



# MIT Open Access Articles

## *Scheduling Advantages of Network Coded Storage in Point-to-Multipoint Networks*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

<b>Citation</b>	Ferner, Ulric J., Sadeghi, Parastoo, Aboutorab, Neda and Medard, Muriel. 2014. "Scheduling Advantages of Network Coded Storage in Point-to-Multipoint Networks."
<b>As Published</b>	10.1109/netcod.2014.6892120
<b>Publisher</b>	Institute of Electrical and Electronics Engineers (IEEE)
<b>Version</b>	Original manuscript
<b>Citable link</b>	<a href="https://hdl.handle.net/1721.1/137823">https://hdl.handle.net/1721.1/137823</a>
<b>Terms of Use</b>	Creative Commons Attribution-Noncommercial-Share Alike
<b>Detailed Terms</b>	<a href="http://creativecommons.org/licenses/by-nc-sa/4.0/">http://creativecommons.org/licenses/by-nc-sa/4.0/</a>

# Scheduling Advantages of Network Coded Storage in Point-to-Multipoint Networks

Ulric J. Ferner<sup>†</sup>, Parastoo Sadeghi<sup>‡</sup>, Neda Aboutorab<sup>‡</sup>, and Muriel Médard<sup>†</sup>

<sup>†</sup>Research Laboratory for Electronics, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

<sup>‡</sup>Research School of Engineering, Australian National University, Canberra ACT 0200, Australia

**Abstract**—We consider scheduling strategies for point-to-multipoint (PMP) storage area networks (SANs) that use network coded storage (NCS). In particular, we present a simple SAN system model, two server scheduling algorithms for PMP networks, and analytical expressions for internal and external blocking probability. We point to select scheduling advantages in NCS systems under *normal* operating conditions, where content requests can be temporarily denied owing to finite system capacity from drive I/O access or storage redundancy limitations. NCS can lead to improvements in throughput and blocking probability due to increased immediate scheduling options, and complements other well documented NCS advantages such as regeneration, and can be used as a guide for future storage system design.

## I. INTRODUCTION

The prolific growth of online content and streaming video makes serving content requests to multiple users simultaneously an important technique for modern storage area networks (SANs). Two fundamental measures of service quality are system external blocking probability, i.e., the probability that a requesting user is denied immediate access to content, as well as system throughput. Under *normal* operating conditions and given perfect scheduling, network coded storage (NCS) has been identified as a promising technique to reduce blocking probability. For instance, [1] used queuing theory to show that network coding can reduce system blocking probability. In this paper we build upon this idea and develop simple and intuitive server scheduling algorithms for such NCS systems. We then explore their impact on both throughput as well as blocking probability. The main contributions of this paper are:

- We introduce a simple storage model for point-to-multipoint (PMP) storage networks that allows direct evaluation of blocking probability and system throughput;
- Using this model, we propose two intuitive scheduling algorithms—one for uncoded storage (UCS) and one for NCS—that can achieve maximal throughput;
- We quantify the blocking probability and throughput savings of NCS over UCS scheduling, showing that a small improvement in throughput translates to a comparatively large improvement in blocking probability.

This paper builds upon and complements existing work in this area. The use of NCS as regenerating codes is a well studied repair technique to enhance SAN reliability [2] in both centralized and distributed systems. This particularly holds in *less common* operating conditions, such as permanent drive failures. In modern systems *traffic-induced* temporary unavailability significantly dominates disk failures [3], and so

like in [1], this paper focusses on *normal* operating conditions and seeks to avoid highly transient and temporary bottlenecks in data liveness. General scheduling for coded storage in point-to-point networks, when users are served sequentially instead of simultaneously, are considered in [4], [5].

Server scheduling is also well studied in matched networks such as cross-bar switches. Throughput-optimal schedules are considered for  $N \times M$  point-to-point cross-bar switches using graph theory and techniques such as the Birkhoff-von Neumann theorem [6]. Switches with multicast and broadcast capabilities with a queuing analysis flavor are considered in [7]. References [8], [9] attempt to map the multicast problem in cross-bar switches to simpler problems such as block-packing games and round-robin based multicast. By characterizing flow conflict graphs and their corresponding stable set polytopes in multicast cross-bar switches, [18] proposed online and offline network coding schedules for enhancing throughput. For general PMP storage networks, developing appropriate storage models, corresponding conflict graphs, and throughput optimal scheduling is an interesting and largely unaddressed area of research. This paper takes a first step towards this by considering a particular kind of PMP network, namely broadcast, and by developing intuitive coded and uncoded leader-based scheduling, which do not explicitly require conflict graph construction. Chunk scheduling problems in uncoded peer-to-peer networks, as opposed to PMPs, are considered in [10], and for star-based broadcast networks in [11]. Note also that unlike classical asynchronous broadcast problems [12], [13], our goal is not to reduce content delivery delay or to optimize caching. Instead, by taking into account intermittent drive availability, we aim to determine the impact of scheduling drive reads and the impact of content storage format on blocking probability and throughput. We expect that by using appropriate caching, system performance can be further improved. However, this is beyond the scope of this work.

The remainder of this paper is organized as follows. Section II details our system model. Section III presents service schemes and Section IV describes numerical results. Section V concludes the paper.

## II. SYSTEM MODEL

Fig. 1 depicts our tree-structured connectivity model made of single server  $S$ , connected to  $R$  drives, that receives user requests for content.

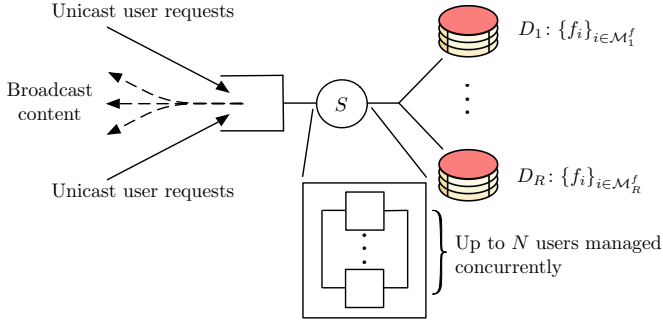


Fig. 1. System model.

### A. Drives

The SAN in Fig. 1 stores a single chunked file  $\mathcal{F} = \{f_1, f_2, \dots, f_T\}$ , where  $T$  is the number of chunks in  $\mathcal{F}$ , and  $\mathcal{F}$  is stored across a set of drives  $\mathcal{D} = \{D_1, \dots, D_R\}$ . If drive  $D_i$  receives a read request for chunk  $f_j$ , and if  $D_i$  stores  $f_j$  and  $D_i$  is available, then it takes one timeslot to read out that content and broadcast to all users. We model the overall effect of drives having finite I/O access bandwidth with parameter  $P_b^D$ , where  $P_b^D$  is the probability that any drive is blocked in timeslot  $t$ . For simplicity, we assume drives are blocked independently of one another and across timeslots.<sup>1</sup>

In UCS, let  $\mathcal{M}_i^d \subseteq \mathcal{D}$  be the collection of drives that hold uncoded file chunk  $f_i$  and conversely, let  $\mathcal{M}_i^f \subset \mathcal{F}$  be the collection of file chunks held by drive  $D_i$ . The only requirement of chunks to drives is that  $R$  drives collectively hold at least one copy of the whole file, i.e.,  $\mathcal{F} = \cup_{i=1}^R \mathcal{M}_i^f$ .

In NCS, the  $r$ th coded file chunk is represented as [1]

$$c_r = \sum_{j=1}^T \alpha_{j,r} f_j \quad (1)$$

where  $\alpha_{j,r}$  is the encoding coefficient of file chunk  $f_j$  and the corresponding encoding vector is

$$\mathbf{k}_r = \sum_{j=1}^T \alpha_{j,r} \mathbf{e}_j \quad (2)$$

In (2),  $\mathbf{e}_j = [e_{j,1}, \dots, e_{j,T}]$  is the unit encoding row vector of length  $T$  with elements  $e_{j,r} = \delta_{j,r}$ . Function  $\delta_{j,r}$  is the Kronecker delta function with  $\delta_{j,r} = 1$  iff  $j = r$ . We assume that a total of  $H$  linearly coded chunks  $c_1$  to  $c_H$  are stored onto drives via some MDS code, such that any  $T$  coded chunks are linearly independent so that the original file chunks can be recovered from them using Gaussian elimination. If encoding coefficients  $\alpha_{j,r}$  are randomly selected from a finite field  $\mathbb{F}_q$  with sufficiently large size  $q$ , this requirement is satisfied with high probability [14].

<sup>1</sup>This blocking model can be applicable where other servers have access to the same drives and therefore, there is some probabilistic traffic-induced blocking observed by  $S$ . More realistic models for traffic-induced drive blocking as well as more general PMP traffic patterns are beyond the scope of this paper and subject of our current research.

### B. Server

We assume server  $S$  has a bounded buffer to manage concurrent user requests. Let  $N$  be the maximum number of users that can be managed and serviced concurrently and suppose  $S$  operates in slotted time. In particular, in any timeslot,  $S$  can serve at most  $N$  active requests for content. A user request for a content is cleared from the buffer when all its requested file chunks have been transmitted by  $S$ .

Any additional request beyond  $N$  for the same content will be *externally blocked*.<sup>2</sup> We will discuss the relation between external and *internal* blocking in Section II-D. This is a similar model to existing drive blocking models [1] and existing practical server experimentation test [15]. When a user request arrives and is not externally blocked, one slot of the server buffer is allocated to manage and service this user request. We make the following additional assumptions about how  $S$  retrieves content from  $\mathcal{D}$ :

- Let the vector  $\mathbf{b}(t)$  of size  $R$  be the drive availability vector, where  $b_i(t) = 0$  means drive  $D_i$  is free for reads and  $b_i(t) = 1$  means it is busy in timeslot  $t$ . We assume that  $\mathbf{b}(t)$  can be obtained by the server at the beginning of timeslot  $t$  with negligible time overhead.
- In timeslot  $t$ , based on  $\mathbf{b}(t)$ ,  $S$  can choose to send a read request to access a single drive and read a single chunk.
- At the end of timeslot  $t$ ,  $S$  broadcasts the received chunk  $x(t)$  to users active in the buffer.

We assume perfect communication so when  $S$  broadcasts content all active users receive that content without error.

### C. Users

We model users with the following key parameters:

- User requests arrives at  $S$  following a Poisson process with rate  $\lambda$ .
- All user requests are for the entire file  $\mathcal{F}$ , so in the long-term there is uniform traffic demand across file chunks.
- Users currently being managed and serviced by  $S$  are referred to as *active users*, which we denote by  $\mathcal{U}_A$ , which is a subset of all serviceable users  $\mathcal{U} = \{u_m\}$ .

Each user  $u_m$  stores the received encoding vectors up to timeslot  $t$  in a buffer (matrix) denoted by  $\mathbf{K}_m(t)$ . This is called the knowledge space of user  $u_m$  at timeslot  $t$ . The rank of knowledge space of user  $u_m$  at timeslot  $t$  is denoted by  $r_m(t) = \text{rank}(\mathbf{K}_m(t))$ .

A user is said to *receive a new degree of freedom (d.o.f.)* if the rank of its knowledge space increases by one after reception of a chunk from  $S$ , that is, if  $r_m(t+1) = r_m(t) + 1$ . A file chunk  $f_j$  is said to *decoded* by user  $u_m$  if the user can obtain the corresponding unit encoding vector  $\mathbf{e}_j$  (possibly after Gaussian elimination) from its knowledge space  $\mathbf{K}_m(t)$ .

An active user  $u_m$  at timeslot  $t$  is a user whose d.o.f. satisfies  $r_m(t) < T$ . User  $u_m$  is said to *depart* the queue at time  $t$  when the rank of its knowledge space becomes  $T$ .

<sup>2</sup> $N$  is an arbitrary, possibly time varying, quantity and hence this model does not limit our analysis.

Throughout the rest of the paper, a user always refers to an active user who has not yet departed from the server's buffer.

References [1], [2] have assumed perfect scheduling by the server, which is not assumed in our model. Somewhat related to this issue is the assumption that the coefficients of a coded chunk are cycled or refreshed to ensure innovative chunks for every drive read. Finally, to be able to apply queuing theoretical arguments in [1], requests for different file chunks of the same content arrive randomly and independently of other chunks at the server. In that paper, the notion of users is abstracted away, which we do not do here.

#### D. Performance Metrics

Let  $\mathcal{U}_n(t) \subset \mathcal{U}_A$  be the subset of *targeted* users who receive an innovative d.o.f. from the broadcast of chunk  $x(t)$  at timeslot  $t$ . We define three throughput metrics in order of strongest to weakest, which are equivalent to those used in cross-bar switch scheduling [16].

*Definition 1. (Throughput optimal)* A scheduling service is *throughput optimal* if every service can guarantee  $\mathcal{U}_n(t) = \mathcal{U}_A$ . That is,  $r_m(t+1) = r_m(t) + 1, \forall u_m \in \mathcal{U}_A, \forall t$ .

Since system constraints may mean that throughput optimality is not feasible, we consider maximum and maximal throughput, which are in general the best any scheduling scheme can do *up to* or *at* any timeslot based on constraints such as drive availability.

*Definition 2. (Maximum throughput)* A service scheme achieves *maximum throughput* if the total number of targeted users up to time  $t$ , denoted by  $\sum_{i=1}^t |\mathcal{U}_n(i)|$  is maximized, across all service schemes for a given data storage allocation.

*Definition 3. (Maximal throughput)* A scheduling service achieves *maximal throughput* if at each timeslot  $t$ , the number of targeted users  $|\mathcal{U}_n(t)|$  is maximized, across all service schemes for a given data storage allocation.

Note that any service scheme that achieves maximal throughput is necessarily a greedy algorithm. In a given timeslot, active users that are not targeted by a scheduling scheme are said to be *internally blocked*. These users are not externally blocked as they are already in the server's buffer, but are held up for service. The better the throughput of a scheduling scheme, the lower its internal blocking probability will be. Intuitively, a lower internal blocking probability should lead to lower external blocking probability as active users are flushed out of the system faster.

### III. DATA SCHEDULING SCHEMES

We introduce the concept of a service *leader* and considers two system types. First, to develop intuition for our problem and to verify expectations, we consider systems in which drives never block, i.e., drives with infinite I/O access bandwidth. Second, we consider systems with traffic-induced drive blocking, i.e., drives with finite I/O access bandwidth. In both systems, we propose service schemes for UCS and NCS. Schemes presented in this section can be formulated as integer linear programs over content demand graphs, similar to those for cross-bar switches [16], [17] and are omitted here.

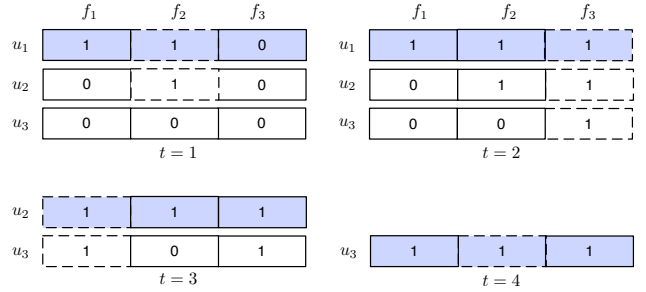


Fig. 2. Consider a system with  $T = 3$ , and the following example sequence of Algorithm 1 showing the evolution of users' decoded chunk vector. Users  $u_1$  and  $u_2$  are already in the system with different demands when user  $u_3$  arrives at  $t = 1$ . During each timeslot  $t$ ,  $S$  transmits the earliest undecoded chunk of the shaded leader.

#### A. Infinite I/O access bandwidth systems

To verify expectations, consider a system in which drives have infinite I/O access bandwidth, so  $P_b^D = 0$ .

1) *Uncoded Scheme:* Consider UCS and the scheme outlined in Algorithm 1. We introduce the following terminology for our leader-based scheme, which will also be used in the finite I/O access bandwidth case. Let  $\mathbf{a}_m(t)$  be a binary valued *decoded chunk vector* of length  $T$  for user  $u_m$  with elements  $a_{m,j}(t)$ . If  $a_{m,j}(t) = 0$  then user  $u_m$  has decoded file chunk  $f_j$  and if  $a_{m,j}(t) = 1$  then file chunk  $f_j$  is yet to be decoded. Upon arrival of user  $u_m$ 's file request,  $a_{m,j}(t) = 1$  for all  $1 \leq j \leq T$  and upon departure  $a_{m,j}(t') = 0$  for all  $1 \leq j \leq T$ .

- The *leader user*  $u_\ell$  at timeslot  $t$  is the user with maximum knowledge space rank. That is,

$$\ell = \operatorname{argmax}_{m:u_m \in \mathcal{U}_A} r_m(t). \quad (3)$$

- The *earliest undecoded chunk* or simply *min chunk* of user  $u_m$  is the chunk for which  $a_{m,j}(t) = 1$  and all  $a_{m,k}(t) = 0$  for  $k < j$ .
- The *earliest undecoded chunk of the leader* or simply *min-max chunk*  $f_{j^*}$  is the chunk for which  $a_{\ell,j^*}(t) = 1$  and  $a_{\ell,k}(t) = 0$  for  $k < j^*$  for the leader user  $u_\ell$ .

See Fig. 2 for an example of the leader-based scheme in Algorithm 1.

---

**Algorithm 1** Leader-based scheduling scheme for UCS with infinite I/O access bandwidth.

---

- 1: **for** timeslot  $t$  **do**
  - 2: Find the leading user  $u_\ell$  among all active users in  $\mathcal{U}_A$ , which has the highest knowledge space rank  $r_\ell(t)$ .
  - 3: Find the leader's earliest undecoded chunk denoted by  $f_{j^*}$ .
  - 4: Read  $f_{j^*}$  from a drive in  $\mathcal{M}_{j^*}^d$  and broadcast  $x(t) = f_{j^*}$  to all active users.
  - 5: All active users get to decode  $f_{j^*}$  and the server updates the decoded chunk list of all active users. That is,  $a_{m,j^*}(t) = 0$  for all active users  $u_m \in \mathcal{U}_A$ .
  - 6: **end for**
-

2) *Uncoded Scheme Analysis*: Intuitively a system with infinite I/O access bandwidth and perfect communications will allow for throughput optimal scheduling by  $S$ , since the scheme is without errors. We now formalize that the leader-based scheme in Algorithm 1 is throughput optimal according to Definition 1.

**Lemma 1.** *The scheduling scheme of Algorithm 1 is throughput optimal.*

*Proof:* We prove optimality by induction.

Base step: Consider an empty server queue. When the first user arrives, it immediately become the leader and the system services uncoded chunks sequentially starting from file chunk  $f_1$ . Therefore, in each timeslot this user will successfully receive a d.o.f. so the scheme is throughput optimal during this time.

Inductive step: Consider a throughput optimal scheme with a server queue comprising  $m$  users in  $\mathcal{U}_A$  where all users have received a d.o.f. in all previous timeslots. The  $(m+1)$ th user arrives. Since the knowledge space rank of the new  $(m+1)$ th user is zero, the leader remains unchanged. Choose uncoded chunk  $x(t) = f_{j^*}$  corresponding to the leader as per Algorithm 1. Then the new user will also receive a d.o.f. as it has received no chunks so far. So all users continue to receive a d.o.f. in every timeslot and the scheme remains throughput optimal. ■

**Lemma 2.** *If a scheduling scheme is throughput optimal, then it also minimizes the blocking probability across all feasible scheduling schemes.*

*Proof:* A throughput optimal scheme means that all users in  $\mathcal{U}_A$  receive an innovative d.o.f. in each timeslot. This means all users are serviced in  $T$  timeslots after their arrival, which is the minimum number possible because only one chunk can be broadcast per timeslot. Hence, the service rate  $\mu = |\mathcal{U}|/T$  is at the maximum for a throughput optimal scheme. Given a fixed arrival rate  $\lambda$  and a fixed buffer size  $N$ , the Erlang B blocking formula monotonically decreases with increasing service rate  $\mu$ . Hence the maximum service rate results in minimum blocking probability. ■

Applying Lemma 2 to Algorithm 1 shows that it is also blocking probability optimal.

**Lemma 3.** *The blocking probability of a throughput optimal scheme is given by*

$$P_b^s = \frac{(\lambda T)^N / N!}{\sum_{i=0}^N (\lambda T)^i / i!}. \quad (4)$$

*Proof:* The arrival process for  $S$  is a Poisson process. Under a throughput optimal scheme, all users immediately begin being serviced upon arrival until a total of  $N$  users are in the server buffer. We can view each active user as being serviced by an individual service unit with deterministic service time  $T$  timeslots. Hence, the average service rate is  $1/T$  for each server and  $S$  is equivalent to an  $M/D/M/M$  queue, where  $D$  is a deterministic service time. The blocking probability for  $S$  is then given by (4). ■

In a system with infinite I/O access bandwidth all drives are always available for read. Then there is no need to store more than one copy of each file chunk. That is,  $|\mathcal{M}_i^d| = 1$  for all  $f_i \in \mathcal{F}$  suffices for throughput optimality.

*Remark 1.* Serving the earliest undecoded chunk of the leader is not essential for the optimality of the algorithm. Selecting any undecoded chunk by the leader will suffice. However, by serving undecoded chunks of the leader in a contiguous way, we promote better in-order delivery to the application.

This verifies the intuitive result that NCS does not provide benefit over UCS in an infinite I/O access bandwidth system. Note that Algorithm 1 can be adjusted to operate with coded storage via simple modifications.

### B. Finite I/O access bandwidth systems

In this subsection we consider systems with drives that can become busy owing to serving other requests, i.e., drives with finite I/O access bandwidth for which  $P_b^D > 0$ . We still assume ideal chunk transport medium with no erasures and broadcast capabilities to all active users, such as TCP for multicast variants, Ethernet, or emulated broadcasting systems.

1) *Uncoded Scheme*: In the finite I/O case, the concept of leaders needs modification depending on what chunks are available for access. We then distinguish between a *true* leader and a *temporary* leader in our modified scheduling algorithm. This is to handle temporary unavailability of drives that store undecoded chunks demanded by the true leader. We modify Algorithm 1 to find undecoded chunks of the true leader that are available for read. If no such undecoded chunk for the true leader is available, then we will limit our search to the *next* leading user and the undecoded chunks of that user, which by the approach of the service scheme must have been all decoded by the excluded leader. We continue until we can find one user who is leading among the remaining users and for whom one of its undecoded chunks is available for read. The modified scheme operates as per Algorithm 2.

2) *Coded Scheme*: The proposed scheme for NCS finite I/O bandwidth systems is similar to Algorithm 2 in terms of finding temporary leaders depending on drive availability. The main difference with Algorithm 2 is the choice of the chunk for service: The scheduler needs to keep track of coded chunks so far received by the users.

For each timeslot  $t$ , we define a binary *coded chunk reception* vector of size  $H$  for user  $u_m$ , denoted by  $\mathbf{q}_m(t)$ , as follows:  $q_{m,r}(t) = 0$  if coded chunk  $c_r$  has been so far received by user  $u_m$  and  $q_{m,r}(t) = 1$  otherwise. Algorithm 3 describes the scheme.

3) *Schemes Analysis and Comparison*: In a finite I/O storage system neither UCS nor NCS can guarantee throughput optimality, since we can always find a drive unavailability pattern with non-zero probability of occurring that would block at least one user (for instance, consider the simple case when all drives are blocked in the same timeslot). We now show simple proofs showing that while both uncoded schedule of Algorithm 2 and coded schedule of Algorithm

**Algorithm 2** Leader-based scheduling scheme for UCS with finite I/O access bandwidth.

- 1: **for** timeslot  $t$  **do**
- 2: Obtain the drive availability vector  $\mathbf{b}(t)$ .
- 3: Create a temporary list of active users, denoted by  $\mathcal{U}_t$ , and initialize it to all active users in the system:  $\mathcal{U}_t = \mathcal{U}_A$ .
- 4: Find the leader from the temporary list of active users from (3).
- 5: Find the leader's set of all undecoded chunks denoted by  $\mathcal{F}_\ell^u \subset \mathcal{F}$ . That is,  $f_j \in \mathcal{F}_\ell^u \leftrightarrow a_{\ell,j}(t) = 1$ .
- 6: If there exists at least one available drive for at least one chunk in  $\mathcal{F}_\ell^u$ , then select one such chunk, denoted by  $f_{j^*}$ , and go to step 7. Otherwise, remove the leader from temporary active users ( $\mathcal{U}_t \leftarrow \mathcal{U}_t \setminus u_\ell$ ) and go to step 4.
- 7: Read the chunk  $f_{j^*}$  from one of the available drives in  $\mathcal{M}_{j^*}^d$  and broadcast  $x(t) = f_{j^*}$  to all active users.
- 8: All users decode the chunk  $f_{j^*}$  and the server updates their decoded chunk list. That is,  $a_{m,j^*}(t) = 0$  for all active users in  $\mathcal{U}_t$ . (Note that the excluded leading users have already decoded  $f_{j^*}$  and hence at the end of this step  $a_{m,j^*}(t) = 0$  for all users  $u_m \in \mathcal{U}_A$ ).
- 9: **end for**

**Algorithm 3** Leader-based scheme for NCS with finite I/O access bandwidth.

- 1: **for** timeslot  $t$  **do**
- 2: Obtain the drive availability vector  $\mathbf{b}(t)$ .
- 3: Create a temporary list of active users,  $\mathcal{U}_t = \mathcal{U}_A$ .
- 4: Find the leader from the temporary list of active users.
- 5: Find the leader's set of all unreceived coded chunks denoted by  $\mathcal{C}_\ell^u$ . That is,  $c_r \in \mathcal{C}_\ell^u \leftrightarrow q_{\ell,r}(t) = 1$ .
- 6: If there exists at least one available coded chunk in  $\mathcal{C}_\ell^u$  for read, then select one such chunk, denoted by  $c_{r^*}$ , and go to step 7. Otherwise, remove the leader from temporary active users ( $\mathcal{U}_t \leftarrow \mathcal{U}_t \setminus u_\ell$ ) and go to step 4.
- 7: Read the chunk  $c_{r^*}$  from its corresponding drive and broadcast  $x(t) = c_{r^*}$  to all active users.
- 8: Update  $q_{m,r^*}(t) = 0$  for all active users  $u_m \in \mathcal{U}_A$ .
- 9: **end for**

3 achieve maximal throughput across their respective data storage formats, that the number of targeted users using NCS with maximal throughput scheduling is at least as high as that in the UCS system.

**Lemma 4.** *Algorithms 2 and 3 achieve maximal throughput across their respective data storage formats.*

*Proof:* The scheme of Algorithm 2 (Algorithm 3) identifies a leading user with maximum rank with available file chunk(s) for read. All other users with smaller or equal ranks will also receive a d.o.f. since no chunks exist that non-

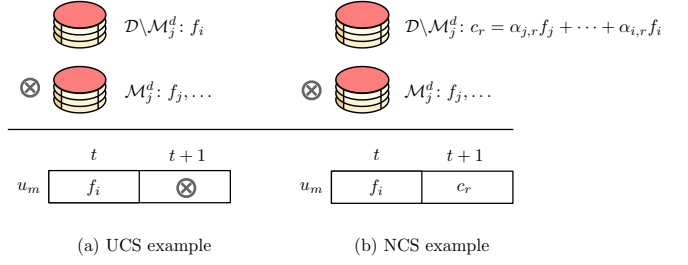


Fig. 3. Given an instance of UCS, an example construction NCS can improve throughput compared to a UCS.

leader users have decoded (received) but the leader has not. Consequently, at any given time, the number of serviced users with a d.o.f. is maximized subject to instantaneous drive availability given the storage format. Therefore, both algorithms achieve maximal throughput. ■

**Lemma 5.** *The number of targeted users  $|\mathcal{U}_n(t)|$  in a finite I/O storage system using NCS with maximal throughput scheduling can be at least as high as that in a UCS system with maximal throughput scheduling.*

*Proof:* Given any instance of UCS, we need to show (1) that no drive blocking patterns exist where the number of targeted users  $|\mathcal{U}_n(t)|$  is higher than that in all instances of NCS, and (2) that there exist drive blocking patterns for which  $|\mathcal{U}_n(t)|$  in NCS is higher than that in the UCS instance.

For (1), consider an instance of NCS which is constructed as follows. Each coded chunk  $c_r$  stored on  $D_i$  is a linear combination of the uncoded chunks stored on  $D_i$  in the UCS instance. Under this scenario, given linear independence from earlier chunks, if any read from  $D_i$  in the UCS instance can target  $|\mathcal{U}_n(t)|$  users, it is clear that a read from  $D_i$  in the NCS counterpart can also provide a new d.o.f. to at least the same number of users. For (2), we proceed by counterexample. We can always consider a single active user  $u_m$  with  $r_m(t) = T - 1$  under the UCS instance. See Fig. 3 for a toy-example, when  $r_1(2) = 1$  and the only missing chunk of user  $u_1$  is  $f_j$ . Assume that all drives in  $\mathcal{M}_j^d$  are blocked during timeslot  $t+1$ . For the UCS system,  $u_m$  cannot be targeted so  $|\mathcal{U}_n(t+1)| = 0$ . However, in the NCS instance of the system, although  $\mathcal{M}_j^d$  is blocked, any unseen coded chunks with  $\alpha_{j,r} \neq 0$  stored on drives in  $\mathcal{D} \setminus \mathcal{M}_j^d$  can still provide a new d.o.f. to user  $u_m$ , so  $|\mathcal{U}_n(t+1)| = 1$ . ■

To further illustrate NCS improved blocking performance, we now focus on the internal true leader's blocking probability. First, consider the following restricted UCS file layout with replication and striping. There are a total of  $R = WT$  drives in the system. The uncoded system stores a single file chunk per drive, where drive  $D_{(w-1)T+i}$  stores the  $w$ th copy of the  $i$ th file chunk for  $w = 1, \dots, W$  and  $i = 1, \dots, T$ . The coded system stores  $H = R = WT$  coded file chunks such as  $c_r$ , one on each drive  $D_r$ . Assuming that each drive becomes unavailable with probability  $P_b^D$  independently of other drives and previous timeslots, the following lemma gives the leader internal blocking probability in each system.



**Lemma 6.** *At timeslot  $t$ , the internal blocking probability of the true leader who has a knowledge space rank of  $r_\ell(t)$  is given by*

$$P_b^c = (P_b^D)^{WT - r_\ell(t)} \quad (5)$$

*in an NCS system and by*

$$P_b^u = (P_b^D)^{WT - Wr_\ell(t)} \quad (6)$$

*in a UCS system.*

*Proof:* If the leader has received  $r_\ell(t)$  coded chunks up to time  $t$ , there remain only  $WT - r_\ell(t)$  useful drives for service and (5) follows. In the uncoded system, if the leader has decoded  $r_\ell(t)$  file chunks up to time  $t$ , there remain only  $WT - Wr_\ell(t)$  useful drives for service and (6) follows. ■

For large  $r_\ell(t)$  or  $W$ , the improvement in leader blocking probability enabled by coded storage can become significant. Next we consider regular  $s$ -striped storage systems [1] with a total of  $R = Ws$  drives and  $T/s$  file chunks in each stripe set which is an integer. The following lemma gives the leader blocking probability in uncoded and coded systems.

**Lemma 7.** *Assume that at timeslot  $t$ , the leader in the uncoded system has completely decoded  $r$  out of  $s$  stripe sets, where  $r = 0, \dots, s-1$ , such that its knowledge space rank satisfies  $rT/s \leq r_\ell(t) < (r+1)T/s$ . Then, its internal blocking probability is given by*

$$P_b^u = (P_b^D)^{Ws - Wr} \quad (7)$$

*Now assume that in the coded system, the leader's knowledge space rank is also  $r_\ell(t)$ . A simple upper bound for the internal blocking probability  $P_{b,ub}^c$  is given by*

$$(P_b^D)^{Ws-r} \leq P_{b,ub}^c = (P_b^D)^{Ws - \lfloor \frac{r_\ell(t)}{T/s} \rfloor} < (P_b^D)^{Ws - (r+1)} \quad (8)$$

*And a simple lower bound  $P_{b,lb}^c$  is given by*

$$P_{b,lb}^c = (P_b^D)^{Ws - \max(0, r_\ell(t) - WT - Ws)} \quad (9)$$

*which will deviate from the best possible blocking probability of  $(P_b^D)^{Ws}$  only when  $W = 1$  and  $T - s < r_\ell(t) < T$ .*

*Proof:* In the uncoded system, if the leader has completely decoded  $r$  stripe sets up to time  $t$ , there only remains  $Ws - Wr$  useful drives for service and (7) follows.

The worst case for the coded system occurs when during  $r_\ell(t)$  previous services of the leader,  $\lfloor \frac{r_\ell(t)}{T/s} \rfloor$  out of  $Ws$  available drives were completely read off and hence are unavailable for further service, in which case (8) follows. The bounds are derived by using the inequalities  $rT/s \leq r_\ell(t) < (r+1)T/s$ .

The best case for the coded system occurs where all previous  $r_\ell(t)$  services of the leader were uniformly read across  $Ws$  available drives. Therefore, one can verify that until the leader's rank reaches  $r_\ell(t) = Ws(T/s - 1) + 1 = WT - Ws + 1$ , none of the drives are completely read off and are all available for service. Hence, we get,  $P_{b,ub}^c = (P_b^D)^{Ws}$  for  $r_\ell(t) < WT - Ws + 1$ . After this point, the drives become sequentially unavailable and  $(P_b^D)^{Ws - \max(0, r_\ell(t) - WT - Ws)}$

follows. One can easily verify the last statement of the lemma using  $r_\ell(t) < T$ , the assumption that  $T/s$  is an integer and  $s \leq T/2$ . ■

*Remark 2.* Lemma 7 demonstrates the importance of drive selection in Algorithms 2 and 3, when more than one drive containing undecoded file chunks of the leader is available for read. One can think about this as memory in the system: Drive service units cease being helpful if all their content has been read. When comparing different variations of Algorithms 2 and 3, we expect that those which temporally spread reads across drives to have better average throughput.

## IV. NUMERICAL RESULTS

Using typical values found in various modern systems, we present Monte Carlo simulation results comparing the performance of the proposed leader-based scheduling scheme for UCS and NCS systems. By using (4), the analytical results for the blocking probability of the proposed leader-based scheduling scheme for uncoded/coded storage with infinite I/O access bandwidth are presented. For all simulations, we use a regular striped mapping of file chunks onto drives.

Fig. 4 illustrates the external blocking probability of the proposed scheduling scheme for both uncoded and coded storage under drives' infinite and finite I/O access bandwidth conditions versus server buffer size,  $N$ . We see that when drives have finite I/O access bandwidth, NCS reduces system blocking probability over UCS and that the gap tends to grow with increasing buffer size.

Fig. 5 shows the average throughput and external blocking probability of the proposed schemes for various drive internal blocking probabilities,  $P_b^D$ . As shown, throughput and external blocking probability are improved in NCS compared with UCS as drives become more overwhelmed. In addition, we see that a small 3% improvement in throughput renders a comparatively large improvement of 150% in external blocking probability.

The internal blocking probability of the true leader versus its knowledge space rank in the uncoded and coded storage for various drive internal blocking probabilities,  $P_b^D$ , is presented in Fig. 6. Here, the internal blocking probability of the true leader is much lower in the coded system compared with the uncoded system as the leader's knowledge space rank increases. This lower blocking probability is one factor explaining lower external blocking probability of the coded system compared to the uncoded system.

## V. CONCLUSIONS

In this paper, we introduced a novel and simple storage model for point-to-multipoint SANs and investigated the impact of scheduling and content storage format on system blocking probability and throughput in PMP networks. We proposed two intuitive drive access scheduling techniques for both UCS and NCS systems, under infinite and finite I/O access bandwidth conditions. In finite I/O access networks, we showed that NCS scheduling flexibility improves blocking

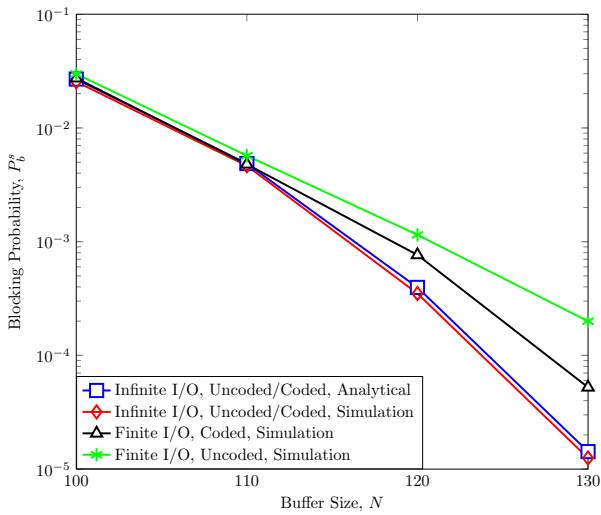


Fig. 4. Blocking probability versus server buffer size  $N$  for  $\lambda = 0.9, T = 100, W = 2, R = 8, s = 4, P_b^D = 0.5$ .

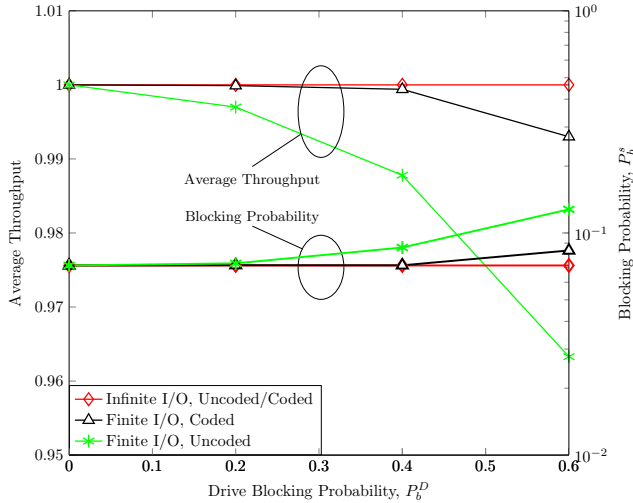


Fig. 5. Average throughput and blocking probability versus  $P_b^D$  for  $\lambda = 0.9, T = 8, W = 2, R = 8, s = 4$ .

probability and throughput over UCS. Our numerical evaluations and simulation results verify these advantages and can be used to guide future storage system design.

## REFERENCES

- [1] U. J. Ferner, M. Médard, and E. Soljanin, "Toward sustainable networking: Storage area networks with network coding," in *Proc. Allerton Conf. on Commun., Control and Computing*, Champaign, IL, Oct. 2012.
- [2] A. G. Dimakis, K. Ramchandran, Y. Wu, and C. Suh, "A survey on network codes for distributed storage," *Proc. IEEE*, vol. 99, no. 3, pp. 476–489, Mar. 2011.
- [3] D. Ford, F. Labelle, F. I. Popovici, M. Stokely, V.-A. Truong, L. Barroso, C. Grimes, and S. Quinlan, "Availability in globally distributed storage systems," in *Proceedings of the 9th USENIX conference on Operating systems design and implementation*, ser. OSDI'10. Berkeley, CA: USENIX Association, 2010, pp. 1–7. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1924943.1924948>
- [4] N. B. Shah, K. Lee, and K. Ramchandran, "The MDS Queue: Analysing latency performance of codes and redundant requests," *CoRR*, <http://arxiv.org/abs/1211.5405>, 2012.
- [5] L. Huang, S. Pawar, Z. Hao, and K. Ramchandran, "Codes can reduce queueing delay in data centers," in *Proc. IEEE Int. Symp. on Inf. Theory*, Jul. 2012, pp. 2766–2770.

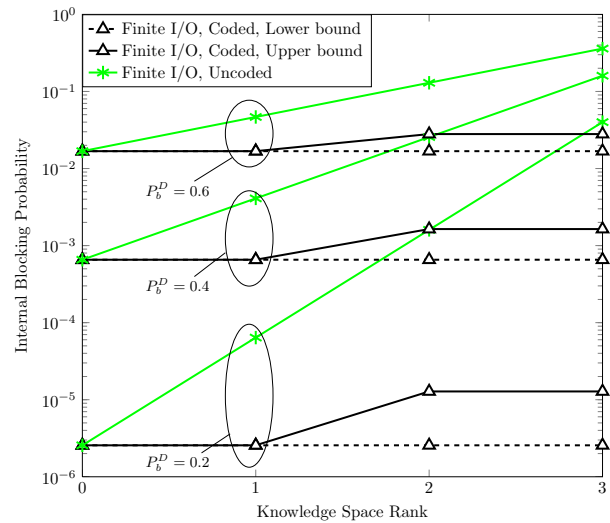


Fig. 6. Blocking probability of the true leader versus drive blocking probability  $P_b^D$  for  $T = 8, W = 2, s = 4$ .

- [6] T. E. Anderson, S. S. Owicki, J. B. Saxe, and C. P. Thacker, "High-speed switch scheduling for local-area networks," *ACM Trans. Comput. Syst.*, vol. 11, no. 4, pp. 319–352, Nov. 1993. [Online]. Available: <http://doi.acm.org/10.1145/161541.161736>
- [7] M. A. Marsan, A. Bianco, P. Giaccone, E. Leonardi, and F. Neri, "Multicast traffic in input-queued switches: optimal scheduling and maximum throughput," *IEEE/ACM Trans. Netw.*, vol. 11, no. 3, pp. 465–477, Jun. 2003. [Online]. Available: <http://dx.doi.org/10.1109/TNET.2003.813048>
- [8] B. Prabhakar, N. McKeown, and R. Ahuja, "Multicast scheduling for input-queued switches," vol. 15, no. 5, pp. 855–866, Jun. 1997.
- [9] H. Yu, S. Ruepp, and M. S. Berger, "Multi-level round-robin multicast scheduling with look-ahead mechanism," in *Proc. IEEE Int. Conf. on Commun.*, Kyoto, Japan, Jun. 2011, pp. 1–5.
- [10] C. Feng and B. Li, *Network coding: Fundamentals and applications*, 1st ed. Academic Press, 2012, ch. Network coding for content distribution and multimedia streaming in peer-to-peer networks.
- [11] G. N. Rouskas and V. Sivaraman, "Packet scheduling in broadcast WDM networks with arbitrary transceiver tuning latencies," *IEEE/ACM Trans. Netw.*, vol. 5, no. 3, pp. 359–370, Jun. 1997.
- [12] D. Aksoy, M. J. Franklin, and S. Zdonik, "Data staging for on-demand broadcast," in *Proc. 27th VLDB Conf.*, Roma, Italy 2001.
- [13] A. Hu, "Video-on-demand broadcasting protocols: A comprehensive study," in *Proc. IEEE Conf. on Computer Commun.*, Anchorage, AK, Apr. 2001, pp. 508–517.
- [14] T. Ho, R. Koetter, M. Médard, M. Effros, J. Shi, and D. Karger, "A random linear network coding approach to multicast," *IEEE Trans. Inf. Theory*, vol. 52, no. 10, pp. 4413–4430, Oct. 2006.
- [15] U. J. Ferner, Q. Long, M. Pedroso, L. Voloch, and M. Médard, "Building a network coded storage testbed for data center energy reduction," in *Proc. IEEE SustainIT*, Palermo, Italy, Oct. 2013.
- [16] C.-S. Chang, W.-J. Chen, and H.-Y. Huang, "On service guarantees for input-buffered crossbar switches: a capacity decomposition approach by Birkhoff and von Neumann," in *Proc. IWQoS*, 1999, pp. 79–86.
- [17] J. Sundararajan, S. Deb, and M. Médard, "Extending the Birkhoff-von Neumann switching strategy for multicast - on the use of optical splitting in switches," *IEEE J. Sel. Areas Commun.*, vol. 25, pp. 36–50, 2007.
- [18] M. Kim, J. K. Sundararajan, M. Médard, A. Eryilmaz, and R. Kotter, "Network coding in a multicast switch," *IEEE Trans. Inf. Theory*, vol. 57, no. 1, pp. 436–460, 2011.