



# Communication Lower Bounds and Optimal Algorithms for Multiple Tensor-Times-Matrix Computation

Hussam Al Daas, Grey Ballard, Laura Grigori, Suraj Kumar, Kathryn Rouse

## ► To cite this version:

Hussam Al Daas, Grey Ballard, Laura Grigori, Suraj Kumar, Kathryn Rouse. Communication Lower Bounds and Optimal Algorithms for Multiple Tensor-Times-Matrix Computation. SIAM Journal on Matrix Analysis and Applications, 2024, 45 (1), pp.450-477. 10.1137/22M1510443 . hal-03950359

**HAL Id: hal-03950359**

**<https://inria.hal.science/hal-03950359>**

Submitted on 21 Jan 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# COMMUNICATION LOWER BOUNDS AND OPTIMAL ALGORITHMS FOR MULTIPLE TENSOR-TIMES-MATRIX COMPUTATION

HUSSAM AL DAAS\*, GREY BALLARD<sup>†</sup>, LAURA GRIGORI<sup>‡</sup>, SURAJ KUMAR<sup>‡</sup>, AND KATHRYN ROUSE<sup>§</sup>

**Abstract.** Multiple Tensor-Times-Matrix (Multi-TTM) is a key computation in algorithms for computing and operating with the Tucker tensor decomposition, which is frequently used in multidimensional data analysis. We establish communication lower bounds that determine how much data movement is required to perform the Multi-TTM computation in parallel. The crux of the proof relies on analytically solving a constrained, nonlinear optimization problem. We also present a parallel algorithm to perform this computation that organizes the processors into a logical grid with twice as many modes as the input tensor. We show that with correct choices of grid dimensions, the communication cost of the algorithm attains the lower bounds and is therefore communication optimal. Finally, we show that our algorithm can significantly reduce communication compared to the straightforward approach of expressing the computation as a sequence of tensor-times-matrix operations.

**Key words.** Communication lower bounds, Multi-TTM, Tensor computations, Parallel algorithms, HBL-inequalities

**MSC codes.** 15A69, 68Q17, 68Q25, 68W10, 68W15, 68W40

**1. Introduction.** The Tucker tensor decomposition is a low-rank representation or approximation that enables significant compression of multidimensional data. The Tucker format consists of a core tensor, which is much smaller than the original data tensor, along with a factor matrix for each mode, or dimension, of the data. Computations involving Tucker-format tensors, such as tensor inner products, often require far fewer operations than with their full-format, dense representations. As a result, the Tucker decomposition is often used as a dimensionality reduction technique before other types of analysis are done, including computing a CP decomposition [8], for example.

A 3-way Tucker-format tensor can be expressed using the tensor notation  $\mathcal{T} = \mathcal{G} \times_1 \mathbf{A}^{(1)} \times_2 \mathbf{A}^{(2)} \times_3 \mathbf{A}^{(3)}$ , where  $\mathcal{G}$  is the 3-way core tensor,  $\mathbf{A}^{(n)}$  is a tall-skinny factor matrix corresponding to mode  $n$ , and  $\times_n$  denotes the tensor-times-matrix (TTM) operation in the  $n$ th mode [17]. Here,  $\mathcal{T}$  is the full-format representation of the tensor that can be constructed explicitly by performing multiple TTM operations. We call this collective operation the Multi-TTM computation, which is the focus of this work.

Multi-TTM is a fundamental computation in the context of Tucker-format tensors. When the Tucker decomposition is used as a data compression tool, Multi-TTM is exactly the decompression operation, which is necessary when the full format is required for visualization [18], for example. In the case of full decompression, the input tensor is small and the output tensor is large. One of the quasi-optimal algorithms for computing the Tucker decomposition is the Truncated Higher-Order SVD algorithm [27, 19], in which each factor matrix is computed as the leading left singular vectors of a matrix unfolding of the tensor. In this algorithm, the smaller core tensor is computed via Multi-TTM involving the larger data tensor and the computed

\*Rutherford Appleton Laboratory, Didcot, Oxfordshire, UK (hussam.al-daas@stfc.ac.uk).

<sup>†</sup>Wake Forest University, Winston-Salem, NC, USA (ballard@wfu.edu).

<sup>‡</sup>Inria Paris, France (laura.grigori@inria.fr, suraj.kumar@inria.fr).

<sup>§</sup>Inmar Intelligence, Winston-Salem, NC, USA (kathryn.rouse@inmar.com).

factor matrices. When the computational costs of the matrix SVDs are reduced using randomization, Multi-TTM becomes the overwhelming bottleneck computation [22, 25].

Since the overall size of multidimensional data grows quickly, there have been many recent efforts to parallelize the computation of the Tucker decomposition and the operations on Tucker-format tensors [2, 9, 21, 11, 4]. There has also been recent progress in establishing lower bounds on the communication costs of parallel algorithms for tensor computations, including the Matricized-Tensor Times Khatri-Rao product (MTTKRP) [5, 6, 28] and symmetric tensor contractions [24]. However, to our knowledge, no communication lower bounds have been previously established for computations relating to Tucker-format tensors. In this work, we prove communication lower bounds for a class of Multi-TTM algorithms. Additionally, we provide a parallel algorithm that attains the lower bound to within a constant factor and is therefore communication optimal.

To minimize the number of arithmetic operations in a Multi-TTM computation, the TTM operations should be performed in sequence, forming temporary intermediate tensors after each step. One of the key observations of this work is that when Multi-TTM is performed in parallel, this approach may communicate more data than necessary, even if communication-optimal algorithms are used for each individual TTM. By considering the Multi-TTM computation as a whole, we can devise parallel algorithms that can communicate less than this TTM-in-Sequence approach, often with negligible increase in computation. Our proposed algorithm provides greatest benefit when the input and output tensors vary greatly in size.

The main contributions of this paper are to

- establish communication lower bounds for the parallel load balanced Multi-TTM computation;
- propose a communication optimal parallel algorithm;
- show that in many typical scenarios, the straightforward approach based on a sequence of TTM operations communicates more than performing Multi-TTM as a whole.

The rest of the paper is organized as follows. Section 2 describes previous work on communication lower bounds for matrix multiplication and some tensor operations. In Section 3, we present our notations and preliminaries for the general Multi-TTM computation. To reduce the complexity of notations, we first focus on 3-dimensional Multi-TTM computation for which we present communication lower bounds and a communication optimal algorithm in Section 4 and Section 5, respectively. In Section 6, we validate the optimality of the proposed algorithm and show that it significantly reduces communication compared to the TTM-in-Sequence approach with negligible increase in computation in many practical cases. We present our general results in Sections 7 and 8, and propose conclusions and perspectives in Section 9.

**2. Related Work.** A number of studies have focused on communication lower bounds for matrix multiplication, starting with the work by Hong and Kung [14] to determine the minimum number of I/O operations for sequential matrix multiplication using red-blue pebble game. Irony et al. [15] extended this work for the parallel case. Demmel et al. [13] studied memory independent communication lower bounds for rectangular matrix multiplication based on aspect ratios of matrices. Recently, Smith et al. [23] and Al Daas et al. [1] have tightened communication lower bounds for matrix multiplication. Ballard et al. [3] extended communication lower bounds of the matrix multiplication for any computations that can be written as 3 nested loops.

Christ et al. [12] generalized the method to prove communication lower bounds of 3 nested loop computations for arbitrary loop nesting. We apply their approach to our Multi-TTM definition.

There is limited work on communication lower bounds for tensor operations. Solomonik et al. [24] proposed communication lower bounds for symmetric tensor contraction algorithms. Ballard et al. [5] proposed communication lower bounds for MTTKRP computation with cubical tensors. This work is extended in [6] to handle varying tensor dimensions. A sequential lower bound for tile-based MTTKRP algorithms is proved by Ziogas et al. [28]. We use some results from [5, 6] to prove communication lower bounds for Multi-TTM.

**3. Notations and Preliminaries.** In this section, we present our notations and basic lemmas for  $d$ -dimensional Multi-TTM computation. In Sections 4 to 6, we focus on  $d = 3$ , i.e.,  $\mathbf{Y} = \mathbf{X} \times_1 \mathbf{A}^{(1)\top} \times_2 \mathbf{A}^{(2)\top} \times_3 \mathbf{A}^{(3)\top}$ . We present our general results in Sections 7 and 8.

We use boldface uppercase Euler script letters to denote tensors ( $\mathbf{X}$ ) and boldface uppercase letters with superscripts to denote matrices ( $\mathbf{A}^{(1)}$ ). We use lowercase letters with subscripts to denote sizes ( $n_1$ ) and add the prime symbol to them to denote the indices ( $n'_1$ ). We use one-based indexing throughout and  $[d]$  to denote the set  $\{1, 2, \dots, d\}$ . To improve the presentation, we denote the product of elements having the same lowercase letter with all subscripts by the lowercase letter only ( $n_1 \cdots n_d$  by  $n$  and  $r_1 \cdots r_d$  by  $r$ ). We denote the product of the  $i$  rightmost terms with the capital letter with subscript  $i$ ,  $N_i = \prod_{j=d-i+1}^d n_j$  and  $R_i = \prod_{j=d-i+1}^d r_j$ , thus  $n = N_d$ , and  $n_d = N_1$ .

Let  $\mathbf{Y} \in \mathbb{R}^{r_1 \times \cdots \times r_d}$  be the  $d$ -mode output tensor,  $\mathbf{X} \in \mathbb{R}^{n_1 \times \cdots \times n_d}$  be the  $d$ -mode input tensor, and  $\mathbf{A}^{(k)} \in \mathbb{R}^{n_k \times r_k}$  be the matrix of the  $k$ th mode. Then the Multi-TTM computation can be represented as  $\mathbf{Y} = \mathbf{X} \times_1 \mathbf{A}^{(1)\top} \cdots \times_d \mathbf{A}^{(d)\top}$ . Without loss of generality and to simplify notation, we consider that the input tensor  $\mathbf{X}$  is larger than the output tensor  $\mathbf{Y}$ , or  $n \geq r$ . This corresponds to computing the core tensor of a Tucker decomposition given computed factor matrices, for example. However, the opposite relationship where the output tensor is larger (e.g.,  $\mathbf{X} = \mathbf{Y} \times_1 \mathbf{A}^{(1)} \cdots \times_d \mathbf{A}^{(d)}$ ) is also an important use case, corresponding to forming an explicit representation of a (sub-)tensor of a Tucker-format tensor. Our results extend straightforwardly to this case. We also assume without loss of generality that the tensor modes are ordered in such a way that  $n_1 r_1 \leq n_2 r_2 \leq \cdots \leq n_d r_d$ .

**DEFINITION 3.1.** Let  $\mathbf{X}$  be an  $n_1 \times \cdots \times n_d$  tensor,  $\mathbf{Y}$  be an  $r_1 \times \cdots \times r_d$  tensor, and  $\mathbf{A}^{(j)}$  be an  $n_j \times r_j$  matrix for  $j \in [d]$ . Multi-TTM computes

$$\mathbf{Y} = \mathbf{X} \times_1 \mathbf{A}^{(1)\top} \cdots \times_d \mathbf{A}^{(d)\top}$$

where for each  $(r'_1, \dots, r'_d) \in [r_1] \times \cdots \times [r_d]$ ,

$$(3.1) \quad \mathbf{Y}(r'_1, \dots, r'_d) = \sum_{\{n'_k \in [n_k]\}_{k \in [d]}} \mathbf{X}(n'_1, \dots, n'_d) \prod_{j \in [d]} \mathbf{A}^{(j)}(n'_j, r'_j).$$

Let us consider an example when  $d = 2$ . In this scenario, the input and output tensors are in fact matrices  $\mathbf{X}$ ,  $\mathbf{Y}$ , and  $\mathbf{Y} = \mathbf{A}^{(1)\top} \mathbf{X} \mathbf{A}^{(2)}$ . As mentioned earlier, Multi-TTM computation can be performed as a sequence of TTM operations, in this case two matrix multiplications. However, we define the Multi-TTM to perform all the products at once for each term of the summation of (3.1). Our definition comes at

greater arithmetic cost, as partial  $(d+1)$ -ary multiplies are not computed and reused, but we will see that this approach can reduce communication cost. We describe how the extra computation can often be reduced to a negligible cost in Subsection 5.1.

We can write pseudocode for the Multi-TTM with the following:

$$\begin{aligned} &\text{for } n'_1 = 1:n_1, \dots, \text{ for } n'_d = 1:n_d, \\ &\quad \text{for } r'_1 = 1:r_1, \dots, \text{ for } r'_d = 1:r_d, \\ &\quad \quad \mathcal{Y}(r'_1, \dots, r'_d) += \mathcal{X}(n'_1, \dots, n'_d) \cdot \mathbf{A}^{(1)}(n'_1, r'_1) \cdot \dots \cdot \mathbf{A}^{(N)}(n'_d, r'_d) \end{aligned}$$

**DEFINITION 3.2.** *A parallel atomic Multi-TTM algorithm computes each term of the summation of (3.1) atomically on a unique processor, but it can distribute the  $nr$  terms over processors in any way.*

Here atomic computation of a single  $(d+1)$ -ary multiplication for a parallel algorithm means that all the multiplications of this operation are performed on only one processor, i.e., all  $d+1$  inputs are accessed on that processor in order to compute the single output value. This assumption is necessary for our communication lower bounds. Processors can reorganize their local atomic operations to reduce computational costs without changing the communication or violating parallel atomicity. However it is reasonable for an algorithm to break this assumption in order to improve arithmetic costs by reusing partial results across processors, and we compare against such algorithms in Section 6.

**3.1. Parallel Computation Model.** We consider that the computation is distributed across  $P$  processors. Each processor has its own local memory and is connected to all other processors via a fully connected network. Every processor can operate on data in its local memory and must communicate to access data of other processors. Hence, communication refers to send and receive operations that transfer data from local memory to the network and vice-versa. Communication cost mainly depends on two factors – the amount of data communicated (bandwidth cost) and the number of messages (latency cost). Latency cost is dominated by bandwidth cost for computations involving large messages, so we focus on bandwidth cost in this work and refer it as communication cost throughout the text. We assume the links of the network are bidirectional and that the communication cost is independent of the number of pairs of processors that are communicating. Each processor can send and receive at most one message at the same time. In our model, the communication cost of an algorithm refers to the cost along the critical path.

**3.2. Existing Results.** Our work relies on two fundamental results. The first, a geometric result on lattices, allows us to relate the volume of computation to the amount of data accessed by determining the maximum data reuse. The result is a specialization of the Hölder-Brascamp-Lieb inequalities [7]. This result has previously been used to derive lower bounds for tensor computations [5, 6, 12, 16] in a similar way to the use of the Loomis-Whitney inequality [20] in derivations of communication lower bounds for linear algebra [3]. The result is proved in [12], but we use the statement from [5, Lemma 4.1]. Here  $\mathbf{1}$  represents a vector of all ones.

**LEMMA 3.3.** *Consider any positive integers  $\ell$  and  $m$  and any  $m$  projections  $\phi_j : \mathbb{Z}^\ell \rightarrow \mathbb{Z}^{\ell_j}$  ( $\ell_j \leq \ell$ ), each of which extracts  $\ell_j$  coordinates  $S_j \subseteq [\ell]$  and forgets the  $\ell - \ell_j$  others. Define  $\mathcal{C} = \{\mathbf{s} \in [0, 1]^m : \mathbf{\Delta} \cdot \mathbf{s} \geq \mathbf{1}\}$ , where the  $\ell \times m$  matrix  $\mathbf{\Delta}$  has entries*

$\Delta_{i,j} = 1$  if  $i \in S_j$  and  $\Delta_{i,j} = 0$  otherwise. If  $[s_1 \cdots s_m]^\top \in \mathcal{C}$ , then for all  $F \subseteq \mathbb{Z}^\ell$ ,

$$|F| \leq \prod_{j \in [m]} |\phi_j(F)|^{s_j}.$$

The second result, a general constrained optimization problem, allows us to cast the communication cost of an algorithm as the objective function in an optimization problem where the constraints are imposed by properties of the computation within the algorithm. A version of the result is proved in [6, Lemma 5.1] and used to derive the general communication lower bound for MTTKRP.

**THEOREM 3.4.** *Consider the constrained optimization problem:*

$$\min \sum_{j \in [d]} x_j$$

such that

$$\frac{nr}{P} \leq \prod_{j \in [d]} x_j \quad \text{and} \quad 0 \leq x_j \leq k_j \quad \text{for all} \quad 1 \leq j \leq d$$

for some positive constants  $k_1 \leq k_2 \leq \cdots \leq k_d$  with  $\prod_{j \in [d]} k_j = nr$ . Then the minimum value of the objective function is

$$I(K_I/P)^{1/I} + \sum_{j \in [d-I]} k_j$$

where we use the notation  $K_I = \prod_{j=d-I+1}^d k_j$  and  $1 \leq I \leq d$  is defined such that

$$\begin{aligned} k_j &< (K_{d-j+1}/P)^{1/(d-j+1)} \quad \text{for } 1 \leq j \leq d-I, \\ k_\ell &\geq (K_{d-\ell+1}/P)^{1/(d-\ell+1)} \quad \text{for } d-I < \ell \leq d. \end{aligned}$$

The minimum is achieved at the point  $\mathbf{x}^*$  defined by  $x_j^* = k_j$  for  $1 \leq j \leq d-I$ ,  $x_\ell^* = (K_I/P)^{1/I}$  for  $d-I < \ell \leq d$ .

While Theorem 3.4 can be straightforwardly derived from the previous work, we provide a proof in Appendix A for completeness. We represent it in this form to be directly applicable to all the constrained optimization problems in this paper. The lower bound and constraint on the products of the upper bounds are derived from the Multi-TTM computation. The additional constraint on the product of the upper bounds implies that there is always a feasible solution to the optimization problem for  $P \geq 1$ .

We can note that the  $d$  conditions are examined to determine the value  $I$ . We calculate the ranges of  $P$  for each  $I$  based on these conditions in Corollaries 4.1, 4.2 and 7.1.

**4. Lower Bounds for 3-dimensional Multi-TTM.** We obtain the lower bound results for 3D tensors in this section, presented as Theorem 4.3. The lower bound is independent of the size of the local memory of each processor, similar to previous results for matrix multiplication [1, 13] and MTTKRP [5, 6], and it varies with respect to the number of processors  $P$  relative to the matrix and tensor dimensions of the problem.

The crux of the proof considers a single processor that performs  $1/P$ th of the computation and has access to  $1/P$ th of the data. Finding the lower bound on the data

that processor must communicate is reduced to solving a constrained optimization problem: we seek to minimize the number of elements of the matrices and tensors that the processor must access in order to execute its computation subject to structure constraints of Multi-TTM. The most important constraint derives from Lemma 3.3, which relates a subset of the computation within a Multi-TTM algorithm to the data it requires. The other constraints provide upper bounds on the data required from each array. The upper bounds are necessary to establish the tightest lower bounds in the cases where  $P$  is small. We show that the optimization problem can be separated into two independent problems, one for the matrix data and one for the tensor data. Corollaries 4.1 and 4.2 state the two constrained optimization problems along with their analytic solutions, both of which follow from Theorem 3.4. That is, setting  $d = 3$ ,  $k_1 = n_1 r_1$ ,  $k_2 = n_2 r_2$  and  $k_3 = n_3 r_3$  in Theorem 3.4, we obtain Corollary 4.1. Similarly, setting  $d = 2$  with  $k_1 = r$  and  $k_2 = n$ , we obtain Corollary 4.2. We recall here that  $r = r_1 r_2 r_3$  and  $n = n_1 n_2 n_3$ .

COROLLARY 4.1. *Consider the following optimization problem:*

$$\min_{x,y,z} x + y + z$$

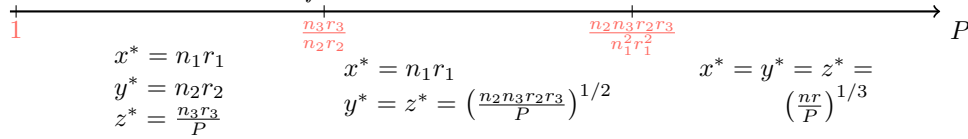
such that

$$\begin{aligned} \frac{nr}{P} &\leq xyz \\ 0 &\leq x \leq n_1 r_1 \\ 0 &\leq y \leq n_2 r_2 \\ 0 &\leq z \leq n_3 r_3, \end{aligned}$$

where  $n_1 r_1 \leq n_2 r_2 \leq n_3 r_3$ , and  $n_1, n_2, n_3, r_1, r_2, r_3, P \geq 1$ . The optimal solution  $(x^*, y^*, z^*)$  depends on the relative values of the constraints, yielding three cases:

1. if  $P < \frac{n_3 r_3}{n_2 r_2}$ , then  $x^* = n_1 r_1$ ,  $y^* = n_2 r_2$ ,  $z^* = \frac{n_3 r_3}{P}$ ;
2. if  $\frac{n_3 r_3}{n_2 r_2} \leq P < \frac{n_2 n_3 r_2 r_3}{n_1^2 r_1^2}$ , then  $x^* = n_1 r_1$ ,  $y^* = z^* = \left( \frac{n_2 n_3 r_2 r_3}{P} \right)^{\frac{1}{2}}$ ;
3. if  $\frac{n_2 n_3 r_2 r_3}{n_1^2 r_1^2} \leq P$ , then  $x^* = y^* = z^* = \left( \frac{nr}{P} \right)^{\frac{1}{3}}$ ;

which can be visualized as follows.



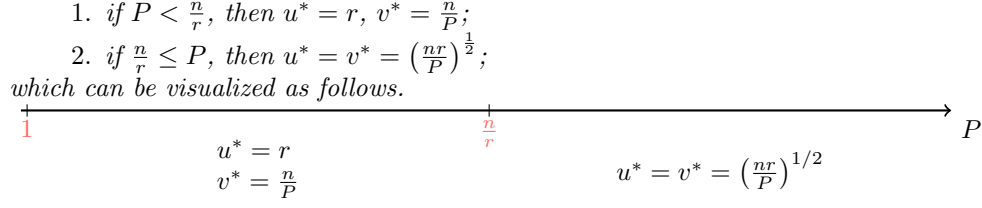
COROLLARY 4.2. *Consider the following optimization problem:*

$$\min_{u,v} u + v$$

such that

$$\begin{aligned} \frac{nr}{P} &\leq uv \\ 0 &\leq u \leq r \\ 0 &\leq v \leq n, \end{aligned}$$

where  $n \geq r$ , and  $n, r, P \geq 1$ . The optimal solution  $(u^*, v^*)$  depends on the relative values of the constraints, yielding two cases:



**4.1. Communication Lower Bounds for Multi-TTM.** We now state the lower bounds for 3-dimensional Multi-TTM. After this, we also present a corollary for cubical tensors.

**THEOREM 4.3.** *Any computationally load balanced atomic Multi-TTM algorithm that starts and ends with one copy of the data distributed across processors involving 3D tensors with dimensions  $n_1, n_2, n_3$  and  $r_1, r_2, r_3$  performs at least  $A + B - \left(\frac{n}{P} + \frac{r}{P} + \sum_{j=1}^3 \frac{n_j r_j}{P}\right)$  sends or receives where*

$$A = \begin{cases} n_1 r_1 + n_2 r_2 + \frac{n_3 r_3}{P} & \text{if } P < \frac{n_3 r_3}{n_2 r_2} \\ n_1 r_1 + 2 \left( \frac{n_2 n_3 r_2 r_3}{P} \right)^{\frac{1}{2}} & \text{if } \frac{n_3 r_3}{n_2 r_2} \leq P < \frac{n_2 n_3 r_2 r_3}{n_1^2 r_1^2} \\ 3 \left( \frac{nr}{P} \right)^{\frac{1}{3}} & \text{if } \frac{n_2 n_3 r_2 r_3}{n_1^2 r_1^2} \leq P \end{cases}$$

$$B = \begin{cases} r + \frac{n}{P} & \text{if } P < \frac{n}{r} \\ 2 \left( \frac{nr}{P} \right)^{\frac{1}{2}} & \text{if } \frac{n}{r} \leq P. \end{cases}$$

*Proof.* Let  $F$  be the set of loop indices associated with the 4-ary multiplication performed by a processor. As we assumed the algorithm is computationally load balanced,  $|F| = nr/P$ . We define  $\phi_{\mathbf{x}}(F)$ ,  $\phi_{\mathbf{y}}(F)$  and  $\phi_j(F)$  to be the projections of  $F$  onto the indices of the arrays  $\mathbf{X}, \mathbf{Y}$ , and  $\mathbf{A}^{(j)}$  for  $1 \leq j \leq 3$  which correspond to the elements of the array that must be accessed by the processor.

We use Lemma 3.3 to obtain a lower bound on the number of array elements that must be accessed by the processor. The computation involves 5 arrays (2 tensors and 3 matrices) with 6 loop indices (see the atomic Multi-TTM definition in Section 3), hence the  $6 \times 5$  matrix corresponding to the projections above is given by

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

Here  $\mathbf{1}_3$  and  $\mathbf{0}_3$  represent the 3-dimensional vectors of all ones and zeros, respectively, and  $\mathbf{I}_{3 \times 3}$  represents the  $3 \times 3$  identity matrix. We recall from Lemma 3.3 that  $\Delta_{i,j} = 1$  if loop index  $i$  is used to access array  $j$  and  $\Delta_{i,j} = 0$  otherwise. The first three columns of  $\Delta$  correspond to matrices and the remaining two columns correspond to tensors. In this case, we have

$$\mathcal{C} = \{\mathbf{s} \in [0, 1]^5 : \Delta \cdot \mathbf{s} \geq \mathbf{1}\}.$$

Recall that  $\mathbf{1}$  represents a vector of all ones. Here  $\Delta$  is not full rank, therefore, we consider all vectors  $\mathbf{v} \in \mathcal{C}$  such that  $\Delta \cdot \mathbf{v} = \mathbf{1}$ . Such a vector  $\mathbf{v}$  is of the form  $[a \ a \ a \ 1-a \ 1-a]$  where  $0 \leq a \leq 1$ . Therefore, we obtain

$$\frac{nr}{P} \leq \left( \prod_{j \in [3]} |\phi_j(F)| \right)^a (|\phi_{\mathbf{x}}(F)| |\phi_{\mathbf{y}}(F)|)^{1-a} \text{ for all } 0 \leq a \leq 1.$$



The above constraint is equivalent to  $\frac{nr}{P} \leq \prod_{j \in [3]} |\phi_j(F)|$  and  $\frac{nr}{P} \leq |\phi_{\mathbf{x}}(F)| |\phi_{\mathbf{y}}(F)|$ . To see this equivalence note that the forward direction is implied by setting  $a = 0$  and  $a = 1$ . For the opposite direction, taking the first of the two constraints to the power  $a$  and the second to the power  $1 - a$  then multiplying the two terms yields the original.

Clearly a projection onto an array cannot be larger than the array itself, thus  $|\phi_{\mathbf{x}}(F)| \leq n$ ,  $|\phi_{\mathbf{y}}(F)| \leq r$ , and  $|\phi_j(F)| \leq n_j r_j$  for  $1 \leq j \leq 3$ .

As the constraints related to projections of matrices and tensors are disjoint, we solve them separately and then sum the results to get a lower bound on the set of elements that must be accessed by the processor. We obtain a lower bound on  $A$ , the number of elements of the matrices that must be accessed by the processor by using Corollary 4.1, and a lower bound on  $B$ , the number of elements of the tensors that must be accessed by the processor by using Corollary 4.2. By summing both, we get the positive terms of the lower bound.

To bound the sends or receives, we consider how much data the processor could have had at the beginning or at the end of the computation. Assuming there is exactly one copy of the data at the beginning and at the end of the computation, there must exist a processor which has access to at most  $1/P$  of the elements of the arrays at the beginning or at the end of the computation. By employing the previous analysis, this processor must access  $A + B$  elements of the arrays, but can only have  $\frac{n}{P} + \frac{r}{P} + \sum_{j \in [3]} \frac{n_j r_j}{P}$  elements of the arrays stored. Thus it must perform the specified amount of sends or receives.  $\square$

We denote the lower bound of Theorem 4.3 by  $LB$  and use it extensively in Subsection 5.2 while analyzing the communication cost of our parallel algorithm.

We also state the result for 3-dimensional Multi-TTM computation with cubical tensors, which is a direct application of Theorem 4.3 with  $n_1 = n_2 = n_3 = n^{\frac{1}{3}}$  and  $r_1 = r_2 = r_3 = r^{\frac{1}{3}}$ .

**COROLLARY 4.4.** *Any computationally load balanced atomic Multi-TTM algorithm that starts and ends with one copy of the data distributed across processors involving 3D cubical tensors with dimensions  $n^{\frac{1}{3}} \times n^{\frac{1}{3}} \times n^{\frac{1}{3}}$  and  $r^{\frac{1}{3}} \times r^{\frac{1}{3}} \times r^{\frac{1}{3}}$  (with  $n \geq r$ ) performs at least*

$$3 \left( \frac{nr}{P} \right)^{\frac{1}{3}} + r - \frac{3(nr)^{\frac{1}{3}} + r}{P}$$

*sends or receives when  $P < \frac{n}{r}$  and at least*

$$3 \left( \frac{nr}{P} \right)^{\frac{1}{3}} + 2 \left( \frac{nr}{P} \right)^{\frac{1}{2}} - \frac{n + 3(nr)^{\frac{1}{3}} + r}{P}$$

*send or receives when  $P \geq \frac{n}{r}$ .*

In particular, we note that the lower bound for cubical atomic Multi-TTM algorithms is smaller than that of a TTM-in-Sequence approach for many typical scenarios in the case  $P < n/r$ , as we discuss further in Section 6.

**5. Parallel Algorithm for 3-dimensional Multi-TTM.** We organize  $P$  processors into a 6-dimensional  $p_1 \times p_2 \times p_3 \times q_1 \times q_2 \times q_3$  logical processor grid. We arrange the grid dimensions such that  $p_1, p_2, p_3, q_1, q_2, q_3$  evenly distribute  $n_1, n_2, n_3, r_1, r_2, r_3$ , respectively. A processor coordinate is represented as  $(p'_1, p'_2, p'_3, q'_1, q'_2, q'_3)$ , where  $1 \leq p'_k \leq p_k$ ,  $1 \leq q'_k \leq q_k$  for  $k = 1, 2, 3$ . To be consistent with our notation, we denote  $p_1 p_2 p_3$  and  $q_1 q_2 q_3$  by  $p$  and  $q$ .

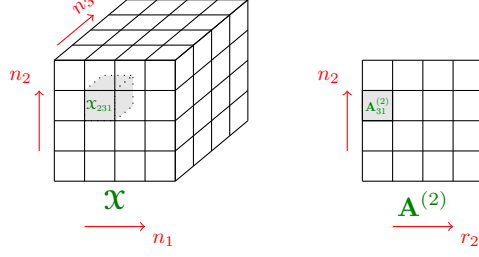


Fig. 1: Subtensor  $\mathcal{X}_{231}$  is distributed evenly among processors  $(2, 3, 1, *, *, *)$ . Similarly, submatrix  $\mathbf{A}_{31}^{(2)}$  is distributed evenly among processors  $(*, 3, *, *, 1, *)$ .

$\mathcal{X}_{p'_1 p'_2 p'_3}$  denotes the subtensor of  $\mathcal{X}$  owned by processors  $(p'_1, p'_2, p'_3, *, *, *)$ . Similarly,  $\mathcal{Y}_{q'_1 q'_2 q'_3}$  denotes the subtensor of  $\mathcal{Y}$  owned by processors  $(*, *, *, q'_1, q'_2, q'_3)$ .  $\mathbf{A}_{p'_1 q'_1}^{(1)}$ ,  $\mathbf{A}_{p'_2 q'_2}^{(2)}$  and  $\mathbf{A}_{p'_3 q'_3}^{(3)}$  denote submatrices of  $\mathbf{A}^{(1)}$ ,  $\mathbf{A}^{(2)}$  and  $\mathbf{A}^{(3)}$  owned by processors  $(p'_1, *, *, q'_1, *, *)$ ,  $(*, p'_2, *, *, q'_2, *)$  and  $(*, *, p'_3, *, *, q'_3)$ , respectively.

We impose that there is one copy of data in the system at the start and end of the computation, and every array is distributed evenly among the sets of processors whose coordinates are different for the corresponding dimensions of the variable. For example,  $\mathcal{X}_{111} = \mathcal{X}(1 : \frac{n_1}{p_1}, 1 : \frac{n_2}{p_2}, 1 : \frac{n_3}{p_3})$  is owned by processors  $(1, 1, 1, *, *, *)$ . Similarly,  $\mathbf{A}_{12}^{(1)} = \mathbf{A}^{(1)}(1 : \frac{n_1}{p_1}, \frac{r_1}{q_1} + 1 : 2\frac{r_1}{q_1})$  is owned by processors  $(1, *, *, 2, *, *)$ . We assume that data inside these sets of processors is also evenly distributed. For example, in the beginning, processor  $(1, 1, 1, 2, 1, 3)$  owns  $\frac{1}{P}$ th portion of each input variable:  $\frac{p_2}{P}$ th portion of  $\mathcal{X}_{111}$ ,  $\frac{p_1 q_1}{P}$ th portion of  $\mathbf{A}_{12}^{(1)}$ ,  $\frac{p_2 q_2}{P}$ th portion of  $\mathbf{A}_{11}^{(2)}$ , and  $\frac{p_3 q_3}{P}$ th portion of  $\mathbf{A}_{13}^{(3)}$ . Figure 1 illustrates examples of our data distribution model for two of the arrays.

---

**Algorithm 5.1** Parallel Atomic 3-dimensional 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  $\mathcal{X}$
  - 3:  $\mathcal{X}_{p'_1 p'_2 p'_3} = \text{All-Gather}(\mathcal{X}, (p'_1, p'_2, p'_3, *, *, *))$
  - 4: //All-gather input matrices
  - 5:  $\mathbf{A}_{p'_1 q'_1}^{(1)} = \text{All-Gather}(\mathbf{A}^{(1)}, (p'_1, *, *, q'_1, *, *))$
  - 6:  $\mathbf{A}_{p'_2 q'_2}^{(2)} = \text{All-Gather}(\mathbf{A}^{(2)}, (*, p'_2, *, *, q'_2, *))$
  - 7:  $\mathbf{A}_{p'_3 q'_3}^{(3)} = \text{All-Gather}(\mathbf{A}^{(3)}, (*, *, p'_3, *, *, q'_3))$
  - 8: //Local computations in a temporary tensor  $\mathcal{T}$
  - 9:  $\mathcal{T} = \text{Local-Multi-TTM}(\mathcal{X}_{p'_1 p'_2 p'_3}, \mathbf{A}_{p'_1 q'_1}^{(1)}, \mathbf{A}_{p'_2 q'_2}^{(2)}, \mathbf{A}_{p'_3 q'_3}^{(3)})$
  - 10: //Reduce-scatter the output tensor in  $\mathcal{Y}_{q'_1 q'_2 q'_3}$
  - 11:  $\text{Reduce-Scatter}(\mathcal{Y}_{q'_1 q'_2 q'_3}, \mathcal{T}, (*, *, *, q'_1, q'_2, q'_3))$
- 

Algorithm 5.1 presents a parallel algorithm to compute 3-dimensional Multi-TTM. When it completes,  $\mathcal{Y}_{q'_1 q'_2 q'_3}$  is distributed evenly among processors  $(*, *, *, *, *, *)$ .

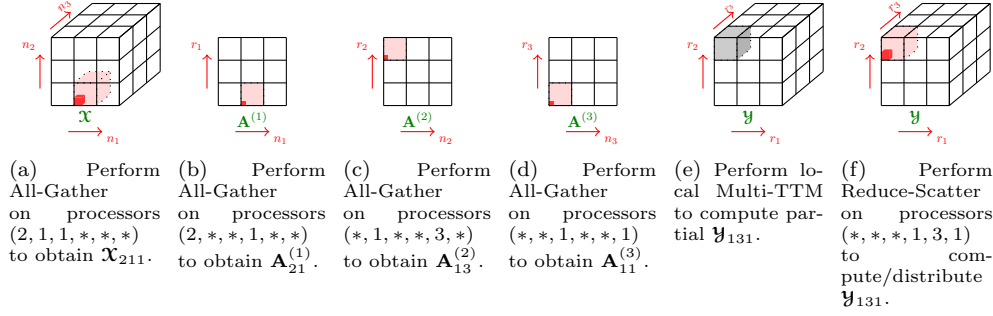


Fig. 2: S-eps-converted-to.pdf of Alg. 5.1 for processor (2, 1, 1, 1, 3, 1), where  $p_1 = p_2 = p_3 = q_1 = q_2 = q_3 = 3$ . Highlighted areas correspond to the data blocks on which the processor is operating. The dark red highlighting represents the input/output data initially/finally owned by the processor, and the light red highlighting corresponds to received/sent data from/to other processors in All-Gather/Reduce-Scatter collectives to compute  $\mathcal{Y}_{131}$ .

$q'_1, q'_2, q'_3$ ). Figure 2 shows the s-eps-converted-to.pdf of the algorithm for a single processor in a  $3 \times 3 \times 3 \times 3 \times 3 \times 3$  grid.

**5.1. Cost Analysis.** Now we analyze computation and communication costs of the algorithm. The dimension of the local tensor  $\mathcal{X}_{p'_1 p'_2 p'_3}$  is  $\frac{n_1}{p_1} \times \frac{n_2}{p_2} \times \frac{n_3}{p_3}$ , the dimension of the local matrix  $\mathbf{A}_{p'_k q'_k}^{(k)}$  is  $\frac{n_i}{p_i} \times \frac{r_i}{q_i}$  for  $i = 1, 2, 3$ , and the dimension of the temporary tensor  $\mathcal{T}$  is  $\frac{r_1}{q_1} \times \frac{r_2}{q_2} \times \frac{r_3}{q_3}$ .

The local Multi-TTM computation in Line 9 can be performed as a sequence of TTM operations to minimize the number of arithmetic operations. Assuming the TTM operations are performed in their order, first with  $\mathbf{A}^{(1)}$ , then with  $\mathbf{A}^{(2)}$ , and in the end with  $\mathbf{A}^{(3)}$ , then each processor performs  $2 \left( \frac{n_1 n_2 n_3 r_1}{p_1 p_2 p_3 q_1} + \frac{n_2 n_3 r_1 r_2}{p_2 p_3 q_1 q_2} + \frac{n_3 r_1 r_2 r_3}{p_3 q_1 q_2 q_3} \right)$  operations to perform the local computation.

Communication occurs only in the All-Gather and Reduce-Scatter collectives in Lines 3, 5, 6, 7 and 11. Each processor is involved in one All-Gather involving the input tensor, three All-Gathers involving input matrices and one Reduce-Scatter involving the output tensor. Lines 3, 5, 6, and 7 specify simultaneous  $\frac{P}{p}$ ,  $\frac{P}{p_1 q_1}$ ,  $\frac{P}{p_2 q_2}$  and  $\frac{P}{p_3 q_3}$  All-Gathers respectively, and Line 11 specifies simultaneous  $\frac{P}{q}$  Reduce-Scatters.

For simplicity of discussion, we consider that the number of processors involved in the collectives is a power of 2. We also assume that communication optimal collective algorithms are used. The optimal latency and bandwidth costs of both collectives on  $Q$  processors are  $\log(Q)$  and  $(1 - \frac{1}{Q})w$ , respectively, where  $w$  denotes the words of data in each processor after All-Gather or before Reduce-Scatter collective. Each processor also performs  $(1 - \frac{1}{Q})w$  computations for the Reduce-Scatter collective. We point the reader to [26, 10] for more details on efficient algorithms for collectives.

Hence the bandwidth costs of Lines 3, 5, 6, 7 in Alg. 5.1 are  $(1 - \frac{P}{p})\frac{n}{p}$ ,  $(1 - \frac{P}{p_1 q_1})\frac{n_1 r_1}{p_1 q_1}$ ,  $(1 - \frac{P}{p_2 q_2})\frac{n_2 r_2}{p_2 q_2}$ ,  $(1 - \frac{P}{p_3 q_3})\frac{n_3 r_3}{p_3 q_3}$  respectively to accomplish All-Gather operations, and the bandwidth cost of performing the Reduce-Scatter operation in Line 11

is  $(1 - \frac{q}{P})\frac{r}{q}$ . Thus the overall bandwidth cost of Alg. 5.1 for each processor is

$$(5.1) \quad \frac{n}{p} + \frac{n_1 r_1}{p_1 q_1} + \frac{n_2 r_2}{p_2 q_2} + \frac{n_3 r_3}{p_3 q_3} + \frac{r}{q} - \left( \frac{n + n_1 r_1 + n_2 r_2 + n_3 r_3 + r}{P} \right).$$

The latency costs of Lines 3, 5, 6, 7, 11 are  $\log(\frac{P}{p})$ ,  $\log(\frac{P}{p_1 q_1})$ ,  $\log(\frac{P}{p_2 q_2})$ ,  $\log(\frac{P}{p_3 q_3})$ ,  $\log(\frac{P}{q})$  respectively. Thus the overall latency cost of Alg. 5.1 for each processor is  $\log(\frac{P}{p}) + \log(\frac{P}{p_1 q_1}) + \log(\frac{P}{p_2 q_2}) + \log(\frac{P}{p_3 q_3}) + \log(\frac{P}{q}) = \log(\frac{P^5}{p^2 q^2}) = 3 \log(P)$ .

Due to the Reduce-Scatter operation, each processor also performs  $(1 - \frac{q}{P})\frac{r}{q}$  computations, which is asymptotically small compared to the computations of Line 9.

**5.2. Selection of  $p_i$  and  $q_i$  in Algorithm 5.1.** We must select the processor dimensions carefully such that Alg. 5.1 is communication optimal.

We attempt to select the processor dimensions  $p_i$  and  $q_i$  in such a way that the terms in the communication cost match the optimal solutions of Corollaries 4.1 and 4.2. In other words, we want to select  $p_i$  and  $q_i$  such that  $\frac{n_1 r_1}{p_1 q_1} = x^*$ ,  $\frac{n_2 r_2}{p_2 q_2} = y^*$ , and  $\frac{n_3 r_3}{p_3 q_3} = z^*$  from Corollary 4.1, and  $\frac{n}{p} = v^*$ ,  $\frac{r}{q} = u^*$  from Corollary 4.2.

While the optimal values are given as a single term, we have two or three processor grid dimensions we need to fix in order to match each term, and each processor grid dimension appears in two equations. In general, we are able to set the processor grid dimensions in a way that is consistent with these equations. However, they are subject to additional constraints that are not imposed by the optimization problem. Specifically, we have  $1 \leq p_i \leq n_i$  and  $1 \leq q_i \leq r_i$  for  $1 \leq i \leq 3$ . The lower bounds are imposed because processor grid dimensions must be at least 1. The upper bounds are imposed to ensure that each processor performs its fair share of the computations. We assume that  $P \leq nr$ , so that every processor has at least one 4-ary multiplication term to compute. For simplicity, we assume that the final grid dimensions are integers and perfectly divide the corresponding input and output dimensions. However, we also discuss how to handle non-integer grid dimensions for a specific set of inputs in Appendix B.1.

In order to define processor grid dimensions, we begin by determining a set of values that match the lower bound terms and denote these by  $\hat{p}_i, \hat{q}_i$  with their products denoted  $\hat{p}$  and  $\hat{q}$ . Then, we will consider how to adapt  $\hat{p}_i$  and  $\hat{q}_i$  so that the additional constraints are met. During the adaption, we maintain the tensor communication costs, modify the matrix communication costs, and then bound the additional costs in terms of communication lower bounds of tensors.

As  $\mathcal{X}$  and  $\mathcal{Y}$  are 3-dimensional tensors, we have  $n_i, r_i \geq 2$  for all  $1 \leq i \leq 3$ . For better readability, we use the notation  $O = \frac{\sum_{j \in [3]} n_j r_j + r + n}{P}$ , the amount of data owned by a single processor at the beginning and end of the algorithm.

**THEOREM 5.1.** *There exist  $p_i, q_i$  with  $1 \leq p_i \leq n_i, 1 \leq q_i \leq r_i$  for  $i = 1, 2, 3$  such that Alg. 5.1 is communication optimal to within a constant factor.*

*Proof.* We break our analysis into 2 scenarios which are further broken down into cases. In each case, we obtain  $\hat{p}_i$  and  $\hat{q}_i$  such that the terms in the communication cost match the corresponding lower bound terms and also satisfy at least one of the two constraints:  $1 \leq \hat{p}_i \leq n_i$  or  $1 \leq \hat{q}_i \leq r_i$  for  $1 \leq i \leq 3$ . We handle all cases from both scenarios together in the end, and adapt these values to get  $p_i$  and  $q_i$  which respect both lower and upper bounds.

• **Scenario I** ( $P < \frac{n}{r}$ ): This scenario corresponds to the first case of the tensor term in LB. Thus, we set  $\hat{p}_i, \hat{q}_i$  in such a way that the tensor terms in the communication

cost match the tensor terms of  $LB$ :

$$(5.2) \quad \hat{p} = P, \hat{q} = 1.$$

This implies  $\hat{q}_i = 1$  for  $1 \leq i \leq 3$ . We break this scenario into 3 cases, each corresponding to a case in the matrix term of  $LB$ .

(Case 1)  $P < \min \left\{ \frac{n_3 r_3}{n_2 r_2}, \frac{n}{r} \right\}$ : Setting the matrix communication costs to the matrix terms in the corresponding case of the lower bound yields

$$(5.3) \quad \frac{n_1 r_1}{\hat{p}_1 \hat{q}_1} = n_1 r_1, \quad \frac{n_2 r_2}{\hat{p}_2 \hat{q}_2} = n_2 r_2, \quad \frac{n_3 r_3}{\hat{p}_3 \hat{q}_3} = \frac{n_3 r_3}{P}.$$

Thus, we set  $\hat{p}_1 = \hat{p}_2 = \hat{q}_1 = \hat{q}_2 = \hat{q}_3 = 1$  and  $\hat{p}_3 = P$  to satisfy (5.2) and (5.3).

(Case 2)  $\frac{n_3 r_3}{n_2 r_2} \leq P < \min \left\{ \frac{n_2 n_3 r_2 r_3}{n_1^2 r_1^2}, \frac{n}{r} \right\}$ : Setting the matrix communication costs to the matrix terms in the corresponding case of the lower bound yields

$$(5.4) \quad \frac{n_1 r_1}{\hat{p}_1 \hat{q}_1} = n_1 r_1, \quad \frac{n_2 r_2}{\hat{p}_2 \hat{q}_2} = \frac{n_3 r_3}{\hat{p}_3 \hat{q}_3} = \left( \frac{n_2 n_3 r_2 r_3}{P} \right)^{1/2}.$$

We set  $\hat{p}_1 = \hat{q}_1 = \hat{q}_2 = \hat{q}_3 = 1$ ,  $\hat{p}_2 = n_2 r_2 \left( \frac{P}{n_2 n_3 r_2 r_3} \right)^{\frac{1}{2}}$ , and  $\hat{p}_3 = n_3 r_3 \left( \frac{P}{n_2 n_3 r_2 r_3} \right)^{\frac{1}{2}}$  to satisfy (5.2) and (5.4).

(Case 3)  $\frac{n_2 n_3 r_2 r_3}{n_1^2 r_1^2} \leq P < \frac{n}{r}$ : Setting the matrix communication costs to match the matrix terms in the corresponding case of the lower bound yields

$$(5.5) \quad \frac{n_1 r_1}{\hat{p}_1 \hat{q}_1} = \frac{n_2 r_2}{\hat{p}_2 \hat{q}_2} = \frac{n_3 r_3}{\hat{p}_3 \hat{q}_3} = \left( \frac{nr}{P} \right)^{1/3}.$$

Thus we set  $\hat{q}_1 = \hat{q}_2 = \hat{q}_3 = 1$ ,  $\hat{p}_1 = n_1 r_1 \left( \frac{P}{nr} \right)^{\frac{1}{3}}$ ,  $\hat{p}_2 = n_2 r_2 \left( \frac{P}{nr} \right)^{\frac{1}{3}}$ , and  $\hat{p}_3 = n_3 r_3 \left( \frac{P}{nr} \right)^{\frac{1}{3}}$  to satisfy (5.2) and (5.5).

Note that in all the cases of this scenario we have  $1 \leq \hat{q}_i = 1 < r_i$  for  $1 \leq i \leq 3$ , but we cannot ensure similar upper bounds for  $\hat{p}_i$ . We will adapt processor grid dimensions for both scenarios in the end as they require the same s-eps-converted-to.pdf.

• **Scenario II** ( $\frac{n}{r} \leq P$ ): This scenario corresponds to the second case of the tensor term in  $LB$ . Thus, we set  $\hat{p}_i, \hat{q}_i$  in such a way that

$$(5.6) \quad \frac{n}{\hat{p}} = \frac{r}{\hat{q}} = \left( \frac{nr}{P} \right)^{1/2}.$$

Again, we break this scenario into 3 cases each corresponding to a case in the matrix term of  $LB$ .

(Case 1)  $\frac{n}{r} \leq P < \frac{n_3 r_3}{n_2 r_2}$ : Setting the matrix communication costs to the matrix terms in the corresponding case of the lower bound yields

$$(5.7) \quad \frac{n_1 r_1}{\hat{p}_1 \hat{q}_1} = n_1 r_1, \quad \frac{n_2 r_2}{\hat{p}_2 \hat{q}_2} = n_2 r_2, \quad \frac{n_3 r_3}{\hat{p}_3 \hat{q}_3} = \frac{n_3 r_3}{P}.$$

Thus we set  $\hat{p}_1 = \hat{q}_1 = \hat{p}_2 = \hat{q}_2 = 1$ ,  $\hat{p}_3 = n \left( \frac{P}{nr} \right)^{1/2}$ , and  $\hat{q}_3 = r \left( \frac{P}{nr} \right)^{1/2}$  to satisfy (5.6) and (5.7). As  $\frac{n}{r} \leq P \leq nr$  and  $r \leq n$ , we have  $1 \leq \hat{p}_3 \leq n$  and  $1 \leq \hat{q}_3 \leq r$ , but cannot ensure  $\hat{p}_3 \leq n_3$  or  $\hat{q}_3 \leq r_3$ . However,  $\hat{p}_3 \hat{q}_3 = P < \frac{n_3 r_3}{n_2 r_2}$  implies that at least one is

satisfied. Therefore, we guarantee that  $1 \leq i \leq 3$ ,  $1 \leq \hat{p}_i$ ,  $1 \leq \hat{q}_i$  and either  $\hat{p}_i \leq n_i$  or  $\hat{q}_i \leq r_i$ .

(Case 2)  $\max \left\{ \frac{n_3 r_3}{n_2 r_2}, \frac{n}{r} \right\} \leq P < \frac{n_2 n_3 r_2 r_3}{n_1^2 r_1^2}$ : Setting the matrix communication costs to the matrix terms in the corresponding case of the lower bound yields

$$(5.8) \quad \frac{n_1 r_1}{\hat{p}_1 \hat{q}_1} = n_1 r_1, \frac{n_2 r_2}{\hat{p}_2 \hat{q}_2} = \frac{n_3 r_3}{\hat{p}_3 \hat{q}_3} = \left( \frac{n_2 n_3 r_2 r_3}{P} \right)^{1/2}.$$

The equations (5.6) and (5.8) do not uniquely determine  $\hat{p}_1, \hat{p}_2, \hat{p}_3, \hat{q}_1, \hat{q}_2$ , and  $\hat{q}_3$ .

The following is one possible solution:  $\hat{p}_1 = \hat{q}_1 = 1$ ,  $\hat{p}_2 = n_2 \left( \frac{n_1 P}{n_2 n_3 r} \right)^{1/4}$ ,  $\hat{p}_3 = n_3 \left( \frac{n_1 P}{n_2 n_3 r} \right)^{1/4}$ ,  $\hat{q}_2 = r_2 \left( \frac{r_1 P}{n r_2 r_3} \right)^{1/4}$ , and  $\hat{q}_3 = r_3 \left( \frac{r_1 P}{n r_2 r_3} \right)^{1/4}$ . Note that  $P < \frac{n_2 n_3 r_2 r_3}{n_1^2 r_1^2}$  implies that  $\hat{p}_2 < n_2$ ,  $\hat{p}_3 < n_3$ ,  $\hat{q}_2 < r_2$ , and  $\hat{q}_3 < r_3$ . We are not able to ensure  $\hat{p}_2, \hat{p}_3, \hat{q}_2, \hat{q}_3$  are all at least 1 in this case. We will handle both Case 2 and Case 3 together as they require the same analysis.

(Case 3)  $\max \left\{ \frac{n_2 n_3 r_2 r_3}{n_1^2 r_1^2}, \frac{n}{r} \right\} \leq P$ : Setting the matrix communication costs to the matrix terms in the corresponding case of the lower bound yields

$$(5.9) \quad \frac{n_1 r_1}{\hat{p}_1 \hat{q}_1} = \frac{n_2 r_2}{\hat{p}_2 \hat{q}_2} = \frac{n_3 r_3}{\hat{p}_3 \hat{q}_3} = \left( \frac{nr}{P} \right)^{\frac{1}{3}}.$$

As in Case 2, the equations (5.6) and (5.9) do not uniquely determine  $\hat{p}_i, \hat{q}_i$  for  $1 \leq i \leq 3$ . We choose a cubical distribution, namely  $\frac{n_1}{\hat{p}_1} = \frac{n_2}{\hat{p}_2} = \frac{n_3}{\hat{p}_3} = \frac{r_1}{\hat{q}_1} = \frac{r_2}{\hat{q}_2} = \frac{r_3}{\hat{q}_3}$  and obtain the following solution,  $\hat{p}_i = n_i \left( \frac{P}{nr} \right)^{1/6}$ ,  $\hat{q}_i = r_i \left( \frac{P}{nr} \right)^{1/6}$  for  $1 \leq i \leq 3$ . As  $P \leq nr$  we have  $\hat{p}_i \leq n_i$  and  $\hat{q}_i \leq r_i$  for  $1 \leq i \leq 3$ . Again we are not able to ensure  $\hat{p}_i$  and  $\hat{q}_i$  are all greater than 1 for  $1 \leq i \leq 3$ .

Now we handle Case 2 and Case 3 of Scenario II here. Communication cost for the obtained set of values matches the lower bound, therefore we have  $1 \leq \hat{p}_i \hat{q}_i \leq n_i r_i$  for  $1 \leq i \leq 3$ ,  $1 \leq \hat{p} \leq n$  and  $1 \leq \hat{q} \leq r$ . We perform the following: for a given  $i$ , at most one of  $\hat{p}_i$  or  $\hat{q}_i$  can be smaller than 1. If either is smaller than one, set it to 1 and update the other so that  $\hat{p}_i \hat{q}_i$  does not change. This, however, might change  $\hat{p}$  and  $\hat{q}$ . The increase factor in one of  $\hat{p}$  and  $\hat{q}$  (say  $\hat{q}$  without loss of generality),  $f = \hat{q}_{\text{new}}/\hat{q}_{\text{orig}}$ , is reflected in a decrease in the other ( $\hat{p}$ ).

As the original  $\hat{q} \geq 1$ , we can factor  $f$ ,  $f = f_1 \cdot f_2 \cdot f_3$  ( $f_i \geq 1$ ) such that  $\hat{q}_i = \hat{q}_i/f_i \geq 1$  and  $\hat{p}_i = \hat{p}_i f_i \geq 1$ , and hence  $\hat{p}$  and  $\hat{q}$  are back to their original values. Note that in the above updates, we always maintain  $\hat{q}_i \leq r_i$ . Therefore, it is guaranteed that  $1 \leq \hat{q}_i \leq r_i$ ,  $1 \leq \hat{p}_i$  for  $1 \leq i \leq 3$ .

If the increase factor was in  $\hat{p}$ , we would have obtained  $1 \leq \hat{p}_i \leq n_i$ ,  $1 \leq \hat{q}_i$  for  $1 \leq i \leq 3$ . Therefore, we have  $1 \leq \hat{p}_i \leq n_i$ ,  $1 \leq \hat{q}_i$  and/or  $1 \leq \hat{q}_i \leq r_i$ ,  $1 \leq \hat{p}_i$  for  $1 \leq i \leq 3$ .

Now for all the cases of both scenarios, it remains to adapt  $\hat{p}_i$  and  $\hat{q}_i$  such that  $\hat{p}_i \leq n_i$  and  $\hat{q}_i \leq r_i$ . We know that one of them is guaranteed from the way they are selected in all the cases. We now obtain  $p_1, p_2, p_3, q_1, q_2, q_3$  from  $\hat{p}_i, \hat{q}_i$  such that both lower and upper bounds are respected, and  $p_1 p_2 p_3 = \hat{p}$  and  $q_1 q_2 q_3 = \hat{q}$ . The intuition is to maintain the tensor communication terms in the lower bound.

Initially we set  $p_i = \hat{p}_i, q_i = \hat{q}_i$  for  $1 \leq i \leq 3$ . If  $\hat{p}_i > n_i$  for some  $i$ , we represent this index with  $l$ , and the other two indices with  $j$  and  $k$ . As  $\hat{p} \leq n$ , therefore  $\hat{p}_j \leq n_j$  or/and  $\hat{p}_k \leq n_k$ . Without loss of generality, we assume that  $\hat{p}_k \leq n_k$ . Now we first update  $p_l$ , and then  $p_j$ , and in the end  $p_k$  with the following expressions:

$p_l = \min \left\{ n_l, \frac{\hat{p}}{p_j p_k} \right\}$ ,  $p_j = \min \left\{ n_j, \frac{\hat{p}}{p_k p_l} \right\}$ ,  $p_k = \min \left\{ n_k, \frac{\hat{p}}{p_l p_j} \right\}$ . The same update can be done if  $\hat{q}_i > r_i$  for some  $i$ .

Now we assess how much additional communication is required for the matrices. If  $\nexists i \in [3]$  such that  $\hat{p}_i > n_i$  or  $\hat{q}_i > r_i$  then  $\sum_{i \in [3]} \frac{n_i r_i}{p_i q_i} = \sum_{i \in [3]} \frac{n_i r_i}{\hat{p}_i \hat{q}_i}$ . We can note that due to our particular selections of  $\hat{p}_i$  and  $\hat{q}_i$ ,  $\nexists i, j \in [3]$  such that  $\hat{p}_i > n_i$  and  $\hat{q}_j > r_j$ . Suppose  $\exists i \in [3]$  such that  $\hat{p}_i > n_i$  then  $\hat{p} > 2$  and

$$\begin{aligned}
\sum_{i \in [3]} \frac{n_i r_i}{p_i q_i} &\leq \sum_{i \in [3]} \max \left\{ \frac{n_i r_i}{\hat{p}_i \hat{q}_i}, \frac{r_i}{\hat{q}_i} \right\} && q_i = \hat{q}_i, \text{ and } p_i \geq \hat{p}_i \text{ or } p_i = n_i \\
&= \sum_{i \in [3]} \left( \frac{n_i r_i}{\hat{p}_i \hat{q}_i} + \frac{r_i}{\hat{q}_i} - \min \left\{ \frac{n_i r_i}{\hat{p}_i \hat{q}_i}, \frac{r_i}{\hat{q}_i} \right\} \right) && \max\{a, b\} = a + b - \min\{a, b\} \\
&< \sum_{i \in [3]} \left( \frac{n_i r_i}{\hat{p}_i \hat{q}_i} + \frac{r_i}{\hat{q}_i} \right) - 2 && \hat{p}_i \hat{q}_i \leq n_i r_i \text{ and } \hat{q}_i \leq r_i \\
&\leq \sum_{i \in [3]} \frac{n_i r_i}{\hat{p}_i \hat{q}_i} + \frac{r}{\hat{q}} && \forall a_i \geq 1, a_1 + a_2 + a_3 - 2 \leq a_1 a_2 a_3 \\
&< \sum_{i \in [3]} \frac{n_i r_i}{\hat{p}_i \hat{q}_i} + 2 \left( \frac{r}{\hat{q}} - \frac{r}{\hat{p} \hat{q}} \right) \\
&= \sum_{i \in [3]} \frac{n_i r_i}{\hat{p}_i \hat{q}_i} + 2 \left( \frac{r}{\hat{q}} - \frac{r}{P} \right).
\end{aligned}$$

Similarly, if  $\exists i \in [3]$  such that  $\hat{q}_i > r_i$  we can obtain  $\sum_{i \in [3]} \frac{n_i r_i}{p_i q_i} < \sum_{i \in [3]} \frac{n_i r_i}{\hat{p}_i \hat{q}_i} + 2 \left( \frac{n}{\hat{p}} - \frac{n}{P} \right)$ .

Therefore, in all situations,  $\sum_{i \in [3]} \frac{n_i r_i}{p_i q_i} + \frac{r}{q} + \frac{n}{p} - O \leq 3 \left( \sum_{i \in [3]} \frac{n_i r_i}{\hat{p}_i \hat{q}_i} + \frac{r}{\hat{q}} + \frac{n}{\hat{p}} - O \right) = 3LB$ .  $\square$

**6. Simulated Evaluation.** In this section, we verify our theoretical claims on particular sets of 3D tensor dimensions using a simulated evaluation. In Subsection 6.1, we demonstrate that the communication cost of Alg. 5.1 matches the lower bound of Theorem 4.3, and we provide intuition for relationships among the communication costs of the individual tensors and matrices. In Subsection 6.2, we compare the approach of Alg. 5.1 for evaluating Multi-TTM with a TTM-in-Sequence approach, demonstrating realistic scenarios when Alg. 5.1 communicates significantly less data and performs a negligible amount of extra computation.

Throughout this section, we restrict to cases where all tensor dimensions and numbers of processors are powers of 2. We vary the number of processors  $P$  from 2 to  $P_{\max} = \min\{n_1 r_1, n_2 r_2, n_3 r_3, n, r\}$ , which ensures that each processor owns some data of every tensor and matrix. The costs of Alg. 5.1 depend on the processor grid, and in these experiments, we perform an exhaustive search for the best configuration. We describe in Appendix B.1 how to adapt the processor grid selection scheme described in Subsection 5.2 to obtain integer-valued processor grid dimensions, and we show that we can obtain nearly optimal configurations without exhaustive search.

**6.1. Verifying Optimality of Algorithm 5.1.** Theorem 5.1 states that Alg. 5.1 attains the communication lower bound to within a constant factor, and in this section we verify the result in a variety of scenarios. Recall from Theorem 4.3 that the

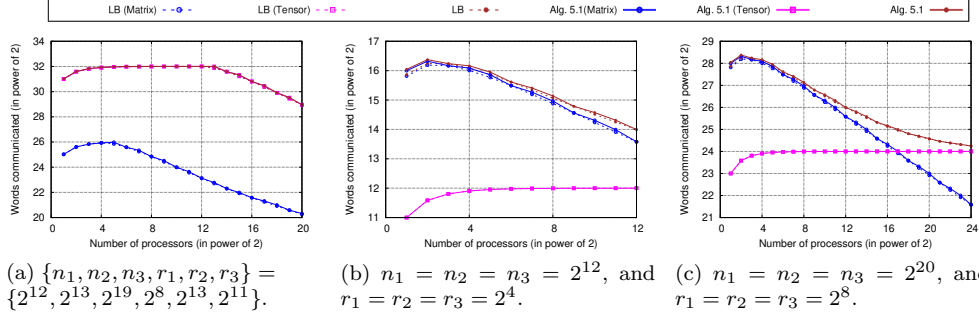


Fig. 3: Matrix and tensor communication costs in LB and Alg. 5.1 for different configurations. The sum of  $LB(Matrix)$  and  $LB(Tensor)$  equals to the lower bound ( $LB$ ), and the sum of Alg. 5.1 (Matrix) and Alg. 5.1 (Tensor) equals to the upper bound (Alg. 5.1). Lower bounds are almost indistinguishable from the corresponding upper bounds.

lower bound is  $A + B - O$ , where

$$A = \begin{cases} n_1 r_1 + n_2 r_2 + \frac{n_3 r_3}{P} & \text{if } P < \frac{n_3 r_3}{n_2 r_2} \\ n_1 r_1 + 2 \left( \frac{n_2 n_3 r_2 r_3}{P} \right)^{\frac{1}{2}} & \text{if } \frac{n_3 r_3}{n_2 r_2} \leq P < \frac{n_2 n_3 r_2 r_3}{n_1 r_1} \\ 3 \left( \frac{n r}{P} \right)^{\frac{1}{3}} & \text{if } \frac{n_2 n_3 r_2 r_3}{n_1 r_1} \leq P \end{cases}$$

$$B = \begin{cases} r + \frac{n}{P} & \text{if } P < \frac{n}{r} \\ 2 \left( \frac{n r}{P} \right)^{\frac{1}{2}} & \text{if } \frac{n}{r} \leq P. \end{cases}$$

$$O = \frac{n_1 r_1 + n_2 r_2 + n_3 r_3 + r + n}{P}.$$

Here,  $A$  corresponds to the matrix entries accessed,  $B$  corresponds to the tensor entries accessed, and  $O$  corresponds to the data owned by a single processor. The costs of Alg. 5.1 are given by (5.1), which we re-write here as

$$\frac{n_1 r_1}{p_1 q_1} + \frac{n_2 r_2}{p_2 q_2} + \frac{n_3 r_3}{p_3 q_3} + \frac{n}{p} + \frac{r}{q} - O,$$

where  $\{p_i\}$  and  $\{q_i\}$  specify the processor grid dimensions. The first three terms correspond to matrix entries and the middle two terms correspond to tensor entries.

Figure 3 shows both components, matrix and tensor communication costs, for three distinct input sizes as we vary the number of processors. In these plots, we show both algorithmic costs (upper bounds) and lower bounds, but they are indistinguishable because the largest differences in overall costs we observe are 9% for Fig. 3a at  $P = 2^{13}$  and 13% for Figures 3b and 3c at  $P = 2$ , verifying Theorem 4.3 for these scenarios.

In Fig. 3a, the input and output tensors have varying dimensions: the input tensor is  $2^{12} \times 2^{13} \times 2^{19}$  and the output is  $2^8 \times 2^{13} \times 2^{11}$ . We choose these dimensions so that all five cases of the values of  $A$  and  $B$  are represented. For these inputs, the tensor communication cost dominates the matrix communication for all values of  $P$  considered. When  $P < 2^4$ , the first cases for  $A$  and  $B$  apply, and the algorithm selects a processor grid such that  $p_3 = P$ , implying that only one tensor and two matrices are



communicated. In this case, both expressions simplify to  $(r + n_1 r_1 + n_2 r_2)(1 - 1/P)$ , which is why we see initial increase as  $P$  increases at the left end of the plot. For  $2^4 \leq P < 2^{12}$ , the second case for  $A$  and the first case for  $B$  apply, and the algorithm selects a processor grid with  $p_2 > 1$  and  $p_3 > 1$ . Here, the matrix communication begins to decrease, but it is dominated by the tensor communication, which is maintained at  $r(1 - 1/P)$ . For  $2^{12} \leq P$ , the second case for  $B$  applies, and we see that tensor communication decreases as  $P$  increases (proportional to  $P^{-1/2}$  as we see from the lower bound). In this regime, the algorithm is selecting grids with both  $p > 1$  and  $q > 1$  and communicating both tensors. Another transition occurs at  $P = 2^{16}$ , switching from the second to third case of  $A$ , but this change in matrix cost has a negligible effect.

Figure 3b demonstrates a scenario where the matrix costs dominate the tensor costs: the input tensor is cubical with dimension  $2^{12}$  and the output tensor is cubical with dimension  $2^4$ . Here we scale  $P$  only up to  $2^{12}$ , the number of entries in the output tensor. Because the tensors are cubical, the lower bounds simplify as in Corollary 4.4, and the algorithm chooses processor grids that are as cubical as possible. For all values of  $P$  in this experiment, the third case of  $A$  and the first case of  $B$  apply, and the algorithm selects  $p_1 \approx p_2 \approx p_3$  and  $q = 1$ . We see that the overall cost is decreasing proportional to  $P^{-1/3}$  until the tensor communication cost starts to contribute more significantly.

Figure 3c considers cubical tensors with larger dimensions to show a more general pattern. For tensor dimensions  $n_i = 2^{20}$  and  $r_i = 2^8$ , we observe a transition point where tensor communication overtakes matrix communication. Similar to the case of Figure 3b, matrix costs dominate for small  $P$  and scale like  $P^{-1/3}$ . However, for  $P \geq 2^{17}$ , the tensor costs dominate the matrix costs and communication costs scale less efficiently as the first case of  $B$  applies. We emphasize that for all three of these experiments, the algorithmic costs match the lower bounds nearly exactly for all values of  $P$ .

**6.2. Comparing Algorithm 5.1 with TTM-in-Sequence.** As mentioned previously, a Multi-TTM computation may be performed as sequence of TTM operations. In this TTM-in-Sequence approach, a single matrix is multiplied with the tensor and an intermediate tensor is computed and stored. For each remaining matrix, single-matrix TTMs are performed in sequence until the final result is computed. This approach can reduce the number of arithmetic operations compared to direct evaluation of atomic expression given in Def. 3.1. The computational cost depends (often significantly) on the order of the TTMs performed. The TTM-in-Sequence approach is parallelized in the TuckerMPI library [4]. We note that Theorem 4.3 does not apply to this parallelization, as it violates the parallel atomicity assumption.

In this section, we provide a comparison between Alg. 5.1 and the TTM-in-Sequence approach to show that our approach can significantly reduce communication in important scenarios without performing too much extra computation. In particular, we observe greatest benefit of Alg. 5.1 when  $r$  is very small relative to  $n$  (or vice versa) and  $P$  is small relative to the ratio of  $n$  and  $r$ . These scenarios occur in the context of computing and using Tucker decompositions for highly compressible tensors that exhibit small multilinear ranks.

The computational cost of TuckerMPI's algorithm with cubical tensors is the same for all possible orderings of the TTMs. In our comparison, we consider that the TTMs are performed in increasing mode order. While no single communication lower bound exists for all parallel TTM-in-Sequence algorithms, we show in Appendix B.2

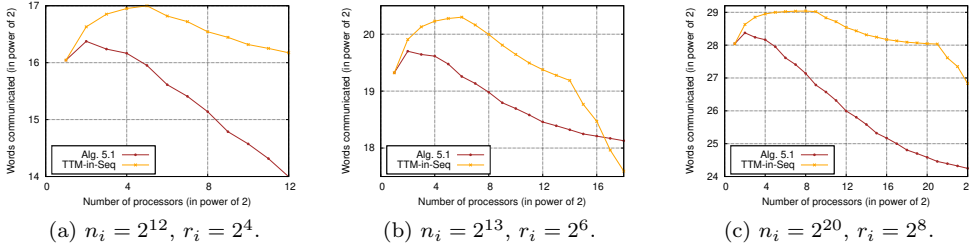
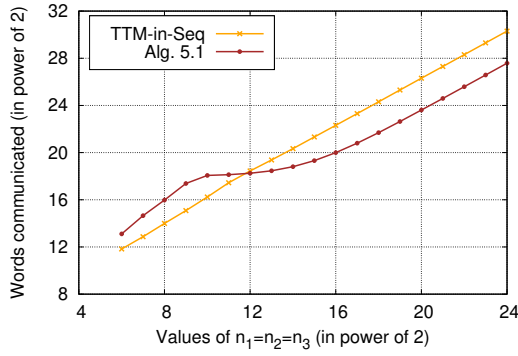


Fig. 4: Communication cost comparison of Alg. 5.1 and TTM-in-Sequence [4].

Fig. 5: Comparison of Alg. 5.1 and the TTM-in-Sequence approach for fixed  $r_1 = r_2 = r_3 = 2^6$  and  $P = 2^{12}$ .

that TuckerMPI's algorithm attains nearly the same cost as tight matrix multiplication lower bounds [1] applied to each TTM it chooses to perform. Thus, no other parallelization of the TTM-in-Sequence approach can reduce communication without breaking the assumptions of the matrix multiplication lower bounds (e.g., using fast matrix multiplication).

The TuckerMPI parallelization uses a 3D logical processor grid with dimensions  $\tilde{p}_1 \times \tilde{p}_2 \times \tilde{p}_3$ . When the TTMs are performed in increasing mode order, the overall communication cost of their algorithm is

$$(6.1) \quad \frac{r_1 n_2 n_3}{\tilde{p}_2 \tilde{p}_3} + \frac{n_1 r_1}{\tilde{p}_1} + \frac{r_1 r_2 n_3}{\tilde{p}_1 \tilde{p}_3} + \frac{n_2 r_2}{\tilde{p}_2} + \frac{r_1 r_2 r_3}{\tilde{p}_1 \tilde{p}_2} + \frac{n_3 r_3}{\tilde{p}_3} - \frac{r_1 n_2 n_3 + r_1 r_2 n_3 + r_1 r_2 r_3 + n_1 r_1 + n_2 r_2 + n_3 r_3}{P},$$

as specified in [4, Section 6.3], though we include the cost of communicating the matrices (their analysis assumes the matrices are already redundantly distributed). We use exhaustive search to determine the processor grid that minimizes the cost of (6.1) in our comparisons.

**6.2.1. Communication Cost.** To compare communications costs, we perform 4 experiments involving cubical tensors. The first three simulated evaluations consider strong scaling and are presented in Fig. 4. Two of these experiments use the same

tensor dimensions as the two cubical examples in Fig. 3. The first experiment involves an input tensor of dimension  $n_i = 2^{12}$  and output dimension  $r_i = 2^4$  (Fig. 4a), the second has dimensions  $n_i = 2^{13}$  and  $r_i = 2^6$  (Fig. 4b), and the third has the largest dimensions  $n_i = 2^{20}$  and  $r_i = 2^8$  (Fig. 4c).

Figure 4a shows that Alg. 5.1 performs less communication than TTM-in-Sequence for  $P \leq 2^{12} < n/r$ . The largest communication reduction occurs at  $P = 2^{12}$  and is approximately  $5\times$ . In the second experiment, we see cases where TTM-in-Sequence performs less communication than Alg. 5.1 and in fact beats the lower bound of Theorem 4.3 (which is possible because it breaks the atomicity assumption). Algorithm 5.1 is more communication efficient for  $P \leq 2^{16}$ , achieving a speedup of up to  $2\times$ , but communicates more for larger  $P$ . In the third experiment with larger tensors, Fig. 4c demonstrates similar qualitative behavior to the first, with Alg. 5.1 outperforming TTM-in-Sequence and a maximum communication reduction of approximately  $12\times$  at  $P = 2^{21}$ .

In the fourth experiment, with results shown in Fig. 5, we fix the output tensor dimension  $r_i = 2^6$  and number of processors  $P = 2^{12}$  and vary the input tensor dimension  $n_i$ . We observe that for  $2^6 \leq n_i < 2^{12}$ , the TTM-in-Sequence approach communicates less data than Alg. 5.1. For  $n_i \geq 2^{12}$ , Alg. 5.1 communicates less data, and the factor of improvement is maintained at approximately  $6\times$  as  $n_i$  scales up.

**6.2.2. Computation Cost.** Assuming TuckerMPI uses increasing mode order, the parallel computational cost is

$$2 \cdot \frac{r_1 n_1 n_2 n_3 + r_1 r_2 n_2 n_3 + r_1 r_2 r_3 n_3}{P} = 2 \left( \frac{r^{1/3} n}{P} + \frac{r^{2/3} n^{2/3}}{P} + \frac{r n^{1/3}}{P} \right),$$

where the right hand side is simplified under the assumption of cubical tensors. In these experiments where  $n \gg r$ , Alg. 5.1 selects a processor grid such that  $q = 1$  and  $p_1 \approx p_2 \approx p_3$ . In this case the computation cost given in Subsection 5.1 simplifies to

$$2 \left( \frac{r^{1/3} n}{P} + \frac{r^{2/3} n^{2/3}}{P^{2/3}} + \frac{r n^{1/3}}{P^{1/3}} \right).$$

Note that this cost is much smaller than  $4nr/P$ , the cost of evaluating (3.1) directly with computational load balance, and it is achieved by performing local computation using a TTM-in-Sequence approach.

While the first terms of the two computational cost expressions match, we observe greater computational cost from Alg. 5.1 in the second and third terms. These terms are lower order when  $P \ll n/r$ , in which case the extra computational cost of Alg. 5.1 is negligible. When  $P = n/r$ , the extra computational cost is no more than  $3\times$ .

In the first three experiments, when our approach reduces communication, the extra computational costs were at most 6%, 30%, and 7%, respectively. The extra computation required for the greatest reductions in communication in those experiments were 6%, 2%, and 7%. For the fourth experiment, the extra computation is approximately 13% at  $n_i = 2^{13}$ , where Alg. 5.1 provides communication reduction, and decreases as  $n_i$  increases.

In all these experiments, we see that when Alg. 5.1 provides a reduction in communication costs, the extra computational costs remain negligible.

**7. Lower Bounds of General Multi-TTM.** We present our lower bound results for  $d$ -dimensional tensors in this section. Similar to the 3-dimensional lower bound proof, we consider a single processor that performs  $1/P$ th of the computation

and has access to  $1/P$ th of the data. We again seek to minimize the number of elements of the matrices and tensors that the processor must access in order to execute its computation subject to the constraints of the structure of Multi-TTM by solving two independent problems, one for the matrix data and one for the tensor data.

**7.1. General Constrained Optimization Problems.** Here we present a generalization of Corollary 4.1 for  $d$  dimensions. As before, this corollary is a direct result of Theorem 3.4. Recall the notation  $N_i = \prod_{j=d-i+1}^d n_j$  and  $R_i = \prod_{j=d-i+1}^d r_j$ .

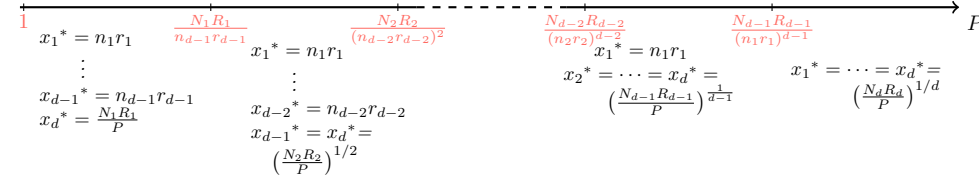
COROLLARY 7.1. *Consider the following optimization problem:*

$$\min_{\mathbf{x}} \sum_{i \in [d]} x_i$$

such that

$$\frac{nr}{P} \leq \prod_{i \in [d]} x_i \quad \text{and} \quad 0 \leq x_i \leq n_i r_i \quad \text{for all} \quad 1 \leq i \leq d,$$

where  $n_i, r_i, P \geq 1$  and  $n_i r_i \leq n_{i+1} r_{i+1}$ . The optimal solution  $\mathbf{x} = [x_1^* \cdots x_d^*]$  depends on the values of constants, yielding  $d$  cases.



- If  $P < \frac{N_1 R_1}{n_{d-1} r_{d-1}}$ , then

$$x_j^* = n_j r_j \quad \text{for} \quad 1 \leq j \leq d-1 \quad \text{and} \quad x_d^* = \frac{N_1 R_1}{P}.$$

- If  $\frac{N_{i-1} R_{i-1}}{(n_{d+1-i} r_{d+1-i})^{i-1}} \leq P < \frac{N_i R_i}{(n_{d-i} r_{d-i})^i}$  for some  $i = 2, \dots, d-1$ , then

$$x_j^* = n_j r_j \quad \text{for} \quad 1 \leq j \leq d-i \quad \text{and} \quad x_{d+1-i}^* = \cdots = x_d^* = (N_i R_i / P)^{1/i}.$$

- If  $\frac{N_{d-1} R_{d-1}}{(n_1 r_1)^{d-1}} \leq P$ , then

$$x_1^* = \cdots = x_d^* = (N_d R_d / P)^{1/d}.$$

**7.2. Communication Lower Bounds.** We now state the lower bounds for the general Multi-TTM computation in the following theorem.

THEOREM 7.2. *Any computationally load balanced atomic Multi-TTM algorithm that starts and ends with one copy of the data distributed across processors and involves  $d$ -dimensional tensors with dimensions  $n_1, n_2, \dots, n_d$  and  $r_1, r_2, \dots, r_d$  performs at least  $A + B - \left( \frac{n}{P} + \frac{r}{P} + \sum_{j=1}^d \frac{n_j r_j}{P} \right)$  sends or receives where*

$$A = \begin{cases} \sum_{j=1}^{d-1} n_j r_j + \frac{N_1 R_1}{P} & \text{if } P < \frac{N_1 R_1}{n_{d-1} r_{d-1}}, \\ \sum_{j=1}^{(d-i)} n_j r_j + i \left( \frac{N_i R_i}{P} \right)^{\frac{1}{i}} & \text{if } \frac{N_{i-1} R_{i-1}}{(n_{d+1-i} r_{d+1-i})^{i-1}} \leq P < \frac{N_i R_i}{(n_{d-i} r_{d-i})^i}, \\ & \text{for some } 2 \leq i \leq d-1, \\ d \left( \frac{N_d R_d}{P} \right)^{\frac{1}{d}} & \text{if } \frac{N_{d-1} R_{d-1}}{(n_1 r_1)^{d-1}} \leq P. \end{cases}$$

$$B = \begin{cases} r + \frac{n}{P} & \text{if } P < \frac{n}{r}, \\ 2 \left( \frac{nr}{P} \right)^{\frac{1}{2}} & \text{if } \frac{n}{r} \leq P. \end{cases}$$

We prove this by applying Corollaries 4.2 and 7.1 and extending the arguments of Theorem 4.3 in a straightforward way (though with more complicated notation). Interested readers can see the detailed proof in Appendix C.1.

**8. Parallel Algorithm for General Multi-TTM.** We present a parallel algorithm for  $d$ -dimensional Multi-TTM computation in this section, which is analogous to Alg. 5.1. We organize  $P$  processors into a  $2d$ -dimensional logical processor grid with dimensions  $p_1 \times \cdots \times p_d \times q_1 \times \cdots \times q_d$ . As before, we consider that  $\forall i \in [d]$ ,  $p_i$  and  $q_i$  evenly divide  $n_i$  and  $r_i$ , respectively. A processor coordinate is represented as  $(p'_1, \dots, p'_d, q'_1, \dots, q'_d)$ , where  $\forall i \in [d]$ ,  $1 \leq p'_i \leq p_i$  and  $1 \leq q'_i \leq q_i$ .

We again impose that there is one copy of data in the system at the beginning and the end of the computation. Algorithm 8.1 presents our proposed parallel algorithm for  $d$ -dimensional Multi-TTM computation.

---

**Algorithm 8.1** Parallel Atomic  $d$ -dimensional Multi-TTM

---

**Require:**  $\mathcal{X}, \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(d)}, p_1 \times \cdots \times p_d \times q_1 \times \cdots \times q_d$  logical processor grid

**Ensure:**  $\mathcal{Y}$  such that  $\mathcal{Y} = \mathcal{X} \times_1 \mathbf{A}^{(1)\top} \cdots \times_d \mathbf{A}^{(d)\top}$

```

1:  $(p'_1, \dots, p'_d, q'_1, \dots, q'_d)$  is my processor id
2: //All-gather input tensor  $\mathcal{X}$ 
3:  $\mathcal{X}_{p'_1 \dots p'_d} = \text{All-Gather}(\mathcal{X}, (p'_1, \dots, p'_d, *, \dots, *))$ 
4: //All-gather all input matrices
5: for  $i = 1, \dots, d$  do
6:    $\mathbf{A}_{p'_i q'_i}^{(i)} = \text{All-Gather}(\mathbf{A}^{(i)}, (*, \dots, *, p'_i, * \dots, *, q'_i, *))$ 
7: end for
8: //Perform local computations in a temporary tensor  $\mathcal{T}$ 
9:  $\mathcal{T} = \text{Local-Multi-TTM}(\mathcal{X}_{p'_1 \dots p'_d}, \mathbf{A}_{p'_1 q'_1}^{(1)}, \dots, \mathbf{A}_{p'_d q'_d}^{(d)})$ 
10: //Reduce-scatter the output tensor in  $\mathcal{Y}_{q'_1 \dots q'_d}$ 
11:  $\text{Reduce-Scatter}(\mathcal{Y}_{q'_1 \dots q'_d}, \mathcal{T}, (*, \dots, *, q'_1, \dots, q'_d))$ 

```

---

The data distribution model and cost analysis of the algorithm are similar to those of Alg. 5.1 and are presented in Appendix C.2.

Theorem 8.1 extends Theorem 5.1 to the  $d$ -dimensional case. A detailed proof can be found in Appendix C.2.2.

**THEOREM 8.1.** *There exist  $p_i, q_i$  with  $1 \leq p_i \leq n_i, 1 \leq q_i \leq r_i$  for  $i = 1, \dots, d$  such that Alg. 8.1 is communication optimal to within a constant factor.*

We present a comparison of the communication costs of Alg. 8.1 with the TTM-in-Sequence approach implemented in TuckerMPI [4] for 4/5/6-dimensional Multi-TTM computations in Appendix C.3. Our results are consistent with what we observe for 3-dimensional Multi-TTM computations in Section 6. When the input tensor is much larger than the output tensor and the number of entries in the output tensor is less than that of the matrices, our algorithm significantly reduces communication compared to the TTM-in-Sequence approach. As in the 3D case, when  $P \ll n/r$ , the extra computation is negligible when the TTM-in-Sequence approach is used locally to reduce computation.

**9. Conclusions.** In this work, we establish communication lower bounds for the parallel Multi-TTM computation and present an optimal parallel algorithm that organizes the processors in a  $2d$ -dimensional grid for  $d$ -dimensional tensors. By judi-

ciously selecting the processor grid dimensions, we prove that our algorithm attains the lower bounds to within a constant factor. To verify the theoretical analysis, we simulate Multi-TTM computations using a variety of values for the number of processors,  $P$ , the dimension,  $d$ , and sizes,  $n_i$  and  $r_i$ ; compute the communication costs of our algorithm corresponding to each simulation; and compute the optimal communication cost provided by the theoretical lower bound. These simulations show that the communication costs of the proposed algorithm are close to optimal. When one of the tensors is much larger than the other tensor, which is typical in compression algorithms based on the Tucker decomposition, our algorithm significantly reduces communication costs over the conventional approach of performing the computation as a sequence of tensor-times-matrix operations.

While applying HBL-inequalities (Lemma 3.3) to compute lower bounds, we encounter a  $\Delta$  matrix that is not full rank. In previous work where  $\Delta$  is square and nonsingular, the tightest constraint is derived from the unique solution to  $\Delta \mathbf{s} = \mathbf{1}$ . Our approach for Multi-TTM, where  $\Delta$  is rectangular and consistent, is to derive constraints for all vectors  $\mathbf{s}$  that satisfy  $\Delta \mathbf{s} = \mathbf{1}$  and then to select only two constraints which represent all constraints. We would like to explore how to select such representative constraints directly for general computations that involve a rank-deficient  $\Delta$  from HBL.

Motivated by the simulated communication cost comparisons, our next goal is to implement the parallel atomic algorithm and verify the performance improvement in practice. Further, because neither the atomic or TTM-in-sequence approach is always superior in terms of communication, we wish to explore hybrid algorithms to account for significant dimension reduction in some modes but modest reduction in others. Given the computation and communication capabilities of a parallel platform, it would also be interesting to study the computation-communication tradeoff for these two approaches and how to minimize the overall execution time in practice. Finally, this work considers that each processor has enough memory. A natural extension is to study communication lower bounds for Multi-TTM computations with limited memory sizes.

#### Appendix A. Proof of Theorem 3.4.

In this section, we prove Theorem 3.4 as it is written, instead of relying on the reader to derive this result from [6, Lemma 5.1]. This proof relies on two additional results. The first, [6, Lemma 2.2], states that the first constraint of the optimization problem is quasiconvex [1]. The second, [1, Lemma 3], states that satisfying the Karush-Kuhn-Tucker (KKT) conditions is sufficient for a solution to the optimization problem to be optimal as the optimization problem minimizes a differentiable convex function and the constraints are all differentiable quasiconvex functions.

*Proof of Theorem 3.4.* To begin we note that the objective and all but the first constraint are affine functions, which are differentiable, convex, and quasiconvex. The first constraint is differentiable, and it is quasiconvex in the positive orthant by [6, Lemma 2.2]. Thus the KKT conditions are sufficient to demonstrate the optimality of any solution by [1, Lemma 3], and we will prove the optimality of the solution  $\mathbf{x}^* = [x_1^* \ x_2^* \ \cdots \ x_d^*]$  by finding dual variables  $\mu_i^*$  for  $0 \leq i \leq d$  such that the KKT conditions are satisfied.

We now convert the problem to standard notation. The minimization objective function is

$$f(\mathbf{x}) = \sum_{j \in [d]} x_j,$$

and the constraints are given by

$$g_0(\mathbf{x}) = \frac{nr}{P} - \prod_{j \in [d]} x_j,$$

$$g_i(\mathbf{x}) = x_i - k_i \text{ for all } i \in [d].$$

Partial derivatives for  $j \in [d]$  are given by

$$\frac{\partial f}{\partial x_j}(\mathbf{x}) = 1,$$

$$\frac{\partial g_0}{\partial x_j}(\mathbf{x}) = - \prod_{\ell \in [d] - \{j\}} x_\ell,$$

$$\frac{\partial g_i}{\partial x_j}(\mathbf{x}) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else.} \end{cases}$$

The KKT conditions of  $(\mathbf{x}^*, \mu^*)$  are:

- *Primal feasibility*:  $g_i(\mathbf{x}^*) \leq 0$ , for  $0 \leq i \leq d$ .
- *Stationarity*:  $\frac{\partial f}{\partial x_j}(\mathbf{x}^*) + \sum_{i=0}^d \mu_i^* \frac{\partial g_i}{\partial x_j}(\mathbf{x}^*) = 0$ , for  $j \in [d]$ .
- *Dual feasibility*:  $\mu_i^* \geq 0$ , for  $0 \leq i \leq d$ .
- *Complementary slackness*:  $\mu_i^* g_i(\mathbf{x}^*) = 0$ , for  $0 \leq i \leq d$ .

Recall from the statement of Theorem 3.4 that  $K_I = \prod_{j=d-I+1}^d k_j$  and  $1 \leq I \leq d$  is defined such that

$$k_j < (K_{d-j+1}/P)^{1/(d-j+1)} \text{ for } 1 \leq j \leq d-I,$$

$$k_\ell \geq (K_{d-\ell+1}/P)^{1/(d-\ell+1)} \text{ for } d-I < \ell \leq d.$$

We claim that the optimal primal solution is

$$x_j^* = \begin{cases} k_j & \text{if } j \leq d-I, \\ (K_I/P)^{1/I} & \text{if } d-I < j \leq d. \end{cases}$$

and the optimal dual solution is

$$\mu_i^* = \begin{cases} (K_I/P)^{1/I} & \text{if } i = 0, \\ \frac{(K_I/P)^{1/I}}{k_i} - 1 & \text{if } 0 < i \leq d-I, \\ 0 & \text{if } d-I < i \leq d. \end{cases}$$

We now check that  $\mathbf{x}^*$  satisfies the primal feasibility condition. By direct verification (and the fact that  $nr = \prod_{j \in [d]} x_j$ ), we have  $g_0(\mathbf{x}^*) = 0$ . Clearly  $g_i(\mathbf{x}^*) = 0$  for all  $i \in [d-I]$ , as  $x_i^* = k_i$  for all  $i \in [d-I]$ . To see that  $g_i(\mathbf{x}^*) \leq 0$  for  $d-I < i \leq d$ , it is sufficient to recall that  $(K_I/P)^{1/I} \leq k_{d-I+1}$  by the definition of  $I$ , and that  $x_i^* = (K_I/P)^{1/I}$  and  $k_{d-I+1} \leq k_i$  for  $d-I < i \leq d$ .

Stationarity follows from direct verification of the condition for  $j \in [d]$ .

To check dual feasibility, we note that all the factors of  $\mu_0^*$  are positive, thus  $\mu_0^* > 0$ . To show that  $\mu_i^* > 0$  for  $i \in [d-I]$ , it is sufficient to show that  $k_{d-I} < (K_I/P)^{1/I}$  as  $k_1 \leq \dots \leq k_{d-I}$ . This is implied by  $k_{d-I} < (K_{I+1}/P)^{1/(I+1)} = (k_{d-I} K_I/P)^{1/I}$ ,

where the first inequality comes from the definition of  $I$ , and the second from the definition of the right products  $K_j$ .

Finally, complementary slackness is satisfied because  $g_i(\mathbf{x}^*) = 0$  for  $0 \leq i \leq d - I$ , and  $\mu_i = 0$  for  $d - I < i \leq d$ .

□

**Appendix B. Details for Simulated Evaluation of 3-Dimensional Multi-TTM.** In this section, we provide more details for the simulated evaluation of Alg. 5.1 and its comparison to the TTM-in-Sequence approach. The analysis of the communication optimality of Alg. 5.1 did not consider integrality constraints on the processor grid dimensions. The simulated evaluation in Section 6 considered all possible processor grid configurations using exhaustive search; we explain in Appendix B.1 a more efficient process for determining an optimal grid when  $P$  is a power of two. In Subsection 6.2, we also compare Alg. 5.1 against an implementation of the TTM-in-Sequence approach as implemented by TuckerMPI [4]. We argue in Appendix B.2 that this implementation is nearly communication optimal given the computation that it performs, validating our comparison against it. Figure 6 presents results relevant to both Appendices B.1 and B.2.

**B.1. Obtaining Integral Processor Grids for Alg. 5.1.** In order to determine the communication cost of Alg. 5.1, one must determine the processor grid. Obtaining  $p_i$  and  $q_i$  from the procedure in Section 5 may yield non-integer values. The following procedure allows us to convert these to integers under our assumption that all parameters are powers of 2. Recall that we consider  $P = pq$  with  $p = p_1 p_2 p_3$  and  $q = q_1 q_2 q_3$ .

If  $\lfloor \log_2(p) + 0.5 \rfloor = \lfloor \log_2(p) \rfloor$ , then we set  $p = 2^{\lfloor \log_2(p) \rfloor}$ , otherwise we set  $p = 2^{\lceil \log_2(p) \rceil}$ , distributing the modification evenly between  $p_1, p_2$ , and  $p_3$ . Now, we keep  $p = p_1 p_2 p_3$  constant, and convert each  $p_i$  to an integer. We set  $p_1 = 2^{\lfloor \log_2(p_1) + 0.5 \rfloor}$  distributing the changes evenly among  $p_2$  and  $p_3$ . To see that our new value of  $p_1$  must still be smaller than  $n_1$ , we note that our original  $p_1$  was less than  $n_1$  which is a power of 2 by our assumption. If we increased  $p$  in our first step, then distributing the modifications evenly between  $p_1, p_2$  and  $p_3$  increased them by at most  $2^{1/6}$ . Thus  $p_1 \leq n_1$  will imply that  $\lfloor \log_2(p_1 \cdot 2^{1/6}) + 0.5 \rfloor \leq \log_2(n_1)$ . Note that this most recent modification to  $p_1$  changes  $p_2$  and  $p_3$ . Then, we set  $p_2 = 2^{\lfloor \log_2(p_2) + 0.5 \rfloor}$  and adapt  $p_3$  accordingly. A similar argument to what is used for  $p_1$  will show that  $p_2$  and  $p_3$  are also not larger than their corresponding dimensions. Having completed our work on the processor dimensions associated with the first tensor, we set  $q = \frac{P}{p}$  distributing the changes evenly among the  $q_i$ , then force each  $q_i$  to be an integer following the same procedure as for the  $p_i$ .

We denote the communication cost of Alg. 5.1 for the grid determined using this method by *Alg. 5.1 (fast)* and the communication cost using exhaustive search by *Alg. 5.1 (best)*. We note that this procedure can increase the total number of accessed elements of any variable at most 4 times, but we see in Fig. 6 that the communication costs of both procedures are exactly the same for the examples we consider. These problems match those presented in Fig. 4.

**B.2. TTM-in-Sequence Lower Bounds.** Here we discuss communication lower bounds for the TTM-in-Sequence approach with cubical tensors. There has not been any proven bound for this approach other than individual bounds for each TTM (a single matrix multiply) computation, assuming the sequence of TTMs has been



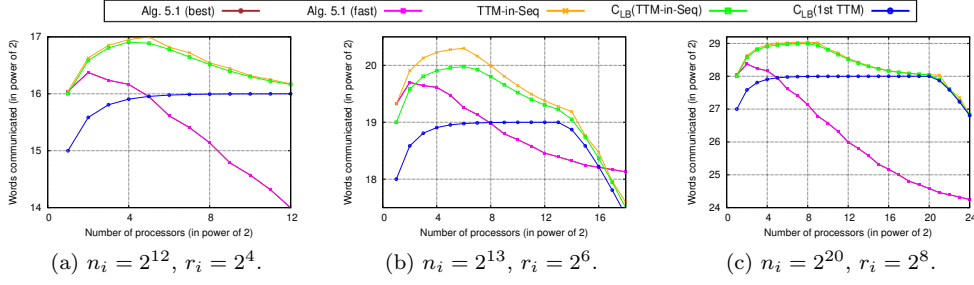


Fig. 6: Communication cost comparison of Alg. 5.1 using best processor grid against fast method and of the *TTM-in-Sequence* approach implemented by TuckerMPI against the lower bounds. *Alg. 5.1 (fast)* and *Alg. 5.1 (best)* are the same for all the configurations.

specified. The sum of individual bounds provides a communication lower bound for this approach. We obtain the tightest (and obtainable) lower bound for each TTM from [1], which depends on the relative matrix dimensions and number of processors, and represent the sum by  $C_{LB}(TTM-in-Seq)$ . We also note that  $C_{LB}(TTM-in-Seq)$  may not be always attainable as data distributions for two successive TTMs may be non-compatible and require extra communication. When the input tensor dimensions are much larger than the output tensor dimensions, most of the computation and communication occurs in the first TTM, so we also consider the communication lower bound of only that matrix multiplication, which also provides a valid lower bound for the entire TTM-in-Sequence computation. Recall that we obtain the algorithmic cost of TTM-in-Sequence by exhaustively searching for the best processor grid configuration given the communication costs specified by (6.1). Figure 6 shows a comparison of *TTM-in-Seq* and  $C_{LB}(TTM-in-Seq)$  for the tensor dimensions presented in Fig. 4. We can see that the communication costs of *TTM-in-Seq* are very close to  $C_{LB}(TTM-in-Seq)$ , the largest differences are 7.9% for Fig. 6a at  $P = 2^5$ , 25% for Fig. 6b at  $P = 2^6$ , and 9.3% for Fig. 6c at  $P = 2^{21}$ . Comparing  $C_{LB}(1st\ TTM)$  and  $C_{LB}(TTM-in-Seq)$ , we see that for these examples at least half the communication of the entire TTM-in-Sequence is required by the first TTM, and it is completely dominated by the first TTM when  $P$  is large.

**Appendix C. General Case.** In this section, we provide results and proofs for general multi-TTM computations. We present the proof of communication lower bounds, Theorem 7.2, in Appendix C.1. In Appendix C.2, we discuss computation and communication costs of Alg. 8.1, and how to select processor grid dimensions such that our algorithm is communication optimal. We present a comparison of the communication costs of our algorithm with a TTM-in-Sequence approach for 3/4/5/6-dimensional Multi-TTM computations in Appendix C.3.

### C.1. Communication Lower Bounds for General Multi-TTM.

*Proof of Theorem 7.2.* Let  $F$  be the set of loop indices associated with the  $(d+1)$ -ary multiplications performed by a processor. As we assumed the algorithm is computationally load balanced,  $|F| = nr/P$ . We define  $\phi_{\mathbf{x}}(F)$ ,  $\phi_{\mathbf{y}}(F)$  and  $\phi_j(F)$  to be the projections of  $F$  onto the indices of the arrays  $\mathbf{x}$ ,  $\mathbf{y}$ , and  $\mathbf{A}^{(j)}$  for  $1 \leq j \leq d$  which correspond to the elements of the arrays that must be accessed by the processor.

We use Lemma 3.3 to obtain a lower bound on the number of array elements that must be accessed by the processor. The matrix corresponding to the projections above is given by

$$\Delta = \begin{bmatrix} \mathbf{I}_{d \times d} & \mathbf{1}_d & \mathbf{0}_d \\ \mathbf{I}_{d \times d} & \mathbf{0}_d & \mathbf{1}_d \end{bmatrix}.$$

Here  $\mathbf{1}_d$  and  $\mathbf{0}_d$  denote the  $d$ -dimensional vectors of all ones and zeros, respectively, and  $\mathbf{I}_{d \times d}$  denotes the  $d \times d$  identity matrix. As before we define

$$\mathcal{C} = \{\mathbf{s} \in [0, 1]^{d+2} : \Delta \cdot \mathbf{s} \geq \mathbf{1}\}.$$

We recall that  $\mathbf{1}$  represents a vector of all ones. As in the proof of Theorem 4.3,  $\Delta$  is not full rank, so we again consider each vector  $\mathbf{v} \in \mathcal{C}$  such that  $\Delta \cdot \mathbf{v} = \mathbf{1}$ . Such a vector  $\mathbf{v}$  is of the form  $[a \ \cdots \ a \ 1-a \ 1-a]$  where  $0 \leq a \leq 1$ . Thus, we obtain

$$\frac{nr}{P} \leq \left( \prod_{j \in [d]} |\phi_j(F)| \right)^a (|\phi_{\mathbf{x}}(F)| |\phi_{\mathbf{y}}(F)|)^{1-a}.$$

Similar to the 3D case, the above constraint is equivalent to  $\frac{nr}{P} \leq \prod_{j \in [d]} |\phi_j(F)|$  and  $\frac{nr}{P} \leq |\phi_{\mathbf{x}}(F)| |\phi_{\mathbf{y}}(F)|$ .

Clearly a projection onto an array can not be larger than the array itself, thus  $|\phi_{\mathbf{x}}(F)| \leq n$ ,  $|\phi_{\mathbf{y}}(F)| \leq r$ , and  $|\phi_j(F)| \leq n_j r_j$  for  $1 \leq j \leq d$ .

As the constraints related to the projections of matrices and tensors are disjoint, we solve them separately and then sum the results to get a lower bound on the set of elements that must be accessed by the processor. We obtain a lower bound on  $A$ , the number of elements of the matrices that must be accessed by the processor by using Corollary 7.1, and a lower bound on  $B$ , the number of elements of the tensors that must be accessed by the processor by using Corollary 4.2. By summing both, we get the positive terms of the lower bound.

To bound the sends or receives, we consider how much data the processor could have had at the beginning or at the end of the computation. Assuming there is exactly one copy of the data at the beginning and at the end of the computation, there must exist a processor which has access to at most  $1/P$  of the elements of the arrays at the beginning or at the end of the computation. By employing the previous analysis, this processor must access  $A + B$  elements of the arrays, but can only have  $\frac{n}{P} + \frac{r}{P} + \sum_{j \in [d]} \frac{n_j r_j}{P}$  elements of the arrays stored. Thus it must perform the specified amount of sends or receives.  $\square$

As we did previously, we denote the lower bound of Theorem 7.2 by  $LB$  and use it extensively while determining the optimal processor grid dimensions for our algorithm in Appendix C.2.2.

**C.2. Parallel Algorithm for General Multi-TTM.** Here we discuss our data distribution model for Alg. 8.1.  $\mathcal{X}_{p'_1 \dots p'_d}$  and  $\mathcal{Y}_{q'_1 \dots q'_d}$  denote the subtensors of  $\mathcal{X}$  and  $\mathcal{Y}$  owned by processors  $(p'_1, \dots, p'_d, *, \dots, *)$  and  $(*, \dots, *, q'_1, \dots, q'_d)$ , respectively.  $\mathbf{A}_{p'_i q'_i}^{(i)}$  denotes the submatrix of  $\mathbf{A}^{(i)}$  owned by processors  $(*, \dots, *, p'_i, *, \dots, *, q'_i, *, \dots, *)$ . We impose that there is one copy of data in the system at the beginning and the end of the computation, and each subarray is distributed evenly among the set of processors which own the data. When Alg. 8.1 completes,  $\mathcal{Y}_{q'_1 \dots q'_d}$  is distributed evenly among processors  $(*, \dots, *, q'_1, \dots, q'_d)$ .

We recall that  $\prod_{i=1}^d p_i$  and  $\prod_{i=1}^d q_i$  are denoted by  $p$  and  $q$ , respectively.

**C.2.1. Cost Analysis.** Now we analyze computation and communication costs of the algorithm. As before, the local Multi-TTM computation in Line 9 can be performed as a sequence of TTM operations to minimize the number of arithmetic operations. Assuming the TTM operations are performed in their order, first with  $\mathbf{A}^{(1)}$ , then with  $\mathbf{A}^{(2)}$ , and so on until the last is performed with  $\mathbf{A}^{(d)}$ , then each processor performs  $\sum_{k=1}^d \left( 2 \prod_{i=1}^k \frac{r_i}{q_i} \prod_{j=k}^d \frac{n_j}{p_j} \right)$  local computations in Line 9. In Line 11, each processor also performs  $(1 - \frac{q}{P}) \frac{r}{q}$  computations due to the Reduce-Scatter operation.

Communication occurs only in All-Gather and Reduce-Scatter collectives in Lines 3, 6, and 11. Line 3 specifies  $\frac{P}{p}$  simultaneous All-Gathers, Line 6 specifies  $\frac{P}{p_i q_i}$  simultaneous All-Gathers in the  $i$ th loop iteration, and Line 11 specifies simultaneous  $\frac{P}{q}$  Reduce-Scatters. Each processor is involved in one All-Gather involving the input tensor,  $d$  All-Gathers involving input matrices and one Reduce-Scatter involving the output tensor.

As before, we assume bandwidth and latency optimal algorithms are used for the All-Gather and Reduce-Scatter collectives. Hence the bandwidth costs of the All-Gather operations are  $(1 - \frac{p}{P}) \frac{n}{p}$  for Line 3, and  $\sum_{i=1}^d (1 - \frac{p_i q_i}{P}) \frac{n_i r_i}{p_i q_i}$  for the  $d$  iterations of Line 6. The bandwidth cost of the Reduce-Scatter operation in Line 11 is  $(1 - \frac{q}{P}) \frac{r}{q}$ . Hence the overall bandwidth cost of Alg. 8.1 for each processor is  $\frac{n}{p} + \frac{r}{q} + \sum_{i=1}^d \frac{n_i r_i}{p_i q_i} - \left( \frac{n+r+\sum_{i=1}^d n_i r_i}{P} \right)$ . The latency costs are  $\log \left( \frac{P}{p} \right)$  and  $\log \left( \frac{P}{q} \right)$  for Lines 3 and 11 respectively, and  $\sum_{i=1}^d \log \left( \frac{P}{p_i q_i} \right)$  for the  $d$  iterations of Line 6. Thus the overall latency cost of Alg. 8.1 is  $\log \left( \frac{P}{p} \right) + \sum_{i=1}^d \log \left( \frac{P}{p_i q_i} \right) + \log \left( \frac{P}{q} \right) = d \log(P)$ .

**C.2.2. Selection of  $p_j$  and  $q_j$  in Alg. 8.1.** Similar to Subsection 5.2, we prove the optimality of Alg. 8.1 by selecting  $p_j$  and  $q_j$  such that the communication cost of Alg. 8.1 matches the communication lower bounds of the Multi-TTM computation (Theorem 7.2), and then determine how much additional communication is required to ensure the selected  $p_j, q_j$  satisfy the additional constraints  $(1 \leq p_j \leq n_j, 1 \leq q_j \leq r_j)$  necessary for processor grid dimensions.

As  $\mathbf{X}$  and  $\mathbf{Y}$  are  $d$ -dimensional tensors, we have  $n_j, r_j \geq 2$  for  $1 \leq j \leq d$ . For better readability, similar to the previous convention, we denote  $\frac{\sum_{j \in [d]} n_j r_j + r + n}{P}$  by  $O$ .

*Proof of Theorem 8.1.* As we did previously, we break our analysis into 2 scenarios which are further broken down into all possible cases.

In each case, we obtain  $\hat{p}_j$  and  $\hat{q}_j$  such that the terms in the communication cost match the corresponding lower bound terms and satisfy at least one of the two sets of constraints,  $1 \leq \hat{p}_j \leq n_j, 1 \leq \hat{q}_j$  or  $1 \leq \hat{q}_j \leq r_j, 1 \leq \hat{p}_j$  for  $1 \leq j \leq d$ . We handle all cases of both scenarios together in the end, and adapt these values to get  $p_j$  and  $q_j$  which respect both lower and upper bounds for all values of  $j$ . Then we determine how much additional communication may be required. We denote  $\prod_{i=1}^d \hat{p}_i$  and  $\prod_{i=1}^d \hat{q}_i$  by  $\hat{p}$  and  $\hat{q}$ .

• **Scenario I** ( $P < \frac{n}{r}$ ): This scenario corresponds to the first case of the tensor term in  $LB$ . Thus, we set  $\hat{p}_j, \hat{q}_j$  in such a way that the tensor terms in the communication cost match the tensor terms of  $LB$ :

$$(C.1) \quad \hat{p} = P, \hat{q} = 1.$$

This implies  $\hat{q}_j = 1$  for  $1 \leq j \leq d$ . We break this scenario into  $d$  cases parameterized

by  $I$ :  $\frac{N_{I-1}R_{I-1}}{(n_{d-I+1}r_{d-I+1})^{I-1}} \leq P < \min \left\{ \frac{N_I R_I}{(n_{d-I}r_{d-I})^I}, \frac{n}{r} \right\}$ . The cases degenerate to  $P < \min \left\{ \frac{N_1 R_1}{n_{d-1}r_{d-1}}, \frac{n}{r} \right\}$ , when  $I = 1$ , and  $\frac{N_{d-1}R_{d-1}}{(n_1r_1)^{d-1}} \leq P < \frac{n}{r}$  when  $I = d$ . Setting the matrix communication costs to the matrix terms of the lower bound in the corresponding cases yields

$$(C.2) \quad \frac{n_j r_j}{\hat{p}_j \hat{q}_j} = n_j r_j \text{ if } 1 \leq j \leq d-I, \quad \frac{n_j r_j}{\hat{p}_j \hat{q}_j} = \left( \frac{N_I R_I}{P} \right)^{\frac{1}{I}} \text{ if } d-I < j \leq d.$$

Thus,  $\hat{q}_j = 1$  for all  $1 \leq j \leq d$ ,  $\hat{p}_j = 1$  if  $1 \leq j \leq d-I$  and  $\hat{p}_j = n_j r_j \left( \frac{P}{N_I R_I} \right)^{\frac{1}{I}}$  if  $d-I < j \leq d$  to satisfy (C.1) and (C.2). Note that when  $I = 1$ ,  $p_d = P \geq 1$ , and for the other values of  $I$ ,  $p_j \geq 1$  because  $n_1 r_1 \leq \dots \leq n_d r_d$  and  $\frac{N_{I-1}R_{I-1}}{(n_{d-I+1}r_{d-I+1})^{I-1}} \leq P$ . Additionally we have that  $1 = \hat{q}_j < r_j$  for  $1 \leq j \leq d$ . However, we are not able to ensure  $\hat{p}_j \leq n_j$  when  $d-I < j \leq d$ . We will handle all cases of both scenarios together as they require the same analysis.

• **Scenario II** ( $\frac{n}{r} \leq P$ ): This scenario corresponds to the second case of the tensor term in  $LB$ . Thus, we set  $\hat{p}_i, \hat{q}_i$  in such a way that

$$(C.3) \quad \frac{n}{\hat{p}} = \frac{r}{\hat{q}} = \left( \frac{nr}{P} \right)^{1/2}.$$

Again, we break this scenario into  $d$  cases parameterized by  $I$ :

$$\max \left\{ \frac{N_{I-1}R_{I-1}}{(n_{d-I+1}r_{d-I+1})^{I-1}}, \frac{n}{r} \right\} \leq P < \frac{N_I R_I}{(n_{d-I}r_{d-I})^I}.$$

The cases degenerate to  $P < \frac{N_1 R_1}{n_{d-1}r_{d-1}}$ , when  $I = 1$ , and  $\max \left\{ \frac{N_{d-1}R_{d-1}}{(n_1r_1)^{d-1}}, \frac{n}{r} \right\} \leq P$  when  $I = d$ .

Setting the matrix communication costs to match the corresponding matrix terms in the lower bound yields

$$(C.4) \quad \frac{n_j r_j}{\hat{p}_j \hat{q}_j} = n_j r_j \text{ if } 1 \leq j \leq d-I, \quad \frac{n_j r_j}{\hat{p}_j \hat{q}_j} = \left( \frac{N_I R_I}{P} \right)^{\frac{1}{I}} \text{ if } d-I < j \leq d.$$

Thus we set  $\hat{q}_j = \hat{p}_j = 1$  for all  $1 \leq j \leq d-I$ . When  $2 \leq I \leq d$ , the equations above do not uniquely determine  $\hat{p}_j, \hat{q}_j$  for  $d-I < j \leq d$ . However, setting  $\hat{p}_j = n_j \left( \frac{nP}{rN_I^2} \right)^{1/2I}$

and  $\hat{q}_j = r_j \left( \frac{rP}{nR_I^2} \right)^{1/2I}$  for  $d-I < j \leq d$  satisfies equations C.3, C.4 in all cases. Note that we cannot ensure lower and upper bounds on  $\hat{p}_j$  and  $\hat{q}_j$ . We now look for new solutions to the equations twice. First, we ensure that all lower bounds are respected, i.e.,  $1 \leq \hat{p}_j$  and  $1 \leq \hat{q}_j$ , and then we guarantee that all upper bounds of  $\hat{p}_j$  or  $\hat{q}_j$  are satisfied, i.e.,  $\hat{p}_j \leq n_j$  or  $\hat{q}_j \leq r_j$ .

As we compared communication cost of each term with its corresponding lower bound to obtain  $\hat{p}_j$  and  $\hat{q}_j$ , we have  $1 \leq \hat{p}_j \hat{q}_j \leq n_j r_j$  for  $1 \leq j \leq d$ ,  $1 \leq \hat{p} \leq n$  and  $1 \leq \hat{q} \leq r$  in all  $d$  cases. However, we may not have  $1 \leq \hat{p}_j$  or  $1 \leq \hat{q}_j$  for some  $j$ . We now seek new solutions that are all greater than 1. First we will increase all  $\hat{q}_j, \hat{p}_j$  that are less than 1 in a way that preserves products  $\hat{p}_j \hat{q}_j$  but does not preserve  $\hat{p}$  and  $\hat{q}$ . Then we will adjust  $\hat{p}_j$  and  $\hat{q}_j$  to force the products  $\hat{p}$  and  $\hat{q}$  back to their initial values.

Let  $q^b$  denote the product of all  $\hat{q}_j$  such that  $\hat{q}_j < 1$ , and  $p^b$  denote the product of all  $\hat{p}_j$  such that  $\hat{p}_j < 1$ . Without loss of generality, if  $q^b \leq p^b$ , set  $a = 1$ ,  $\hat{p}^{orig} = \hat{p}$ , and  $\hat{q}^{orig} = \hat{q}$ . We perform the following updates:

Looping over the index  $j$  from 1 to  $d$ , if  $\hat{q}_j < 1$  then set  $a = a \cdot \hat{q}_j$ ,  $\hat{p}_j = \hat{p}_j \hat{q}_j$ ,  $\hat{q}_j = 1$ ; else if  $\hat{p}_j < 1$  then set  $a = a / \hat{p}_j$ ,  $\hat{q}_j = \hat{p}_j \hat{q}_j$ ,  $\hat{p}_j = 1$ . This step preserves all products  $\hat{p}_j \hat{q}_j$  and enforces  $1 \leq \hat{p}_j, 1 \leq \hat{q}_j$  for  $1 \leq j \leq d$ , but it does not preserve  $\hat{p}, \hat{q}$ . At the end of this step, we have  $a = q^b / p^b < 1$ ,  $\hat{p} = a \cdot \hat{p}^{orig}$ , and  $\hat{q} = \hat{q}^{orig} / a$ . In order to force  $\hat{p}$  and  $\hat{q}$  to match their initial values, we decrease some  $\hat{q}_j$  in such a way that  $\hat{q}$  is decreased by a factor of  $a$ . This is possible because  $1 \leq \hat{q}^{orig} = a \cdot \hat{q}$ . Looping over the index  $j$  from 1 to  $d$ , if  $\hat{q}_j > 1$  then set  $\hat{q}_j^{prev} = \hat{q}_j$ ,  $\hat{q}_j = \max(1, a \cdot \hat{q}_j)$ ,  $a = a \left( \frac{\hat{q}_j^{prev}}{\hat{q}_j} \right)$ ,  $\hat{p}_j = \hat{p}_j \left( \frac{\hat{q}_j^{prev}}{\hat{q}_j} \right)$ . At the end of this step,  $\hat{q}$  has been decreased by a factor of  $q^b / p^b$ ,  $\hat{p}_j \hat{q}_j$  were all preserved, and thus,  $\hat{p}$  has been increased by a factor of  $p^b / q^b$ , hence  $\hat{q} = \hat{q}^{orig}$  and  $\hat{p} = \hat{p}^{orig}$ . After the above updates, we have  $1 \leq \hat{p}_j, 1 \leq \hat{q}_j$  for  $1 \leq j \leq d$ , and the products match the initial products thus are valid solutions to the original equations. If  $p_b < q_b$  we would perform the same process, but changing the actions on the  $\hat{q}_j$  to be performed on the  $\hat{p}_j$  and vice versa.

We now handle upper bounds of  $\hat{p}_j$  and  $\hat{q}_j$ . If  $\exists j, k$  such that  $\hat{p}_j > n_j$  and  $\hat{q}_k > r_k$ , we again seek new solutions such that all  $\hat{p}_j$  or all  $\hat{q}_j$  satisfy upper bounds while respecting lower bounds of all variables. Let  $p^t$  and  $q^t$  denote the products of all  $\hat{p}_j$  and  $\hat{q}_j$ , respectively, such that  $\hat{p}_j \hat{q}_j \neq 1$ . As  $\forall j, 1 \leq \hat{p}_j \hat{q}_j \leq n_j r_j$ , therefore  $p^t$  is not more than the product of the corresponding  $n_j$  and/or  $q^t$  is not more than the product of the corresponding  $r_j$ . If the first constraint is satisfied, then we perform the following updates:

Looping over the index  $j$  from 1 to  $d$ , if  $p^t > 1$  and  $\hat{p}_j \hat{q}_j \neq 1$  then set  $a = \hat{p}_j \hat{q}_j$ ,  $\hat{p}_j = \min(p^t, n_j)$ ,  $\hat{q}_j = \frac{a}{\hat{p}_j}$ ,  $p^t = \frac{p^t}{\hat{p}_j}$ . After the above updates, we have  $1 \leq \hat{p}_j \leq n_j, 1 \leq \hat{q}_j$  for  $1 \leq j \leq d$ , and  $\hat{p}_j \hat{q}_j, \hat{p}$  and  $\hat{q}$  are back to their original values. If the first constraint is not satisfied, we would perform the same process on  $\hat{q}_j$  instead of  $\hat{p}_j$ .

Now for all cases of both scenarios, we know that  $1 \leq \hat{p}_j$  and  $1 \leq \hat{q}_j$  for  $1 \leq j \leq d$ , and either  $\hat{p}_j \leq n_j$  for  $1 \leq j \leq d$  or  $\hat{q}_j \leq r_j$  for  $1 \leq j \leq d$ . It remains to adapt  $\hat{p}_j$  and  $\hat{q}_j$  such that both  $\hat{p}_j \leq n_j$  and  $\hat{q}_j \leq r_j$  for  $1 \leq j \leq d$ . We obtain  $p_1, \dots, p_d, q_1, \dots, q_d$  from  $\hat{p}_j$  and  $\hat{q}_j$  such that  $p_1 \cdots p_d = \hat{p}$  and  $q_1 \cdots q_d = \hat{q}$ . The intuition is to maintain the tensor communication terms in the lower bound.

Initially, we set  $p_j = \hat{p}_j$  and  $q_j = \hat{q}_j$  for  $1 \leq j \leq d$ . If  $1 \leq p_j \leq n_j$  and  $1 \leq q_j \leq r_j$  for  $1 \leq j \leq d$ , then  $\sum_{j \in [d]} \frac{n_j r_j}{p_j q_j} = \sum_{j \in [d]} \frac{n_j r_j}{\hat{p}_j \hat{q}_j}$  and the communication cost exactly matches the lower bound. Otherwise, we adapt the values and determine the effects on the matrix communication costs. We recall that due of our particular selections of  $\hat{p}_j$  and  $\hat{q}_j$ ,  $\nexists j, \ell \in [d]$  such that  $\hat{p}_j > n_j$  and  $\hat{q}_\ell > r_\ell$ . If  $\hat{p}_j > n_j$  for some  $j \in [d]$ , then we iterate over the index  $j$  from  $d$  to 1 setting  $p_j = \min \left\{ n_j, \frac{\hat{p}}{\prod_{\ell \in [d] - \{j\}} p_\ell} \right\}$ . We iterate again from  $d$  to 1 with the same expression. Iterating twice ensures that all updates are visible to all  $p_j$ .

Now we assess how much additional communication is required for the matrices.

As  $\hat{p}_j > n_j$  for some  $j$ , it must be the case that  $\hat{p} \geq 2$ . Thus

$$\begin{aligned}
\sum_{j \in [d]} \frac{n_j r_j}{p_j q_j} &\leq \sum_{j \in [d]} \max \left\{ \frac{n_j r_j}{\hat{p}_j \hat{q}_j}, \frac{r_j}{\hat{q}_j} \right\} \\
&= \sum_{j \in [d]} \left( \frac{n_j r_j}{\hat{p}_j \hat{q}_j} + \frac{r_j}{\hat{q}_j} - \min \left\{ \frac{n_j r_j}{\hat{p}_j \hat{q}_j}, \frac{r_j}{\hat{q}_j} \right\} \right) \\
&< \sum_{j \in [d]} \left( \frac{n_j r_j}{\hat{p}_j \hat{q}_j} + \frac{r_j}{\hat{q}_j} \right) - (d-1) \\
&\leq \sum_{j \in [d]} \frac{n_j r_j}{\hat{p}_j \hat{q}_j} + \frac{r}{\hat{q}} \\
&< \sum_{j \in [d]} \frac{n_j r_j}{\hat{p}_j \hat{q}_j} + 2 \left( \frac{r}{\hat{q}} - \frac{r}{\hat{p} \hat{q}} \right) \\
&= \sum_{j \in [d]} \frac{n_j r_j}{\hat{p}_j \hat{q}_j} + 2 \left( \frac{r}{\hat{q}} - \frac{r}{P} \right).
\end{aligned}$$

Similarly, if  $\exists j \in [d]$  such that  $\hat{q}_j > r_j$ , the same update can be performed to the  $q_j$ , and we obtain  $\sum_{j \in [d]} \frac{n_j r_j}{p_j q_j} < \sum_{j \in [d]} \frac{n_j r_j}{\hat{p}_j \hat{q}_j} + 2 \left( \frac{n}{\hat{p}} - \frac{n}{P} \right)$ .

Therefore,  $\sum_{j \in [d]} \frac{n_j r_j}{p_j q_j} + \frac{r}{q} + \frac{n}{p} - O \leq 3 \left( \sum_{j \in [d]} \frac{n_j r_j}{\hat{p}_j \hat{q}_j} + \frac{r}{\hat{q}} + \frac{n}{\hat{p}} - O \right) = 3LB$ .  $\square$

**C.3. Simulated Evaluation.** Similar to Section 6, we compare communication costs of our algorithm and a TTM-in-Sequence approach implemented in the TuckerMPI library. We again restrict to cases where all dimensions are powers of 2, and vary the number of processors  $P$  from 2 to  $P_{\max}$  in multiples of 2, where  $P_{\max} = \min\{n_1 r_1, \dots, n_d r_d, n, r\}$ .

Like Section 6, we look at all possible processor grid dimensions and represent the minimum communication costs of our algorithm and TuckerMPI algorithm by *Alg. 8.1 (best)* and *TTM-in-Seq*, respectively. The TTM-in-Sequence approach described in [4] organizes  $P$  in a  $d$ -dimensional  $\tilde{p}_1 \times \dots \times \tilde{p}_d$  logical processor grid. Assuming TTMs are performed in increasing mode order, the overall communication cost of this algorithm is

$$\begin{aligned}
&\frac{r_1 n_2 \dots n_d}{\frac{P}{\tilde{p}_1}} + \frac{r_1 r_2 n_3 \dots n_d}{\frac{P}{\tilde{p}_2}} + \dots + \frac{r_1 r_2 \dots r_d}{\frac{P}{\tilde{p}_d}} - \frac{r_1 n_2 \dots n_d + r_1 r_2 n_3 \dots n_d + \dots + r_1 r_2 \dots r_d}{P} \\
&\quad + \frac{n_1 r_1}{\tilde{p}_1} + \dots + \frac{n_d r_d}{\tilde{p}_d} - \frac{n_1 r_1 + \dots + n_d r_d}{P}.
\end{aligned}$$

The first line corresponds to tensor communication and the second line corresponds to matrix communication. As mentioned earlier, the TTM-in-Sequence approach forms a tensor after each TTM computation. Each positive term of the first line corresponds to the number of entries of such a tensor accessed by a processor in TuckerMPI.

We again look at cases where the input tensors are large and the output tensors are small. Figure 7 shows comparison of *Alg. 8.1 (best)* and *TTM-in-Seq* with our communication lower bounds (*LB*) for 3/4/5-dimensional Multi-TTM computations. For  $P = 2$ , both approaches perform the same amount of communication. After that, the total number of accessed elements in both approaches decreases, however the rate of owned elements decreases at the faster rate. Hence we see slight increase in both curves. This behavior continues roughly till  $2^{n_i - r_i}$  processors for *TTM-in-Seq* curve. In this region, TuckerMPI selects  $\tilde{p}_1 = \dots = \tilde{p}_{d-1} = 1$

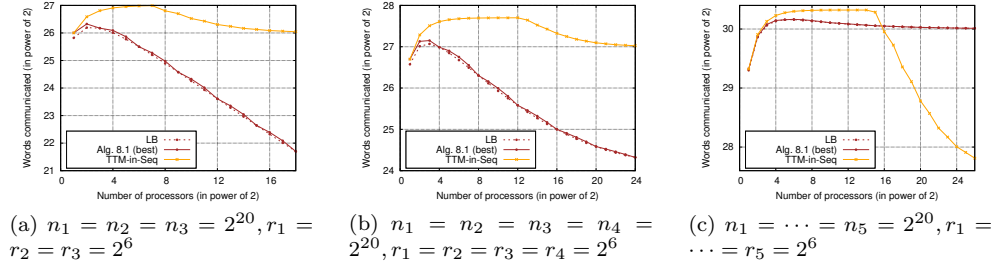


Fig. 7: Communication cost comparison of Alg. 8.1 and the TTM-in-Sequence approach implemented by the TuckerMPI library. Note that  $LB$  is a communication lower bound for atomic Multi-TTM algorithms, not for the TTM-in-Sequence approach. Communication cost of our approach (*Alg. 8.1 (best)*) is very close to the lower bound ( $LB$ ).

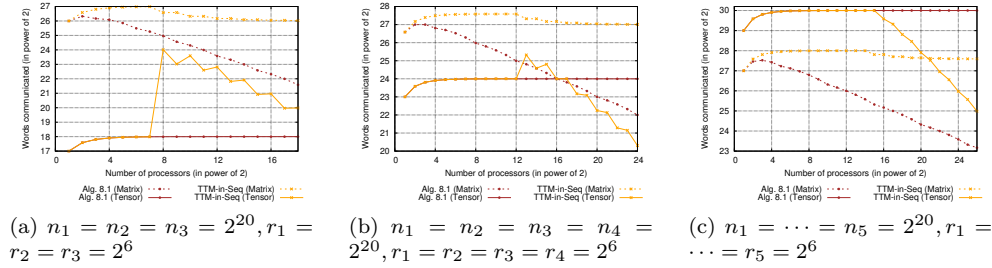


Fig. 8: Matrix and Tensor communication costs in Alg. 8.1 and the TTM-in-Sequence approach.

and  $\tilde{p}_d = P$ . Our algorithm selects  $p_1 \approx \dots \approx p_d$  and  $q_1 = \dots = q_d = 1$ . These processor grid dimensions result in the same tensor communication cost for both approaches. However our approach reduces matrix communication cost roughly

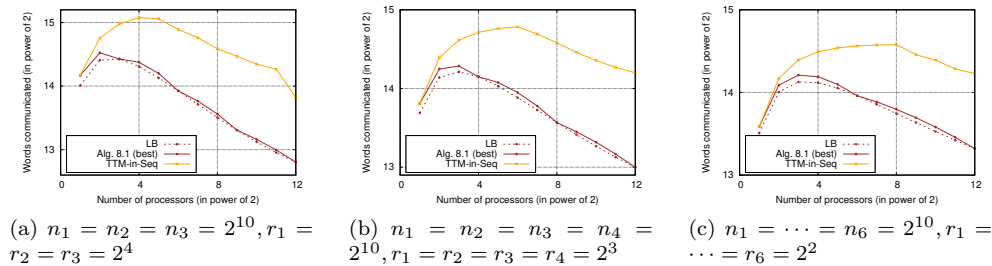


Fig. 9: Communication cost comparison of Alg. 8.1 and the TTM-in-Sequence approach for 3/4/6-dimensional Multi-TTM computations.

$(1 - \frac{1}{d})P^{\frac{1}{d}}$  times, hence it is better than the TTM-in-Sequence approach. Fig. 8 shows the distribution of matrix and tensor communication costs in both approaches. In general, our approach significantly minimizes the matrix communication costs in all the plots and is better when the number of the entries in the output tensor is less than that of the matrices. When the communication cost is dominated by the output tensor, our approach is outperformed by the TTM-in-Sequence approach, which is the case in Fig. 7c.

Now we consider a different set of experiments. Here the number of entries in the input tensor is  $(2^{10})^d$  for  $d$ -dimensional computation. We fix the number of entries in the output tensor to  $2^{12}$  and present comparisons of the considered approaches in Figure 9 for 3/4/6-dimensional Multi-TTM computations. These dimensions allow both tensors to be cubical. As the number of entries in the matrices are greater than the number of entries in the output tensor, our approach is always superior to the TTM-in-Sequence approach.

**Acknowledgments.** This work is supported by the National Science Foundation under Grant No. CCF-1942892 and OAC-2106920. This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (Grant agreement No. 810367).

## REFERENCES

- [1] H. AL DAAS, G. BALLARD, L. GRIGORI, S. KUMAR, AND K. ROUSE, *Tight memory-independent parallel matrix multiplication communication lower bounds*, 2022, <https://arxiv.org/abs/2205.13407>.
- [2] W. AUSTIN, G. BALLARD, AND T. G. KOLDA, *Parallel tensor compression for large-scale scientific data*, in Proceedings of the 30th IEEE International Parallel and Distributed Processing Symposium, May 2016, pp. 912–922, <https://doi.org/10.1109/IPDPS.2016.67>.
- [3] G. BALLARD, E. CARSON, J. DEMMEL, M. HOEMMEN, N. KNIGHT, AND O. SCHWARTZ, *Communication lower bounds and optimal algorithms for numerical linear algebra*, Acta Numerica, 23 (2014), pp. 1–155, <https://doi.org/10.1017/S0962492914000038>.
- [4] G. BALLARD, A. KLINEX, AND T. G. KOLDA, *TuckerMPI: A parallel C++/MPI software package for large-scale data compression via the tucker tensor decomposition*, ACM Trans. Math. Softw., 46 (2020), <https://doi.org/10.1145/3378445>.
- [5] G. BALLARD, N. KNIGHT, AND K. ROUSE, *Communication lower bounds for matricized tensor times Khatri-Rao product*, in 2018 IEEE International Parallel and Distributed Processing Symposium (IPDPS), 2018, pp. 557–567, <https://doi.org/10.1109/IPDPS.2018.00065>.
- [6] G. BALLARD AND K. ROUSE, *General memory-independent lower bound for MTTKRP*, in Proceedings of the 2020 SIAM Conference on Parallel Processing for Scientific Computing (PP), 2020, pp. 1–11, <https://doi.org/10.1137/1.9781611976137.1>.
- [7] J. BENNETT, A. CARBERY, M. CHRIST, AND T. TAO, *Finite bounds for Hölder-Brascamp-Lieb multilinear inequalities*, Mathematical Research Letters, 17 (2010), pp. 647–666, <https://doi.org/10.4310/MRL.2010.v17.n4.a6>.
- [8] R. BRO AND C. A. ANDERSSON, *Improving the speed of multiway algorithms: Part II: Compression*, Chemometrics and Intelligent Laboratory Systems, 42 (1998), pp. 105–113, [https://doi.org/10.1016/S0169-7439\(98\)00011-2](https://doi.org/10.1016/S0169-7439(98)00011-2).
- [9] V. T. CHAKRAVARTHY, J. W. CHOI, D. J. JOSEPH, X. LIU, P. MURALI, Y. SABHARWAL, AND D. SREEDHAR, *On optimizing distributed Tucker decomposition for dense tensors*, in 2017 IEEE International Parallel and Distributed Processing Symposium (IPDPS), May 2017, pp. 1038–1047, <https://doi.org/10.1109/IPDPS.2017.86>.
- [10] E. CHAN, M. HEIMLICH, A. PURKAYASTHA, AND R. VAN DE GEIJN, *Collective communication: theory, practice, and experience*, Concurrency and Computation: Practice and Experience, 19 (2007), pp. 1749–1783, <https://doi.org/10.1002/cpe.1206>.
- [11] J. CHOI, X. LIU, AND V. CHAKRAVARTHY, *High-performance dense Tucker decomposition on GPU clusters*, in Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis, SC ’18, Piscataway, NJ, USA, 2018, IEEE Press, pp. 42:1–42:11, <http://dl.acm.org/citation.cfm?id=3291656.3291712>.



- [12] M. CHRIST, J. DEMMEL, N. KNIGHT, T. SCANLON, AND K. A. YELICK, *Communication lower bounds and optimal algorithms for programs that reference arrays - part 1*, Tech. Report UCB/EECS-2013-61, EECS Department, University of California, Berkeley, May 2013, <http://www2.eecs.berkeley.edu/Pubs/TechRpts/2013/EECS-2013-61.html>.
- [13] J. DEMMEL, D. ELIAHU, A. FOX, S. KAMIL, B. LIPSHITZ, O. SCHWARTZ, AND O. SPILLINGER, *Communication-optimal parallel recursive rectangular matrix multiplication*, in 2013 IEEE 27th International Symposium on Parallel and Distributed Processing, 2013, pp. 261–272, <https://doi.org/10.1109/IPDPS.2013.80>.
- [14] J.-W. HONG AND H. T. KUNG, *I/O complexity: The red-blue pebble game*, in Proceedings of the Thirteenth Annual ACM Symposium on Theory of Computing, STOC '81, New York, NY, USA, 1981, Association for Computing Machinery, pp. 326–333, <https://doi.org/10.1145/800076.802486>.
- [15] D. IRONY, S. TOLEDO, AND A. TISKIN, *Communication lower bounds for distributed-memory matrix multiplication*, Journal of Parallel and Distributed Computing, 64 (2004), pp. 1017–1026, <https://doi.org/10.1016/j.jpdc.2004.03.021>.
- [16] N. KNIGHT, *Communication-Optimal Loop Nests*, PhD thesis, EECS Department, University of California Berkeley, Aug 2015, <http://www2.eecs.berkeley.edu/Pubs/TechRpts/2015/EECS-2015-185.html>.
- [17] T. G. KOLDA AND B. W. BADER, *Tensor decompositions and applications*, SIAM Review, 51 (2009), pp. 455–500, <https://doi.org/10.1137/07070111X>.
- [18] H. KOLLA, K. ADITYA, AND J. H. CHEN, *Higher Order Tensors for DNS Data Analysis and Compression*, Springer International Publishing, Cham, 2020, ch. 6, pp. 109–134, [https://doi.org/10.1007/978-3-030-44718-2\\_6](https://doi.org/10.1007/978-3-030-44718-2_6).
- [19] L. D. LATHAUWER, B. D. MOOR, AND J. VANDEWALLE, *A multilinear singular value decomposition*, SIAM Journal on Matrix Analysis and Applications, 21 (2000), pp. 1253–1278, <https://doi.org/10.1137/S0895479896305696>.
- [20] L. H. LOOMIS AND H. WHITNEY, *An inequality related to the isoperimetric inequality*, Bulletin of the American Mathematical Society, 55 (1949), pp. 961 – 962, <https://doi.org/bams/1183514163>.
- [21] L. MA AND E. SOLOMONIK, *Accelerating alternating least squares for tensor decomposition by pairwise perturbation*, Numerical Linear Algebra with Applications, (2022), pp. 1–33, <https://doi.org/10.1002/nla.2431>.
- [22] R. MINSTER, A. K. SAIBABA, AND M. E. KILMER, *Randomized algorithms for low-rank tensor decompositions in the Tucker format*, SIAM Journal on Mathematics of Data Science, 2 (2020), pp. 189–215, <https://doi.org/10.1137/19M1261043>.
- [23] T. M. SMITH, B. LOWERY, J. LANGOU, AND R. A. VAN DE GEIJN, *A tight I/O lower bound for matrix multiplication*, 2019, <https://arxiv.org/abs/1702.02017>.
- [24] E. SOLOMONIK, J. DEMMEL, AND T. HOEFLER, *Communication lower bounds of bilinear algorithms for symmetric tensor contractions*, SIAM Journal on Scientific Computing, 43 (2021), pp. A3328–A3356, <https://doi.org/10.1137/20M1338599>.
- [25] Y. SUN, Y. GUO, C. LUO, J. TROPP, AND M. UDELL, *Low-rank Tucker approximation of a tensor from streaming data*, SIAM Journal on Mathematics of Data Science, 2 (2020), pp. 1123–1150, <https://doi.org/10.1137/19M1257718>.
- [26] R. THAKUR, R. RABENSEIFNER, AND W. GROPP, *Optimization of collective communication operations in MPICH*, The International Journal of High Performance Computing Applications, 19 (2005), pp. 49–66, <https://doi.org/10.1177/1094342005051521>.
- [27] L. R. TUCKER, *Some mathematical notes on three-mode factor analysis*, Psychometrika, 31 (1966), pp. 279–311, <https://doi.org/10.1007/BF02289464>.
- [28] A. N. ZIOGAS, G. KWASNIEWSKI, T. BEN-NUN, T. SCHNEIDER, AND T. HOEFLER, *Deinsum: Practically I/O optimal multilinear algebra*, tech. report, arXiv, 2022, <https://doi.org/10.48550/ARXIV.2206.08301>.