for all non-empty  $\mathcal{S} \subseteq [K]$  where  $|\mathcal{S}| \leq s$  for some integer  $s \in [K-1]$ .

Case  $|\mathcal{S}| = s+1$ : Having assumed that the lemma holds for all  $\mathcal{S} \subseteq [K]$  where  $|\mathcal{S}| \leq s$ , we aim to show that for any set  $\mathcal{J} \subseteq [K]$  where  $|\mathcal{J}| = s+1$ , we have

$$H(X_{\mathcal{J}}^t | Y_{[K] \setminus \mathcal{J}}^t) \geq \sum_{m=1}^{|\mathcal{J}|} \sum_{k \in \mathcal{V}_{\mathcal{J}}} \sum_{i \in \mathcal{A}_k^t} \sum_{\substack{\mathcal{W} \subseteq (\mathcal{J} \setminus \{k\}): \\ |\mathcal{W}| = m, u_i^{t-1} \in \mathcal{W}}} \frac{|F_{i, \mathcal{W}}|}{m}. \quad (83)$$

From the independence bound on entropy we have

$$\begin{aligned} H(X_{\mathcal{J}}^t | Y_{[K] \setminus \mathcal{J}}^t) &= \frac{1}{|\mathcal{J}|} \sum_{k \in \mathcal{J}} \left( H(X_{\mathcal{J} \setminus \{k\}}^t | X_k^t, Y_{[K] \setminus \mathcal{J}}^t) + H(X_k^t | Y_{[K] \setminus \mathcal{J}}^t) \right) \end{aligned} \quad (84a)$$

$$\geq \frac{1}{|\mathcal{J}|} \left( \sum_{k \in \mathcal{J}} H(X_{\mathcal{J} \setminus \{k\}}^t | X_k^t, Y_{[K] \setminus \mathcal{J}}^t) + H(X_{\mathcal{J}}^t | Y_{[K] \setminus \mathcal{J}}^t) \right), \quad (84b)$$

and thus

$$(|\mathcal{J}| - 1)H(X_{\mathcal{J}}^t | Y_{[K] \setminus \mathcal{J}}^t) \geq \sum_{k \in \mathcal{J}} H(X_{\mathcal{J} \setminus \{k\}}^t | X_k^t, Y_{[K] \setminus \mathcal{J}}^t) \quad (84c)$$

$$\geq \sum_{k \in \mathcal{J}} H(X_{\mathcal{J} \setminus \{k\}}^t | X_k^t, Y_{[K] \setminus \mathcal{J}}^t, Z_k^{t-1}) \quad (84d)$$

$$= \sum_{k \in \mathcal{J}} H(X_{\mathcal{J} \setminus \{k\}}^t, \{F_i : i \in \mathcal{A}_k^t\} | X_k^t, Y_{[K] \setminus \mathcal{J}}^t, Z_k^{t-1}) \quad (84e)$$

$$= \sum_{k \in \mathcal{J}} H(\{F_i : i \in \mathcal{A}_k^t\} | Z_k^{t-1}, Y_{[K] \setminus \mathcal{J}}^t) \quad (84f)$$

$$+ \sum_{k \in \mathcal{J}} H(X_{\mathcal{J} \setminus \{k\}}^t | \{F_i : i \in \mathcal{A}_k^t\}, Z_k^{t-1}, Y_{[K] \setminus \mathcal{J}}^t) \quad (84f)$$

$$= \sum_{k \in \mathcal{J}} H(\{F_i : i \in \mathcal{A}_k^t\} | Z_k^{t-1}, Y_{[K] \setminus \mathcal{J}}^t) \quad (84g)$$

$$+ \sum_{k \in \mathcal{J}} H(X_{\mathcal{J} \setminus \{k\}}^t | Y_{([K] \setminus \mathcal{J}) \cup \{k\}}^t), \quad (84g)$$

where (84d) follows because we added  $Z_k^{t-1}$  in the conditioning, and conditioning cannot increase entropy, where (84e) follows because  $\{F_i : i \in \mathcal{A}_k^t\}$  is a function of  $(Z_k^{t-1}, X^t)$  by the decoding constraint in (5) (note that the knowledge of  $(Y_{[K] \setminus \mathcal{J}}^t, Z_k^{t-1})$  implies the knowledge of  $Z_{\{k\} \cup [K] \setminus \mathcal{J}}^{t-1}$  and thus of  $X_{\{k\} \cup [K] \setminus \mathcal{J}}^t$  by the encoding constraint in (4)), where (84f) follow because  $X_k^t$  is a function of  $Z_k^{t-1}$  (see the encoding constraint in (4)), and (84g) from the definition in (81).

Next we bound the first term of (84g) by using the independence of the sub-blocks, and the second term of (84g) by the induction assumption. More precisely,

- First term of (84g). For each  $k \in \mathcal{J}$ , if  $k \notin \mathcal{V}_{\mathcal{J}}$ , we have  $\{F_i : i \in \mathcal{A}_k^t\} \subseteq Y_{[K] \setminus \mathcal{J}}^t$ . So for each  $k \in \mathcal{J}$ , by independence of sub-blocks, we have (85) and thus we rewrite the first term of (84g) as

$$\sum_{k \in \mathcal{J}} H(\{F_i : i \in \mathcal{A}_k^t\} | Z_k^{t-1}, Y_{[K] \setminus \mathcal{J}}^t)$$

$$= \sum_{k \in \mathcal{V}_{\mathcal{J}}} \sum_{m \in [|\mathcal{J}|]} \sum_{i \in \mathcal{A}_k^t} \sum_{\substack{\mathcal{W} \subseteq (\mathcal{J} \setminus \{k\}): \\ |\mathcal{W}| = m, u_i^{t-1} \in \mathcal{W}}} |F_{i, \mathcal{W}}|. \quad (86)$$

- Second term of (84g). By the induction assumption,

$$\begin{aligned} &\sum_{k \in \mathcal{J}} H(X_{\mathcal{J} \setminus \{k\}}^t | Y_{([K] \setminus \mathcal{J}) \cup \{k\}}^t) \\ &= \sum_{k \in \mathcal{J}} H(X_{\mathcal{J} \setminus \{k\}}^t | Y_{[K] \setminus (\mathcal{J} \setminus \{k\})}^t) \\ &\geq \sum_{k \in \mathcal{J}} \sum_{u \in \mathcal{V}_{\mathcal{J} \setminus \{k\}}} \sum_{m=1}^{|\mathcal{J}|-1} \sum_{i \in \mathcal{A}_u^t} \sum_{\substack{\mathcal{W} \subseteq (\mathcal{J} \setminus \{k, u\}): \\ |\mathcal{W}| = m, u_i^{t-1} \in \mathcal{W}}} \frac{|F_{i, \mathcal{W}}|}{m}. \end{aligned} \quad (87)$$

In order to combine (86) with (87), both terms need to have the summations in the same form. Let us focus on one worker  $u' \in \mathcal{V}_{\mathcal{J}}$  and one sub-block  $F_{i', \mathcal{W}'}$ , where  $i' \in \mathcal{A}_{u'}^t$ ,  $\mathcal{W}' \subseteq \mathcal{J} \setminus \{u'\}$ ,  $|\mathcal{W}'| = m$ , and  $u_{i'}^{t-1} \in \mathcal{W}'$ . On the RHS of (87), for each  $k \in \mathcal{J} \setminus (\mathcal{W}' \cup \{u'\})$ , it can be seen that  $F_{i', \mathcal{W}'}$  appears once in the sum

$$\sum_{m \in [|\mathcal{J}|-1]} \sum_{u \in \mathcal{V}_{\mathcal{J} \setminus \{k\}}} \sum_{i \in \mathcal{A}_u^t} \sum_{\substack{\mathcal{W} \subseteq (\mathcal{J} \setminus \{k, u\}): \\ |\mathcal{W}| = m, u_i^{t-1} \in \mathcal{W}}} \frac{|F_{i, \mathcal{W}}|}{m}, \quad (88)$$

hence, the coefficient of  $F_{i', \mathcal{W}'}$  on the RHS of (87) is  $(|\mathcal{J}| - m - 1)/m$ . Thus, from (87), we have

$$\sum_{k \in \mathcal{J}} H(X_{\mathcal{J} \setminus \{k\}}^t | Y_{([K] \setminus \mathcal{J}) \cup \{k\}}^t) \quad (89a)$$

$$\geq \sum_{u' \in \mathcal{V}_{\mathcal{J}}} \sum_{m \in [|\mathcal{J}|-1]} \sum_{i' \in \mathcal{A}_{u'}^t} \sum_{\substack{\mathcal{W}' \subseteq (\mathcal{J} \setminus \{u'\}): \\ |\mathcal{W}'| = m, u_{i'}^{t-1} \in \mathcal{W}'}} \frac{|F_{i', \mathcal{W}'}| (|\mathcal{J}| - m - 1)}{m} \quad (89b)$$

$$= \sum_{u' \in \mathcal{V}_{\mathcal{J}}} \sum_{m \in [|\mathcal{J}|]} \sum_{i' \in \mathcal{A}_{u'}^t} \sum_{\substack{\mathcal{W}' \subseteq (\mathcal{J} \setminus \{u'\}): \\ |\mathcal{W}'| = m, u_{i'}^{t-1} \in \mathcal{W}'}} \frac{|F_{i', \mathcal{W}'}| (|\mathcal{J}| - m - 1)}{m}. \quad (89c)$$

We next take (86) and (89c) into (84g) to obtain,

$$\begin{aligned} &H(X_{\mathcal{J}}^t | Y_{[K] \setminus \mathcal{J}}^t) \\ &\geq \frac{1}{|\mathcal{J}| - 1} \sum_{k \in \mathcal{V}_{\mathcal{J}}} \sum_{m=1}^{|\mathcal{J}|} \sum_{i \in \mathcal{A}_k^t} \sum_{\substack{\mathcal{W} \subseteq (\mathcal{J} \setminus \{k\}): \\ |\mathcal{W}| = m, u_i^{t-1} \in \mathcal{W}}} |F_{i, \mathcal{W}}| \\ &+ \frac{1}{|\mathcal{J}| - 1} \sum_{k \in \mathcal{V}_{\mathcal{J}}} \sum_{m=1}^{|\mathcal{J}|} \sum_{i \in \mathcal{A}_k^t} \sum_{\substack{\mathcal{W} \subseteq (\mathcal{J} \setminus \{k\}): \\ |\mathcal{W}| = m, u_i^{t-1} \in \mathcal{W}}} \frac{|F_{i, \mathcal{W}}| (|\mathcal{J}| - m - 1)}{m} \end{aligned} \quad (90a)$$

$$= \sum_{k \in \mathcal{V}_{\mathcal{J}}} \sum_{m \in [|\mathcal{J}|]} \sum_{i \in \mathcal{A}_k^t} \sum_{\substack{\mathcal{W} \subseteq (\mathcal{J} \setminus \{k\}): \\ |\mathcal{W}| = m, u_i^{t-1} \in \mathcal{W}}} \frac{|F_{i, \mathcal{W}}|}{m}, \quad (90b)$$

which proves Lemma 1.  $\blacksquare$

$$\begin{aligned}
& H(\{F_i : i \in \mathcal{A}_k^t\} | Z_k^{t-1}, Y_{[K] \setminus \mathcal{J}}^t) \\
& = \begin{cases} \sum_{m=1}^{|\mathcal{J}|} \sum_{i \in \mathcal{A}_k^t} \sum_{\mathcal{W} \subseteq (\mathcal{J} \setminus \{k\}) : |\mathcal{W}|=m, u_i^{t-1} \in \mathcal{W}} |F_i, \mathcal{W}|, & k \in \mathcal{V}_{\mathcal{J}} \\ 0 & \text{otherwise} \end{cases} \quad (85)
\end{aligned}$$

## REFERENCES

- [1] J. Chung, K. Lee, R. Pedarsani, D. Papailiopoulos, and K. Ramchandran, "Ubershuffle: Communication-efficient data shuffling for sgd via coding theory," in *NIPS 2017, ML Systems Workshop*.
- [2] M. Gurbuzbalaban, A. Ozdaglar, and P. Parrilo, "Why random reshuffling beats stochastic gradient descent," available at *arXiv:1510.08560*, Oct. 2015.
- [3] 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, Mar. 2018.
- [4] F. A. Tobagi and V. B. Hunt, "Performance analysis of carrier sense multiple access with collision detection," *Computer Networks*, vol. 4, pp. 245–259, 1980.
- [5] ANSI/IEEE Standard 802.5, "Token ring access method and physical layer specifications," *IEEE Press*, 1985.
- [6] M. A. Attia and R. Tandon, "Information theoretic limits of data shuffling for distributed learning," in *IEEE Global Commun. Conf. (GLOBECOM)*, Dec. 2016.
- [7] —, "On the worst-case communication overhead for distributed data shuffling," in *54th Annual Allerton Conf. on Commun., Control, and Computing (Allerton)*, Sep. 2016.
- [8] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Trans. Infor. Theory*, vol. 60, no. 5, pp. 2856–2867, May 2014.
- [9] K. Wan, D. Tuninetti, and P. Piantanida, "On the optimality of uncoded cache placement," in *IEEE Infor. Theory Workshop*, Sep. 2016.
- [10] Q. Yu, M. A. Maddah-Ali, and S. Avestimehr, "The exact rate-memory tradeoff for caching with uncoded prefetching," *IEEE Trans. Infor. Theory*, vol. 64, pp. 1281 – 1296, Feb. 2018.
- [11] M. A. Attia and R. Tandon, "Near optimal coded data shuffling for distributed learning," vol. 65, no. 11, pp. 7325–7349, Nov. 2019.
- [12] A. Elmahdy and S. Mohajer, "On the fundamental limits of coded data shuffling for distributed learning systems," *arXiv:1807.04255, a preliminary version is in IEEE Int. Symp. Inf. Theory 2018*, Jul. 2018.
- [13] M. Ji, G. Caire, and A. Molisch, "Fundamental limits of caching in wireless d2d networks," *IEEE Trans. Inf. Theory*, vol. 62, no. 1, pp. 849–869, 2016.
- [14] 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.
- [15] Y. H. Ezzeldin, M. Karmoose, and C. Fragouli, "Communication vs distributed computation: An alternative trade-off curve," in *IEEE Inf. Theory Workshop (ITW)*, Nov. 2017.
- [16] N. Woolsey, R. Chen, and M. Ji, "A new combinatorial design of coded distributed computing," in *IEEE Int. Symp. Inf. Theory*, Jun. 2018.
- [17] K. Konstantinos and A. Ramamoorthy, "Leveraging coding techniques for speeding up distributed computing," available at *arXiv:1802.03049*, Feb. 2018.
- [18] S. R. Srinivasavaradhan, L. Song, and C. Fragouli, "Distributed computing trade-offs with random connectivity," in *IEEE Int. Symp. Inf. Theory*, June. 2018.
- [19] S. Prakash, A. Reiszadeh, R. Pedarsani, and S. Avestimehr, "Coded computing for distributed graph analytics," in *IEEE Int. Symp. Inf. Theory*, June. 2018.
- [20] B. Guler, A. S. Avestimehr, and A. Ortega, "A topology-aware coding framework for distributed graph processing," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2019.
- [21] L. Songze, Q. Yu, M. A. Maddah-Ali, and A. S. Avestimehr, "A scalable framework for wireless distributed computing," *IEEE/ACM Transactions on Networking*, vol. 25, no. 5, pp. 2643–2654, Oct. 2017.
- [22] Q. Yan, S. Yang, and M. Wigger, "Storage, computation, and communication: A fundamental tradeoff in distributed computing," *arXiv:1806.07565*, Jun. 2018.
- [23] N. Woolsey, R. Chen, and M. Ji, "Cascaded coded distributed computing on heterogeneous networks," *arXiv:1901.07670*, Jan. 2019.
- [24] —, "Coded distributed computing with heterogeneous function assignments," *arXiv:1902.10738*, Feb. 2019.
- [25] P. Sadeghi, F. Arbabjolfaci, and Y. Kim, "Distributed index coding," in *IEEE Inf. Theory Workshop 2016*, Sep. 2016.
- [26] Y. Liu, P. Sadeghi, F. Arbabjolfaci, and Y. Kim, "Capacity theorems for distributed index coding," *arXiv:1801.09063*, Jan. 2018.
- [27] Y. Birk and T. Kol, "Informed source coding on demand (iscod) over broadcast channels," in *Proc. IEEE Conf. Comput. Commun.*, pp. 1257–1264, 1998.
- [28] A. Porter and M. Wootters, "Embedded index coding," available at *arXiv:1904.02179*, Apr. 2019.
- [29] C. Thapa, L. Ong, S. J. Johnson, and M. Li, "Structural characteristics of two-sender index coding," *arXiv:1711.08150v2*, Jun. 2019.
- [30] A. E. Gamal and Y.-H. Kim, *Network Information Theory*. Cambridge, UK: Cambridge University Press, 2011.

**Kai Wan** (S '15 – M '18) received the M.Sc. and Ph.D. degrees in Communications from Université Paris Sud-CentraleSupélec, France, in 2014 and 2018. He is currently a post-doctoral researcher with the Communications and Information Theory Chair (CommIT) at Technische Universität Berlin, Berlin, Germany. His research interests include coded caching, index coding, distributed storage, wireless communications, and distributed computing.

**Daniela Tuninetti** (M '98 – SM '13) is currently a Professor within the Department of Electrical and Computer Engineering at the University of Illinois at Chicago (UIC), which she joined in 2005. Dr. Tuninetti got her Ph.D. in Electrical Engineering in 2002 from ENST/Télécom ParisTech (Paris, France, with work done at the Eurecom Institute in Sophia Antipolis, France), and she was a postdoctoral research associate at the School of Communication and Computer Science at the Swiss Federal Institute of Technology in Lausanne (EPFL, Lausanne, Switzerland) from 2002 to 2004. Dr. Tuninetti is a recipient of a best paper award at the European Wireless Conference in 2002, of an NSF CAREER award in 2007, and named University of Illinois Scholar in 2015. Dr. Tuninetti was the editor-in-chief of the IEEE Information Theory Society Newsletter from 2006 to 2008, an editor for IEEE COMMUNICATION LETTERS from 2006 to 2009, for IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS from 2011 to 2014; and for IEEE TRANSACTIONS ON INFORMATION THEORY from 2014 to 2017. She is currently a distinguished lecturer for the Information Theory society. Dr. Tuninetti's research interests are in the ultimate performance limits of wireless interference networks (with special emphasis on cognition and user cooperation), coexistence between radar and communication systems, multi-relay networks, content-type coding, cache-aided systems and distributed private coded computing.

**Mingyue Ji** (S '09 – M '15) received the B.E. in Communication Engineering from Beijing University of Posts and Telecommunications (China), in 2006, the M.Sc. degrees in Electrical Engineering from Royal Institute of Technology (Sweden) and from University of California, Santa Cruz, in 2008 and 2010, respectively, and the PhD from the Ming Hsieh Department of Electrical Engineering at University of Southern California in 2015. He subsequently was a Staff II System Design Scientist with Broadcom Corporation (Broadcom Limited) in 2015-2016. He is now an Assistant Professor of Electrical and Computer Engineering Department and an Adjunct Assistant Professor of School of Computing at the University of Utah. He received the IEEE Communications Society Leonard G. Abraham Prize for the best IEEE JSAC paper in 2019, the best paper award in IEEE ICC 2015 conference, the best student paper award in IEEE European Wireless 2010 Conference and USC Annenberg Fellowship from 2010 to 2014. He is interested in the broad area of information theory, coding theory, concentration of measure and statistics with the applications of caching networks, wireless communications, distributed computing and storage, security and privacy and (statistical) signal processing.

**Giuseppe Caire** (S '92 – M '94 – SM '03 – F '05) was born in Torino in 1965. He received the B.Sc. in Electrical Engineering from Politecnico di Torino in 1990, the M.Sc. in Electrical Engineering from Princeton University in 1992, and the Ph.D. from Politecnico di Torino in 1994. He has been a post-doctoral research fellow with the European Space Agency (ESTEC, Noordwijk, The Netherlands) in 1994-1995, Assistant Professor in Telecommunications at the Politecnico di Torino, Associate Professor at the University of Parma, Italy, Professor with the Department of Mobile Communications at the Eurecom Institute, Sophia-Antipolis, France, a Professor of Electrical Engineering with the Viterbi School of Engineering, University of Southern California, Los Angeles, and he is currently an Alexander von Humboldt Professor with the Faculty of Electrical Engineering and Computer Science at the Technical University of Berlin, Germany.

He received the Jack Neubauer Best System Paper Award from the IEEE Vehicular Technology Society in 2003, the IEEE Communications Society & Information Theory Society Joint Paper Award in 2004 and in 2011, the Leonard G. Abraham Prize for best IEEE JSAC paper in 2019, the Okawa Research Award in 2006, the Alexander von Humboldt Professorship in 2014, the Vodafone Innovation Prize in 2015, and an ERC Advanced Grant in 2018. Giuseppe Caire is a Fellow of IEEE since 2005. He has served in the Board of Governors of the IEEE Information Theory Society from 2004 to 2007, and as officer from 2008 to 2013. He was President of the IEEE Information Theory Society in 2011. His main research interests are in the field of communications theory, information theory, channel and source coding with particular focus on wireless communications.

**Pablo Piantanida** (SM '16) received both B.Sc. in Electrical Engineering and the M.Sc. (with honors) from the University of Buenos Aires (Argentina) in 2003, and the Ph.D. from Université Paris-Sud (Orsay, France) in 2007. Since October 2007 he has joined the Laboratoire des Signaux et Systèmes (L2S), at CentraleSupélec together with CNRS (UMR 8506) and Université Paris-Sud, as an Associate Professor of Network Information Theory. He is currently associated with Montreal Institute for Learning Algorithms (Mila) at Université de Montréal, Quebec, Canada. He is an IEEE Senior Member and serves as Associate Editor for IEEE Transactions on Information Forensics and Security. He served as General Co-Chair of the 2019 IEEE International Symposium on Information Theory (ISIT). His research interests lie broadly in information theory and its interactions with other fields, including multi-terminal information theory, Shannon theory, machine learning and representation learning, statistical inference, cooperative communications, communication mechanisms for security and privacy.