

On the Upload versus Download Cost for Secure and Private Matrix Multiplication

Wei-Ting Chang Ravi Tandon

Department of Electrical and Computer Engineering

University of Arizona, Tucson, AZ, USA

E-mail: {wchang, tandonr}@email.arizona.edu

Abstract—In this paper, we study the problem of secure and private distributed matrix multiplication. Specifically, we focus on a scenario where a user wants to compute the product of a confidential matrix A , with a matrix B_θ , where $\theta \in \{1, \dots, M\}$. The set of candidate matrices $\{B_1, \dots, B_M\}$ are public, and available at all the N servers. The goal of the user is to distributedly compute AB_θ , such that (a) no information is leaked about the matrix A to any server; and (b) the index θ is kept private from each server. Our goal is to understand the fundamental tradeoff between the upload vs download cost for this problem. Our main contribution is to show that the lower convex hull of following (upload, download) pairs: $(U, D) = (N/(K-1), (K/(K-1))(1 + (K/N) + \dots + (K/N)^{M-1}))$ for $K = 2, \dots, N$ is achievable. The scheme improves upon state-of-the-art existing schemes for this problem, and leverages ideas from secret sharing and coded private information retrieval.

Keywords – Distributed Matrix Multiplication, Secure Matrix Multiplication, Private Information Retrieval.

I. INTRODUCTION

Distributed computing aims at speeding up computationally intensive operations by dividing and outsourcing the computation to multiple servers. Distributed processing often comes with additional communication overhead, when compared to centralized processing. Furthermore, when processing sensitive information in a distributed manner, security and privacy are also critical concerns. Hence, simultaneously utilizing distributed processing, while satisfying security/privacy constraints are critical and of great interest.

A well known issue in distributed systems is the *straggler effect*, where slower machines in the system slow down the entire computation. Recently, error control codes have been successfully adopted to carefully add computational redundancy in order to tackle the stragglers problem. For instance, entangled polynomial codes [1] and PolyDot codes [2] have been proposed for the stragglers problem, within the context of distributed matrix multiplication.

Although mitigating stragglers is important, it is crucial to ensure that the data, if confidential, is not leaked to unauthorized servers. Recently, another line of research has focused on the security aspects for the distributed matrix multiplication problem. Secure matrix-vector multiplication was considered in [3]. In [4], the authors considered a variation of the problem where the matrices are kept secure against the servers as well as the user. The authors in [5] proposed the idea of Lagrange

This work was supported by NSF Grants CAREER 1651492, and CNS 1715947.

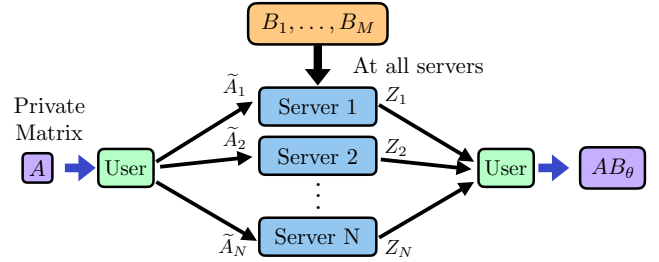


Fig. 1: System model for secure and private matrix multiplication: multiply a confidential matrix A with B_θ , while keeping the index θ private from all N servers.

coded computing that simultaneously deals with straggler, colluding servers and malicious servers. Later, [6], [7] studied the problem of secure distributed matrix multiplication. Both works proposed schemes based on polynomial codes, where the authors of [7] studied the non-colluding version of the problem and the authors of [6] studied the colluding version (also see recent improvements and variants in [8]–[11]).

Another issue that needs to be addressed is privacy of the user. Consider a scenario where in addition to the confidential data, the queries sent for computation are also sensitive, i.e., the queries can leak private information about the user. The authors of [12] considered both security and privacy for the problem of distributed matrix multiplication. There are several other works that have focused primarily on the privacy aspect. In particular, starting from the work of [13], private information retrieval (PIR) has been studied extensively within an information theoretic framework (such as PIR from coded databases [14] and private computation [15]–[18]).

Main Contributions: In this paper, we focus on the model where a user wants to compute the product of a confidential matrix A , with a matrix B_θ , where $\theta \in \{1, \dots, M\}$. The set of candidate matrices $\{B_1, \dots, B_M\}$ are public, and available at all the N servers. The goal of the user is to distributedly compute AB_θ , such that (a) no information is leaked about the matrix A to any server; and (b) the index θ is kept private from each server. Our main contribution is to show that the lower convex hull of following (upload, download) pairs: $(U, D) = (N/(K-1), (K/(K-1))(1 + (K/N) + \dots + (K/N)^{M-1}))$ for $K = 2, \dots, N$ is achievable. Our scheme also improves upon the scheme of [12].

II. PROBLEM FORMULATION

We consider a distributed system with one user and N non-colluding servers. The user has a confidential matrix $A \in \mathbb{F}^{d_1 \times d_2}$, and all servers have access to M matrices $B_m \in \mathbb{F}^{d_2 \times d_3}$, $m = 1, \dots, M$, for some integers d_1, d_2 and d_3 , and a sufficiently large field \mathbb{F} . We assume that the matrices A, B_1, \dots, B_M are all independent of each other, and each with i.i.d. entries from \mathbb{F} . The goal of the user is to compute the product of a confidential matrix A with B_θ , where the index θ is private. Each server is connected to the user through a separate link. We assume that all servers are honest but curious (i.e., all servers correctly follow the protocols, however, they are interested in learning about A and θ). To ensure security, the user securely encodes A using encoding functions $\mathbf{f} = (f_1, \dots, f_N)$, where f_n is the individual encoding function for server n . We denote the encoded version of A that is sent to server n by \tilde{A}_n , i.e., $\tilde{A}_n = f_n(A)$. Along with \tilde{A}_n , the user also sends a query $Q_n^{(\theta)}$ to server n . Without any prior knowledge of the stored data, the queries sent by the user are independent of all the B_m 's, i.e., $I(Q_1^{(\theta)}, \dots, Q_N^{(\theta)}; B_1, \dots, B_M) = 0$. Each server n uses a computing function $Z_n^{(\theta)} : \{B_1, \dots, B_M\} \times \tilde{A}_n \times Q_n^{(\theta)} \rightarrow Z_n^{(\theta)}$ for the assigned task. Once servers finish the requested computations, servers return their answers $Z_n^{(\theta)}$, $n = 1, \dots, N$, to the user. Next, to preserve the privacy of the user, the strategy the user used to download the desired result should not reveal θ . That is, for any θ and θ' , the query and the answer from server n should be statistically identical. Upon receiving all the $Z_n^{(\theta)}$'s, the user decodes the desired result by using decoding function, $AB_\theta = g(Z_1^{(\theta)}, Q_1^{(\theta)}, \dots, Z_N^{(\theta)}, Q_N^{(\theta)})$. A scheme is secure, private and reliable if it satisfies the following constraints:

Security Constraint:

$$I(A; \tilde{A}_n, Q_n^{(\theta)}, B_1, \dots, B_M) = 0, \quad n = 1, \dots, N. \quad (1)$$

Privacy Constraint:

$$(\tilde{A}_n, Q_n^{(\theta)}, Z_n^{(\theta)}, B_{[1:M]}) \sim (\tilde{A}_n, Q_n^{(\theta')}, Z_n^{(\theta')}, B_{[1:M]}), \quad \forall n \quad (2)$$

Decodability Constraint:

$$H(AB_\theta | Z_1^{(\theta)}, Q_1^{(\theta)}, \dots, Z_N^{(\theta)}, Q_N^{(\theta)}) = 0. \quad (3)$$

The performance of a scheme is determined by the normalized upload cost (U) and download cost (D), defined as:

$$U = \frac{\sum_{n=1}^N H(\tilde{A}_n)}{H(AB_\theta)}, \quad \text{and} \quad D = \frac{\sum_{n=1}^N H(Z_n^{(\theta)})}{H(AB_\theta)}. \quad (4)$$

Our aim is to understand the tradeoff between the upload and download costs, i.e., the set of all feasible (U, D) pairs. The optimal download cost for a fixed upload cost is defined as:

$$D^*(U) \triangleq \min\{D : (U, D) \text{ is feasible}\}. \quad (5)$$

The following Lemma shows that the optimal download cost is a convex function of the upload cost.

Lemma 1. *The optimal download cost $D^*(U)$ is a convex function of the upload cost U .*

Proof: To prove this Lemma, we show that for any two upload costs U_1, U_2 and their corresponding optimal download costs $D^*(U_1), D^*(U_2)$, there is a scheme with upload cost $\bar{U} = \alpha U_1 + (1 - \alpha)U_2$ that achieves a download cost of $\bar{D} = \alpha D^*(U_1) + (1 - \alpha)D^*(U_2)$, $0 \leq \alpha \leq 1$. This can be shown by memory sharing argument where we partition A into two parts, $A^{(\alpha)} \in \mathbb{F}^{\alpha d_1 \times d_2}$ and $A^{(1-\alpha)} \in \mathbb{F}^{(1-\alpha)d_1 \times d_2}$, respectively. $A^{(\alpha)}$ is sent to servers securely using scheme 1 with upload cost U_1 . The answers are returned to the user using the optimal scheme that corresponds to U_1 , hence, achieves the download cost of $D^*(U_1)$. Similarly, $A^{(1-\alpha)}$ is sent to servers and the answers are returned to the user using scheme 2 that achieves $(U_2, D^*(U_2))$. The desired results obtained from both parts are of size $\alpha d_1 \times d_3$ and $(1 - \alpha)d_1 \times d_3$. Since these are per-bit costs, the total upload and download costs are

$$U_{\text{total}} = \alpha U_1 d_1 d_2 + (1 - \alpha)U_2 d_1 d_2 = \bar{U} d_1 d_2, \quad (6)$$

$$D_{\text{total}} = \alpha D^*(U_1) d_1 d_3 + (1 - \alpha)D^*(U_2) d_1 d_3 = \bar{D} d_1 d_3. \quad (7)$$

Clearly, the optimal $D^*(\bar{U})$ is upper bounded by $\bar{D}(\bar{U})$ completing the proof of Lemma 1.

In the next proposition, we show that even without any privacy constraints, the minimum upload cost is $N/(N - 1)$.

Proposition 1. *The minimum value of upload cost is $U_{\min} = N/(N - 1)$.*

Proof: We start with the following sequence of inequalities:

$$\begin{aligned} \sum_{n=1}^N H(\tilde{A}_n) &\geq H(\tilde{A}_1, \dots, \tilde{A}_N) \\ &\stackrel{(a)}{=} H(\tilde{A}_1, \dots, \tilde{A}_N | B_1, \dots, B_M) \\ &\stackrel{(b)}{=} H(AB_\theta, \tilde{A}_1, \dots, \tilde{A}_N | B_1, \dots, B_M) \\ &= H(AB_\theta | B_1, \dots, B_M) \\ &\quad + H(\tilde{A}_1, \dots, \tilde{A}_N | AB_\theta, B_1, \dots, B_M) \\ &\stackrel{(c)}{\geq} H(AB_\theta) + H(\tilde{A}_n | AB_\theta, B_1, \dots, B_M) \\ &\stackrel{(d)}{=} H(AB_\theta) + H(\tilde{A}_n | B_1, \dots, B_M) \\ &\stackrel{(e)}{=} H(AB_\theta) + H(\tilde{A}_n) \end{aligned} \quad (8)$$

where (a) follows from the fact that \tilde{A}_n 's are independent of B_m 's for all m, n ; (b) follows from decodability constraint; (c) is due to the fact that $H(AB_\theta | B_1, \dots, B_M) = H(AB_\theta)$ by security constraint; (d) also follows from security constraint and the fact that A, B_1, \dots, B_M are all independent of each other so that AB_θ can be removed from the conditioning; (e) follows due to independence of \tilde{A}_n and $\{B_1, \dots, B_M\}$. Summing up (8) for all n , we have $N(\sum_{n=1}^N H(\tilde{A}_n)) \geq NH(AB_\theta) + \sum_{n=1}^N H(\tilde{A}_n)$. By rearranging, we obtain

$$U = \frac{\sum_{n=1}^N H(\tilde{A}_n)}{H(AB_\theta)} \geq \frac{N}{N - 1}, \quad (9)$$

completing the proof of Proposition 1.

III. MAIN RESULTS AND DISCUSSION

Theorem 1. For the secure and private matrix multiplication problem with N non-colluding servers, where each server has access to M public matrices, the lower convex hull of

$$(U, D) = \left(\frac{N}{K-1}, \frac{K}{K-1} \left(1 + \frac{K}{N} + \dots + \left(\frac{K}{N} \right)^{M-1} \right) \right)$$

for $K = 2, \dots, N$ is achievable.

Remark 1. In Theorem 1, by varying the parameter K , we can trade off upload cost for download cost. We note that for minimum upload cost, i.e., when $K = N$ and $U = N/(N-1)$, and $M = 1$ (i.e., no privacy constraint), the download cost of $D = N/(N-1)$ is information-theoretically optimal as shown in our previous work [6].

Remark 2. We compare our scheme to the scheme in [12]. For any two integers (m_1, m_2) , the scheme in [12] uses $N = (m_1 + 1)(m_2 + 1)$ servers to achieve a lower convex hull of $(U, D) = (N/m_1, N/m_1 m_2)$. In the simulation, we let $N = 12$, $M = 6$ and $K = 2, \dots, 12$ for our proposed scheme. We used $(m_1, m_2) = (1, 5), (2, 3), (3, 2), (5, 1)$ for the scheme in [12] to ensure that $N = 12$. The scheme in [12] essentially achieves $(U, D) = (N/(K-1), (K/(K-1))(N/(N-K)))$ for $K = 2, 3, 4, 6$ in this example. In Fig. 2, we compare the proposed scheme of this paper with that in [12]. We next make the following observations: for the same upload cost, the download cost of the proposed scheme is smaller. Furthermore, there are several (U, D) pairs that the scheme in [12] cannot achieve, particularly for smaller values of upload cost.

A. Illustrative Example: ($N = 4, M = 2, K = 3$)

We next demonstrate the proposed achievable scheme through an example to illustrate the main ideas. Suppose there are $N = 4$ non-colluding servers and each server has access to all $M = 2$ matrices, namely B_1 and B_2 , whose dimensions are all $d_2 \times d_3$. Assume that the user wants to compute AB_1 , and the index $\theta = 1$ must be kept private. The user first partitions the matrix A into $A = [A_1^T A_2^T]^T$, where each A_i is of size $(d_1/2) \times d_2$, $i = 1, 2$. The desired computation AB_1 becomes $AB_1 = [(A_1 B_1)^T (A_2 B_1)^T]^T$. To provide security, A_i 's are encoded using a secure $(N, K) = (4, 3)$ MDS code as follows:

$$\tilde{A}_n = A_1 + A_2 x_n + R x_n^2, \quad n = 1, \dots, 4, \quad (10)$$

where $R \in \mathbb{F}^{(d_1/2) \times d_2}$ is a random matrix (independent of A), whose entries are i.i.d. uniform random variables, and x_n is a distinct element in \mathbb{F} assigned to server n . The user then sends \tilde{A}_n to server n and instructs them to multiply their respective \tilde{A}_n with all B_m 's. Specifically, each server n computes

$$\tilde{A}_n B_m = A_1 B_m + A_2 B_m x_n + R B_m x_n^2, \quad m = 1, 2. \quad (11)$$

To simplify the notation, we let $h_n^T W^{(m)} \triangleq \tilde{A}_n B_m$, where

$$h_n = \begin{bmatrix} 1 \\ x_n \\ x_n^2 \end{bmatrix}, \quad W^{(m)} = \begin{bmatrix} A_1 B_m \\ A_2 B_m \\ R B_m \end{bmatrix}. \quad (12)$$

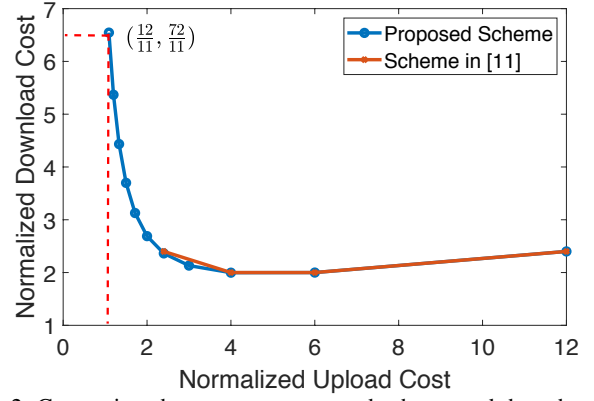


Fig. 2: Comparison between our proposed scheme and the scheme in [12], with $N = 12$, $M = 6$ and $K = 2, 3, \dots, 12$.

Due to the properties of the $(4, 3)$ MDS code, we need at least three different $h_n^T W^{(m)}$ (viewed as three evaluations of a polynomial) to decode the desired result $W^{(1)}$ whose components can then be used to recover AB_1 (via polynomial interpolation). However, only downloading $h_n^T W^{(1)}$ from any three servers will clearly violate the privacy constraint. Thus, to retrieve the answers while preserving privacy, we adopt a similar downloading strategy for the PIR problem with coded databases in [14]. The user asks each server n to divide each $h_n^T W^{(m)}$, $m = 1, 2$ into $N^M = 4^2 = 16$ blocks vertically, i.e.,

$$h_n^T W^{(m)} = \begin{bmatrix} h_n^T W_1^{(m)} \\ \vdots \\ h_n^T W_{16}^{(m)} \end{bmatrix}, \quad \forall n, m. \quad (13)$$

Note that each block is composed of rows of the results, each one of which is coded using $(4, 3)$ MDS code.

The privacy preserving download scheme is broken into repetitions and rounds. Within each repetition, there are $M = 2$ rounds, and there are a total of $K = 3$ repetitions. In Repetition 1, Round 1, the user downloads 3 blocks of the desired computation $W^{(1)}$ from each server, i.e., $h_1^T W_{[1:3]}^{(1)}$ from server 1, $h_2^T W_{[4:6]}^{(1)}$ from server 2 and so on. To maintain privacy, the user needs to download equivalent amount of blocks of the undesired computation $W^{(2)}$. In order to utilize the undesired blocks, the user needs to be able to decode the undesired blocks. Hence, the user downloads $h_1^T W_1^{(2)}, h_2^T W_1^{(2)}$ and $h_3^T W_1^{(2)}$ from server 1, 2 and 3, respectively; downloads $h_4^T W_2^{(2)}, h_1^T W_2^{(2)}$ and $h_2^T W_2^{(2)}$ from server 4, 1 and 2 respectively, and so on, until all $W_{[1:4]}^{(2)}$ can be decoded (see Table 3).

In Round 2, we pair up a new block of $W^{(1)}$ and a undesired block of $W^{(2)}$ decoded from Round 1 and let the user download the sum of them. Since the user has not downloaded $W_4^{(2)}$ from server 1, the user can ask server 1 to send $h_1^T W_4^{(2)} + h_1^T W_{13}^{(1)}$. Similarly, the user can ask server 2 to send $h_2^T W_3^{(2)} + h_2^T W_{14}^{(1)}$ and so on. Clearly, the user can use $W_{[1:4]}^{(2)}$ from Round 1 as side information to obtain $h_1^T W_{13}^{(1)}, h_2^T W_{14}^{(1)}, h_3^T W_{15}^{(1)}$ and $h_4^T W_{16}^{(1)}$. In Repetition 1, the user ends up downloading all the blocks of $W^{(1)}$ once (12

		Server			
		1	2	3	4
Repetition 1	Round 1	$h_1 W_1^{(1)}$	$h_2 W_4^{(1)}$	$h_3 W_7^{(1)}$	$h_4 W_{10}^{(1)}$
		$h_1 W_2^{(1)}$	$h_2 W_5^{(1)}$	$h_3 W_8^{(1)}$	$h_4 W_{11}^{(1)}$
		$h_1 W_3^{(1)}$	$h_2 W_6^{(1)}$	$h_3 W_9^{(1)}$	$h_4 W_{12}^{(1)}$
		$h_1 W_1^{(2)}$	$h_2 W_1^{(2)}$	$h_3 W_1^{(2)}$	$h_4 W_2^{(2)}$
		$h_1 W_2^{(2)}$	$h_2 W_2^{(2)}$	$h_3 W_3^{(2)}$	$h_4 W_3^{(2)}$
		$h_1 W_3^{(2)}$	$h_2 W_4^{(2)}$	$h_3 W_4^{(2)}$	$h_4 W_4^{(2)}$
	Round 2	$h_1 W_{13}^{(1)}$ +	$h_2 W_{14}^{(1)}$ +	$h_3 W_{15}^{(1)}$ +	$h_4 W_{16}^{(1)}$ +
		$h_1 W_4^{(2)}$	$h_2 W_3^{(2)}$	$h_3 W_2^{(2)}$	$h_4 W_1^{(2)}$
	Round 1	$h_1 W_{10}^{(1)}$	$h_2 W_1^{(1)}$	$h_3 W_4^{(1)}$	$h_4 W_7^{(1)}$
		$h_1 W_{11}^{(1)}$	$h_2 W_2^{(1)}$	$h_3 W_5^{(1)}$	$h_4 W_8^{(1)}$
		$h_1 W_{12}^{(1)}$	$h_2 W_3^{(1)}$	$h_3 W_6^{(1)}$	$h_4 W_9^{(1)}$
		$h_1 W_5^{(2)}$	$h_2 W_5^{(2)}$	$h_3 W_5^{(2)}$	$h_4 W_6^{(2)}$
		$h_1 W_6^{(2)}$	$h_2 W_6^{(2)}$	$h_3 W_7^{(2)}$	$h_4 W_7^{(2)}$
		$h_1 W_7^{(2)}$	$h_2 W_8^{(2)}$	$h_3 W_8^{(2)}$	$h_4 W_8^{(2)}$
Repetition 2	Round 1	$h_1 W_{16}^{(1)}$ +	$h_2 W_{13}^{(1)}$ +	$h_3 W_{14}^{(1)}$ +	$h_4 W_{15}^{(1)}$ +
		$h_1 W_8^{(2)}$	$h_2 W_7^{(2)}$	$h_3 W_6^{(2)}$	$h_4 W_5^{(2)}$
	Round 2	$h_1 W_7^{(1)}$	$h_2 W_{10}^{(1)}$	$h_3 W_1^{(1)}$	$h_4 W_4^{(1)}$
		$h_1 W_8^{(1)}$	$h_2 W_{11}^{(1)}$	$h_3 W_2^{(1)}$	$h_4 W_5^{(1)}$
		$h_1 W_9^{(1)}$	$h_2 W_{12}^{(1)}$	$h_3 W_3^{(1)}$	$h_4 W_6^{(1)}$
		$h_1 W_9^{(2)}$	$h_2 W_9^{(2)}$	$h_3 W_9^{(2)}$	$h_4 W_{10}^{(2)}$
Repetition 3	Round 1	$h_1 W_{10}^{(2)}$	$h_2 W_{10}^{(2)}$	$h_3 W_{11}^{(2)}$	$h_4 W_{11}^{(2)}$
		$h_1 W_{11}^{(2)}$	$h_2 W_{12}^{(2)}$	$h_3 W_{12}^{(2)}$	$h_4 W_{12}^{(2)}$
	Round 2	$h_1 W_{15}^{(1)}$ +	$h_2 W_{16}^{(1)}$ +	$h_3 W_{13}^{(1)}$ +	$h_4 W_{14}^{(1)}$ +
		$h_1 W_{12}^{(2)}$	$h_2 W_{11}^{(2)}$	$h_3 W_{10}^{(2)}$	$h_4 W_9^{(2)}$

Fig. 3: Downloading strategy for $(N = 4, M = 2, K = 3)$.

from Round 1, and 4 from Round 2). In order to decode $W^{(1)}$, the user needs to download all 16 blocks two more times with different linear combinations. Thus, the user can follow the same downloading pattern with a right shift of indices in Repetition 2 and one more shift in Repetition 3 with new side information. The normalized upload and download costs of the scheme are given by $U = 4/(3 - 1) = 2$ and $D = (3 \times 28)/(2 \times 16) = 21/8$ per bit. We next describe the general scheme in the following section.

IV. PROOF OF THEOREM 1

We now present the achievable scheme, which can be broken into two phases: (1) secure upload; and (2) private download.

Phase 1 (Secure Upload): Assume that the user wants to compute $AB_\theta, \theta \in \{1, \dots, M\}$. For a given parameter K , the user first divides the input matrix A into $K - 1$ partitions vertically into $A = [A_1^T \dots A_{K-1}^T]^T$, where each A_i is of size $d_1/(K - 1) \times d_2$. The goal can be written as

$$AB_\theta = \begin{bmatrix} A_1 B_\theta \\ \vdots \\ A_{K-1} B_\theta \end{bmatrix}. \quad (14)$$

The confidential matrix A is encoded as follows:

$$\tilde{A}_n = \sum_{i=1}^{K-1} A_i x_n^{i-1} + R x_n^{K-1} \quad (15)$$

for server n , where R is a random matrix with the same dimension as any A_i , with uniformly distributed i.i.d. entries, and x_n is a distinct element in \mathbb{F} assigned to server n . The user uploads \tilde{A}_n to server n and instructs each server to multiply their respective \tilde{A}_n with all the B_m 's, i.e.,

$$\tilde{A}_n B_m = \sum_{i=1}^{K-1} A_i B_m x_n^{i-1} + R B_m x_n^{K-1}, \quad \forall m, \quad (16)$$

which can be written as $h_n^T W^{(m)}$, where

$$h_n = \begin{bmatrix} 1 \\ x_n \\ \vdots \\ x_n^{K-1} \end{bmatrix}, \quad W^{(m)} = \begin{bmatrix} A_1 B_m \\ \vdots \\ A_{K-1} B_m \\ R B_m \end{bmatrix}. \quad (17)$$

It is clear that total normalized upload cost is $U = N/(K - 1)$.

Phase 2 (Private Download): We use a similar downloading technique that was proposed in [14]. Thus, before the user downloads anything from servers, the user first asks each server n to partition the results $\tilde{A}_n B_m$ into blocks as follows

$$h_n^T W^{(m)} = \begin{bmatrix} h_n^T W_1^{(m)} \\ \vdots \\ h_n^T W_{NM}^{(m)} \end{bmatrix}, \quad \forall m, \quad (18)$$

where each block is of size $d_1/(N^M(K - 1)) \times d_3$. The user organizes the download into repetitions and rounds. Within each repetition, there are M rounds, and there are a total of K repetitions. Since the download pattern in each repetition is a cyclic shift of the previous repetition, we focus on describing the details of rounds.

In Round $i = 1$ of Repetition 1, the user downloads K^{M-1} desired blocks from each server, where the user downloads $h_n^T W_{(n-1)K^{M-1}+1}^{(\theta)}, \dots, h_n^T W_{nK^{M-1}}^{(\theta)}$ from server $n, \forall n$. A total of NK^{M-1} distinct desired blocks are downloaded in this step. However, in order to ensure privacy, the user needs to enforce message symmetry. In other words, the user needs to download equal number of blocks of each undesired computation $h_n^T W^{(m)}, m = 1, \dots, M, m \neq \theta$ from server n . While undesired, these blocks can be used as side information in the next round. To ensure that the undesired blocks can be used as side information, the user needs to be able to decode those undesired blocks, i.e., obtain $W_j^{(m)}$ for some j and $m \neq \theta$. Since each block is encoded using (N, K) MDS code, the user needs to download each $h_n^T W_j^{(m)}$ from at least K different servers, for a fixed $m \neq \theta$ and j . The user can choose

to download $h_n^T W_{j_p}^{(m)}$ from server $n = K(p-1) + 1, \dots, Kp \bmod N$ and $p = 1, \dots, NK^{M-2}$ for all $m \neq \theta$. Note that if $n \bmod N = 0$, we set that n to N . Since each $h_n^T W_{j_p}^{(m)}$ is downloaded from K servers, the user is able to decode $W_{j_p}^{(m)}$. Hence, the user can cancel the contribution of any $W_{j_p}^{(m)}$ in the subsequent rounds. In general, in any given Round i , the user downloads $h_{n'}^T W_{\ell_1}^{(m_1)} + h_{n'}^T W_{\ell_2}^{(m_2)} + \dots + h_{n'}^T W_{\ell_i}^{(m_i)}$, where $n' = n, \dots, n + K - 1$, $m_1, \dots, m_i \in \{1, \dots, M\} \setminus \theta$ and $\ell_1, \dots, \ell_i \in \{1, \dots, N^M\}$. Thus, the amount of side information available at the user is $N \binom{M-1}{i} K^{M-i-1} (N-K)^{i-1}$.

Since for a particular block of side information, there are $N-K$ servers who did not download that block, we can pair that block with a fresh desired block at those $N-K$ servers. Starting from Round $i+1$, the user downloads the sums of $i+1$ blocks from different computations with the help of side information from last round. The user downloads a fresh block of $h_n^T W^{(\theta)}$ and i blocks of side information decoded from last round, i.e., $h_{n_1}^T W_{r_1}^{(\theta)} + h_{n_1}^T W_{\ell_1}^{(m_1)} + \dots + h_{n_i}^T W_{\ell_i}^{(m_i)}$, where r_1 is the new block of $W^{(\theta)}$ from server n . This allows the user to obtain $N \binom{M-1}{i} K^{M-i-1} (N-K)^i$ fresh desired blocks.

The user can repeat this process K times except the indices for desired blocks are shifted right once in every repetition and fresh indices are chosen for side information in every repetition. Since each desired block is downloaded K times from K different source, the user can easily decode $W^{(\theta)}$, hence, recover AB_θ . This allows the user to download a total of $K \sum_{i=1}^{M-1} N \binom{M-1}{i} K^{M-i-1} (N-K)^{i-1}$ undesired blocks and a total of $K \sum_{i=0}^{M-1} N \binom{M-1}{i} K^{M-i-1} (N-K)^i$ desired blocks. Note that each desired block has K terms, however, only $K-1$ of those terms are actually useful due to the addition of the random matrix (for ensuring confidentiality of the matrix A). Hence, the download cost for this scheme is

$$\begin{aligned} D &= \frac{\text{total downloaded}}{\frac{K-1}{K} \times \text{total desired}} = \frac{K}{K-1} \left(1 + \frac{\text{total undesired}}{\text{total desired}} \right) \\ &= \frac{K}{K-1} \left(1 + \frac{K \sum_{i=1}^{M-1} N \binom{M-1}{i} K^{M-i-1} (N-K)^{i-1}}{K \sum_{i=0}^{M-1} N \binom{M-1}{i} K^{M-i-1} (N-K)^i} \right) \\ &= \frac{K}{K-1} \left(1 + \frac{\frac{K}{N-K} \sum_{i=1}^{M-1} \binom{M-1}{i} K^{M-1-i} (N-K)^i}{N^{M-1}} \right) \\ &= \frac{K}{K-1} \left(1 + \frac{\frac{K}{N-K} (N^{M-1} - K^{M-1})}{N^{M-1}} \right) \\ &= \frac{K}{K-1} \left(1 + \frac{K}{N-K} \left(1 - \left(\frac{K}{N} \right)^{M-1} \right) \right) \\ &= \frac{K}{K-1} \left(\frac{N - K \left(\frac{K}{N} \right)^{M-1}}{N-K} \right) = \frac{K}{K-1} \left(\frac{1 - \left(\frac{K}{N} \right)^M}{1 - \frac{K}{N}} \right), \end{aligned}$$

which leads to the expression of the normalized download cost D in Theorem 1. This completes the proof of Theorem 1.

V. CONCLUSIONS

In this paper, we studied the problem of secure and private matrix multiplication. We proposed a new scheme that com-

bines ideas from secret sharing and coded PIR. We showed that the lower convex hull of $(U, D) = (N/(K-1), (K/(K-1))(1 + (K/N) + \dots + (K/N)^{M-1}))$ for $K = 2, \dots, N$ is achievable. We also show that our scheme performs better than the scheme proposed in [12] in terms of upload and download cost. An interesting future direction is to obtain converse results for the tradeoff between upload/download cost for this problem.

REFERENCES

- [1] Q. Yu, M. A. Maddah-Ali, and A. S. Avestimehr, "Straggler Mitigation in Distributed Matrix Multiplication: Fundamental Limits and Optimal Coding," *CoRR*, vol. abs/1801.07487, 2018. [Online]. Available: <http://arxiv.org/abs/1801.07487>
- [2] S. Dutta, M. Fahim, F. Haddadpour, H. Jeong, V. R. Cadambe, and P. Grover, "On the optimal recovery threshold of coded matrix multiplication," *CoRR*, vol. abs/1801.10292, 2018. [Online]. Available: <http://arxiv.org/abs/1801.10292>
- [3] R. Bitar, P. Parag, and S. E. Rouayheb, "Minimizing latency for secure distributed computing," in *2017 IEEE International Symposium on Information Theory (ISIT)*, Jun. 2017, pp. 2900–2904.
- [4] H. A. Nodehi and M. A. Maddah-Ali, "Limited-sharing multi-party computation for massive matrix operations," in *2018 IEEE International Symposium on Information Theory (ISIT)*, Jun. 2018, pp. 1231–1235.
- [5] Q. Yu, N. Raviv, J. So, and A. S. Avestimehr, "Lagrange Coded Computing: Optimal Design for Resiliency, Security and Privacy," *CoRR*, vol. abs/1806.00939, 2018. [Online]. Available: <http://arxiv.org/abs/1806.00939>
- [6] W. Chang and R. Tandon, "On the capacity of secure distributed matrix multiplication," in *2018 IEEE Global Communications Conference (GLOBECOM)*, Dec. 2018, pp. 1–6.
- [7] H. Yang and J. Lee, "Secure distributed computing with straggling servers using polynomial codes," *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 1, pp. 141–150, Jan 2019.
- [8] J. Kakar, S. Ebadifar, and A. Sezgin, "Rate-efficiency and straggler-robustness through partition in distributed two-sided secure matrix computation," *CoRR*, vol. abs/1810.13006, 2018. [Online]. Available: <http://arxiv.org/abs/1810.13006>
- [9] R. G. L. D'Oliveira, S. E. Rouayheb, and D. A. Karpuk, "GASP codes for secure distributed matrix multiplication," *CoRR*, vol. abs/1812.09962, 2018. [Online]. Available: <http://arxiv.org/abs/1812.09962>
- [10] M. Aliasgari, O. Simeone, and J. Kliewer, "Distributed and private coded matrix computation with flexible communication load," *CoRR*, vol. abs/1901.07705, 2019. [Online]. Available: <http://arxiv.org/abs/1901.07705>
- [11] W. Chang and R. Tandon, "On the Capacity of Secure Distributed Fast Fourier Transform," in *2018 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, Nov. 2018, pp. 653–657.
- [12] M. Kim and J. Lee, "Private secure coded computation," *CoRR*, vol. abs/1902.00167, 2019. [Online]. Available: <http://arxiv.org/abs/1902.00167>
- [13] H. Sun and S. A. Jafar, "The capacity of private information retrieval," *IEEE Transactions on Information Theory*, vol. 63, no. 7, pp. 4075–4088, Jul. 2017.
- [14] K. Banawan and S. Ulukus, "The capacity of private information retrieval from coded databases," *IEEE Transactions on Information Theory*, vol. 64, no. 3, pp. 1945–1956, March 2018.
- [15] H. Sun and S. A. Jafar, "The capacity of private computation," *CoRR*, vol. abs/1710.11098, 2017. [Online]. Available: <http://arxiv.org/abs/1710.11098>
- [16] M. Mirmohseni and M. A. Maddah-Ali, "Private function retrieval," in *2018 Iran Workshop on Communication and Information Theory (IWCIT)*, Apr. 2018, pp. 1–6.
- [17] D. Karpuk, "Private computation of systematically encoded data with colluding servers," in *2018 IEEE International Symposium on Information Theory (ISIT)*, Jun. 2018, pp. 2112–2116.
- [18] S. A. Obead and J. Kliewer, "Achievable rate of private function retrieval from MDS coded databases," in *2018 IEEE International Symposium on Information Theory (ISIT)*, Jun. 2018, pp. 2117–2121.