

# Network Topologies and Inevitable Contention

Grey Ballard  
Wake Forest University  
ballard@wfu.edu

Benjamin Lipshitz\*  
UC Berkeley  
lipshitz@cs.berkeley.edu

James Demmel  
UC Berkeley  
demmel@cs.berkeley.edu

Yishai Oltchik  
The Hebrew University  
yishai.oltchik@cs.huji.ac.il

Andrew Gearhart  
UC Berkeley  
agearh@cs.berkeley.edu

Oded Schwartz  
The Hebrew University  
odedsc@cs.huji.ac.il

Sivan Toledo  
Tel-Aviv University  
stoledo@tau.ac.il

## ABSTRACT

Network topologies can have significant effect on the execution costs of parallel algorithms due to inter-processor communication. For particular combinations of computations and network topologies, costly network contention may inevitably become a bottleneck, even if algorithms are optimally designed so that each processor communicates as little as possible. We obtain novel contention lower bounds that are functions of the network and the computation graph parameters. For several combinations of fundamental computations and common network topologies, our new analysis improves upon previous per-processor lower bounds which only specify the number of words communicated by the busiest individual processor. We consider torus and mesh topologies, universal fat-trees, and hypercubes; algorithms covered include classical matrix multiplication and direct numerical linear algebra, fast matrix multiplication algorithms, programs that reference arrays,  $N$ -body computations, and the FFT. For example, we show that fast matrix multiplication algorithms (e.g., Strassen’s) running on a 3D torus will suffer from contention bottlenecks. On the other hand, this network is likely sufficient for a classical matrix multiplication algorithm. Our new lower bounds are matched by existing algorithms only in very few cases, leaving many open problems for network and algorithmic design.

## Categories and Subject Descriptors

F.2.1 [Analysis of Algorithms and Problem Complexity]: Numerical Algorithms and Problems—*Computations on matrices*

\*Current affiliation: Google Inc.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.  
Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

## General Terms

Algorithms, Design, Performance.

## Keywords

Network topology, Communication-avoiding algorithms, Strong scaling, Communication costs, Matrix Multiplication, Numerical Linear Algebra, FFT

## 1. INTRODUCTION

Good connectivity of the inter-processor network is necessary for fast execution of parallel algorithms. Insufficient connectivity provably slows down specific parallel algorithms that are communication intensive. Parallel algorithms that ignore network topology can suffer from congestion along network links, and for particular combinations of computations and network topologies, costly network contention may be inevitable even for optimally designed algorithms. In this paper we obtain novel lower bounds on such *contention costs*, and point out cases where this cost is a performance bottleneck for certain, ubiquitous algorithms and network topologies.

In our model a  $\langle P, M, G_{Net} \rangle$ -machine has  $P$  identical processors, each with local memory of size  $M$ , connected with interprocessor network  $G_{Net}$ .<sup>1</sup> The network  $G_{Net}$  may have  $P$  vertices where each vertex is both a processor and a router (direct networks like tori and hypercubes) or may have more than  $P$  vertices, where some vertices represent routers (indirect networks like fat-trees).

Edges of  $G_{Net}$  are network links with weights corresponding to bandwidth. All weights are typically the same in a torus or hypercube, but not in a fat-tree. We ignore processor injection rates in this model, assuming processors can communicate data as fast as their network links allow.

Most previous communication cost lower bounds for parallel algorithms utilize *per-processor* analysis. That is, the lower bounds establish that some processor must communicate a given amount of data. These include classical matrix multiply, direct and iterative linear algebra algorithms,

<sup>1</sup>This model is a variant of the distributed-memory communication model (cf. [7, 14, 17]), where all-to-all connectivity is assumed, and the bandwidth-cost of an algorithm is proportional to the number of words communicated by the worst processor.

FFT, Strassen and Strassen-like fast algorithms, graph related algorithms,  $N$ -body, sorting, programs that reference arrays and others (cf. [1, 4, 5, 7, 13, 19, 24, 26, 29, 32, 40]).

We demonstrate the usefulness of our novel analysis by applying it to a spectrum of algorithm and network combinations, including many of the algorithms above, with networks such as meshes and tori of any dimension, universal fat-trees, and hypercubes. We note that nine out of the twenty fastest supercomputers in the world have either torus or fat-tree network topologies (as of June 2016) [36]. By considering the network graphs, we introduce communication lower bounds for certain computations and networks that are tighter than what was previously known. We also show that often, but not always, the worst contention is expected across the network bisection, supporting a trend that appears in various network designs for decades, namely the importance of considering a network's bisection.

By analyzing the network graphs, we introduce communication lower bounds for certain computations and networks that are tighter than what was previously known. In this work, we bound from below the number of words communicated between a subset of processors and the rest of the processors for a given parallel algorithm (defined by a computation graph and work assignment to the processors), and divide it by the number of links connecting that subset to the rest of the graph. This relates to the *contention cost* of the algorithm, which we specify in Definition 2.2.

Applying the main theorem, we improve (i.e., increase) communication cost lower bounds for several combinations of fundamental computations on common network topologies. Note that we inherit any assumptions made in the original per-processor lower bounds, e.g., assuming no re-computation. Our new lower bounds are matched by existing algorithms only in very few cases, leaving many open problems for network and algorithmic design such as efficient scheduling of heavily utilized computation kernels on (a subset of) a supercomputer.

## 2. CONTENTION LOWER BOUNDS

In this section we state our main result, which translates per-processor communication cost lower bounds to contention cost lower bounds. The following definitions differentiate these costs.

**DEFINITION 2.1.** *Let a parallel algorithm be run on a distributed-memory machine with  $P$  processors. The per-processor bandwidth cost  $W_{proc}$  is the maximum over processors  $1 \leq p \leq P$  of the number of words sent/received by processor  $p$ .*

For ease of notation in the above definition (and the next one) the algorithm in question is omitted, and understood from context. Observe that for many algorithms, there exist two types of per-processor lower bounds: memory-independent  $W_{proc}(P, N)$  (cf. [5]) and memory-dependent  $W_{proc}(P, M, N)$  (cf. [6, 7, 8, 19, 29]) where  $N$  is the input and output size.

**DEFINITION 2.2.** *Consider a  $\langle P, M, G_{Net} \rangle$ -machine running an algorithm. The contention cost  $W_{link}$  is the maximum over edges  $e \in E(G_{Net})$  of the number of words communicated along  $e$  during the execution of the algorithm.*

Recall that the small set expansion  $h_s(G)$  of a  $D$ -regular graph  $G = (V, E)$  is the minimum ratio of edges leaving a

set of vertices of size at most  $s$  [27]. Formally, for  $s \leq |V|/2$ , we have

$$h_s(G) = \min_{S \subseteq V, |S| \leq s} \frac{|E(S, V \setminus S)|}{|E(S, V \setminus S)| + |E(S, S)|}$$

where  $E(S, S)$  is the set of edges that have both endpoints in vertex subset  $S$  and  $E(S, V \setminus S)$  is the set of edges with one endpoint in  $S$  and one endpoint in  $V \setminus S$ .

We present contention lower bounds for two cases: one where  $G_{Net}$  admits an equi-partition such that each part induces a minimal cut, and another where each processor communicates at least  $W_{proc}(P, M, N)$  and  $W_{proc}(P, N)$  words.

**THEOREM 2.3 (MAIN THEOREM).** *Consider some computation with combined input and output data size of  $N$  and per-processor communication lower bounds  $W_{proc}(P, M, N)$  and  $W_{proc}(P, N)$  being performed by a  $\langle P, M, G_{Net} \rangle$ -machine. Consider the following conditions:*

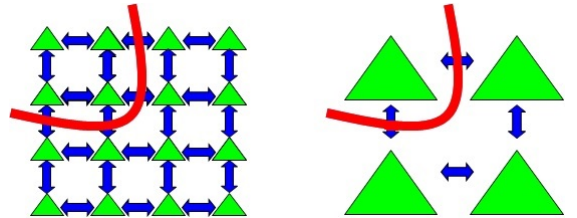
1. *Every  $p \in P$  performs a fraction of  $\Omega(1/P)$  of the flops.*
2. *Every  $p \in P$  stores  $O(N/P)$  of the input and output.*
3. *For any  $t$  that divides  $|V|$ ,  $G_{Net}$  admits an equi-partition  $\Pi$  of  $V$  such that  $|\Pi| = \frac{P}{t}$  and  $\forall A \in \Pi, |E(A, V \setminus A)| = \min_{S \subseteq V, |S|=t} |E(S, V \setminus S)|$ .*

*If at least two of the above conditions are met, then for any  $1 \leq t \leq \frac{|V|}{2}$  that maintains  $\frac{|V|}{t} \in \mathbb{N}$ , the contention cost is bounded below by:*

$$W_{link}(P, M, N) \geq \max_{\substack{S \subseteq V \\ |S|=t}} \frac{W_{proc}(P/t, M \cdot t, N)}{|E(S, V \setminus S)|}$$

and

$$W_{link}(P, N) \geq \max_{\substack{S \subseteq V \\ |S|=t}} \frac{W_{proc}(P/t, N)}{|E(S, V \setminus S)|}.$$



**Figure 1: Computation of  $t = 4$  processors on a 16-processor machine can be emulated as the computation of one processor on a 4-processor machine.**

**PROOF.** Given some graph  $G$  and a partition  $\Pi$  of  $V(G)$ , define  $G_{\Pi}$  to be the result of contracting each part of  $\Pi$  into a single vertex. The proof idea is as follows: given some  $\langle P, M, G_{Net} \rangle$ -machine, identify a partition  $\Pi$  such that at least one vertex  $u \in V(G_{Net})$  is connected to the rest of the graph by  $\min_{\substack{S \subseteq P \\ |S|=t}} (|E(S, V \setminus S)|)$  vertices. Then, emulate the computation of the  $\langle P, M, G_{Net} \rangle$ -machine using such a  $\langle P/t, M \cdot t, G_{\Pi} \rangle$ -machine (see Figure 1), and identify some meta-processor  $p'$  of the emulating machine so

its corresponding processors perform at least  $\Omega(t/P)$  of the flops and store a fraction  $O(t/P)$  of the input and output. By definition of the per-processor communication bounds, the number of words communicated by  $p'$  is bounded below by  $W_{\text{proc}}(P/t, M \cdot t, N)$  and  $W_{\text{proc}}(P/t, N)$ . Therefore, if the meta-processor  $p'$  corresponds to the processor  $u$ , the contention lower bounds are obtained by dividing the total words communicated by the  $\min_{\substack{S \subset V \\ |S|=t}} |E(S, V \setminus S)|$  edges connecting  $p'$  to the rest of the graph.

Recall the three conditions: asymptotic load balance of input and output storage, asymptotic load balance of flops, and the existence of a sparse equi-partition of the network graph. If two of the above are maintained for every meta-processor participating in the computation, then there exists a meta-processor that also maintains the third condition. Therefore, the contention bound will apply to a  $\langle P/t, M \cdot t, G_{\Pi} \rangle$ -machine that emulates the computation.

- (i) Suppose 1 and 2 hold. Let  $S \subset V$  be a subset with  $|S| = t$  and  $|E(S, V \setminus S)| = \min_{\substack{A \subset V \\ |A|=t}} |E(A, V \setminus A)|$ .

By assumption,  $S$  performs a fraction  $\Omega(t/P)$  of the flops and stores  $O(N/P)$  of the input and output. Let  $\Pi$  be a partition of  $G_{Net}$  into subsets of size  $t$  such that  $S \in \Pi$ . Then, it is possible to emulate the computation using a  $\langle \frac{P}{t}, M \cdot t, G_{\Pi} \rangle$ -machine. Let  $u \in V(G_{\Pi})$  be the vertex corresponding to  $S$ . Then  $u$  has a total of  $M \cdot t$  local memory, performs a ratio of  $\Omega(t/P)$  of the computation and has local access to a ratio of  $O(t/P)$  of the input and output. Therefore,  $u$  must communicate a total of  $W_{\text{proc}}(P/t, M \cdot t, N)$  words. Since  $u$  is connected to the rest of  $G_t$  by exactly  $|E(S, V \setminus S)|$  edges, at least one edge must communicate at least  $\frac{W_{\text{proc}}(P/t, M \cdot t, N)}{|E(S, V \setminus S)|}$  words.

- (ii) Suppose 2 and 3 hold. Then, there exists a partition  $\Pi$  such that each part of  $\Pi$  is connected to the rest of the graph by exactly  $\min_{\substack{S \subset V \\ |S|=t}} |E(S, V \setminus S)|$  edges.

Emulate the computation by a  $\langle \frac{P}{t}, M \cdot t, G_{\Pi} \rangle$ -machine. Since 2 holds, every vertex in  $G_{\Pi}$  has stores a fraction  $O(t/P)$  of the input and output data. Additionally,  $\exists v \in V(G_{\Pi})$  such that the processors corresponding to  $v$  perform at least a ratio of  $\Omega(t/P)$  of the computation. Therefore  $v$  must perform communication of  $W_{\text{proc}}(P/t, M \cdot t, N)$  words with the other processors, and at least one edge must communicate at least  $\frac{W_{\text{proc}}(P/t, M \cdot t, N)}{|E(S, V \setminus S)|}$  words.

- (iii) Suppose 1 and 3 hold. The proof is analogous to that of (ii); an averaging argument implies the existence of some subset  $S$  of size  $t$  that stores at most  $O(N/P)$  of the input and output, and therefore the contention bound applies to some  $\langle \frac{P}{t}, M \cdot t, G_{\Pi} \rangle$ -machine.

The proof for the memory-independent bound  $W_{\text{link}}(P, N)$  is analogous.  $\square$

We thus obtain a tool for combining per-processor communication lower bounds and network topologies into a contention lower bound.

**COROLLARY 2.4.** *If in addition to the conditions of Theorem 2.3  $G_{Net}$  is also  $D$ -regular, then the contention cost is*

bounded below by:

$$W_{\text{link}}(P, M, N) \geq \max_{t \in T} \frac{W_{\text{proc}}(P/t, M \cdot t, N)}{D \cdot t \cdot h_t(G_{Net})}$$

and

$$W_{\text{link}}(P, N) \geq \max_{t \in T} \frac{W_{\text{proc}}(P/t, N)}{D \cdot t \cdot h_t(G_{Net})},$$

where

$$T = \{t \mid 1 \leq t \leq P/2, \exists S \subseteq V \text{ such that } |S| = t \text{ and}$$

$$h_t(G) = \frac{|E(S, V \setminus S)|}{|E(S, S)| + |E(S, V \setminus S)|}\}$$

**PROOF.** Recall that for  $D$ -regular graphs,  $D|S| = 2|E(S, S)| + |E(S, V \setminus S)|$  for any  $S \subset V$  with  $|S| \leq \frac{|V|}{2}$ . Therefore, a subset  $S \subset V$  minimizes  $|E(S, V \setminus S)|$  if and only if it minimizes  $\frac{|E(S, V \setminus S)|}{|E(S, S)| + |E(S, V \setminus S)|}$ . Therefore,  $T$  is the set of processor-subset sizes for which sparse equi-partitions (as needed by Theorem 2.3) exist in  $G$ .  $\square$

### 3. PRELIMINARIES

#### 3.1 Per-Processor Lower Bounds

Before deriving the lower bounds on link contention, we review the per-processor communication bounds for several classes of algorithms.

##### *Classical Linear Algebra.*

Most classical direct linear algebra computations can be specified by three nested loops, and for dense  $n \times n$  matrices, the number of flops performed is  $\Theta(n^3)$ .<sup>2</sup> Informally, such computations, which include matrix multiplication, Cholesky and LU decompositions, and many others, can be defined by

$$C_{ij} = f_{ij}(\{g_{ijk}(A_{ik}, B_{kj})\}_{1 \leq k \leq n}) \text{ for } 1 \leq i, j \leq n \quad (1)$$

where  $f$  and  $g$  are sets of functions particular to the computation. For example, in the case of classical matrix multiplication,  $f_{ij}$  is a summation and  $g_{ijk}$  is a scalar multiplication for all  $i, j, k$ . For a more formal definition, see [3, Definition 4.1]. In such a case, we have the following lower bound:

**THEOREM 3.1** ([7],[29]). *Consider an algorithm performing a computation of the form given by equation (1) on  $P$  processors, each with local memory of size  $M$ , and assume one copy of the input data is initially distributed across processors and the computation is load balanced. Then the number of words some processor must communicate is at least*

$$W_{\text{proc}}(P, M, N) = \Omega\left(\frac{n^3}{PM^{1/2}}\right) = \Omega\left(\frac{N^{3/2}}{PM^{1/2}}\right).$$

Note that the local memory size  $M$  appears in the denominator of the expression above, which is why we refer to it as the memory-dependent bound. Additionally, such computations also inherit a memory-independent lower bound:

**THEOREM 3.2** ([5]). *Consider an algorithm performing a computation of the form given by equation (1) on  $P$  processors, and assume just one copy of the input data is initially distributed across processors and the computation is*

<sup>2</sup>For matrix computations, we denote the size of the input/output to be  $N = \Theta(n^2)$ .

load balanced. Then the number of words some processor must communicate is at least

$$W_{\text{proc}}(P, N) = \Omega\left(\frac{n^2}{P^{2/3}}\right) = \Omega\left(\frac{N}{P^{2/3}}\right).$$

### Strassen-like Fast Matrix Multiplication.

Similar lower bounds exist for Strassen’s matrix multiplication and similar algorithms, though the proof techniques differ substantially. Informally, we use the term “Strassen-like” to refer to algorithms that recursively multiply matrices according to a base-case computation. For square algorithms, this corresponds to multiplying  $n_0 \times n_0$  matrices with  $m_0$  scalar multiplications, where  $n_0$  and  $m_0$  are constants. Using recursion, this results in a square matrix multiplication flop count of  $\Theta(n^{\omega_0})$  where  $\omega_0 = \log_{n_0} m_0$ . Note that additional technical assumptions are required for the communication lower bounds to apply and that Strassen-like algorithms may have a rectangular base case; see [8, Section 5.1] for more details. The memory-dependent communication lower bound for Strassen-like algorithms is:

**THEOREM 3.3** ([8, COROLLARY 1.5]). *Consider a Strassen-like matrix multiplication algorithm that requires  $\Theta(n^{\omega_0})$  total flops. Suppose a parallel algorithm performs the computation using  $P$  processors (each with local memory of size  $M$ ) load balances the flops. Then the number of words some processor must communicate is at least*

$$W_{\text{proc}}(P, M, N) = \Omega\left(\frac{n^{\omega_0}}{PM^{\omega_0/2-1}}\right) = \Omega\left(\frac{N^{\omega_0/2}}{PM^{\omega_0/2-1}}\right).$$

Additionally, such computations also inherit a memory-independent lower bound:

**THEOREM 3.4.** *Suppose a parallel algorithm performs a Strassen-like matrix multiplication algorithm requiring  $\Theta(n^{\omega_0})$  flops load balances the computation across  $P$  processors. Then under some technical assumptions (see [8]) the number of words some processor must communicate is at least*

$$W_{\text{proc}}(P, N) = \Omega\left(\frac{n^2}{P^{2/\omega_0}}\right) = \Omega\left(\frac{N}{P^{2/\omega_0}}\right).$$

The proof is identical to [5, Theorem 2.1], with  $\omega_0$  replacing  $\log_2 7$ .

Recently, the per-processor communication lower bounds have been shown also for Strassen’s matrix multiplication even if recomputation is used [11]. In the more general case of Strassen-like algorithms, the previously-known bounds have been extended for any algorithm that does not recompute nontrivial linear combinations [39].

### Programs Referencing Arrays.

The model defined in Equation (1) encompasses most direct linear algebra computations, but lower bounds can be obtained for a considerably more general set of computations. In particular, Christ et al. [19] consider programs of the following form:

$$\begin{aligned} &\text{for all } \mathcal{I} \in \mathcal{Z} \subseteq \mathbb{Z}^d, \text{ in some order,} \\ &\quad \text{inner\_loop}(\mathcal{I}, (A_1, \dots, A_m), (\phi_1, \dots, \phi_m)) \end{aligned} \quad (2)$$

where  $\mathbb{Z}^d$  is the  $d$ -dimensional space of integers and  $\text{inner\_loop}()$  represents a computation involving arrays  $A_1, \dots, A_m$  of dimensions  $d_1, \dots, d_m$  that are referenced by the corresponding subscripts  $\phi_1(\mathcal{I}), \dots, \phi_m(\mathcal{I})$  where  $\phi_i$  are affine maps  $\phi_j$ :

$\mathbb{Z}^d \rightarrow \mathbb{Z}^{d_j}$  for iteration  $\mathcal{I} = (i_1, \dots, i_d)$ . For example, matrix-matrix multiplication has  $(A_1, A_2, A_3) = (A, B, C)$ ,  $\phi_1(\mathcal{I}) = \phi_1(i_1, i_2, i_3) = (i_1, i_3)$ ,  $\phi_2(\mathcal{I}) = \phi_2(i_1, i_2, i_3) = (i_3, i_2)$ ,  $\phi_3(\mathcal{I}) = \phi_3(i_1, i_2, i_3) = (i_1, i_2)$  and the function  $\text{inner\_loop}()$  is defined as  $A_3(\phi_3(\mathcal{I})) = A_3(\phi_3(\mathcal{I})) + A_1(\phi_1(\mathcal{I})) * A_2(\phi_2(\mathcal{I}))$ .

Because the work inside the loop is currently defined as a general function, the space of potential executions of  $\text{inner\_loop}()$  must be restricted in a manageable manner, or to “legal parallel executions” as defined in [19]. To express the lower bounds, we define a set of linear constraints on a vector of unknown scalars  $(s_1, \dots, s_m)$

$$\text{rank}(H) \leq \sum_{j=1}^m s_j \text{rank}(\phi_j(H)), \quad (3)$$

for all subgroups  $H$  of  $\mathbb{Z}^d$ , where  $\text{rank}(H)$  is the cardinality of any maximal subset of Abelian group  $H$  that is linearly independent.<sup>3</sup> For such computations we have the following lower bound:

**THEOREM 3.5** ([19]). *Consider an algorithm performing a computation of the form given by (2) on  $P$  processors, each with local memory of size  $M$ , and assume the input data is initially evenly distributed across processors. Then for any legal parallel execution and sufficiently large  $|\mathcal{Z}|/P$ , the number of words some processor must communicate is at least*

$$W_{\text{proc}}(P, M, N) = \Omega\left(\frac{|\mathcal{Z}|}{PM^{s_{\text{HBL}}-1}}\right),$$

where  $s_{\text{HBL}}$  is the minimum value of  $\sum_{i=1}^m s_i$  subject to Inequality (3), assuming that this linear program is feasible (see [19]).

We restate the memory-independent bound from [19] for such computations (note that the formal proof has not yet appeared). For legal parallel executions of computations of the form (2) on  $P$  processors, some processor must move

$$W_{\text{proc}}(P, N) = \Omega\left(\left(\frac{|\mathcal{Z}|}{P}\right)^{1/s_{\text{HBL}}}\right) \quad (4)$$

words where  $s_{\text{HBL}}$  is defined as in Theorem 3.5.

Note that Theorem 3.5 generalizes Theorem 3.1. For example, matrix multiplication satisfies both forms (1) and (2), where in the latter case  $|\mathcal{Z}| = n^3$  and  $s_{\text{HBL}} = 3/2$ .

Theorem 3.5 also applies to, for example,  $N$ -body computations where all pairs of interactions are computed [22]. In this case,  $|\mathcal{Z}| = \Theta(N^2)$  and  $s_{\text{HBL}} = 2$ , yielding lower bounds of  $W_{\text{proc}}(P, M, N) = \Omega(N^2/(PM))$  and  $W_{\text{proc}}(P, N) = \Omega(N/P^{1/2})$ . We also note that Theorem 3.5 applies to  $N$ -body computations that use a distance cutoff to reduce the number of neighbor interactions, i.e.  $|\mathcal{Z}| \ll N^2$ .

### FFT/Sorting.

We next discuss per-processor communication cost bounds for the FFT and comparison sorts.

A sequential communication cost lower bound of  $\Omega(n \log n / \log M)$  was given by Hong and Kung [26]. A parallel memory-independent per-processor bound has been proven for the LPRAM model [1] and the BSP model [13].

<sup>3</sup>The rank of an Abelian group is analogous to the concept of the dimension of a vector space.

**THEOREM 3.6** ([1, 13]). *The per-processor communication cost of FFT algorithm on input of size  $N$ , run on a  $\langle P, M, G_{Net} \rangle$ -machine with no recomputation is bounded below by*

$$W_{proc}(P, N) = \Omega\left(\frac{N \log N}{P \log(N/P)}\right).$$

We are unaware of a previous memory-dependent per-processor bound for FFT, although this is a straightforward extension of the existing sequential lower bounds in [26] or the one in [37] for memory hierarchy model. We provide it here for completeness.

**THEOREM 3.7.** *The per-processor communication cost of FFT algorithm on input of size  $N$ , run on a  $\langle P, M, G_{Net} \rangle$ -machine with no recomputation is*

$$W_{proc}(P, M, N) = \Omega\left(\frac{N \log N}{P \log M} - M\right).$$

**PROOF.** The bound follows from the same logic as the partitioning argument in Hong-Kung [26], applied to one of the processors that performs  $\Theta((N \log N)/P)$  of the computations. Since (as in [26]) each  $M$ -partition may have no more than  $M \log M$  vertices, the number of  $M$ -partitions is at least  $N \log N / (PM \log M)$ . Thus, the total per-processor bandwidth is  $\Omega(M \lfloor (N \log N) / (PM \log M) \rfloor)$ , and the theorem follows.  $\square$

Note that this memory dependent bound is dominated by the memory independent one, since we assume  $MP \geq N$ .

## 3.2 Torus, Mesh and Hypercube Networks

### 3.2.1 Torus Networks

**DEFINITION 3.8.** *A  $D$ -torus graph  $G = \langle V, E \rangle$  is the Cartesian product of  $D$  cycle graphs  $C^1, \dots, C^D$ . That is,  $V = \mathbb{Z}_{|C^1|} \times \dots \times \mathbb{Z}_{|C^D|}$  and  $(u, v) \in E$  if and only if  $u, v$  disagree on a single index  $i$  and  $|u_i - v_i| = 1 \pmod{C_i}$ . A  $D$ -torus is said to be of the form  $[n]^D$  if  $C^1 = \dots = C^D = C_n$ .*

Torus networks are common topologies among current supercomputers (six of the current top twenty [36], for example). Cray's XK7 [20] and IBM's Blue Gene/P [28] machines utilize 3D tori, Blue Gene/Q a 5-dimensional torus [18], and the K computer in Japan has a 6-dimensional network topology [2]. Intel Xeon Phi coprocessors rely on a ring-based (a 1-dimensional torus) on-chip communication network between cores [30]. In this section, we derive a tight bound on the network small set expansion for this class of networks.

From [15], if  $G$  is a torus graph of the form  $[n]^D$  then for any  $S \subset V$  with  $|S| = t \leq |V|/2$  we have the tight inequality:

$$E(S, V \setminus S) \geq \min_{r \in [D]} \{2rt^{1-1/r}n^{(D/r)-1}\}$$

**LEMMA 3.9.** *For a torus graph  $G$  of the form  $[n]^D$  and any  $t \in \{1, \dots, \frac{n^D}{2}\}$ , the small set expansion is*

$$h_t(G) = \begin{cases} \frac{2}{t^{D/r} + 1} & t \leq \left(\frac{n}{e}\right)^D \\ \frac{2 \ln \frac{n^D}{t}}{\frac{Dn}{e} + \ln \frac{n^D}{t}} & t > \left(\frac{n}{e}\right)^D \end{cases}$$

**PROOF.** As a  $D$ -dimensional torus is  $2D$  regular, any subset  $S$  also maintains  $2D|S| = 2|E(S, S)| + |E(S, V \setminus S)|$ . From this follows that minimizing  $|E(S', V \setminus S')|$  is equivalent to minimizing  $\frac{|E(S', V \setminus S')|}{|E(S', V \setminus S')| + |E(S', S')|}$ .  $\square$

Given some  $t$  that divides  $n^D$ , we next describe an equi-partition  $\Pi_t$  of  $G$  such that every part  $A \in \Pi_t$  maintains  $|E(A, V \setminus A)| = \min_{r \in [D]} \{2r|A|^{1-1/r}n^{(D/r)-1}\}$ .

### Equi-Partitions of Torus Networks.

The construction of an expanding equi-partition  $\Pi$  of tori is straightforward. Consider  $t$  that divides  $|V|$ . Partition  $G$  into sub-tori of the form  $[n]^{D-r} \times [a]^r$  where  $a = \frac{t^{1/r}}{n^{D/r-1}}$  and  $r \in [D]$ . Intuitively,  $r$  denotes the dimensions of  $G$  that are not wholly contained by the sub-torus, and these partially-filled dimensions are  $r$ -dimensional cubes on  $\frac{t}{n^{D-r}}$  vertices.

If  $t \leq \left(\frac{n}{e}\right)^D$  then the minimum cut is attained by the sub-torus with  $r = D$ , and as  $t$  divides  $n^D$ ,  $t^{1/D}$  divides  $n$ . Therefore,  $\frac{|E(S, V \setminus S)|}{|E(S, S)| + |E(S, V \setminus S)|} = \frac{2}{t^{D/r} + 1}$  for any  $S \in \Pi$ . Otherwise, the minimum cut is attained by a sub-torus with  $r = \ln\left(\frac{n^D}{t}\right)$  and  $a = \frac{n}{e}$ . In this case, every part  $S$  of  $\Pi$

$$\text{maintains } \frac{|E(S, V \setminus S)|}{|E(S, S)| + |E(S, V \setminus S)|} = \frac{2 \ln\left(\frac{n^D}{t}\right)}{\frac{Dn}{e} + \ln\left(\frac{n^D}{t}\right)}.$$

It is assumed above that  $t^{1/D} \in \mathbb{N}$ . If not, we define a nearly-equal partition  $\Pi$  into sub-tori of the form  $[n]^D - r \times [a_1]^r - k \times [a_2]^k$  where  $a_1 = \lfloor a \rfloor$  and  $a_2 = \lceil a \rceil$ . Different parts of  $\Pi$  may have different values of  $k$ ; this only affects the analysis by at most a constant factor.

### 3.2.2 Mesh Networks

**DEFINITION 3.10.** *A  $D$ -dimensional mesh graph  $G = \langle V, E \rangle$  is the Cartesian product of  $D$  path graphs  $P^1, \dots, P^D$ . That is,  $V = \mathbb{Z}_{|P^1|} \times \dots \times \mathbb{Z}_{|P^D|}$  and  $(u, v) \in E$  if and only if  $u, v$  disagree on a single index  $i$  and  $|u_i - v_i| = 1$ . A  $D$ -dimensional mesh is said to be of the form  $[n]^D$  if  $P^1 = \dots = P^D = P_n$ .*

Unlike other direct topologies discussed in this paper, meshes are not  $D$ -regular. However, all degrees in a  $D$ -dimensional mesh are between  $\frac{D}{2}$  and  $D$ . In that case, we refer instead to the metric  $\frac{|E(S, V \setminus S)|}{D|S|}$ . Observe that for a  $D$ -torus,  $\frac{|E(S, V \setminus S)|}{D|S|} \leq h_s(G) \leq \frac{2|E(S, V \setminus S)|}{D|S|}$ . Therefore the results for grid mesh networks are within factor 2 of the results for torus networks, with the minor caveat that a sparse equi-partition  $\Pi$  in a mesh grid can have  $\max_{S_1 \in \Pi} (|E(S_1, V \setminus S_1)|) = 2 \cdot \min_{S_2 \in \Pi} (|E(S_2, V \setminus S_2)|)$  instead of being exactly equal for any  $S \in \Pi$ .

### 3.2.3 Hypercube Networks

**DEFINITION 3.11.** *A hypercube graph on  $2^n$  vertices  $G = \langle V, E \rangle$  is the Cartesian product of  $n$  path graphs of length 2. That is,  $V = \underbrace{[C_2] \times \dots \times [C_2]}_{n \text{ times}}$  and  $(u, v) \in E$  if and only if  $u, v$  have Hamming distance of 1.*

**LEMMA 3.12.** *The small-set expansion of a hypercube graph  $G$  on  $2^k$  vertices is  $h_t(G) = \frac{2(k - \log t)}{2k - \log t}$*

**PROOF.** Given a hypercube graph  $G = \langle V, E \rangle$  on  $2^k$  vertices and  $t$  divides  $2^k$ , consider a sub-hypercube  $G_t = \langle V_t, E_t \rangle$  on  $t = 2^m$  vertices,  $m < k$ . As a sub-hypercube attains the sparsest cut in a hypercube [34], and therefore characterizing its expansion is sufficient. Selecting a  $2^m$  sub-hypercube

is equivalent to fixing  $m$  coordinates of the vertices, (enumerating them in standard binary fashion), and then each vertex in  $G_t$  contributes exactly  $(k - m)$  edges to the cut. Therefore  $|E(V_t, V \setminus V_t)| = 2^m \cdot (k - m)$ . Since an  $n$ -hypercube has  $n \cdot 2^{n-1}$  edges,  $|E(V_t, V_t)| = m \cdot 2^{m-1}$ . Therefore:

$$h_t(G) = \frac{2^m(k - m)}{2^m(k - m) + m \cdot 2^{m-1}} = \frac{2(k - \log t)}{2k - \log t} \quad \square$$

### Equi-Partitions of Hypercube Networks.

To construct an equi-partition of a hypercube  $G$  on  $2^k$  vertices into subsets of size  $t = 2^m$ , we decompose  $G$  into disjoint sub-hypercubes on  $t$  vertices each.

### 3.3 Universal Fat-Tree Networks

We also consider universal fat-trees [33], which are indirect networks consisting of processors connected by a binary tree of router nodes. As of June 2016 [36], three of the world's top twenty fastest supercomputers have this topology. Also, the world's second fastest supercomputer is China's National University of Defense Technology Tianhe-2, which has a custom interconnect with a fat-tree topology [21]. In a fat-tree network, the bandwidth capacity of links between routers varies, increasing from links connecting the processor leaves to the "fattest" links connected to the root node. Universal fat-trees are parametrized by the root links' capacity:  $w$  words of data per unit time, where the capacity of the leaf links (to processors) is normalized to 1. Given the parameter  $w$ , the capacity of each link at level  $0 \leq i \leq \log P$  in the tree is  $\min\{P/2^i, w/2^{2i/3}\}$ . This implies that the capacities of the links between subsequent levels of the tree differs by either a factor of 2 or  $2^{2/3}$ . Because the fat-tree network graph consists of both processor and router nodes (an indirect network) and link capacities are variable, we may not be able to straightforwardly define meaningful equi-partitions of  $G_{Net}$ 's vertices for Theorem 2.3 and will instead apply stronger load-balancing requirements to the computations.

## 4. APPLICATIONS

In this section, we apply the general contention lower bounds of Theorem 2.3 to pairs of computations and network topologies. That is, we combine the previously known per-processor bounds of a computation with properties of the network topology to obtain novel contention bounds that are tighter in some scenarios. We first derive contention lower bounds for all the computations on  $D$ -dimensional tori (or meshes) in Section 4.1, and then in Section 4.2 we compare the bounds to determine in which scenarios each lower bound dominates. We also derive contention bounds for the computations on fat-tree topologies in Section 4.4; we discuss the comparisons among bounds for that topology more briefly, as the relationships are qualitatively similar. Table 1 presents the communication lower bounds for each of the computations on both  $D$ -dimensional tori and fat-trees with root capacity  $w$ .

### 4.1 Contention Lower Bounds for Tori

In this section, we derive contention lower bounds by plugging the memory-dependent and memory-independent per-processor lower bounds [5, 8, 19, 29] into Theorem 2.3 and using the properties of  $D$ -dimensional tori. Table 1 summarizes these results.

### Direct Linear Algebra, Strassen-like, and $N$ -body.

We apply Theorem 2.3 to the relevant per-processor bounds given in Section 3.1. Let  $F$  denote the number of work operations (e.g. flops or loop iterations) of the different computations. The per-processor memory-dependent bound is thus:

$$W_{\text{proc}}(P, M, N) = \Omega\left(\frac{F}{PM^{\alpha-1}}\right) \quad (5)$$

where  $\alpha = 3/2$  for direct dense linear algebra,  $\alpha = \omega_0/2$  for Strassen-like matrix multiplication,  $\alpha = 2$  for the  $O(N^2)$   $N$ -body problem. By Lemma 3.9, the memory-dependent contention bound is:

$$W_{\text{link}}^{\text{tor}}(P, M, N) = \max_{t \in T} \Omega\left(\frac{F \cdot D}{PM^{\alpha-1}} \cdot t^{1-\alpha+1/D}\right).$$

Note that  $t^{1-\alpha+1/D}$  is monotonic (in the given range), but that the exponent can be positive, negative or zero. If the exponent of  $t$  is negative or zero, then the expression is maximized at  $t = 1$ , reproducing the per-processor bound (up to a constant factor). If the exponent is positive, namely  $D \leq D_1 = 1/(\alpha - 1)$ , then the expression is maximized at  $t = P/2$ ,<sup>4</sup> and we obtain a new and tighter bound:

$$W_{\text{link}}^{\text{tor}}(P, M, N) = \Omega\left(\frac{F \cdot D}{P^{\alpha-1/D} M^{\alpha-1}}\right).$$

The per-processor memory-independent bound is

$$W_{\text{proc}}(P, N) = \Omega\left(\frac{N}{P^{1/\alpha}}\right), \quad (6)$$

and we can apply Theorem 2.3 to obtain:

$$W_{\text{link}}^{\text{tor}}(P, N) = \max_{t \in T} \Omega\left(\frac{N \cdot D}{P^{1/\alpha}} \cdot t^{1/\alpha-1+1/D}\right).$$

Again,  $t^{1/\alpha-1+1/D}$  is monotonic and may be positive, negative or zero. If the exponent of  $t$  is negative or zero, then the expression is maximized at  $t = 1$ , reproducing the per-processor bound (up to a constant factor). If the exponent is positive, e.g.  $D \leq D_2 = \alpha/(\alpha - 1)$ , then the expression is maximized at  $t = P/2$ , and we obtain a tighter new bound:

$$W_{\text{link}}^{\text{tor}}(P, N) = \Omega\left(\frac{N \cdot D}{P^{1-1/D}}\right).$$

### Programs that Reference Arrays.

According to Theorem 3.5, the memory-dependent per-processor bandwidth lower bound for programs defined by (2) is  $W_{\text{proc}}(P, M, N) = \Omega(|Z|/(PM^{s_{\text{HBL}}-1}))$ . Similar to the derivation for the previous problems (albeit with  $\alpha = s_{\text{HBL}}$ ), the contention bound becomes

$$W_{\text{link}}^{\text{tor}}(P, M, N) = \max_{1 \leq t \leq P/2} \Omega\left(\frac{|Z|}{PM^{s_{\text{HBL}}-1}} \cdot t^{1-s_{\text{HBL}}+1/D}\right)$$

which is maximized at either  $t = 1$  (the per-processor bound), or  $t = P/2$  (see Footnote 4). So, we obtain

$$W_{\text{link}}^{\text{tor}}(P, M, N) = \Omega\left(\frac{|Z|}{P^{s_{\text{HBL}}-1/D} M^{s_{\text{HBL}}-1}}\right)$$

<sup>4</sup> Note that there may not be a subset of the vertices of  $G_{Net}$  that attains the small set expansion  $h_t(G_{Net})$  of size *exactly*  $P/2$ . However, the small set expansion of tori and meshes is attained for small sets of size  $P/c$  for some constant  $c \geq 2$  (e.g. consider a sub-torus), hence the following contention analysis holds up to a constant factor.

as a memory-dependent lower bound on contention. In a similar manner, we can derive a memory-independent contention lower bound from Equation (4). Observing that the contention bound is maximized at either  $t = 1$  or  $t = P/2$ , we derive the memory-independent lower bound on contention at  $t = P/2$ :  $W_{\text{link}}^{\text{tor}}(P, N) = \Omega(|\mathcal{Z}|^{1/s_{\text{HBL}}}/(P^{1-1/D}))$ .

### FFT/Sorting.

As with the previous algorithms, we apply Theorem 2.3 to the relevant per-processor bounds given in Section 3.1. As we mentioned earlier, the memory-independent per-processor bound always dominates the memory-dependent bound because we assume that  $M \geq N/P$ .

The per-processor memory-independent bound is

$$W_{\text{proc}}(P, N) = \Omega\left(\frac{N \log(N)}{P \log(N/P)}\right),$$

and we can apply this bound to Theorem 2.3 to obtain:

$$\begin{aligned} W_{\text{link}}^{\text{tor}}(P, N) &= \max_{t \in T} \Omega\left(\frac{N \log(N) \cdot D}{P \log(Nt/P)t^{-1/D}}\right) \\ &= \frac{N \log(N) \cdot D}{P} \max_{t \in T} \Omega\left(\frac{t^{1/D}}{\log(Nt/P)}\right). \end{aligned} \quad (7)$$

Again, when  $t = 1$  we obtain the original per-processor bound. Equation (7) has a stationary point at  $t = P2^D/N$ , but via consideration of the second derivative with respect to  $t$ , it can be shown that this point is a minimum for all relevant values of  $N$ ,  $P$ , and  $D$ . Thus, we can derive a memory-independent contention bound by setting  $t = P/2$ :

$$W_{\text{link}}^{\text{tor}}(P, N) = \Omega\left(\frac{N \cdot D}{P^{1-1/D}}\right).$$

## 4.2 Analysis and Interpretation for Tori

### Which bound dominates?

Our first observation is that, for these computations, the memory-independent contention bound dominates the memory-dependent contention bound for many algorithms. In the cases of direct linear algebra, Strassen and Strassen-like, and the  $O(N^2)$  N-body problem we prove this by contradiction: if the memory-dependent contention bound dominates, then the problem is too large to be distributed across all the processors' local memories. Thus, if

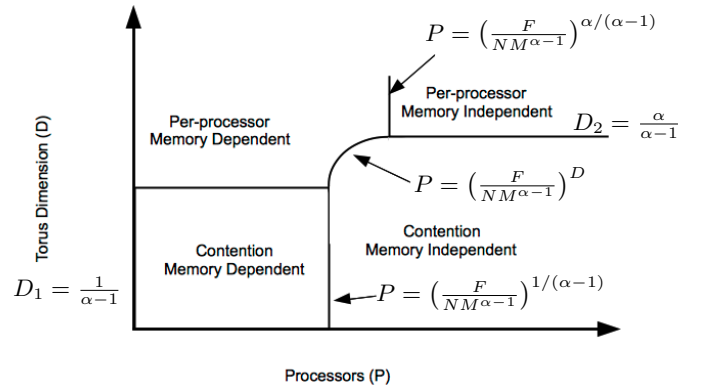
$$\frac{F}{P^{\alpha-1/D}M^{\alpha-1}} > \frac{N}{P^{1-1/D}}$$

then, as  $F = \theta(N^\alpha)$ , we have  $N^{\alpha-1} > P^{\alpha-1}M^{\alpha-1}$  which is a contradiction as we assumed that  $N \leq PM$ . For programs that reference arrays, the proof requires a bit more of the theoretical apparatus from [19] and is proven in Appendix A. We note that in practice the value of constants may result in the memory-dependent contention bound being dominant, despite the asymptotic result.

For direct linear algebra, Strassen, Strassen-like and  $O(N^2)$  N-body algorithms, Figure 2 illustrates the relationships between the four types of communication lower bounds for a fixed computation, fixed problem size  $N$ , and fixed local memory size  $M$ , varying the number of processors  $P$  and the torus dimension  $D$ . See Appendix B for the derivation of the expressions used in Figure 2.

		Memory Dependent	Memory Independent
Direct Linear Algebra	$W_{\text{proc}}$	$\Omega\left(\frac{N^{\frac{3}{2}}}{PM^{\frac{1}{2}}}\right)$	$\Omega\left(\frac{N}{P^{\frac{2}{3}}}\right)$
	$W_{\text{link}}^{\text{tor}}$	$\Omega\left(\frac{N^{\frac{3}{2}} \cdot D}{P^{\frac{3}{2}} - \frac{1}{D} M^{\frac{1}{2}}}\right)$	$\Omega\left(\frac{N \cdot D}{P^{1-\frac{1}{D}}}\right)$
	$W_{\text{link}}^{\text{f-t}}$	$\Omega\left(\frac{N^{\frac{3}{2}}}{w(MP)^{\frac{1}{2}}}\right)$	$\Omega\left(\frac{N}{w}\right)$
Strassen and Strassen-like	$W_{\text{proc}}$	$\Omega\left(\frac{N^{\frac{\omega_0}{2}}}{PM^{\frac{\omega_0}{2}-1}}\right)$	$\Omega\left(\frac{N}{P^{\frac{2}{\omega_0}}}\right)$
	$W_{\text{link}}^{\text{tor}}$	$\Omega\left(\frac{N^{\frac{\omega_0}{2}} \cdot D}{P^{\frac{\omega_0}{2}} - \frac{1}{D} M^{\frac{\omega_0}{2}-1}}\right)$	$\Omega\left(\frac{N \cdot D}{P^{1-\frac{1}{D}}}\right)$
	$W_{\text{link}}^{\text{f-t}}$	$\Omega\left(\frac{N^{\frac{\omega_0}{2}}}{w(MP)^{\frac{\omega_0}{2}-1}}\right)$	$\Omega\left(\frac{N}{w}\right)$
N-body	$W_{\text{proc}}$	$\Omega\left(\frac{N^2}{PM}\right)$	$\Omega\left(\frac{N}{P^{1/2}}\right)$
	$W_{\text{link}}^{\text{tor}}$	$\Omega\left(\frac{N^2 \cdot D}{P^2 - \frac{1}{D} M}\right)$	$\Omega\left(\frac{N \cdot D}{P^{1-\frac{1}{D}}}\right)$
	$W_{\text{link}}^{\text{f-t}}$	$\Omega\left(\frac{N^2}{wMP}\right)$	$\Omega\left(\frac{N}{w}\right)$
Programs Referencing Arrays	$W_{\text{proc}}$	$\Omega\left(\frac{F}{PM^{\tilde{s}-1}}\right)$	$\Omega\left(\left(\frac{F}{P}\right)^{\frac{1}{\tilde{s}}}\right)$
	$W_{\text{link}}^{\text{tor}}$	$\Omega\left(\frac{D \cdot F}{P^{\tilde{s}} - \frac{1}{D} M^{\tilde{s}-1}}\right)$	$\Omega\left(\frac{D \cdot F^{\frac{1}{\tilde{s}}}}{P^{1-\frac{1}{D}}}\right)$
	$W_{\text{link}}^{\text{f-t}}$	$\Omega\left(\frac{F}{w(MP)^{\tilde{s}-1}}\right)$	$\Omega\left(\frac{F^{\frac{1}{\tilde{s}}}}{w}\right)$
FFT	$W_{\text{proc}}$	$\Omega\left(\frac{N \log(N)}{P \log(M)}\right)$	$\Omega\left(\frac{N \log(N)}{P \log\left(\frac{N}{P}\right)}\right)$
	$W_{\text{link}}^{\text{tor}}$	$\Omega\left(\frac{N \log(N) \cdot D}{P^{1-\frac{1}{D}} \log(MP)}\right)$	$\Omega\left(\frac{N \cdot D}{P^{1-\frac{1}{D}}}\right)$
	$W_{\text{link}}^{\text{f-t}}$	$\Omega\left(\frac{N \log N}{w \log(MP)}\right)$	$\Omega\left(\frac{N}{w}\right)$

**Table 1: Per-processor bounds ( $W_{\text{proc}}$ ) vs. the new contention bounds ( $W_{\text{link}}$ ) on a  $D$ -dimensional torus and fat-trees with root capacity  $w$  for classical (dense) linear algebra, fast matrix multiplication,  $O(N^2)$  N-body, a general set of programs that reference arrays, and Fast Fourier Transform (FFT). For readability, we denote  $\tilde{s} = s_{\text{HBL}}$ .**



**Figure 2: Relationship between the per-processor and contention communication lower bounds for direct linear algebra, Strassen/Strassen-like and the  $O(N^2)$  N-body problems.**

Depending on the dimension of the torus and number of processors, the tightest bound may be one of the previously known per-processor bounds or the memory-independent contention bound. We first consider subdividing the vertical axis of Figure 2, which corresponds to the torus dimension  $D$ . Intuitively speaking, the smaller  $D$  is, the more likely contention will dominate communication costs. For a given algorithm, we let  $D = \lfloor 1/(\alpha - 1) \rfloor = \lfloor D_1 \rfloor$  is the maximum torus dimension such that the communication cost is dominated by contention for all input and machine parameters. Similarly, we let  $D = \lceil \alpha/(\alpha - 1) \rceil = \lceil D_2 \rceil$  be the minimum torus dimension so that the communication cost is not dominated by the contention (at least not by the bound proved here). Note that for a combination of an algorithm and a  $D$ -dimensional torus such that  $D_1 < D < D_2$ , either the per-processor memory-dependent bound or the memory-independent contention bound may dominate. See Table 2 for values of  $D_1$  and  $D_2$  for various matrix multiplication algorithms. In particular, note that for the classical algorithm, a 2D torus is not sufficient to avoid contention. While Cannon’s algorithm [16] does not suffer from contention on a 2D torus network, it is also not communication-optimal. The more communication-efficient “3D” algorithms [10, 1, 35, 41], which utilize extra memory and have the ability to strong scale perfectly, require a 3D torus to attain the per-processor lower bounds. For matrix multiplication algorithms with smaller exponents, the torus dimension requirements for remaining contention-free are even larger.

### Range of perfect strong scaling

We next consider subdividing the horizontal axis of Figure 2, which corresponds to the number of processors  $P$ . Because Figure 2 shows a fixed problem size, increasing  $P$  (moving to the right) corresponds to “strong scaling.” We differentiate between whether or not the computation has the possibility of strong scaling perfectly: that is, for a fixed problem size, increasing the number of processors by a constant factor reduces the communication costs (and running time) by the same constant factor. Note that of the bounds, the memory-dependent per-processor bound (see Equation (5) for example) exhibits this possibility of perfect strong scaling, as  $P$  appears in the denominator with an exponent of 1. However, as  $P$  increases, one of the memory-independent bounds eventually dominates and perfect strong scaling is no longer possible. See [5] for a discussion of this behavior given only per-processor bounds.

For direct linear algebra, Strassen-like algorithms and the  $O(N^2)$   $N$ -body problem, if  $D \geq D_2$  and  $P \leq (F/NM^{\alpha-1})^{\frac{\alpha}{\alpha-1}}$ , then the memory-dependent per-processor bound dominates. When this happens, we have a perfect strong scaling range. For values of  $P$  beyond this range, the communication cost is dominated by the memory-independent per-processor bound (see [5] for further discussion). When  $D_1 < D < D_2$ , a smaller strong-scaling ranges exists for  $P \leq (F/NM^{\alpha-1})^D$ ; for values of  $P$  beyond this range, the communication cost bound is dominated by contention. If  $D \leq D_1$ , then the contention bounds always dominate and there is no strong-scaling range. A similar analysis can demonstrate such a region of perfect strong scaling in runtime for programs that reference arrays.

Figure 3 shows this behavior for Strassen’s matrix multiplication (where  $\alpha = (\log_2 7)/2$ ) given the relevant torus dimensions. For Strassen,  $F/NM^{\alpha-1} = (N/M)^{\alpha-1} = P_{min}^{\alpha-1}$ ,

where  $P_{min}$  is the minimum number of processors required to store the problem as  $F = O(n^\alpha)$ . Note that the lower subfigure in Figure 3 is a log-log scale, while the upper subfigure’s y-axis is linear. For a sufficiently good network ( $D \geq 4$ ), the perfect strong scaling range is  $P_{min} < P < P_{min}^{(\log_2 7)^{1/2}} \approx P_{min}^{1.40}$ . For a 3D torus, the perfect strong scaling range shrinks to  $P_{min} < P < P_{min}^{3(\log_2 7-2)/2} \approx P_{min}^{1.21}$ . On 2D torus, perfect strong scaling is impossible. These three regions of network dimension ( $D \geq D_2$ ,  $D \leq D_1$  and  $D_1 < D < D_2$ ) are illustrated in Figure 2 as being the points of transition between dominance of the various bounds. The upper portion of Figure 3 demonstrates the regions of dominance for the various network dimensions in the case of Strassen’s algorithm.

Algorithm	$\omega_0$	$\lfloor D_1 \rfloor$	$\lceil D_2 \rceil$
Classical	3	2	3
Strassen (1969) [42]	$\approx 2.81$	2	4
Schönhage (1981) [38]	$\approx 2.55$	3	5
Strassen (1987) [43]	$\approx 2.48$	4	6
Le Gall (2014) [23]	$\approx 2.3729$	5	7

**Table 2: Torus dimensions so that communication cost is either always ( $D \leq \lfloor D_1 \rfloor$ ) or never ( $D \geq \lceil D_2 \rceil$ ) contention-bound, for several matrix multiplication algorithms. Assertions for the last three algorithms are under some technical assumptions, see [8].**

### 4.3 Contention Lower Bounds for Hypercubes

On hypercubes our analysis does not show meaningful contention bounds for the algorithms discussed in this paper. It is easy to verify that the per-processor communication lower bounds always dominate the contention lower bounds by plugging them into Theorem 2.3 and comparing the resulting bounds to the respective per-processor bounds. For a complete analysis, see Appendix C.

### 4.4 Contention Lower Bounds for Fat-Trees

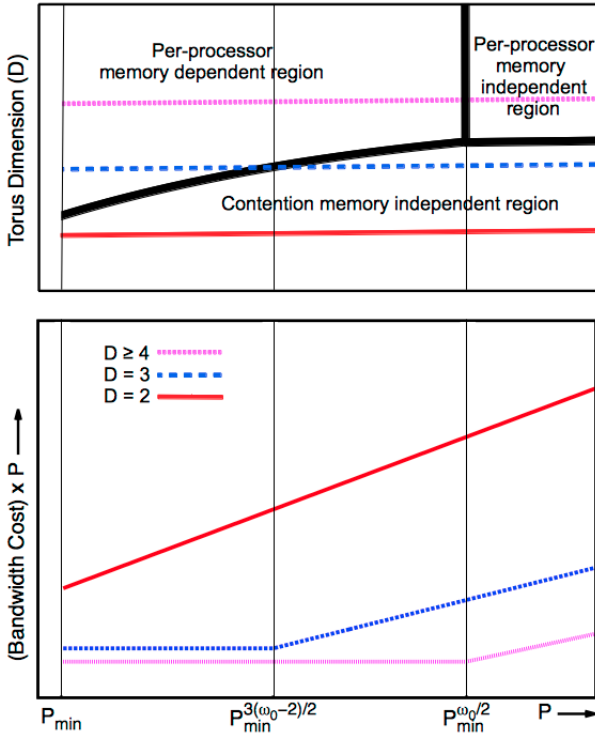
To obtain contention bounds for fat-trees, we define a set of partitions  $\{\mathcal{P}_i\}$  and compute the corresponding aggregate bandwidths  $\{L_i\}$ . We consider partitioning the processors into complete subtrees of sizes ranging from one processor to  $P/2$  processors. Then, each subset of processors has one external link from its root node to the next level in the tree. Let  $\mathcal{P}_i$  be a partition of the processors into  $2^i$  subsets, where each subset consists of a complete binary subtree of  $P/2^i$  processors. The minimum aggregate bandwidth is given by the capacity of the link connected to the root of the subtree:

$$L_i = \min \left\{ \frac{P}{2^i}, \frac{w}{2^{2i/3}} \right\}. \quad (8)$$

#### Direct Linear Algebra, Strassen-like, and $N$ -body.

We first consider direct linear algebra, Strassen-like matrix multiplication, and  $O(N^2)$   $N$ -body algorithms with per-processor lower bounds given by Equations (5) and (6) where  $1 < \alpha \leq 2$ , depending on the algorithm. The aggregate bandwidth of a complete subtree of size  $P/2^i$  is given by Equation (8), so by Theorem 2.3 our memory-dependent





**Figure 3: Communication bounds for Strassen's algorithm on  $D$ -dimensional tori. The lower plot is log-log, while the upper is linear on the y-axis. Horizontal lines in the lower plot correspond to perfect strong scaling.**

contention bound is

$$\begin{aligned} W_{\text{link}}^{\text{f-t}}(P, M, N) &= \max_{1 \leq i \leq \log P} \Omega \left( \frac{\frac{F}{2^i (MP/2^i)^{\alpha-1}}}{\min\{P/2^i, w/2^{2i/3}\}} \right) \\ &= \max_{1 \leq i \leq \log P} \Omega \left( \frac{F \cdot 2^{i(\alpha-1)}}{P^\alpha M^{\alpha-1}} + \frac{F \cdot 2^{i(\alpha-4/3)}}{w (MP)^{\alpha-1}} \right). \end{aligned}$$

The first term is an increasing function of  $i$  because  $\alpha > 1$ . The second term can be either increasing or decreasing (it will be decreasing for a fast matrix multiplication algorithm with exponent  $\omega_0 < 2.66$ ). The maximum is therefore achieved at either  $i = 1$  or  $i = \log P$ , so in order to obtain a new lower bound we evaluate the expression at  $i = 1$  to obtain  $W_{\text{link}}^{\text{f-t}}(P, M, N) = \Omega(F/(w(MP)^{\alpha-1}))$ . Likewise, the memory-independent contention bound is

$$\begin{aligned} W_{\text{link}}^{\text{f-t}}(P, N) &= \max_{1 \leq i \leq \log P} \Omega \left( \frac{N/(2^i)^{1/\alpha}}{\min\{P/2^i, w/2^{2i/3}\}} \right) \\ &= \max_{1 \leq i \leq \log P} \Omega \left( \frac{N}{P} \cdot 2^{i(1-1/\alpha)} + \frac{N}{w} \cdot 2^{i(2/3-1/\alpha)} \right). \end{aligned}$$

As before, the first term is always increasing but the second term can be decreasing (this time for fast matrix multiplication with exponent  $\omega_0 < 3$ ), so the maximum could be achieved at either endpoint. By plugging in  $i = 1$  we obtain a new bound of  $W_{\text{link}}^{\text{f-t}}(P, N) = \Omega(N/w)$ .

As in the case of torus networks, of the two contention bounds, the memory-independent one dominates the memory-dependent one assuming  $N < MP$ , or that the data fits

across all processors' memories. Relationships among the two per-processor bounds and memory-independent contention bound can be derived as in Section 4.2 for torus networks (see Figures 2 and 3). We note that the memory-independent contention bound is the tightest bound for Strassen-like matrix multiplication when, for example, the fat-tree parameter  $w$  falls in the range  $P^{2/3} \leq w < P^{2/\omega_0} = P^{1/\alpha}$  and  $P = \Omega((N/M)^{\omega_0/2})$ . That is, the bound  $w > P^{1/\alpha}$  for a fat-tree is analogous to the value  $D > D_2 = \lceil \alpha/(\alpha-1) \rceil$  for a torus. We do not obtain any tighter bounds for direct linear algebra or the  $N$ -body problem. This analysis suggests that for those computations, a fat-tree with the cheapest choice of  $w = P^{2/3}$  is sufficient to avoid contention becoming the communication bottleneck, though we point out that tighter contention bounds may exist that contradict this conjecture.

### Programs Referencing Arrays.

The analysis for programs referencing arrays follows that of direct linear algebra, matrix-multiplication, and  $N$ -body computations. Replacing  $F$  with  $|\mathcal{Z}|$ ,  $N$  with  $|\mathcal{Z}|^{1/s_{\text{HBL}}}$ , and  $\alpha$  with  $s_{\text{HBL}}$ , we obtain the contention bounds  $W_{\text{link}}^{\text{f-t}}(P, M, N) = \Omega(|\mathcal{Z}|/(w(MP)^{s_{\text{HBL}}-1}))$  and  $W_{\text{link}}^{\text{f-t}}(P, N) = \Omega(|\mathcal{Z}|^{1/s_{\text{HBL}}}/w)$ . The dominance of the memory-independent bound follows from Claim A.1, and we can also conclude that choosing  $w > P^{1/s_{\text{HBL}}}$  guarantees that these contention bounds do not dominate the per-processor bounds.

### FFT/Sorting.

To obtain a contention bound for the FFT on a fat-tree, we combine (via Theorem 2.3) the per-processor bounds given in Theorems 3.6 and 3.7 with the aggregate bandwidth defined by Equation (8). This yields a memory-independent contention bound of

$$\begin{aligned} W_{\text{link}}^{\text{f-t}}(P, N) &= \max_{1 \leq i \leq \log P} \Omega \left( \frac{\frac{N \log N}{2^i \log(N/2^i)}}{\min\{P/2^i, w/2^{2i/3}\}} \right) \\ &= \max_{1 \leq i \leq \log P} \Omega \left( \frac{N \log N}{P \log(N/2^i)} + \frac{N \log N}{w \log(N/2^i) 2^{i/3}} \right). \end{aligned}$$

Again, this function is maximized at either endpoint, so to obtain a new bound we choose  $i = 1$ , which evaluates to  $W_{\text{link}}^{\text{f-t}}(P, N) = \Omega(N/w)$ . The memory-dependent contention bound can be derived similarly, evaluating to

$$W_{\text{link}}^{\text{f-t}}(P, M, N) = \Omega(N \log N / (w \log(MP))).$$

As in the per-processor case, the memory-independent contention bound dominates the memory-dependent contention bound due to the assumption of  $N < MP$ . However, either of the two memory-independent bounds may dominate; the contention bound dominates when  $w < P(1 - (\log P / \log N))$ . That is, for sufficiently small  $N$  ( $N$  close to  $P$ ), contention is not a bottleneck even for  $w = P^{2/3}$ ; for sufficiently large  $N$  ( $N = P^C$  for some constant  $C$ ), then contention will bottleneck the computation unless  $w$  is chosen to be  $\Omega(P)$ .

## 5. DISCUSSION AND FUTURE RESEARCH

### Is it always about the bisection bandwidth?

For the algorithms discussed in this paper on torus and fat-tree networks, the contention lower bound, when it differs from the per-processor bound, is maximized for  $t = P/2$ ;

that is, the contention bound corresponds to a network bisection cut. Is this always the case, or do we expect to have combinations of algorithms and machines where contention bounds dominate, but the constricting cut is not balanced? An example might be when  $h_s(G_{Net})$  is minimized by some  $2 < s_0 < \frac{|G_{Net}|}{2}$ ; for example, two networks of processors, a large and a small one, where each of them is well connected, but the connection between the large and the small one is narrow (e.g., two racks with one router each, connected with a narrow link one to the other, where the racks contain different numbers of processors).

### *Applicability*

Immediate applications of the technique include combinations of other networks (e.g. dragonfly networks [31]) and other classes of algorithms for which per-processor lower bounds are known. A network may have expansion sufficiently large to preclude the use of our contention bound on a given computation, yet the contention may still dominate the communication cost. This calls for further study on how well computations and networks match each other. Similar questions have been addressed by Leiserson and others [9, 25, 33], and had a large impact on the design of supercomputer networks. In particular, a parallel computer that uses a fat tree communication network can simulate any other routing network, at the cost of at most polylogarithmic slowdown. [12] considers similar issues from a VLSI view.

### *Contention-Efficient Algorithms*

Some parallel algorithms attaining per-processor communication lower bounds have also been tuned to particular topologies (cf. [41] for classical matrix multiplication on 3D torus). Algorithmic analysis of the contention costs will likely show that the contention bounds for these and related computations are attainable. Many other algorithms are communication optimal when all-to-all connectivity is assumed, but their performance on other topologies has not yet been studied. Additionally, a parallel algorithm can have

limited knowledge of the physical layout of the processors it runs on. Even with such knowledge, node failure or migration may replace a well-positioned node with a poorly-positioned one. Are there algorithms that attain the communication lower bounds for any realistic network graph, either by automatic tuning or topology-oblivious tools?

## **Acknowledgments**

We thank Guy Kindler for pointing us to [15]. Research partially funded by DARPA Award Number HR0011-12-2-0016, the Center for Future Architecture Research, a member of STARnet, a Semiconductor Research Corporation program sponsored by MARCO and DARPA, and ASPIRE Lab industrial sponsors and affiliates Intel, Google, Nokia, NVIDIA, Oracle, MathWorks and Samsung. Research is also supported by U.S. Department of Energy Office of Science, Office of Advanced Scientific Computing Research, Applied Mathematics program DE-SC0004938, DE-SC0005136, DE-SC0003959, DE-SC0008700, DE-SC0008699, DE-SC0010200 and DOE AC02-05CH11231. Research is supported by grants 1878/14, 1901/14, and 1045/09 from the Israel Science Foundation (founded by the Israel Academy of Sciences and Humanities). This research is supported by grant 3-10891 from the Ministry of Science and Technology, Israel. This paper is supported by the Intel Collaborative Research Institute for Computational Intelligence (ICRI-CI). This research was supported by a grant from the United States-Israel Binational Science Foundation (BSF), Jerusalem, Israel. This work was supported by the HUJI Cyber Security Research Center in conjunction with the Israel National Cyber Bureau in the Prime Minister's Office. Research is supported by the Einstein Foundation. Research is supported by the Minerva Foundation. This research was supported in part by an appointment to the Sandia National Laboratories Truman Fellowship in National Security Science and Engineering, sponsored by Sandia Corporation (a wholly owned subsidiary of Lockheed Martin Corporation) as Operator of Sandia National Laboratories under its U.S. Department of Energy Contract No. DE-AC04-94AL85000.

## 6. REFERENCES

- [1] A. Aggarwal, A. K. Chandra, and M. Snir. Communication complexity of PRAMs. *Theoretical Computer Science*, 71(1):3–28, 1990.
- [2] Y. Ajima, S. Sumimoto, and T. Shimizu. Tofu: A 6D mesh/torus interconnect for exascale computers. *Computer*, 42(11):36–40, Nov 2009.
- [3] G. Ballard. *Avoiding Communication in Dense Linear Algebra*. PhD thesis, EECS Department, University of California, Berkeley, Aug 2013.
- [4] G. Ballard, A. Buluç, J. Demmel, L. Grigori, B. Lipshitz, O. Schwartz, and S. Toledo. Communication optimal parallel multiplication of sparse random matrices. In *SPAA '13: Proceedings of the 25rd ACM Symposium on Parallelism in Algorithms and Architectures*, 2013.
- [5] G. Ballard, J. Demmel, O. Holtz, B. Lipshitz, and O. Schwartz. Brief announcement: strong scaling of matrix multiplication algorithms and memory-independent communication lower bounds. In *Proceedings of the 24th ACM Symposium on Parallelism in Algorithms and Architectures*, SPAA '12, pages 77–79, New York, NY, USA, 2012. ACM.
- [6] G. Ballard, J. Demmel, O. Holtz, B. Lipshitz, and O. Schwartz. Graph expansion analysis for communication costs of fast rectangular matrix multiplication. In G. Even and D. Rawitz, editors, *Design and Analysis of Algorithms*, volume 7659 of *Lecture Notes in Computer Science*, pages 13–36. Springer Berlin Heidelberg, 2012.
- [7] G. Ballard, J. Demmel, O. Holtz, and O. Schwartz. Minimizing communication in numerical linear algebra. *SIAM Journal on Matrix Analysis and Applications*, 32(3):866–901, 2011.
- [8] G. Ballard, J. Demmel, O. Holtz, and O. Schwartz. Graph expansion and communication costs of fast matrix multiplication. *Journal of the ACM*, 59(6):32:1–32:23, Dec. 2012.
- [9] P. Bay and G. Bilardi. Deterministic on-line routing on area-universal networks. In *Proceedings of the 31st Annual Symposium on the Foundations of Computer Science (FOCS)*, pages 297–306, 1990.
- [10] J. Berntsen. Communication efficient matrix multiplication on hypercubes. *Parallel Computing*, 12(3):335 – 342, 1989.
- [11] G. Bilardi and L. De Stefani. The I/O complexity of strassen’s matrix multiplication with recomputation. *arXiv preprint arXiv:1605.02224*, 2016.
- [12] G. Bilardi and F. P. Preparata. Area-time lower-bound techniques with applications to sorting. *Algorithmica*, 1(1-4):65–91, 1986.
- [13] G. Bilardi, M. Scquizzato, and F. Silvestri. A lower bound technique for communication on BSP with application to the FFT. In *Euro-Par 2012 Parallel Processing*, pages 676–687. Springer, 2012.
- [14] L. S. Blackford, J. Choi, A. Cleary, E. D’Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R. C. Whaley. *ScaLAPACK Users’ Guide*. SIAM, Philadelphia, PA, USA, May 1997. Also available from <http://www.netlib.org/scalapack/>.
- [15] B. Bollobás and I. Leader. Edge-isoperimetric inequalities in the grid. *Combinatorica*, 11(4):299–314, 1991.
- [16] L. Cannon. *A cellular computer to implement the Kalman filter algorithm*. PhD thesis, Montana State University, Bozeman, MN, 1969.
- [17] E. Chan, M. Heimlich, A. Purkayastha, and R. Van De Geijn. Collective communication: theory, practice, and experience. *Concurrency and Computation: Practice and Experience*, 19(13):1749–1783, 2007.
- [18] D. Chen, N. Easley, P. Heidelberger, R. Senger, Y. Sugawara, S. Kumar, V. Salapura, D. Satterfield, B. Steinmacher-Burow, and J. Parker. The IBM BG/Q Interconnection Fabric. *IEEE Micro*, 32(1):32–43, 2012.
- [19] M. Christ, J. Demmel, N. Knight, T. Scanlon, and K. Yelick. Communication lower bounds and optimal algorithms for programs that reference arrays - part 1. Technical Report UCB/EECS-2013-61, EECS Department, University of California, Berkeley, 2013.
- [20] Cray. Cray XK7 brochure, 2011.
- [21] J. Dongarra. Visit to the National University for Defense Technology Changsha, China, June 2013.
- [22] M. Driscoll, E. Georganas, P. Koanantakool, E. Solomonik, and K. Yelick. A communication-optimal N-body algorithm for direct interactions. In *Proceedings of IPDPS '13*, 2013.
- [23] F. L. Gall. Powers of tensors and fast matrix multiplication. In *Proceedings of the 39th International Symposium on Symbolic and Algebraic Computation (ISSAC 2014)*, pages 296–303, 2014.
- [24] M. T. Goodrich. Communication-efficient parallel sorting. *SIAM J. Computing*, 29(2):416–432, 1999.
- [25] R. I. Greenberg and C. E. Leiserson. Randomized routing on fat-tress. In *Proceedings of the 26th Annual Symposium on the Foundations of Computer Science (FOCS)*, pages 241–249, 1985.
- [26] J. W. Hong and H. T. Kung. I/O complexity: The red-blue pebble game. In *Proc. 14th STOC*, pages 326–333, New York, NY, USA, 1981. ACM.
- [27] S. Hoory, N. Linial, and A. Wigderson. Expander graphs and their applications. *Bulletin of the AMS*, 43(4):439–561, 2006.
- [28] IBM Blue Gene Team. Overview of the IBM Blue Gene/P project. *IBM Journal of Research and Development*, 52(1.2):199–220, Jan 2008.
- [29] D. Irony, S. Toledo, and A. Tiskin. Communication lower bounds for distributed-memory matrix multiplication. *J. Parallel Distrib. Comput.*, 64(9):1017–1026, 2004.
- [30] J. Jeffers, J. Jeffers, and J. Reinders. *Intel Xeon Phi Coprocessor High Performance Programming*. Elsevier Science & Technology Books, 2013.
- [31] J. Kim, W. Dally, S. Scott, and D. Abts. Technology-driven, highly-scalable dragonfly topology. In *Computer Architecture, 2008. ISCA '08. 35th International Symposium on*, pages 77–88, June 2008.
- [32] N. Knight, E. Carson, and J. Demmel. Exploiting data sparsity in parallel matrix powers computations. In *Proceedings of PPAM '13*, Lecture Notes in Computer Science. Springer (to appear), 2013.
- [33] C. E. Leiserson. Fat-trees: Universal networks for hardware-efficient supercomputing. *IEEE*

- Transactions on Computers*, C-34(10):892–901, 1985.
- [34] J. H. Lindsey. Assignment of numbers to vertices. *The American Mathematical Monthly*, 71(5):508–516, 1964.
- [35] W. McColl and A. Tiskin. Memory-efficient matrix multiplication in the BSP model. *Algorithmica*, 24(3-4):287–297, 1999.
- [36] H. Meuer, E. Strohmaier, J. Dongarra, and H. Simon. Top500 Supercomputer Sites, 2016. [www.top500.org](http://www.top500.org).
- [37] J. E. Savage. Extending the Hong-Kung model to memory hierarchies. In *COCOON*, pages 270–281, 1995.
- [38] A. Schönhage. Partial and total matrix multiplication. *SIAM J. Computing*, 10(3):434–455, 1981.
- [39] J. Scott, O. Holtz, and O. Schwartz. Matrix multiplication I/O-complexity by path routing. In *Proceedings of the 27th ACM symposium on Parallelism in Algorithms and Architectures*, pages 35–45. ACM, 2015.
- [40] M. Scquizzato and F. Silvestri. Communication lower bounds for distributed-memory computations. In E. W. Mayr and N. Portier, editors, *31st International Symposium on Theoretical Aspects of Computer Science (STACS 2014)*, volume 25 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 627–638, Dagstuhl, Germany, 2014. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
- [41] E. Solomonik and J. Demmel. Communication-optimal parallel 2.5D matrix multiplication and LU factorization algorithms. In E. Jeannot, R. Namyst, and J. Roman, editors, *Euro-Par 2011 Parallel Processing*, volume 6853, pages 90–109. Springer Berlin Heidelberg, 2011.
- [42] V. Strassen. Gaussian elimination is not optimal. *Numer. Math.*, 13:354–356, 1969.
- [43] V. Strassen. Relative bilinear complexity and matrix multiplication. *Journal für die reine und angewandte Mathematik (Crelles Journal)*, 1987(375–376):406–443, 1987.

## APPENDIX

### A. DOMINANCE OF MEMORY-INDEPENDENT CONTENTION BOUND

CLAIM A.1. *Let Alg be an algorithm performing a computation of the form given by (2) on  $P$  processors, each with local memory of size  $M$ , and assume the input data is initially evenly distributed across processors. Then,*

$$\frac{|\mathcal{Z}|^{1/s_{\text{HBL}}}}{M} \leq \sum_{j=1}^m \frac{|\phi_j(\mathcal{Z})|}{M}.$$

*As the minimum number of processors required to hold the problem is the right-hand side of this inequality, we conclude that the memory-independent contention bound dominates the memory-dependent contention bound as the two bounds are equivalent when  $P = |\mathcal{Z}|^{1/s_{\text{HBL}}}/M$ .*

PROOF. To begin a proof, the HBL bound discussed in Christ et al. [19], states (with certain assumptions) that

$$|\mathcal{Z}| \leq \prod_{j=1}^m |\phi_j(\mathcal{Z})|^{s_j}.$$

To detail an argument from Section 2 of [19], we present several greater upper bounds on  $|\mathcal{Z}|$  that will allow us to demonstrate the desired result:

$$\begin{aligned} |\mathcal{Z}| &\leq \prod_{j=1}^m |\phi_j(\mathcal{Z})|^{s_j} \leq \prod_{j=1}^m \left( \max_{j=1}^m |\phi_j(\mathcal{Z})| \right)^{s_j} \\ &= \left( \max_{j=1}^m |\phi_j(\mathcal{Z})| \right)^{\sum_{j=1}^m s_j} = \left( \max_{j=1}^m |\phi_j(\mathcal{Z})| \right)^{s_{\text{HBL}}} \end{aligned}$$

As  $\max_{j=1}^m x_j \leq \sum_{j=1}^m x_j$  if all  $x_j \geq 0$ ,

$$|\mathcal{Z}| \leq \left( \max_{j=1}^m |\phi_j(\mathcal{Z})| \right)^{s_{\text{HBL}}} \leq \left( \sum_{j=1}^m |\phi_j(\mathcal{Z})| \right)^{s_{\text{HBL}}}$$

which proves the desired inequality if we take  $s_{\text{HBL}}$ th root of both sides and divide by  $M$ .  $\square$

### B. DERIVATION OF EXPRESSIONS IN FIGURE 2

- **Equivalence point for per-processor bounds**

We set the per-processor bounds equal to each other, and solve for  $P$ :

$$\frac{F}{PM^{\alpha-1}} = \Theta \left( \frac{N}{P^{1/\alpha}} \right)$$

$$P = \Theta \left( \frac{F}{NM^{\alpha-1}} \right)^{\alpha/(\alpha-1)}$$

- **Equivalence point for contention bounds**

We set the contention bounds equal to each other, and solve for  $P$ :

$$\frac{F}{P^{\alpha-1/D} M^{\alpha-1}} = \Theta \left( \frac{N}{P^{1-1/D}} \right)$$

$$P = \Theta \left( \frac{F}{NM^{\alpha-1}} \right)^{1/(\alpha-1)}$$

- **Equivalence point for the memory-dependent per-processor and memory-independent contention bounds**

We set the memory-dependent per-processor and memory-independent contention bounds equal to each other, and solve for  $P$  as a function of  $D$ :

$$\frac{F}{PM^{\alpha-1}} = \Theta \left( \frac{N}{P^{1-1/D}} \right)$$

$$P = \Theta \left( \frac{F}{NM^{\alpha-1}} \right)^D$$

### C. HYPERCUBE NETWORK CONTENTION LOWER BOUNDS

Our analysis reveals no meaningful new bounds for direct linear algebra, fast matrix multiplication,  $O(N^2)N$ -body and the FFT. We show this here in detail.

*Linear algebra, fast matrix multiplication and  $N$ -body.*

We show the results for direct linear algebra only. For fast matrix multiplication and  $N$ -body, the analysis is essentially the same. Recall that for direct linear algebra,

$$W_{\text{proc}}(P, M, N) = \Omega \left( \frac{N^{\frac{3}{2}}}{PM^{\frac{1}{2}}} \right) \text{ and } W_{\text{proc}}(P, N) = \Omega \left( \frac{N}{P^{\frac{2}{3}}} \right).$$

Assigning these per-processor lower bounds into Theorem 2.3, we have the following contention bounds:

$$W_{\text{link}}(P, M, N) \geq \Omega \left( \max_{t \in T} \frac{N^{\frac{3}{2}}}{PM^{\frac{1}{2}} t^{\frac{1}{2}} (\log P - \log t)} \right)$$

$$W_{\text{link}}(P, N) \geq \Omega \left( \max_{t \in T} \frac{N}{P^{\frac{2}{3}} t^{\frac{1}{3}} (\log P - \log t)} \right)$$

Deriving both expressions by  $t$  shows that  $W_{\text{link}}(P, M, N)$  and  $W_{\text{link}}(P, N)$  are maximized by  $t = 1$ . It is therefore impossible for our contention lower bounds to dominate over the per-processor lower bounds.

*FFT and sorting.* Recall that  $W_{\text{proc}}(P, M, N) = \Omega \left( \frac{N \log N}{P \log M} \right)$

and  $W_{\text{proc}}(P, N) = \Omega \left( \frac{N \log N}{P \log \frac{N}{P}} \right)$ . As above, using Theorem 2.3 to obtain the following contention lower bounds:

$$W_{\text{link}}(P, M, N) \geq \Omega \left( \max_{t \in T} \left( \frac{N \log N}{P (\log M + \log t) (\log P - \log t)} \right) \right)$$

$$W_{\text{link}}(P, N) \geq \Omega \left( \max_{t \in T} \left( \frac{N \log N}{P (\log N + \log t - \log P) (\log P - \log t)} \right) \right)$$

For the memory-dependent bound, considering the derivative by  $t$  of the bounds for  $W_{\text{link}}(P, M, N)$  reveals the maximal contention bound is at  $t = \frac{P}{2}$ . That is:

$$W_{\text{link}}(P, M, N) \geq \Omega \left( \frac{N \log N}{P \log(MP)} \right)$$

and so the memory-dependent per-processor communication bound dominates over the memory-dependent contention bound.

Similarly for the memory-independent bound, considering the derivative by  $t$  of the bounds for  $W_{link}(P, N)$  reveals the maximal contention bound is at either  $t = 1$  or at  $t = \frac{P}{2}$ . For  $t = 1$ , the maximal bound is  $\Omega\left(\frac{N \log N}{P \log P (\log N - \log P)}\right)$ . For  $t = \frac{P}{2}$ , the maximal bound is  $\Omega\left(\frac{N}{P}\right)$ . The memory-independent per-processor communication bound dominates over both potential memory-independent contention bounds.

Therefore, our analysis reveals no meaningful new bounds for FFT when computed over a hypercube network.