

# CS 598: Provably Efficient Algorithms for Numerical and Combinatorial Problems

## Part 2: Algorithm Representation

Edgar Solomonik

University of Illinois at Urbana-Champaign

# Straight Line Programs

- ▶ Often, we want to quantify the efficiency of an algorithm that solves any problem of size  $n$  in  $f(n)$  iterations, i.e., it is a *straight line program*
  - ▶ *Numerical algorithms often fit this form, including for matrix multiplication, convolution, matrix factorizations, direct solvers,  $k$ -iterations of an sparse iterative method*
  - ▶ *Programs that have branches or conditional loop bounds may not be described a straight line program*
- ▶ The completed execution of a program for a particular problem may always be described by a straight line program
  - ▶ *Ultimately a sequence of instructions are executed*
  - ▶ *These instructions may be independent and so can be reordered, thus we have a set of operations which need to be executed according to a partial order*

# Algorithms as Directed Acyclic Graphs

- ▶ A directed acyclic graph (DAG) describes a straight line program in terms of elementwise operations (addition, multiplication, etc.)
  - ▶ *This DAG has a vertex for each scalar value input to or computed within the program*
  - ▶ *Computed values are vertices with in-degree one or two in the DAG*
- ▶ Assuming an algorithm is a straight line program, we may ask questions regarding parallelism and communication cost
  - ▶ *Depth of DAG gives lower bound on parallel execution time*
  - ▶ *Expansion properties describe communication cost*

# Schedules of an Algorithm

- ▶ A *schedule* assigns the vertices of a straight-line program to instructional units and manages associated communication
  - ▶ *Schedule depends on architectural model*
    - ▶ *Sequential with explicitly managed bounded fast memory (cache/register file), schedule needs to provide reads/writes/discards*
    - ▶ *Sequential with implicitly managed bounded fast memory (cache/register file) via a given caching protocol*
    - ▶ *Parallel shared memory with no fast memory (PRAM), variants regarding how to handle concurrent reads/writes to same locations*
    - ▶ *Distributed-memory parallel (BSP/ $\alpha$ - $\beta$ ), schedule may manage initial data layout, communication and synchronization*

# Parameterization of Algorithms

- ▶ Oftentimes, we may want to parameterize the algorithm (and not just the schedule) depending on the architecture
  - ▶ *For example, we may use a flat reduction tree on a single processor as opposed to a binary reduction tree with many processors*
  - ▶ *Reorganization of DAG may minimize expansion with respect to a parameter controlling subset size (which may be correlated with fast memory size)*
- ▶ An algorithm may also be designed to be *oblivious* to a parameter, i.e., to minimize execution time for any choice of a particular parameter
  - ▶ *This notion is most important for fast memory (cache) size, with corresponding algorithms referred to as *cache-oblivious**
  - ▶ *Other examples include network-oblivious algorithms*

## Matrix Multiplication as a DAG

Lets consider the matrix multiplication problem: compute  $C$  such that  $C = AB$  with  $A, B, C \in \mathbb{R}^{n \times n}$

- ▶ Loop-nest can be used to describe algorithm/DAG (for  $i$ , for  $j$ , for  $k$ ,  $c_{ij}^{(k)} = c_{ij}^{(k-1)} + a_{ik}b_{kj}$  with  $c_{ij}^{(0)} = 0$  and  $c_{ij} = c_{ij}^{(n)}$ )
  - ▶ *As stated, loop nest describes flat reduction tree with depth  $O(n)$*
  - ▶ *Using associativity of addition can obtain DAGs with depth  $O(\log(n))$*
- ▶ Recursive formulation describes another algorithm/DAG

$$\begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$

- ▶ *Performing eight block products recursively yields cost*

$$T(n) = 8T(n/2) + O(n^2) = O(n^3)$$

- ▶ *Recursive calls can be done in parallel and algorithm is cache-oblivious*

# Family of Classical Matrix Multiplication Algorithms

- ▶ The nested-loop and recursive formulations are two instances of a family of classical matrix multiplication algorithms
  - ▶ *They yield different partial orders*
  - ▶ *Both compute products  $a_{ik}b_{kj}$  and sum along  $k$*
  - ▶ *Consequently, associativity of addition implies they are equivalent in exact arithmetic*
- ▶ Can describe family of DAGs as a hypergraph
  - ▶ *Define hyperedges  $h_{ij}^{(C)} = (\{a_{ik}b_{kj} : k \in \{1, \dots, n\}\}, c_{ij})$  to denote reduction/sum*
  - ▶ *Could further abstract DAG edges  $(a_{ik}, a_{ik}b_{kj})$  as ‘broadcast’ hyperedges  $h_{ik}^{(A)} = (a_{ik}, \{a_{ik}b_{kj} : k \in \{1, \dots, n\}\})$  and  $h_{kj}^{(B)}$  similarly*
  - ▶ *Expansion properties of this hypergraph imply expansion properties for DAGs arising from any summation order*

## Surface Area to Volume Ratio in Hypergraphs

- ▶ We can analyze the hypergraph to determine communication cost bounds
  - ▶ *Would like to consider any partial order equivalent via associativity*
  - ▶ *For simplicity, often want to assume no intermediate values are recomputed*
- ▶ The *Loomis-Whitney* is a *volumetric inequality* that provides a way to bound expansion
  - ▶ *Consider a set of products  $S \subseteq \{a_{ik}b_{kj} : i, j, k \in \{1, \dots, n\}\}$*
  - ▶ *Let  $H_S = H_S^{(A)} \cup H_S^{(B)} \cup H_S^{(C)}$  be the set of hyperedges adjacent to  $S$ ,*

$$|S| \leq \left( |H_S^{(A)}| \cdot |H_S^{(B)}| \cdot |H_S^{(C)}| \right)^{1/2}$$

- ▶ *Consequently, we can derive communication lower bounds by inferring that e.g.,*

$$|H_S| \geq (1/3)^{1/3} |S|^{2/3}$$



## Compression and Recomputation

- ▶ Our previous discussion of communication assumed that each hypergraph edge requires communication of a matrix entry
  - ▶ *Perhaps it is possible to communicate less information by transforming a set of entries, e.g. taking linear combinations there of (sending  $a_{ij} + a_{i'j'}$  instead of both individually)?*
  - ▶ *Could lower bound information flow and compression via linear combinations by considering rank*
- ▶ A method that computes bilinear products  $a_{ik}b_{kj}$  may take arbitrary linear combinations of entries of  $A$ ,  $B$ , or partial sums for  $C$ 
  - ▶ *Given a vector of input linear combinations  $s^{(A)}$  of entries of  $A$ ,  $s^{(B)}$  of entries of  $B$ , output some set of linear combinations  $s^{(C)}$  of products  $a_{ik}b_{kj}$ , we have that for some  $A^{(S)}$ ,  $B^{(S)}$ ,  $C^{(S)}$ ,*

$$s^{(C)} = C^{(S)} \left[ (A^{(S)T} s^{(A)}) \odot (B^{(S)T} s^{(B)}) \right]$$

- ▶ *To lower bound the dimensions of  $s^{(A)}$ ,  $s^{(B)}$ , and  $s^{(C)}$ , need to relate  $A^{(S)}$ ,  $B^{(S)}$ , and  $C^{(S)}$  to overall computation*

# Bilinear Algorithms

A bilinear algorithm (V. Pan, 1984)  $\Lambda = (\mathbf{F}^{(A)}, \mathbf{F}^{(B)}, \mathbf{F}^{(C)})$  computes

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

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

$$\begin{bmatrix} \mathbf{c} \end{bmatrix} = \begin{bmatrix} \begin{matrix} x & & x & & x & & x \\ x & x & & x & x & & x \\ & x & & x & & x & x \\ x & & x & & & x & x \\ & x & x & & x & x & x \\ x & & x & & x & x & x \\ & x & x & & x & x & x \end{matrix} \end{bmatrix} \left[ \left( \begin{bmatrix} \begin{matrix} x & x & x & x & x & x & x \\ x & x & x & x & x & x & x \\ & x & & & x & x & x \\ x & x & x & x & x & x & x \\ & x & x & x & x & x & x \\ x & x & x & x & x & x & x \end{matrix} \end{bmatrix}^T \begin{bmatrix} \mathbf{a} \end{bmatrix} \right) \circ \left( \begin{bmatrix} \begin{matrix} x & x & x & x & x & x & x \\ x & & x & x & x & x & x \\ & x & x & x & x & x & x \\ x & x & & x & x & x & x \\ & x & x & x & x & x & x \\ x & x & x & x & x & x & x \end{matrix} \end{bmatrix}^T \begin{bmatrix} \mathbf{b} \end{bmatrix} \right) \right]$$

# Bilinear Algorithms as Tensor Factorizations

- ▶ A bilinear algorithm corresponds to a CP tensor decomposition

$$\begin{aligned}c_i &= \sum_{r=1}^R f_{ir}^{(C)} \left( \sum_j f_{jr}^{(A)} a_j \right) \left( \sum_k f_{kr}^{(B)} b_k \right) \\&= \sum_j \sum_k \left( \sum_{r=1}^R f_{ir}^{(C)} f_{jr}^{(A)} f_{kr}^{(B)} \right) a_j b_k \\&= \sum_j \sum_k t_{ijk} a_j b_k \quad \text{where} \quad t_{ijk} = \sum_{r=1}^R f_{ir}^{(C)} f_{jr}^{(A)} f_{kr}^{(B)}\end{aligned}$$

- ▶ For multiplication of  $n \times n$  matrices, we can define a *matrix multiplication tensor* and consider algorithms with various bilinear rank
  - ▶  $\mathbf{T}$  is  $n^2 \times n^2 \times n^2$
  - ▶ Classical algorithm has rank  $R = n^3$
  - ▶ Strassen's algorithm has rank  $R \approx n^{\log_2(7)}$

# Strassen's Algorithm

$$\text{Strassen's algorithm } \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \cdot \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$

$$M_1 = (A_{11} + A_{22}) \cdot (B_{11} + B_{22})$$

$$C_{11} = M_1 + M_4 - M_5 + M_7$$

$$M_2 = (A_{21} + A_{22}) \cdot B_{11}$$

$$C_{21} = M_2 + M_4$$

$$M_3 = A_{11} \cdot (B_{12} - B_{22})$$

$$C_{12} = M_3 + M_5$$

$$M_4 = A_{22} \cdot (B_{21} - B_{11})$$

$$C_{22} = M_1 - M_2 + M_3 + M_6$$

$$M_5 = (A_{11} + A_{12}) \cdot B_{22}$$

$$M_6 = (A_{21} - A_{11}) \cdot (B_{11} + B_{12})$$

$$M_7 = (A_{12} - A_{22}) \cdot (B_{21} + B_{22})$$

By performing the nested calls recursively, Strassen's algorithm achieves cost,

$$T(n) = 7T(n/2) + O(n^2) = O(7^{\log_2 n}) = O(n^{\log_2 7})$$

## Expansion in Bilinear Algorithms

- ▶ The communication cost of a bilinear algorithm depends on the amount of data needed to compute subsets of the bilinear products.
  - ▶ *A schedule may involve computations of parts of the bilinear algorithm, of the form,  $\Lambda = (\mathbf{F}^{(A)}, \mathbf{F}^{(B)}, \mathbf{F}^{(C)})$ ,  $\Lambda_{\text{sub}} \subseteq \Lambda$  if for some projection matrix  $\mathbf{P}$ ,*

$$\Lambda_{\text{sub}} = (\mathbf{F}^{(A)} \mathbf{P}, \mathbf{F}^{(B)} \mathbf{P}, \mathbf{F}^{(C)} \mathbf{P}).$$

- ▶ *The projection matrix extracts  $\# \text{cols}(\mathbf{P})$  columns of each matrix.*
- ▶ A bilinear algorithm  $\Lambda$  can be associated expansion bound  $\mathcal{E}_{\Lambda} : \mathbb{N}^3 \rightarrow \mathbb{N}$ 
  - ▶ *Expansion bounds holds if for all*

$$\Lambda_{\text{sub}} := (\mathbf{F}_{\text{sub}}^{(A)}, \mathbf{F}_{\text{sub}}^{(B)}, \mathbf{F}_{\text{sub}}^{(C)}) \subseteq \Lambda$$

*we have*  $\text{rank}(\Lambda_{\text{sub}}) \leq \mathcal{E}_{\Lambda} \left( \text{rank}(\mathbf{F}_{\text{sub}}^{(A)}), \text{rank}(\mathbf{F}_{\text{sub}}^{(B)}), \text{rank}(\mathbf{F}_{\text{sub}}^{(C)}) \right)$

- ▶ *For matrix mult., Loomis-Whitney inequality  $\rightarrow \mathcal{E}_{\text{MM}}(x, y, z) = \sqrt{xyz}$*