# Parallel Numerical Algorithms

Chapter 3 – Dense Linear Systems
Section 3.1 – Vector and Matrix Products

## Michael T. Heath and Edgar Solomonik

Department of Computer Science
University of Illinois at Urbana-Champaign

## CS 554 / CSE 512

# Outline

1. **BLAS**

2. **Inner Product**

3. **Outer Product**

4. **Matrix-Vector Product**

5. **Matrix-Matrix Product**

## Basic Linear Algebra Subprograms

- Basic Linear Algebra Subprograms (BLAS) are building blocks for many other matrix computations

- BLAS encapsulate basic operations on vectors and matrices so they can be optimized for particular computer architecture while high-level routines that call them remain portable

- BLAS offer good opportunities for optimizing utilization of memory hierarchy

- Generic BLAS are available from netlib, and many computer vendors provide custom versions optimized for their particular systems

## Examples of BLAS

| Level | Work | Examples | Function |
|-------|------|----------|----------|
| 1 | $\mathcal{O}(n)$ | daxpy | Scalar $\times$ vector $+$ vector |
|   |      | ddot | Inner product |
|   |      | dnrm2 | Euclidean vector norm |
| 2 | $\mathcal{O}(n^2)$ | dgemv | Matrix-vector product |
|   |      | dtrsv | Triangular solve |
|   |      | dger | Outer-product |
| 3 | $\mathcal{O}(n^3)$ | dgemm | Matrix-matrix product |
|   |      | dtrsm | Multiple triangular solves |
|   |      | dsyrk | Symmetric rank-$k$ update |

$$\underbrace{\gamma_1}_{\text{BLAS 1 effective sec/flop}} > \underbrace{\gamma_2}_{\text{BLAS 2 effective sec/flop}} \gg \underbrace{\gamma_3}_{\text{BLAS 3 effective sec/flop}}$$

BLAS
**Inner Product**
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Scalability

## Inner Product

- Inner product of two $n$-vectors $x$ and $y$ given by

$$\boldsymbol{x}^T \boldsymbol{y} = \sum_{i=1}^{n} x_i \, y_i$$

- Computation of inner product requires $n$ multiplications and $n-1$ additions

$$M_1 = \Theta(n), \quad Q_1 = \Theta(n), \quad T_1 = \Theta(\gamma \, n)$$

- Effectively as hard as scalar reduction, can be done via binary or binomial tree summation
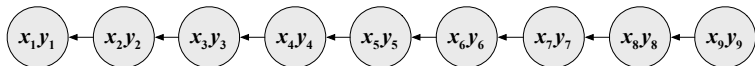
BLAS
**Inner Product**
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Scalability

## Parallel Algorithm

### *Partition*

- For $i = 1, \ldots, n$, fine-grain task $i$ stores $x_i$ and $y_i$, and computes their product $x_i \, y_i$

### *Communicate*

- Sum reduction over $n$ fine-grain tasks

$$\left(x_1 y_1\right) \leftarrow \left(x_2 y_2\right) \leftarrow \left(x_3 y_3\right) \leftarrow \left(x_4 y_4\right) \leftarrow \left(x_5 y_5\right) \leftarrow \left(x_6 y_6\right) \leftarrow \left(x_7 y_7\right) \leftarrow \left(x_8 y_8\right) \leftarrow \left(x_9 y_9\right)$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Scalability

## Fine-Grain Parallel Algorithm

$z_i = x_i y_i$                                           { local scalar product }

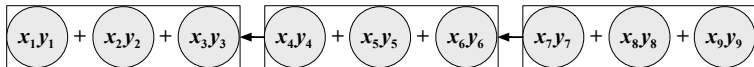reduce $z_i$ across all tasks $i = 1, ..., n$      { sum reduction }

BLAS
**Inner Product**
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

**Parallel Algorithm**
Scalability

# Agglomeration and Mapping

*Agglomerate*

- Combine $k$ components of both $x$ and $y$ to form each coarse-grain task, which computes inner product of these subvectors

- Communication becomes sum reduction over $n/k$ coarse-grain tasks

*Map*

- Assign $(n/k)/p$ coarse-grain tasks to each of $p$ processors, for total of $n/p$ components of $x$ and $y$ per processor

$$\boxed{x_1y_1 + x_2y_2 + x_3y_3} \leftarrow \boxed{x_4y_4 + x_5y_5 + x_6y_6} \leftarrow \boxed{x_7y_7 + x_8y_8 + x_9y_9}$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Scalability

# Coarse-Grain Parallel Algorithm

$z_i = \boldsymbol{x}_{[i]}^T \boldsymbol{y}_{[i]}$  { local inner product }

reduce $z_i$ across all processors $i = 1, ..., p$  { sum reduction }

$\left[ \boldsymbol{x}_{[i]} - \text{subvector of } \boldsymbol{x} \text{ assigned to processor } i \right]$

BLAS
**Inner Product**
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Scalability

## Performance

The parallel costs $(L_p, W_p, F_p)$ for the inner product are given by

- Computational cost $F_p = \Theta(n/p)$ regardless of network
- The latency and bandwidth costs depend on network:
  - 1-D mesh: $L_p, W_p = \Theta(p)$
  - 2-D mesh: $L_p, W_p = \Theta(\sqrt{p})$
  - hypercube: $L_p, W_p = \Theta(\log p)$
- For a hypercube or fully-connected network time is

$$T_p = \alpha L_p + \beta W_p + \gamma F_p = \Theta(\alpha \log(p) + \gamma n/p)$$

- Efficiency and scaling are the same as for binary tree sum

BLAS
Inner Product
**Outer Product**
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Outer Product

- Outer product of two $n$-vectors $x$ and $y$ is $n \times n$ matrix $Z = xy^T$ whose $(i, j)$ entry $z_{ij} = x_i y_j$

- For example,

$$
\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}^T = \begin{bmatrix} x_1y_1 & x_1y_2 & x_1y_3 \\ x_2y_1 & x_2y_2 & x_2y_3 \\ x_3y_1 & x_3y_2 & x_3y_3 \end{bmatrix}
$$

- Computation of outer product requires $n^2$ multiplications,

$$
M_1 = \Theta(n^2), \quad Q_1 = \Theta(n^2), \quad T_1 = \Theta(\gamma n^2)
$$

(in this case, we should treat $M_1$ as output size or define the problem as in the BLAS: $Z = Z_{\text{input}} + xy^T$)

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
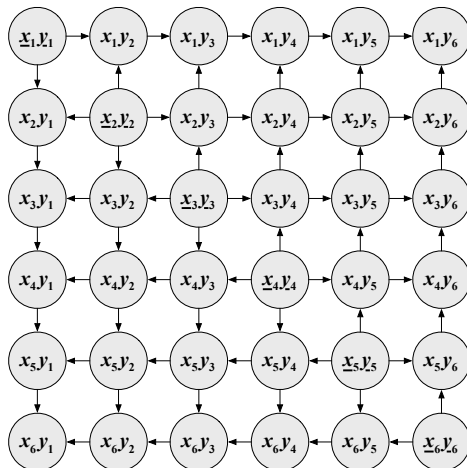Agglomeration Schemes
Scalability

## Parallel Algorithm

*Partition*

- For $i, j = 1, \ldots, n$, fine-grain task $(i, j)$ computes and stores $z_{ij} = x_i\, y_j$, yielding 2-D array of $n^2$ fine-grain tasks

- Assuming no replication of data, at most $2n$ fine-grain tasks store components of $x$ and $y$, say either
  - for some $j$, task $(i, j)$ stores $x_i$ and task $(j, i)$ stores $y_i$, or
  - task $(i, i)$ stores both $x_i$ and $y_i$, $i = 1, \ldots, n$

*Communicate*

- For $i = 1, \ldots, n$, task that stores $x_i$ broadcasts it to all other tasks in $i$th task row

- For $j = 1, \ldots, n$, task that stores $y_j$ broadcasts it to all other tasks in $j$th task column

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Fine-Grain Tasks and Communication

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Fine-Grain Parallel Algorithm

broadcast $x_i$ to tasks $(i, k), \ k = 1, \dots, n$     { horizontal broadcast }

broadcast $y_j$ to tasks $(k, j), \ k = 1, \dots, n$     { vertical broadcast }

$z_{ij} = x_i y_j$     { local scalar product }

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Agglomeration

## *Agglomerate*

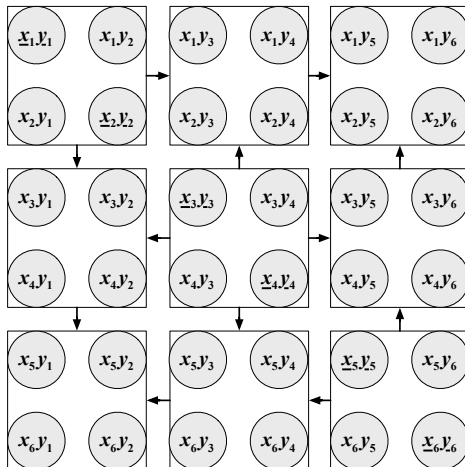With $n \times n$ array of fine-grain tasks, natural strategies are

- 2-D: Combine $k \times k$ subarray of fine-grain tasks to form each coarse-grain task, yielding $(n/k)^2$ coarse-grain tasks

- 1-D column: Combine $n$ fine-grain tasks in each column into coarse-grain task, yielding $n$ coarse-grain tasks

- 1-D row: Combine $n$ fine-grain tasks in each row into coarse-grain task, yielding $n$ coarse-grain tasks

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## 2-D Agglomeration

- Each task that stores portion of $x$ must broadcast its subvector to all other tasks in its task row

- Each task that stores portion of $y$ must broadcast its subvector to all other tasks in its task column
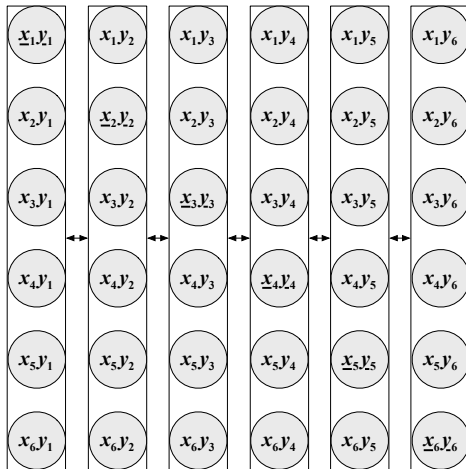
BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## 2-D Agglomeration

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product
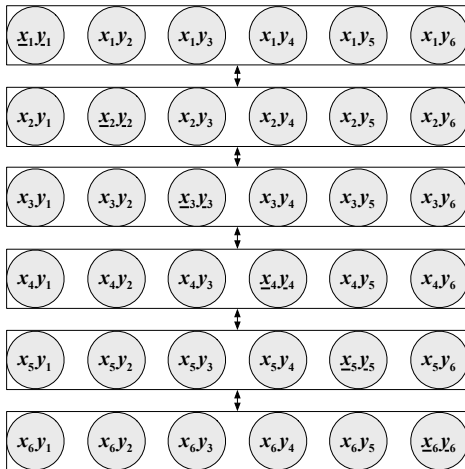
Parallel Algorithm
Agglomeration Schemes
Scalability

# 1-D Agglomeration

- If either $x$ or $y$ stored in one task, then broadcast required to communicate needed values to all other tasks

- If either $x$ or $y$ distributed across tasks, then multinode broadcast required to communicate needed values to other tasks

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# 1-D Column Agglomeration

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# 1-D Row Agglomeration

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

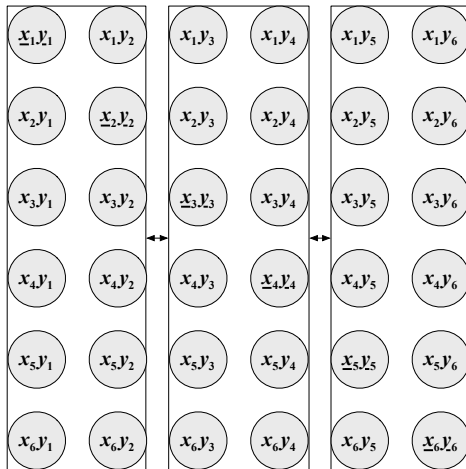Parallel Algorithm
Agglomeration Schemes
Scalability

# Mapping

*Map*

- 2-D: Assign $(n/k)^2/p$ coarse-grain tasks to each of $p$ processors using any desired mapping in each dimension, treating target network as 2-D mesh

- 1-D: Assign $n/p$ coarse-grain tasks to each of $p$ processors using any desired mapping, treating target network as 1-D mesh
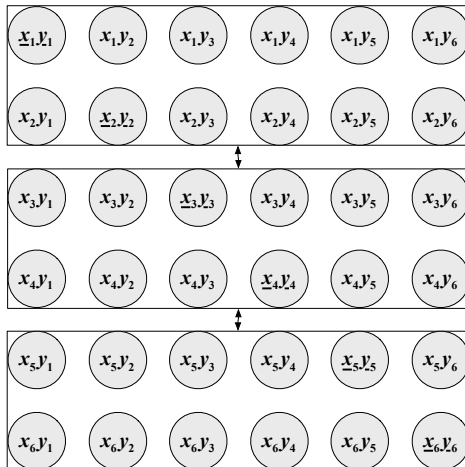
BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## 2-D Agglomeration with Block Mapping

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# 1-D Column Agglomeration with Block Mapping

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# 1-D Row Agglomeration with Block Mapping

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Coarse-Grain Parallel Algorithm

broadcast $x_{[i]}$ to $i$th process row      { horizontal broadcast }

broadcast $y_{[j]}$ to $j$th process column      { vertical broadcast }

$Z_{[i][j]} = x_{[i]} y_{[j]}^T$      { local outer product }

$[\,Z_{[i][j]}$ means submatrix of $Z$ assigned to process $(i, j)$ by mapping $]$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Performance

The parallel costs $(L_p, W_p, F_p)$ for the outer product are

- Computational cost $F_p = \Theta(n^2/p)$ regardless of network
- The latency and bandwidth costs can be derived from the cost of broadcast/allgather
  - 1-D agglomeration: $L_p = \Theta(\log p), W_p = \Theta(n)$
  - 2-D agglomeration: $L_p = \Theta(\log p), W_p = \Theta(n/\sqrt{p})$
- For 1-D agglomeration, execution time is

$$T_p^{\text{1-D}} = T_p^{\text{allgather}}(n) + \Theta(\gamma n^2/p) = \Theta(\alpha \log(p) + \beta n + \gamma n^2/p)$$

- For 2-D agglomeration, execution time is

$$T_p^{\text{2-D}} = 2T_{\sqrt{p}}^{\text{bcast}}(n/\sqrt{p}) + \Theta(\gamma n^2/p) = \Theta(\alpha \log(p) + \beta n/\sqrt{p} + \gamma n^2/p)$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Outer Product Strong Scaling

- 1-D agglomeration is strongly scalable to

$$p_s = \Theta(\min((\gamma/\alpha)n^2/\log((\gamma/\alpha)n^2), (\gamma/\beta)n))$$

processors, since

$$E_{p_s}^{\text{1-D}} = \Theta(1/(1 + (\alpha/\gamma)\log(p_s)p_s/n^2 + (\beta/\gamma)p_s/n))$$

- 2-D agglomeration is strongly scalable to

$$p_s = \Theta(\min((\gamma/\alpha)n^2/\log((\gamma/\alpha)n^2), (\gamma/\beta)^2 n^2))$$

processors, since

$$E_{p_s}^{\text{2-D}} = \Theta(1/(1 + (\alpha/\gamma)\log(p_s)p_s/n^2 + (\beta/\gamma)\sqrt{p_s}/n))$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Outer Product Weak Scaling

- 1-D agglomeration is weakly scalable to

$$p_w = \Theta(\min(2^{(\gamma/\alpha)n^2}, (\gamma/\beta)^2 n^2))$$

processors, since

$$E_{p_w}^{\text{1-D}}(\sqrt{p_w}n) = \Theta(1/(1 + (\alpha/\gamma)\log(p_w)/n^2 + (\beta/\gamma)\sqrt{p_w}/n))$$

- 2-D agglomeration is weakly scalable to

$$p_w = \Theta(2^{(\gamma/\alpha)n^2})$$

processors, since

$$E_{p_w}^{\text{2-D}}(\sqrt{p_w}n) = \Theta(1/(1 + (\alpha/\gamma)\log(p_w)/n^2 + (\beta/\gamma)/n))$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Memory Requirements

- The memory requirements are dominated by storing $Z$, which in practice means the outer-product is a poor primitive (local *flop-to-byte ratio* of $1$)

- If possible, $Z$ should be operated on as it is computed, e.g. if we really need

$$v_i = \sum_j f(x_i y_j) \quad \text{for some scalar function } f$$

- If $Z$ does not need to be stored, vector storage dominates

- Without further modification, 1-D algorithm requires $M_p = \Theta(np)$ overall storage for vector

- Without further modification, 2-D algorithm requires $M_p = \Theta(n\sqrt{p})$ overall storage for vector

BLAS
Inner Product
Outer Product
**Matrix-Vector Product**
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Matrix-Vector Product

- Consider matrix-vector product

$$y = Ax$$

where $A$ is $n \times n$ matrix and $x$ and $y$ are $n$-vectors

- Components of vector $y$ are given by inner products:

$$y_i = \sum_{j=1}^{n} a_{ij} x_j, \quad i = 1, \ldots, n$$

- The sequential memory, work, and time are

$$M_1 = \Theta(n^2), \quad Q_1 = \Theta(n^2), \quad T_1 = \Theta(\gamma n^2)$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
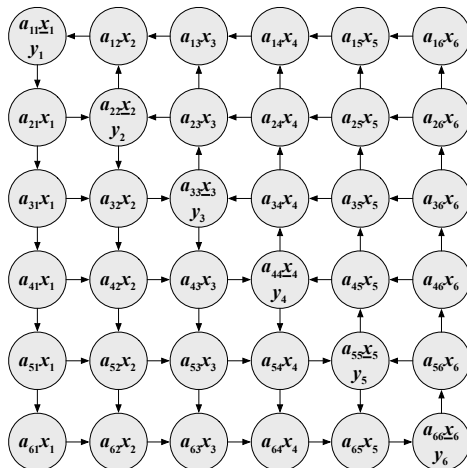Agglomeration Schemes
Scalability

## Parallel Algorithm

*Partition*

- For $i, j = 1, \ldots, n$, fine-grain task $(i, j)$ stores $a_{ij}$ and computes $a_{ij}\, x_j$, yielding 2-D array of $n^2$ fine-grain tasks
- Assuming no replication of data, at most $2n$ fine-grain tasks store components of $x$ and $y$, say either
    - for some $j$, task $(j, i)$ stores $x_i$ and task $(i, j)$ stores $y_i$, or
    - task $(i, i)$ stores both $x_i$ and $y_i$, $i = 1, \ldots, n$

*Communicate*

- For $j = 1, \ldots, n$, task that stores $x_j$ broadcasts it to all other tasks in $j$th task column
- For $i = 1, \ldots, n$, sum reduction over $i$th task row gives $y_i$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Fine-Grain Tasks and Communication

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Fine-Grain Parallel Algorithm

broadcast $x_j$ to tasks $(k, j), \ k = 1, \ldots, n$    { vertical broadcast }

$y_i = a_{ij} x_j$                                         { local scalar product }

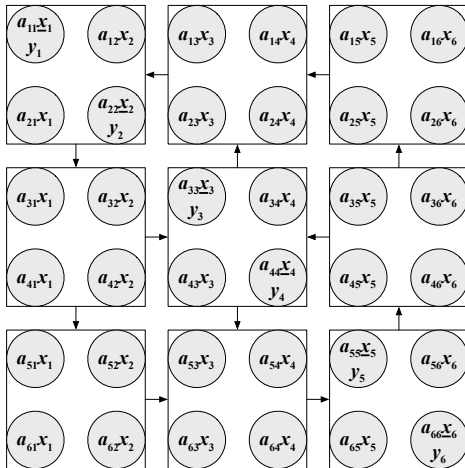reduce $y_i$ across tasks $(i, k), \ k = 1, \ldots, n$    { horizontal sum reduction }

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Agglomeration

### *Agglomerate*

With $n \times n$ array of fine-grain tasks, natural strategies are

- 2-D: Combine $k \times k$ subarray of fine-grain tasks to form each coarse-grain task, yielding $(n/k)^2$ coarse-grain tasks

- 1-D column: Combine $n$ fine-grain tasks in each column into coarse-grain task, yielding $n$ coarse-grain tasks

- 1-D row: Combine $n$ fine-grain tasks in each row into coarse-grain task, yielding $n$ coarse-grain tasks

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## 2-D Agglomeration

- Subvector of $x$ broadcast along each task column

- Each task computes local matrix-vector product of submatrix of $A$ with subvector of $x$

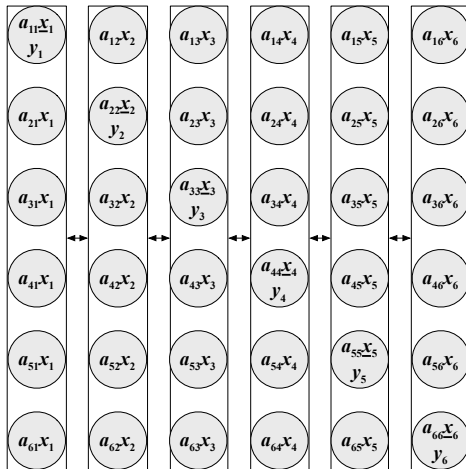- Sum reduction along each task row produces subvector of result $y$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# 2-D Agglomeration

BLAS
Inner Product
Outer Product
**Matrix-Vector Product**
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# 1-D Agglomeration

1-D column agglomeration

- Each task computes product of its component of $x$ times its column of matrix, with no communication required
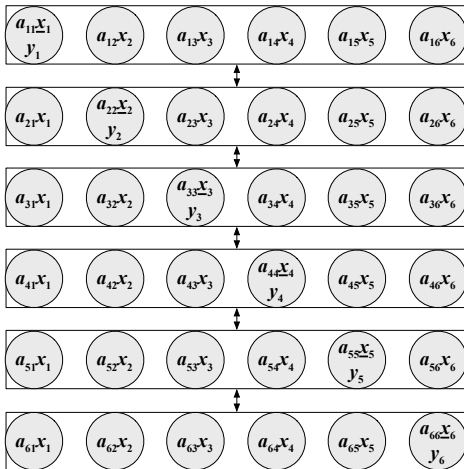
- Sum reduction across tasks then produces $y$

1-D row agglomeration

- If $x$ stored in one task, then broadcast required to communicate needed values to all other tasks

- If $x$ distributed across tasks, then multinode broadcast required to communicate needed values to other tasks

- Each task computes inner product of its row of $A$ with *entire* vector $x$ to produce its component of $y$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# 1-D Column Agglomeration

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# 1-D Row Agglomeration

BLAS
Inner Product
Outer Product
**Matrix-Vector Product**
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## 1-D Agglomeration

Column and row algorithms are dual to each other

- Column algorithm begins with communication-free local vector scaling (`daxpy`) computations combined across processors by a reduction
- Row algorithm begins with broadcast followed by communication-free local inner-product (`ddot`) computations

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
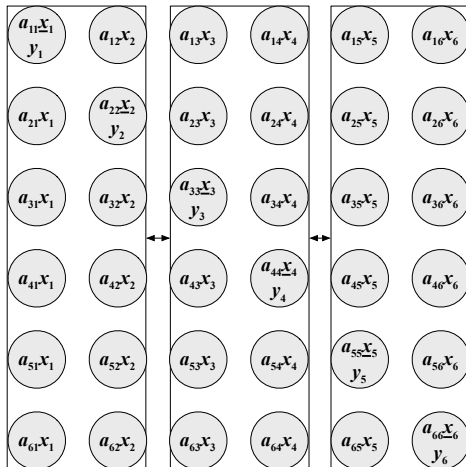Agglomeration Schemes
Scalability

# Mapping

## *Map*

- 2-D: Assign $(n/k)^2/p$ coarse-grain tasks to each of $p$ processes using any desired mapping in each dimension, treating target network as 2-D mesh

- 1-D: Assign $n/p$ coarse-grain tasks to each of $p$ processes using any desired mapping, treating target network as 1-D mesh
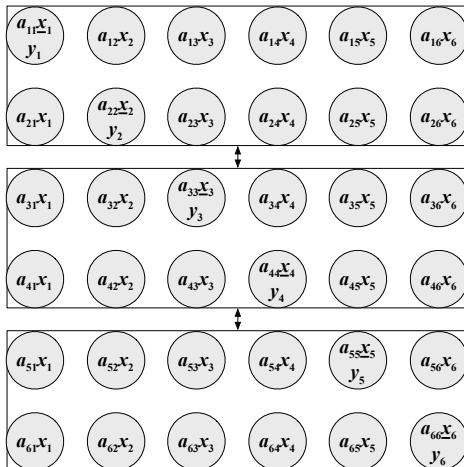
BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## 2-D Agglomeration with Block Mapping

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# 1-D Column Agglomeration with Block Mapping

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# 1-D Row Agglomeration with Block Mapping

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Coarse-Grain Parallel Algorithm

broadcast $x_{[j]}$ to $j$th process column     { vertical broadcast }

$y_{[i]} = A_{[i][j]} x_{[j]}$     { local matrix-vector product }

reduce $y_{[i]}$ across $i$th process row     { horizontal sum reduction }

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Performance

The parallel costs $(L_p, W_p, F_p)$ for the matrix-vector product are

- Computational cost $F_p = \Theta(n^2/p)$ regardless of network
- Communication costs can be derived from the cost of collectives
  - 1-D agglomeration: $L_p = \Theta(\log p), W_p = \Theta(n)$
  - 2-D agglomeration: $L_p = \Theta(\log p), W_p = \Theta(n/\sqrt{p})$
- For 1-D row agglomeration, perform allgather,

$$T_p^{\text{1-D}} = T_p^{\text{allgather}}(n) + \Theta(\gamma n^2/p) = \Theta(\alpha \log(p) + \beta n + \gamma n^2/p)$$

- For 2-D agglomeration, perform broadcast and reduction,

$$T_p^{\text{2-D}} = T_{\sqrt{p}}^{\text{bcast}}(n/\sqrt{p}) + T_{\sqrt{p}}^{\text{reduce}}(n/\sqrt{p}) + \Theta(\gamma n^2/p)$$
$$= \Theta(\alpha \log(p) + \beta n/\sqrt{p} + \gamma n^2/p)$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Matrix-Matrix Product

- Consider matrix-matrix product

$$C = A\,B$$

where $A$, $B$, and result $C$ are $n \times n$ matrices

- Entries of matrix $C$ are given by

$$c_{ij} = \sum_{k=1}^{n} a_{ik}\,b_{kj}, \quad i,j = 1,\ldots,n$$

- Each of $n^2$ entries of $C$ requires $n$ multiply-add operations, so model serial time as

$$T_1 = \gamma\,n^3$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Matrix-Matrix Product

- Matrix-matrix product can be viewed as

  - $n^2$ inner products, or
  - sum of $n$ outer products, or
  - $n$ matrix-vector products
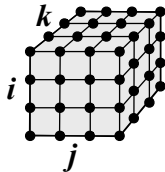
  and each viewpoint yields different algorithm

- One way to derive parallel algorithms for matrix-matrix product is to apply parallel algorithms already developed for inner product, outer product, or matrix-vector product

- However, considering the problem as a whole yields the best algorithms

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Parallel Algorithm

### *Partition*

- For $i, j, k = 1, \ldots, n$, fine-grain task $(i, j, k)$ computes product $a_{ik} b_{kj}$, yielding 3-D array of $n^3$ fine-grain tasks

- Assuming no replication of data, at most $3n^2$ fine-grain tasks store entries of $A$, $B$, or $C$, say task $(i, j, j)$ stores $a_{ij}$, task $(i, j, i)$ stores $b_{ij}$, and task $(i, j, k)$ stores $c_{ij}$ for $i, j = 1, \ldots, n$ and some fixed $k$

- We refer to subsets of tasks along $i$, $j$, and $k$ dimensions as rows, columns, and layers, respectively, so $k$th column of $A$ and $k$th row of $B$ are stored in $k$th layer of tasks

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Parallel Algorithm

*Communicate*

- Broadcast entries of $j$th column of $A$ horizontally along each task row in $j$th layer

- Broadcast entries of $i$th row of $B$ vertically along each task column in $i$th layer

- For $i, j = 1, \ldots, n$, result $c_{ij}$ is given by sum reduction over tasks $(i, j, k)$, $k = 1, \ldots, n$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Fine-Grain Algorithm

broadcast $a_{ik}$ to tasks $(i, q, k)$, $q = 1, \ldots, n$      { horizontal broadcast }

broadcast $b_{kj}$ to tasks $(q, j, k)$, $q = 1, \ldots, n$      { vertical broadcast }

$c_{ij} = a_{ik} b_{kj}$      { local scalar product }

reduce $c_{ij}$ across tasks $(i, j, q)$, $q = 1, \ldots, n$      { lateral sum reduction }

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Agglomeration

### *Agglomerate*

With $n \times n \times n$ array of fine-grain tasks, natural strategies are

- 3-D: Combine $q \times q \times q$ subarray of fine-grain tasks

- 2-D: Combine $q \times q \times n$ subarray of fine-grain tasks, eliminating sum reductions

- 1-D column: Combine $n \times 1 \times n$ subarray of fine-grain tasks, eliminating vertical broadcasts and sum reductions

- 1-D row: Combine $1 \times n \times n$ subarray of fine-grain tasks, eliminating horizontal broadcasts and sum reductions

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Mapping
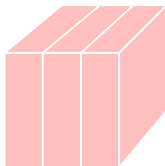
### *Map*

Corresponding mapping strategies are

- 3-D: Assign $(n/q)^3/p$ coarse-grain tasks to each of $p$ processors using any desired mapping in each dimension, treating target network as 3-D mesh

- 2-D: Assign $(n/q)^2/p$ coarse-grain tasks to each of $p$ processors using any desired mapping in each dimension, treating target network as 2-D mesh

- 1-D: Assign $n/p$ coarse-grain tasks to each of $p$ processors using any desired mapping, treating target network as 1-D mesh

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Agglomeration with Block Mapping

1-D row      1-D col      2-D      3-D

BLAS
Inner Product
Outer Product
Matrix-Vector Product
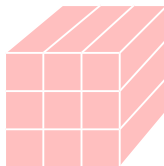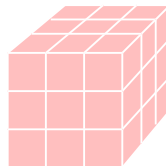Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Coarse-Grain 3-D Parallel Algorithm

broadcast $A_{[i][k]}$ to $i$th processor row { horizontal broadcast }

broadcast $B_{[k][j]}$ to $j$th processor column { vertical broadcast }

$C_{[i][j]} = A_{[i][k]} B_{[k][j]}$ { local matrix product }

reduce $C_{[i][j]}$ across processor layers { lateral sum reduction }

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Agglomeration with Block Mapping



**2-D:**
$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{bmatrix}$$

**1-D column:**
$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{bmatrix}$$

**1-D row:**
$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{bmatrix}$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Coarse-Grain 2-D Parallel Algorithm

allgather $A_{[i][j]}$ in $i$th processor row      { horizontal broadcast }
allgather $B_{[i][j]}$ in $j$th processor column      { vertical broadcast }
$C_{[i][j]} = 0$
**for** $k = 1, \ldots, \sqrt{p}$
     $C_{[i][j]} = C_{[i][j]} + A_{[i][k]} B_{[k][j]}$      { sum local products }
**end**

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# SUMMA Algorithm

- Algorithm just described requires excessive memory, since each process accumulates $\sqrt{p}$ blocks of both $A$ and $B$

- One way to reduce memory requirements is to

  - broadcast blocks of $A$ successively across processor rows
  - broadcast blocks of $B$ successively across processor cols

$C_{[i][j]} = 0$
**for** $k = 1, \ldots, \sqrt{p}$
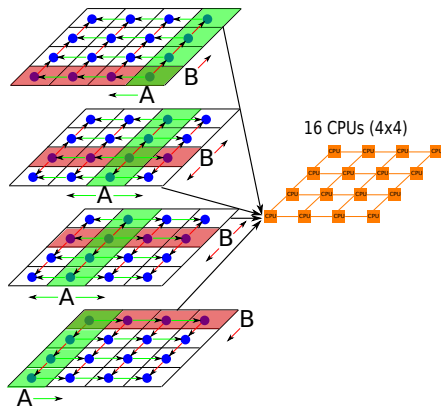  broadcast $A_{[i][k]}$ in $i$th processor row     { horizontal broadcast }
  broadcast $B_{[k][j]}$ in $j$th processor column     { vertical broadcast }
  $C_{[i][j]} = C_{[i][j]} + A_{[i][k]} B_{[k][j]}$     { sum local products }
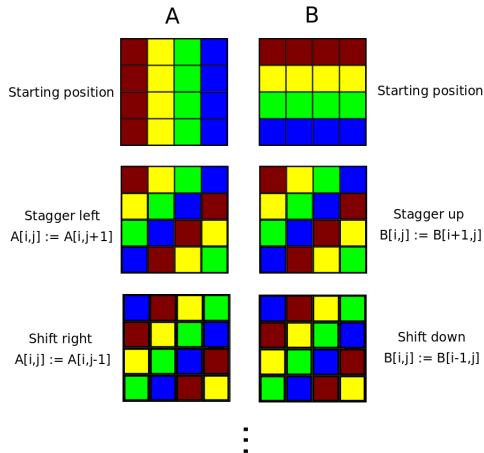**end**

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# SUMMA Algorithm

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Cannon Algorithm

- Another approach, due to Cannon (1969), is to circulate blocks of $B$ vertically and blocks of $A$ horizontally in ring fashion

- Blocks of both matrices must be initially aligned using circular shifts so that correct blocks meet as needed

- Requires less memory than SUMMA and replaces line broadcasts with shifts, lowering the number of messages

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Cannon Algorithm

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Fox Algorithm

- It is possible to mix techniques from SUMMA and Cannon's algorithm:

  - circulate blocks of $B$ in ring fashion vertically along processor columns step by step so that each block of $B$ comes in conjunction with appropriate block of $A$ broadcast at that same step

- This algorithm is due to Fox et al.

- Shifts, especially in Cannon's algorithm, are harder to generalize to nonsquare processor grids than collectives in algorithms like SUMMA

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Execution Time for 3-D Agglomeration

- For 3-D agglomeration, computing each of $p$ blocks $C_{[i][j]}$ requires matrix-matrix product of two $(n/\sqrt[3]{p}) \times (n/\sqrt[3]{p})$ blocks, so

$$F_p = (n/\sqrt[3]{p})^3 = n^3/p$$

- On 3-D mesh, each broadcast or reduction takes time

$$T^{\text{bcast}}_{p^{1/3}}((n/p^{1/3})^2) = O(\alpha \log p + \beta n^2/p^{2/3})$$

- Total time is therefore

$$T_p = O(\alpha \log p + \beta\, n^2/p^{2/3} + \gamma\, n^3/p)$$

- However, overall memory footprint is

$$M_p = \Theta(p(n/p^{1/3})^2) = \Theta(p^{1/3}n^2)$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Execution Time for 2-D Agglomeration

- For 2-D agglomeration (SUMMA), computation of each block $C_{[i][j]}$ requires $\sqrt{p}$ matrix-matrix products of $(n/\sqrt{p}) \times (n/\sqrt{p})$ blocks, so

$$F_p = \sqrt{p}\,(n/\sqrt{p})^3 = n^3/p$$

- For broadcast among rows and columns of processor grid, communication time is

$$2\sqrt{p}\,T_{\sqrt{p}}^{\text{bcast}}(n^2/p) \;\; = \Theta(\alpha\sqrt{p}\log(p) + \beta n^2/\sqrt{p})$$
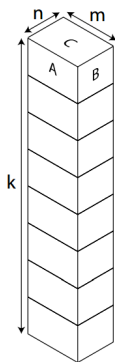
- Total time is therefore

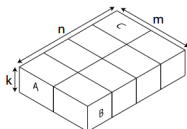$$T_p = O(\alpha\sqrt{p}\log(p) + \beta n^2/\sqrt{p} + \gamma\,n^3/p)$$

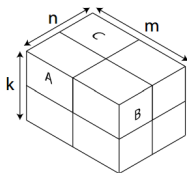- The algorithm is memory-efficient, $M_p = \Theta(n^2)$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# Rectangular Matrix Multiplication



(b) Two large dimensions

(a) One large dimension

(c) Three large dimensions

If $C$ is $m \times n$, $A$ is $m \times k$, and $B$ is $k \times n$, choosing a 3D grid that optimizes volume-to-surface-area ratio yields bandwidth cost...

$$W_p(m,n,k) = O\left( \min_{p_1 p_2 p_3 = p} \left[ \frac{mk}{p_1 p_2} + \frac{kn}{p_1 p_3} + \frac{mn}{p_2 p_3} \right] \right)$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

## Communication vs. Memory Tradeoff

- Communication cost for 2-D algorithms for matrix-matrix product is optimal, assuming no replication of storage
- If explicit replication of storage is allowed, then lower communication volume is possible via 3-D algorithm
- Generally, we assign $n/p_1 \times n/p_2 \times n/p_3$ computation subvolume to each processor
- The largest face of the subvolume gives communication cost, the smallest face gives minimal memory usage
  - can keep smallest face local while successively computing slices of subvolume

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
**Scalability**

# Leveraging Additional Memory in Matrix Multiplication

- Provided $\bar{M}$ total available memory, 2-D and 3-D algorithms achieve bandwidth cost

$$W_p(n, \bar{M}) = \begin{cases} \infty & : \bar{M} < n^2 \\ n^2/\sqrt{p} & : \bar{M} = \Theta(n^2) \\ n^2/p^{2/3} & : \bar{M} = \Theta(n^2 p^{1/3}) \end{cases}$$

- For general $\bar{M}$, possible to pick processor grid to achieve

$$W_p(n, \bar{M}) = O(n^3/(\sqrt{p}\bar{M}^{1/2}) + n^2/p^{2/3})$$

- and an overall execution time of

$$T_p(n, \bar{M}) = O(\alpha(\log p + n^3\sqrt{p}/\bar{M}^{3/2}) + \beta W_p(n, \bar{M}) + \gamma n^3/p)$$

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# References

- R. C. Agarwal, S. M. Balle, F. G. Gustavson, M. Joshi, and P. Palkar, A three-dimensional approach to parallel matrix multiplication, *IBM J. Res. Dev.*, 39:575-582, 1995

- J. Berntsen, Communication efficient matrix multiplication on hypercubes, *Parallel Comput.* 12:335-342, 1989

- J. Demmel, A. Fox, S. Kamil, B. Lipshitz, O. Schwartz, and O. Spillinger, Communication-optimal parallel recursive rectangular matrix multiplication, IPDPS, 2013

- J. W. Demmel, M. T. Heath, and H. A. van der Vorst, Parallel numerical linear algebra, *Acta Numerica* 2:111-197, 1993

BLAS
Inner Product
Outer Product
Matrix-Vector Product
Matrix-Matrix Product

Parallel Algorithm
Agglomeration Schemes
Scalability

# References

- R. Dias da Cunha, A benchmark study based on the parallel computation of the vector outer-product $A = uv^T$ operation, *Concurrency: Practice and Experience* 9:803-819, 1997

- G. C. Fox, S. W. Otto, and A. J. G. Hey, Matrix algorithms on a hypercube I: matrix multiplication, *Parallel Comput.* 4:17-31, 1987

- D. Irony, S. Toledo, and A. Tiskin, Communication lower bounds for distributed-memory matrix multiplication, *J. Parallel Distrib. Comput.* 64:1017-1026, 2004.

- S. L. Johnsson, Communication efficient basic linear algebra computations on hypercube architectures, *J. Parallel Distrib. Comput.* 4(2):133-172, 1987

BLAS
Inner Product
Outer Product
Matrix-Vector Product
**Matrix-Matrix Product**

Parallel Algorithm
Agglomeration Schemes
**Scalability**

## References

- S. L. Johnsson, Minimizing the communication time for matrix multiplication on multiprocessors, *Parallel Comput.* 19:1235-1257, 1993

- B. Lipshitz, Communication-avoiding parallel recursive algorithms for matrix multiplication, Tech. Rept. UCB/EECS-2013-100, University of California at Berkeley, May 2013.

- O. McBryan and E. F. Van de Velde, Matrix and vector operations on hypercube parallel processors, *Parallel Comput.* 5:117-126, 1987

- R. A. Van De Geijn and J. Watts, SUMMA: Scalable universal matrix multiplication algorithm, *Concurrency: Practice and Experience* 9(4):255-274, 1997