#### Communication Lower Bounds for Bilinear Algorithms

Edgar Solomonik

L P. N A @CS@Illinois

Department of Computer Science University of Illinois at Urbana-Champaign

SIAM Annual Meeting

# Laboratory for Parallel Numerical Algorithms

Group research foci

- tensor computations
- numerical linear algebra
- numerical optimization
- communication cost analysis\*
- paralallel algorithms
- quantum simulation
- graph algorithms



http://lpna.cs.illinois.edu

#### Overview

- General framework for communication lower bounds based on rank structure of bilinear algorithm
  - Communication lower bounds of bilinear algorithms for symmetric tensor contractions
    E.S., James Demmel, and Torsten Hoefler (SISC 2021, arXiv:1707.04618)
- Communication lower bounds for nested bilinear algorithms
  - Communication lower bounds for nested bilinear algorithms Caleb Ju, Yifan Zhang, and E.S. (arXiv:2107.09834)



#### **Bilinear Algorithms**

A bilinear algorithm<sup>1</sup>  $\Lambda = (\mathbf{F}^{(A)}, \mathbf{F}^{(B)}, \mathbf{F}^{(C)})$  computes

$$\boldsymbol{c} = \boldsymbol{F}^{(C)}[(\boldsymbol{F}^{(A)T}\boldsymbol{a}) \odot (\boldsymbol{F}^{(B)T}\boldsymbol{b})]$$

where a and b are inputs and  $\odot$  is the Hadamard (pointwise) product



<sup>1</sup>Victor Pan, How to multiply matrices faster. 1984.

LPNA

Lower Bounds for Bilinear Algorithms

### Example: Convolution Bilinear Algorithm

Given  $oldsymbol{a},oldsymbol{b}\in\mathbb{R}^n$ , compute

$$c_i = \sum_k a_k b_{i-k}$$

• Equivalent to product of polynomials, compute via interpolation

$$\boldsymbol{c} = \boldsymbol{V}^{-1}[(\boldsymbol{V}^T \boldsymbol{a}) \odot (\boldsymbol{V}^T \boldsymbol{b})]$$

where V is a Vandermonde matrix

- Known as Toom-Cook algorithm; with complex roots of unity as nodes, V is DFT matrix, enabling use of FFT algorithm
- Other bilinear algorithms and problem variants possible<sup>1</sup>

<sup>&</sup>lt;sup>1</sup>Caleb Ju and E.S., Derivation and analysis of fast bilinear algorithms for convolution. SIAM Review 2020.

#### Bilinear Algorithms as Tensor Factorizations

A bilinear algorithm corresponds to a CP tensor decomposition

$$\begin{split} c_{i} &= \sum_{r=1}^{R} f_{ir}^{(C)} \bigg( \sum_{j} f_{jr}^{(A)} a_{j} \bigg) \bigg( \sum_{k} f_{kr}^{(B)} b_{k} \bigg) \\ &= \sum_{j} \sum_{k} \bigg( \sum_{r=1}^{R} f_{ir}^{(C)} f_{jr}^{(A)} f_{kr}^{(B)} \bigg) a_{j} b_{k} \\ &= \sum_{j} \sum_{k} t_{ijk} a_{j} b_{k} \quad \text{where} \quad t_{ijk} = \sum_{r=1}^{R} f_{ir}^{(C)} f_{jr}^{(A)} f_{kr}^{(B)} \end{split}$$

For example, the convolution problem is defined by

$$t_{ijk} = 1$$
 iff  $i = j + k$ 

so that

$$c_i = \sum_{j,k} t_{ijk} a_j b_k = \sum_k a_k b_{i-k}$$

LPNA

Lower Bounds for Bilinear Algorithms

#### **Bilinear Problems**

- convolution of *n*-dimensional vectors
  - optimal rank<sup>1</sup>:  $R = \Theta(n)$
- matrix multiplication of  $n \times n$  matrices
  - best currently known algorithm<sup>2</sup>:  $R = O(n^{2.37286})$

• symmetric tensor contraction with  $n \times \cdots \times n$  tensors

• for symmetry preserving algorithms<sup>3</sup>, e.g., AB + BA via

$$c_{ij} = \sum_{k} (a_{ij} + a_{ik} + a_{jk})(b_{ij} + b_{ik} + b_{jk}) - \dots$$

 $R = \frac{n^d}{d!} + O(n^{d-1}),$  where d is total number of indices in contraction

<sup>1</sup>Ke Ye and Lek-Heng Lim, FoCM 2018.

 $^2 {\sf Josh}$  Alman and V. V. Williams, A refined laser method and faster matrix multiplication. SODA 2020.

<sup>3</sup>E.S. and James Demmel, Fast bilinear algorithms for symmetric tensor contractions. CMAM 2021.

# Communication Lower Bounds for Bilinear Algorithms

How much communication does a bilinear algorithm require?

- bilinear algorithm specifies only products, not order of additions, need to consider various partial sums
- measure two types of communication:
  - ${\scriptstyle \bullet }$  vertical words communicated between memory and cache of size H
  - $\bullet\,$  horizontal words communicated by some processor of p
- assume load-balanced input/output and no replication of work



## Execution DAG

Model computation by dependency graph (execution DAG)

- vertex input/output/computed value
- edge input to operation computing a value
- parallelization graph partitioning
- cache-blocking pebbling game on DAG

Communication lower bounds characterize minimal comm. cost of one or a family of execution DAG

- various orderings of a sum permitted by expressing it as a hyperedge<sup>1</sup>
- columns of bilinear algorithm matrices encode such hyperedges



<sup>&</sup>lt;sup>1</sup>E.S., E. Carson, N. Knight, Trade-offs between synchronization, communication, and computation in parallel linear algebra computations. TOPC 2016.

#### Some Known Lower Bounds

Matrix multiplication of  $n \times n$  matrices

- classical  $O(n^3)$  algorithm requires  $Q = \Theta(n^3/\sqrt{H})$  vertical<sup>1</sup> and  $W = \Theta(n/p^{2/3})$  horizontal communication<sup>2</sup>
- Strassen's  $O(n^{\log_2 7})$  algorithm requires  $Q = \Theta(n^{\log_2 7}H/H^{\log_4 7})$  vertical and  $W = \Theta(n^2/p^{1/\log_4 7})$  horizontal communication<sup>3</sup>
- bounds for many numerical linear algebra algorithms (LU, QR, SVD, eigensolvers) depend on above<sup>4</sup>

Convolution of n-dimensional vectors

• FFT requires  $Q = \Theta(n \log n / \log H)$  vertical<sup>1</sup> and  $W = \Theta(n \log n / \log(n/p))$  horizontal comm.<sup>5</sup>

Given computational complexity f(n), these bounds share a common form

$$Q = \Theta(f(n) \cdot H/\tilde{f}(H)) \qquad W = \Theta(\tilde{f}^{-1}(f(n)/p))$$

 $<sup>^1</sup>$ Jia-Wei Hong and H.T. Kung, I/O complexity: The red-blue pebble game. STOC 1981

<sup>&</sup>lt;sup>2</sup>D. Irony, S. Toledo, A. Tiskin, Communication lower bounds for distributed-memory matrix multiplication. JPDC 2004.

<sup>&</sup>lt;sup>3</sup>G. Ballard, J. Demmel, O. Holtz, O. Schwartz, Graph expansion and comm. costs of fast mat. mult., JACM 2013.

<sup>&</sup>lt;sup>4</sup>G. Ballard, J. Demmel, O. Holtz, O. Schwartz. Minimizing communication in numerical linear algebra. SIMAX 2011.

<sup>&</sup>lt;sup>5</sup>A. Aggarwal, A.K. Chandra, M. Snir, Communication complexity of PRAMs, 1990.

# Lower Bounds for Bilinear Algorithms

Communication lower bounds based on the rank  ${\it R}$  of the bilinear algorithm

- ${\ensuremath{\bullet}}$  lower bound information needed to compute any subset of the R products
  - consider inputs and outputs symmetrically
  - products are columns in bilinear matrices, so seek lower bound on rank of submatrices
- $\bullet$  above expansion bound  $\tilde{f}$  yields horiz. and vertical comm. lower bounds
  - lower bound on amount of useful work with  ${\cal O}({\cal H})$  input/output words yields vertical communication lower bound

$$Q = \Omega(R \cdot H/\tilde{f}(H))$$

- some processor computes at least R/p products, use  $\bar{f}$  to bound information it must send/receive

$$W = \Omega(\tilde{f}^{-1}(R/p) - (\#\mathsf{inputs} + \#\mathsf{outputs})/p)$$

### Bilinear Algorithm Expansion Bound



• generally, for bilinear algorithm  $\Lambda$  seek function  $\mathcal{E}_{\Lambda}(x,y,z)$ , s.t., for any subset of  $k\leq R$  products

$$k \leq \mathcal{E}_{\Lambda}\left(\mathrm{rank}(\boldsymbol{F}_{\mathsf{sub}}^{(A)}), \mathrm{rank}(\boldsymbol{F}_{\mathsf{sub}}^{(B)}), \mathrm{rank}(\boldsymbol{F}_{\mathsf{sub}}^{(C)})\right)$$

where  $F_{sub}^{(A)}$  is the submatrix of  $F^{(A)}$  obtained from extracting k columns, and similar for respective columns of  $F^{(B)}$  and  $F^{(C)}$ 

•  $\mathcal{E}_{\Lambda}(x,y,z)$  describes the amount of useful work we can do with m information

$$\tilde{f}(m) = \max_{x+y+z \le m} \mathcal{E}_{\Lambda}(x, y, z)$$

# Rank Expansion of Matrices

• Let  $\sigma_{F^{(X)}}(k)$  be the minimum rank of a subset k columns of  ${\pmb F}^{(X)}$  , then

$$\mathcal{E}_{\Lambda}(x, y, z) \le \max_{x+y+z \le m} \min(\sigma_{F^{(A)}}^{-1}(x), \sigma_{F^{(B)}}^{-1}(y), \sigma_{F^{(C)}}^{-1}(z))$$

- $\bullet$  so we may simply consider the rank expansion  $\sigma(k)$  of each bilinear algorithm matrix
- stronger bounds are sometimes possible by considering all three matrices simultaneously (since the subset of columns must be matching among the 3 matrices)
  - e.g., for classical matrix multiplication, this is needed to use the surface-to-volume-ratio argument (Loomis-Whitney inequality)
  - in other cases, it suffices to consider  $\sigma$  for each bilinear matrix separately, e.g., compare

$$c_{ij} \leftarrow a_{ik}b_{kj} \quad \text{to} \quad \{c_{ij}, c_{ik}, c_{jk}\} \leftarrow (a_{ij} + a_{ik} + a_{jk})(b_{ij} + b_{ik} + b_{jk})$$

# Rank Expansion of Matrices

We are then interested in characterizing by  $\sigma(k)$ , the minimum rank of a subset of k columns of a given matrix C

- $\sigma_C(k) = k$  if C is full rank
- the first kink in σ<sub>C</sub> is determined by the Kruskal rank of C (largest k such that any subset of k columns of C are linearly independent)
- if C is sparse (e.g., banded), lower bounds on  $\sigma_C(k)$  may be derived based on sparsity
- nested (recursive) bilinear algorithms, (e.g., Strassen's algorithm) result in matrices with Kronecker product structure

$$C = A \otimes B$$

since rank(C) = rank(A) rank(B), can we characterize σ<sub>C</sub> in terms of σ<sub>A</sub> and σ<sub>B</sub>?

# Rank Expansion of Kronecker-Product-Structured Matrices

• suppose we select  $k_A$  columns of  $A(A_1)$  and  $k_B$  columns of  $B(B_1)$ , then

$$\operatorname{rank}(\boldsymbol{A}_1 \otimes \boldsymbol{B}_1) = \operatorname{rank}(\boldsymbol{A}_1) \operatorname{rank}(\boldsymbol{B}_1)$$

consequently, we expect that

$$\sigma_{A\otimes B}(k) \ge \min_{k_A k_B \ge k} \sigma_A(k_A) \sigma_B(k_B)$$

• we formalize this intuition into a proof and quantify when it holds<sup>1</sup>



 $<sup>^1 \</sup>mbox{Caleb}$  Ju, Yifan Zhang, and E.S., Communication lower bounds for nested bilinear algorithms. arXiv 2021.

# Proof of Rank Expansion Bound under Kronecker Product

We associate a given subset of columns of  $m{C}$  with grid points  $m{a}_i \otimes m{b}_j$ 



then we transform and shrink the set of grid points in a way that can only decrease rank of the associated submatrix of  ${\boldsymbol C}$ 



# Proof of Rank Expansion Bound under Kronecker Product

We then prove that the compactified grid has a minimal basis (subset of grid points spanning original set) with a Kronecker product structure



Having identified a nicely-structured minimal basis, we can upper bound the total number of columns it can span using  $\sigma_A$  and  $\sigma_B$ 



# Proof of Rank Expansion Bound under Kronecker Product

Finally, we can transform the shape of the geometric region (set of columns spanned by basis) by removing 'stairs' without decreasing area



the optimal region for general concave  $\sigma_A$  and  $\sigma_B$  is L-shaped

LPNA

Lower Bounds for Bilinear Algorithms

## General Kronecker-Product Rank-Expansion Bound

#### Theorem

Let 
$$f(x) = \sigma_A^{-1}(x)$$
,  $g(x) = \sigma_B^{-1}(x)$ . If

$$\hat{f}(x) = \frac{f(r_A) - f(r_A - x)}{xf'(r_A - x)}, \ \hat{g}(x) = \frac{g(r_B) - g(r_B - x)}{xg'(r_B - x)},$$

are increasing on  $(0, r_A)$  and  $(0, r_B)$ , respectively, then

$$\sigma_C(k) = \min_{\substack{k_A \in [d_A, n_A], \ k_B \in [d_B, n_B] \\ k_A \in \{ d_A, n_A \} \text{ or } k_B \in \{ d_B, n_B \}, \\ k_A k_B \ge k}} \sigma_A(k_A) \cdot \sigma_B(k_B)$$

is a rank expansion lower bound of  $C = A \otimes B$ .

The conditions are satisfied e.g., when f and g are monomial or exponential

# Applications and Conclusion

We apply the nested rank expansion bounds to derive communciation lower bounds for various nested bilinear algorithms

Previous (V)	Previous (H)	Our Work (V)	Our Work (H)
$n^{\log_2(7)}$	$n^2$	$n^{\log_2(7)}$	$n^{\log_3(7)}$
$\frac{1}{H^{\log_4(7)-1}}$	$\overline{p^{\log_7(4)}}$	$\overline{H^{\log_2(3)-1}}$	$\overline{p^{\log_3(2)}}$
$\frac{n^{\log_k(2k-1)}}{H^{\log_k(2k-1)-1}}$	-	Match Prev	$\frac{n}{n^{\log_{2k-1}(k)}}$
	$\frac{n^{\log_2(7)}}{H^{\log_4(7)-1}} \frac{n^{\log_2(7)}}{n^{\log_k(2k-1)}}$	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $

For contraction among symmetric tensors of order s+v and v+t over v modes

- $\bullet\,$  asymptotically more communication is required by symmetry-preserving techniques when s,t,v are unequal
- comm. lower bounds for symmetry preserving algorithms on some contractions of partially-symmetric tensors follow from nested bound

For further details, see arXiv:2107.09834

#### Backup slides

#### Symmetric Matrix Vector Product

- Consider computing  $oldsymbol{c} = oldsymbol{A} oldsymbol{b}$  with  $oldsymbol{A} = oldsymbol{A}^T$ 
  - Typically requires  $n^2$  multiplications since  $a_{ij}b_j \neq a_{ji}b_i$  and  $n^2-n$  additions
  - Instead can compute

$$v_i = \sum_{j=1}^{i-1} u_{ij} + \sum_{j=i+1}^n u_{ji}$$
 where  $u_{ij} = a_{ij}(b_i + b_j)$ 

using n(n-1)/2 multiplications (since we only need  $u_{ij}$  for i > j) and about  $3n^2/2$  additions, then

$$c_i = (2a_{ii} - \sum_{j=1}^n a_{ij})b_i + v_i$$

using  $n \mbox{ more multiplications and } n^2 \mbox{ additions }$ 

- Beneficial when multiplying elements of  ${m A}$  and  ${m b}$  costs more than addition
- ${\ensuremath{\, \rm o}}$  This technique yields a bilinear algorithm with rank n(n+1)/2

LPNA

Lower Bounds for Bilinear Algorithms

July 23, 2021 23 / 21

# Partially-Symmetric Tensor Times Matrix (TTM)

- Can use symmetric mat-vec algorithm to accelerate TTM with partially symmetric tensor from  $2n^4$  operations to  $(3/2)n^4 + O(n^3)$ 
  - Given  $A \in \mathbb{R}^{n \times n \times n}$  with symmetry  $a_{ijk} = a_{jik}$  and  $B \in \mathbb{R}^{n \times n}$ , we compute

$$c_{ikl} = \sum_{j} a_{ijk} b_{jl}$$

• We can think of this as a set of symmetric matrix-vector products  ${m c}^{(k,l)}-{m A}^{(k)}{m b}^{(l)}$ 

and apply the fast bilinear algorithm

$$\begin{aligned} v_{ikl} &= \sum_{j=1}^{i-1} u_{ijkl} + \sum_{j=i+1}^{n} u_{ijkl} \quad \text{where} \quad u_{ijkl} = a_{ijk} (b_{il} + b_{jl}) \\ c_{ikl} &= (2a_{iik} - \sum_{j=1}^{n} a_{ijk}) b_{il} + v_{ikl} \end{aligned}$$

using about  $n^4/2$  multiplications and  $n^4 + O(n^3)$  additions (need only  $n^3$  distinct sums of elements of B) to compute  $\mathcal{V}$ , then  $O(n^3)$ operations to get  $\mathcal{C}$  from  $\mathcal{V}$ Lower Bounds for Bilinear Algorithms July 23, 2021 24/21

# Permutational Symmetry in Tensor Contractions



New contraction algorithms reduce cost via permutational symmetry<sup>1</sup>

- Symmetry is hard to use in contraction e.g.  $m{y}=Am{x}$  with A symmetric
- For contraction of order s + v and v + t tensors to produce an order s + t tensor, previously known approaches reduce cost by s!t!v!
- New algorithm reduces number of *products* by  $\omega$ ! where  $\omega = s + t + v$ , leads to same reduction in *cost* for partially-symmetric contractions

$$\boldsymbol{C} = \boldsymbol{A}\boldsymbol{B} + \boldsymbol{B}\boldsymbol{A} \Rightarrow c_{ij} = \sum_{k} [(a_{ij} + a_{ik} + a_{jk}) \cdot (b_{ij} + b_{ik} + b_{jk})] - \dots$$

<sup>1</sup>E.S, J. Demmel, CMAM 2020

LPNA