

Communication Optimal Algorithms for Matrix and Tensor Computations

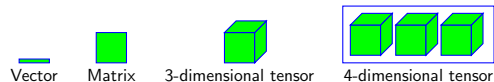
Suraj KUMAR

ROMA team
Inria & ENS Lyon

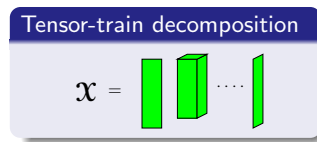
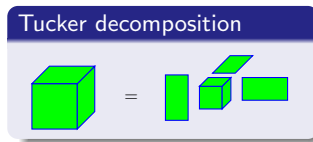
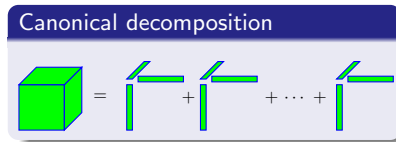
Project committee meeting
Jan 17, 2025

Tensors and their uses

- **Neuroscience:** Neuron \times Time \times Trial
- **Media:** User \times Movie \times Time
- **Ecommerce:** User \times Product \times Time
- **Social-Network:** Person \times Person \times Time \times Type



- High dimensional tensors: Neural network, Molecular simulation, Quantum computing
- People work with low dimensional structure (decomposition) of tensors

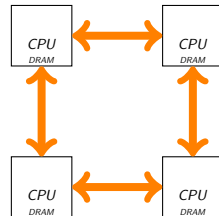


Importance of communication in high performance computing

- Gaps between computation and communication costs growing exponentially

	Annual improvements
Time-per-operation	59 %
Network-bandwidth	26 %
Network-latency	15 %

Source: Getting up to speed: The future of supercomputing (observed from 2004)



- **Goal:** Scalable and communication optimal tools for tensor computations

Communication Optimal Algorithms for Matrix and Tensor Computations

Hussam AL DAAS¹, Grey BALLARD², Laura GRIGORI³, Suraj KUMAR⁴, Kathryn ROUSE⁵, and Mathieu VÉRITÉ³

¹Rutherford Appleton Laboratory, UK

²Wake Forest University, USA

³EPFL, Switzerland

⁴Inria and ENS Lyon, France

⁵Inmar Intelligence, USA

Project committee meeting (Jan 17, 2025)

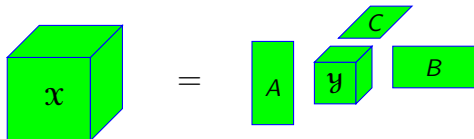
Outline

- 1 Parallel Multiple Tensor-Times-Matrix (Multi-TTM) computation
- 2 Parallel Nystrom approximation with random matrices

Our approach:

- Obtain communication lower bounds for each computation
- Design algorithms based on the communication lower bounds

Higher-order SVD (HOSVD) to compute Tucker decomposition

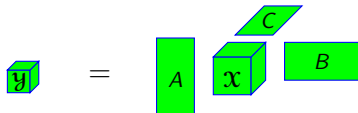


Algorithm 1 3-dimensional HOSVD Algorithm(\mathcal{X})

- 1: Obtain factor matrices A, B and C from the matrix representations of the input tensor \mathcal{X}
 - 2: $\mathcal{Y} = \mathcal{X} \times_1 A^T \times_2 B^T \times_3 C^T$
 - 3: Return \mathcal{Y}, A, B, C
-

- \mathcal{X}, \mathcal{Y} : 3-dimensional input and output tensors (or arrays) & A, B, C : matrices
- \times_i : tensor contraction along the i th dimension (similar to matrix multiplication)
- Multiple Tensor-Times-Matrix (Multi-TTM) computation: $\mathcal{Y} = \mathcal{X} \times_1 A^T \times_2 B^T \times_3 C^T$
- When A, B and C are obtained using randomized approaches, Multi-TTM becomes the bottleneck

Multi-TTM: $\mathcal{Y} = \mathcal{X} \times_1 A^T \times_2 B^T \times_3 C^T$



- Focus on communication cost of Multi-TTM on a parallel homogeneous machine
- Multi-TTM is also the bottleneck computation for Sequentially Truncated HOSVD

Settings

- P number of processors
- Each processor performs (asymptotically) equal amount of operations
- One copy of data is in the system
- Focus on bandwidth cost (volume of data transfers)

H. Al Daas, G. Ballard, L. Grigori, S. Kumar, and K. Rouse, *Communication lower bounds and optimal algorithms for multiple tensor-times-matrix computation*, SIMAX, 2024.

1 Parallel Multiple Tensor-Times-Matrix (Multi-TTM) computation

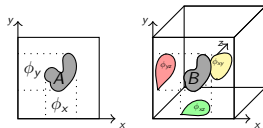
- Our approach to compute communication lower bounds
 - For Traditional Matrix Matrix Multiplication
 - For Multi-TTM Computation
- Communication Optimal Algorithm and Simulated Experiments
- Conclusion

2 Parallel Nyström approximation with random matrices

- Symmetric Nyström approximation
- Conclusion

Approach to obtain communication lower bounds

- Loomis-Whitney inequality: for $d - 1$ dimensional projections
 - For the 2d object A , $\phi_x \phi_y \geq \text{Area}(A)$
 - For the 3d object B , $(\phi_{xy} \phi_{yz} \phi_{xz})^{\frac{1}{2}} \geq \text{Volume}(B)$
- Hölder-Brascamp-Lieb (HBL) inequality – generalization for arbitrary dimensional projections
 - Provide exponent for each projection



Constraints for parallel load balanced matrix matrix multiplication

- $C = AB$ with $A \in \mathbb{R}^{n_1 \times n_2}$, $B \in \mathbb{R}^{n_2 \times n_3}$, and $C \in \mathbb{R}^{n_1 \times n_3}$

for $i = 1:n_1$, for $k = 1:n_2$, for $j = 1:n_3$

$$C[i][j] += A[i][k] * B[k][j]$$

- ϕ_A, ϕ_B, ϕ_C : projections of computations on arrays A, B, C
- From Loomis-Whitney/HBL inequality: $\phi_A^{\frac{1}{2}} \phi_B^{\frac{1}{2}} \phi_C^{\frac{1}{2}} \geq \text{number of multiplications per processor} = \frac{n_1 n_2 n_3}{P}$
- Extra constraints: $\frac{n_1 n_2}{P} \leq \phi_A \leq n_1 n_2$, $\frac{n_2 n_3}{P} \leq \phi_B \leq n_2 n_3$, $\frac{n_1 n_3}{P} \leq \phi_C \leq n_1 n_3$

Optimization problem and communication lower bounds

Minimize $\phi_A + \phi_B + \phi_C$ s.t.

$$\phi_A^{\frac{1}{2}} \phi_B^{\frac{1}{2}} \phi_C^{\frac{1}{2}} \geq \frac{n_1 n_2 n_3}{P}$$

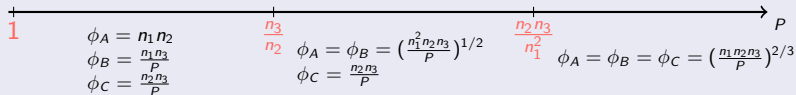
$$\frac{n_1 n_2}{P} \leq \phi_A \leq n_1 n_2$$

$$\frac{n_2 n_3}{P} \leq \phi_B \leq n_2 n_3$$

$$\frac{n_1 n_3}{P} \leq \phi_C \leq n_1 n_3$$

Amount of array accesses = $\phi_A + \phi_B + \phi_C$

- Estimate the solution and prove optimality by showing Karush–Kuhn–Tucker (KKT) conditions are satisfied
- For $n_1 \leq n_2 \leq n_3$,



- Communication lower bound = $\phi_A + \phi_B + \phi_C - \text{data owned by the processor} = \phi_A + \phi_B + \phi_C - \frac{n_1 n_2 + n_2 n_3 + n_1 n_3}{P}$

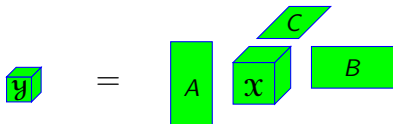
1 Parallel Multiple Tensor-Times-Matrix (Multi-TTM) computation

- Our approach to compute communication lower bounds
 - For Traditional Matrix Matrix Multiplication
 - For Multi-TTM Computation
- Communication Optimal Algorithm and Simulated Experiments
- Conclusion

2 Parallel Nyström approximation with random matrices

- Symmetric Nyström approximation
- Conclusion

3-dimensional Multi-TTM computation



- $\mathcal{Y} = \mathcal{X} \times_1 A^T \times_2 B^T \times_3 C^T$
- \mathcal{X}, \mathcal{Y} : 3-dimensional input and output tensors
- A, B, C : matrices
- \times_i : analogous to matrix multiplication

- TTM-in-Sequence approach (used in TuckerMPI library): $\mathcal{Y} = ((\mathcal{X} \times_1 A^T) \times_2 B^T) \times_3 C^T$
- Our All-at-Once definition with $\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, $\mathcal{Y} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$, $A \in \mathbb{R}^{n_1 \times r_1}$, $B \in \mathbb{R}^{n_2 \times r_2}$, $C \in \mathbb{R}^{n_3 \times r_3}$

for $\{n'_1, n'_2, n'_3, r'_1, r'_2, r'_3\} = 1:\{n_1, n_2, n_3, r_1, r_2, r_3\}$

$$\mathcal{Y}(r'_1, r'_2, r'_3) = \mathcal{X}(n'_1, n'_2, n'_3) \cdot A(n'_1, r'_1) \cdot B(n'_2, r'_2) \cdot C(n'_3, r'_3)$$

Solving optimization problems to compute lower bounds

- Each processor performs $\frac{n_1 r_1 n_2 r_2 n_3 r_3}{P}$ amount of 4 – *array* operations
- After applying lower and upper bounds for each projection, we need to solve the following optimization problem

Minimize $\phi_x + \phi_y + \phi_1 + \phi_2 + \phi_3$ s.t.

$$\phi_x^{1-a} \phi_y^{1-a} \phi_1^a \phi_2^a \phi_3^a \geq \frac{n_1 r_1 n_2 r_2 n_3 r_3}{P}$$

$$\frac{n_1 n_2 n_3}{P} \leq \phi_x \leq n_1 n_2 n_3$$

$$\frac{r_1 r_2 r_3}{P} \leq \phi_y \leq r_1 r_2 r_3$$

$$\frac{n_1 r_1}{P} \leq \phi_1 \leq n_1 r_1$$

$$\frac{n_2 r_2}{P} \leq \phi_2 \leq n_2 r_2$$

$$\frac{n_3 r_3}{P} \leq \phi_3 \leq n_3 r_3$$

$$0 \leq a \leq 1$$

Divide the problem into two parts

Matrix part

Minimize $\phi_1 + \phi_2 + \phi_3$ s.t.

$$\phi_1 \phi_2 \phi_3 \geq \frac{n_1 r_1 n_2 r_2 n_3 r_3}{P}$$

$$\frac{n_1 r_1}{P} \leq \phi_1 \leq n_1 r_1$$

$$\frac{n_2 r_2}{P} \leq \phi_2 \leq n_2 r_2$$

$$\frac{n_3 r_3}{P} \leq \phi_3 \leq n_3 r_3$$

Tensor part

Minimize $\phi_x + \phi_y$ s.t.

$$\phi_x \phi_y \geq \frac{n_1 r_1 n_2 r_2 n_3 r_3}{P}$$

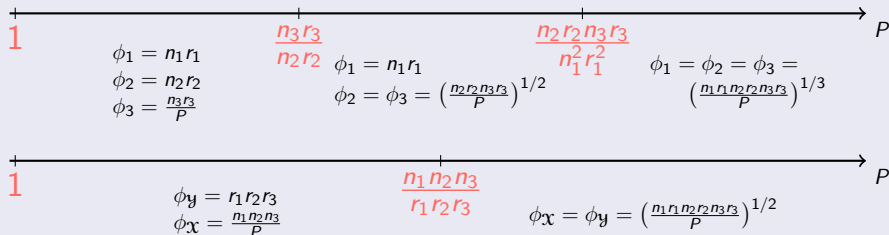
$$\frac{n_1 n_2 n_3}{P} \leq \phi_x \leq n_1 n_2 n_3$$

$$\frac{r_1 r_2 r_3}{P} \leq \phi_y \leq r_1 r_2 r_3$$

Amount of accesses and lower bounds

- We assume $n_1 r_1 \leq n_2 r_2 \leq n_3 r_3$ and $r_1 r_2 r_3 \leq n_1 n_2 n_3$
- Estimate the solutions and prove optimality by showing KKT conditions are satisfied

Amount of accesses = $\phi_x + \phi_y + \phi_1 + \phi_2 + \phi_3$



$$\text{Communication lower bound} = \phi_x + \phi_y + \phi_1 + \phi_2 + \phi_3 - \frac{n_1 n_2 n_3 + r_1 r_2 r_3 + n_1 r_1 + n_2 r_2 + n_3 r_3}{P}$$

1 Parallel Multiple Tensor-Times-Matrix (Multi-TTM) computation

- Our approach to compute communication lower bounds
 - For Traditional Matrix Matrix Multiplication
 - For Multi-TTM Computation
- Communication Optimal Algorithm and Simulated Experiments
- Conclusion

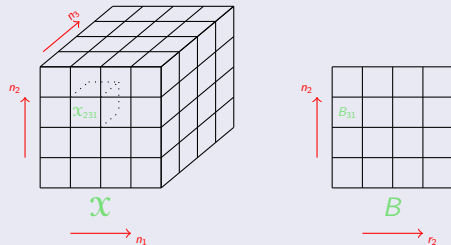
2 Parallel Nyström approximation with random matrices

- Symmetric Nyström approximation
- Conclusion

Design of communication optimal algorithms

Data Distribution (P is organized into a $p_1 \times p_2 \times p_3 \times q_1 \times q_2 \times q_3$ grid)

- p_1, p_2, p_3, q_1, q_2 , and q_3 evenly distribute n_1, n_2, n_3, r_1, r_2 , and r_3
- Each processor has $\frac{1}{p}$ th amount of input and output variables
- Subtensor $\mathcal{X}_{231} = \mathcal{X}(\frac{n_1}{p_1} + 1 : 2\frac{n_1}{p_1}, 2\frac{n_2}{p_2} + 1 : 3\frac{n_2}{p_2}, 1 : \frac{n_3}{p_3})$ is distributed evenly among processors $(2, 3, 1, *, *, *)$
- Submatrix $B_{31} = B(2\frac{n_2}{p_2} + 1 : 3\frac{n_2}{p_2}, 1 : \frac{r_2}{q_2})$ is distributed evenly among processors $(*, 3, *, *, 1, *)$



6-dimensional algorithm to compute Multi-TTM

Algorithm 1 3-dimensional Parallel Atomic Multi-TTM

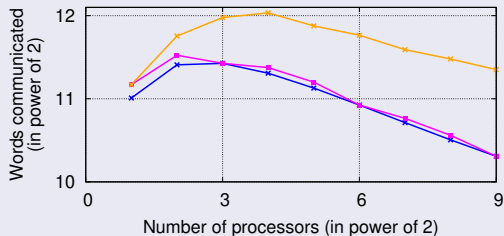
Require: \mathcal{X} , A , B , C , $p_1 \times p_2 \times p_3 \times q_1 \times q_2 \times q_3$ logical processor grid

Ensure: \mathcal{Y} such that $\mathcal{Y} = \mathcal{X} \times_1 A^\top \times_2 B^\top \times_3 C^\top$

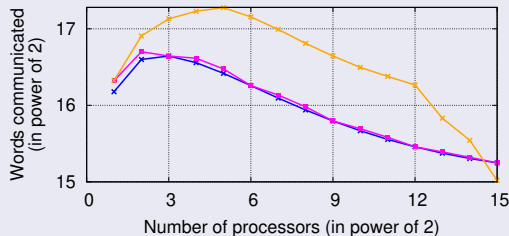
- 1: $(p'_1, p'_2, p'_3, q'_1, q'_2, q'_3)$ is my processor id
 - 2: //All-gather input tensor and matrices
 - 3: $\mathcal{X}_{p'_1 p'_2 p'_3} = \text{All-Gather}(\mathcal{X}, (p'_1, p'_2, p'_3, *, *, *))$
 - 4: $A_{p'_1 q'_1} = \text{All-Gather}(A, (p'_1, *, *, q'_1, *, *))$
 - 5: $B_{p'_2 q'_2} = \text{All-Gather}(B, (*, p'_2, *, *, q'_2, *))$
 - 6: $C_{p'_3 q'_3} = \text{All-Gather}(C, (*, *, p'_3, *, *, q'_3))$
 - 7: //Perform local Multi-TTM computation in a temporary tensor \mathcal{T}
 - 8: $\mathcal{T} = \text{Local-Multi-TTM}(\mathcal{X}_{p'_1 p'_2 p'_3}, A_{p'_1 q'_1}, B_{p'_2 q'_2}, C_{p'_3 q'_3})$
 - 9: //Reduce-scatter the output tensor in $\mathcal{Y}_{q'_1 q'_2 q'_3}$
 - 10: $\text{Reduce-Scatter}(\mathcal{Y}_{q'_1 q'_2 q'_3}, \mathcal{T}, (*, *, *, q'_1, q'_2, q'_3))$
-

Performance comparison (simulated experiments)

$$n_1 = n_2 = n_3 = 2^8, r_1 = r_2 = r_3 = 2^3$$



$$n_1 = n_2 = n_3 = 2^{11}, r_1 = r_2 = r_3 = 2^5$$



Lower Bound — x —

Our Approach — ■ —

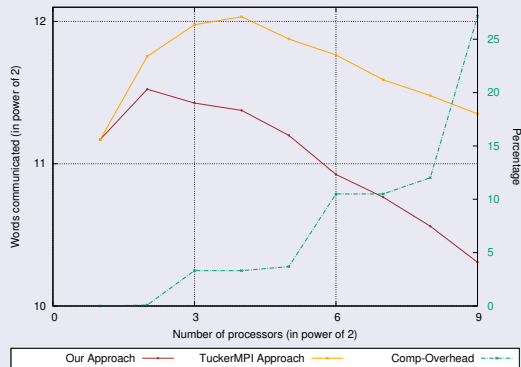
TuckerMPI Approach — x —

- Typical scenarios in data compression problems
- For small P , our approach communicates much less than TuckerMPI approach

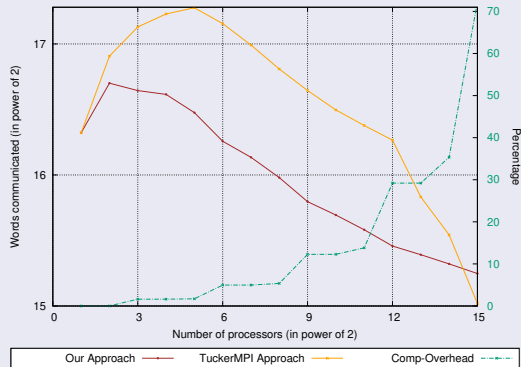
Implementation of the proposed approach: R. Minster, Z. Li, and G. Ballard, *Parallel Randomized Tucker Decomposition Algorithms*, SISC, 2024.

Extra computation in our approach

$$n_1 = n_2 = n_3 = 2^8, r_1 = r_2 = r_3 = 2^3$$



$$n_1 = n_2 = n_3 = 2^{11}, r_1 = r_2 = r_3 = 2^5$$



1 Parallel Multiple Tensor-Times-Matrix (Multi-TTM) computation

- Our approach to compute communication lower bounds
 - For Traditional Matrix Matrix Multiplication
 - For Multi-TTM Computation
- Communication Optimal Algorithm and Simulated Experiments
- Conclusion

2 Parallel Nyström approximation with random matrices

- Symmetric Nyström approximation
- Conclusion

Conclusion and future work

Conclusion

- Communication lower bounds and optimal algorithms for All-at-Once Multi-TTM
- Our algorithm communicates much less data than TTM-in-Sequence for small P

Future Work

- Detailed study of what scenarios are favorable for our approach
- Combine both All-at-Once and TTM-in-Sequence approaches

Outline

- 1 Parallel Multiple Tensor-Times-Matrix (Multi-TTM) computation
- 2 Parallel Nyström approximation with random matrices

Nonsymmetric Nyström approximation

$\tilde{A} = (AV)(U^T AV)^\dagger (A^T U)^T$ with $A \in \mathbb{R}^{n_1 \times n_2}$, $U \in \mathbb{R}^{n_1 \times r_1}$, and $V \in \mathbb{R}^{n_2 \times r_2}$.

- U and V are random matrices
 - can be generated on any processor without any extra communication costs
- Need to focus on $D = U^T A$, $B = AV$ and $C = U^T AV$ computations together

We will focus mainly on symmetric Nyström approximation today.

1 Parallel Multiple Tensor-Times-Matrix (Multi-TTM) computation

- Our approach to compute communication lower bounds
 - For Traditional Matrix Matrix Multiplication
 - For Multi-TTM Computation
- Communication Optimal Algorithm and Simulated Experiments
- Conclusion

2 Parallel Nyström approximation with random matrices

- Symmetric Nyström approximation
- Conclusion

Symmetric Nystrom approximation

$\tilde{A} = (AV)(V^T AV)^\dagger (AV)^T$ with $A \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{n \times r}$.

- V is a random matrix and $A = A^T$
- Need to focus on $B = AV$ and $C = V^T AV$ computations together

A naive way to minimize communication cost for $B = AV$

$$\begin{pmatrix} B_1 \\ B_2 \\ B_3 \end{pmatrix} = \begin{pmatrix} A_1 \\ A_2 \\ A_3 \end{pmatrix} \cdot V$$

Structure of the algorithm for each processor i :

- owns A_i row block ($1/P$ th portion of A), generates random matrix V , and performs $B_i = A_i \cdot V$

When P is small, communication cost of the algorithm is 0 (**Optimal**).

What about when P is large?

Solving optimization problems to compute lower bounds for $B = AV$

- Each processor performs $\frac{n^2 r}{P}$ multiplications
- After applying lower and upper bounds for each projection, we need to solve the following optimization problem

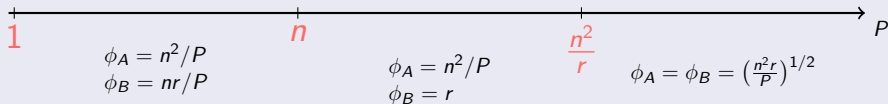
Minimize $\phi_A + \phi_B$ s.t.

$$\phi_A \phi_B \geq \text{ncols}(\phi_A) \phi_B \geq \frac{n^2 r}{P}$$

$$\phi_A^{1/2} \phi_B^{1/2} \phi_V^{1/2} \geq \frac{n^2 r}{P}$$

$$\frac{n^2}{P} \leq \phi_A \leq n^2, \quad \frac{nr}{P} \leq \phi_B \leq nr, \quad \frac{nr}{P} \leq \phi_V \leq nr, \quad r \leq n$$

Communication lower bound = $\phi_A + \phi_B - \frac{nr+n^2}{P}$



Our algorithm to compute $B = AV$ computation

Algorithm 2 $B = AV$ computation

Require: A , $p_1 \times p_2 \times p_3$ logical processor grid

Ensure: B such that $B = AV$

- 1: (p'_1, p'_2, p'_3) is my processor id
 - 2: $A_{p'_1 p'_2} = \text{All-Gather}(A, (p'_1, p'_2, *))$ //All-gather input matrix A
 - 3: $V_{p'_2 p'_3} = \text{GenerateRequiredRandomMatrix}()$ //Generate random matrix of size $\frac{n}{p_2} \times \frac{r}{p_3}$
 - 4: //Perform local computation in a temporary matrix T
 - 5: $T = A_{p'_1 p'_2} \cdot V_{p'_2 p'_3}$
 - 6: $\text{Reduce-Scatter}(B_{p'_1 p'_3}, T, (p'_1, *, p'_3))$ //Reduce-scatter the output matrix in $B_{p'_1 p'_3}$
-

Our algorithm is communication optimal when p_1 , p_2 & p_3 are chosen based on lower bounds.

Compute lower bounds for $B = AV$ and $C = V^T AV$

$$B = AV \quad \text{and} \quad C = V^T AV$$

- Each processor performs $\frac{n^2 r}{P}$ multiplications to compute B and $\frac{nr^2}{P}$ multiplications to compute C
- Assume that same entries of B are accessed on a processor for both computations
- After combining iteration spaces of both computations in a 6-dimensional lattice, we need to solve the following optimization problem

Minimize $\phi_A + \phi_B + \phi_C$ s.t.

$$\phi_A \phi_B \phi_C \geq \frac{n^2 r}{P} \cdot \frac{nr^2}{P}$$

$$\frac{n^2}{P} \leq \phi_A \leq n^2, \quad \frac{nr}{P} \leq \phi_B \leq nr, \quad \frac{r^2}{P} \leq \phi_C \leq r^2, \quad r \leq n$$

- Can obtain communication lower bounds by solving the above optimization problem
- Similar to the previous algorithms, we can design communication optimal algorithms for a $p_1 \times p_2 \times p_3$ logical processor grid

1 Parallel Multiple Tensor-Times-Matrix (Multi-TTM) computation

- Our approach to compute communication lower bounds
 - For Traditional Matrix Matrix Multiplication
 - For Multi-TTM Computation
- Communication Optimal Algorithm and Simulated Experiments
- Conclusion

2 Parallel Nyström approximation with random matrices

- Symmetric Nyström approximation
- Conclusion

Conclusion

- Communication lower bounds and optimal algorithms for the computations of symmetric Nyström approximation
- An approach to obtain communication lower bounds for a set of computations

Future Work

- Implementation of the proposed algorithms for real datasets
- Simplify and refine the lower bound computations

Thank You!