

Rateless Codes for Near-Perfect Load Balancing in Distributed Matrix-Vector Multiplication

ANKUR MALLICK*, Carnegie Mellon University
 MALHAR CHAUDHARI†, Oracle Corporation
 UTSAV SHETH†, Automation Anywhere
 GANESH PALANIKUMAR†, Apple Inc.
 GAURI JOSHI, Carnegie Mellon University

Large-scale machine learning and data mining applications require computer systems to perform massive matrix-vector and matrix-matrix multiplication operations that need to be parallelized across multiple nodes. The presence of straggling nodes – computing nodes that unpredictably slowdown or fail – is a major bottleneck in such distributed computations. Ideal load balancing strategies that dynamically allocate more tasks to faster nodes require knowledge or monitoring of node speeds as well as the ability to quickly move data. Recently proposed fixed-rate erasure coding strategies can handle unpredictable node slowdown, but they ignore partial work done by straggling nodes thus resulting in a lot of redundant computation. We propose a *rateless fountain coding* strategy that achieves the best of both worlds – we prove that its latency is asymptotically equal to ideal load balancing, and it performs asymptotically zero redundant computations. Our idea is to create linear combinations of the m rows of the matrix and assign these encoded rows to different worker nodes. The original matrix-vector product can be decoded as soon as slightly more than m row-vector products are collectively finished by the nodes. We conduct experiments in three computing environments: local parallel computing, Amazon EC2, and Amazon Lambda, which show that rateless coding gives as much as $3\times$ speed-up over uncoded schemes.

CCS Concepts: • **Mathematics of computing** → **Queueing theory**; *Coding theory*; • **Computer systems organization** → **Redundancy**; *Reliability*.

Additional Key Words and Phrases: Erasure coded Computing, Rateless Fountain Codes, Large-scale Parallel Computing

ACM Reference Format:

Ankur Mallick, Malhar Chaudhari, Utsav Sheth, Ganesh Palanikumar, and Gauri Joshi. 2019. Rateless Codes for Near-Perfect Load Balancing in Distributed Matrix-Vector Multiplication. *Proc. ACM Meas. Anal. Comput. Syst.* 3, 3, Article 58 (December 2019), 40 pages. <https://doi.org/10.1145/3366706>

1 INTRODUCTION

Matrix-vector multiplications form the core of a plethora of scientific computing and machine learning applications that include solving partial differential equations [3], forward and back propagation in neural networks [7], computing the PageRank of graphs [48] etc. In the age of Big Data, most of these applications involve multiplying extremely large matrices and vectors and

*Correspondence Author

†Work done while at CMU

Authors' addresses: Ankur Mallick, Carnegie Mellon University, Pittsburgh, PA, amallic1@andrew.cmu.edu; Malhar Chaudhari, Oracle Corporation, Redwood City, CA, malharchaudhari@gmail.com; Utsav Sheth, Automation Anywhere, San Jose, CA, utsavsheth1994@gmail.com; Ganesh Palanikumar, Apple Inc., Cupertino, CA, ganeshpkumar93@gmail.com; Gauri Joshi, Carnegie Mellon University, Pittsburgh, PA, gaurij@andrew.cmu.edu.

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM.

This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, <https://doi.org/10.1145/3366706>.

the computations cannot be performed efficiently on a single machine. This has motivated the development of several algorithms [40], [18] that seek to speed-up matrix-vector multiplication by distributing the computation across multiple computing nodes. The individual nodes (the *workers*) perform their respective tasks in parallel while a central node (the *master*) aggregates the output of all these workers to complete the computation.

The Problem of Stragglers. Unfortunately, large-scale distributed computation jobs are often bottlenecked by tasks that are run on unpredictably slow or unresponsive workers called *stragglers* [10]. Since the job is complete only when all its parallel tasks are executed, this problem is aggravated for jobs with a large number of parallel tasks. Even a tiny probability of a node slowing down and becoming a straggler can cause a big increase in the expected latency of the job. As pointed out in [10, Table 1], the latency of executing many parallel tasks could be significantly larger (140 ms) than the median latency of a single task (1 ms). Straggling of nodes is widely observed in cloud infrastructure [10] and it is the norm rather than an exception.

1.1 Previous Solution Approaches

Load Balancing Strategies. An obvious solution to overcome the bottleneck of waiting for slow nodes is to move tasks from busy or slow nodes to idle or fast nodes. Such work stealing or dynamic load balancing strategies are often implemented in shared and distributed memory settings [12, 13, 26]. This approach involves establishing a protocol for continually monitoring workers and moving tasks from slow to fast workers. It entails considerable centralized control over the computing environment, which may not be feasible in cloud systems where the nodes can unpredictably slow down due to background processes, network outages etc. There may also be concerns regarding data privacy, and the communication cost of moving data between nodes in a distributed system spread over a large geographical area. Thus it is desirable to develop a principled and easy-to-implement straggler-mitigation approach that does not involve moving data between workers.

Task Replication. Existing systems like MapReduce [11] and Spark [65] generally deal with the problem of stragglers by launching replicas of straggling tasks, which are referred to as *back-up* tasks. This strategy of task replication has many variants such as [4, 5], and has been theoretically analyzed in [58–60] where schemes for adding redundant copies based on the tail of the runtime distribution at the workers are proposed. In the area of queueing theory there is a line of interesting recent works analyzing the effect of task replication on queueing delays in multi-server systems [19–21, 29, 30, 54, 55]. For distributed matrix-vector multiplication, which is the focus of this work, a simple replication strategy is to divide matrix A into p/r (where r divides the number of workers p) sub-matrices and replicate each sub-matrix at r workers. Then the master waits for the fastest worker from each set of r to finish multiplying its sub-matrix with the vector x in order to recover the overall result $b = Ax$.

Erasured Coded Matrix-vector Multiplication. From a coding-theoretic perspective, task replication is a special case of more general *erasure codes* that overcome loss or erasure of data and recover the message from a subset of the transmitted bits. Erasure codes were first employed to overcome stragglers in the context of fast content download from distributed storage [28, 31, 32]. A file that is divided into k chunks and encoded using a (p, k) maximum-distance-separable (MDS) code (for example a Reed-Solomon code), can be recovered by downloading any k out of p encoded chunks. Queueing models to analyze the latency of coded content download jobs were proposed and analyzed in [32, 34, 35, 42, 50].

Unlike distributed storage, erasure coding of computing jobs is not straightforward. A job with n parallel tasks needs to be designed such that the execution of any k out of n tasks is sufficient to complete the job. However, this is possible for linear computations such as matrix-vector multiplication. The usage of codes to provide error-resilience in computation has its origins

in works on algorithmic fault tolerance [27]. Recent works such as [14, 41, 44] have employed Maximum Distance Separable (MDS) codes to speed up the computation of matrix vector products in a distributed setting. For example, suppose that we want to multiply a matrix \mathbf{A} with vector \mathbf{x} using 3 worker nodes and a (3, 2) MDS code. Then we split \mathbf{A} along rows into two matrices \mathbf{A}_1 and \mathbf{A}_2 such that $\mathbf{A} = [\mathbf{A}_1^T \ \mathbf{A}_2^T]^T$. The worker nodes store matrices \mathbf{A}_1 , \mathbf{A}_2 and $\mathbf{A}_1 + \mathbf{A}_2$ respectively, and each node multiplies its matrix with \mathbf{x} . Results from any two worker nodes are sufficient to obtain \mathbf{Ax} , and thus the system is tolerant to 1 straggling node.

Erasure Coding for Other Linear Computations. A natural generalization of matrix-vector multiplication is matrix-matrix multiplication, considered in [16, 61, 63]. There are also works dealing with coded gradient descent [24, 43, 56], coded convolution [15], coded Fourier Transform [64], Page Rank [62], and coded distributed optimization [36] – in general any distributed *linear* computation can essentially be expressed as a matrix multiplication/addition operation and can be made straggler-proof using erasure codes. Most of these works use MDS codes as the core idea and modify it for the specific computation in question. In this work, we focus on the original problem of coded distributed matrix-vector multiplication, but we consider a rateless-coded approach in contrast to MDS codes. We expect that the underlying principles of our work can be extended to matrix-matrix multiplication and other linear computations like gradient descent in the future.

1.2 Rateless Coding Approach and its Benefits

The replication or MDS coding strategies used for matrix-vector multiplication are fixed-rate strategies, that is, they fix a redundancy rate k/p when encoding matrix \mathbf{A} , use the results of the fastest k out of p worker nodes. The key drawbacks of this approach are that: 1) it cannot perform load balancing within the fastest k nodes and account for variabilities in their speeds, and 2) it discards all partial work done by the $p - k$ straggling workers. We address both these issues by proposing the use of *rateless fountain codes*, specifically Luby Transform (LT) codes [33, 45, 52] which are known to be scalable and form the basis for practical erasure coding schemes in wireless communication standards [53].

The rateless coded matrix-vector multiplication algorithm generates $m_e = \alpha m$ ($\alpha > 1$) coded linear combinations of the m rows of matrix \mathbf{A} and distributes them equally across p worker nodes. Each of these linear combinations is generated by choosing d of the m rows uniformly at random and adding them. For example, if $d = 2$, and we choose rows \mathbf{a}_1 and \mathbf{a}_3 of \mathbf{A} , then the encoded row is $\mathbf{a}_1 + \mathbf{a}_3$, as shown in Fig. 5a. The value d , referred to as the degree of the linear combination is an i.i.d. realization of a carefully chosen degree distribution $\rho(d)$. For LT codes, $\rho(d)$ is the Robust Soliton distribution (given in (4) below). Each worker receives m_e/p encoded rows of matrix \mathbf{A} and a copy of the vector \mathbf{x} . It computes row-vector products, for example $(\mathbf{a}_1 + \mathbf{a}_3)^T \mathbf{x} = b_1 + b_3$ for the encoded row $(\mathbf{a}_1 + \mathbf{a}_3)$, and sends them back to the master node. Due to the carefully chosen degree distribution $\rho(d)$, the master node can use an iterative peeling decoder [45] (illustrated in Fig. 5b) to recover each element of the product vector $\mathbf{b} = \mathbf{Ax}$ with a low decoding complexity of $O(\log m)$. Overall, it needs to wait for any $M' = m(1 + \epsilon)$ row-vector products to be completed across *all* the nodes, where ϵ is a small overhead; $\epsilon \rightarrow 0$ as $m \rightarrow \infty$).

Rateless codes offer the following key benefits over previously proposed coding techniques based on MDS codes.

(1) **Near-Ideal Load Balancing.** In order to adjust to varying speeds of worker nodes and minimize the overall time to complete the multiplication \mathbf{Ax} , one can use an *ideal load-balancing scheme that dynamically assigns one row-vector product computation task to each worker node as soon as the node finishes its current task*. Thus, faster nodes complete more tasks than slower nodes, and the final product $\mathbf{b} = \mathbf{Ax}$ is obtained when the p nodes collectively finish m row-vector products. Our

rateless coding strategy achieves nearly the same load balancing benefit without the communication overhead of dynamically allocating the tasks one row-vector product at a time. In our strategy, the nodes need to collectively finish $M' = m(1 + \epsilon)$ row-vector products, for small ϵ that goes to zero as $m \rightarrow \infty$. In contrast, MDS coding strategies do not adjust to different degrees of node slowdown; they use the results from k nodes, and ignore the remaining $p - k$ nodes. As a result rateless codes achieve a much lower delay than MDS coding strategies.

(2) **Negligible Redundant Computation.** A major drawback of MDS coding is that if there is no straggling, the worker nodes collectively perform mp/k row-vector products, instead of m . With the rateless coding strategy, the nodes collectively perform a maximum of $M' = m(1 + \epsilon)$ row-vector products where, the overhead ϵ goes to zero as m , the number of rows in the matrix \mathbf{A} increases.

(3) **Maximum straggler tolerance.** A (p, k) MDS coded distributed computation is robust to $p - k$ straggling nodes, for $k \in [1, 2, \dots, p]$. Reducing k increases straggler tolerance but also adds more redundant computation. The rateless coding scheme can tolerate up to $p - 1$ stragglers, with negligible redundant computation overhead.

(4) **Low Decoding Complexity.** One may argue that MDS coding approaches can also use partial computations and achieve near-perfect load balancing if we construct an (m_e, m) MDS code (for a given amount of redundancy m_e/m) to encode a $m \times n$ matrix. The decoding complexity of such a code is $O(m^3)$ which is unacceptable for large m in practical implementations. Rateless codes offer a low decoding complexity: $O(m \log m)$ for LT codes [45], and $O(m)$ for Raptor codes [52].

Difference from [49, 61] on LT-coded Matrix-vector Multiplication. The use of Luby Transform (LT) codes for matrix-vector multiplication has been recently proposed in [49, 61]. *However, these works do not utilize the ‘rateless’ property of LT codes and instead use them in a fixed-rate setting.* For example, the algorithm in [49] generates m_e LT-coded rows from an m -row matrix using LT codes, and it allocates each row to ηq workers for some $\frac{1}{p} \leq \eta \leq 1$. Each worker completes the entire set of row-vector product tasks assigned to them, and the master waits for the fastest q workers to finish. Partial computations performed by slow workers are discarded. The scheme proposed in [61] also uses LT codes in this fixed-rate setting and focuses on using the sparsity of LT codes to reduce the decoding complexity of coded matrix multiplication. Thus, although these works use LT codes, they are similar in spirit to fixed-rate MDS-coding approaches.

To the best of our knowledge, our work is the first to exploit the *rateless* nature of LT codes to perform load-balancing in distributed matrix computations and utilize all the partial work done by slow workers. We also provide the first theoretical analysis of the latency achieved by this strategy with ideal load balancing and show that it asymptotically achieves near-perfect latency and computation cost. Previous works [49, 61] do not present such analyses. Moreover, we present extensive experimental results on 3 different computing environments: local parallel computing, distributed computing on Amazon EC2 and serverless computing on Amazon Lambda.

1.3 Main Theoretical and Experimental Results

Besides proposing the rateless coding strategy, one of the main contributions of our work is to theoretically analyze and compare it with ideal load balancing. In particular, we consider two performance metrics: 1) latency T , which is the time until $\mathbf{b} = \mathbf{Ax}$ can be recovered by the master, and 2) number of computations C , which is the number of row-vector product tasks completed by the p workers until \mathbf{b} can be recovered. We consider a simple delay model where worker i has an initial delay of X_i after which it spends a constant time τ per row-vector product task.

Comparison with Ideal Load Balancing. In the ideal load balancing strategy, the m row-vector product tasks (which comprise the job of multiplying the $m \times n$ size \mathbf{A} with vector \mathbf{x}) are kept in a central queue at the master and dynamically allocated to idle workers one task at a time. The job

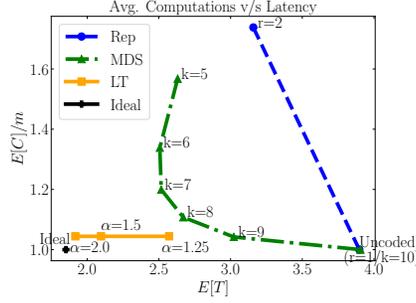


Fig. 1. The expected latency ($\mathbb{E}[T]$) of the LT-Coded approach smoothly decays on adding redundancy (increasing α) and approaches the Ideal approach without any increase in computational overhead ($\mathbb{E}[C]/m$). Previous approaches - Replication (Rep) and MDS coding not only have a higher latency for the same task but also perform far more redundant computations. The simulation parameters are $m = 10000$, number of worker nodes $p = 10$ and delay model parameters $\mu = 1.0$, $\tau = 0.001$.

is complete when m tasks are collectively finished by the workers. The rateless coding strategy differs from this ideal policy in two ways due to which its latency is larger: 1) each worker gets $m_e/p = \alpha m/p$ encoded rows and thus a fast worker may run out of rows before the master is able to recover $\mathbf{b} = \mathbf{A}\mathbf{x}$ and 2) the workers collectively need to finish $m(1 + \epsilon)$ tasks where ϵ is a small overhead that diminishes as $m \rightarrow \infty$. Our main theoretical result stated in the following (informal) theorem compares the two latencies.

THEOREM 1. *The latency T_{LT} and computations C_{LT} of our LT coded distributed matrix-vector multiplication strategy in computing the product of a $m \times n$ matrix \mathbf{A} with a $n \times 1$ vector \mathbf{x} satisfy the following for large m :*

$$\Pr(T_{LT} > T_{ideal}) = p \exp\left(-\frac{\mu\tau m(\alpha - 1)}{p^2}\right) \quad (1)$$

$$\frac{\mathbb{E}[T_{LT}] - \mathbb{E}[T_{ideal}]}{\mathbb{E}[T_{ideal}]} = O\left(\exp\left(-\frac{\tau m(\alpha - 1)}{p^2}\right)\right) \quad (2)$$

$$\frac{\mathbb{E}[C_{LT}]}{\mathbb{E}[C_{ideal}]} = \frac{m(1 + \epsilon)}{m} \text{ where } \epsilon \rightarrow 0 \text{ as } m \rightarrow \infty. \quad (3)$$

where $m_e = \alpha m$ (for $\alpha > 1$) is the number of encoded rows, the initial delay at each worker is $X_i \sim \exp(\mu)$ and τ is the time taken to complete each row-vector product task. Due to the inherent design of LT codes, the overhead $\epsilon \rightarrow 0$ as $m \rightarrow \infty$.

This results shows that as long as the number of encoded rows m_e is sufficiently larger than m , despite not performing dynamic task assignment, the rateless coding strategy can seamlessly adapt to varying initial delays at the workers. Its runtime T_{LT} and C_{LT} asymptotically converge to the ideal strategy. The exact results are derived in Theorem 3 and Theorem 4.

Comparison with MDS and Replication Strategies. Unlike our rateless coding strategy, MDS-coded and replication-based strategies give strictly worse latency and cost than the ideal scheme and the gap does not go to zero. In Section 4 we analyze the expected latency and computations of these strategies. Fig. 1 shows simulation plots of the latency-computation trade-off of these strategies clearly demonstrating the superiority of using rateless LT codes.

Experimental Results. Fig. 2 shows the results of implementing our rateless coded strategy for a real distributed matrix-vector multiplication task on a cluster of 70 EC2 [1] workers deployed using Kubernetes [22]. The computation involves multiplying a 11760×9216 matrix \mathbf{A} extracted from the

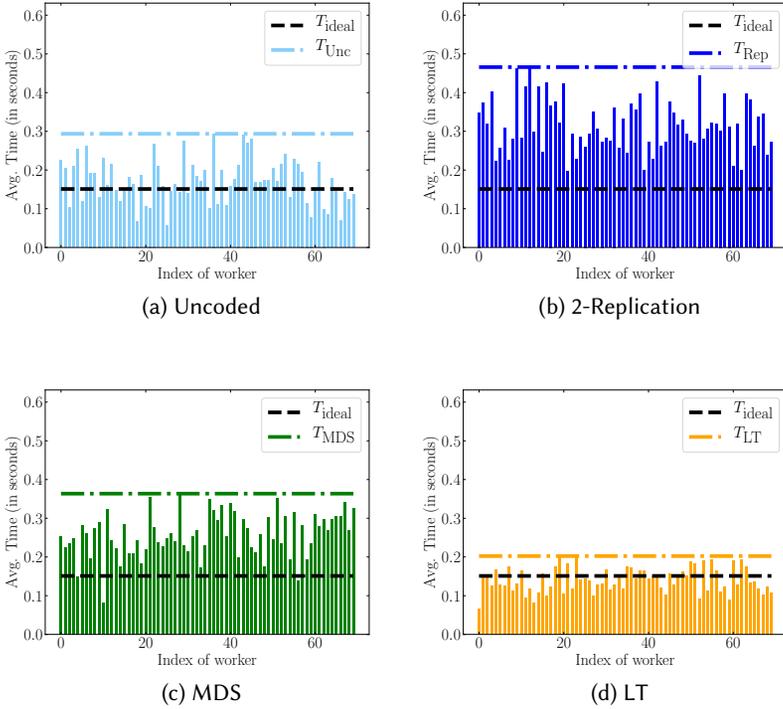


Fig. 2. Comparison of load balancing across different matrix-vector multiplication approaches. The rows of a 11760×9216 matrix A are encoded and distributed among 70 EC2 workers. The height of the bar plot for each worker indicates the time spent by the worker computing row-vector products either until it finishes its assigned tasks or is terminated by the master because the final matrix-vector product Ax has been successfully decoded. The dash-dot line indicates the overall latency (time at which matrix-vector product Ax can be successfully decoded) in each case, and the black dashed line is the latency of ideal load balancing. The LT coded approach exhibits near-ideal load balancing, and has lower latency than other approaches.

STL-10 dataset [6] with vectors extracted from the same dataset, and is implemented using Dask [8] a popular framework for parallel computing in Python. The proposed rateless coded strategy significantly outperforms the uncoded ($3 \times$ speedup) and MDS coded ($2 \times$ speedup) approaches. The plots in Fig. 2 also show that the variability in individual worker times is significantly lower for our rateless coded strategy (Fig. 2d) than for other approaches as fast nodes perform more tasks than slow nodes under our approach leading to much better load balancing. The latency of each approach is also compared to T_{ideal} , the latency of the ideal load-balancing strategy, approximated in this case as the minimum time required by the workers to compute 11760 encoded row-vector products in total. Observe that T_{LT} is closest to T_{ideal} . We also obtain similar improvements with LT coding in parallel computing using Python’s Multiprocessing library [17] library on a single machine, and in serverless computing on AWS Lambda [2] as described in Section 6.

1.4 Organization

The rest of the paper is organized as follows. Section 2 presents the system model, performance criteria and comparison benchmarks. Section 3 describes our rateless fountain coding strategy

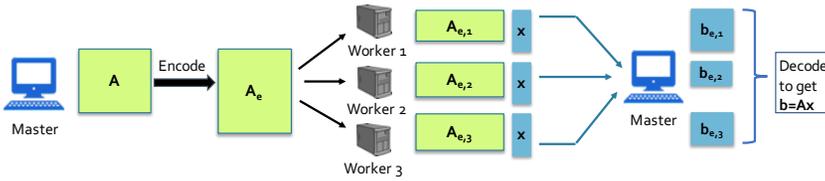


Fig. 3. The system model for coded distributed matrix vector multiplication with a master-worker framework. The master generates the encoded matrix A_e by applying a coding scheme to the rows of A . Worker i stores a submatrix of A_e denoted by $A_{e,i}$ and sends encoded row-vector products $b_{e,i}$ to the master ($i = 1, \dots, p$). Different $b_{e,i}$'s may have different sizes. The master decodes the encoded row-vector products in $b_{e,i}$, $i = 1, \dots, p$ to recover $\mathbf{b} = \mathbf{A}\mathbf{x}$

for distributed matrix-vector multiplication. Section 4 shows theoretical analyses and a latency-cost comparison of rateless coding with other strategies. Section 5 extends these results to the queuing setting where vectors $\mathbf{x}_1, \mathbf{x}_2, \dots$ that need to be multiplied with matrix A arrive at rate λ . Experimental results are presented in Section 6. All proofs are deferred to the Appendix.

2 PROBLEM FORMULATION

2.1 System Model

Consider the problem of multiplying a $m \times n$ matrix A with a $n \times 1$ vector \mathbf{x} using p worker nodes and a master node as shown in Fig. 3. The worker nodes can only communicate with the master, and cannot directly communicate with other workers. The goal is to compute the result $\mathbf{b} = \mathbf{A}\mathbf{x}$ in a distributed fashion and mitigate the effect of unpredictable node slowdown or straggling. The rows of A are encoded using an error correcting code to give the $m_e \times n$ encoded matrix A_e , where $m_e \geq m$. We denote the amount of redundancy added by the parameter $\alpha = m_e/m$. Matrix A_e is split along its rows to give p submatrices $A_{e,1}, \dots, A_{e,p}$ of equal size such that worker i stores submatrix $A_{e,i}$. To compute the matrix-vector product $\mathbf{b} = \mathbf{A}\mathbf{x}$, the vector \mathbf{x} is communicated to the workers such that Worker i is tasked with computing the product $A_{e,i}\mathbf{x}$.

To complete the assigned task, each worker needs to compute a sequence of row vector products of the form $\mathbf{a}_{e,j}\mathbf{x}$ where $\mathbf{a}_{e,j}$ is the j^{th} row of A_e . The time taken by a worker node to finish computing one or more row-vector products may be random due to variability in the node speed or variability in the amount of computation assigned to it. The master node aggregates the computations of all, or a subset of, the workers into the vector \mathbf{b}_e , which is then decoded to give the final result $\mathbf{b} = \mathbf{A}\mathbf{x}$. If \mathbf{b}_e is not decodable, the master waits until workers compute more row-vector products.

2.2 Performance Criteria

We use the following metrics to compare different distributed matrix-vector multiplication schemes via theoretical analysis and associated simulations (Section 4), and experiments in parallel, distributed, and serverless environments (Section 6).

DEFINITION 1 (LATENCY (T)). *The latency T is the time required by the system to complete enough computations so that $\mathbf{b} = \mathbf{A}\mathbf{x}$ can be successfully decoded from worker computations aggregated in \mathbf{b}_e .*

DEFINITION 2 (COMPUTATIONS (C)). *The number of computations C is defined as the total number of row-vector products $\mathbf{a}_{e,j}\mathbf{x}$ performed collectively by the worker nodes until $\mathbf{b} = \mathbf{A}\mathbf{x}$ is decoded.*

For any strategy we always have $C \geq m$ where m is the number of rows of A or the number of elements in \mathbf{b} .

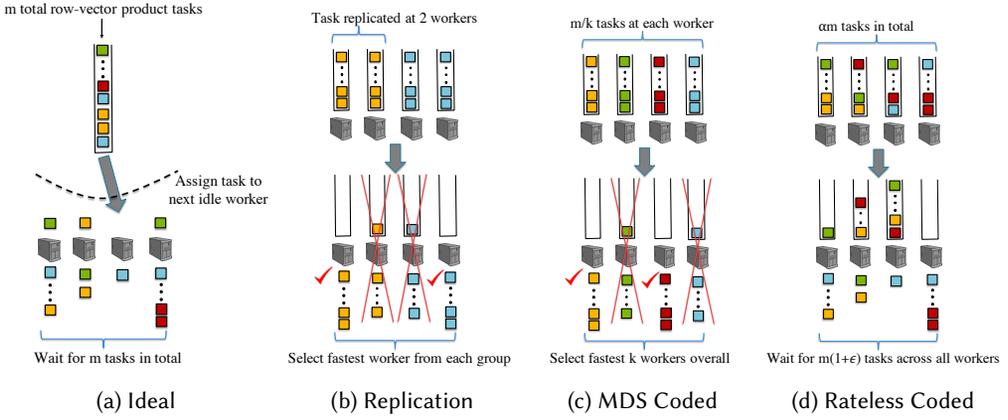


Fig. 4. Each square represents one row-vector product task out of a total of m tasks to be completed by the p workers. In the ideal scheme we have a central queue of m tasks and each worker is assigned a new task as soon as it becomes idle until all m tasks are completed. In the replication scheme, the master waits for the fastest worker for each sub-matrix. With MDS coding, the master needs to wait for k out of p workers, but each worker has to complete m/k tasks. The rateless coded strategy requires waiting for only $m(1 + \epsilon)$ tasks across *all* workers.

2.3 Benchmarks for Comparison

We compare the performance of the proposed rateless coded strategy with three benchmarks: ideal load balancing, r -replication, and the (p, k) MDS-coded strategy, which are described formally below. Fig. 4 illustrates the differences in the way row-vector product tasks are assigned to and collected from workers in each strategy.

Ideal Load Balancing. The multiplication of the $m \times n$ matrix A with the $n \times 1$ vector x can be treated as a job with m tasks, where each task corresponding to one row-vector product. In the ideal load balancing strategy, the master node maintains a central queue of these m tasks. It dynamically assigns one task to each of the p workers as soon as a worker finishes its previous task. The matrix-vector multiplication is complete when exactly m tasks are collectively finished by the workers. This strategy seamlessly adapts to varying worker speeds without performing any redundant computation ($C = m$); hence it gives the optimal latency-computation trade-off. This strategy may be impractical due to the constant communication between the master and the worker nodes. Nevertheless, it serves as a good theoretical benchmark for comparison with the rateless, replication and MDS strategies.

The r -Replication Strategy. A simple distributed multiplication strategy is to split A along its rows into p/r submatrices $A_1, \dots, A_{p/r}$, with rm/p rows each (assume that p/r divides m) and multiply each submatrix with x in parallel on r distinct worker nodes. The master collects the results from the fastest of the r nodes that have been assigned the task of computing the product $A_i x$ for all i . The computed products are aggregated into the $m \times 1$ vector b . *Setting $r = 1$ corresponds to the naive or uncoded strategy where A is split into p sub-matrices and each worker node computes the corresponding submatrix-vector product.* While this approach performs the least number of computations it is susceptible to straggling nodes or node failures. Increasing the number of replicas provides greater straggler tolerance at the cost of redundant computations. Real distributed computing frameworks like MapReduce [11] and Spark [65] often use $r = 2$ i.e. each computation is assigned to 2 different worker nodes for added reliability and straggler tolerance.

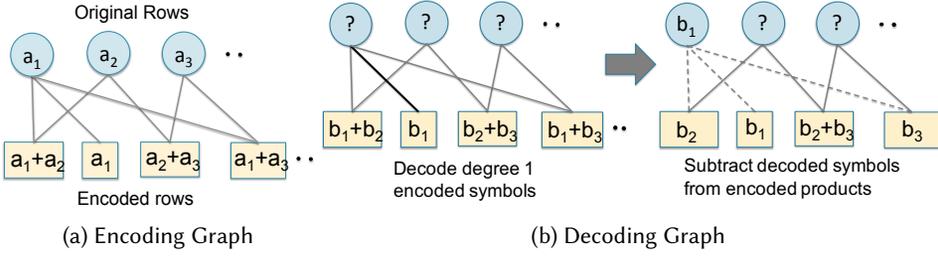


Fig. 5. (a) Bipartite graph representation of the encoding of the rows a_1, a_2, \dots, a_m of matrix A . Each encoded row is the sum of d rows of A chosen uniformly at random, where d is drawn from the Robust Soliton degree distribution given by (4). (b) In each step of the iterative decoding process, a single degree one encoded symbol is decoded directly, and is subtracted from all sums in which it participates.

The (p, k) MDS Coded Strategy. Recent works like [41, 44] have applied MDS coding to overcome the problem of stragglers in the uncoded strategy. The strategy involves pre-multiplying A at the central node with a suitable encoding matrix F denoting the MDS codes. For encoding using a (p, k) MDS code, the matrix A is split along its rows into k matrices A_1, \dots, A_k , each having m/k rows. The MDS code adds $p - k$ redundant matrices A_{k+1}, \dots, A_p which are independent linear combinations of the matrices A_1, \dots, A_k . Worker i computes the product $A_i x$. Thus the system is robust to $p - k$ stragglers. However this strategy adds a significant computation overhead. When none of the nodes are slow, the system performs mp/k row-vector products (as opposed to m row-vector products in the uncoded case).

3 PROPOSED RATELESS CODED STRATEGY

We describe how rateless codes, specifically LT codes [45], can be applied to perform coded matrix vector multiplication, and then propose a distributed implementation of this scheme for straggler mitigation in computing the matrix-vector product $b = Ax$ using the master-worker framework of Section 2.1.

3.1 LT-Coded Matrix-vector Multiplication

Luby Transform (LT) codes proposed in [45] are a class of erasure codes that can be used to generate a limitless number of encoded symbols from a finite set of source symbols. We apply LT codes to matrix-vector multiplication by treating the m rows of the matrix A as source symbols. Each encoded symbol is the sum of d source symbols chosen uniformly at random from the matrix rows. Thus if $S_d \subseteq \{1, 2, \dots, m\}$ is the set of d row indices, the corresponding encoded row is $a_e = \sum_{i \in S_d} a_i$.

The number of original rows in each encoded row, or the degree d , is chosen according to the Robust Soliton degree distribution

$$\rho(d) = \begin{cases} \frac{R}{dm} + \frac{1}{m} & \text{for } d = 1 \\ \frac{R}{dm} + \frac{1}{m(m-1)} & \text{for } d = 2, \dots, m/R - 1 \\ \frac{R \ln(R/\delta)}{m} + \frac{1}{m(m-1)} & \text{for } d = m/R \\ \frac{1}{m(m-1)} & \text{for } d = m/R + 1, \dots, m \end{cases} \quad (4)$$

where $R = c \log(m/\delta) \sqrt{m}$ for some $c > 0$ and $\delta \in [0, 1]$, with c and δ being design parameters. Some guidelines for choosing c and δ can be found in [46]. The probability of choosing $d = d_0$ is equal to $\rho(d_0) / \sum_{i=1}^m \rho(i)$. Once the degree d is chosen, encoding is performed by choosing d source symbols

uniformly at random (this determines \mathcal{S}_d) and adding them to generate an encoded symbol. The encoding process is illustrated in Fig. 5a.

Once the rows of the encoded matrix \mathbf{A}_e are generated, we can compute the encoded matrix vector product $\mathbf{b}_e = \mathbf{A}_e \mathbf{x}$. To decode the desired matrix vector product $\mathbf{b} = \mathbf{A} \mathbf{x}$ from a subset of M' symbols of \mathbf{b}_e we use the *iterative peeling decoder* described in [45, 52, 53]. If $\mathbf{b} = [b_1, b_2, \dots, b_m]$, the decoder may receive symbols $b_1 + b_2 + b_3, b_2 + b_4, b_3, b_4$, and so on since each row of \mathbf{A}_e is a sum of some rows of \mathbf{A} . Decoding is performed in an iterative fashion. In each iteration, the decoder finds a degree one encoded symbol, covers the corresponding source symbol, and subtracts the symbol from all other encoded symbols connected to that source symbols. This decoding process is illustrated in Fig. 5b.

Since the encoding uses a random bipartite graph, the number of symbols required to decode the m source symbols successfully is a random variable M' which we call the decoding threshold,

DEFINITION 3 (DECODING THRESHOLD (M')). *The decoding threshold M' is the number of encoded symbols required to decode a set of m source symbols using the rateless coding strategy.*

For the Robust Soliton distribution, [45] gives the following high probability bound on M' .

LEMMA 1 (THEOREMS 12 AND 17 IN [45]). *The original set of m source symbols can be recovered from a set of any $M' = m + O(\sqrt{m} \ln^2(m/\delta))$ with probability at least $1 - \delta$.*

REMARK 1. While \mathbf{A} can be encoded using any random linear code to ensure successful decoding of \mathbf{b} from m symbols of \mathbf{b}_e with a high probability, the key benefit of using LT codes is the low decoding complexity owing to the careful design of the Robust Soliton distribution. The complexity of LT decoding is $O(m \ln m)$ while for any other random linear code it would be $O(m^3)$ which is unacceptable for large m . (see Appendix A)

3.2 Distributed Implementation

The $m \times n$ matrix \mathbf{A} is encoded to generate an $m_e \times n$ encoded matrix \mathbf{A}_e where $m_e = \alpha m$. Each row of \mathbf{A}_e is the sum of a random subset of rows of \mathbf{A} as described in Section 3.1. The knowledge of the mapping between the rows of \mathbf{A} and the rows of \mathbf{A}_e is crucial for successful decoding as illustrated in Figures 5a and 5b. Hence this mapping is stored at the master. The encoding step can be treated as a pre-processing step in that it is only performed initially.

The αm rows of the encoded matrix are distributed equally among the p worker nodes as illustrated in Fig. 3. To multiply \mathbf{A} with a vector \mathbf{x} , the master communicates \mathbf{x} to the workers. Each worker multiplies \mathbf{x} with each row of \mathbf{A}_e stored in its memory and returns the product (a scalar) to the master. The master collects row-vector products of the form $\mathbf{a}_{e,j} \mathbf{x}$ (elements of \mathbf{b}_e) from the workers until it has enough elements to be able to recover \mathbf{b} . If a worker node completes all the $\alpha m/p$ row-vector products assigned to it before the master is able to decode \mathbf{b} , it will remain idle, while the master collects more row-vector products from other workers.

Once the master has collected a sufficient number of coded row-vector products from the workers it can recover the desired matrix vector product $\mathbf{b} = \mathbf{A} \mathbf{x}$ from the subset of the elements of $\mathbf{b}_e = \mathbf{A}_e \mathbf{x}$ that it has collected using the iterative peeling decoder. Once the master decodes all elements of the product vector $\mathbf{b} = \mathbf{A} \mathbf{x}$, it sends a *done* signal to all workers nodes to stop their local computation.

The following modifications can make the current implementation even more efficient in real systems:

(1) **Blockwise Communication:** To truly monitor the partial work done by each worker the master needs to receive each encoded row-vector product $\mathbf{a}_{e,j} \mathbf{x}$ from the workers. However this imposes a large communication overhead which may increase latency in a slow network. To prevent this, in our distributed computing experiments, we communicate submatrix-vector products $\mathbf{A}_{e_i}^j \mathbf{x}$

where $\mathbf{A}_{e_i}^j$ is the j^{th} part of the encoded submatrix \mathbf{A}_{e_i} stored at worker i , and each part corresponds to approximately 10% of the total rows of the submatrix. Note that if \mathbf{A} is very large then it will not be feasible for worker i to read the entire submatrix \mathbf{A}_{e_i} from memory at once. As a result $\mathbf{A}_{e_i}\mathbf{x}$ needs to be computed in parts for *any* coding scheme.

(2) **Using Raptor Codes:** Despite their ease of implementation and fast decoding, LT codes [45] are sub-optimal in practice due to the overhead of $M' - m$ extra symbols required to decode the original m source symbols. In our experiments we observe that for a matrix \mathbf{A} with $m = 11760$ rows, we need to wait for 12500 encoded row-vector products to decode $\mathbf{b} = \mathbf{A}\mathbf{x}$ with 99% probability. Advanced rateless codes like Raptor Codes [52] can decode m source symbols from $m(1 + \epsilon)$ symbols for any *constant* ϵ even for *finite values* of m . Since Raptor Codes are the rateless codes used in practical wireless standards [53] we expect them to be used in practical implementations of our coded distributed matrix vector multiplication strategy to improve efficiency.

(3) **Using Systematic Rateless Codes:** We can entirely avoid decoding (in the absence of significant straggling) by using Systematic LT/Raptor Codes [52] where the m source rows $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m$ form a subset of the encoded rows in \mathbf{A}_e . The overall scheme can be designed so that each worker first computes the row-vector products corresponding to the systematic symbols $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m$ and then computes other encoded products (in the event of node slowdown). This would preclude the need for decoding if there is no/little straggling thereby reducing the overall latency.

4 PERFORMANCE ANALYSIS

In this section we theoretically analyze the performance of LT coding and the three benchmark strategies – ideal load balancing, (p, k) -MDS, and r -Replication – in terms of latency (Definition 1) and computations (Definition 2). Our results are summarized in Table 1 and the proofs of the theoretical results are contained in Appendix C. We begin by describing our delay model.

4.1 Delay Model

We assume that worker i requires time Y_i to perform B_i row-vector product computations where

$$Y_i = X_i + \tau B_i, \quad \text{for all } i = 1, \dots, p \quad (5)$$

Thus, the delay involves the sum of two components: 1) a random variable X_i that includes initial setup time at the worker before it actually begins performing the computations, and 2) a shift that is linear in the number of computations performed at the worker. This delay model is motivated by the observations of [10] where it is noted that the variability in latency arises largely from delays due to background tasks running at worker nodes and that once a request actually begins execution, the variability is considerably lower. When X_i is exponentially distributed with rate μ , the time taken by worker i to perform b computations is distributed as

$$\Pr(Y_i \leq t) = 1 - \exp(-\mu(t - \tau b)). \quad (6)$$

While this follows the shifted exponential delay models used in [41], [14] and [15], the key difference is that the shift is parameterized by the number of computations at each worker. We believe this is a more realistic model as it captures the effect of increasing the amount of computations on the delay – if a worker is assigned more computations, there is larger delay. Moreover, unlike previous works, the decay rate μ of the exponential part of the delay does not change with the number of computations performed by that worker. Fig. 6 illustrates our delay model.

Also, in our analysis, we use $X_{k:p}$ to denote the k^{th} order statistic i.e. the k^{th} smallest of p random variables X_1, \dots, X_p and we define $U_l = X_{l+1:p} - X_{l:p}$, $l = 1, \dots, p-1$ as the difference of consecutive order statistics. We also use the notation $H_j = \sum_{v=1}^j 1/v$, for the j^{th} Harmonic number.

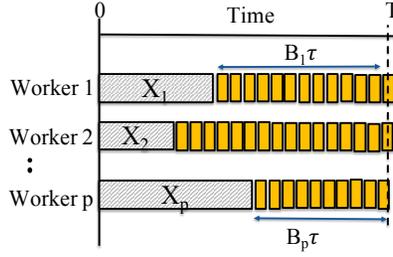


Fig. 6. Worker i has a random initial delay X_i , after which it completes row-vector product tasks (denoted by the small rectangles), taking time τ per task. The latency T is the time until enough tasks have been completed for the product $\mathbf{b} = \mathbf{A}\mathbf{x}$ to be recovered.

Strategy	Latency	# of Comp	Complexity
Ideal	$\frac{\tau m}{p} + \frac{1}{\mu}$	m	$O(m)$
LT (large α)	$\frac{\tau m(1+\epsilon)}{p} + \frac{1}{\mu}$	$m(1 + \epsilon)$	$O(m \log m)$
r -Replication	$\frac{\tau r m}{p} + \frac{1}{r\mu} \log \frac{p}{r}$	rm	$O(m)$
(p, k) MDS	$\frac{\tau m}{k} + \frac{1}{\mu} \log \frac{p}{p-k}$	mp/k	$O(mk + k^3)$

Table 1. Comparison of different strategies to multiply a $m \times n$ matrix \mathbf{A} with vector \mathbf{x} using p worker nodes. The latency values are approximate, and number of computations values are for the case when none of the nodes slowdown.

4.2 Ideal Load Balancing Strategy

Recall that in the ideal load balancing strategy, we have a central queue at the master and tasks being allocated to a worker as soon as it becomes idle (either immediately after the initial delay or after it completes the current task) as illustrated in Fig. 6. Thus it computes exactly $C = m$ row-vector products in total when an $m \times n$ matrix is multiplied with $n \times 1$ vector and performs zero redundant computations. Theorem 2 below proves the optimality of ideal load balancing in terms of latency and Lemma 2 and Corollary 1 give bounds on the expected latency.

THEOREM 2 (OPTIMALITY OF IDEAL LOAD BALANCING). *For any distributed matrix-vector multiplication scheme, for the delay model of (5), the latency T is no less than the latency of ideal load balancing, denoted by T_{ideal} . In other words, for any scheme,*

$$T \geq T_{ideal}. \quad (7)$$

LEMMA 2 (LATENCY OF IDEAL LOAD BALANCING). *The latency for the ideal load balancing strategy with p workers has the following upper and lower bounds.*

$$\frac{\tau m}{p} + X_{1:p} \leq \mathbb{E}[T_{ideal}] \leq \frac{\tau m}{p} + \frac{1}{p} \sum_{i=1}^p X_i + \tau. \quad (8)$$

COROLLARY 1. *The expected latency for the ideal strategy with $X_i \sim \exp(\mu)$ for all workers $i = 1, \dots, p$, has the following upper and lower bounds.*

$$\frac{\tau m}{p} + \frac{1}{p\mu} \leq \mathbb{E}[T_{ideal}] \leq \frac{\tau m}{p} + \frac{1}{\mu} + \tau. \quad (9)$$

Note that the ideal load balancing scheme is not exactly realizable in practice. Approaches like work stealing [13, 26] can potentially approximate this strategy by physically moving tasks from

busy workers to idle workers. However implementing such approaches may not be feasible in all settings, for e.g. when the communication latency between workers is too large, or the data is restricted to lie on a particular worker due to privacy concerns. In this work we aim to show that it is possible to *algorithmically* achieve near-ideal latency performance for matrix-vector multiplication by using the rateless coded computing strategy described in Section 3 which does not require physically moving data between workers.

4.3 Rateless Coded Strategy

We make the following assumption for analyzing the latency of the proposed rateless coded strategy.

ASSUMPTION 1. *The decoding threshold M' (Definition 3) of the LT coded strategy satisfies $M' \simeq m$.*

We believe the above assumption is reasonable because the problem of distributed matrix vector multiplication arises only when m (the number of rows of \mathbf{A}) is large and the high probability bound of Lemma 1 can be used to show that $\mathbb{E}[M'] = m(1 + \epsilon)$, where $\epsilon \rightarrow 0$ as $m \rightarrow \infty$. Note that this assumption is only to facilitate a better *theoretical* comparison between the LT coded and ideal strategies. In our experiments in Section 6 we choose a value of M' according to Lemma 1 that is slightly larger than m and ensures that the original matrix-vector product $\mathbf{b} = \mathbf{A}\mathbf{x}$ can be recovered with high ($> 99\%$) probability.

REMARK 2. The Rateless coded computing strategy described in Section 3 is identical to the ideal load balancing strategy described above for large values of m and infinite redundancy i.e. $\alpha = m_e/m \rightarrow \infty$. This is because both the rateless coded strategy and ideal load balancing strategies are based on collecting a pre-determined number of computations across all workers by greedily picking the next available task at each worker.

However in practice, we cannot set $\alpha = m_e/m \rightarrow \infty$ owing to limitations in computation power and memory at workers. Instead the amount of redundancy in the LT coded strategy is fixed initially by choosing the number of encoded rows $m_e = \alpha m/p$ for some $\alpha > 1$. Computations are divided equally among the p workers and thus each worker can perform a maximum of $\alpha m/p$ computations.

THEOREM 3 (RATELESS V/S IDEAL). *The latency of the proposed rateless coded strategy, T_{LT} decreases on increasing α and approaches the latency of the ideal strategy T_{ideal} . This is quantified by the following probabilistic upper bound:*

$$\Pr(T_{LT} > T_{ideal}) \leq \sum_{j=2}^p \Pr\left(\sum_{l=1}^{j-1} U_l \geq \frac{\tau m(\alpha - 1)}{p - 1}\right), \quad (10)$$

where $U_l = X_{l+1:p} - X_{l:p}$.

REMARK 3. The effect of straggling is captured through the term $U_l = X_{l+1:p} - X_{l:p}$ in the above expression. High straggling, implies a high variability in the initial worker delays X_i due to which U_l is large and the probability of $T_{LT} > T_{ideal}$ is also higher.

Thus as α (and consequently m_e) increases, T_{LT} is equal to T_{ideal} with a high probability. A cleaner result is obtained for the case when $X_i \sim \exp(\mu)$, as given below.

COROLLARY 2. *If $X_i \sim \exp(\mu)$ for all workers $i = 1, \dots, p$ then*

$$\Pr(T_{LT} > T_{ideal}) \leq p \exp\left(-\frac{\mu \tau m(\alpha - 1)}{p^2}\right). \quad (11)$$

We also derive an upper bound on the difference between the expected latencies of the rateless and the ideal strategies. We only state the result for $X_i \sim \exp(\mu)$ over here and defer the (more complicated) general result and its proof to Appendix C.

THEOREM 4 (LATENCY OF THE RATELESS CODED STRATEGY). *If $X_i \sim \exp(\mu)$ for all workers $i = 1, \dots, p$ then*

$$\mathbb{E}[T_{LT}] - \mathbb{E}[T_{ideal}] \leq \left(\tau \alpha m p^2 + \frac{p^2}{\mu} + \tau p \right) \exp\left(-\frac{\mu \tau m (\alpha - 1)}{p^2}\right). \quad (12)$$

The second term decays exponentially and dominates the first term, which is a polynomial. The rate of decay increases as the amount of redundancy α increases. In other words, $\mathbb{E}[T_{LT}]$ approaches $\mathbb{E}[T_{ideal}]$ exponentially fast as redundancy is increased.

REMARK 4. Another important advantage of the rateless coded strategy is that the number of computations performed by the workers, C_{LT} , is always equal to M' (the decoding threshold, defined in Definition 3), and does not increase on increasing redundancy (increasing α) unlike for the MDS and Replication strategies. Moreover since $\mathbb{E}[M'] = m(1 + \epsilon)$ and $\epsilon \rightarrow 0$ as $m \rightarrow \infty$, $\mathbb{E}[C_{LT}]$ asymptotically approaches the minimum number of computations (m) required to recover a m -dimensional matrix-vector product.

In the following subsections (and in Appendix E) we show that the latency of the MDS and Replication strategies is much larger than that of ideal load balancing and does not converge to T_{ideal} on increasing redundancy. We also show that the number of computations performed by both replication and MDS coding in computing $\mathbf{b} = \mathbf{A}\mathbf{x}$ is much larger than m .

4.4 MDS Coded Strategy

Recall that for the (p, k) MDS coded strategy, the encoded submatrices $\mathbf{A}_{\epsilon_1}, \dots, \mathbf{A}_{\epsilon_p}$ are generated by applying a (p, k) MDS Code to submatrices $\mathbf{A}_1, \dots, \mathbf{A}_k$. The master then waits for the fastest k workers to complete all the tasks assigned to them.

LEMMA 3 (LATENCY OF THE MDS CODED STRATEGY). *The latency of the (p, k) MDS-coded strategy, T_{MDS} , is given by*

$$T_{MDS} = X_{k:p} + \tau \frac{m}{k}. \quad (13)$$

COROLLARY 3. *The expected latency of the (p, k) MDS-coded strategies with $X_i \sim \exp(\mu)$ for all workers $i = 1, \dots, p$ is*

$$\mathbb{E}[T_{MDS}] = \frac{\tau m}{k} + \frac{1}{\mu} (H_p - H_{p-k}) \approx \frac{\tau m}{k} + \frac{1}{\mu} \log \frac{p}{p-k}. \quad (14)$$

Observe that in (14) above, adding redundancy (reducing k) leads to an increase in the first term (more computation at each node) and decrease in the second term (less delay due to stragglers). Thus, straggler mitigation comes at the cost of additional computation at the workers which might even lead to an increase in latency. This is in contrast to Theorems 3 and 4 which indicates that the expected latency of the rateless coded strategy always decreases on adding redundancy (increasing α). Moreover, the presence of the log-factor in the second term causes T_{MDS} to always be larger than T_{ideal} since there is no log-factor in the term containing $1/\mu$ in the upper bound on T_{ideal} (Lemma 2).

We now analyze the number of computations performed by the (p, k) MDS coding. The following result shows that with a high probability, the number of computations performed by the MDS Coded strategy is very close to the worst-case number of computations (mp/k) i.e. when all the workers perform all the tasks assigned to them in time T_{MDS} .

LEMMA 4 (TAIL OF COMPUTATIONS FOR MDS CODING). *The tail of the number of computations of the MDS coded strategy, C_{MDS} , with p workers and the delay model of (5) is bounded as*

$$\Pr(C_{MDS} \geq \frac{mp}{k} - C_0) \geq 1 - \Pr\left(\sum_{l=k}^{p-1} U_l \geq \frac{\tau C_0}{p-k} - \tau\right). \quad (15)$$

When $X_i \sim \exp(\mu) \forall i = 1, \dots, p$ this reduces to

$$\Pr(C_{MDS} \geq \frac{mp}{k} - C_0) \geq 1 - \exp\left(-\mu \left(\frac{\tau C_0}{(p-k)^2} - \frac{\tau}{p-k}\right)\right). \quad (16)$$

Even for a small value of C_0 in the above expression, $\Pr(C_{MDS} \geq \frac{mp}{k} - C_0)$ can be very large. Thus the overhead $C_{MDS} - m$ is quite large (we only need m computations in the uncoded case to reconstruct the m -dimensional matrix-vector product).

4.5 Replication Strategy

The r -Replication strategy involves replicating each of the p/r submatrices $A_1, \dots, A_{p/r}$ at r distinct workers and selecting the result of the fastest worker for each submatrix.

LEMMA 5 (LATENCY OF THE REPLICATION STRATEGY). *The latency of the r -Replication strategy, T_{rep} is given by*

$$T_{rep} = \max_{1 \leq i \leq p/r} \min_{1 \leq j \leq r} X_{(i-1)r+j} + \frac{\tau mr}{p}. \quad (17)$$

COROLLARY 4. *The expected latency of the r -Replication with $X_i \sim \exp(\mu)$ for all workers $i = 1, \dots, p$ is*

$$\mathbb{E}[T_{rep}] = \frac{\tau mr}{p} + \frac{1}{\mu} H_{p/r} \approx \frac{\tau mr}{p} + \frac{1}{\mu} \log \frac{p}{r}. \quad (18)$$

Once again we see that adding redundancy (increasing r) leads to an increase in the first term (more computation at each node) and decrease in the second term (less delay due to stragglers). Thus the extra computation at the workers may lead to an increase in latency even in this case. Moreover the log-factor in the second term causes T_{Rep} to always be larger than T_{ideal} , just like T_{MDS} .

LEMMA 6 (TAIL OF COMPUTATIONS FOR REPLICATION). *The tail of the number of computations of the replication strategy, C_{rep} , with p workers and the delay model of (5) is bounded as*

$$\Pr(C_{rep} \geq mr - C_0) \geq 1 - \Pr\left(\sum_{i=1}^{p/r} \sum_{j=1}^{r-1} (V_{j+1:r}^i - V_{j:r}^i) \geq \frac{\tau C_0}{r-1} - \frac{\tau p}{r}\right), \quad (19)$$

where $V_j^i = X_{(i-1)r+j}$ and $V_{j:r}^i$ are the corresponding order statistics. When $X_i \sim \exp(\mu) \forall i = 1, \dots, p$ this reduces to

$$\Pr(C_{rep} \geq mr - C_0) \geq 1 - \sum_{i=0}^{p/r-1} \frac{1}{i!} \exp(-\mu\theta)(\mu\theta)^i, \quad (20)$$

$$\text{where } \theta = \frac{\tau C_0}{(r-1)^2} - \frac{\tau p}{r(r-1)}. \quad (21)$$

Thus with a high probability, the number of computations performed is very close to the worst-case number (mr) i.e. when all the workers perform all the tasks assigned to them in time T_{rep} .

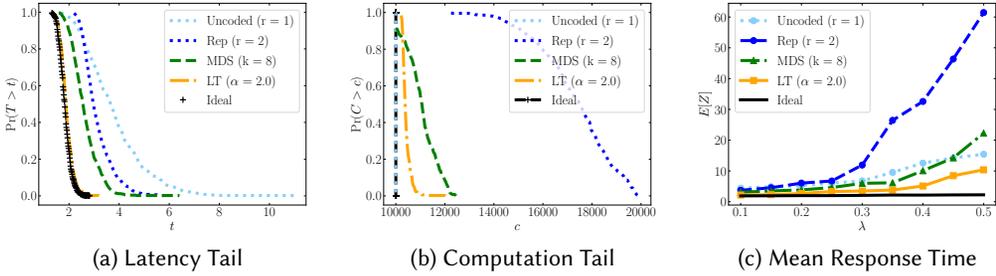


Fig. 7. The tail probability of the latency is the highest for the replication schemes. MDS codes perform better in terms of latency but they perform a large number of redundant computations. The latency tail of LT codes is the minimum among all the schemes. Moreover the LT coded schemes performs significantly fewer redundant computations than MDS Codes or replication. When there are multiple jobs in the queue, the mean response time is least for the LT Coded setting under all values of arrival rate λ . All simulations are performed for a distributed matrix-vector multiplication task with $m = 10000$ matrix rows, $p = 10$ worker nodes, and delay model parameters $\mu = 1.0$, $\tau = 0.001$.

REMARK 5. While the benefits of using partial work from all workers can be obtained by using any random linear code on the rows of \mathbf{A} , the key strength of LT codes is their low $O(m \ln m)$ decoding complexity. Using an (m_e, m) MDS code on the rows of \mathbf{A} has $O(m^3)$ decoding complexity which is unacceptable for large m .

We simulate the MDS, replication, and LT-coded schemes under our delay model (5) for distributed matrix-vector multiplication with $m = 10000$ matrix rows, $p = 10$ workers and delay model parameters $\mu = 1.0$, $\tau = 0.001$ (Fig. 7). We limit the amount of redundancy to $\alpha = m_e/m \leq 2.0$ since this is the amount of redundancy in the basic 2-replication scheme. Observe that the LT coded strategy ($\alpha = 2.0$) clearly outperforms MDS coding (with $k = 8$) in that it not only exhibits near-ideal latency (Fig. 7a) but also performs fewer total computations (Fig. 7b) than MDS coding. Changing k does not improve the performance of MDS coding much. Specifically, increasing redundancy (reducing k) in MDS coding leads to higher latency after a point, as illustrated in Fig. 1 (and as expected from Lemma 3). On the other hand, the latency of LT coding converges to that of the Ideal scheme on increasing α , without any increase in computations. Additional simulations for $X_i \sim \text{Pareto}(1, 3)$ given in Fig. 11 in Appendix F also show similar improvements with LT coding.

5 QUEUEING ANALYSIS

In most real applications of matrix-vector multiplication in machine learning and data analytics, the matrix representing the model is fixed, while vectors representing the data that need to be multiplied with this matrix arrive as a real-time stream. Prior works on coded computing like [14, 41] do not consider the effect of multiple incoming jobs which could lead to queueing delays at the master, in addition to straggling at workers. We analyze the latency with queueing for the proposed LT coded strategy as well as the MDS coded and Replication strategies. The LT coding results are presented below, while MDS and replication results are given in Appendix D.

Suppose that vectors $\mathbf{x}_1, \mathbf{x}_2, \dots$ arrive according to a Poisson process with rate λ and are broadcast by the master to the p workers. Worker i multiplies each vector it receives with the sub-matrix \mathbf{A}_{e_i} stored in its memory (where \mathbf{A}_e is generated according to the corresponding encoding strategy) and communicates the corresponding elements of $\mathbf{b}_e = \mathbf{A}_e \mathbf{x}$ to the master. Once the master has

enough elements of \mathbf{b}_e to successfully decode \mathbf{b} , the remaining tasks at all the workers are cancelled. Then the mean response time $\mathbb{E}[Z]$ (waiting time in queue plus service time) of a matrix-vector multiplication job is as follows.

THEOREM 5 (LATENCY OF LT CODING WITH QUEUEING OF JOBS). *For large m_e i.e. $\alpha = m_e/m \rightarrow \infty$, the mean response time of the LT coded scheme Z_{LT} when vectors \mathbf{x} are arriving at rate λ according to a Poisson process is*

$$\mathbb{E}[Z_{LT}] = \mathbb{E}[T_{LT}] + \frac{\lambda \mathbb{E}[(T_{LT})^2]}{2(1 - \lambda \mathbb{E}[T_{LT}])}, \quad (22)$$

where $\mathbb{E}[T_{LT}]$ is bounded as described in Lemma 2 and bounds on $\mathbb{E}[(T_{LT})^2]$ are derived in Appendix D.

The proof is given in Appendix D. The key idea used in the proof is that for large α this system becomes equivalent to an M/G/1 queue with service time T_{LT} . Then we simply apply the Pollaczek-Khinchine formula [25] for the mean response time of M/G/1 queues. When α is small the analysis becomes very difficult – it is a generalization of the fork-join queueing system, whose response time is notoriously hard to analyze [37, 47, 57]. This is an open question for future research.

For the MDS and replication strategies, we reduce the queueing system to a fork-join queueing system with redundancy, and then use previous results [32, 35] to obtain bounds on the mean response time. The results are presented in Appendix D.

REMARK 6 (INSIGHTS FROM MDS AND REPLICATION QUEUEING ANALYSES). In the MDS and replication strategies, increasing redundancy (lower k and higher r) reduces the number of workers that need to complete their tasks. However, the added redundancy increases the number of computations that each worker needs to perform due to which the waiting time for incoming jobs at the master increases, thus increasing the overall mean response time (Z_{MDS} and Z_{rep} respectively). On the other hand for the rateless coded (and ideal) strategies we just need to wait for M' (or m) computations across all workers for each job. Moreover, for LT coding, adding redundancy (increasing m_e) always reduces the service time for each job (Theorem 3) and thus the overall queueing delay Z_{LT} always decreases on adding redundancy.

Fig. 7c shows simulation results of mean response time Z under our delay model with $X \sim \exp(1)$ and $\tau = 0.001$ for a distributed matrix-vector multiplication task with $m = 10000$ matrix rows using $p = 10$ worker nodes. The mean response time is averaged over 10 trials with 100 jobs in each trial. Jobs arrive according to a Poisson process with rate $\lambda \in (0.1, 0.6)$. The results illustrate that the benefits of our LT coded strategy over previous approaches are further enhanced when there is queueing of jobs.

6 EXPERIMENTAL RESULTS

We demonstrate the effectiveness of rateless codes in speeding up distributed matrix-vector multiplication in parallel, distributed and serverless environments. No artificial delays/background tasks were added to induce straggling in any experiments.

6.1 Parallel Computing Experiments

We consider multiplication of a 10000×10000 matrix \mathbf{A} of random integers with a 10000×1 vector \mathbf{x} of random integers on an iMac Desktop with 8 GB of RAM and a 3.6 GHz Intel i7 Processor. This computation is parallelized over 100 processes using Python's Multiprocessing Library [17]. We compare the uncoded, 2-replication, MDS coding ($k = 80, 50$) and LT coding ($\alpha = 1.25, 2.0$) approaches. For fair comparison, we consider instances of MDS and LT codes that have the same number of encoded rows m_e . The encoded matrix \mathbf{A}_e was divided equally among the $p = 100$

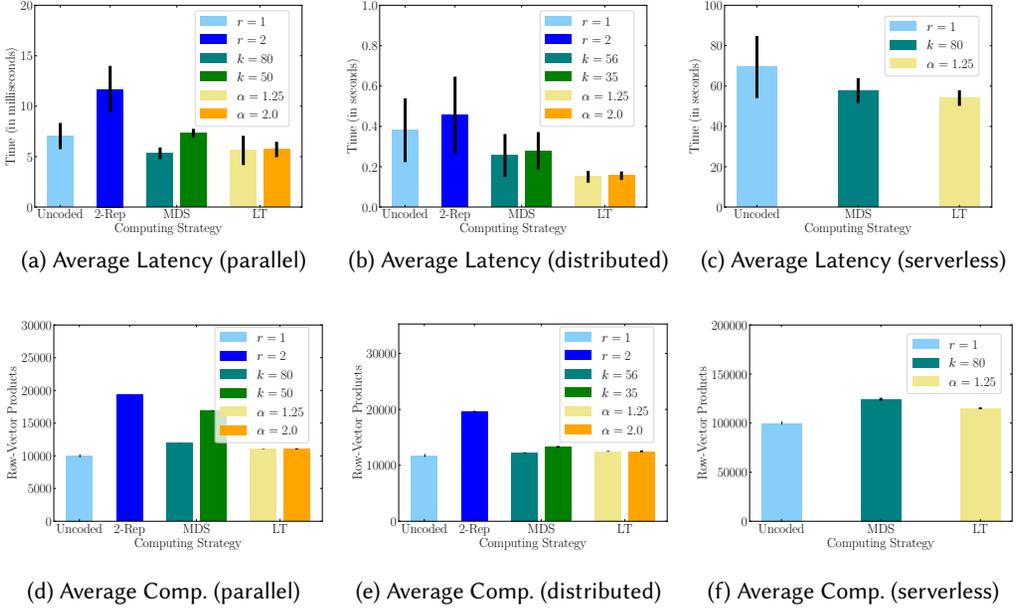


Fig. 8. Experiments on coded distributed matrix vector multiplication in parallel (Python Multiprocessing [17]), distributed (AWS EC2 [1]) and serverless (AWS Lambda [2]) settings show that the LT Coded strategy has lower average latency than all other approaches ($1.2 \times$ – to $3 \times$ – improvement across scenarios) and performs fewer total computations than Replication or MDS Coding. Each error bar corresponds to 1 standard deviation.

processes and the experiment was repeated 10 times with a different random x each time. The processes multiply the rows of A_c with x in parallel and we record the average latency (time required to collect enough row-vector products for successful decoding) and total computations. Results of average latency (Fig. 8a) show that LT coded and MDS coded approaches are clearly faster (about at least $1.2 \times$ –) than the Uncoded and 2-Replication approaches while Fig. 8d shows that the LT coded approaches also perform fewer total computations than the MDS or 2-replication strategies thus leading to more efficient resource utilization. Note that while MDS coding with $k = 80$ has latency comparable to that of LT coding (both for $\alpha = 1.25$ and $\alpha = 2.0$), both latency and total computations with MDS coding increase on increasing k to 50 due to the higher computational load at each node (as discussed in Section 4). Recall that k corresponds to the number of "fast" workers in the system. In most real systems the number of "fast" workers is transient and thus, unpredictable. Our experiments show that MDS coding is highly sensitive to the choice of k with incorrect choices leading to *higher* latency. LT Coding on the other hand is not only fast, but is also insensitive to the amount of redundancy (α) in that, the system designer can choose α to be as large as permitted by memory constraints without a risk of loss in performance (unlike MDS).

6.2 Distributed Computing Experiments

We created a cluster of 70 t2.small workers on AWS EC2 [1] using Kubernetes [22]. Each worker was allocated 1 GB of memory. Computations were performed using Dask [8], a popular framework for parallel computing in Python. A 11760×9216 matrix A was extracted from the STL-10 [6]

dataset. Once again we compared the uncoded, 2-replication, MDS coding ($k = 56, 35$) and LT coding ($\alpha = 1.25, 2.0$) approaches. The encoded matrix A_e is divided equally among the $p = 70$ workers and multiplied with 5 different vectors, each of length 9216, also extracted from the STL-10 [6] dataset. Fig. 8b shows the average latency of the different approaches. Both LT coded approaches are almost $2 \times$ – faster than the MDS coded approaches in this setting and almost $3 \times$ – faster than the uncoded approaches. Each worker computes approximately 14 row-vector products at a time before communicating the results to the master. This corresponds to approximately 10% of the data stored in the workers memory thus balancing excessive communication (if results were communicated after every row-vector product computation) and memory limitations (if submatrices at the workers are too large to be communicated as a single chunk). Fig. 8e shows that the LT coded strategies also perform fewer total computations than MDS or 2-replication. Additional experiments in Appendix F show that LT coding also demonstrates greater resilience to node failures than Replication or MDS coding in this setting.

6.3 Serverless Computing Experiments

We also performed experiments in the serverless computing environment AWS Lambda [2]. Serverless computing eschews the master-worker set-up in favor of only workers (resources) which read data from storage, perform computations on the data, and write it back to storage. Any further computations (like decoding) can be performed as and when desired by re-reading data from storage. As described in [23], there is typically significant variability (straggling) across workers in this setting and therefore we expect to obtain speedups through coding. We use Numpywren [51] for performing linear algebra on AWS Lambda. We multiply a 100000×10000 matrix A (approximately $10 \times$ – larger than A in the previous two experiments) with a 10000×1 vector x in this setting. We compare the uncoded, MDS-Coded ($k = 80$) and LT-coded ($\alpha = 2.0$) approaches. As per the requirements of [51], the encoding is performed over blocks of 10 rows instead of individual rows. Our results, averaged over 5 trials, are presented in Figures 8c and 8f. They clearly show that the LT coded approach is faster than previous approaches, and performs fewer computations than MDS coding. We note that the experimental nature of current serverless computing frameworks makes it challenging to perform fine-grained logging of task times and to use larger encoded matrices (larger α). Thus we expect even better results once the limitations of current frameworks have been resolved.

7 CONCLUDING REMARKS

We propose an erasure coding strategy based on *rateless fountain codes* to speed up distributed matrix-vector multiplication in the presence of slow nodes (stragglers). For a matrix with m rows, our strategy requires the nodes to *collectively* finish slightly more than m row-vector products. Thus, it seamlessly adapts to varying node speeds and achieves near-perfect load balancing. Moreover, it has a small overhead of redundant computations (asymptotically zero), and low decoding complexity. Theoretical analysis and experiments show that our approach strikes a better latency-computation trade-off than existing uncoded, replication and maximum-distance-separable (MDS) coding approaches.

Going forward, we plan to extend our approach to other linear computations like sparse matrix-vector multiplication (SpMV), Matrix-Matrix multiplication, and Fourier Transforms. Previous work [61] has used fixed-rate variants of LT codes to speed-up SpMV; we expect even better performance by exploiting the rateless properties of fountain codes and utilizing partial work as described in this paper. Since erasure codes are inherently linear, extending coding techniques to speed-up distributed *non-linear* computations such as neural network inference is difficult. Recently [38, 39] propose the use of neural networks to learn the encoder and decoder to handle non-linear

computations. Coming up with a principled rateless coding approach in this setting remains an open problem.

8 ACKNOWLEDGEMENTS

The authors are grateful to Pulkit Grover, Sanghamitra Dutta, Yaoqing Yang, Haewon Jeong, Rashmi Vinayak, and Jack Kosaian for helpful discussions. Author Joshi also sincerely thanks Emina Soljanin, Alyson Fox, Fiona Knoll and Nadia Kazemi for fruitful initial discussions during the Women in Data Science and Mathematics (WiSDM) Research Collaboration Workshop held at Brown University in July 2017. This project was supported in part by the CMU Dean’s fellowship, Qualcomm Innovation Fellowship, NSF CCF grant no. 1850029 and an Amazon Credits for Research Grant.

REFERENCES

- [1] Amazon. 2006. Amazon Web Services EC2. <https://aws.amazon.com/ec2/>.
- [2] Amazon. 2014. Amazon Web Services Lambda. <https://aws.amazon.com/lambda/>.
- [3] William F Ames. 2014. *Numerical Methods for Partial Differential Equations*. Academic Press.
- [4] Ganesh Ananthanarayanan, Ali Ghodsi, Scott Shenker, and Ion Stoica. 2013. Effective Straggler Mitigation: Attack of the Clones.. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, Vol. 13. 185–198.
- [5] Ganesh Ananthanarayanan, Srikanth Kandula, Albert G Greenberg, Ion Stoica, Yi Lu, Bikas Saha, and Edward Harris. 2010. Reining in the Outliers in Map-Reduce Clusters using Mantri.. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, Vol. 10. 24.
- [6] Adam Coates, Andrew Ng, and Honglak Lee. 2011. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. 215–223.
- [7] William Dally. 2015. High-performance Hardware for Machine Learning. *NIPS Tutorial* (2015).
- [8] Dask Development Team. 2016. Dask: Library for dynamic task scheduling. <https://dask.org>
- [9] H. A. David and H. N. Nagaraja. 2003. *Order statistics*. John Wiley, Hoboken, NJ.
- [10] Jeffrey Dean and Luiz André Barroso. 2013. The tail at scale. *Commun. ACM* 56, 2 (2013), 74–80.
- [11] Jeffrey Dean and Sanjay Ghemawat. 2008. MapReduce: simplified data processing on large clusters. *Commun. ACM* 51, 1 (2008), 107–113.
- [12] James Dinan, D Brian Larkins, Ponnuswamy Sadayappan, Sriram Krishnamoorthy, and Jarek Nieplocha. 2009. Scalable work stealing. In *High Performance Computing Networking, Storage and Analysis, Proceedings of the Conference on*. IEEE, 1–11.
- [13] James Dinan, Stephen Olivier, Gerald Sabin, Jan Prins, P Sadayappan, and Chau-Wen Tseng. 2007. Dynamic load balancing of unbalanced computations using message passing. In *Parallel and Distributed Processing Symposium, 2007. IPDPS 2007. IEEE International*. IEEE, 1–8.
- [14] Sanghamitra Dutta, Viveck Cadambe, and Pulkit Grover. 2016. Short-dot: Computing large linear transforms distributively using coded short dot products. In *Advances In Neural Information Processing Systems*. 2100–2108.
- [15] Sanghamitra Dutta, Viveck Cadambe, and Pulkit Grover. 2017. Coded convolution for parallel and distributed computing within a deadline. In *IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2403–2407.
- [16] Sanghamitra Dutta, Mohammad Fahim, Farzin Haddadpour, Haewon Jeong, Viveck Cadambe, and Pulkit Grover. 2018. On the Optimal Recovery Threshold of Coded Matrix Multiplication. *arXiv preprint arXiv:1801.10292* (2018).
- [17] Python Software Foundation. [n. d.]. Multiprocessing. <https://docs.python.org/3/library/multiprocessing.html>.
- [18] Geoffrey C. Fox, Steve W. Otto, and Anthony JG. Hey. 1987. Matrix Algorithms on a Hypercube I: Matrix Multiplication. *Parallel Comput.* 4, 1 (1987), 17–31.
- [19] K. Gardner, M. Harchol-Balter, and A. Scheller-Wolf. 2016. A Better Model for Job Redundancy: Decoupling Server Slowdown and Job Size. In *Proceedings of IEEE MASCOTS*.
- [20] K. Gardner, S. Zbarsky, S. Doroudi, M. Harchol-Balter, E. Hyttiä, and A. Scheller-Wolf. 2015. Reducing Latency via Redundant Requests: Exact Analysis. In *Proceedings of the ACM SIGMETRICS*.
- [21] K. Gardner, S. Zbarsky, M. Harchol-Balter, and A. Scheller-Wolf. 2015. Analyzing Response Time in the Redundancy-d System. In *CMU-CS-15-141 archive*.
- [22] Google. 2015. Kubernetes. <https://kubernetes.io>.
- [23] Vipul Gupta, Shusen Wang, Thomas Courtade, and Kannan Ramchandran. 2018. Oversketch: Approximate matrix multiplication for the cloud. In *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 298–304.
- [24] Wael Halbawi, Navid Azizan-Ruhi, Fariborz Salehi, and Babak Hassibi. 2017. Improving distributed gradient descent using reed-solomon codes. *arXiv preprint arXiv:1706.05436* (2017).

- [25] Mor Harchol-Balter. 2013. *Performance modeling and design of computer systems: queuing theory in action*. Cambridge University Press.
- [26] Aaron Harlap, Henggang Cui, Wei Dai, Jinliang Wei, Gregory R Ganger, Phillip B Gibbons, Garth A Gibson, and Eric P Xing. 2016. Addressing the straggler problem for iterative convergent parallel ML. In *Proceedings of the Seventh ACM Symposium on Cloud Computing*. ACM, 98–111.
- [27] Kuang-Hua Huang et al. 1984. Algorithm-based Fault Tolerance for Matrix Operations. *IEEE Trans. Comput.* 100, 6 (1984), 518–528.
- [28] Longbo Huang, S. Pawar, Hao Zhang, and K. Ramchandran. 2012. Codes can reduce queueing delay in data centers. In *IEEE International Symposium on Information Theory Proceedings (ISIT)*. 2766–2770.
- [29] Gauri Joshi. 2017. Boosting Service Capacity via Adaptive Task Replication. *SIGMETRICS Perform. Eval. Rev.* 45, 2 (Oct. 2017), 9–11. <https://doi.org/10.1145/3152042.3152046>
- [30] Gauri Joshi. 2018. Synergy via Redundancy: Boosting Service Capacity with Adaptive Replication. *SIGMETRICS Performance Evaluation Review* 45, 3 (March 2018), 21–28. <http://doi.acm.org/10.1145/3199524.3199530>
- [31] Gauri Joshi, Yanpei Liu, and Emina Soljanin. 2012. Coding for fast content download. In *Allerton Conference on Communication, Control, and Computing*. IEEE, 326–333.
- [32] Gauri Joshi, Yanpei Liu, and Emina Soljanin. 2014. On the Delay-Storage Trade-Off in Content Download from Coded Distributed Storage Systems. *IEEE Journal on Selected Areas of Communications* 32, 5 (May 2014), 989–997.
- [33] Gauri Joshi, Joong Bum Rhim, John Sun, and Da Wang. 2010. Fountain codes. In *Global telecommunications conference (GLOBECOM 2010)*. 7–12.
- [34] Gauri Joshi, Emina Soljanin, and Gregory Wornell. 2015. Queues with redundancy: Latency-cost analysis. *ACM SIGMETRICS Performance Evaluation Review* 43, 2 (2015), 54–56.
- [35] Gauri Joshi, Emina Soljanin, and Gregory Wornell. 2017. Efficient Redundancy Techniques for Latency Reduction in Cloud Systems. *ACM Transactions on Modeling and Performance Evaluation of Computing Systems* 2, 12 (may 2017).
- [36] Can Karakus, Yifan Sun, and Suhas Diggavi. 2017. Encoded distributed optimization. In *IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2890–2894.
- [37] C. Kim and A. K. Agrawala. 1989. Analysis of the Fork-Join Queue. *IEEE Trans. Comput.* 38, 2 (Feb. 1989), 250–255.
- [38] Jack Kosaian, K. V. Rashmi, and Shivaram Venkataraman. 2018. Learning a Code: Machine Learning for Approximate Non-Linear Coded Computation. *CoRR abs/1806.01259* (2018). arXiv:1806.01259 <http://arxiv.org/abs/1806.01259>
- [39] Jack Kosaian, K. V. Rashmi, and Shivaram Venkataraman. 2019. Parity Models: A General Framework for Coding-Based Resilience in ML Inference. *CoRR abs/1905.00863* (2019). arXiv:1905.00863 <http://arxiv.org/abs/1905.00863>
- [40] Vipin Kumar, Ananth Grama, Anshul Gupta, and George Karypis. 1994. *Introduction to Parallel Computing: Design and Analysis of Algorithms*. Vol. 400. Benjamin/Cummings Redwood City.
- [41] Kangwook Lee, Maximilian Lam, Ramtin Pedarsani, Dimitris Papailiopoulos, and Kannan Ramchandran. 2017. Speeding Up Distributed Machine Learning Using Codes. *IEEE Transactions on Information Theory* (2017).
- [42] Kangwook Lee, Nihar B. Shah, Longbo Huang, and Kannan Ramchandran. 2017. The MDS Queue: Analysing the Latency Performance of Erasure Codes. *IEEE Transactions on Information Theory* 63, 5 (May 2017), 2822–2842.
- [43] Songze Li, Seyed Mohammadreza Mousavi Kalan, A Salman Avestimehr, and Mahdi Soltanolkotabi. 2017. Near-Optimal Straggler Mitigation for Distributed Gradient Methods. *arXiv preprint arXiv:1710.09990* (2017).
- [44] Songze Li, Mohammad Ali Maddah-Ali, and A Salman Avestimehr. 2016. A Unified Coding Framework for Distributed Computing with Stragglers. In *IEEE Global Communications Conference (GLOBECOM) Workshops*. IEEE, 1–6.
- [45] Michael Luby. 2002. LT codes. In *null*. IEEE, 271.
- [46] David JC MacKay. 2003. *Information theory, inference and learning algorithms*. Cambridge university press.
- [47] R. Nelson and A. Tantawi. 1988. Approximate Analysis of Fork/Join Synchronization in Parallel Queues. *IEEE Trans. Comput.* 37, 6 (Jun. 1988), 739–743.
- [48] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. 1999. *The PageRank Citation Ranking: Bringing order to the Web*. Technical Report. Stanford InfoLab.
- [49] Albin Severinson, Alexandre Graell i Amat, and Eirik Rosnes. 2017. Block-Diagonal and LT Codes for Distributed Computing With Stragglers. *arXiv preprint arXiv:1712.08230* (dec 2017).
- [50] Nihar B. Shah, Kangwook Lee, and Kannan Ramchandran. 2016. When Do Redundant Requests Reduce Latency? *IEEE Transactions on Communications* 64, 2 (Feb 2016), 715–722.
- [51] Vaishaal Shankar, Karl Krauth, Qifan Pu, Eric Jonas, Shivaram Venkataraman, Ion Stoica, Benjamin Recht, and Jonathan Ragan-Kelley. 2018. numpywren: serverless linear algebra. *arXiv preprint arXiv:1810.09679* (2018).
- [52] Amin Shokrollahi. 2006. Raptor codes. *IEEE/ACM Transactions on Networking (TON)* 14, SI (2006), 2551–2567.
- [53] Amin Shokrollahi, Michael Luby, et al. 2011. Raptor codes. *Foundations and trends® in communications and information theory* 6, 3–4 (2011), 213–322.
- [54] Yin Sun, Can Emre Koksak, and Ness B. Shroff. 2016. On Delay-Optimal Scheduling in Queueing Systems with Replications. *arXiv:1603.07322* (March 2016).

- [55] Yin Sun, Zizhan Zheng, Can Emre Koksal, Kyu-Han Kim, and Ness B. Shroff. 2015. Provably Delay Efficient Data Retrieving in Storage Clouds. In *Proceedings of the IEEE Conference on Computer Communications (INFOCOM)*.
- [56] Rashish Tandon, Qi Lei, Alexandros G Dimakis, and Nikos Karampatziakis. 2017. Gradient Coding: Avoiding Stragglers in Synchronous Gradient Descent. *stat* 1050 (2017), 8.
- [57] Elizabeth Varki, Arif Merchant, and Hui Chen. 2008. The M/M/1 fork-join queue with variable sub-tasks. *unpublished, available online* (2008).
- [58] Da Wang, Gauri Joshi, and Gregory Wornell. 2014. Efficient Task Replication for Fast Response times in Parallel Computation. In *ACM SIGMETRICS Performance Evaluation Review*, Vol. 42. ACM, 599–600.
- [59] Da Wang, Gauri Joshi, and Gregory Wornell. 2015. Using Straggler Replication to Reduce Latency in Large-scale Parallel Computing. *ACM SIGMETRICS Performance Evaluation Review* 43, 3 (2015), 7–11.
- [60] Da Wang, Gauri Joshi, and Gregory W. Wornell. 2019. Efficient Straggler Replication in Large-Scale Parallel Computing. *ACM Trans. Model. Perform. Eval. Comput. Syst.* 4, 2, Article 7 (April 2019), 23 pages. <http://doi.acm.org/10.1145/3310336>
- [61] Sinong Wang, Jiashang Liu, and Ness Shroff. 2018. Coded Sparse Matrix Multiplication. *arXiv preprint arXiv:1802.03430* (2018).
- [62] Yaoqing Yang, Pulkit Grover, and Soumya Kar. 2017. Coded Distributed Computing for Inverse Problems. In *Advances in Neural Information Processing Systems*. 709–719.
- [63] Qian Yu, Mohammad Maddah-Ali, and Salman Avestimehr. 2017. Polynomial codes: an optimal design for high-dimensional coded matrix multiplication. In *Advances in Neural Information Processing Systems*. 4406–4416.
- [64] Qian Yu, Mohammad Ali Maddah-Ali, and A Salman Avestimehr. 2017. Coded Fourier Transform. *arXiv preprint arXiv:1710.06471* (2017).
- [65] Matei Zaharia, Mosharaf Chowdhury, Michael J Franklin, Scott Shenker, and Ion Stoica. 2010. Spark: Cluster computing with working sets. *HotCloud* 10, 10-10 (2010), 95.

A PROPERTIES OF LT CODES

Fig. 9 shows simulation results for the number of symbols decoded successfully for each encoded symbol received. For this we perform LT-Coded multiplication of a randomly generated $10,000 \times 10,000$ matrix with a $10,000 \times 1$ vector. The matrix \mathbf{A} is encoded using an LT code with parameters c and δ chosen according to the guidelines of [46]. We generate a single row of the encoded matrix \mathbf{A}_e at a time which is then multiplied with the $10,000 \times 1$ size vector \mathbf{x} to give a single element of the encoded matrix vector product \mathbf{b}_e . The process is repeated until we have enough symbols for successfully decoding the entire $10,000 \times 1$ size vector \mathbf{b} using the peeling decoder. The plots of Fig. 9 correspond to different choices of c and δ . In each case we observe an avalanche behavior wherein very few symbols are decoded up to a point (approximately up to 10,000 encoded symbols received) after which the decoding proceeds very rapidly to completion. This effectively illustrates the fact that the computation overhead of the proposed LT coded matrix vector multiplication strategy is very small ($m_d = m(1 + \epsilon)$). The theoretical encoding and decoding properties of LT codes are summarized in the following lemmas:

LEMMA 7 (THEOREM 13 IN [45]). *For any constant $\delta > 0$, the average degree of an encoded symbol is $O(\log(m/\delta))$ where m is the number of source symbols.*

COROLLARY 5. *Each encoding symbol can be generated using $O(\log m)$ symbol operations on average.*

LEMMA 8 (THEOREM 17 IN [45]). *For any constant $\delta > 0$ and for a source block with m source symbols, the LT decoder can recover all the source symbols from a set of $M' = m + O(\sqrt{m} \log^2(m/\delta))$ with probability at least $1 - \delta$.*

COROLLARY 6. *The expected decoding threshold $\mathbb{E}[M']$ is given by $\mathbb{E}[M'] = m(1 + \epsilon)$ where $\epsilon \rightarrow 0$ as $m \rightarrow \infty$*

COROLLARY 7. *Since the average degree of an encoded symbol is $O(\log(m/\delta))$ the decoding requires $O(m \log m)$ symbol operations on average.*

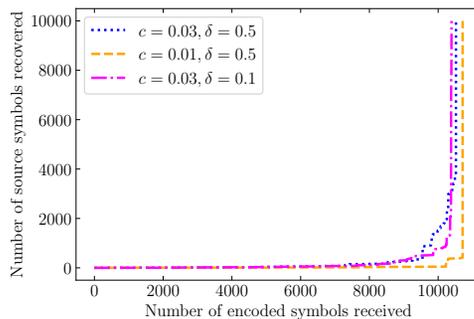


Fig. 9. The number of decoded symbols is almost constant until $m = 10,000$ encoded symbols are received after which it increases rapidly.

B ON THE ORDER STATISTICS OF EXPONENTIAL RANDOM VARIABLES

We first state some standard results [9] on order statistics of exponential random variables to aid the understanding of the latency analysis presented subsequently. If X_1, X_2, \dots, X_p are exponential random variables with rate μ , their k^{th} order statistic is denoted by $X_{k:p}$. Thus, $X_{1:p} =$

$\min(X_1, X_2, \dots, X_p)$, and $X_{p:p} = \max(X_1, X_2, \dots, X_p)$. The expected value of $X_{k:p}$ is given by

$$\mathbb{E}[X_{k:p}] = \frac{1}{\mu} \left(\frac{1}{p} + \dots + \frac{1}{p-k+1} \right) = \frac{H_p - H_{p-k}}{\mu}, \quad (23)$$

where H_p is the p^{th} Harmonic number

$$H_p \triangleq \begin{cases} \sum_{i=1}^p \frac{1}{i} & \text{for } p = 1, 2, \dots \\ 0 & \text{for } p = 0 \end{cases} \quad (24)$$

For large p , $H_p = \log p + \gamma$, where γ is the Euler-Mascheroni constant and thus we can use the approximation $H_p \approx \log p$ for large p .

Also the difference of consecutive order statistics of i.i.d exponential random variables is also exponentially distributed. Specifically, $U_l = X_{l+1:p} - X_{l:p} \sim \exp((p-l)\mu)$ in this case.

C PROOF OF DELAY ANALYSIS RESULTS

C.1 Ideal Load Balancing Strategy

PROOF OF THEOREM 2. The result follows from Lemma 9 and Lemma 10 for non-redundant and redundant task allocation policies respectively given below. Redundant policies refer to policies where the same task can be allocated to multiple workers (such as the r -Replication policy). \square

LEMMA 9. *For any distributed matrix-vector multiplication scheme following the delay model of (5) without any redundancy in task allocation, the latency T is no less than T_{ideal} .*

PROOF OF LEMMA 9. Consider any scheme other than the ideal load balancing scheme. Let $W_{i,i+1}$ be the time elapsed between completion of the i^{th} and $(i+1)^{\text{th}}$ computation in this scheme and let $W_{i,i+1}^{\text{ideal}}$ be the time elapsed between completion of the i^{th} and $(i+1)^{\text{th}}$ computation in the ideal scheme. By definition, the ideal load balancing is work-conserving, that is, no worker is idle while there are pending computations at other workers. Hence for the other scheme there must be some computation \hat{i} after which at least one worker is idle even though there are pending computations at the other workers. Therefore $W_{i,\hat{i}+1}^{\text{ideal}} \leq W_{i,\hat{i}+1}$. Moreover since the tasks are allocated by the master initially in our setting, it means that if a worker is idle after computation \hat{i} is completed by the system, then it is idle for *all* subsequent computations. Thus, $W_{i,\hat{i}+1}^{\text{ideal}} \leq W_{i,i+1}$, $\forall i \geq \hat{i}$. Therefore if T^{ideal} and T^{other} are the times taken by the ideal and any other scheme respectively to complete m tasks,

$$T_{\text{ideal}} = W_{0,1}^{\text{ideal}} + \dots + W_{\hat{i},\hat{i}+1}^{\text{ideal}} + \dots + W_{m-1,m}^{\text{ideal}}, \quad (25)$$

$$T = W_{0,1} + \dots + W_{\hat{i},\hat{i}+1} + \dots + W_{m-1,m}, \quad (26)$$

and $T^{\text{ideal}} \leq T$ for any other task allocation strategy. \square

LEMMA 10. *For any distributed matrix-vector multiplication scheme following the delay model of (5) with potentially redundant task allocations, the latency T is no less than T_{ideal} .*

PROOF OF LEMMA 10. Consider the a scheme with redundancy where task j is allocated to r distinct workers. Let V_1, \dots, V_r be the time instant at which each of the r workers start working on the task. If W^j is the earliest time instant at which the task is received at the master then,

$$W^j = \min_{1 \leq i \leq r} (V_i + \tau), \quad (27)$$

$$= \tau + \min_{1 \leq i \leq r} V_i. \quad (28)$$

Thus the time at which a task is completed is equal to the sum of the earliest time instant ($\min_{1 \leq i \leq r} V_i$) at which one of the r worker is available to process the task and the time (τ) required by any worker to process the task. Since the ideal load balancing scheme greedily assigns tasks to the first of the p workers that are available to process it and $r \leq p$, therefore the ideal scheme is essentially equivalent to a scheme with *maximum* redundancy and thus cannot be outperformed by any task allocation strategy with redundancy. \square

PROOF OF LEMMA 2. As per our model, the time taken by worker i to perform B_i computations is given by

$$Y_i = X_i + \tau B_i, \text{ for } i = 1, \dots, p. \quad (29)$$

The latency T_{ideal} is the earliest time when $\sum_{i=1}^p B_i = m$, as illustrated in Fig. 10a. We note that, in this case it is not necessary that each worker has completed at least 1 computation. Specifically, if $T_{\text{ideal}} - X_i \leq \tau$ for any i then it means that worker i has not performed even a single computation in the time that the system as a whole has completed m computations (owing to the large initial delay X_i). Therefore we define

$$\mathcal{W}_{\text{ideal}} := \{i : T_{\text{ideal}} - X_i \geq \tau\}. \quad (30)$$

Here $\mathcal{W}_{\text{ideal}}$ is the set of workers for which $B_i > 0$. Thus

$$T_{\text{ideal}} = \max_{i \in \mathcal{W}_{\text{ideal}}} Y_i = \max_{i \in \mathcal{W}_{\text{ideal}}} (X_i + \tau B_i), \quad (31)$$

$$\geq \min_{i \in \{1, \dots, p\}} X_i + \tau \max_{i \in \mathcal{W}_{\text{ideal}}} B_i, \quad (32)$$

$$\geq X_{1:p} + \frac{\tau m}{p}, \quad (33)$$

where to obtain (32), we replace each X_i in (31) by $\min_{i \in [1, \dots, p]} X_i$ and then we can bring it outside the maximum. To obtain (33), we observe that in order for the p workers to collectively finish m computations, the maximum number of computations completed by a worker has to be at least m/p .

To derive the upper bound, we note that

$$T_{\text{ideal}} \leq X_i + \tau(B_i + 1), \text{ for all } i = 1, \dots, p \quad (34)$$

This is because at time T_{ideal} each of the workers $1, \dots, p$, have completed B_1, \dots, B_p row-vector product tasks respectively, but they may have partially completed the next task. The 1 added to each B_i accounts for this edge effect, which is also illustrated in Fig. 6. Summing over all i on both sides, we get

$$\sum_{i=1}^p T_{\text{ideal}} \leq \sum_{i=1}^p X_i + \sum_{i=1}^p \tau (B_i + 1), \quad (35)$$

$$pT_{\text{ideal}} \leq \sum_{i=1}^p X_i + \tau (m + p), \quad (36)$$

$$T_{\text{ideal}} \leq \frac{1}{p} \sum_{i=1}^p X_i + \frac{\tau m}{p} + \tau \quad (37)$$

\square

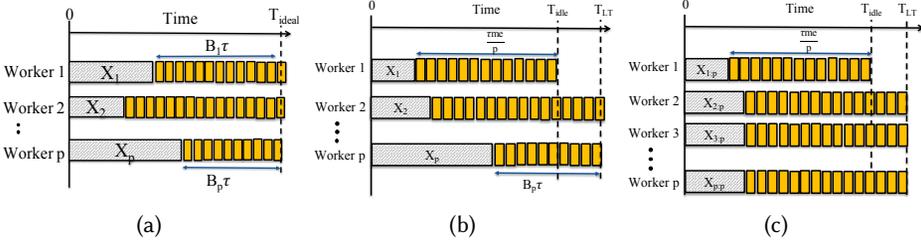


Fig. 10. (a) Worker i has a random exponential initial delay X_i , after which it completes row-vector product tasks taking time τ per task. In the ideal case the latency T_{ideal} is simply the time to complete m tasks in total. (b) A general scenario Worker 1 runs out of computations at T_{idle} before the system completes enough computations necessary for successful decoding. (c) The specific scenario (event E_2) where Worker 2 starts so late that Worker 1 runs out of computations even though workers $3, \dots, p$ start at the same time as 2.

PROOF OF COROLLARY 1. If $X_i \sim \exp(\mu)$ then taking expectation on both sides of (33) gives

$$\mathbb{E}[T_{\text{ideal}}] \geq \mathbb{E}[X_{1:p}] + \frac{\tau m}{p}, \quad (38)$$

$$= \frac{1}{p\mu} + \frac{\tau m}{p}. \quad (39)$$

where the lower bound in (39) follows from the result (23) on order statistics of exponential random variables. Likewise for the upper bound we can compute expectation on both sides of (37) to get,

$$\mathbb{E}[T_{\text{ideal}}] \leq \frac{1}{p} \sum_{i=1}^p \mathbb{E}[X_i] + \frac{\tau m}{p} + \tau \quad (40)$$

$$\mathbb{E}[T_{\text{ideal}}] \leq \frac{1}{\mu} + \frac{\tau m}{p} + \tau. \quad (41)$$

□

C.2 Rateless Coded Strategy

PROOF OF THEOREM 3. Recall that we make the following assumption for analysing the latency of the proposed rateless coded strategy.

ASSUMPTION 2. *The decoding threshold M' (Definition 3) of the LT coded strategy satisfies $M' \simeq m$.*

We believe the above assumption is reasonable because the problem of distributed matrix vector multiplication arises only when m (the number of rows of \mathbf{A}) is large and the high probability bound of Lemma 1 can be used to show that $\mathbb{E}[M'] = m(1 + \epsilon)$, where $\epsilon \rightarrow 0$ as $m \rightarrow \infty$. To analyze the probability $\Pr(T_{\text{LT}} > T_{\text{ideal}})$ in light of the above assumption, let us first understand when T_{LT} will exceed T_{ideal} . This situation will arise when the fastest node runs out of computations and becomes idle before M' computations are collectively collectively completed by all the workers as illustrated in Fig. 10b.

Recall that according to our delay model, the time Y_i required by worker i to perform B_i computations is

$$Y_i = X_i + \tau B_i, \quad \text{for all } i = 1, \dots, p \quad (42)$$

where X_i is an initial delay at worker i . Without loss of generality, we assume that nodes are ordered in increasing order of initial delays. Thus, $X_i = X_{i:p}$ for $i = 1, \dots, p$, where $X_{i:p}$ denotes the i^{th} order statistic of p i.i.d. random variables with distribution F_X .

Now let us determine the events that cause T_{LT} to be greater than T_{Ideal} for a given set of realizations of the initial delays. Since each node is assigned m_e/p row-vector product tasks, the time T_{idle} at which the fastest node, node 1, becomes idle is

$$T_{\text{idle}} = X_{1:p} + \frac{\tau m_e}{p} \quad (43)$$

Suppose that the initial delays of the fastest j workers are $X_{1:p}, X_{2:p}, \dots, X_{j:p}$ respectively. Then let C_j denote the number of computations collectively performed by the p workers until time T_{idle} in the event that all the remaining $p - j$ workers start at the earliest possible instant and have initial delays equal to $X_{j:p}$.

Observe that $T_{\text{LT}} > T_{\text{Ideal}}$ if node 1 becomes idle before all workers collectively perform m computations. This situation (depicted in Fig. 10b) arises if $X_{j:p}$ is so large for some node j that $C_j < m$. We distill this event of node 1 becoming idle before the result is recovered into sub-events E_1, E_2, \dots, E_p , where we refer to E_j as the event of idling due to node j , which is defined as follows.

DEFINITION 4 (EVENT E_j : IDLING DUE TO NODE j). *Given $X_{1:p}, X_{2:p}, \dots, X_{j-1:p}, E_j$ occurs if $C_j < m$.*

The event E_2 is depicted in Fig. 10c. Therefore,

$$\Pr(T_{\text{LT}} > T_{\text{ideal}}) = \Pr(E_2 \cup \dots \cup E_p) \quad (44)$$

$$\leq \sum_{j=2}^p \Pr(E_j) \quad (45)$$

$$= \sum_{j=2}^p \Pr(C_j < m) \quad (46)$$

where (45) follows from the union bound. Note that the event of idling due to node 1, E_1 is not defined since node 1 is the fastest node by definition.

Now we derive a lower bound on C_j in terms of $X_{1:p}, \dots, X_{j:p}$ as follows.

$$C_j = \sum_{l=1}^{j-1} \left(\left\lfloor \frac{T_{\text{idle}} - X_{l:p}}{\tau} \right\rfloor \right)^+ + (p - j + 1) \left(\left\lfloor \frac{T_{\text{idle}} - X_{j:p}}{\tau} \right\rfloor \right)^+ \quad (47)$$

$$= \sum_{l=1}^{j-1} \left(\left\lfloor \frac{m_e}{p} - \frac{(X_{l:p} - X_{i:p})}{\tau} \right\rfloor \right)^+ + (p - j + 1) \left(\left\lfloor \frac{m_e}{p} - \frac{(X_{j:p} - X_{1:p})}{\tau} \right\rfloor \right)^+ \quad (48)$$

$$\geq m_e - \sum_{l=1}^{j-1} \left(\frac{X_{l:p} - X_{i:p}}{\tau} \right) - (p - j + 1) \left(\frac{X_{j:p} - X_{1:p}}{\tau} \right) \quad (49)$$

$$\geq m_e - \frac{p-1}{\tau} \sum_{l=1}^{j-1} (X_{l+1:p} - X_{l:p}) \quad (50)$$

$$= m_e - \frac{p-1}{\tau} \sum_{l=1}^{j-1} U_l \quad (51)$$

where in (49) we use the fact that for any $a, b \geq 0$, $[a - b]_+ \geq (a - b)$.

As a consequence of the stochastic dominance implied by the above bound we have.

$$\Pr(E_j) = \Pr(C_j < M') \quad (52)$$

$$\leq \Pr\left(\sum_{l=1}^{j-1} U_{l;p} \geq \tau \frac{m_e - m}{p-1}\right) \quad (53)$$

Recall from (45) that,

$$\Pr(T_{LT} > T_{ideal}) \leq \sum_{j=2}^p \Pr(E_j) \quad (54)$$

$$\leq \sum_{j=2}^p \Pr\left(\sum_{l=1}^{j-1} U_l \geq \frac{\tau m(\alpha-1)}{p-1}\right) \quad (55)$$

□

PROOF OF COROLLARY 2. If $X_i \sim \exp(\mu)$ we can use the results from Appendix B to simplify (53) further,

$$\Pr(E_j) \leq \Pr\left(\sum_{l=1}^{j-1} U_l \geq \frac{\tau m(\alpha-1)}{p-1}\right) \quad (56)$$

$$\leq \Pr\left((j-1)U_{j-1} \geq \frac{\tau m(\alpha-1)}{p-1}\right) \quad (57)$$

$$= \Pr\left(U_{j-1} \geq \frac{\tau m(\alpha-1)}{(p-1)(j-1)}\right) \quad (58)$$

$$= \exp\left(-\frac{\mu \tau m(\alpha-1)(p-j+1)}{(p-1)(j-1)}\right) \quad (59)$$

where (57) is obtained from the fact that $\Pr(U_{j-1} \geq u) \geq \Pr(U_l \geq u)$ for $l = 1, \dots, j-1$ for any u since $U_l \sim \exp((p-l)\mu)$. Lastly (59) is obtained from the expression for the tail distribution of an exponential random variable.¹

Since $(p-j+1)/(j-1) > 1/(p-1) \forall j = 2, \dots, p$, we obtain the final result

$$\Pr(T_{LT} > T_{ideal}) \leq \sum_{j=2}^p \Pr(E_j) \quad (60)$$

$$\leq \sum_{j=2}^p \Pr\left(\sum_{l=1}^{j-1} U_l \geq \frac{\tau m(\alpha-1)}{p-1}\right) \quad (61)$$

$$\leq (p-1) \exp\left(-\frac{\mu \tau m(\alpha-1)}{(p-1)^2}\right) \quad (62)$$

$$\leq p \exp\left(-\frac{\mu \tau m(\alpha-1)}{p^2}\right) \quad (63)$$

□

¹For a two-server system ($p = 2$), $\Pr(T_{LT} > T_{ideal}) = \Pr(E_2) = \Pr(C_2 < m) \leq \Pr(X_{2;2} - X_{1;2} \geq \tau(m_e - m))$ and we can see that increasing redundancy (increasing m) will decrease the probabilistic upper bound. Actually for any $m_e > 2m$, $\Pr(T_{LT} > T_{ideal}) = 0$ since in this case the entire work can be done by the faster node but in this case our upper bound, which is a little loose, has a small positive value.

PROOF OF THEOREM 4. We can write the expectation of T_{LT} as follows.

$$\mathbb{E}[T_{LT}] = \Pr(T_{LT} = T_{ideal})\mathbb{E}[T_{LT} \mid T_{LT} = T_{ideal}] + \Pr(T_{LT} > T_{ideal})\mathbb{E}[T_{LT} \mid T_{LT} > T_{ideal}] \quad (64)$$

$$\begin{aligned} &= \Pr(T_{LT} = T_{ideal})\mathbb{E}[T_{ideal} \mid T_{LT} = T_{ideal}] + \Pr(T_{LT} > T_{ideal})\mathbb{E}[T_{ideal} \mid T_{LT} > T_{ideal}] \\ &+ \Pr(T_{LT} > T_{ideal})\mathbb{E}[T_{LT} - T_{ideal} \mid T_{LT} > T_{ideal}] \end{aligned} \quad (65)$$

$$= \mathbb{E}[T_{ideal}] + \Pr(T_{LT} > T_{ideal})\mathbb{E}[T_{LT} - T_{ideal} \mid T_{LT} > T_{ideal}] \quad (66)$$

$$\leq \mathbb{E}[T_{ideal}] + \Pr(T_{LT} > T_{ideal})\mathbb{E}[T_{LT} \mid T_{LT} > T_{ideal}], \quad (67)$$

where (67) follows from the fact that $T_{LT} - T_{ideal} < T_{LT}$. Since we have already derived an upper bound for $\Pr(T_{LT} > T_{ideal})$ in Theorem 3, we will now derive an upper bound for $\mathbb{E}[T_{LT} \mid T_{LT} > T_{ideal}]$. To do so, we first define \mathcal{W}_{LT} as the set of workers that have *not* completed all the $\alpha m/p$ computations assigned to them i.e. $\mathcal{W}_{LT} := \{i : B_i < \frac{\alpha m}{p}\}$. The set \mathcal{W}_{LT} does not include the workers which have completed all the $\alpha m/p$ tasks assigned to them and are idle at T_{LT} since they have completed their tasks at some (unknown earlier) time and thus do not increase the upper bound for T_{LT} . For workers in \mathcal{W}_{LT} ,

$$T_{LT} \leq X_i + \tau(B_i + 1) \text{ for all } i \in \mathcal{W}_{LT}, \quad (68)$$

$$\leq \sum_{i=1}^p X_i + \frac{\tau \alpha m}{p} + \tau, \quad (69)$$

$$= X_{1:p} + X_{2:p} + \dots + X_{p:p} + \frac{\tau \alpha m}{p} + \tau \quad (70)$$

$$= pX_{1:p} + (p-1)U_1 + \dots + U_{p-1} + \frac{\tau \alpha m}{p} + \tau, \quad (71)$$

where (69) follows from the fact that $X_i \leq \sum_{i=1}^p X_i$ for any worker i . In (70) we express the sum of X_i 's in terms of the order statistics $X_{l:p}$. In (71), $U_l = X_{l+1:p} - X_{l:p}$, $l = 1, \dots, p-1$.

We can then use (71) to upper bound $\mathbb{E}[T_{LT} \mid T_{LT} > T_{ideal}]$ as

$$\mathbb{E}[T_{LT} \mid T_{LT} > T_{ideal}] \leq p\mathbb{E}[X_{1:p}] + \mathbb{E}\left[\sum_{l=1}^{p-1} (p-l)U_l \mid T_{LT} > T_{ideal}\right] + \frac{\tau \alpha m}{p} + \tau, \quad (72)$$

where in the first term of (72), we do not condition by $T_{LT} > T_{ideal}$ since the initial delay $X_{1:p}$ at the fastest worker $X_{1:p}$ is independent of the event $T_{LT} > T_{ideal}$. This is because as seen in the proof of Theorem 3, the event $T_{LT} > T_{ideal}$ depends on the values of U_1, \dots, U_{p-1} but not on $X_{1:p}$.

Now let us define $\bar{U}_{l-1} = [U_1, \dots, U_{l-1}]$ and $S_j = \sum_{l=1}^j (p-l)U_l$ for all $j = 1, \dots, p-1$. Using Lemma 11 (stated and proved below) we can simplify the second term in the right-hand-side of (72) as,

$$\mathbb{E}_{\bar{U}_{p-1}}[S_j \mid T_{LT} > T_{ideal}] = \mathbb{E}_{\bar{U}_{p-2}}[\mathbb{E}_{U_{p-1}}[U_{p-1} + S_{p-2} \mid T_{LT} > T_{ideal}, \bar{U}_{p-2}]] \quad (73)$$

$$\leq \mathbb{E}_{\bar{U}_{p-3}}[R_{p-1} + \mathbb{E}_{U_{p-2}}[2U_{p-2} + S_{p-3} \mid T_{LT} > T_{ideal}, \bar{U}_{p-3}]] \quad (74)$$

where $R_l = \mathbb{E}[(p-l)U_l \mid U_l > \tau \alpha m/p]$.

Repeatedly applying Lemma 11 gives the final upper bound

$$\mathbb{E}\left[\sum_{l=1}^{p-1} (p-l)U_l \mid T_{LT} > T_{ideal}\right] \leq \sum_{l=1}^{p-1} R_l \quad (75)$$

Substituting (75) in (72) gives

$$\mathbb{E}[T_{LT} \mid T_{LT} > T_{ideal}] \leq p\mathbb{E}[X_{1:p}] + \sum_{l=1}^{p-1} R_l + \frac{\tau\alpha m}{p} + \tau \quad (76)$$

Thus finally we have

$$\mathbb{E}[T_{LT}] - \mathbb{E}[T_{ideal}] \leq \Pr(T_{LT} > T_{ideal})\mathbb{E}[T_{LT} \mid T_{LT} > T_{ideal}] \quad (77)$$

$$\leq \left(\sum_{j=2}^p \Pr \left(\sum_{l=1}^{j-1} U_{l:p} \geq \tau \frac{m_e - M'}{p-1} \right) \right) \left(p\mathbb{E}[X_{1:p}] + \sum_{l=1}^{p-1} R_l + \frac{\tau\alpha m}{p} + \tau \right). \quad (78)$$

If $X_i \sim \exp(\mu)$ we can use the results from Appendix B to simplify the above results further by using the bounds on $\Pr(T_{LT} > T_{ideal})$ for this case which is given by

$$\Pr(T_{LT} > T_{ideal}) \leq p \exp \left(-\frac{\mu\tau m(\alpha-1)}{p^2} \right), \quad (79)$$

and the fact that $\mathbb{E}[X_{1:p}] = 1/p\mu$ and

$$\begin{aligned} R_l &= \mathbb{E}[(p-l)U_l \mid U_l > \frac{\tau\alpha m}{p}] \\ &= \frac{\tau\alpha m}{p}(p-l) + (p-l)\frac{1}{(p-l)\mu} \\ &= \frac{\tau\alpha m}{p}(p-l) + \frac{1}{\mu}, \end{aligned} \quad (80)$$

since $U_l \sim \exp((p-l)\mu)$ which implies that

$$p\mathbb{E}[X_{1:p}] + \sum_{l=1}^{p-1} R_l = p\frac{1}{p\mu} + \frac{\tau\alpha m}{p} \frac{p(p-1)}{2} + \frac{p-1}{\mu}, \quad (81)$$

$$\leq \frac{\tau\alpha mp}{2} + \frac{p}{\mu}. \quad (82)$$

Thus,

$$\mathbb{E}[T_{LT}] - \mathbb{E}[T_{ideal}] \leq \Pr(T_{LT} > T_{ideal})\mathbb{E}[T_{LT} \mid T_{LT} > T_{ideal}] \quad (83)$$

$$\leq p \exp \left(-\frac{\mu\tau m(\alpha-1)}{p^2} \right) \left(\frac{\tau\alpha mp}{2} + \frac{p}{\mu} + \frac{\tau\alpha m}{p} + \tau \right) \quad (84)$$

$$\leq \left(\tau\alpha mp^2 + \frac{p^2}{\mu} + \tau p \right) \exp \left(-\frac{\mu\tau m(\alpha-1)}{p^2} \right) \quad (85)$$

where the last expression is true for $p \geq 2$ which is always the case in distributed settings. \square

LEMMA 11. *Given the values of U_1, \dots, U_{l-1} we have the following upper bound*

$$\mathbb{E}[(p-l)U_l \mid T_{LT} > T_{ideal}, \bar{U}_{l-1}] \leq \mathbb{E}[(p-l)U_l \mid U_l > \frac{\tau\alpha m}{p}] \quad (86)$$

where $\bar{U}_{l-1} = [U_1, \dots, U_{l-1}]$.

PROOF. Observe that given the values of U_1, \dots, U_{l-1} , the condition $U_l > \tau\alpha m/p$ is always sufficient to guarantee that $T_{LT} > T_{ideal}$ even if that may not be guaranteed by the values of $U_{1:p}, \dots, U_{l-1:p}$. This is because $U_l = X_{l+1:p} - X_{l:p}$ and thus $U_l > \tau\alpha m/p$ implies that worker $l+1$ starts computing only *after* worker l has completed all $\alpha m/p$ computations assigned to it. For eg. If $U_2 > \tau\alpha m/p$ then the second fastest worker (and all subsequent workers) will only start computing

after the fastest worker has completed all of its $\alpha m/p$ computations. Based on this we can show that,

$$\begin{aligned} & \mathbb{E}[(p-l)U_l \mid T_{LT} > T_{ideal}, \bar{U}_{l-1}] \\ &= \Pr(U_l < \frac{\tau \alpha m}{p}) \mathbb{E}[(p-l)U_l \mid T_{LT} > T_{ideal}, U_l < \frac{\tau \alpha m}{p}, \bar{U}_{l-1}] + \\ & \Pr(U_l > \frac{\tau \alpha m}{p}) \mathbb{E}[(p-l)U_l \mid T_{LT} > T_{ideal}, U_l > \frac{\tau \alpha m}{p}, \bar{U}_{l-1}] \end{aligned} \quad (87)$$

Since $\mathbb{E}[(p-l)U_l \mid T_{LT} > T_{ideal}, U_l < \frac{\tau \alpha m}{p}, \bar{U}_{l-1}] < \mathbb{E}[(p-l)U_l \mid T_{LT} > T_{ideal}, U_l > \frac{\tau \alpha m}{p}, \bar{U}_{l-1}]$, it follows that

$$\begin{aligned} & \mathbb{E}[(p-l)U_l \mid T_{LT} > T_{ideal}, \bar{U}_{l-1}] \\ & \leq \mathbb{E}[(p-l)U_l \mid T_{LT} > T_{ideal}, U_l > \frac{\tau \alpha m}{p}, \bar{U}_{l-1}] \end{aligned} \quad (88)$$

However since $U_l > \frac{\tau \alpha m}{p}$ is sufficient to guarantee $T_{LT} > T_{ideal}$ and $(p-l)U_l$ does not depend on $\bar{U}_{l-1} = [U_1, \dots, U_{l-1}]$, the above expression reduces to

$$\mathbb{E}[(p-l)U_l \mid T_{LT} > T_{ideal}, \bar{U}_{l-1}] \leq \mathbb{E}[(p-l)U_l \mid U_l > \frac{\tau \alpha m}{p}] \quad (89)$$

□

C.3 MDS Coded Strategy

PROOF OF LEMMA 3. The latency in the MDS-coded case is $T_{MDS} = Y_{k:p}$, where $Y_{k:p}$ is the k^{th} order statistic of the individual worker latencies Y_1, Y_2, \dots, Y_p since we only wait for the fastest k workers to finish the task assigned to them. In this case, each of the fastest k workers performs $\frac{m}{k}$ computations and thus the overall latency is given by

$$T_{MDS} = Y_{k:p} = X_{k:p} + \tau \frac{m}{k} \quad (90)$$

□

PROOF OF COROLLARY 3. If $X_i \sim \exp(\mu)$, the expected overall latency is given by

$$\mathbb{E}[T_{MDS}] = \mathbb{E}[X_{k:p}] + \tau \frac{m}{k}, \quad (91)$$

$$= \frac{\tau m}{k} + \frac{1}{\mu} (H_p - H_{p-k}), \quad (92)$$

$$\simeq \frac{\tau m}{k} + \frac{1}{\mu} \log \frac{p}{p-k}. \quad (93)$$

where (92) and (93) follow from the exponential order statistics results in (23) and (24). □

PROOF OF LEMMA 4. As per our model, we represent the number of computations at worker i by the random variable B_i . We also use the random variable C_{MDS} to denote the total number of computations performed by all p workers until T_{MDS} , which is the time when the master collects enough computations to be able to recover the matrix-vector product $\mathbf{b} = \mathbf{A}\mathbf{x}$. Thus

$$C_{MDS} = B_1 + B_2 + \dots + B_p \quad (94)$$

$$= B_{1:p} + B_{2:p} + \dots + B_{p:p}, \quad (95)$$

where the second expression is simply the right-hand side of the first expression written in terms of the corresponding order statistics. We note that under our model the time spent by worker i in performing B_i computations is $Y_i = X_i + \tau B_i$ where X_i denotes setup/initial delay and τ is a

constant denoting the time taken to perform a single computation. Thus $B_{1:p}$ corresponds to the worker that performs the least number of computations which is also the worker with the largest value of setup time i.e $X_{p:p}$ since all workers stop computing at the same time (T_{MDS}). Thus for a given C_0 , the tail of the total number of computations performed in the MDS Coded strategy is given by

$$\Pr\left(C_{\text{MDS}} \leq \frac{mp}{k} - C_0\right) = \Pr\left(\sum_{i=1}^p B_{i:p} \leq \frac{mp}{k} - C_0\right) \quad (96)$$

$$= \Pr\left(\sum_{i=1}^{p-k} B_{i:p} + \frac{m}{k} \times k \leq \frac{mp}{k} - C_0\right) \quad (97)$$

$$= \Pr\left(\sum_{i=1}^{p-k} B_{i:p} \leq \frac{m(p-k)}{k} - C_0\right) \quad (98)$$

$$\leq \Pr\left((p-k) B_{1:p} \leq \frac{m(p-k)}{k} - C_0\right) \quad (99)$$

$$= \Pr\left(B_{1:p} \leq \frac{m}{k} - \frac{C_0}{p-k}\right) \quad (100)$$

where (97) follows from the fact that the fastest k workers correspond to $B_{p-k+1:p}, B_{p-k+2:p}, \dots, B_{p:p}$ and must perform all the tasks assigned to them i.e. m/k computations each, while (99) follows from the fact that $B_{2:p}, \dots, B_{p:p}$ are always larger than $B_{1:p}$ by definition.

At this point we note that the worker which performs $B_{1:p}$ computations has setup time $X_{p:p}$. There can be two possibilities – either $T_{\text{MDS}} > X_{p:p}$, or $T_{\text{MDS}} \leq X_{p:p}$. If $T_{\text{MDS}} > X_{p:p}$ then

$$T_{\text{MDS}} \leq X_{p:p} + \tau (B_{1:p} + 1) \quad (101)$$

where the added 1 accounts for the edge effect of partial computations at the nodes. If $T_{\text{MDS}} \leq X_{p:p}$ then also the upper bound (101) holds. Thus overall (by rearranging terms in (101)) we obtain,

$$B_{1:p} \geq \frac{T_{\text{MDS}} - X_{p:p}}{\tau} - 1. \quad (102)$$

Thus we can write

$$\Pr\left(C_{\text{MDS}} \leq \frac{mp}{k} - C_0\right) \leq \Pr\left(\frac{T_{\text{MDS}} - X_{p:p}}{\tau} - 1 \leq \frac{m}{k} - \frac{C_0}{p-k}\right) \quad (103)$$

$$= \Pr\left(X_{p:p} - X_{k:p} \geq \frac{\tau C_0}{p-k} - \tau\right) \quad (104)$$

$$= \Pr\left(\sum_{l=k}^{p-1} (X_{l+1:p} - X_{l:p}) \geq \frac{\tau C_0}{p-k} - \tau\right), \quad (105)$$

where (104) follows from the fact that $T_{\text{MDS}} = X_{k:p} + \tau m/k$. If $X_i \sim \exp(\mu)$ we can use the result on the difference of consecutive order statistics of exponential random variables from Appendix B to

simplify the above expression further,

$$\Pr\left(C_{\text{MDS}} \leq \frac{mp}{k} - C_0\right) \leq \Pr\left(\sum_{l=k}^{p-1} U_l \geq \frac{\tau C_0}{p-k} - \tau\right) \quad (106)$$

$$\leq \Pr\left((p-k)U_{p-1} \geq \frac{\tau C_0}{p-k} - \tau\right) \quad (107)$$

$$= \Pr\left(U_{p-1} \geq \frac{\tau C_0}{(p-k)^2} - \frac{\tau}{p-k}\right) \quad (108)$$

$$= \exp\left(-\mu\left(\frac{\tau C_0}{(p-k)^2} - \frac{\tau}{p-k}\right)\right) \quad (109)$$

where (107) is obtained from the fact that $\Pr(U_{p-1} \geq u) \geq \Pr(U_l \geq u)$ for $l = k, \dots, p-1$ for any u since $U_l \sim \exp((p-l)\mu)$. Lastly (109) is obtained from the expression for the tail distribution of an exponential random variable. \square

C.4 Replication Strategy

PROOF OF LEMMA 5. In the r -replication strategy each submatrix A_i is replicated at r workers and we wait for the fastest of these r workers. Without loss of generality, we assume that submatrix A_1 is stored at workers $1, 2, \dots, r$, submatrix A_2 is stored at workers $r+1, r+2, \dots, 2*r$ and so on. More generally submatrix A_i is stored at workers $(i-1)r+1, \dots, ir$. Thus the time taken to compute the product $A_i \mathbf{x}$ is given by

$$V_i = \min(Y_{(i-1)r+1}, Y_{(i-1)r+2}, \dots, Y_{ir}) \quad (110)$$

$$= \min(X_{(i-1)r+1} + \tau B_{(i-1)r+1}, \dots, X_{ir} + \tau B_{ir}) \quad (111)$$

$$= \min(X_{(i-1)r+1}, \dots, X_{ir}) + \frac{\tau mr}{p} \quad (112)$$

$$= W_i + \frac{\tau mr}{p}. \quad (113)$$

This is because the fastest of the r workers that store A_i corresponds to $\min(X_{(i-1)r+1}, \dots, X_{ir})$ and this worker must perform $\frac{mr}{p}$ computations to compute the product $A_i \mathbf{x}$.

The latency T_{rep} is the time at which the product $A_i \mathbf{x}$ is computed for all $i = 1, \dots, p/r$ since A is split into p/r submatrices. Thus

$$T_{\text{rep}} = \max(V_1, V_2, \dots, V_{p/r}), \quad (114)$$

$$= \max(W_1, W_2, \dots, W_{p/r}) + \frac{\tau mr}{p}, \quad (115)$$

\square

PROOF OF COROLLARY 4. If $X_j \sim \exp(\mu)$ then observe that $W_i = \min(X_{(i-1)r+1}, \dots, X_{ir})$ is an $\exp(r\mu)$ random variable since it is the minimum of r $\exp(\mu)$ random variables. Thus taking expectation on both sides of (115),

$$\mathbb{E}[T_{\text{rep}}] = \frac{\tau mr}{p} + \mathbb{E}[\max(W_1, W_2, \dots, W_{p/r})], \quad (116)$$

$$= \frac{\tau mr}{p} + \frac{1}{r\mu} H_{p/r}, \quad (117)$$

$$\simeq \frac{\tau mr}{p} + \frac{1}{r\mu} \log \frac{p}{r}, \quad (118)$$

where (117) and (118) follow from (23) and (24). \square

PROOF OF LEMMA 6. As per our model, we represent the number of computations at worker i by the random variable B_i . We also use the random variable C_{rep} to denote the total number of computations performed by all p workers until T_{rep} , which is the time when the master collects enough computations to be able to recover the matrix-vector product $\mathbf{b} = \mathbf{Ax}$. Thus

$$C_{\text{rep}} = B_1 + B_2 + \dots + B_p \quad (119)$$

$$= \sum_{i=1}^{p/r} \sum_{j=1}^r B_{(i-1)r+j}, \quad (120)$$

where the term inside the summation in the second expression represents the number of computations performed by each worker that store a copy of the submatrix \mathbf{A}_i (for a given i). In what follows, we use the shorthand notation $D_j^i = B_{(i-1)r+j}$ and use $D_{j:r}^i$ to denote the order statistics of D_1^i, \dots, D_r^i . Rewriting the above expression in terms of the order statistics we get,

$$C_{\text{rep}} = \sum_{i=1}^{p/r} \sum_{j=1}^r D_{j:r}^i, \quad (121)$$

and the tail bound,

$$\Pr(C_{\text{rep}} \leq mr - C_0) = \Pr\left(\sum_{i=1}^{p/r} \sum_{j=1}^r D_{j:r}^i \leq mr - C_0\right) \quad (122)$$

$$= \Pr\left(\sum_{i=1}^{p/r} \sum_{j=1}^{r-1} D_{j:r}^i \leq m(r-1) - C_0\right) \quad (123)$$

$$\leq \Pr\left((r-1) \sum_{i=1}^{p/r} D_{1:r}^i \leq m(r-1) - C_0\right) \quad (124)$$

$$= \Pr\left(\sum_{i=1}^{p/r} D_{1:r}^i \leq m - \frac{C_0}{r-1}\right) \quad (125)$$

where (123) follows from the fact that for any given submatrix \mathbf{A}_i , $i = 1, \dots, p/r$, the fastest worker that stores a copy of that submatrix, which corresponds to $D_{r:r}^i$ (fastest worker performs the most computations) must perform all the tasks assigned to it i.e. mr/p computations each, while (124) follows from the fact that $D_{2:r}^i, \dots, D_{r:r}^i$ are always larger than $D_{1:r}^i$ by definition.

At this point we introduce the shorthand notation $V_j^i = X_{(i-1)r+j}$ for the setup time of the worker that stores the j^{th} copy of submatrix \mathbf{A}_i and note that the worker which performs $D_{1:r}^i$ computations has setup time $V_{r:r}^i$ ($V_{j:r}^i$ are the order statistics of V_1^i, \dots, V_r^i). There can be two possibilities – either $T_{\text{rep}} > V_{r:r}^i$, or $T_{\text{rep}} \leq V_{r:r}^i$. If $T_{\text{rep}} > V_{r:r}^i$ then

$$T_{\text{rep}} \leq V_{r:r}^i + \tau(D_{1:r}^i + 1) \quad (126)$$

where the added 1 accounts for the edge effect of partial computations at the nodes. If $T_{\text{rep}} \leq V_{r:r}^i$ then also the upper bound (126) holds. Thus overall (by rearranging terms in (126)) we obtain,

$$D_{1:r}^i \geq \frac{T_{\text{rep}} - V_{r:r}^i}{\tau} - 1 \quad (127)$$

Thus we can write

$$\Pr(C_{\text{rep}} \leq mr - C_0) \leq \Pr\left(\sum_{i=1}^{p/r} \left(\frac{T_{\text{rep}} - V_{r:r}^i}{\tau} - 1\right) \leq m - \frac{C_0}{r-1}\right) \quad (128)$$

$$= \Pr\left(\sum_{i=1}^{p/r} (V_{r:r}^i - W_{\text{rep}}) \geq \frac{\tau C_0}{r-1} - \frac{\tau p}{r}\right) \quad (129)$$

where $W_{\text{rep}} = \max_i V_{1:r}^i = \max_{1 \leq i \leq p/r} \min_{1 \leq j \leq r} X_{(i-1)r+j}$ and (129) follows from the fact that $T_{\text{rep}} = \max_{1 \leq i \leq p/r} \min_{1 \leq j \leq r} X_{(i-1)r+j} + \tau mr/p$. From our definition of W_{rep} we see that,

$$V_{r:r}^i - W_{\text{rep}} \leq V_{r:r}^i - V_{1:r}^i \quad (130)$$

and the consequent stochastic dominance can be used to get an upper bound on (129) as,

$$\Pr(C_{\text{rep}} \leq mr - C_0) \leq \Pr\left(\sum_{i=1}^{p/r} (V_{r:r}^i - V_{1:r}^i) \geq \frac{\tau C_0}{r-1} - \frac{\tau p}{r}\right) \quad (131)$$

$$= \Pr\left(\sum_{i=1}^{p/r} \sum_{j=1}^{r-1} (V_{j+1:r}^i - V_{j:r}^i) \geq \frac{\tau C_0}{r-1} - \frac{\tau p}{r}\right) \quad (132)$$

If $X_i \sim \exp(\mu)$ we can use the result from Appendix B on the difference of consecutive order statistics of exponential random variables to simplify the above expression further (since $V_j^i = X_{(i-1)r+j}$ are also exponentially distributed and thus $U_j^i = (V_{j+1:r}^i - V_{j:r}^i) \sim \exp((r-j)\mu)$),

$$\Pr(C_{\text{rep}} \leq mr - C_0) \leq \Pr\left(\sum_{i=1}^{p/r} \sum_{j=1}^{r-1} U_j^i \geq \frac{\tau C_0}{r-1} - \frac{\tau p}{r}\right) \quad (133)$$

$$\leq \Pr\left((r-1) \sum_{i=1}^{p/r} U_{r-1}^i \geq \frac{\tau C_0}{r-1} - \frac{\tau p}{r}\right) \quad (134)$$

$$= \Pr\left(\sum_{i=1}^{p/r} U_{r-1}^i \geq \frac{\tau C_0}{(r-1)^2} - \frac{\tau p}{r(r-1)}\right) \quad (135)$$

$$= \sum_{i=0}^{p/r-1} \frac{1}{i!} \exp(-\mu\theta)(\mu\theta)^i. \quad (136)$$

where (134) is obtained from the fact that $\Pr(U_{r-1}^i \geq u) \geq \Pr(U_j^i \geq u)$ for $j = 1, \dots, r-2$ for any u since $U_j \sim \exp((r-j)\mu)$. Lastly (136) is obtained from the expression for the tail distribution of an Erlang random variable which is the sum of p/r exponential random variables with rate μ and

$$\theta = \frac{\tau C_0}{(r-1)^2} - \frac{\tau p}{r(r-1)} \quad (137)$$

□

D THEORETICAL RESULTS FOR QUEUEING ANALYSIS

PROOF OF THEOREM 5. For large m_e i.e. $\alpha = m_e/m \rightarrow \infty$, under the proposed rateless coded strategy we just need to wait for M' computations to be performed by the workers in total. Hence the p workers can be treated as a single server with service time equal to T_{LT} . According to the

Pollaczek Khinchine formula [25] for Poisson arrivals with rate λ , the expected total processing time (time in queue + service time) for a job in such a queue is given by

$$\mathbb{E}[Z_{LT}] = \mathbb{E}[T_{LT}] + \frac{\lambda \mathbb{E}[(T_{LT})^2]}{2(1 - \lambda \mathbb{E}[T_{LT}])} \quad (138)$$

Recall that for $\alpha = m_e/m \rightarrow \infty$ the LT and Ideal schemes are identical by definition. Therefore we can use the bounds derived for T_{ideal} in Lemma 2 to bound T_{LT} as,

$$T_{LT} \leq \tau \left(\frac{M'}{p} + 1 \right) + \frac{1}{p} \sum_{i=1}^p X_i \quad (139)$$

$$T_{LT} \geq \tau \frac{M'}{p} + X_{1:p} \quad (140)$$

To analyze the second moment $\mathbb{E}[(T_{LT})^2]$, let $\tau' = \tau/p$ and let $\bar{X} = (1/p) \sum_{i=1}^p X_i$. Then,

$$(T_{LT})^2 \leq (\tau + \tau' M')^2 + \bar{X}^2 + 2(\tau + \tau' M')\bar{X} \quad (141)$$

$$(T_{LT})^2 \geq (\tau' M')^2 + \bar{X}^2 + 2\tau' M'\bar{X} \quad (142)$$

Taking expectations and noting that M' and \bar{X} are independent,

$$\begin{aligned} \mathbb{E}[(T_{LT})^2] &\leq \tau^2 + 2\tau\tau' m_d + (\tau')^2 \mathbb{E}[(M')^2] + \mathbb{E}[\bar{X}^2] + \\ &\quad 2(\tau + \tau' m_d) \mathbb{E}[\bar{X}] \end{aligned} \quad (143)$$

$$\mathbb{E}[(T_{LT})^2] \geq (\tau')^2 \mathbb{E}[(M')^2] + \mathbb{E}[\bar{X}^2] + 2\tau' m_d \mathbb{E}[\bar{X}] \quad (144)$$

□

LEMMA 12 (MULTIPLE JOBS WITH MDS CODING). *The expected latency of the MDS coded scheme Z_{MDS} when a stream of vectors $\mathbf{x}_1, \mathbf{x}_2, \dots$ need to be multiplied with the same matrix \mathbf{A} (assuming Poisson arrivals with rate λ for the vectors) is given by*

$$\mathbb{E}[Z_{MDS}] \leq \mathbb{E}[Y_{k:p}] + \frac{\lambda((\mathbb{E}[Y_{k:p}])^2 + \text{Var}[Y_{k:p}])}{2(1 - \lambda \mathbb{E}[Y_{k:p}])} \quad (145)$$

$$\mathbb{E}[Z_{MDS}] \geq \mathbb{E}[Y_{k:p}] + \frac{\lambda((\mathbb{E}[Y_{1:p}])^2 + \text{Var}[Y_{1:p}])}{2(1 - \lambda \mathbb{E}[Y_{1:p}])} \quad (146)$$

where $Y_i = X_i + \tau m/k$ is the service time at worker i , $i = 1, \dots, p$ for the MDS coded case.

PROOF OF LEMMA 12. When a stream of incoming vectors $\mathbf{x}_1, \mathbf{x}_2, \dots$ need to be multiplied with the matrix \mathbf{A} over p workers, the resulting system is a (p, k) fork-join queue since the task of computing matrix vector products of the form $\mathbf{A}\mathbf{x}$ is forked to the p workers and we need to wait for k workers to complete the tasks assigned to them. The expression for latency Z_{MDS} follows from Theorem 4 of [35] which gives bounds on the latency of (p, k) fork-join queues assuming Poisson arrivals. □

LEMMA 13 (MULTIPLE JOBS WITH REPLICATION). *The expected latency of the replication scheme Z_{rep} when a stream of vectors $\mathbf{x}_1, \mathbf{x}_2, \dots$ need to be multiplied with the same matrix \mathbf{A} (assuming Poisson arrivals with rate λ for the vectors) is given by*

$$\mathbb{E}[Z_{rep}] \leq \mathbb{E}[V_{p/r:p/r}] + \frac{\lambda((\mathbb{E}[V_{p/r:p/r}])^2 + \text{Var}[V_{p/r:p/r}])}{2(1 - \lambda \mathbb{E}[V_{p/r:p/r}])} \quad (147)$$

$$\mathbb{E}[Z_{rep}] \geq \mathbb{E}[V_{p/r:p/r}] + \frac{\lambda((\mathbb{E}[V_{1:p/r}])^2 + \text{Var}[V_{1:p/r}])}{2(1 - \lambda \mathbb{E}[V_{1:p/r}])} \quad (148)$$

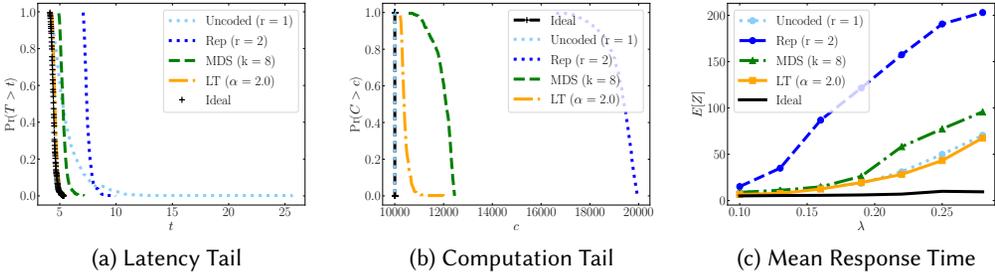


Fig. 11. The tail probability of the latency is the highest for the replication schemes. MDS codes perform better in terms of latency but they perform a large number of redundant computations. The latency tail of LT codes is the minimum among all the schemes. Moreover the LT coded schemes performs significantly fewer redundant computations than MDS Codes or replication. When there are multiple jobs in the queue, the mean response time is least for the LT Coded setting under all values of arrival rate λ . All simulations are performed with $m = 10000$ matrix rows, $p = 10$ worker nodes, $\tau = 0.001$, and a Pareto (1,3) distribution on the initial delays $X_i, i = 1, \dots, p$

where $W_i = \min_j X_{(i-1)r+j} + \tau mr/p, i = 1, \dots, p/r$.

PROOF OF LEMMA 13. When a stream of incoming vectors $\mathbf{x}_1, \mathbf{x}_2, \dots$ need to be multiplied with the matrix \mathbf{A} over p workers, the resulting system is a $(p/r, p/r)$ fork-join queue. This is because each submatrix $\mathbf{A}_1, \dots, \mathbf{A}_{p/r}$ is replicated at r workers and we need to wait for the fastest worker for each submatrix. Thus i^{th} group of r workers has an effective service time of V_i as defined in (113). Since we need to wait for all p/r groups of r workers in this fashion it is equivalent to a $(p/r, p/r)$ fork-join queue where the i^{th} node has service time V_i . Once again the expression for latency Z_{rep} follows from Theorem 4 of [35] assuming Poisson arrivals. \square

E ADDITIONAL THEORETICAL RESULTS

The following theorem compares the latency of the MDS coded, and ideal load balancing strategies. It indicates that T_{MDS} approaches T_{ideal} only when the fastest k workers start computing at approximately the same time while the stragglers do not start until much later. In this rare situation neglecting partial computations by stragglers does not adversely impact latency performance. However in all other cases discarding work done by stragglers causes T_{MDS} to be larger than T_{ideal} .

THEOREM 6 (MDS V/S IDEAL). *The latency of the MDS coded strategy, T_{MDS} is larger than the latency of the ideal strategy T_{ideal} with a high probability. Specifically $\Pr(T_{\text{MDS}} > T_{\text{ideal}}) = 1 - \delta_{\text{MDS}}$ where*

$$\delta_{\text{MDS}} = \Pr(X_{k:p} - X_{1:p} \leq \tau, X_{k+1:p} - X_{k:p} > \tau \left(\frac{m}{k} - 1 \right)) \quad (149)$$

PROOF. To compare the latency of the MDS coded strategy to that of the ideal scheme we note that the latency T_{ideal} of the ideal scheme, is the earliest time when $\sum_{i=1}^p B_i = m$, as illustrated in Fig. 6. We note that, in this case it is not necessary that each worker has completed at least 1 computation. Specifically, if $T_{\text{ideal}} - X_i \leq \tau$ for any i then it means that worker i has not performed even a single computation in the time that the system as a whole has completed m computations (owing to the large initial delay X_i). Therefore we define

$$\mathcal{W}_{\text{ideal}} := \{i : T_{\text{ideal}} - X_i \geq \tau\} \quad (150)$$

Here $\mathcal{W}_{\text{ideal}}$ is the set of workers for which $B_i > 0$ in time up to T_{ideal} . We also note that

$$T_{\text{ideal}} < X_i + \tau(B_i + 1), \quad \text{for all } i = 1, \dots, p \quad (151)$$

This is because at time T_{ideal} each of the workers $1, \dots, p$, have completed B_1, \dots, B_p row-vector product tasks respectively, but they may have partially completed the next task. The 1 added to each B_i accounts for this edge effect, which is also illustrated in Fig. 6.

We will compare T_{MDS} and T_{ideal} for the following three cases assuming the same realizations of initial delay X_i , $i = 1, \dots, p$ for both schemes (we will also assume without loss of generality that $X_1 < X_2 < \dots < X_p$ i.e $X_i = X_{i:p} \forall i$):

- Case 1: If node $k \notin \mathcal{W}_{\text{ideal}}$, $B_k = 0$ in the ideal scheme. Thus if we use the ideal scheme the latency is $T_{\text{ideal}} < X_{k:p} + \tau$ (from (151)) whereas if we use the MDS coded scheme the latency is $T_{\text{MDS}} = X_{k:p} + \tau m/k > T_{\text{ideal}}$.
- Case 2: If node $k \in \mathcal{W}_{\text{ideal}}$ but $X_k \neq \max_{i \in \mathcal{W}_{\text{ideal}}} X_i$ then nodes $1, \dots, k-1$ must also lie in $\mathcal{W}_{\text{ideal}}$ (since $X_1 < X_2 < \dots < X_k$) and at least one of nodes $k+1, \dots, p$ must also lie in $\mathcal{W}_{\text{ideal}}$ (since $X_k \neq \max_{i \in \mathcal{W}_{\text{ideal}}} X_i$). If we use the ideal scheme and assume that node k performs m/k computations then nodes $1, \dots, k-1$ must perform m/k or more computations each (since $X_1 < X_2 < \dots < X_k$). Together nodes $1, \dots, p$ perform at least $k \times m/k = m$ computations and since at least one of nodes $k+1, \dots, p$ must also lie in $\mathcal{W}_{\text{ideal}}$, the total number of computations performed is greater than m which is not possible since the number of computations performed in the ideal scheme is m by definition. Thus node k can perform at most $(m/k - 1)$ computations and thus from (151) we have that $T_{\text{ideal}} \leq X_{k:p} + \tau(m/k - 1) < T_{\text{MDS}}$.
- Case 3: If node $k \in \mathcal{W}_{\text{ideal}}$ and $X_k = \max_{i \in \mathcal{W}_{\text{ideal}}} X_i$ then there are exactly k workers in $\mathcal{W}_{\text{ideal}}$. If $X_k - X_1 > \tau$ then node 1 performs at least 1 more computation than node k since it takes time τ for a node to perform a computation. In that case node k performs fewer than m/k computations since total number of computations performed by nodes $1, \dots, k$ must be m . Thus in this case $T_{\text{ideal}} < X_k + \tau m/k = T_{\text{MDS}}$ (from (151)). Likewise for nodes $2, \dots, k-1$ as well. Thus $T_{\text{ideal}} = T_{\text{MDS}}$ only when $X_1, \dots, X_{k-1} \geq X_k - \tau$

Thus overall it is only in Case 3 above under specific circumstances that we can have $T_{\text{ideal}} = T_{\text{MDS}}$. In all other cases $T_{\text{ideal}} < T_{\text{MDS}}$. Hence probabilistically $\Pr(T_{\text{MDS}} > T_{\text{ideal}}) = 1 - \delta_{\text{MDS}}$ where

$$\delta_{\text{MDS}} = \Pr(X_{k:p} - X_{1:p} \leq \tau, X_{k+1:p} - X_{k:p} > \tau \left(\frac{m}{k} - 1 \right)) \quad (152)$$

where the first condition is equivalent to $X_{1:p}, \dots, X_{k-1:p} \leq X_{k:p}$ by the definition of the order statistics and the second condition ensures that $X_{k:p} = \max_{i \in \mathcal{W}_{\text{ideal}}} X_i$ (nodes $k+1, \dots, p$ do not start until node k has completed m/k computations). \square

The following theorem shows that except in rare cases T_{Rep} is larger than T_{ideal} .

THEOREM 7 (REPLICATION V/S IDEAL). *The latency of the replication coded strategy, T_{rep} is larger than the latency of the ideal strategy T_{ideal} with a high probability. Specifically $\Pr(T_{\text{rep}} > T_{\text{ideal}}) = 1 - \delta_{\text{rep}}$ where*

$$\delta_{\text{rep}} = \Pr \left(\min_{1 \leq i \leq p/r} X_{2:r}^{(i)} - \max_{1 \leq i \leq p/r} X_{1:r}^{(i)} > \tau \left(\frac{mr}{p} - 1 \right) \right) \quad (153)$$

and $X_j^{(i)} = X_{(i-1)r+j}$.

PROOF. In the r -replication scheme the matrix \mathbf{A} is split into p/r submatrices $\mathbf{A}_1, \dots, \mathbf{A}_{p/r}$. Each submatrix is replicated at r distinct workers and we wait for the fastest worker for each submatrix. There are p/r such groups of workers, with all workers in group i computing the product $\mathbf{A}_i \mathbf{x}$.

Following similar arguments to the MDS coded case, we can conclude that $T_{\text{rep}} = T_{\text{ideal}}$ only when the workers in each group other than the fastest worker do not start computing until the fastest workers of *all* groups have finished their computations. The fastest worker group i is the one that takes time V_i (in (113)) to compute $A_i \mathbf{x}$.

This is because in any other scenario the workers in the replications scheme that are not the fastest in their respective groups perform redundant computations and hence some work is wasted as compared to the ideal scheme where some computations could have been transferred from the fastest workers to the slower workers in each group thus reducing the total time taken to compute m row-vector products. However if the slower workers in the groups do not start their computations until after *all* the fastest workers across *all* groups complete their computations then the replication and ideal schemes are essentially identical. This is because in this case even in the ideal scheme no work could have been transferred from the fastest workers in each group to the slow workers.

The following condition quantifies the above explanation,

$$\min_{1 \leq i \leq p/r} X_{2:r}^{(i)} - \max_{1 \leq i \leq p/r} X_{1:r}^{(i)} > \tau \left(\frac{mr}{p} - 1 \right) \quad (154)$$

This is because $X_{1:r}^{(i)} + \tau(mr/p - 1)$ is the time at which the fastest worker in group i (with initial delay corresponding to $X_{1:r}^{(i)}$) completes $mr/p - 1$ computations and starts working on the last $((mr/p)^{\text{th}})$ computation. Thus even if the second slowest worker of the group (corresponding to $X_{2:r}^{(i)}$) and any other worker(s) starts after this time they cannot perform any computations before the fastest worker completes all the computations assigned to it. Thus no computations are transferred to the slower workers even in the ideal scheme in this case.

Here $\max_{1 \leq i \leq p/r} X_{1:r}^{(i)} + \tau \left(\frac{mr}{p} - 1 \right)$ is the time taken by the slowest of the fastest workers across all groups to complete $mr/p - 1$ computations and $\min_{1 \leq i \leq p/r} X_{2:r}^{(i)}$ is the initial delay of the worker with the least initial delay among the second fastest workers of all groups.

Therefore $\Pr(T_{\text{rep}} > T_{\text{ideal}}) = 1 - \delta_{\text{rep}}$ where

$$\delta_{\text{rep}} = \Pr \left(\min_{1 \leq i \leq p/r} X_{2:r}^{(i)} - \max_{1 \leq i \leq p/r} X_{1:r}^{(i)} > \tau \left(\frac{mr}{p} - 1 \right) \right) \quad (155)$$

and $X_j^{(i)} = X_{(i-1)r+j}$ □

F ADDITIONAL SIMULATIONS AND EXPERIMENTS

Additional Simulations: We simulate the MDS, replication, and LT-coded schemes under our delay model ((5)) with initial delays X_i distributed according to a Pareto (1,3) distribution. The results are summarised in Fig. 11. Once again we observe that LT coding ($\alpha = 2.0$) clearly outperforms MDS coding ($k = 8$) both in terms of latency (Fig. 11a), number of redundant computations (Fig. 11b), and mean response time averaged over 10 trials with 100 jobs per trial and Poisson (λ) arrivals (Fig. 11c).

Additional Experiments: We study the effect of worker failures in the different coded computing strategies using an EC2 [1] cluster 10 t2.micro workers. Note that redundancy/coding is essential if workers fail as the naive uncoded approach cannot handle the resulting loss of data. A is a 10000×10000 identity matrix and is encoded using replication ($r = 2$), MDS coding ($k = 5$), and LT coding ($\alpha = 2.0$). Results in Fig. 12 clearly show that the LT coded approach is more robust multiple node failures than the other approaches.

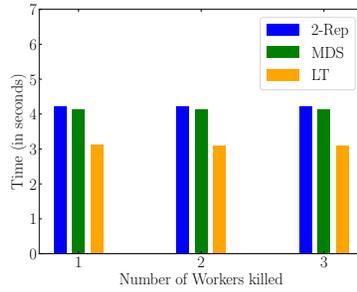


Fig. 12. Worker i has a random initial delay X_i , after which it completes row-vector product tasks (denoted by the small rectangles), taking time τ per task. The latency T is the time until enough tasks have been completed for the product $\mathbf{b} = \mathbf{Ax}$ to be recovered.