

GENERAL CONSTRUCTION AND ANALYSIS OF PARALLEL MATRIX MULTIPLICATION ALGORITHMS

WENJING CUN, LIHUA PEI, AND YUEFAN DENG

ABSTRACT. This paper proposes a general vision onto the construction of the algorithms of parallel matrix multiplication in distributed-memory system. Via the analysis of the communication forest of the algorithms, this paper presents a map of the parallel algorithms for universal matrix multiplication, based on 3-dimensional hypercube algorithm with computational complexity of $O(N^3)$, and presents the experiments results of several algorithms on clusters to compare the performance.

1. INTRODUCTION

High-performance algorithms of linear operations are critical for many of the computing tasks, and to promote the performance of computing in computers, the topology-based algorithms of operations on specific fields are more valuable for research, in order to promote the flexibility and robustness more regardless of the values involved.

And as one of the most essential linear operations, matrix multiplication is a good start point to launch our research on abstract operation. This paper will mainly discuss about the general method of constructing algorithms for matrix multiplication, and the parallelization of it, with several classical methods as well as analysis of their performances for instance. In this research, we have discovered a general map to organize the family of the algorithms, and with the metric of parallel efficiency of the lattice, it is now possible for mathematicians to evaluate the upper bound of the algorithms depending on any given dimensions of the matrices, Furtherly, this paper proposes a formulated problem to reduce the computational complexity of the algorithm of matrix multiplication, which is still not resolved, yet we will figure it as a future work to research on, to lead the way to optimize the distributed algorithms with computational complexity lower than $O(N^3)$.

The organization of this paper is as the following: Section 2 will present the general idea of designing the algorithms of matrix multiplication on 3D hypercube, based the maps between sets of coordinates. In Section 3, we will propose the problem of constructing the communication forest among the computation nodes, while proceeding the parallel algorithms, and organize the family of the parallel algorithms based on 3D hypercube. Then Section 4 will present the analysis of several classical parallel algorithms for matrix multiplication and propose a general algorithm for constructing the distributed plan for universal matrix multiplication.

Date: March 29, 2018.

Key words and phrases. Parallel Computation, Matrix Theory.

The first author was supported in part by NSF Grant #000000.

Support information for the second author.

Algorithm 1 (Naive Algorithm)**Input:** $A \in \mathbb{R}^{m \times l}, B \in \mathbb{R}^{l \times n}$ **Output:** $C \in \mathbb{R}^{m \times n}$ **for** $i = 1 : m$ **do** **for** $i = 1 : l$ **do** **for** $i = 1 : n$ **do** $C_{ik} = C_{ik} + A_{ij} \cdot B_{jk}$ **Return** C

Finally this paper will propose a problem of reducing the computational complexity, and summarize the construction of algorithms.

2. CONSTRUCTION OF MATRIX MULTIPLICATION IN 3D HYPERCUBE

Firstly, we shall review the naïve algorithm of matrix multiplication.

Naïve algorithm is the direct definition of matrix multiplications for any given pair of matrices: $A \in R^{(m \times l)}, B \in R^{(l \times n)}$, where m, l and n are some positive integers. And if we just focus on one final result in $C = A \cdot B$, namely C_{ij} , then we have:

$$(2.1) \quad C_{ik} = \sum_{j=1}^l (A_{ij} \cdot B_{jk})$$

Then regardless of the actual values of A_{ij} and B_{jk} , formula (1) reveals the information that C_{ik} only results from the values in the coordinates depending on three indices i, j and k , namely (i, j, k) . Then group all (i, j, k) corresponding to the coordinates in A and B together, we may expect there is a well-defined map from such a set of (i, j, k) to that of (i, k) in C .

Sourced from such an idea, and in order to make the algorithms totally independent of the values involved, so that can be flexible and robust, we will define such a map from a set of coordinates to another one in the first subsection.

2.1. Coordinate Set and Operation Map. Maps between sets of values with associate coordinates are often independent of the values involved, but more like the relationship among the indices of coordinates. So to abstract such relationships, we proposed the concept called operation map, of which the definition is as the following:

Definition 2.1. An operation map is a map from one set A to another set B , where there are two injective maps $L_1 : C \rightarrow A, L_2 : C \rightarrow B$, and C is a set with a well-defined operation function τ .

For convenience, we denote such an operation map f as: $f_{C,\tau} : A \rightarrow B$, and name L_1, L_2 as the coordinate maps. Here A and B are the topologies from which the sets of values in C abstracted into, and sourcing from the most intuitive concept in geometry, we call A and B the corresponding coordinate sets. And an abstract operation map has a property sourced from the abstraction of the value set to coordinate sets:

Theorem 2.2. Given an operation map $f_{C,\tau} : C_1 \rightarrow C_2$, if S is closed under a well-defined operation τ , or say G is an algebraic group, then for any subset $S \subseteq G$, $f_{C,\tau} : C_1 \rightarrow C_2$ is also an operation map. Vice Versa.

Such a property shows the stability of an operation map regardless of the values when the inputs in conservative of the algebraic operation.

For a series of operation maps, they will follow the chain rule if they follow Theorem 2.1:

Theorem 2.3. *Given a set G that is closed under a series of operations τ_1, \dots, τ_n , a series of sets C_1, \dots, C_n , also a series of subsets of G , S_1, \dots, S_n , there exist at least n operation maps:*

$$(2.2) \quad (f_1)_{S_1, \tau_1}, \dots, (f_n)_{S_n, \tau_n}$$

And $(f_1)_{S_1, \tau_1}, \dots, (f_n)_{S_n, \tau_n}$ is also an operation map.

Following the chain rule in the opposite direction, if an abstract operation map $f(G, \tau) : C_1 \rightarrow C_2$ can be decomposed into a series of abstract operation maps, then we say $f(G, \tau)$ is separable.

The third property we raise here shows one additional great value of abstract operations, that is the flexibility, aka the scalability in computational engineering, if the operation map can be conservative in topology:

Theorem 2.4. *Given two homomorphic algebraic structures G_1 and G_2 , closed under two operations τ_1 and τ_2 respectively, if $S_1 \subseteq G_1$ and $S_2 \subseteq G_2$ are two subsets, and an operation map $f(S_1, \tau_1) : C_1 \rightarrow C_2$, then $f(S_2, \tau_2) : C_1 \rightarrow C_2$ is also an operation map.*

In the next subsection, we will present the ideas of abstracting the process of matrix multiplication, based on the naive algorithm 2.1, into a separable operation map composed by several independent stages, and shows the methodology to parallelize the algorithm.

2.2. Algorithms Based on 3-Dimensional Hypercube. To abstract the naïve algorithm 2.1 into the operation maps into operations maps mentioned in the previous subsection, the first step is to split all the multiplications between a pair of scalars out, then for the 3D-hypercube algorithm, this step means an abstract operation map from a coordinate set C_{11} to the other one C_{12} . Since the original coordinate set represents the positions of all the entries of A and B , then we denote the set of the coordinate as:

$$(2.3) \quad C_A = \mathbb{Z}/m\mathbb{Z} \times \mathbb{Z}/l\mathbb{Z}, C_B = \mathbb{Z}/l\mathbb{Z} \times \mathbb{Z}/n\mathbb{Z}$$

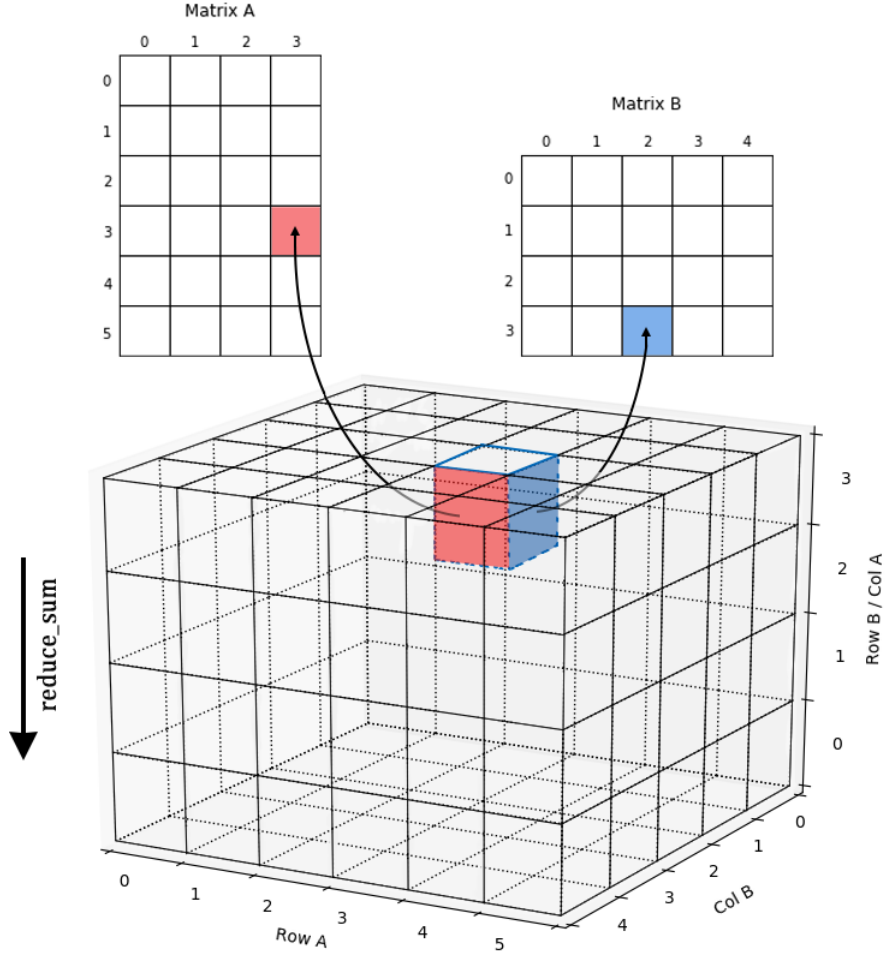
And define the following abstract operation map:

$$(2.4) \quad (f_1)_{\mathbb{R}, \times} : C_A \times C_B \rightarrow C_{mult3D}$$

where $C_{mult} = \mathbb{Z}/m\mathbb{Z} \times \mathbb{Z}/l\mathbb{Z} \times \mathbb{Z}/n\mathbb{Z}$, " \times " is the scalar multiplication and:

$$(2.5) \quad (f_1)_{\mathbb{R}, \times}((r_a, c_a), (r_b, c_b)) = \begin{cases} (r_a, c_b, c_a), & \text{if } c_a = r_b \\ \emptyset, & \text{otherwise} \end{cases}$$

This abstract operation map represents the process: $A_{ij} \cdot B_{jk}$, and from C_{mult} , it can be observed that all the mapping can be described as the following figure:

FIGURE 1. Operation map $(f_1)_{(\mathbb{R}, \times)} : C_A \times C_B \rightarrow C_{\text{mult}}$

The next process to be considered is the addition part, here since it is known that only the scalar operation is remained, and the result is a matrix $C \in \mathbb{R}^{m \times n}$, thus we define the third coordinate set as:

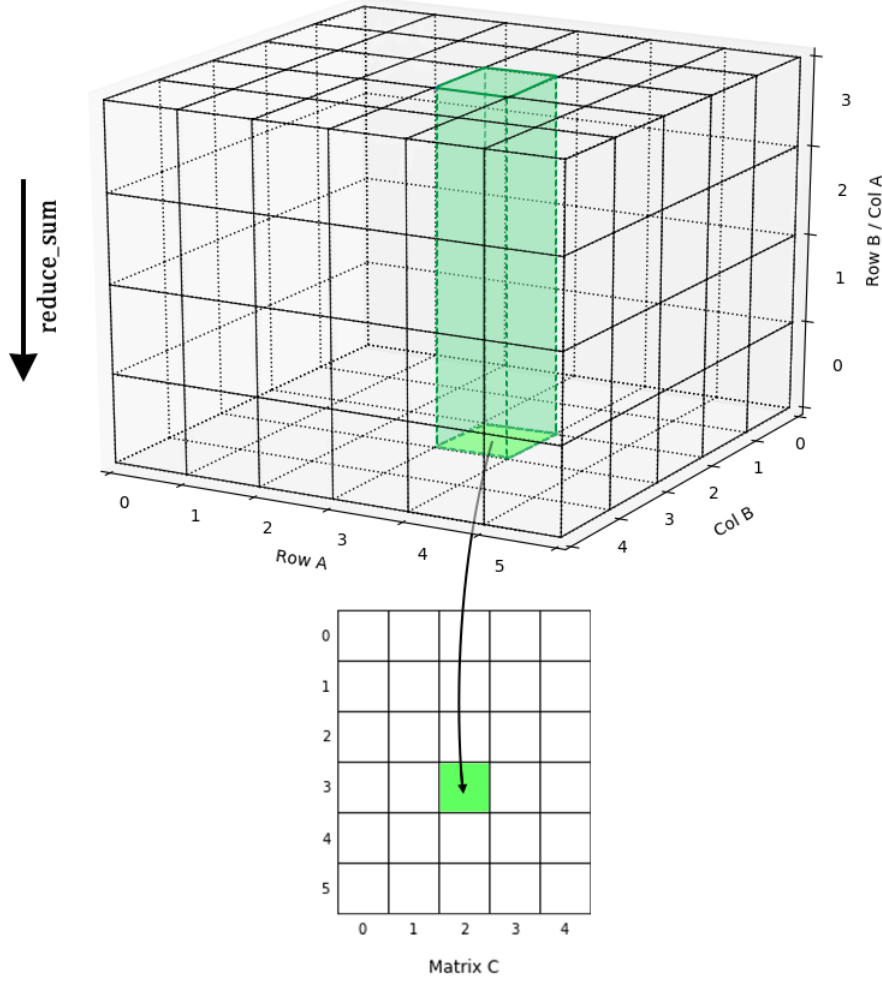
$$(2.6) \quad C_C = \mathbb{Z}/m\mathbb{Z} \times \mathbb{Z}/n\mathbb{Z}$$

And the corresponding abstract operation map is defined as:

$$(2.7) \quad \begin{aligned} (f_2)_{R,+} &: C_{\text{mult}3D} \rightarrow C_C \\ (f_2)_{R,+}((r_a, c_b, c_a)) &= (r_a, c_b) \end{aligned}$$

Here "+" is the scalar addition.

Observed from formula(2.7), all the 3D "cubes" in the Figure 2.1 with the same "z-value" are reduced in to one in the "xy-plane", which can be plotted as the following:


 FIGURE 2. Operation map $(f_2)_{F^{m \times l \times n}, \text{sum}}$

Now the whole process has been abstracted into a chain of abstract operation maps, for the matrix multiplication operation: $F^{m \times l} \times F^{l \times n} \rightarrow F^{m \times n}$, we can define the chain of abstract operation maps as the following:

$$(2.8) \quad f_{\mathbb{R}^{m \times l} \times \mathbb{R}^{l \times n}, \text{matrix multiply}} = (f_1)_{\mathbb{R}, \times} \circ (f_2)_{\mathbb{R}, +}$$

where we define $f_{\mathbb{R}^{m \times l} \times \mathbb{R}^{l \times n}, \text{matrix multiply}}(x) = x$, is a trivial abstract operation map directly representing that: $C = AB$.

For a local-memory system, there seems nothing to modify (2.7) for such a fixed chain, but to a distributed-memory cluster, we can add one more abstract operation map to each computing node p , $g_{R,=}$, where "=" simply means the identity operation, but can filter the coordinates for each node, and we define it as:

$$(2.9) \quad \begin{aligned} g_{R,=} &: C_{mult3D} \rightarrow C_{mult3D} \\ g_{R,=}(r_a, c_b, c_a) &= \begin{cases} (r_a, c_b, c_a), & \text{if } (r_a, c_b, c_a) \in U_p \\ \emptyset, & \text{otherwise} \end{cases} \end{aligned}$$

Here U_p is a custom set chosen for each node p , called a filtering set, represents the “blocks” chosen for node p to calculate. Thus, the chain of abstract operation maps for each node p is written as:

$$(2.10) \quad (f_p)_{\mathbb{R}^{m \times l} \times \mathbb{R}^l \times n, \text{matrix multiply}} = (f_1)_{\text{mathbb{R} \times \mathbb{R}, \times} \circ g_{\text{mathbb{R} \times \mathbb{R}, =} \circ (f_2)_{\text{mathbb{R} \times \mathbb{R}, +}$$

This formula represents all the value-independent parallel algorithms based on Naïve Algorithm.

Now we take such a case for example: if there are totally 4 nodes for computing $A \cdot B$, where $A \in \mathbb{R}^{6 \times 4}$, $B \in \mathbb{R}^{4 \times 6}$, then choose the four filtering sets as:

$$(2.11) \quad \begin{aligned} U_1 &= \{4, 5\} \times \{2, 3, 4\} \times (\mathbb{Z}/4\mathbb{Z}) \\ U_2 &= \{3\} \times \{2, 3, 4\} \times (\mathbb{Z}/4\mathbb{Z}) \\ U_3 &= \{3, 4, 5\} \times \{0, 1\} \times (\mathbb{Z}/4\mathbb{Z}) \\ U_4 &= \{0, 1, 2\} \times \{0, 1, 2, 3, 4\} \times (\mathbb{Z}/4\mathbb{Z}) \end{aligned}$$

So the “blocks” in C_{mult} allotted to each node are shown as the following figure:

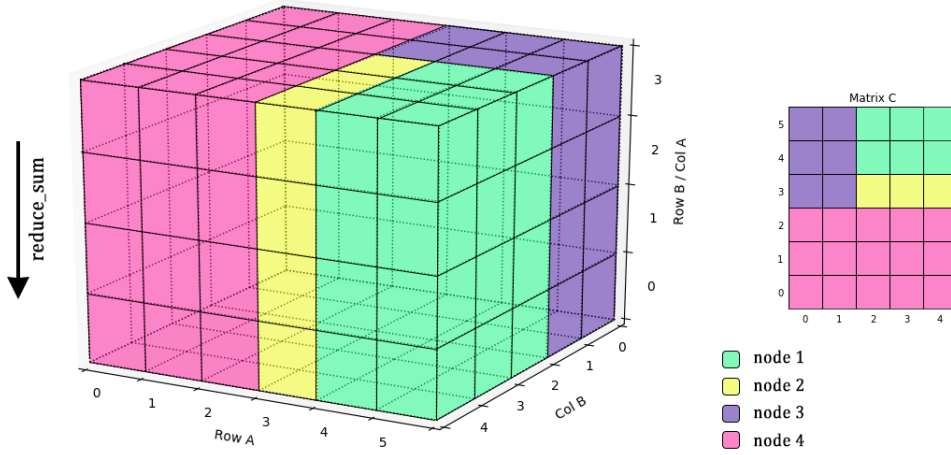


FIGURE 3. Parallel 3D-Hypercube

In fact, the strategy to choose filtering sets as single “cubes” that all span over z-axis, is similar to one of the popular general algorithm for parallel universal matrix multiplication, called BMR[1], if not consider the communication among the nodes. Specially, if the blocks are not finished computing at one communicating step, it can deduce many interesting algorithms, where Cannon’s algorithm[2] and SUMMA[3] are two of the famous ones.

So only consider the computational part of the whole process of matrix multiplication, all the algorithms based on Naïve Algorithm can be deduced by selecting the filtering sets $\{U_p\}$, so that the algorithm is presented as the following pseudocode:

Algorithm 2 (3D-Hyperblock Matrix Multiplication)**Input:** $A \in \mathbb{R}^{m \times l}, B \in \mathbb{R}^{l \times n}$ **Output:** $C \in \mathbb{R}^{m \times n}$ **if** at rank p **do** choose $U_p = \{(i_1, j_1, k_1), (i_2, j_2, k_2), \dots, (i_K, j_K, k_K)\}$ **for** i, j, k in U_p **do** $C_{ik} = C_{ik} + A_{ij} \cdot B_{jk}$ **Return** C

The members in the 3D-Hyperblock family often have the best potential in minimizing the time in communication, however, while dealing with large-scale matrices, the buffer size often becomes a problem for the cluster to limit the total size of data to be deployed onto each computation node. Therefore, each such optimization in communication involves deploying the best subset $V_p \subseteq U_p$, and find the best communication tree to pass the data among the nodes, in order to finish the tasks undertaken by U_p efficiently, which is sophisticated while dealing with such P cubes simultaneously.

In the next section, we will discuss the influence of limiting buffer size on designing the communication part for a parallel algorithm, and find out a general map for designing the algorithm for parallel matrix multiplication.

3. COMMUNICATION IN 3D-HYPERCUBE

Allocation of blocks of operations to each computational node is straightforward for 3D hypercube algorithms, however, while dealing with large-scale matrix multiplications, the restricted buffer size will limit the size of data stored onside each node.

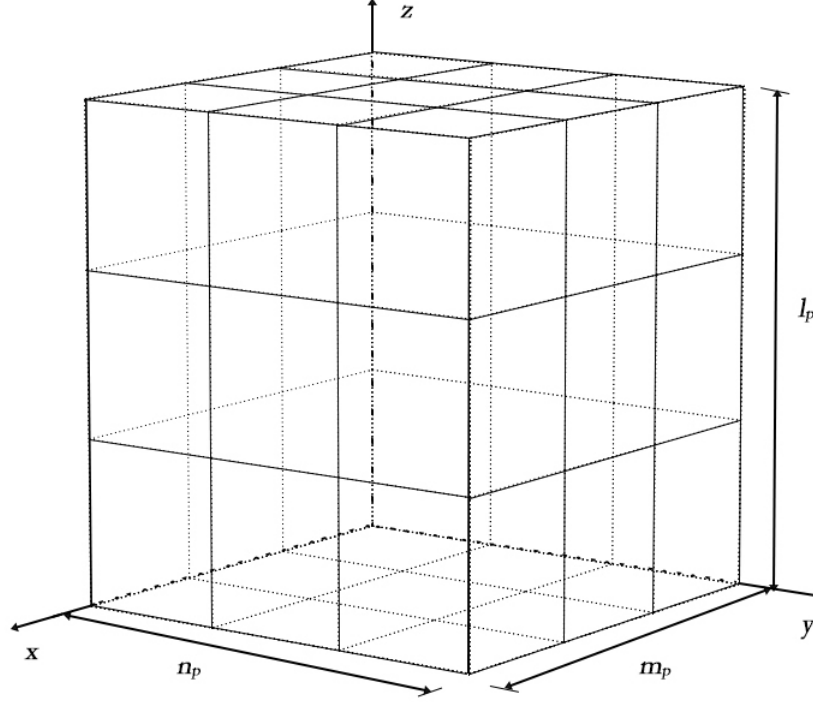
Generally, when a parallel program allows all the nodes to own all the necessary data already, in terms of the algorithm itself, the cost of communication can be minimized. Otherwise, one node lacking necessary data needs to fetch those needed from the other ones. So for parallel matrix multiplication, referring to algorithm 2.2, often we can only deploy part of U_p onto node p , then following the communication tree, at each step, one need to finish as much as computation tasks with its data owing, release or send the data no longer needed, and require and receive the data needed for next step of computation.

3.1. Communication Rules and Cost. For convenience, we describe a “cube” assigned to node p with:

- m_p, l_p, n_p , the lengths of the cube assigned to node p , along with x-axis, y-axis and z-axis respectively;

For many of the algorithms in 3D-hypercube family, we may assume that such P communication trees are similar, or out of the geometric vision, such P “cubes” are similar in the following properties:

- Volume, i.e. the number of multiplication between entries;
- Total projection area onto xz-plane and yz-plane, i.e. the total size of data needed for the computation tasks;
- Maximum area projective onto xz-plane and yz-plane can be covered at each communication step, i.e. the maximum buffer size assigned to each core, denoted as Buffer_p .

FIGURE 4. Cube Onside Node p

The three properties above in fact determined the performance of the part of scalar (or even block) multiplication. As for the part of scalar (or even block) addition, there is only one additional property needed to be considered:

- Total projection area onto xy-plane, i.e. the size of data to be reduced onto one node.

Then to design a smart algorithm, we may design following the rules below:

Rule 3.1 The total volume of the P cubes is exactly mln , i.e. the total volume of the whole 3D-hypercube. i.e. there would be no repetitive computational tasks among the nodes:

$$(3.1) \quad \sum_{p=1}^P m_p l_p n_p = mln$$

Rule 3.2 The entries of two matrices can be covered by the data stored in the P nodes at each communication step, so that it requires:

$$(3.2) \quad \sum_{p=1}^P \text{Buffer}_p \geq ml + ln$$

Rule 3.3 Any entry (or block) stored onside one node can only be released if and only if all the scalar (block) multiplications involving it have been finished. So that

the total cost of receiving data on node p is:

$$(3.3) \quad T_{p,recv} = C \cdot (m_p l_p + l_p n_p - \text{Buffer}_p)$$

From Rule 3.3, since the total number of sending is always conservative with that of receiving, so that we may estimate the average cost of communication on send-receive for each node would follow:

$$(3.4) \quad \begin{cases} T_{comm,send-recv} \leq \sum_{p=1}^P T_{p,recv} & \text{send/require one by one} \\ T_{comm,send-recv} \geq \max_p \{T_{p,recv}\} & \text{communication parallelized} \end{cases}$$

where the lower bound comes from the strategy that minimizing the waiting time, where we let the one with the largest demand to send nothing out.

Next, based on the analysis of cost on send-recv, we are going to discuss the strategy of designing the cubes and communication tree for each node.

3.2. Algorithm for Parallel Matrix Multiplication. In practical cases, since where all the computation nodes are in the similar conditions (in CPU, GPU, RAM, etc.), we often divide the whole 3D hypercube into P similar (not necessary to be identical) cubes, and assigned with similar limits of buffer sizes.

For each cube assigned to each node, due to the restriction to buffer size, our communication objective is to minimizing the time cost for communication among the nodes.

So numerate $U_p = \{(x_{pi}, y_{pi}, z_{pi})\}$ by taking the operation map (2.9) mentioned in subsection 2.1, which is source from $A_{pk} \times B_{pk}$ at each communication step k , where it should satisfy that:

$$(3.5) \quad |C_{pAk}| + |C_{pBk}| \leq \text{Buffer}_p$$

So that we can see, to parallelize the matrix multiplication by adding the communication part in 3D-hypercube, the form of the operation maps can be totally conserved, and so construct a union of a series of operation maps in the same forms:

$$(3.6) \quad \begin{aligned} & (f_{p1})_{\mathbb{R}, \times} : C_{pA} \times C_{pB} \rightarrow C_{p,mult3D} \\ \Rightarrow & C_{p,mult3D} = \bigcup_{k=1}^{\text{total steps } K} C_{p,mult3D,k} = \bigcup_{k=1}^K C_{pAk} \times C_{pBk} \end{aligned}$$

Now to minimize the communication cost on each node with fixed limit of buffer size, we want to prove that:

Theorem 3.1. *To numerate $U_p = \{(x_{pi}, y_{pi}, z_{pi})\}$, for $\{C_{pAk}\}$ and $\{C_{pBk}\}$ satisfying (3.5), there can exist $\{C_{pAk}\}$ and $\{C_{pBk}\}$ so that:*

- For $\Delta k > 1, d \in C_{pAk} \cap C_{pA,(k+\Delta k)}$ only if $d \in C_{pAk} \cap C_{pA,(k+\Delta k-1)}$.
- For $\Delta k > 1, d \in C_{pBk} \cap C_{pB,(k+\Delta k)}$ only if $d \in C_{pBk} \cap C_{pB,(k+\Delta k-1)}$.

if given that $\text{Buffer}_p \geq \min\{m_p, n_p\}$.

This theorem states the availability of Rule 3.3, so that no data (entry or block) would be sent to any one node redundantly. And the given condition can propose a general strategy for send-receive and so prove the existence.

This strategy with the condition $\text{Buffer}_p \geq \min\{m_p, n_p\}$ can be stated as the following:

Algorithm 3 (3D-Hyperblock Algorithm with Communication)**Input:** $A \in \mathbb{R}^{m \times l}$, $B \in \mathbb{R}^{l \times n}$, m_p, l_p, n_p , $\text{Buffer}_p \in N^+$, ($m_p \leq n_p$)**Output:** $C \in \mathbb{R}^{m \times n}$ **if** at rank p **do**

Choose $U_p = \{(i_1, j_1, k_1), (i_2, j_2, k_2), \dots, (i_K, j_K, k_K)\}$
 $\subset (\mathbb{Z}_{m_p} \times \mathbb{Z}_{l_p} \times \mathbb{Z}_{n_p}) \oplus (i_{row}, j_{row}, k_{row})$

Let $C_{pA} = (Z(m_p) \times Z(l_p)) \oplus (i_{low}, j_{low})$ and

$C_{pB} = (Z(l_p) \times Z(n_p)) \oplus (j_{low}, k_{low})$.

Sort U_p, C_{pA} and C_{pB} increasingly with j .

Let $\text{StepSize}_A = \lfloor \text{Buffer}_p / m_p \rfloor$, and initialize the buffer as:

$C_{pA,0} = C_{pA}[0 : \text{StepSize}_A]$, $C_{pB,0} = C_{pB}[0 : \text{Buffer}_p - \text{StepSize}_A]$

Set numbers of steps associated to two sets as $k_A = k_B = 0$.

for i, j, k in U_p ,

if $(i, j) \in C_{pA, k_A}$ and $(j, k) \in C_{pB, k_B}$ **do**

$C_{ik} = C_{ik} + A_{ij} \cdot B_{jk}$

ready to send $(i, j) \in C_{pA, k_A}$ and $(j, k) \in C_{pB, k_B}$

else if $(i, j) \notin C_{pA, k_A}$ **do**

require (i, j) from another node,

pop the front of C_{pA, k_A} ,

push (i, j) into the back of C_{pA, k_A} ,

$k_A = k_A + 1$.

else if $(j, k) \notin C_{pB, k_B}$ **do**

require (j, k) from another node,

pop the front of C_{pB, k_B} ,

push (j, k) into the back of C_{pB, k_B} ,

$k_B = k_B + 1$.

Gather the results onto the root node.

Return C

This algorithm designs a communication tree for each node, with which we can construct an algorithm with only the following two given conditions:

- m_p, l_p, n_p , i.e. a cube assigned to each node p ;
- Buffer_p .

And the algorithm designed based on algorithm 3.1 is called the **Buffer Adaptive Matrix Multiplication Algorithm (BAMMA)**, which can fit to pair of matrices with random dimensions and make the best usage of buffer onside each node.

Next we will evaluate the performance of BAMMA, based on the variables raised in section 3.1, and show the general prediction of parallel efficiency of the algorithms for parallel matrix multiplication in 3D-Hypercube family.

3.3. Evaluation of Performance. To evaluate the performance, or say the total cost of time of a parallel matrix multiplication algorithm, we need to evaluate both the computational cost and the communication cost, i.e.:

$$(3.7) \quad T = \max_p \{T_{p,comm} + T_{p,comp} + T_{p,idle}\}$$

where the $T_{p,idle}$ is the waiting time on node p .

If we assigned similar cubes to p nodes with similar limits of buffer sizes, then we may assume that there is a way to minimize $T_{p,idle}$ by building static routines for the whole communication forest. So that we may remove the maximum in (3.7) to estimate the total cost:

$$(3.8) \quad T_{total,P} \approx T_{p,comp} + T_{p,idle}$$

And considering the practical cases in computational clusters, we assume that the average time cost on one pair of send-receive process on one floating number (notated as $\hat{t}_{send-recv}$), and that of one arithmetical operation, (notated as \hat{t}_{op}), will satisfy that:

$$(3.9) \quad \gamma := \frac{\hat{t}_{send-recv}}{\hat{t}_{op}}$$

and for practical cases, we assume τ is a constant value associated with particular machine, cluster topology, to be used to estimate the real fraction while running the program.

Now firstly, for the computational cost $T_{p,comp}$, by the operation map defined as (2.10), it sources from two parts: scalar (or block) multiplication and addition, so that we have:

$$(3.10) \quad T_{p,comp} = T_{p,mult} + T_{p,add} = (2 \cdot m_p l_p n_p - m_p n_p) \cdot \hat{t}_{op} \approx 2m_p l_p n_p \cdot \hat{t}_{op}$$

And for those similar p cubes:

$$(3.11) \quad m_p l_p n_p \approx \frac{m l n}{P}$$

where P is the total number of nodes. So in total:

$$(3.12) \quad T_{comp} = T_{p,comp} \approx \frac{2m l n}{P} \cdot \hat{t}_{op}$$

Now consider the next part of communication, it consists of two parts: send-receive in order to complete all the computations onside each cube, and the reduction with summation:

$$(3.13) \quad T_{comm} = T_{send-recv} + T_{reduce}$$

For reduction, since the total cost is just the total number of entries (or blocks) need to be sent to the root node, along the z-axis, or say it is just the total projection area onto the xy-plane from the P nodes. And since for each step, it involves one addition and one send/receive action, thus:

$$(3.14) \quad T_{send-recv} \geq C \cdot \max_p \{m_p l_p + n_p l_p - \text{Buffer}_p\} \cdot \hat{t}_{send-recv}$$

where generally C is a constant, usually equals or larger than 2 for machines assigned with similar conditions, since one node can only do send or receive action at one single step.

So in summary, we add all of the costs of parts respectively together, and it deduces the prediction formula of the total cost of time:

$$(3.15) \quad T_{total,P} \geq \frac{2m l n}{P} \cdot \hat{t}_{op} + (\hat{t}_{send-recv} + \hat{t}_{op}) \cdot \sum_{p=1}^P m_p n_p + \hat{t}_{send-recv} \cdot 2 \cdot \max_p \{m_p l_p + n_p l_p - \text{Buffer}_p\}$$

For the convenience to construct the space to include all the algorithms in 3D-Hypercube family, we define two new variables: **Z-variance** and **Maximum Saturability**, denoted by A_z and \bar{S} respectively, as the following:

$$(3.16) \quad A_z := \frac{\sum_{p=1}^P m_p n_p}{mn}, \bar{S} := 1 - \max_p \left\{ \frac{m_p l_p + n_p l_p - \text{Buffer}_p}{l(m+n)} \right\}$$

here $A_z \geq 1$ and $S \in (0, 1]$.

And by parallelizing the reducing process, the whole process to “compress” the cubes can be optimized into a flipped binary tree, so that it would take much shorter time than that in (3.13):

$$(3.17) \quad A_z := \frac{\sum_{p=1}^P m_p n_p}{mn}, T_{reduce} = (\hat{t}_{send-recv} + \hat{t}_{op}) \cdot \log_2 A_z \cdot mn \cdot \frac{A_z}{P}$$

with considering that the columns involving different cores can reduce parallelly.

And based on the assumptions of similar cubes we made above, we may deduce that:

$$(3.18) \quad \bar{S} := 1 - \frac{\bar{m}_p \bar{l}_p + \bar{n}_p \bar{l}_p - \text{Buffer}_p}{l(m+n)}$$

where $\bar{m}_p, \bar{l}_p, \bar{n}_p$ and Buffer_p are set to be the similar for every node.

Then consider the total cost of the matrix multiplication on a single node is always:

$$(3.19) \quad T_{single} = (2mln - mn) \cdot \hat{t}_{op}$$

So inside the $A_z - \bar{S}$ plane, we can now sketch the parallel efficiency:

$$(3.20) \quad \begin{aligned} \tau(m, l, n, P) &\leq \frac{1}{P} \cdot \frac{(2mln - mn) \cdot \hat{t}_{op}}{\frac{2mln - mn}{P} \cdot \hat{t}_{op} + (\hat{t}_{send-recv} + \hat{t}_{op}) \cdot \log_2 A_z \cdot mn \cdot \frac{A_z}{P} + \hat{t}_{send-recv} \cdot 2 \cdot l(m+n)(1 - \bar{S})} \\ &\approx \frac{mln}{mln + (\frac{1+\gamma}{2} \cdot \log_2 A_z \cdot mn \cdot \frac{A_z}{P} + P \cdot \gamma \cdot l(m+n)(1 - \bar{S}))} \end{aligned}$$

because in most of the cases, $l \gg 1$ and simplify (3.18) we can finally get:

$$(3.21) \quad \tau(m, l, n, P) \leq \frac{1}{1 + \frac{1+\gamma}{2l} \cdot A_z \cdot \log_2 A_z + P \cdot \gamma \cdot (\frac{1}{m} + \frac{1}{n})(1 - \hat{S})}$$

Here to normalize the space, we set $\hat{A}_z = 1 - \frac{\log_2 A_z}{\log_2 P} = 1 - \log_P A_z \in (0, 1]$, since each node can at most own projection area mn onto xy -plane, and for any one column along the z -axis, there are at most P segments belong to different nodes. So that:

$$(3.22) \quad \tau(m, l, n, P) \leq \frac{1}{1 + \frac{1+\gamma}{2l} \cdot \log_2 P \cdot P^{1-\hat{A}_z} \cdot (1 - \hat{A}_z) + P \cdot \gamma \cdot (\frac{1}{m} + \frac{1}{n})(1 - \hat{S})}$$

Now with formula (3.22), we can easily predict the upper bound of the parallel efficiency by given certain $m, l, n, P, \hat{m}_p, \hat{l}_p, \hat{n}_p$ and $\text{hattertextBuffer}_p$, and we can also locate the algorithms inside the space, with the metric of parallel efficiency.

3.4. 3D-Hypercube Family. Out of vision of the algorithm designer, it cannot be directly indicated from formula (3.22) that how to choose $\{m_p, l_p, n_p\}$ for given m, l, n and Buffer_p , but with which we may estimate the performance of a parallel algorithm for matrices with certain dimensions, and so support us to promote the algorithm.

For classical method, such as Cannon's method[2], Fox method[1], and SUMMA[3], they all focus on dividing the hypercube into cubes with special shapes following some certain rules, as shown in the following figure for example:

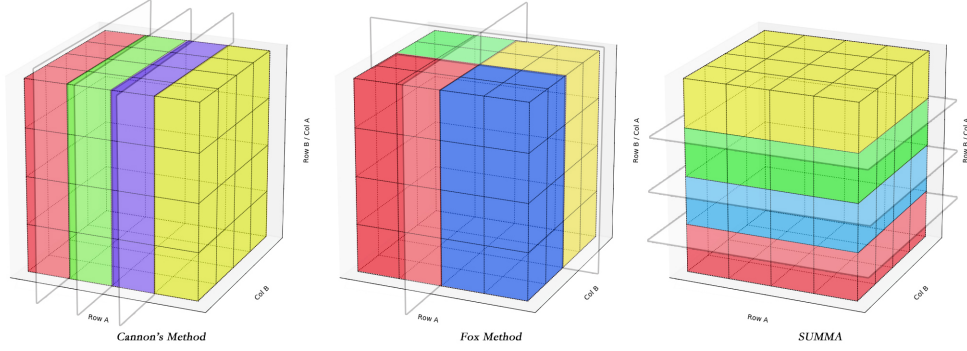


FIGURE 5. Some Classical Methods

But due to different dimensions of matrices involved in the multiplication, one certain strategy to design the shape of cubes can comparatively result in different performances. Here we may plot three different color maps to show the distributions of 3D-Hypercube algorithms for parallel matrix multiplication, with the metric of parallel efficiency $\tau_{nogather}(m, l, n, P)$ defined in (3.22), under 3 different sets of dimensions, namely m, l and n :

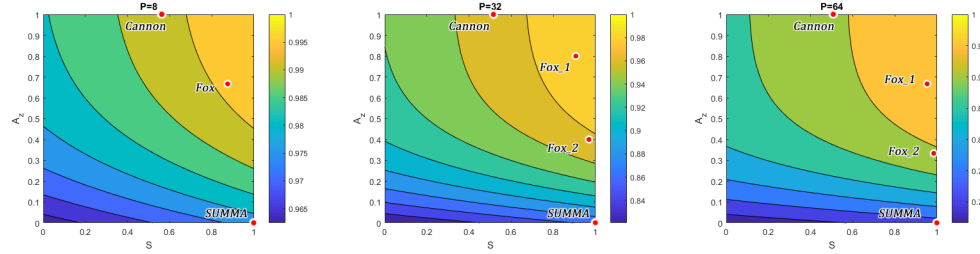
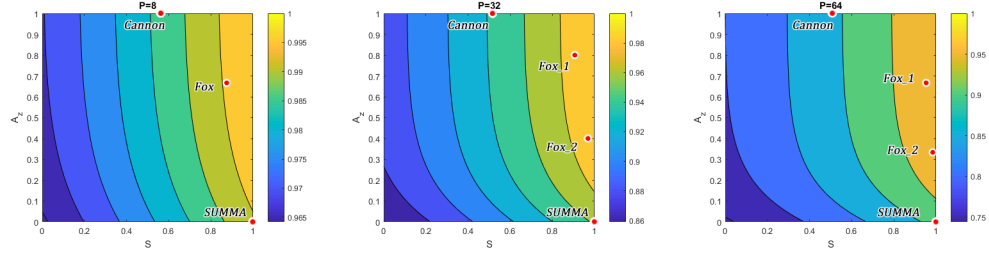
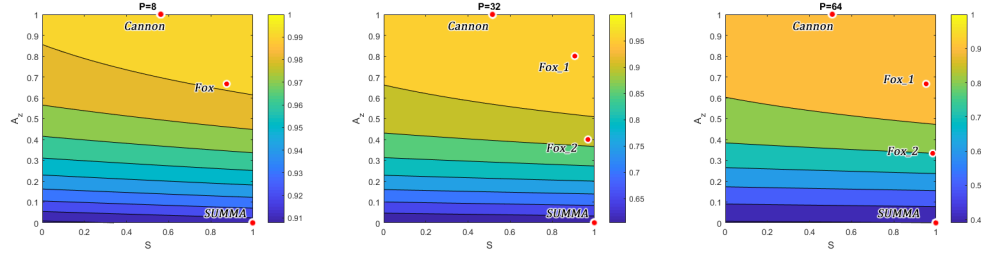


FIGURE 6. Efficiency Distribution: Square Matrices ($m=l=n=1024, \tau=1$)

So from the colormaps above, we may conclude some general ideas for constructing the algorithms for parallel matrix multiplication based on 3D-Hypercube:

- It would be hard to increase \hat{A}_z and \hat{S} with limited buffer size;
- When l is larger than m and n , it would be better to be prior to slice the whole hypercube along z-axis, which means smaller \hat{S} can compromise to smaller \hat{A}_z ;
- When l is smaller than m and n , it would be better to let each cube cover more blocks along the z-axis;

FIGURE 7. Efficiency Distribution: Tall Matrices ($m=n=512, l=4096, \tau=1$)FIGURE 8. Efficiency Distribution: Flat Matrices ($m=n=2048, l=256, \tau=1$)

And we have also raised some observations from the defined variables up to now:

- If a cube on node P is more regular, i.e. m_p, l_p and n_p are closer to each other, then with fixed volume $m_p l_p n_p$, $m_p l_p + n_p l_p$ can be smaller so that it can be more possible to increase \hat{S} ;
- The larger the limit of buffer size Buffer_p is, the more that both \hat{A}_z and \hat{S} can increase;
- For matrices with larger scale, the proportion of computations would be larger, so that the distribution of such colormap would be closer to 1.

Based on those strategies, though we cannot reduce the total cost of computation, we can estimate and help to search out the best 3D-Hypercube algorithm for parallel matrix multiplication, adaptive to certain matrices dimensions, limited buffer size and some other specific conditions.

4. APPLICATION OF 3D-HYPERCUBE ALGORITHM

In this section, we designed a series of experiments to verify our estimations of parallel efficiency using formula (3.22).

From formula (3.22), we want to find out the relationship between the parallel efficiency of a constructed 3D-hypercube algorithm and a few variables, namely the dimensions of matrices, buffer size, and the number of nodes.

We designed the experiments consist of two parts: one for **estimating** τ and the other one for **measuring time cost for particular algorithms**, all the experiments were tested in the SeaWulf Clusters¹, and run with the same one testing program, which is developed with MPI for C/C++ and is capable of parallel communication, and parallel matrix multiplication of complexity $O(N^3)$.

¹Seawulf Clusters: <https://it.stonybrook.edu/help/kb/understanding-seawulf>

4.1. Estimating the Value of τ . As mentioned in Subsection 3.3 and referring to definition (18), τ is presumed as a constant value, associated with particular clusters with P cores. We ran the testing program to test two jobs: one for multiplying two matrices, both of dimensions 1024×1024 , and the other ones for broadcasting an 1024×1024 matrix from each of P cores one by one, so that to measuring the average costs, on unit send-receive communication, and floating operation.

For taking the average, we have run two jobs both for 10 times, in 8, 32 and 64 cores assigned by SeaWulf respectively. And the results are shown as below:

P		Individual Multiplication	Broadcast
8	Mean	5.2372	0.1120
	StdDev	0.0014	0.0014
32	Mean	5.2732	2.0944
	StdDev	0.0049	0.0457
64	Mean	5.2729	8.4876
	StdDev	0.0059	0.1488

So with these data, we can calculate γ respective to each P following the deduction below:

$$(4.1) \quad \frac{T_{bcast}}{T_{mult}} = \frac{P(P-1) \cdot 1024^2 \cdot \hat{t}_{sendrecv}}{(C_{program} \cdot 1024^3 - 1024^2) \cdot \hat{t}_{op}} = \frac{P(P-1)}{C_{program} \cdot 1024 - 1} \cdot \gamma$$

Here $C_{program}$ means the account of floating operations executed per step in a practical program for matrix multiplication. However, in practical cases, it is tough to measure $C_{program}$ stably and precisely, so compromising to experimental cases, formula (31) is modified as the following:

$$(4.2) \quad \begin{aligned} \tau(m, l, n, P) &= \frac{1}{P} \cdot \frac{m \ln \cdot T_{mult_per_step}}{\frac{m \ln}{P} \cdot T_{mult_per_step} + (\hat{t}_{sendrecv} + \hat{t}_{op}) \cdot \log_2 A_z \cdot mn \cdot \frac{A_z}{P} + \hat{t}_{sendrecv} \cdot 2 \cdot l} \\ &= \frac{1}{1 + \frac{\hat{t}_{sendrecv} + \hat{t}_{op}}{T_{mult_per_step}} \cdot \frac{\log_2 A_z}{l} \cdot \hat{A}_z + \frac{2 \cdot \hat{t}_{sendrecv}}{T_{mult_per_step}} \cdot P \cdot (\frac{1}{m} + \frac{1}{n})(1 - \bar{S})} \end{aligned}$$

Here $T_{mult_per_step}$ is the time cost per step in the loop of a matrix multiplication in programming.

Then comparing to formula (32), it is reasonable to obtain the following relationship:

$$(4.3) \quad \gamma \approx \frac{2 \cdot \hat{t}_{sendrecv}}{T_{mult_per_step}} = 2 \cdot \frac{T_{bcast}/P(P-1) \cdot 1024^2}{T_{mult}/(1024^3)} = \frac{2048}{P(P-1)} \cdot \frac{T_{bcast}}{T_{mult}}$$

So the estimated value of each γ is calculated as the following:

P	8	32	64
γ	0.7822	0.8195	0.8170

It has to be admitted that there values are less than the real ones due to those not considered, and the communication performance is usually more complicated

in parallel program, than that analyzed in theory so that hard to obtain the correct γ .

4.2. Experiments of Particular Algorithms. For each cases with 8, 32 and 64 cores respectively, we designed 7, 10 and 10 different strategies to split the 3D hypercube and allocate the buffers in different way, so that constructed various algorithms for testing the performances of each.

To evaluate the performance, the formula for calculating the parallel efficiency is applied here:

$$(4.4) \quad \tau = \frac{1}{P} \cdot \frac{T_{one_core}}{T_{P_Cores}}$$

The plans with results are shown as the following subsections:

4.3.1 $P = 8$ (total 7 cases)

(1) Cube Dimension $(m_p, l_p, n_p) = (128, 1024, 1024)$, $\bar{A}_z = 1$:

Buffer size on A	128 * 1024	128 * 1024	128 * 1024
Buffer size on B	128 * 1024	256 * 1024	1024 * 1024
\bar{S}	0.5625	0.625	1
Efficiency	0.9712	0.9787	0.9899
Expected (sec)	0.9795	0.9947	0.9954
Error (%)	0.8346	1.5963	0.5575

(2) Cube Dimension $(m_p, l_p, n_p) = (512, 512, 512)$, $\bar{A}_z = 0.6667$:

Buffer size on A	512 * 256	512 * 256	512 * 512
Buffer size on B	512 * 256	512 * 512	512 * 512
\bar{S}	0.875	0.9375	1
Efficiency	1.0007	0.9942	0.9948
Expected (sec)	1.0000	0.9967	0.9975
Error (%)	0.0678	0.2491	0.2662

(3) Cube Dimension $(m_p, l_p, n_p) = (1024, 128, 1024)$, $\bar{A}_z = 0$:

Buffer size on A	1024 * 128
Buffer size on B	1024 * 128
\bar{S}	1
Efficiency	0.9951
Expected (sec)	0.9983
Error (%)	0.3151

4.3.2 $P = 32$ (total 10 cases)

(1) Cube Dimension $(m_p, l_p, n_p) = (32, 1024, 1024)$, $\bar{A}_z = 1$:

Buffer size on A	32 * 1024	32 * 1024	32 * 1024	32 * 1024
Buffer size on B	32 * 1024	256 * 1024	512 * 1024	1024 * 1024
\bar{S}	0.515625	0.625	0.75	1
Efficiency	0.9514	0.9809	0.9880	1.0054
Expected (sec)	0.9761	0.9814	0.9875	1.0000
Error (%)	2.4684	0.0546	0.0447	0.5386

(2) Cube Dimension $(m_p, l_p, n_p) = (256, 512, 256)$, $\bar{A}_z = 0.8$:

Buffer size on A	256 * 128	256 * 256	256 * 512
Buffer size on B	256 * 128	256 * 256	256 * 512
\bar{S}	0.90625	0.9375	1
Efficiency	0.9872	0.9894	0.9902
Expected (sec)	0.9951	0.9982	0.9817
Error (%)	0.7892	0.8836	0.8493

(3) Cube Dimension $(m_p, l_p, n_p) = (512, 128, 512)$, $\bar{A}_z = 0.4$:

Buffer size on A	512 * 64	512 * 128
Buffer size on B	512 * 64	512 * 128
\bar{S}	0.96875	1
Efficiency	0.9729	0.9743
Expected (sec)	0.9777	0.9792
Error (%)	0.4849	0.4980

(4) Cube Dimension $(m_p, l_p, n_p) = (1024, 32, 1024)$, $\bar{A}_z = 0$:

Buffer size on A	1024 * 32
Buffer size on B	1024 * 32
\bar{S}	1
Efficiency	0.8438
Expected (sec)	0.8762
Error (%)	3.2378

4.3.2 $P = 64$ (total 10 cases)

(1) Cube Dimension $(m_p, l_p, n_p) = (32, 1024, 1024)$, $\bar{A}_z = 1$:

Buffer size on A	16 * 1024	16 * 1024	16 * 1024	16 * 1024
Buffer size on B	16 * 1024	64 * 1024	256 * 1024	1024 * 1024
\bar{S}	0.5078125	0.53125	0.625	1
Efficiency	0.9015	0.9141	0.9539	0.9959
Expected (sec)	0.9525	0.9547	0.9634	1.0000
Error (%)	5.1028	4.0593	0.9562	0.4088

(2) Cube Dimension $(m_p, l_p, n_p) = (256, 512, 256)$, $\bar{A}_z = 0.6667$:

Buffer size on A	256 * 64	256 * 128	256 * 256
Buffer size on B	256 * 64	256 * 128	256 * 256
\bar{S}	0.953125	0.96875	1
Efficiency	0.9732	0.9746	0.9748
Expected (sec)	0.9883	0.9899	0.9930
Error (%)	1.5095	1.5263	1.8142

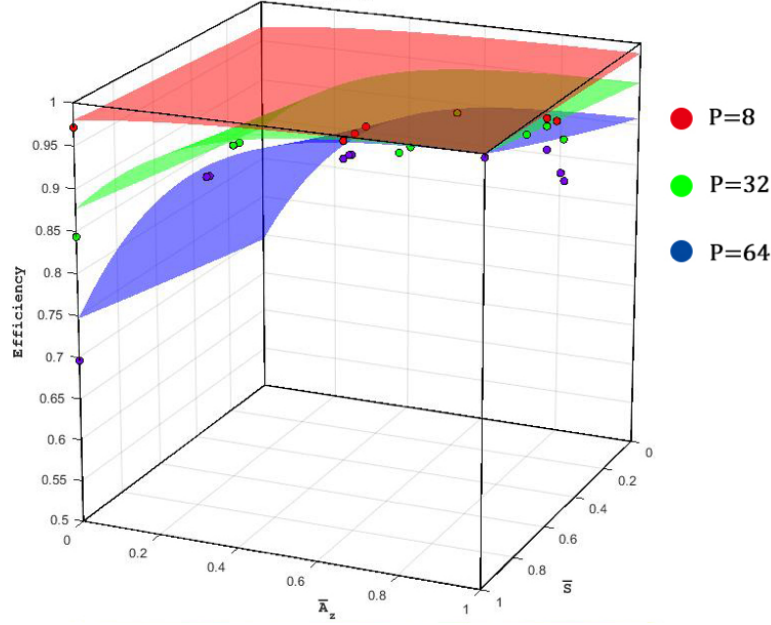
(3) Cube Dimension $(m_p, l_p, n_p) = (512, 128, 512)$, $\bar{A}_z = 0.3333$:

Buffer size on A	512 * 32	512 * 64
Buffer size on B	512 * 32	512 * 64
\bar{S}	0.984375	1
Efficiency	0.9329	0.9342
Expected (sec)	0.9450	0.9464
Error (%)	1.2093	1.2239

(4) Cube Dimension $(m_p, l_p, n_p) = (1024, 32, 1024)$, $\bar{A}_z = 0$:

Buffer size on A	1024 * 16
Buffer size on B	1024 * 16
\bar{S}	1
Efficiency	0.7002
Expected (sec)	0.7465
Error (%)	4.6290

Therefore, finally we can plot the experimental results inside a 3D space together with the expected surfaces sketched from formula (32), which is shown as the following:



Comparing to the analytic surfaces, we found that the communication parts of the algorithms are indeed affected by more other elements not discussed in this paper, which to be searched out is part of the future work on BAMMA.

5. CONCLUSIONS AND FUTURE WORK

In this paper, we have discussed about the construction of algorithms for parallel matrix multiplication, based on 3D-Hypercube so that with total computational complexity $O(N^3)$. We have observed that the construction of such algorithm is influenced by the number of computational nodes, the dimensions of matrices and the limit of the buffer size. Concluded from both the theoretical analysis and the experiments, it is learnt that often it is better to construct the parallel algorithms considering adaptive to certain cases, described by the three variables mentioned above, rather than primly follow the rules defined by classical methods.

All of the theoretical analysis is based on the assumption that all the cubes and communication trees onside the P nodes are designed similarly. So as another part of our future work on parallel matrix multiplication, we are going to dig out the possibility to design and evaluate the algorithm, with more flexibility on each specific node with specific condition respectively, so that it can adapt to more

practical cases with the best performance. As one of the most latest coming work, we are going to optimize the algorithms on clusters with particular or even special topologies, using the ideas raised in this paper.

And as mentioned in Section 3.2, another one of our research plans in the future is to optimize and finish the program of BAMMA, so that it can satisfy the requirement for parallel universal matrix multiplication adaptive to specific conditions, expected with better performance comparable to the classical method. And referring to a popular branch of algorithms of parallel matrix multiplication, the recursive algorithms, for instance CARMA[4] raised by the team led by James Demmel in 2013, we may need to furtherly research on the influence of sequential recursion, on the performance of a parallel algorithm of parallel matrix multiplication.

Finally, as the final target to optimize the algorithm for general matrix multiplications, our group has proposed a formulated problem in 4D Hypercube, with the concepts of operation map mentioned in Section 2.1, to reduce the computational complexity of matrix multiplication, of which is lower than $O(N^3)$, inspired by the ideas[5], raised by Strassen in 1969, as well as Coppersmith and Winograd[6] in 1987. Our next target is to find a general strategy to solve to problem so that find the limit of recursive algorithms to reduce the total complexity of matrix multiplications, and furtherly construct stable and scalable parallel algorithm for universal matrix multiplication, with lower computational complexity.

DEPARTMENT OF APPLIED MATHEMATICS, STONYBROOK UNIV. SUNY, STONY BROOK, NEW YORK 11790

E-mail address: wenjing.cun@stonybrook.edu

DEPARTMENT OF APPLIED MATHEMATICS, STONYBROOK UNIV. SUNY, STONY BROOK, NEW YORK 11790

E-mail address: lihua.pei@stonybrook.edu

DEPARTMENT OF APPLIED MATHEMATICS, STONYBROOK UNIV. SUNY, STONY BROOK, NEW YORK 11790

E-mail address: yuefan.deng@stonybrook.edu