

Received November 18, 2019, accepted January 29, 2020, date of publication February 3, 2020, date of current version February 10, 2020.

Digital Object Identifier 10.1109/ACCESS.2020.2971036

High-Performance Homomorphic Matrix Completion on Multiple GPUs

TAO ZHANG¹, (Member, IEEE), HAN LU¹, AND XIAO-YANG LIU², (Student Member, IEEE)

¹School of Computer Engineering and Science, Shanghai University, Shanghai 200444, China

²Department of Electrical Engineering, Columbia University, New York, NY 10027, USA

Corresponding author: Xiao-Yang Liu (xl2427@columbia.edu)

The work of Tao Zhang was supported in part by the Natural Science Foundation of Shanghai under Grant 17ZR1409800.

ABSTRACT In various applications such as trajectory tracking in mobile social networks and online recommendation systems, the massive raw data are often incomplete due to various unpredictable or unavoidable reasons. Matrix completion algorithms are effective for reconstructing two-dimensional data, but sending raw data containing personal, sensitive information to cloud computing nodes for matrix completion may lead to privacy exposure issue. The homomorphic matrix completion is a promising approach to perform matrix completion while preserving privacy. However, CPU-based homomorphic matrix completion has low performance, making it impractical to process multiple or large-scale data completion tasks in real-time. In this paper, we propose a high-performance homomorphic matrix completion scheme by exploiting commodity GPUs (Graphics Processing Units) that are widely available in HPC servers and cloud computing nodes. First, we design and implement a baseline GPU-based homomorphic matrix completion, and propose techniques to optimize memory accesses, GPU utilization, and communications. Second, we propose a shard mode for large-scale matrices exceeding GPU memory capacity. Third, we propose a multi-GPU mode to fully utilize multiple GPUs in computing nodes. Experiment results show that the proposed scheme is both fast and accurate. On matrices of varying sizes, the proposed scheme running on a single Tesla V100 GPU achieves up to $116.23\times$ speedups over the CPU MATLAB implementation running on dual Xeon CPUs. The multi-GPU mode achieves up to $1.84\times$ speedups on two GPUs versus on a single GPU. For large-scale matrices, the shard mode achieves up to $174.92\times$ speedups on a single GPU over the CPU MATLAB implementation on two CPUs, and further achieves up to $1.35\times$ speedups when running on two GPUs using the multi-GPU mode.

INDEX TERMS GPU, homomorphic matrix completion, least squares minimization.

I. INTRODUCTION

In many applications such as video and image processing [2], out-door and in-door localization [3], [4], and recommendation systems [5], the massive data are often incomplete owing to various reasons in data acquisition and transmission [6]. For instance, the UC Berkeley INTEL project [7], [8] reported 40% data missings and 8% data errors. Many scientific researches such as geoexploration [9] depend on complete data to draw accurate conclusion, since data analysis on incomplete data leads to inaccurate and even wrong conclusions. Therefore, various data completion methods [8], [10], [11] have been designed to recover the incomplete data. For two-dimensional data (i.e., matrix),

The associate editor coordinating the review of this manuscript and approving it for publication was Tomás F. Peña¹.

matrix completion approaches [8] [11] [12] [13] were proposed to recover the missing elements by exploiting the latent low-rank property of matrices. Alternating minimization [7] [11] [14] is an effective algorithm to address the matrix completion problem. Matrix completion approaches have been widely adopted in diverse applications including link prediction in social networks [15], image recovery [16], image classification [17], and action detection [18].

Existing matrix completion approaches have limitations in the cloud computing paradigm. Plain-data (i.e., not encrypted) based matrix completion on cloud computing nodes may have privacy issues. In cloud computing, users send requests with plain data from their mobile devices to cloud computing nodes, then cloud computing nodes execute the matrix completion algorithms, and finally responses are sent back from cloud computing nodes to users [4] [19].

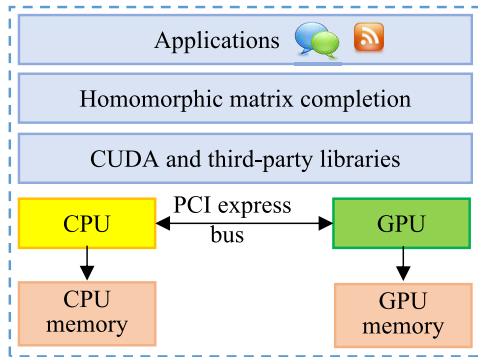


FIGURE 1. Overview of the homomorphic matrix completion on GPU.

In this process, users' sensitive information may be exposed to cloud computing nodes, such as locations and trajectories, preferences, and website navigation histories. Privacy exposure may lead to economic losses and even lawsuits. For example, the recommendation contest Netflix prize was canceled after the privacy lawsuit [20]. To address the privacy issue, the homomorphic matrix completion approach [19] was proposed to perform matrix completion while preserving privacy. Homomorphic matrix completion [19] adopts a three-step procedure and conduct matrix completion on encrypted data instead of plain data to preserve privacy. However, the major computation of the homomorphic matrix completion is tens of iterative least square minimizations, whose running time grows rapidly with the matrix size. As a result, it is impractical for a cloud computing node to process multiple or large-scale homomorphic matrix completion tasks in real-time.

In this paper, we propose high-performance and large-scale homomorphic matrix completion on GPU, as shown in Fig. 1. We design, implement, and optimize fast and accurate homomorphic matrix completion on single and multiple GPUs. First, we design and implement a baseline homomorphic matrix completion scheme on GPU. This baseline scheme is full functional, but its performance is unoptimized. We then investigate the bottlenecks of the baseline scheme, and optimize the memory accesses, GPU utilizations, and CPU-GPU communications to obtain an optimized scheme. Second, we design a multi-GPU mode to distribute the computations among multiple GPUs on a computing node. Third, we propose a shard mode to complete large-scale matrices exceeding the GPU memory capacity.

Our contributions are summarized as follows.

- We design, implement and optimize a GPU-based homomorphic matrix completion scheme to achieve high performance and accuracy. We propose techniques to optimize memory accesses, GPU utilizations, and CPU-GPU communications. These optimization techniques improve performance significantly.
- We propose a multi-GPU mode to fully utilize multiple GPUs in cloud computing nodes. This mode achieves up to $1.84\times$ speedups on two GPUs versus on a single GPU.

We further propose a shard mode to complete large-scale matrices exceeding the GPU memory capacity.

- We perform extensive experiments to evaluate the performance of the GPU-based homomorphic matrix completion. For small- or medium-sized matrices, the proposed scheme running on a Tesla V100 GPU achieves up to $116.23\times$ speedup versus the CPU MATLAB implementation [19] running on dual Xeon CPUs. With the multi-GPU mode, the proposed scheme achieves up to $1.84\times$ speedups on two GPUs versus on a single GPU. With the shard mode, the proposed scheme achieves up to $174.92\times$ speedups on a Tesla V100 GPU versus the CPU MATLAB implementation [19] running on dual Xeon CPUs for large matrices, and further achieves up to $1.35\times$ speedups when running on two GPUs using the multi-GPU mode. The proposed scheme achieves similar recovery error with the CPU MATLAB implementation [19]. The project is available at: https://github.com/hust512/Homomorphic_Matrix_Completion_Multiple_GPU.

The remainder of this paper is organized as follows. In Section II, we describe the notations and the homomorphic matrix completion algorithm [19]. Section III presents the design, implementation and optimization of the GPU-based homomorphic matrix completion scheme. Section IV presents the shard mode and the multi-GPU mode. In Section V, we evaluate the performance of the proposed scheme. Section VI discusses the related works. The conclusions are drawn in Section VII.

II. THE HOMOMORPHIC MATRIX COMPLETION ALGORITHM

We introduce the notations, summarize the homomorphic matrix completion algorithm [19], and provide parallel acceleration analysis for this algorithm.

A. NOTATIONS

We use lowercase and uppercase boldface letters to denote vectors (e.g., $\mathbf{x} \in \mathbb{R}^J$) and matrices (e.g., $\mathbf{X} \in \mathbb{R}^{I \times J}$), respectively. We use i, j to index the rows and columns of a matrix, respectively. We use $[I]$ to denote the set $\{1, 2, \dots, I\}$. For a matrix $\mathbf{X} \in \mathbb{R}^{I \times J}$, \mathbf{X}_j or $\mathbf{X}(:, j)$ denotes the j -th column of \mathbf{X} , the (i, j) -th element is $X(i, j)$ or simply X_{ij} , where $i \in [I]$ and $j \in [J]$. The transposed matrix of \mathbf{X} is denoted as \mathbf{X}^T . We use \mathbf{X}^i to denote the i -th matrix in a series of matrices.

1) FROBENIUS NORM OF A MATRIX

The Frobenius norm of a matrix \mathbf{X} is $\|\mathbf{X}\|_F = \sqrt{\sum_{i=1}^I \sum_{j=1}^J |X_{ij}|^2}$.

2) INDEX MATRIX

The index matrix Φ indicates whether an element $X_{ij} \in \mathbf{X}$ is missing. It is defined as follows:

$$\Phi_{ij} = \begin{cases} 0, & \text{if } X_{ij} \text{ is missing,} \\ 1, & \text{otherwise.} \end{cases}$$

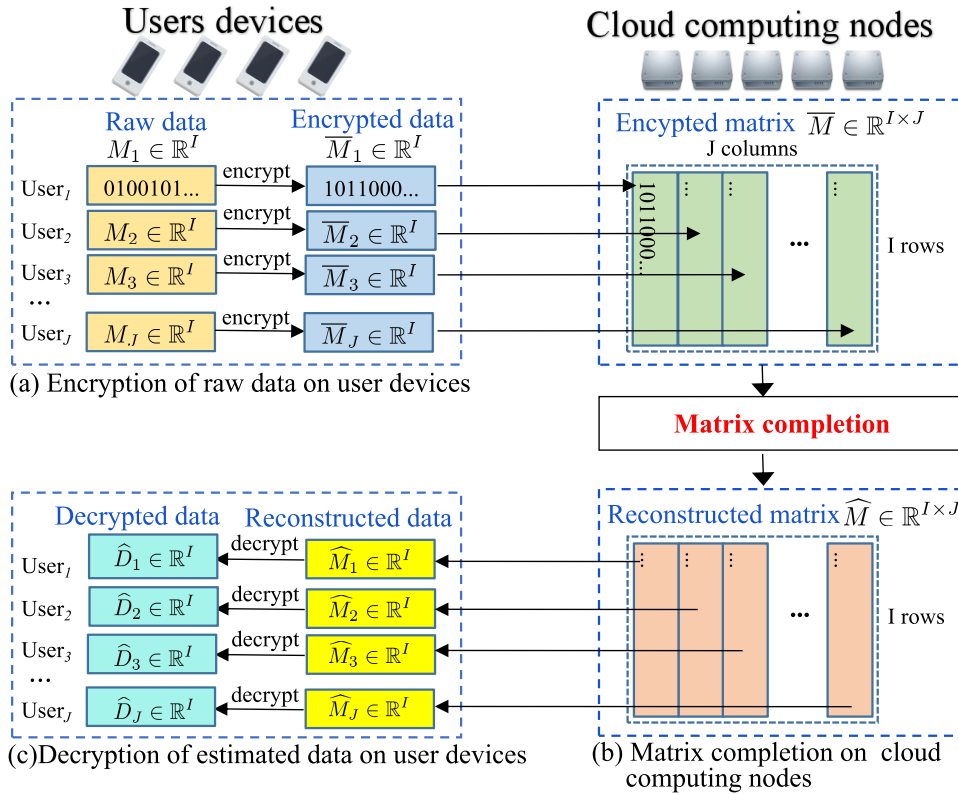


FIGURE 2. Overview of the homomorphic matrix completion algorithm [19].

Algorithm 1 Homomorphic Matrix Completion Algorithm on User Devices j , for $j \in [J]$

Input: The raw vector $M_j \in \mathbb{R}^I$, index vector $\Phi_j \in \mathbb{R}^I$.

- 1) **Initialize:** Randomly generate $\psi_0, \psi_1, \dots, \psi_K \in [0, 1]$ and $\sum_{i=0}^K \psi_i = 1$,
- 2) Receive public matrix P from the cloud server (line 2 in Alg. 2),
- 3) Encrypt M_j to \bar{M}_j through $\bar{M}_j = (\psi_0 M_j + \psi_1 P_1 + \dots + \psi_K P_K) \circ \Phi_j$,
- 4) Upload the encrypted vector \bar{M}_j to the cloud server,
- 5) Receive the recovered vector \hat{M}_j from the cloud server (line 8 in Alg. 2),
- 6) Decrypt \hat{M}_j to obtain \hat{D}_j through $\hat{D}_j = (\hat{M}_j - \psi_1 P_1 - \dots - \psi_K P_K) / \psi_0$,

Output: \hat{D}_j .

3) HADAMARD PRODUCT (ELEMENT-WISE PRODUCT)

The Hadamard product Z of $X \in \mathbb{R}^{I \times J}$ and $Y \in \mathbb{R}^{I \times J}$ is $Z = X \circ Y$, where $Z_{ij} = X_{ij} Y_{ij}$.

4) RAW MATRIX

The raw (incomplete) matrix M can be represented as the Hadamard product of the true (complete) matrix $D \in \mathbb{R}^{I \times J}$ and the index matrix $\Phi \in \mathbb{R}^{I \times J}$, i.e., $M = D \circ \Phi \in \mathbb{R}^{I \times J}$.

B. OVERVIEW OF THE HOMOMORPHIC MATRIX COMPLETION ALGORITHM

The homomorphic matrix completion algorithm [19] consists of three steps, as shown in Fig. 2:

- 1) Encryption of raw data on user devices as in Fig. 2(a): each user uses the public matrix $P \in \mathbb{R}^{I \times K}$ from cloud computing nodes to encrypt her data vector on user devices such as mobile phones. These encrypted vectors are sent to a cloud computing node and form an encrypted matrix.
- 2) Matrix completion on cloud computing nodes as in Fig. 2(b): the encrypted matrix is recovered targeting a low rank R using the matrix completion algorithm on the cloud computing node.
- 3) Decryption of reconstructed data on user devices as in Fig. 2(c): the reconstructed matrix \hat{M} is returned to users. Each user obtains her own reconstructed data vector by decrypting her column in the reconstructed matrix using her private keys on a user device.

In the first and third steps, homomorphic encryption and decryption are used to encrypt the raw (incomplete) data and decrypt the reconstructed data, so as to protect data privacy. In the second step, the alternating minimization algorithm [7] [11] is utilized for matrix completion. The parts of the homomorphic matrix completion algorithm on user devices and on cloud computing nodes are described in Algs. 1 and 2, respectively.

Algorithm 2 Homomorphic Matrix Completion Algorithm on Cloud Computing Nodes

Input: Index matrix $\Phi \in \mathbb{R}^{I \times J}$, target low rank R , number of iterations L .

- 1: **Initialize:** Public matrix $P \in \mathbb{R}^{I \times K}$, matrix $X^0 \in \mathbb{R}^{I \times R}$,
- 2: Broadcast P to all J users,
- 3: Gather all J encrypted vectors and form the encrypted matrix $\bar{M} \in \mathbb{R}^{I \times J}$,
- 4: **for** $\ell = 1$ to L **do**
- 5: $Y^\ell = \arg \min_{Y^\ell \in \mathbb{R}^{R \times J}} \|\bar{M} - \Phi \circ (X^{\ell-1} Y^\ell)\|_F^2$,
- 6: $X^\ell = \arg \min_{X^\ell \in \mathbb{R}^{I \times R}} \|\bar{M} - \Phi \circ (X^\ell Y^\ell)\|_F^2$,
- 7: **end for**
- 8: Compute the reconstructed matrix $\hat{M} = X^L Y^L$,
- 9: Broadcast the reconstructed vector \hat{M}_j to the j -th user, for $j \in [J]$,

Output: \hat{M} .

1) ENCRYPTION OF RAW DATA

The encryption of raw data on user devices is shown in lines 2–4 of Alg. 1. Each user uses the public matrix $P \in \mathbb{R}^{I \times K}$ from a cloud computing node to encrypt her own raw data vector $M_1 \in \mathbb{R}^I$ and obtain $\bar{M}_1 \in \mathbb{R}^I$ on a user device. Since user devices only have limited computation power, memory capacity and power budget, the homomorphic matrix completion algorithm [19] utilizes a light-weight encryption scheme as follows:

$$\bar{M}_j = (\psi_0 M_j + \psi_1 P_1 + \dots + \psi_K P_K) \circ \Phi_j \in \mathbb{R}^I, \quad (1)$$

where $\psi_0, \psi_1, \dots, \psi_K \in [0, 1]$ are private keys generated on user devices and $\sum_{i=0}^K \psi_i = 1$. Each user generates her own private keys ψ and encrypt her own data vector. Therefore, user j encrypts a raw data vector $M_j \in \mathbb{R}^I$ and produces an encrypted data vector $\bar{M}_j \in \mathbb{R}^I$. The encrypted data vectors \bar{M}_j from J users are sent to a cloud computing node and form an encrypted matrix $\bar{M} \in \mathbb{R}^{I \times J}$, where each encrypted data vector becomes a column of the encrypted matrix \bar{M} . Let K denote the rank of the public matrix P , and $R - K$ denote the rank of the raw (incomplete) matrix before encryption, then the rank of the encrypted matrix \bar{M} is denoted by R .

2) MATRIX COMPLETION ON CLOUD COMPUTING NODES

The matrix completion on cloud computing nodes is shown in Alg. 2. First, the cloud computing node generates a public matrix $P \in \mathbb{R}^{I \times K}$ and an initial matrix $X^0 \in \mathbb{R}^{I \times R}$ (in line 1), where R is the target rank of the encrypted matrix. The public matrix P is constructed by selecting K public vectors from all available public vectors or generated randomly. Public vectors are subsidiary data that are complete and without encryption, which are provided by some users [19]. The initial X^0 is generated randomly for the specified size. Second, the cloud computing node broadcasts the public matrix P to all J users, then gathers all J encrypted vectors and form an encrypted matrix $\bar{M} \in \mathbb{R}^{I \times J}$ (in lines 2–3). Third, the cloud computing node performs the compute-intensive matrix completion (in lines 4–8). We adopt the alternating minimization

algorithm [7] [11] [14] to solve matrix completion. In each iteration, we fix matrix $X \in \mathbb{R}^{I \times R}$ to solve $Y \in \mathbb{R}^{R \times J}$ using (2) and then fix Y to solve X using (3). The core operation is least squares minimization. The reconstructed matrix $\hat{M} \in \mathbb{R}^{I \times J}$ is computed as the matrix multiplication of X^L and Y^L (line 8).

$$Y = \arg \min_{Y \in \mathbb{R}^{R \times J}} \|\bar{M} - \Phi \circ (XY)\|_F^2, \quad (2)$$

$$X = \arg \min_{X \in \mathbb{R}^{I \times R}} \|\bar{M} - \Phi \circ (XY)\|_F^2. \quad (3)$$

Finally, the cloud computing node broadcast the reconstructed vector \hat{M}_j to the j -th user, for $j \in [J]$.

3) DECRYPTION OF RECONSTRUCTED DATA

The decryption of reconstructed data on user devices is described in line 6 of Alg. 1. The homomorphic matrix completion algorithm [19] utilizes a light-weight decryption scheme as follows:

$$\hat{D}_j = (\hat{M}_j - \psi_1 P_1 - \dots - \psi_K P_K) / \psi_0 \in \mathbb{R}^I. \quad (4)$$

Each user utilizes her own private keys $\psi_0, \psi_1, \dots, \psi_K$ to perform decryption on her reconstructed vector \hat{M}_j to obtain decrypted data column \hat{D}_j .

C. PARALLEL ACCELERATION ANALYSIS

We profiled the CPU MATLAB implementation of the homomorphic matrix completion algorithm [19] and found that the second step took more than 95% of the total running time for most matrix sizes. Therefore, the efficiency of the algorithm can be improved by accelerating the second step.

In Alg. 2, the algorithm solves Y^ℓ and X^ℓ alternately within the for-loop (lines 4–7), then computes the reconstructed matrix (line 8). This process involves intensive computations, and we propose to offload these computations from CPU to high-performance GPU. GPUs are effective for accelerating massively parallel computations such as matrix operations. The computations in lines 4–8 involve matrix Hadamard product, matrix multiplication, matrix transpose, and least square minimization, which can be accelerated on GPU. Thus, we focus on the design, implementation and optimizations of the matrix completion (lines 4–8) implementation on GPU.

III. EFFICIENT HOMOMORPHIC MATRIX COMPLETION ON GPU

We parallelize the homomorphic matrix completion algorithm [19] in Alg. 2 on GPU. Then we propose techniques to optimize memory accesses, GPU utilizations and CPU-GPU communications.

A. DESIGN AND IMPLEMENTATION OF THE BASELINE GPU HOMOMORPHIC MATRIX COMPLETION

1) DATA STORAGE

We store the data in the column-major format because of two reasons. First, the homomorphic matrix completion algorithm

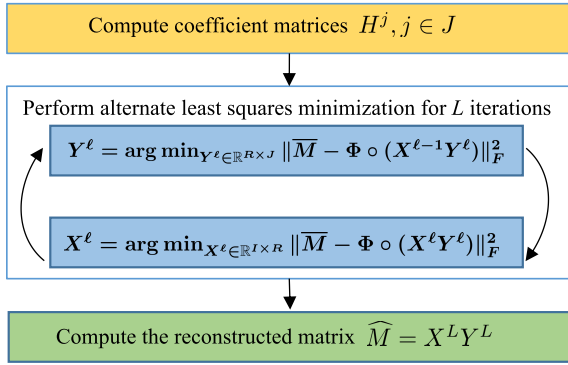


FIGURE 3. The three-step baseline of homomorphic matrix completion on GPU.

in Alg. 2 always accesses each column of matrices, therefore the column-major storage guarantees continuous memory accesses. Second, the GPU-based homomorphic matrix completion utilizes some routines in the CUDA libraries that use the column-major format.

2) PARALLELIZATION OF THE HOMOMORPHIC MATRIX COMPLETION

In Alg. 2, the major computations are the for-loop in lines 4–7 for alternate least square minimizations and the matrix multiplications in line 8. Lines 5 and 6 inside the for-loop perform the least squares minimizations in (2) and (3), respectively. Since the processes to obtain Y (line 5) and X (line 6) are quite similar, we only discuss the process to obtain Y in the following sections for conciseness. We use the following equation to solve the least squares minimization in (2) and get a column Y_j of Y :

$$Y_j = \arg \min_{Y_j \in \mathbb{R}^R} \|\bar{M}_j - H^j Y_j\|_F^2, \quad (5)$$

where $H^j \in \mathbb{R}^{I \times R}$ is the j -th matrix in a series of coefficient matrices and R is the target rank of the encrypted matrix. H^j is calculated as the Hadamard product of Φ_j and each column of matrix X as follows:

$$H^j(i, r) = \Phi(i, j) \circ X(i, r). \quad (6)$$

We compute lines 4–8 in Alg. 2 on GPU in the following three steps, as shown in Fig. 3:

- 1) Computing coefficient matrices H^j : Compute J coefficient matrices H^1, H^2, \dots, H^J using (6). In the baseline GPU implementation, we compute these J coefficient matrices in parallel on GPU. We allocate multiple GPU threads and use one GPU thread to calculate each $H^j \in \mathbb{R}^{I \times R}$.
- 2) Using (5) to solve the least squares minimizations in lines 5–6 of Alg. 2. We perform QR factorization on each H^j to obtain an orthogonal matrix and an upper triangular matrix. Therefore, a single least squares minimization is solved and a column Y_j of Y is obtained. For an encrypted matrix $\bar{M} \in \mathbb{R}^{I \times J}$, the algorithm computes J least squares minimizations for Y and J

least squares minimizations for X in each iteration. Therefore, the algorithm perform $(I + J)L$ least squares minimizations for L iterations.

- 3) Performing matrix multiplication of X^L and Y^L (in line 8 of of Alg. 2) to obtain the reconstructed matrix \hat{M} . We utilize the high-performance matrix multiplication routine in the cuBLAS library [21].

B. OPTIMIZATIONS

We performed preliminary experiments to evaluate the baseline GPU homomorphic matrix completion scheme, and found that it achieved merely similar performance with the CPU MATLAB implementation. We profiled the baseline scheme using the profiler in CUDA and observed three major bottlenecks: low memory access efficiency and big memory footprint, low utilization of GPU compute units, and high communication cost between CPU and GPU. Therefore, we optimize the baseline scheme on these aspects.

1) OPTIMIZING MEMORY ACCESSES

We optimize the memory access efficiency in GPU-based computation. The first step of lines 4–8 in Alg. 2 is to compute a series of coefficient matrices H^j using (6). In the baseline GPU implementation, Φ and X are stored in the GPU global memory which has the largest capacity on GPU. During the computing, each column of Φ in the GPU global memory is accessed J times. However, the access latency of the GPU global memory is high and introduce high time cost. To improve memory access efficiency, we store Φ into low-latency, small shared memory inside each GPU streaming multiprocessor (SM). To compute J coefficient matrices, we launch J thread blocks and copy Φ_j from the GPU global memory to the shared memory on j -th block. In this way, the algorithm accesses the shared memory J times to calculate H^j , which is much faster than accessing the GPU global memory.

We reduce the memory footprint in GPU computation. The matrix form of (5) is $H^j Y_j = \bar{M}_j$, where $H^j \in \mathbb{R}^{I \times R}$ and $\bar{M}_j \in \mathbb{R}^{I \times 1}$. To solve all J columns of Y , the space complexity is $O(IRJ)$, which increases rapidly as the grow of encrypted matrix size $I \times J$ and imposes pressure on limited GPU memory capacity. We multiply the transposed matrix $(H^j)^T$ to the matrix form of (5) and derive the following equation:

$$H^j Y_j = \bar{M}_j \Rightarrow (H^j)^T H^j Y_j = (H^j)^T \bar{M}_j. \quad (7)$$

In (7), $(H^j)^T H^j$ is $R \times R$ and $(H^j)^T \bar{M}_j$ is $R \times 1$, so the space complexity becomes $O(R^2 J)$ for all J columns of Y . Since the encrypted matrix \bar{M} holds the low-rank property, the target rank R is much smaller than I and J . Therefore, the memory consumption is significantly reduced. Since we store (H^j) in the column major format, we can obtain $(H^j)^T$ through fetch the elements of (H^j) without explicit matrix transpose.

2) IMPROVING GPU UTILIZATION

In (7), there are two matrix multiplications to calculate $(H^j)^T H^j$ and $(H^j)^T \bar{M}_j$, and a least square minimization.

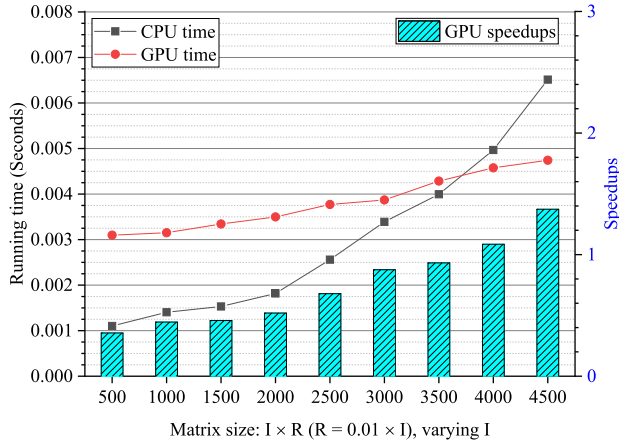


FIGURE 4. The performance of a single least squares minimization on two CPUs and the Tesla V100 GPU, respectively.

So there are $2J$ matrix multiplications and J least square minimizations to calculate all J columns of Y . Similarly, there are multiple matrix multiplications and least square minimizations to obtain X . In regular GPU implementation, these matrix multiplications or least square minimizations are computed one by one on GPU. However, since GPUs have thousands of cores, a single matrix multiplication or least squares minimization is often inadequate to fully utilize all GPU cores. We performed preliminary experiment to evaluate the performance of a single least squares minimization in (5). The performance under varying matrix size of $I \times R$ on two Xeon CPUs or a Tesla V100 GPU is shown in Fig. 4. The least squares minimization on the Tesla V100 GPU poorly achieved an average of $0.75\times$ and up to $1.37\times$ speedups versus two Xeon CPUs, due to the fact that a single least squares minimization is inadequate to fully utilize all GPU resources.

We improve GPU utilization by organizing and batching multiple matrix computation or the least squares minimizations to run parallelly on GPU. Since there are no dependency between these computations, we organize and compute multiple matrix computation or the least squares minimizations in batches. We compute $(H^j)^T H^j$, $j \in [J]$ in a batch and $(H^j)^T \bar{M}_j$, $j \in [J]$ in another batch by exploiting the batched matrix multiplication routine in the cuBLAS library [21]. After these batched matrix multiplications, we obtain J pairs of data from $[(H^1)^T H^1, (H^1)^T \bar{M}_1]$ to $[(H^J)^T H^J, (H^J)^T \bar{M}_J]$. Each pair of data is used for solving a single least squares minimization in (7). We organize and solve these J least squares minimizations parallelly in a batched by exploiting the batched Lower-Upper (LU) decomposition routine in the MAGMA library [22], as illustrated in Fig. 5.

The pseudo code of the optimized GPU matrix completion algorithm is described in Alg. 3. The most outer for-loop in lines 2–15 executes the matrix completion for L iterations. There are two parfor-loop in lines 3–6 and 8–11, in which the computations are in parallel (i.e., batched). The function `getCoefMatrix(·)` in lines 4 and 9 calculates a coefficient

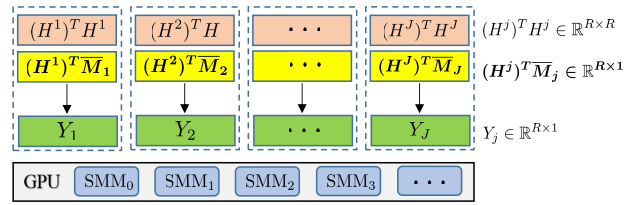


FIGURE 5. Batched least squares minimizations.

Algorithm 3 Pseudo code of the optimized GPU matrix completion

Input: Encrypted matrix $\bar{M} \in \mathbb{R}^{I \times J}$, target rank R , index matrix $\Phi \in \mathbb{R}^{I \times J}$, maximum iterations L .

- 1: **Initialize:** Matrix $X^0 \in \mathbb{R}^{I \times R}$,
- 2: **for** $\ell = 1$ to L **do**
- 3: **parfor** $j = 1$ to J **do**
- 4: $H^j = \text{getCoefMatrix}(X^{\ell-1}, \Phi_j)$,
- 5: $Y_j^\ell \leftarrow \text{LSM}((H^j)^T H^j Y_j^\ell = (H^j)^T \bar{M}_j)$,
- 6: **end parfor**
- 7: $(Y^\ell)^T = \text{Transpose}(Y^\ell)$,
- 8: **parfor** $i = 1$ to I **do**
- 9: $U^i = \text{getCoefMatrix}((Y^\ell)^T, \Phi_i^T)$,
- 10: $(X^\ell)_i^T \leftarrow \text{LSM}((U^i)^T U^i (X^\ell)_i^T = (U^i)^T \bar{M}_i^T)$,
- 11: **end parfor**
- 12: $X^\ell = \text{Transpose}((X^\ell)^T)$,
- 13: **end for**
- 14: $\hat{M} = X^L Y^L$,

Output: \hat{M} .

matrix H^j , and all J coefficient matrices are calculated parallelly on GPU. The function `LSM(·)` in lines 5 and 10 is the least squares minimization function. The `transpose(·)` function in lines 7 and 12 performs matrix transposition.

3) REDUCING COMMUNICATIONS

The profiling result on the baseline GPU implementation showed a high communication cost between CPU and GPU due to a large amount of data transfers. In the optimized GPU matrix completion implementation in Alg. 3, the communication could be a bottleneck if not optimized. We reduce communication cost by eliminating data transfers of intermediate results. We transfer X^0 , Φ and \bar{M} from CPU memory to GPU memory and continuously calculate lines 4–6 of Alg. 3 on the GPU without transferring intermediate results back to the CPU memory. Similar method is applied to lines 10–12.

IV. LARGE-SCALE AND MULTI-GPU HOMOMORPHIC MATRIX COMPLETION

A. LARGE-SCALE HOMOMORPHIC MATRIX COMPLETION ON A SINGLE GPU

Although external GPUs have much higher computation capability than integrated GPUs, their GPU memory capacity is limited. For large matrices, the memory required for matrix completion exceeds GPU memory capacity, limiting

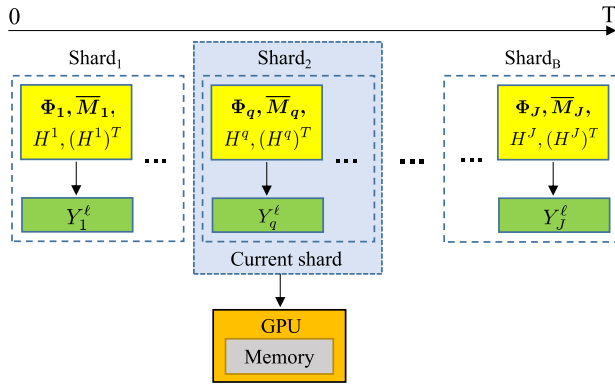


FIGURE 6. Shard mode for large-scale homomorphic matrix completion.

the applicability of external GPUs to small- or medium-sized matrices. By exploiting the separability of the matrix completion algorithm, we design a shard mode to perform large-scale homomorphic matrix completion on a single GPU.

1) SHARDS

As shown in Alg. 3, the parfor-loop in lines 3–6 to solve all columns of Y^ℓ and the parfor-loop in lines 8–11 to solve all columns of $(X^\ell)^T$ consume most of the GPU global memory during computing. Take Y^ℓ as example, J coefficient matrices $H^j \in \mathbb{R}^{I \times R}$ are computed parallelly (in line 4). Then there are two batched multiplications to compute J parallel $(H^j)^T H^j$ and $(H^j)^T \bar{M}_j$, and a batched least squares minimization to obtain all columns of Y^ℓ (in line 5). Since the computations to obtain different column of Y^ℓ are independent of each other (i.e., computing Y_i^ℓ is independent of computing Y_j^ℓ , for $i \neq j$). Therefore, we split the computations and data into shards (i.e., partitions) and each partition computes several columns of Y^ℓ or $(X^\ell)^T$. The shards to compute columns of Y^ℓ are called Y^ℓ shards while those to compute columns of X^ℓ are called X^ℓ shards. Fig. 6 shows example Y^ℓ shards.

Suppose the GPU has S bytes of available memory and the required memory for computing Y^ℓ or $(X^\ell)^T$ on an $I \times J$ encrypted matrix is a function of I and J denoted by $f(I, J)$, then the number of Y^ℓ shards or X^ℓ shards B is determined as $B = f(I, J)/S$.

2) LARGE-SCALE HOMOMORPHIC MATRIX COMPUTING WITH THE SHARD MODE

We propose a shard mode to perform large-scale homomorphic matrix completion exceeding the GPU memory capacity. Using this mode, we perform the computation in the ℓ -th iteration (in lines 2–13 of Alg.3) as follows:

- Computing Y^ℓ shards one by one on the GPU to obtain Y^ℓ , as shown in Fig. 6;
- Transposing Y^ℓ to obtain $(Y^\ell)^T$;
- Computing X^ℓ shards one by one on the GPU to obtain $(X^\ell)^T$;
- Transposing $(X^\ell)^T$ to obtain X^ℓ ;

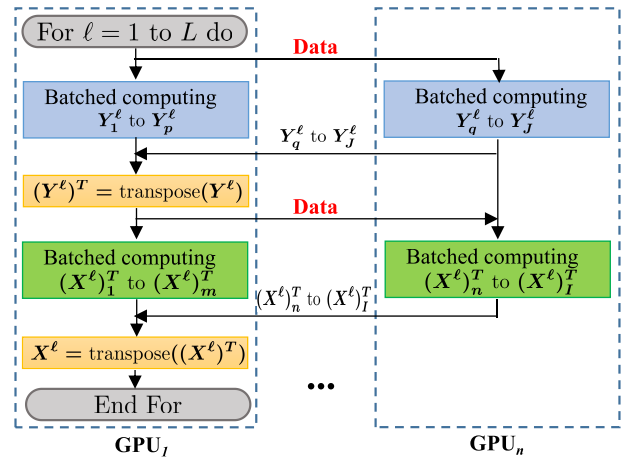


FIGURE 7. Multi-GPU scheme for homomorphic matrix completion.

B. MULTI-GPU HOMOMORPHIC MATRIX COMPLETION

Cloud computing nodes are often equipped with multiple GPUs, therefore we design a multi-GPU scheme to fully utilize multiple GPUs in a node for homomorphic matrix completion, as illustrated in Fig. 7. Since the computing of columns of Y^ℓ or $(X^\ell)^T$ is independent of each other, we can distributed the computation among multiple GPUs to compute in parallel. Assuming there are n GPUs, denoted by GPU₁, ..., GPU_n, in a computing node and the size of the encrypted matrix is $I \times J$, we perform the computation in the ℓ -th iteration (in lines 2–13 of Alg. 3) as follows:

- GPU_j computes C_j continuous columns of Y^ℓ , where C_j is proportional to the performance of the GPU (e.g., GFLOPs when performing matrix multiplication) and $\sum_{j=1}^n C_j = J$, where n is the number of GPUs. GPU₁ sends the data for computing to other GPUs using direct GPU to GPU communication. Each GPU computes its assigned columns of Y^ℓ using batched matrix multiplications and batched least square minimization. Upon completion, all other GPUs send their results to GPU₁ using direct GPU to GPU communication to form Y^ℓ ;
- GPU₁ performs matrix transpose on Y^ℓ to obtain $(Y^\ell)^T$;
- GPU_i computes C_i continuous columns of $(X^\ell)^T$, where C_i is proportional to the performance of the GPU (e.g., GFLOPs when performing matrix multiplication) and $\sum_{i=1}^n C_i = I$, where n is the number of GPUs. GPU₁ sends the data for computing to other GPUs using direct GPU to GPU communication. Each GPU computes its assigned columns of $(X^\ell)^T$ using batched matrix multiplications and batched least square minimization. Upon completion, all other GPUs send their results to GPU₁ using direct GPU to GPU communication to form $(X^\ell)^T$;
- GPU₁ performs matrix transpose on $(X^\ell)^T$ to obtain X^ℓ .

In this scheme, the direct GPU to GPU communication is implemented by exploiting the asynchronous, GPU to GPU memory copying routine in the CUDA library [21], which supports the communication between two and more GPUs.

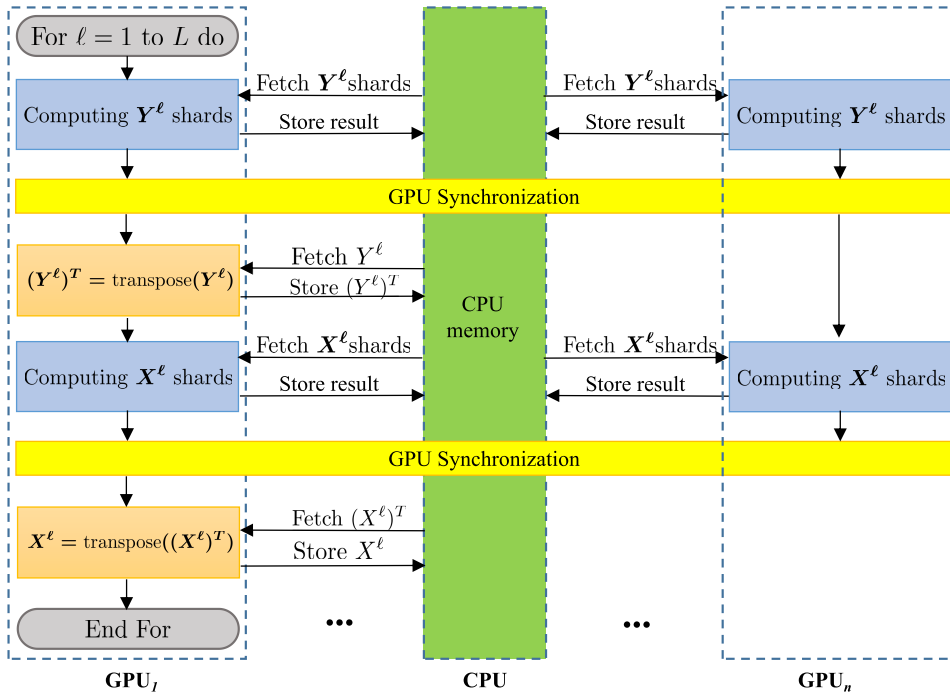


FIGURE 8. Large-scale homomorphic matrix completion on multiple GPUs.

C. LARGE-SCALE HOMOMORPHIC MATRIX COMPLETION ON MULTIPLE GPUS

We further design a scheme to perform large-scale homomorphic matrix completion exceeding GPU memory capacity on Multiple GPUs, as illustrated in Fig. 8. The basic idea is to exploit large CPU memory to store all the data, then distribute shards among multiple GPUs, and utilize GPU synchronization to ensure correct algorithm flow. Assuming there are n GPUs, denoted by GPU₁, ..., GPU_n, in a computing node, we perform the computation in the l -th iteration (in lines 2–13 of Alg. 3) as follows:

- Each GPU fetches a Y^ℓ shard from CPU memory, computes it and stores results (i.e., columns of Y^ℓ) into CPU memory. Then the GPU that finishes computing continues to fetch the next shard if there are remaining Y^ℓ shards;
- Performing GPU synchronization across all GPUs to ensure the computing of all Y^ℓ shards are completed;
- GPU₁ fetches Y^ℓ from CPU, performs matrix transpose on Y^ℓ to obtain $(Y^\ell)^T$, and store $(Y^\ell)^T$ back to CPU;
- Each GPU fetches a X^ℓ shard from CPU memory, compute it and store results (i.e., columns of $(X^\ell)^T$) into CPU memory. Then the GPU that finishes computing continues to fetch the next shard if there are remaining X^ℓ shards;
- Performing GPU synchronization across all GPUs to ensure the computing of all X^ℓ shards are completed;
- GPU₁ fetches $(X^\ell)^T$ from CPU, performs matrix transpose on $(X^\ell)^T$ to obtain X^ℓ , and store X^ℓ back to CPU.

V. PERFORMANCE EVALUATION

A. EVALUATION SETTINGS

The experiment platform has two Intel Xeon E5-2640V4 CPUs, an NVIDIA Tesla V100 GPU (Volta architecture), a Quadro RTX6000 GPU (Turin architecture), and 80 GB DDR memory. Each E5-2640V4 CPU has 10 cores that support 20 threads with hyperthreading technology. The V100 GPU has 5120 CUDA cores with 32 GB DDR memory, achieving 14 TFLOPs single-precision performance and 900 GB/s memory bandwidth. The RTX6000 GPU has 4608 CUDA cores with 24 GB DDR memory, achieving 16.3 TFLOPs single-precision performance and 672 GB/s memory bandwidth. The two GPUs are running on driver 430.40. The operating system is Ubuntu 18.04 64bit. The GPU implementation and the CPU implementation are executed on CUDA 10.1 and MATLAB 2017b, respectively.

We generate low-rank matrices of varying size $(I \times J)$ for performance evaluation. The rank of matrices R is set to $0.01 \times \min(I, J)$. The data missing rate of matrices is set to 50%. For single GPU performance, we run experiments on the Tesla V100 GPU. For multi-GPU performance, we run experiments on Tesla V100 and Quadro RTX6000 parallelly. Both CPU and GPU implementations execute the homomorphic matrix completion algorithm for $L = 10$ iterations. All experiments are executed for five times and we report the average results.

We adopt two metrics for performance evaluation: running time and recovery error. We compare the running time of the CPU MATLAB implementation and the GPU implementation on varying matrix sizes, and calculate speedups as

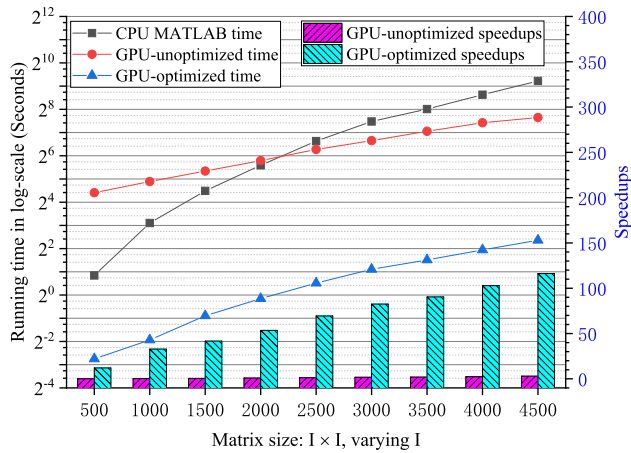


FIGURE 9. Running time and speedups of the homomorphic matrix completion algorithm on the Tesla V100 GPU and two Xeon CPUs for small- or medium-scale matrices.

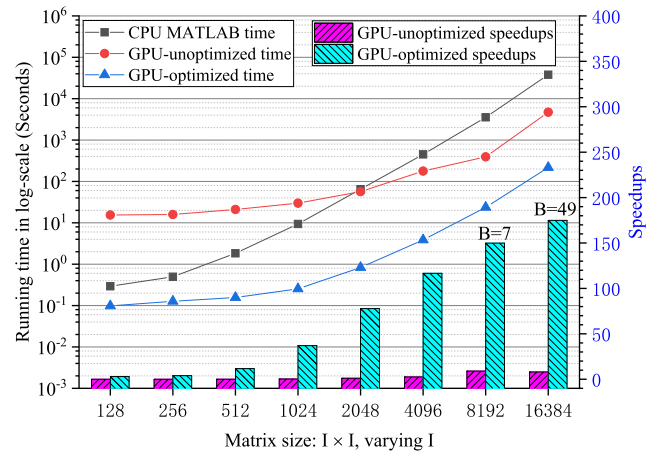


FIGURE 11. Running time and speedups of the large-scale GPU homomorphic matrix completion on a single GPU.

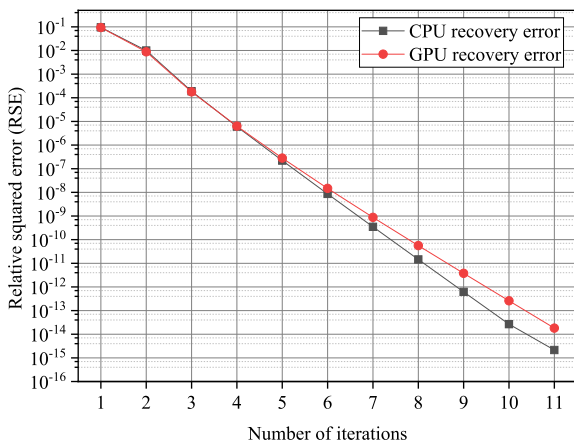


FIGURE 10. Recovery error of the homomorphic matrix completion algorithm on the Tesla V100 GPU and two Xeon CPUs, respectively.

(longer running time / shorter running time). The recovery error is defined as $RSE = \|D - \hat{D}\|_F / \|D\|_F$, where \hat{D} is given in (4).

B. BASIC PERFORMANCE OF THE GPU HOMOMORPHIC MATRIX COMPLETION

Fig. 9 shows the performance of the baseline (i.e., unoptimized) and the optimized GPU homomorphic matrix completion running on the Tesla V100 GPU and the CPU MATLAB implementation running on two Xeon CPUs, respectively. The optimized GPU homomorphic matrix completion achieved up to 116.23× speedups versus two Xeon CPUs. In contrast, the baseline GPU homomorphic matrix completion achieved up to 2.98× speedups versus two Xeon CPUs. The baseline GPU implementation adopts none of the optimization techniques in Section III-B. As a result, it has poor resource utilization and high CPU-GPU communication cost. The optimized GPU homomorphic matrix completion achieved higher speedups on bigger matrix sizes. However,

it supports up to 4,500 × 4,500 matrix due to limited GPU memory capacity.

Fig. 10 shows the recovery error of the optimized GPU homomorphic matrix completion on the Tesla V100 GPU and the CPU MATLAB implementation on two Xeon GPUs, respectively. The matrix is 1,000 × 1,000 with a 50% data missing rate. Both GPU and CPU MATLAB implementations used double precision float. The GPU and CPU implementation achieved similar recovery errors under different number of iterations. The recovery error keep decreasing with the number of iterations till 10e – 15 which is the limit of double precision float in MATLAB.

C. PERFORMANCE OF THE LARGE-SCALE AND MULTI-GPU HOMOMORPHIC MATRIX COMPLETION

Fig. 11 shows the performance on large-scale matrices using the shard mode. Due to the long running time of the CPU MATLAB implementation for comparison, we tested matrix sizes up to 16,384 × 16,384. The memory requirements for performing homomorphic matrix completion on two largest matrices 8,192 × 8,192 and 16,384 × 16,384 exceed the memory capacity of Tesla V100 GPU, therefore they are split into 7 and 49 shards for computing, respectively. On all matrix sizes, the optimized homomorphic matrix completion achieved up to 174.92× speedups versus the CPU MATLAB implementation. In contrast, the unoptimized GPU implementation poorly achieved up to 9.03× speedups over the CPU MATLAB implementation.

Fig. 12 shows the performance of the multi-GPU GPU homomorphic matrix completion. There are three running time curves to contrast the performance on single Tesla V100 GPU, single Quadro RTX6000 GPU, and two GPUs parallelly. In the two-GPU experiment, we set V100 as GPU₁ and RTX6000 as GPU₂, and the multi-GPU homomorphic matrix completion distributed the computation onto two GPUs and computed parallelly. On matrices of 1,000 × 1,000 or bigger, the GPU implementation running on RTX6000 was on average 1.13× faster than

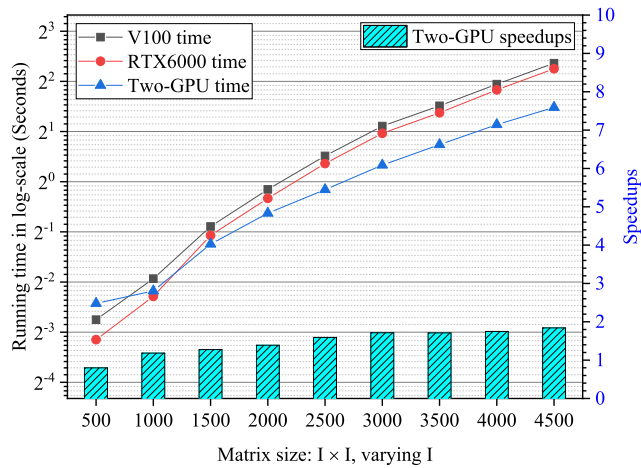


FIGURE 12. Running time and speedups of the multi-GPU homomorphic matrix completion.

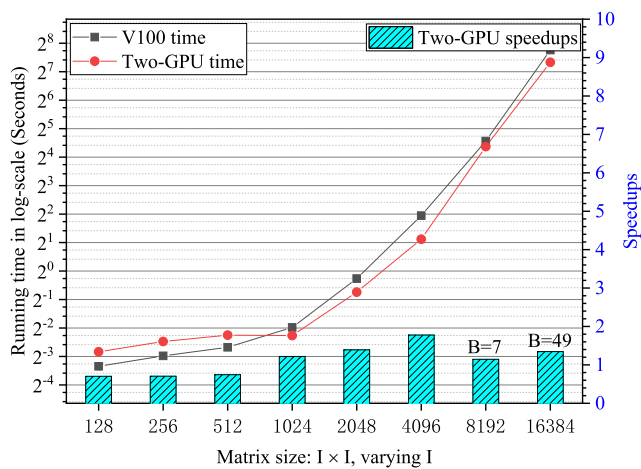


FIGURE 13. Running time and speedups of the large-scale GPU homomorphic matrix completion on multiple GPUs.

the GPU implementation running on V100. The multi-GPU homomorphic matrix completion achieved up to $1.84\times$ speedups against the single-GPU homomorphic matrix completion running on V100. On small matrix of 500×500 , the multi-GPU homomorphic matrix completion is slower than the single-GPU homomorphic matrix completion running on V100, because the running time is small and the cost to partition computation and the communication between multiple GPUs exceeds the performance benefit from multiple GPUs. Therefore, the multi-GPU scheme is more advantageous for medium and large matrices.

Fig. 13 compares the performance of the large-scale GPU homomorphic matrix completion on multiple GPUs. Two largest matrices $8,192 \times 8,192$ and $16,384 \times 16,384$ exceed the memory capacity of Tesla V100 GPU, therefore they are split into shards for computing. In the two-GPU experiment, we set V100 as GPU₁ and RTX6000 as GPU₂, and the large-scale, multi-GPU homomorphic matrix completion distributed all shards onto two GPUs and computed parallelly. On two matrices $8,192 \times 8,192$ and $16,384 \times 16,384$, the homomorphic matrix completion achieved $1.14\times$ and

$1.35\times$ speedups on two GPUs versus on a single V100 GPU, respectively. These speedups are lower than the multi-GPU speedups in Fig. 12 because of two major reasons. First, the large-scale, multi-GPU scheme requires two synchronizations between GPUs in each iteration to ensure the completion of computing, which introduces time overhead. Second, there can be workload imbalance between GPUs when GPUs process different number of shards and lead to GPU waiting time overhead.

VI. RELATED WORKS

Homomorphic encryption [23], [24] is a form of encryption that allows computation on ciphertexts, generating an encrypted result which, when decrypted, matches the result of the operations as if they had been performed on the plaintext. Homomorphic encryption approaches have been applied to diverse fields including secure voting systems, private information retrieval schemes, and collision-resistant hash functions [25], [26]. Kong *et al.* [19] proposed a homomorphic encryption approach for low-rank matrices. Using this approach, incomplete data from multiple users can be reconstructed on cloud computing nodes while still preserving user privacy.

Matrix completion methods are effective in reconstructing two dimensional data. Low rank is often a necessary hypothesis for multi-dimensional data completion to avoid being an undetermined and intractable problem. Low-rank, incomplete matrices can be reconstructed using the alternating minimization algorithm [7] [11] [14]. Candès and Recht [13] proposed a convex optimization method for matrix completion and obtained the sampling lower bound. For data reconstruction of three-dimension and higher-dimension, researchers proposed various tensor completion algorithms, including decomposition-based algorithms, nuclear (trace) norm-based algorithms, and other variants [27].

Data completion algorithms are compute-intensive. Many data completion algorithms are iterative, whose computation complexity grows rapidly with the size and dimension of the data. Therefore, researchers proposed to utilize high-performance GPUs for acceleration. Shah and Majumdar [28] proposed an efficient matrix completion implementation based on stochastic gradient descent (SGD) on GPUs. Later, researchers designed parallel SGD [29] and multi-stream SGD [30] on GPUs that can be used for matrix factorization and completion. For the completion of third-dimensional data, we proposed efficient GPU tensor completion implementations [31], [32]. Different with these works, we propose a high-performance GPU homomorphic matrix completion based on the alternating minimization.

The GPU homomorphic matrix completion in this paper involves matrix operations including matrix multiplication, matrix decomposition (for solving least squares minimization), and matrix transpose. Because matrix algebra is fundamental and widely used, designing high-performance matrix routines could benefit diverse applications. Tao *et al.* [33] accelerated Sparse matrix-vector

multiplication and matrix-transpose vector multiplication. Deveci *et al.* [34] designed efficient GPU sparse matrix-matrix multiplication. Tomov *et al.* [35] accelerated dense linear algebra on hybrid GPU systems. Our work considers how to organize the computation and exploiting matrix routines in batched, large-scale or multi-GPU conditions.

VII. CONCLUSION

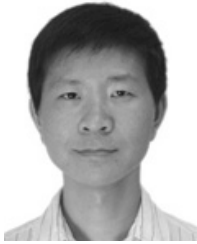
In this paper, we propose a GPU-based homomorphic matrix completion scheme for high-performance matrix completion while preserving user privacy. We optimized memory accesses, GPU utilizations, and CPU-GPU communications. The optimized scheme achieved up to $116.23\times$ speedups over the CPU MATLAB implementation running on two Xeon CPUs. To fully utilize multiple GPUs within a computing node, we designed a multi-GPU mode by exploiting the separability of the homomorphic matrix completion algorithm. To compute large-scale matrices exceeding GPU memory capacity, we designed a shard mode by exploiting large CPU memories. The multi-GPU mode achieves up to $1.84\times$ speedups on two GPUs versus on a single GPU. For large-scale matrices, the shard mode achieves up to $174.92\times$ speedups on a single GPU over the CPU MATLAB implementation on two CPUs, and further achieves up to $1.35\times$ speedups when running on two GPUs using the multi-GPU mode. The proposed scheme can be employed in diverse applications.

ACKNOWLEDGMENT

A conference version [1] of this article was presented at the 2019 IEEE International Conference on High Performance Computing and Communications.

REFERENCES

- [1] H. Lu, T. Zhang, and X.-Y. Liu, "High-performance homomorphic matrix completion on GPUs," in *Proc. IEEE 21st Int. Conf. High Perform. Comput. Commun., IEEE 17th Int. Conf. Smart City, IEEE 5th Int. Conf. Data Sci. Syst. (HPCC/SmartCity/DSS)*, Aug. 2019, pp. 1627–1634.
- [2] H. Wu, X.-Y. Liu, L. Fu, and X. Wang, "Energy-efficient and robust tensor-encoder for wireless camera networks in Internet of Things," *IEEE Trans. Netw. Sci. Eng.*, vol. 6, no. 4, pp. 646–656, Oct. 2019.
- [3] X.-Y. Liu, S. Aeron, V. Aggarwal, X. Wang, and M.-Y. Wu, "Adaptive sampling of RF fingerprints for fine-grained indoor localization," *IEEE Trans. Mobile Comput.*, vol. 15, no. 10, pp. 2411–2423, Oct. 2016.
- [4] S. Rallapalli, L. Qiu, Y. Zhang, and Y.-C. Chen, "Exploiting temporal stability and low-rank structure for localization in mobile networks," in *Proc. 16th Annu. Int. Conf. Mobile Comput. Netw.*, 2010, pp. 161–172.
- [5] Y. Koren, R. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, vol. 42, no. 8, pp. 30–37, Aug. 2009.
- [6] X.-Y. Liu, Y. Zhu, L. Kong, C. Liu, Y. Gu, A. V. Vasilakos, and M.-Y. Wu, "CDC: Compressive data collection for wireless sensor networks," *IEEE Trans. Parallel Distrib. Syst.*, vol. 26, no. 8, pp. 2188–2197, Aug. 2014.
- [7] L. Kong, M. Xia, X.-Y. Liu, G. Chen, Y. Gu, M.-Y. Wu, and X. Liu, "Data loss and reconstruction in wireless sensor networks," *IEEE Trans. Paralle. Distrib. Syst.*, vol. 25, no. 11, pp. 2818–2828, Nov. 2014.
- [8] X.-Y. Liu and X. Wang, "LS-decomposition for robust recovery of sensory big data," *IEEE Trans. Big Data*, vol. 4, no. 4, pp. 542–555, Dec. 2018.
- [9] S. Liao, X.-Y. Liu, F. Qian, M. Yin, and G.-M. Hu, "Tensor super-resolution for seismic data," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2019, pp. 8598–8602.
- [10] M. L. Stein, *Interpolation of Spatial Data: Some Theory for Kriging*. New York, NY, USA: Springer, 2012.
- [11] L. Kong, M. Xia, X.-Y. Liu, M.-Y. Wu, and X. Liu, "Data loss and reconstruction in sensor networks," in *Proc. IEEE INFOCOM*, Apr. 2013, pp. 1654–1662.
- [12] B. Recht, "A simpler approach to matrix completion," *J. Mach. Learn. Res.*, vol. 12, pp. 3413–3430, Jan. 2011.
- [13] E. Candès and B. Recht, "Exact matrix completion via convex optimization," *Commun. ACM*, vol. 55, no. 6, p. 111, Jun. 2012.
- [14] P. Jain, P. Netrapalli, and S. Sanghavi, "Low-rank matrix completion using alternating minimization," in *Proc. 45th Annu. ACM Symp. Symp. Theory Comput. (STOC)*, 2013, pp. 665–674.
- [15] C.-J. Hsieh, N. Natarajan, and I. S. Dhillon, "PU learning for matrix completion," in *Proc. Int. Conf. Mach. Learn.*, vol. 2015, pp. 2445–2453.
- [16] W. Li, L. Zhao, Z. Lin, D. Xu, and D. Lu, "Non-local image inpainting using low-rank matrix completion," *Comput. Graph. Forum*, vol. 34, no. 6, pp. 111–122, Sep. 2015.
- [17] R. Cabral, F. D. L. Torre, J. P. Costeira, and A. Bernardino, "Matrix completion for weakly-supervised multi-label image classification," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 1, pp. 121–135, Jan. 2015.
- [18] E. Adeli-Mosabbeh and M. Fathy, "Non-negative matrix completion for action detection," *Image Vis. Comput.*, vol. 39, pp. 38–51, Jul. 2015.
- [19] L. Kong, L. He, X.-Y. Liu, Y. Gu, M.-Y. Wu, and X. Liu, "Privacy-preserving compressive sensing for crowdsensing based trajectory recovery," in *Proc. IEEE 35th Int. Conf. Distrib. Comput. Syst.*, Jun. 2015, pp. 31–40.
- [20] S. Lohr, "Netflix cancels contest after concerns are raised about privacy," *New York Times*, Mar. 13, 2010, p. B3.
- [21] Nvidia Corporation. (2019). *NVIDIA CUDA SDK 10.1, NVIDIA CUDA Software Download*. [Online]. Available: <https://developer.nvidia.com/cuda-downloads>
- [22] A. Abdelfattah, A. Haidar, S. Tomov, and J. Dongarra, "Performance, design, and autotuning of batched GEMM for GPUs," in *Proc. Int. Conf. High Perform. Comput. Cham, Switzerland: Springer*, 2016, pp. 21–38.
- [23] C. Gentry, "Fully homomorphic encryption using ideal lattices," in *Proc. 41st Annu. ACM Symp. Symp. Theory Comput. (STOC)*, 2009, pp. 169–178.
- [24] M. Van Dijk, C. Gentry, S. Halevi, and V. Vaikuntanathan, "Fully homomorphic encryption over the integers," in *Proc. Annu. Int. Conf. Theory Appl. Cryptograph. Techn. Berlin, Germany: Springer*, 2010, pp. 24–43.
- [25] B. Chor, O. Goldreich, E. Kushilevitz, and M. Sudan, "Private information retrieval," in *Proc. IEEE 36th Annu. Found. Comput. Sci.*, 1995, pp. 41–50.
- [26] H. Sun and S. A. Jafar, "The capacity of private information retrieval," *IEEE Trans. Inf. Theory*, vol. 63, no. 7, pp. 4075–4088, Jul. 2017.
- [27] Q. Song, H. Ge, J. Caverlee, and X. Hu, "Tensor completion algorithms in big data analytics," *ACM Trans. Knowl. Discov. Data*, vol. 13, no. 1, pp. 1–48, Jan. 2019.
- [28] A. Shah and A. Majumdar, "Accelerating low-rank matrix completion on GPUs," in *Proc. Int. Conf. Adv. Comput., Commun. Informat. (ICACCI)*, Sep. 2014, pp. 182–187.
- [29] X. Xie, W. Tan, L. L. Fong, and Y. Liang, "CuMF_SGD: Parallelized stochastic gradient descent for matrix factorization on GPUs," in *Proc. 26th Int. Symp. High-Perform. Parallel Distrib. Comput.*, 2017, pp. 79–92.
- [30] H. Li, K. Li, J. An, and K. Li, "MSGD: A novel matrix factorization approach for large-scale collaborative filtering recommender systems on GPUs," *IEEE Trans. Parallel Distrib. Syst.*, vol. 29, no. 7, pp. 1530–1544, Jul. 2018.
- [31] T. Zhang, X.-Y. Liu, X. Wang, and A. Walid, "CuTensor-tubal: Efficient primitives for tubal-rank tensor learning operations on GPUs," *IEEE Trans. Parallel Distrib. Syst.*, vol. 31, no. 3, pp. 595–610, Mar. 2020.
- [32] H. Li, T. Zhang, R. Zhang, and X.-Y. Liu, "High-performance tensor decoder on GPUs for wireless camera networks in IoT," in *Proc. IEEE 21st Int. Conf. High Perform. Comput. Commun., IEEE 17th Int. Conf. Smart City, IEEE 5th Int. Conf. Data Sci. Syst. (HPCC/SmartCity/DSS)*, Aug. 2019, pp. 1619–1626.
- [33] Y. Tao, Y. Deng, S. Mu, Z. Zhang, M. Zhu, L. Xiao, and L. Ruan, "GPU accelerated sparse matrix-vector multiplication and sparse matrix-transpose vector multiplication," *Concurrency Comput., Pract. Exper.*, vol. 27, no. 14, pp. 3771–3789, Sep. 2015.
- [34] M. Deveci, C. Trott, and S. Rajamanickam, "Multithreaded sparse matrix-matrix multiplication for many-core and GPU architectures," *Parallel Comput.*, vol. 78, pp. 33–46, Oct. 2018.
- [35] S. Tomov, J. Dongarra, and M. Baboulin, "Towards dense linear algebra for hybrid GPU accelerated manycore systems," *Parallel Comput.*, vol. 36, nos. 5–6, pp. 232–240, Jun. 2010.



TAO ZHANG (Member, IEEE) received the bachelor's and master's degrees in computer science and technology from Xidian University, China, in 2001 and 2006, respectively, the Ph.D. degree in computer engineering from The University of New Mexico, USA, and the Ph.D. degree in computer science from Shanghai Jiao Tong University, China, in 2015.

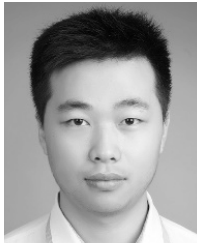
He is currently a Researcher and a Lecturer with the School of Computer Engineering and Science, Shanghai University, China. His research interests include big data processing, GPU heterogeneous computing, machine learning algorithms, and applications.



XIAO-YANG LIU (Student Member, IEEE) received the B.Eng. degree in computer science from the Huazhong University of Science and Technology, China, in 2010. He is currently pursuing the Ph.D. degree with the Department of Electrical Engineering, Columbia University. He graduated from the Ph.D. Program in the Department of Computer Science and Engineering, Shanghai Jiao Tong University, in 2017.

His research interests include tensor theory and high performance tensor computation, deep learning, optimization algorithms, and big data analysis and data privacy.

...



HAN LU received the B.Eng. degree in computer science and technology from the Hebei University of Engineering, China, in 2017. He is currently a Graduate Student with the School of Computer Engineering and Science, Shanghai University, China.

His research interests include GPU parallel computing, tensor computing, and optimization acceleration of algorithm.