

Communication Optimal Algorithms for Multiple Tensor-Times-Matrix Computation

Hussam AL DAAS¹, Grey BALLARD², Laura GRIGORI³, Suraj KUMAR³, and Kathryn ROUSE⁴

¹Rutherford Appleton Laboratory, UK

²Wake Forest University, USA

³Inria Paris, France

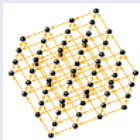
⁴Inmar Intelligence, USA

ROMA Working Group
December 7, 2021

My Past Experience

Parallelization in Polyhedral Model

- Linked-list operations
- Improved spatial locality
- Parallelization using OpenMP



Seismic Imaging on GPUs

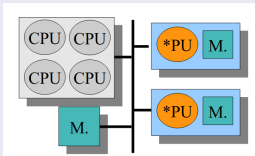
$$H_1 = \sin^2 \theta \cos^2 \phi \frac{\partial^2}{\partial x^2} + \sin^2 \theta \sin^2 \phi \frac{\partial^2}{\partial y^2} + \cos^2 \theta \frac{\partial^2}{\partial z^2} + \sin^2 \theta \sin 2\phi \frac{\partial^2}{\partial x \partial y} + \sin 2\theta \sin \phi \frac{\partial^2}{\partial y \partial z} + \sin 2\theta \cos \phi \frac{\partial^2}{\partial x \partial z}$$
$$H_2 = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2} - H_1$$

Schedulers for Blue Gene Supercomputers

- GASNET API
- Unbalanced Tree Search benchmark
- Comparison to Charm++

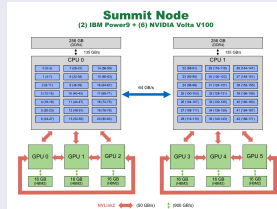


Scheduling on Heterogeneous Platforms



Molecular Simulations on Supercomputers

- NWChemEx Project
- TAMM library
- Hartree Fock and CCSD applications



Communication Optimal Algorithms for Multiple Tensor-Times-Matrix Computation

Hussam AL DAAS¹, Grey BALLARD², Laura GRIGORI³, Suraj KUMAR³, and Kathryn ROUSE⁴

¹Rutherford Appleton Laboratory, UK

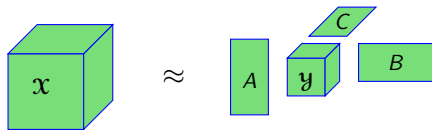
²Wake Forest University, USA

³Inria Paris, France

⁴Inmar Intelligence, USA

ROMA Working Group
December 7, 2021

Higher-order SVD (HOSVD) to compute Tucker decomposition



Algorithm 1 HOSVD Algorithm(\mathcal{X} , R_1 , R_2 , R_3)

- 1: $A \leftarrow R_1$ left singular vectors of $\mathcal{X}_{(1)}$
 - 2: $B \leftarrow R_2$ left singular vectors of $\mathcal{X}_{(2)}$
 - 3: $C \leftarrow R_3$ left singular vectors of $\mathcal{X}_{(3)}$
 - 4: $\mathcal{Y} = \mathcal{X} \times_1 A^T \times_2 B^T \times_3 C^T$
 - 5: Return \mathcal{Y} , A , B , C
-

- \mathcal{X} , \mathcal{Y} : 3-dimensional input and output tensors (or arrays) & A , B , C : matrices
- $\mathcal{X}_{(i)}$: matricization of \mathcal{X} (i th dimension represents rows and remaining dimensions represent columns)
- Multiple Tensor-Times-Matrix (Multi-TTM) computation: $\mathcal{Y} = \mathcal{X} \times_1 A^T \times_2 B^T \times_3 C^T$
 - To obtain full tensor, $\mathcal{X} = \mathcal{Y} \times_1 A \times_2 B \times_3 C$

Outline (Lower Bounds and Communication Optimal Algorithms)

- 1 For Matrix Matrix Multiplications
- 2 For Multi-TTM Computation
 - Simulated Experiments
 - Conclusion

Assumptions

- P number of processors
- Each processor performs (asymptotically) equal amount of operations
- No redundant operations
- One copy of data is in the system
 - $1/P$ th amount of inputs (before the computation) and output (after the computation) on each processor
- Each processor has enough memory

Outline (Lower Bounds and Communication Optimal Algorithms)

1 For Matrix Matrix Multiplications

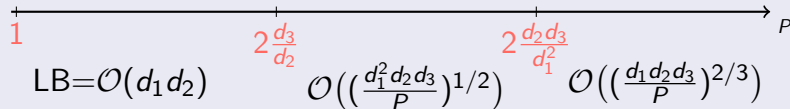
2 For Multi-TTM Computation

- Simulated Experiments
- Conclusion

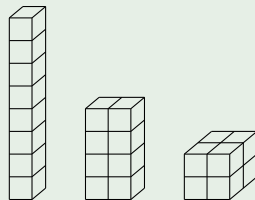
Existing Lower Bounds for Matrix Matrix Multiplications

- $C = AB$, where $A \in \mathbb{R}^{n_1 \times n_2}$, $B \in \mathbb{R}^{n_2 \times n_3}$, and $C \in \mathbb{R}^{n_1 \times n_3}$
- Let $d_1 = \min(n_1, n_2, n_3) \leq d_2 = \text{median}(n_1, n_2, n_3) \leq d_3 = \max(n_1, n_2, n_3)$

Existing Communication Lower Bounds (CARMA [IPDPS 2013])



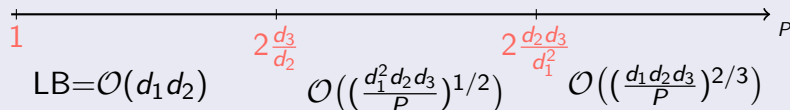
Arrangements of 8 processors



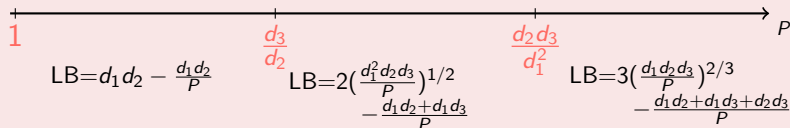
Existing Lower Bounds for Matrix Matrix Multiplications

- $C = AB$, where $A \in \mathbb{R}^{n_1 \times n_2}$, $B \in \mathbb{R}^{n_2 \times n_3}$, and $C \in \mathbb{R}^{n_1 \times n_3}$
- Let $d_1 = \min(n_1, n_2, n_3) \leq d_2 = \text{median}(n_1, n_2, n_3) \leq d_3 = \max(n_1, n_2, n_3)$

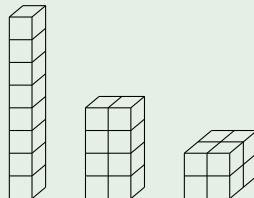
Existing Communication Lower Bounds (CARMA [IPDPS 2013])



Our Communication Lower Bounds



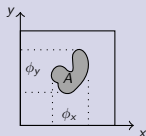
Arrangements of 8 processors



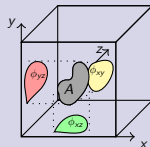
Loomis-Whitney & Hölder-Brascamp-Lieb inequalities

Size of $d - 1$ dimensional projections (Loomis-Whitney inequality)

- 2-dimensional object A and its 1-dimensional projections ϕ_x, ϕ_y
- $\phi_x \phi_y \geq \text{Area}(A)$



- 3-dimensional object A and its 2-dimensional projections: $\phi_{xy}, \phi_{yz}, \phi_{xz}$
- $(\phi_{xy} \phi_{yz} \phi_{xz})^{\frac{1}{3-1}} \geq \text{Volume}(A)$



Hölder-Brascamp-Lieb (HBL) inequality – Generalization of Loomis-Whitney inequality

$$\Delta = \begin{matrix} & A & B & C \\ \begin{matrix} i \\ j \\ k \end{matrix} & \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 0 \end{pmatrix} \end{matrix}$$

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

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

- Find $\mathbf{x} = [x_1 \ x_2 \ x_3]^T$ such that $\Delta \cdot \mathbf{x} \geq \mathbf{1}$, $\mathbf{1}$ is vector of all ones
- ϕ_A, ϕ_B, ϕ_C : projections of computations on arrays A, B, C
- HBL inequality: $\phi_A^{x_1} \phi_B^{x_2} \phi_C^{x_3} \geq \text{Amount of computations}$
- To make inequality tight select \mathbf{x} such that $\mathbf{1}^T \mathbf{x}$ is minimum $\Rightarrow x_1 = x_2 = x_3 = \frac{1}{2}$

Constraints for Matrix Multiplications

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

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

- Total number of multiplications = $n_1 n_2 n_3$
- Consider a processor which performs $\frac{n_1 n_2 n_3}{P}$ amount of multiplications
- Optimization problem:

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

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

Extra constraints (our contributions)

- Each element of A (resp. B) is involved in n_3 (resp. n_1) multiplications
 - To perform at least $\frac{n_1 n_2 n_3}{P}$ multiplications: $\phi_A \geq \frac{n_1 n_2}{P}, \phi_B \geq \frac{n_2 n_3}{P}$
- Each element of C is the sum of n_2 multiplications, therefore $\phi_C \geq \frac{n_1 n_3}{P}$
- Projections can be at max the size of the arrays: $\phi_A \leq n_1 n_2, \phi_B \leq n_2 n_3, \phi_C \leq n_1 n_3$

Optimization Problem to Compute Communication Lower Bounds

- Projections (ϕ_A, ϕ_B, ϕ_C) indicate the amount of array access
- Communication lower bound = $\phi_A + \phi_B + \phi_C$ – data owned by the processor

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

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

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

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

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

Generalized version (in terms of d_1, d_2, d_3)

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

$$\phi_1^{\frac{1}{2}} \phi_2^{\frac{1}{2}} \phi_3^{\frac{1}{2}} \geq \frac{d_1 d_2 d_3}{P}$$

$$\frac{d_1 d_2}{P} \leq \phi_1 \leq d_1 d_2$$

$$\frac{d_1 d_3}{P} \leq \phi_2 \leq d_1 d_3$$

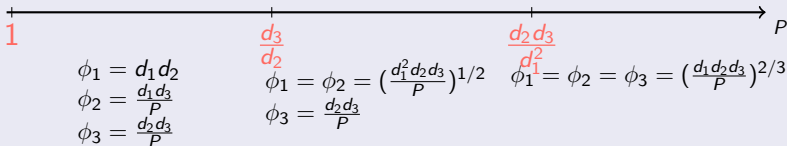
$$\frac{d_2 d_3}{P} \leq \phi_3 \leq d_2 d_3$$

$$d_1 \leq d_2 \leq d_3$$

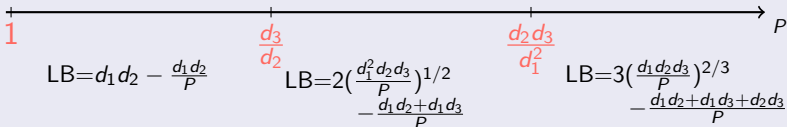
Amount of Accesses and Communication Lower bounds

- Estimate the solution based on Lagrange multipliers
- Prove optimality using all KarushKuhnTucker (KKT) conditions are satisfied

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



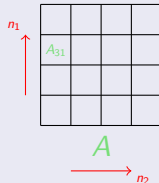
Communication Lower Bounds (Amount of Data Transfers)



Design of Communication Optimal Algorithms

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

- p_1, p_2 , and p_3 evenly distribute n_1, n_2 , and n_3
- Each processor has $\frac{1}{p}$ th amount of input and output variables
- $A_{31} = A(2\frac{n_1}{p_1} + 1 : 3\frac{n_1}{p_1}, 1 : \frac{n_2}{p_2})$ is evenly distributed among $(3, 1, *)$ processors
- $B_{12} = B(1 : \frac{n_2}{p_2}, \frac{n_3}{p_3} + 1 : 2\frac{n_3}{p_3})$ is evenly distributed among $(*, 1, 2)$ processors



Algorithm 2 $C = AB$ Matrix Multiplication Algorithm

- 1: (p'_1, p'_2, p'_3) is my processor id
- 2: //All-gather input matrices A and B
- 3: $A_{p'_1 p'_2} = \text{All-Gather}(A, (p'_1, p'_2, *))$
- 4: $B_{p'_2 p'_3} = \text{All-Gather}(B, (*, p'_2, p'_3))$
- 5: $T = \text{Local-Matrix-Multiplication}(A_{p'_1 p'_2}, B_{p'_2 p'_3})$ // Local matrix multiplication in a temporary
- 6: $\text{Reduce-Scatter}(C_{p'_1 p'_3}, T, (p'_1, *, p'_3))$ // Reduce-scatter the output

Cost Analysis and Open Questions

Cost Analysis

- Total amount of multiplications per processor = $\frac{n_1 n_2 n_3}{p_1 p_2 p_3} = \frac{n_1 n_2 n_3}{P}$
- Total data transfers = $\frac{n_1 n_2}{p_1 p_2} + \frac{n_2 n_3}{p_2 p_3} + \frac{n_1 n_3}{p_1 p_3} - \frac{n_1 n_2 + n_2 n_3 + n_1 n_3}{P}$

Open Questions

- How to select p_1, p_2, p_3 ?
- Are communication lower bounds are achievable for all matrix dimensions?

1 For Matrix Matrix Multiplications

2 For Multi-TTM Computation

- Simulated Experiments
- Conclusion

Multi-TTM Computation

- $\mathcal{X} \in \mathbb{R}^{n_1 \times \dots \times n_d}$: input tensor, $\mathcal{Y} \in \mathbb{R}^{r_1 \times \dots \times r_d}$: output tensor
- $\mathbf{A}^{(k)} \in \mathbb{R}^{n_k \times r_k}$: factor matrix of the k th mode

2-dimensional Multi-TTM ($\mathcal{Y} = \mathcal{X} \times_1 \mathbf{A}^{(1)\top} \times_2 \mathbf{A}^{(2)\top}$)

TTM-in-Sequence approach (used in Tucker-MPI): $\mathbf{Y} = \mathbf{A}^{(1)\top} \mathbf{X} \mathbf{A}^{(2)}$

All-at-Once approach (our contribution):

for $n'_1 = 1:n_1$, for $n'_2 = 1:n_2$

for $r'_1 = 1:r_1$, for $r'_2 = 1:r_2$

$$\mathcal{Y}(r'_1, r'_2) = \mathcal{Y}(r'_1, r'_2) + \left(\mathcal{X}(n'_1, n'_2) * \mathbf{A}^{(1)}(n'_1, r'_1) * \mathbf{A}^{(2)}(n'_2, r'_2) \right)$$

3-dimensional Multi-TTM ($\mathcal{Y} = \mathcal{X} \times_1 \mathbf{A}^{(1)\top} \times_2 \mathbf{A}^{(2)\top} \times_3 \mathbf{A}^{(3)\top}$)

All-at-Once 3-dimensional Multi-TTM Computation

for $n'_1 = 1:n_1$, for $n'_2 = 1:n_2$, for $n'_3 = 1:n_3$

for $r'_1 = 1:r_1$, for $r'_2 = 1:r_2$, for $r'_3 = 1:r_3$

$\mathcal{Y}(r'_1, r'_2, r'_3) = \mathcal{Y}(r'_1, r'_2, r'_3)$

$+ \left(\mathcal{X}(n'_1, n'_2, n'_3) * \mathbf{A}^{(1)}(n'_1, r'_1) * \mathbf{A}^{(2)}(n'_2, r'_2) * \mathbf{A}^{(3)}(n'_3, r'_3) \right)$

$$\Delta = \begin{bmatrix} \mathbf{I}_{3 \times 3} & \mathbf{1}_3 & \mathbf{0}_3 \\ \mathbf{I}_{3 \times 3} & \mathbf{0}_3 & \mathbf{1}_3 \end{bmatrix}$$

- Total number of 4 – array operations = $n_1 r_1 n_2 r_2 n_3 r_3$
- Δ is not full rank
 - Consider each vector \mathbf{x} such that $\Delta \cdot \mathbf{x} = \mathbf{1}$, \mathbf{x} is of the form $[a \ a \ a \ 1-a \ 1-a]^\top$ and $0 \leq a \leq 1$
- $\phi_{\mathcal{X}}, \phi_{\mathcal{Y}}$: tensor projections & ϕ_1, ϕ_2, ϕ_3 : matrix projections
- From HBL, $\phi_{\mathcal{X}}^{1-a} \phi_{\mathcal{Y}}^{1-a} \phi_1^a \phi_2^a \phi_3^a \geq \text{Amount of computations}$

Solving Optimization Problem to Compute Lower Bounds

- Select a processor which 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.

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

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

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

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

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

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

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.

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

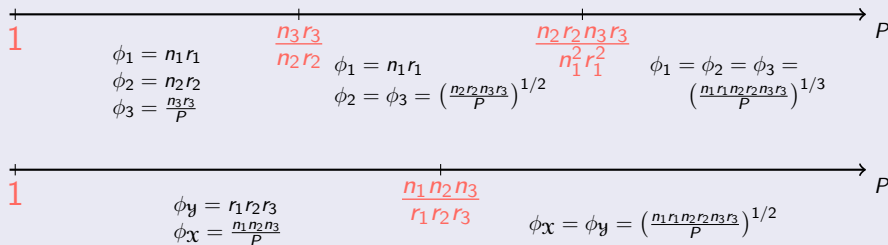
$$\frac{n_1 n_2 n_3}{P} \leq \phi_x \leq n_1 n_2 n_3$$

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

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 solutions for both parts using Lagrange multipliers (optimality can be proven using KKT conditions)

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

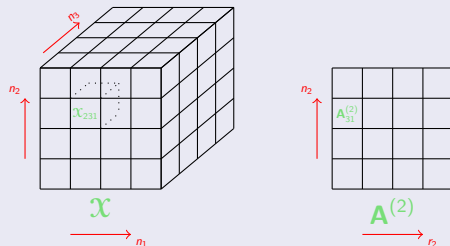


$$\text{Communication lower bound} = \phi_x + \phi_y + \phi_1 + \phi_2 + \phi_3 - \frac{n_1 n_2 n_3 + r_1 r_2 r_3 + n_1 r_1 + n_2 r_2 + n_3 r_3}{P}$$

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 $\mathbf{A}_{31}^{(2)} = \mathbf{A}^{(2)}(2\frac{n_2}{p_2} + 1 : 3\frac{n_2}{p_2}, 1 : \frac{r_2}{q_2})$ 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} , $\mathbf{A}^{(1)}$, $\mathbf{A}^{(2)}$, $\mathbf{A}^{(3)}$, $p_1 \times p_2 \times p_3 \times q_1 \times q_2 \times q_3$ logical processor grid

Ensure: \mathcal{Y} such that $\mathcal{Y} = \mathcal{X} \times_1 \mathbf{A}^{(1)\top} \times_2 \mathbf{A}^{(2)\top} \times_3 \mathbf{A}^{(3)\top}$

- 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: $\mathbf{A}_{p'_1 q'_1}^{(1)} = \text{All-Gather}(\mathbf{A}^{(1)}, (p'_1, *, *, q'_1, *, *))$
- 5: $\mathbf{A}_{p'_2 q'_2}^{(2)} = \text{All-Gather}(\mathbf{A}^{(2)}, (*, p'_2, *, *, q'_2, *))$
- 6: $\mathbf{A}_{p'_3 q'_3}^{(3)} = \text{All-Gather}(\mathbf{A}^{(3)}, (*, *, p'_3, *, *, q'_3))$
- 7: //Perform local Multi-TTM computation in a temporary tensor \mathcal{T}
- 8: $\mathcal{T} = \text{Local-Multi-TTM}(\mathcal{X}_{p'_1 p'_2 p'_3}, \mathbf{A}_{p'_1 q'_1}^{(1)}, \mathbf{A}_{p'_2 q'_2}^{(2)}, \mathbf{A}_{p'_3 q'_3}^{(3)})$
- 9: //Reduce-scatter the output tensor in $\mathcal{Y}_{q'_1 q'_2 q'_3}$
- 10: $\text{Reduce-Scatter}(\mathcal{Y}_{q'_1 q'_2 q'_3}, \mathcal{T}, (*, *, *, q'_1, q'_2, q'_3))$

Cost Analysis of our Algorithm

- Total amount of 4-array operations per processor = $\frac{n_1 r_1 n_2 r_2 n_3 r_3}{p_1 q_1 p_2 q_2 p_3 q_3} = \frac{n_1 r_1 n_2 r_2 n_3 r_3}{P}$
- Data transfers happen only in All-Gather and Reduce-Scatter collective operations
 - Cost on Q processors is $(1 - \frac{1}{Q})w$, w is the amount of total data after All-Gather or before Reduce-Scatter operation
- Total data transfers = $\frac{n_1 n_2 n_3}{p_1 p_2 p_3} + \frac{r_1 r_2 r_3}{q_1 q_2 q_3} + \frac{n_1 r_1}{p_1 q_1} + \frac{n_2 r_2}{p_2 q_2} + \frac{n_3 r_3}{p_3 q_3} - \frac{n_1 n_2 n_3 + r_1 r_2 r_3 + n_1 r_1 + n_2 r_2 + n_3 r_3}{P}$
- p_1, p_2, p_3, q_1, q_2 , and q_3 can be obtained based on lower bounds (Not today)

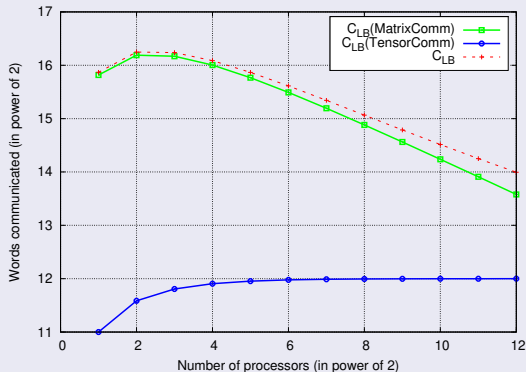
1 For Matrix Matrix Multiplications

2 For Multi-TTM Computation

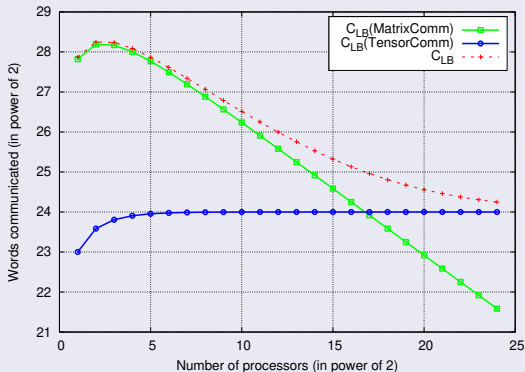
- Simulated Experiments
- Conclusion

Communication Lower Bound (C_{LB}) distributions

$$n_1 = n_2 = n_3 = 2^{12}, r_1 = r_2 = r_3 = 2^4$$



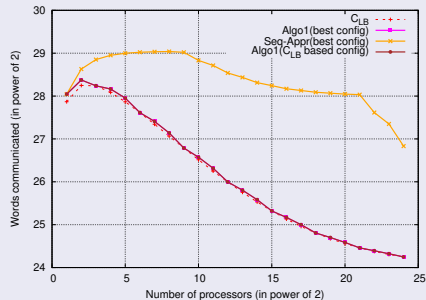
$$n_1 = n_2 = n_3 = 2^{20}, r_1 = r_2 = r_3 = 2^8$$



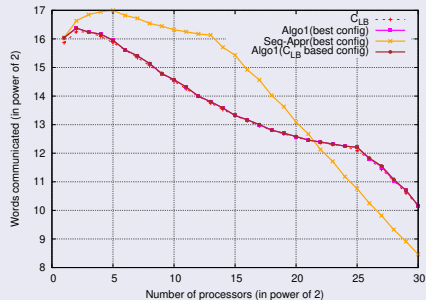
- Matrix communication costs dominate when P is much less than $\frac{n_1 n_2 n_3}{r_1 r_2 r_3}$

Performance Comparison of Our Algorithm

$$n_1 = n_2 = n_3 = 2^{20}, r_1 = r_2 = r_3 = 2^8$$



$$n_1 = n_2 = n_3 = 2^{12}, r_1 = r_2 = r_3 = 2^4$$



- C_{LB} : Communication lower bound, $Algo1(best\ config)$: Our algorithm with the best partition, $Algo1(C_{LB}\ based\ config)$: Our algorithm with the partition based on the lower bound, $Seq-Appr(best\ config)$: Multi-TTM computation used in Tucker-MPI with the best partition
- Our algorithm communicates much less than the approach in Tucker-MPI

- 1 For Matrix Matrix Multiplications
- 2 For Multi-TTM Computation
 - Simulated Experiments
 - Conclusion

Conclusion and Future Work

Conclusion

- Lower bounds for All-at-Once Multi-TTM
- Communication Optimal Algorithm for Multi-TTM
- Comparison of our algorithm with the TTM-in-Sequence Algorithm
- Our algorithm communicates much less data than TTM-in-Sequence for $P < \frac{n_1 n_2 n_3}{r_1 r_2 r_3}$

Future Work

- Detailed study of what scenarios are favorable for our algorithm
- Combine both All-at-Once and TTM-in-Sequence approaches to minimize communication lower bounds
- Extend our work for the full Tucker-decomposition algorithms