
CS 267

Dense Linear Algebra: History and Structure, Parallel Matrix Multiplication

James Demmel

www.cs.berkeley.edu/~demmel

Quick review of earlier lecture

- What do you call
 - A program written in PyGAS, a Global Address Space language based on Python...
 - That uses a Monte Carlo simulation algorithm to approximate π ...
 - That has a race condition, so that it gives you a different funny answer every time you run it?

Monte - π - thon

Outline

- History and motivation
 - What is dense linear algebra?
 - Why minimize communication?
 - Lower bound on communication
- Parallel Matrix-matrix multiplication
 - Attaining the lower bound
- Other Parallel Algorithms (next lecture)

Outline

- History and motivation
 - What is dense linear algebra?
 - Why minimize communication?
 - Lower bound on communication
- Parallel Matrix-matrix multiplication
 - Attaining the lower bound
- Other Parallel Algorithms (next lecture)

Motifs

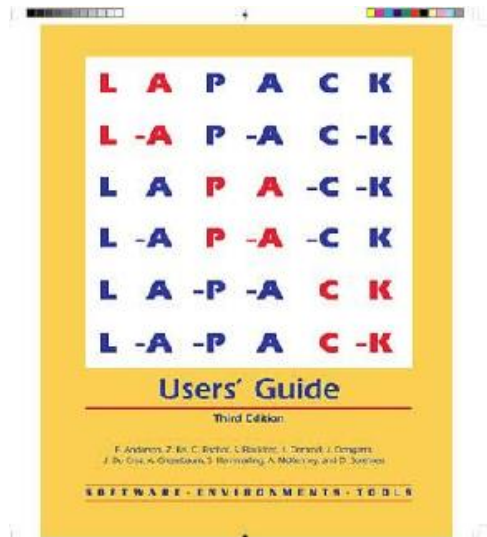
The Motifs (formerly “Dwarfs”) from
“The Berkeley View” (Asanovic et al.)
Motifs form key computational patterns

	Embed	SPEC	DB	Games	ML	HPC
Finite State Mach.	Red	Red	Red	Yellow	Yellow	Light Blue
Circuits	Red	Light Blue	Green	Light Blue	Green	Light Blue
Graph Algorithms	Red	Yellow	Yellow	Yellow	Red	Light Blue
Structured Grid	Red	Red	Light Blue	Yellow	Light Blue	Red
Dense Matrix	Red	Red	Yellow	Red	Red	Red
Sparse matrix	Yellow	Yellow	Light Blue	Red	Red	Red
Spectral (FFT)	Yellow	Light Blue	Light Blue	Yellow	Yellow	Red
Dynamic Prog	Yellow	Light Blue	Red	Light Blue	Red	Light Blue
N-Body	Light Blue	Yellow	Light Blue	Yellow	Light Blue	Red
Backtrack/ B&B	Light Blue	Light Blue	Yellow	Light Blue	Red	Light Blue
Graphical Models	Light Blue	Light Blue	Yellow	Light Blue	Red	Light Blue
Unstructured Grid	Light Blue	Light Blue	Light Blue	Yellow	Yellow	Red

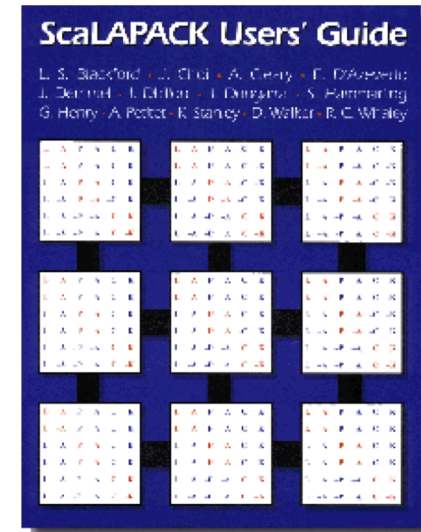
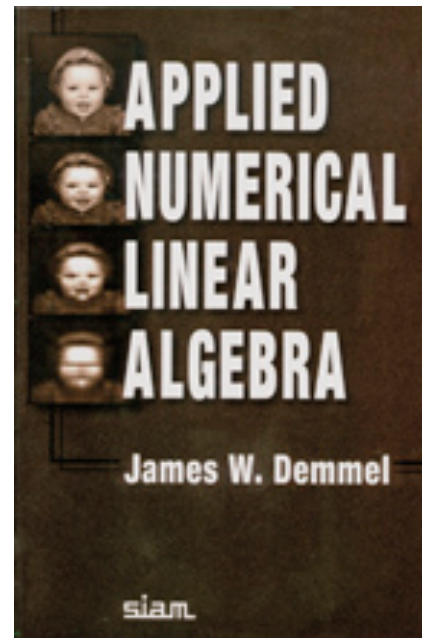
What is dense linear algebra?

- Not just matmul!
- Linear Systems: $Ax=b$
- Least Squares: choose x to minimize $\|Ax-b\|_2$
 - Overdetermined or underdetermined; Unconstrained, constrained, or weighted
- Eigenvalues and vectors of Symmetric Matrices
 - Standard ($Ax = \lambda x$), Generalized ($Ax = \lambda Bx$)
- Eigenvalues and vectors of Unsymmetric matrices
 - Eigenvalues, Schur form, eigenvectors, invariant subspaces
 - Standard, Generalized
- Singular Values and vectors (SVD)
 - Standard, Generalized
- Different matrix structures
 - Real, complex; Symmetric, Hermitian, positive definite; dense, triangular, banded ...
 - 27 types in LAPACK (and growing...)
- Level of detail
 - Simple Driver (" $x=A \backslash b$ ")
 - Expert Drivers with error bounds, extra-precision, other options
 - Lower level routines ("apply certain kind of orthogonal transformation", matmul...)

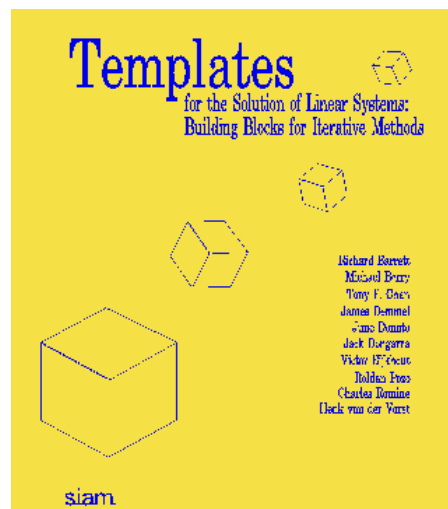
Organizing Linear Algebra – in books



www.netlib.org/lapack

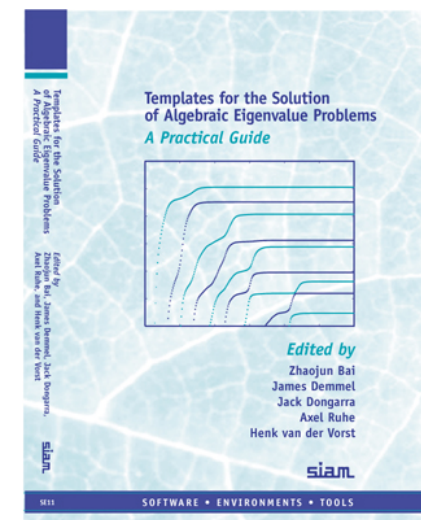


www.netlib.org/scalapack



www.netlib.org/templates

gams.nist.gov



www.cs.utk.edu/~dongarra/etemplates

A brief history of (Dense) Linear Algebra software (1/7)

- In the beginning was the do-loop...
 - Libraries like EISPACK (for eigenvalue problems)
- Then the BLAS (1) were invented (1973-1977)
 - Standard library of 15 operations (mostly) on vectors
 - “AXPY” ($y = \alpha \cdot x + y$), dot product, scale ($x = \alpha \cdot x$), etc
 - Up to 4 versions of each (S/D/C/Z), 46 routines, 3300 LOC
 - Goals
 - Common “pattern” to ease programming, readability
 - Robustness, via careful coding (avoiding over/underflow)
 - Portability + Efficiency via machine specific implementations
 - Why BLAS 1 ? They do $O(n^1)$ ops on $O(n^1)$ data
 - Used in libraries like LINPACK (for linear systems)
 - Source of the name “LINPACK Benchmark” (not the code!)

Current Records for Solving Dense Systems (11/2019)

- Linpack Benchmark
- Fastest machine overall (www.top500.org)
 - Summit at Oak Ridge National Lab (ORNL)
 - 149 Petaflops out of 201 Petaflops peak, $n = O(1M)$
 - 4,608 nodes, each with
 - 2 IBM Power9 CPUs/node, (22 cores/Power9) and
 - 6 NVIDIA Tesla GPUS (80 SMs/GPU) => 2.4M cores
 - 10.1 MWatts of power, 14.6 Gflops/Watt
 - 2.8 Pbytes of memory
- Historical data (www.netlib.org/performance)
 - Palm Pilot III
 - 1.69 Kiloflops
 - $n = 100$

A brief history of (Dense) Linear Algebra software (2/7)

- But the BLAS-1 weren't enough
 - Consider AXPY ($y = \alpha \cdot x + y$): $2n$ flops on $3n$ read/writes
 - Computational intensity = flops/words = $(2n)/(3n) = 2/3$
 - Too low to run near peak speed (read/write dominates)
 - Hard to vectorize (“SIMD’ize”) on supercomputers of the day (1980s)
- So the BLAS-2 were invented (1984-1986)
 - Standard library of 25 operations (mostly) on matrix/vector pairs
 - “GEMV”: $y = \alpha \cdot A \cdot x + \beta \cdot x$, “GER”: $A = A + \alpha \cdot x \cdot y^T$, $x = T^{-1} \cdot x$
 - Up to 4 versions of each (S/D/C/Z), 66 routines, 18K LOC
 - Why BLAS 2 ? They do $O(n^2)$ ops on $O(n^2)$ data
 - So computational intensity still just $\sim (2n^2)/(n^2) = 2$
 - OK for vector machines, but not for machine with caches

A brief history of (Dense) Linear Algebra software (3/7)

- The next step: BLAS-3 (1987-1988)
 - Standard library of 9 operations (mostly) on matrix/matrix pairs
 - “GEMM”: $C = \alpha \cdot A \cdot B + \beta \cdot C$, $C = \alpha \cdot A \cdot A^T + \beta \cdot C$, $B = T^{-1} \cdot B$
 - Up to 4 versions of each (S/D/C/Z), 30 routines, 10K LOC
 - Why BLAS 3 ? They do $O(n^3)$ ops on $O(n^2)$ data
 - So computational intensity $(2n^3)/(4n^2) = n/2$ – big at last!
 - Good for machines with caches, other mem. hierarchy levels
- How much BLAS1/2/3 code so far (all at www.netlib.org/blas)
 - Source: 142 routines, 31K LOC, Testing: 28K LOC
 - Reference (unoptimized) implementation only
 - Ex: 3 nested loops for GEMM
 - Lots more optimized code (eg Homework 1)
 - Motivates “automatic tuning” of the BLAS
 - Part of standard math libraries (eg AMD ACML, Intel MKL)

Level 1 BLAS

	dim	scalar	vector	vector	scalars	5-element array		prefixes
SUBROUTINE xROTG (A, B, C, S)		Generate plane rotation	S, D
SUBROUTINE xROTMG(D1, D2, A, B, C, S)	PARAM)	Generate modified plane rotation	S, D
SUBROUTINE xROT (N,			X, INCX, Y, INCY,		C, S)		Apply plane rotation	S, D
SUBROUTINE xROTM (N,			X, INCX, Y, INCY,			PARAM)	Apply modified plane rotation	S, D
SUBROUTINE xSWAP (N,			X, INCX, Y, INCY)				$x \leftrightarrow y$	S, D, C, Z
SUBROUTINE xSCAL (N,	ALPHA,		X, INCX)				$x \leftarrow \alpha x$	S, D, C, Z, CS, ZD
SUBROUTINE xCOPY (N,			X, INCX, Y, INCY)				$y \leftarrow x$	S, D, C, Z
SUBROUTINE xAXPY (N,	ALPHA,		X, INCX, Y, INCY)				$y \leftarrow \alpha x + y$	S, D, C, Z
FUNCTION xDOT (N,			X, INCX, Y, INCY)				$dot \leftarrow x^T y$	S, D, DS
FUNCTION xDOTU (N,			X, INCX, Y, INCY)				$dot \leftarrow x^T y$	C, Z
FUNCTION xDOTC (N,			X, INCX, Y, INCY)				$dot \leftarrow x^H y$	C, Z
FUNCTION xxDOT (N,			X, INCX, Y, INCY)				$dot \leftarrow \alpha + x^T y$	SDS
FUNCTION xNRM2 (N,			X, INCX)				$nrm2 \leftarrow \ x\ _2$	S, D, SC, DZ
FUNCTION xASUM (N,			X, INCX)				$asum \leftarrow \ re(x)\ _1 + \ im(x)\ _1$	S, D, SC, DZ
FUNCTION IxAMAX(N,			X, INCX)				$amax \leftarrow 1^{st} k \ni re(x_k) + im(x_k) $ $= \max(re(x_i) + im(x_i))$	S, D, C, Z

Level 2

xGEMV (
xGBMV (
xHEMV (
xHBMV (
xHPMV (
xSYMV (
xSBMV (
xSPMV (
xTRMV (
xTBMV (
xTPMV (
xTRSV (
xTBSV (
xTPSV (

xGER (
xGERU (
xGERC (
xHER (
xHPR (
xHER2 (
xHPR2 (
xSYR (
xSPR (
xSYR2 (
xSPR2 (

Level 3

xGEMM (
xSYMM (
xHEMM (
xSYRK (
xHERK (
xSYR2K(
xHER2K(

xTRMM (SIDE, UPLO, TRANSA, DIAG, M, N, ALPHA, A, LDA, B, LDB)
xTRSM (SIDE, UPLO, TRANSA, DIAG, M, N, ALPHA, A, LDA, B, LDB)

$B \leftarrow \alpha op(A)B, B \leftarrow \alpha Bop(A), op(A) = A, A^T, A^H, B - m \times n$ S, D, C, Z
 $B \leftarrow \alpha op(A^{-1})B, B \leftarrow \alpha Bop(A^{-1}), op(A) = A, A^T, A^H, B - m \times n$ S, D, C, Z

BLAS Standards Committee started meeting again in May 2016:

Batched BLAS: many independent BLAS operations at once

Reproducible BLAS: getting bitwise identical answers from

run-to-run, despite nonassociative floating point, and dynamic

scheduling of resources (bebop.cs.berkeley.edu/reproblas)

(new instruction in 2019 IEEE 754 Floating Point Standard

to accelerate this)

Low-Precision BLAS: eg 16 bit floating point (which one?)

See www.netlib.org/blas/blast-forum/ for previous extension attempt

New functions, Sparse BLAS, Extended Precision BLAS

GraphBLAS also underway

A brief history of (Dense) Linear Algebra software (4/7)

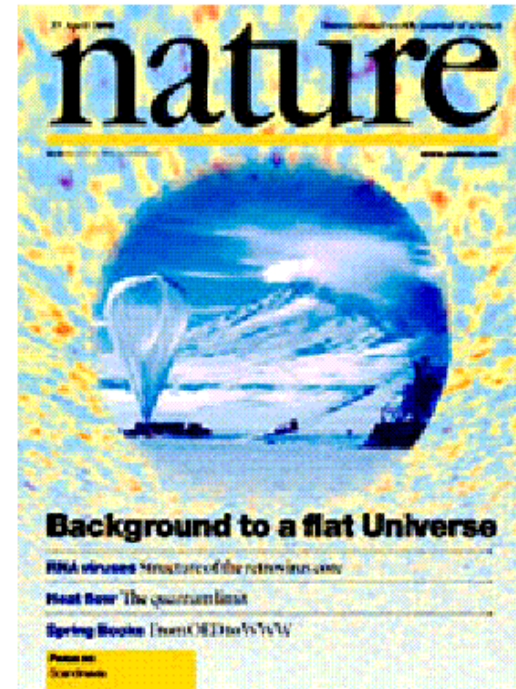
- LAPACK – “Linear Algebra PACKage” - uses BLAS-3 (1989 – now)
 - Ex: Obvious way to express Gaussian Elimination (GE) is adding multiples of one row to other rows – BLAS-1
 - How do we reorganize GE to use BLAS-3 ? (details later)
 - Contents of LAPACK (summary)
 - Algorithms that are (nearly) 100% BLAS 3
 - Linear Systems: solve $Ax=b$ for x
 - Least Squares: choose x to minimize $\|Ax-b\|_2$
 - Algorithms that are only $\approx 50\%$ BLAS 3
 - Eigenproblems: Find λ and x where $Ax = \lambda x$
 - Singular Value Decomposition (SVD)
 - Generalized problems (eg $Ax = \lambda Bx$)
 - Error bounds for everything
 - Lots of variants depending on A 's structure (banded, $A=A^T$, etc)
 - How much code? (Release 3.9.0, Nov 2019) (www.netlib.org/lapack)
 - Source: 1982 routines, 827K LOC, Testing: 1210 routines, 545K LOC
 - Ongoing development (at UCB and elsewhere) (class projects!)

A brief history of (Dense) Linear Algebra software (5/7)

- Is LAPACK parallel?
 - Only if the BLAS are parallel (possible in shared memory)
- ScaLAPACK – “Scalable LAPACK” (1995 – now)
 - For distributed memory – uses MPI
 - More complex data structures, algorithms than LAPACK
 - Only subset of LAPACK’s functionality available
 - Details later (class projects!)
 - All at www.netlib.org/scalapack

Success Stories for Sca/LAPACK (6/7)

- Widely used
 - Adopted by Mathworks, Cray, Fujitsu, HP, IBM, IMSL, Intel, NAG, NEC, SGI, ...
- New Science discovered through the solution of dense matrix systems
 - Nature article on the flat universe used ScaLAPACK
 - Other articles in Physics Review B that also use it
 - 1998 Gordon Bell Prize
 - www.nersc.gov/assets/NewsImages/2003/newNERSCresults050703.pdf



**Cosmic Microwave Background
Analysis, BOOMERanG
collaboration, MADCAP code (Apr.
27, 2000).**

A brief future look at (Dense) Linear Algebra software (7/7)

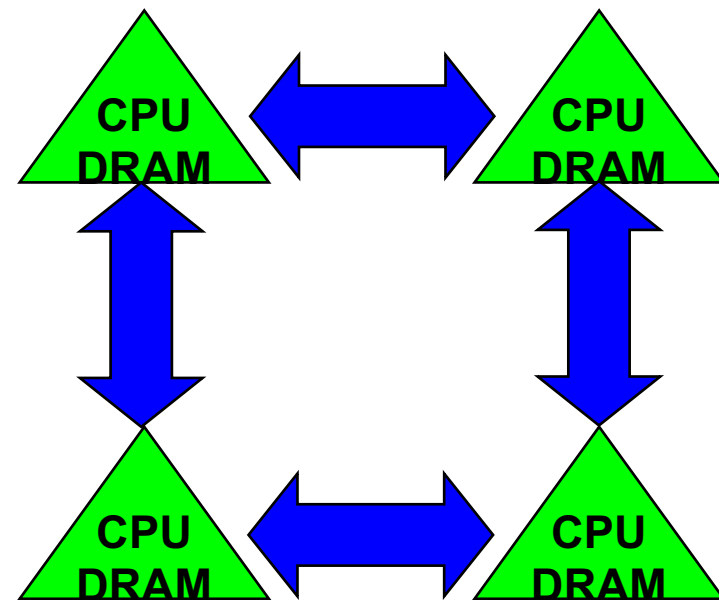
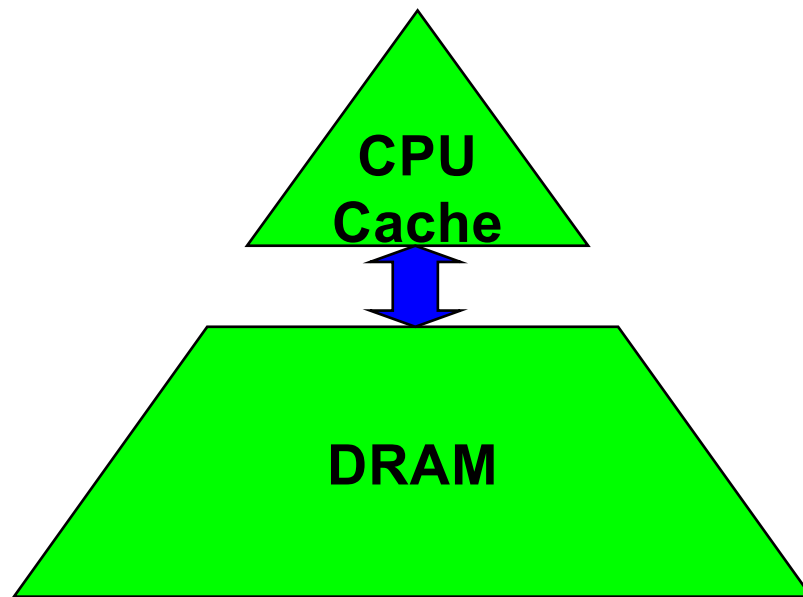
- PLASMA, DPLASMA and MAGMA (now)
 - Ongoing extensions to Multicore/GPU/Heterogeneous
 - Can one software infrastructure accommodate all algorithms and platforms of current (future) interest?
 - How much code generation and tuning can we automate?
 - Details later (icl.cs.utk.edu/{{d}}plasma,magma)
- Other related projects
 - SLATE (icl.utk.edu/slate) – accelerators, data layouts, ...
 - Elemental (libelemental.org)
 - Distributed memory dense linear algebra
 - “Balance ease of use and high performance”
 - FLAME (www.cs.utexas.edu/users/flame/web/index.html)
 - Formal Linear Algebra Method Environment
 - Attempt to automate code generation across multiple platforms
- So far, none of these libraries minimize communication in all cases (not even matmul!)

Back to basics:

Why avoiding communication is important (1/3)

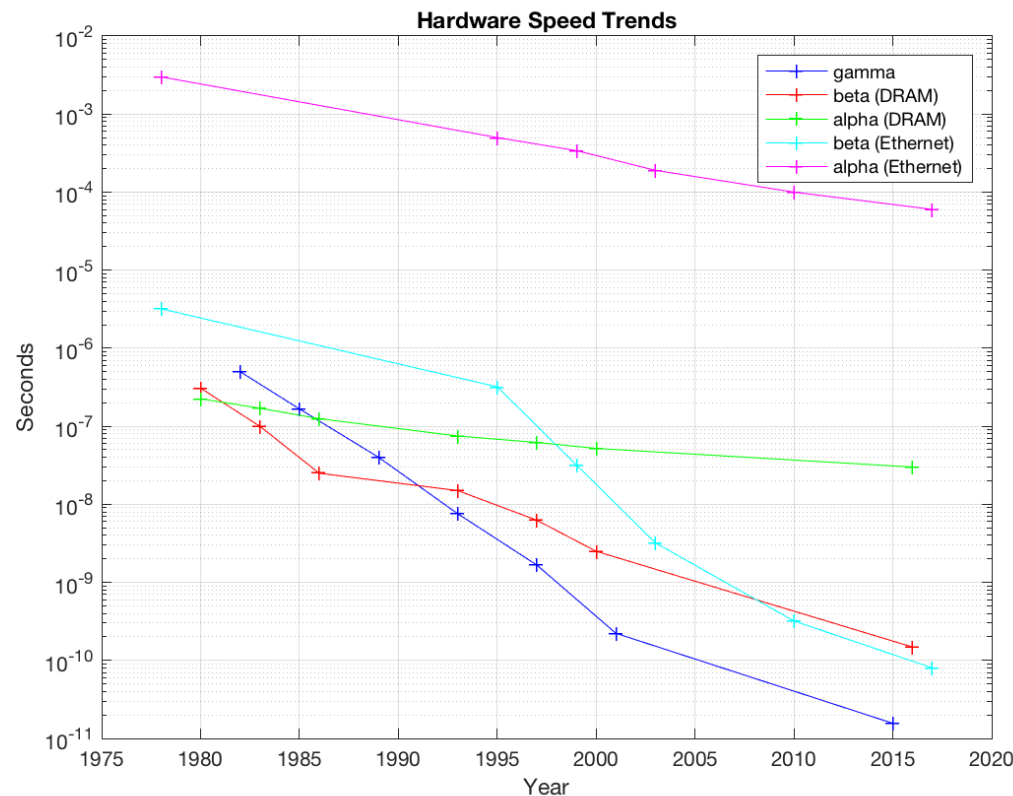
Algorithms have two costs:

1. Arithmetic (FLOPS)
2. Communication: moving data between
 - levels of a memory hierarchy (sequential case)
 - processors over a network (parallel case).



Why avoiding communication is important (2/3)

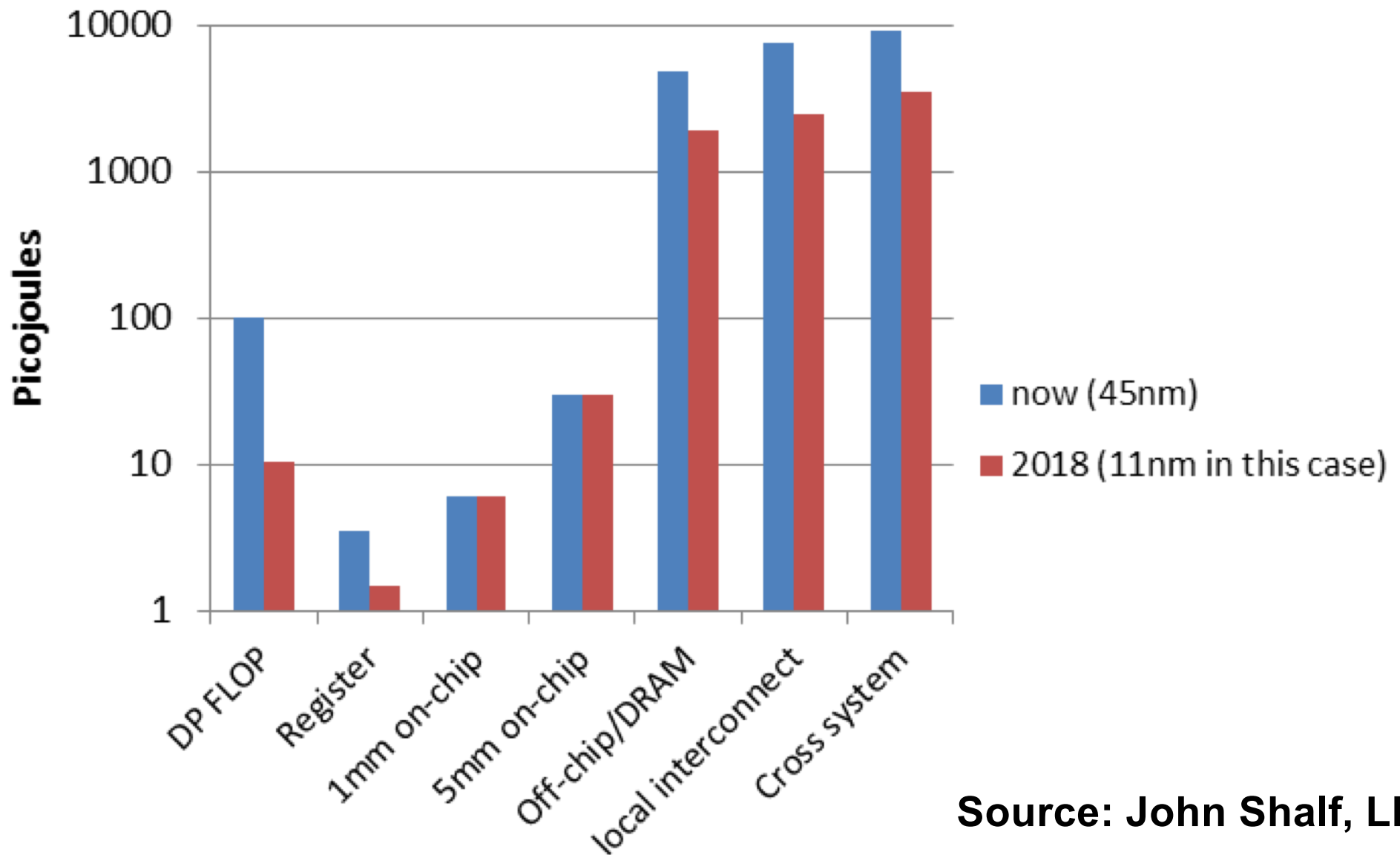
- Running time of an algorithm is sum of 3 terms:
 - # flops * time_per_flop
 - # words moved / bandwidth
 - # messages * latency
- } **communication**



- Time_per_flop
 $\ll 1/\text{bandwidth}$
 $\ll \text{latency}$
- Gaps growing exponentially with time

- Minimize communication to save time**

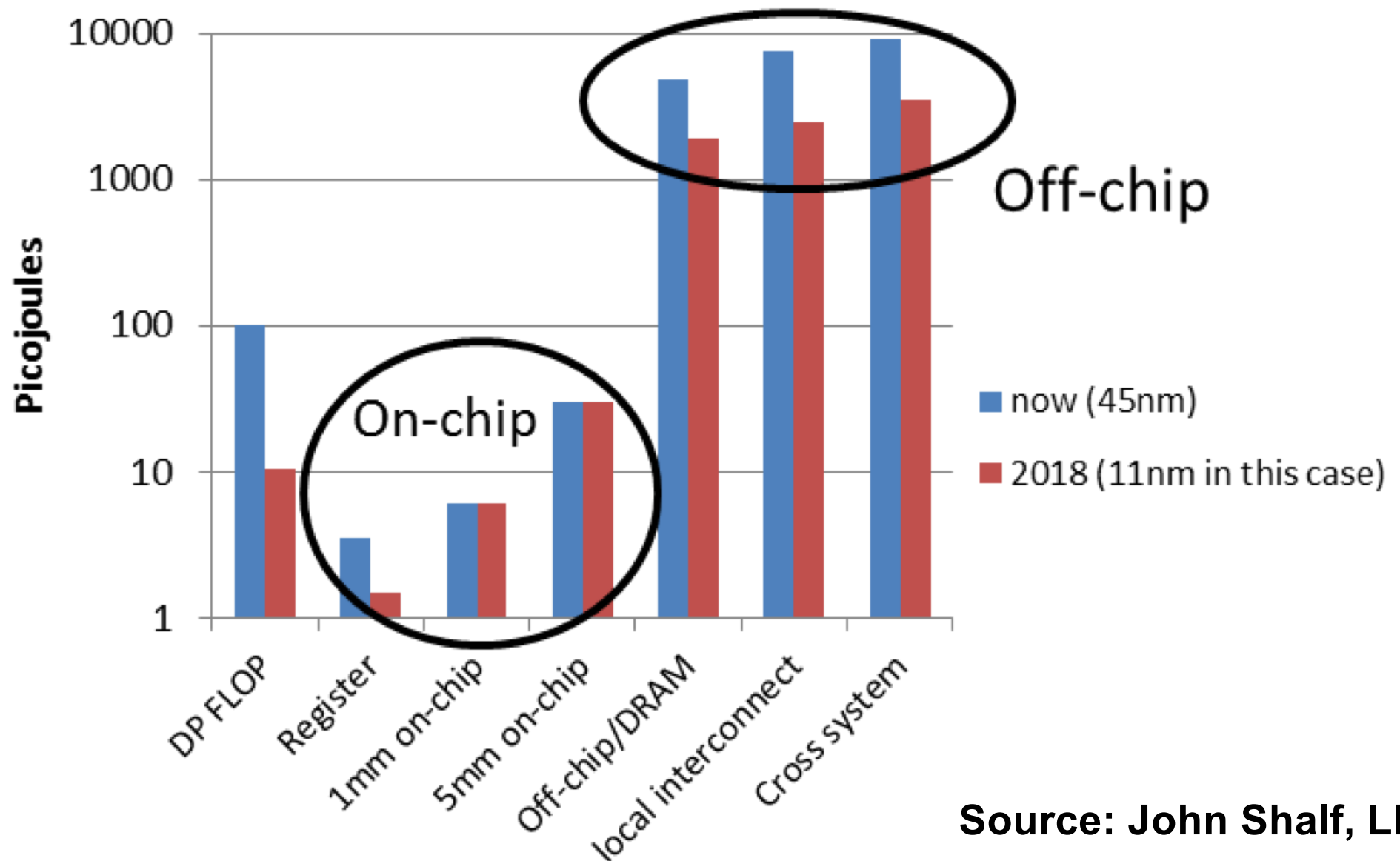
Why Minimize Communication? (3/3)



Source: John Shalf, LBL

Why Minimize Communication? (3/3)

Minimize communication to save energy



Source: John Shalf, LBL

Goal:

Organize Linear Algebra to Avoid Communication

- Between all memory hierarchy levels
 - $L1 \leftrightarrow L2 \leftrightarrow \text{DRAM} \leftrightarrow \text{network, etc}$
- Not just *hiding* communication (overlap with arithmetic)
 - $\text{Speedup} \leq 2x$
- Arbitrary speedups/energy savings possible
- Later: Same goal for other computational patterns
 - Lots of open problems

Review: Blocked Matrix Multiply

- Blocked Matmul $C = A \cdot B$ breaks A , B and C into blocks with dimensions that depend on cache size

... Break $A^{n \times n}$, $B^{n \times n}$, $C^{n \times n}$ into $b \times b$ blocks labeled $A(i,j)$, etc

... b chosen so 3 $b \times b$ blocks fit in cache

for $i = 1$ to n/b , for $j=1$ to n/b , for $k=1$ to n/b

$C(i,j) = C(i,j) + A(i,k) \cdot B(k,j)$... $b \times b$ matmul, $4b^2$ reads/writes

- When $b=1$, get “naïve” algorithm, want b larger ...
- $(n/b)^3 \cdot 4b^2 = 4n^3/b$ reads/writes altogether
- Minimized when $3b^2 = \text{cache size} = M$, yielding $O(n^3/M^{1/2})$ reads/writes
- What if we had more levels of memory? (L1, L2, cache etc)?
 - Would need 3 more nested loops per level
 - Recursive (cache-oblivious algorithm) also possible

Communication Lower Bounds: Prior Work on Matmul

- Assume n^3 algorithm (i.e. not Strassen-like)
- Sequential case, with fast memory of size M
 - Lower bound on #words moved to/from slow memory = $\Omega(n^3 / M^{1/2})$ [Hong, Kung, 81]
 - Attained using blocked or cache-oblivious algorithms
- Parallel case on P processors:
 - Let M be memory per processor; assume load balanced
 - Lower bound on #words moved
= $\Omega((n^3 / p) / M^{1/2})$ [Irony, Tiskin, Toledo, 04]
 - If $M = 3n^2/p$ (one copy of each matrix), then
lower bound = $\Omega(n^2 / p^{1/2})$
 - Attained by SUMMA, Cannon's algorithm

New lower bound for all “direct” linear algebra

Let M = “fast” memory size per processor
= cache size (sequential case) or $O(n^2/p)$ (parallel case)
#flops = number of flops done per processor

$$\text{\#words_moved per processor} = \Omega(\text{\#flops} / M^{1/2})$$

$$\text{\#messages_sent per processor} = \Omega(\text{\#flops} / M^{3/2})$$

- Holds for
 - Matmul, BLAS, LU, QR, eig, SVD, tensor contractions, ...
 - Some whole programs (sequences of these operations, no matter how they are interleaved, eg computing A^k)
 - Dense *and* sparse matrices (where $\text{\#flops} \ll n^3$)
 - Sequential *and* parallel algorithms
 - Some graph-theoretic algorithms (eg Floyd-Warshall)
- Generalizations (Strassen-like algorithms, loops accessing arrays)

New lower bound for all “direct” linear algebra

Let M = “fast” memory size per processor
= cache size (sequential case) or $O(n^2/p)$ (parallel case)
#flops = number of flops done per processor

$$\text{\#words_moved per processor} = \Omega(\text{\#flops} / M^{1/2})$$

$$\text{\#messages_sent per processor} = \Omega(\text{\#flops} / M^{3/2})$$

- Sequential case, dense $n \times n$ matrices, so $O(n^3)$ flops
 - #words_moved = $\Omega(n^3 / M^{1/2})$
 - #messages_sent = $\Omega(n^3 / M^{3/2})$
- Parallel case, dense $n \times n$ matrices
 - Load balanced, so $O(n^3/p)$ flops processor
 - One copy of data, load balanced, so $M = O(n^2/p)$ per processor
 - #words_moved = $\Omega(n^2 / p^{1/2})$
 - #messages_sent = $\Omega(p^{1/2})$

SIAM Linear Algebra Prize, 2012

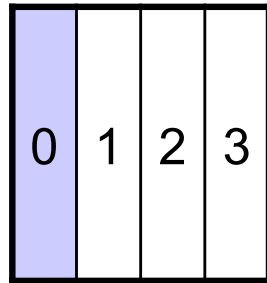
Can we attain these lower bounds?

- Do conventional dense algorithms as implemented in LAPACK and ScaLAPACK attain these bounds?
 - Mostly not yet, work in progress
- If not, are there other algorithms that do?
 - Yes
- Goals for algorithms:
 - Minimize #words_moved
 - Minimize #messages_sent
 - Need new data structures
 - Minimize for multiple memory hierarchy levels
 - Cache-oblivious algorithms would be simplest
 - Fewest flops when matrix fits in fastest memory
 - Cache-oblivious algorithms don't always attain this
- Attainable for nearly all dense linear algebra
 - Just a few prototype implementations so far (class projects!)
 - Only a few sparse algorithms so far (eg Cholesky, sparse-dense matmul)

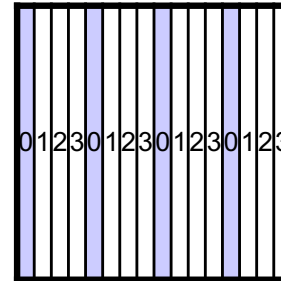
Outline

- History and motivation
 - What is dense linear algebra?
 - Why minimize communication?
 - Lower bound on communication
- **Parallel Matrix-matrix multiplication**
 - **Attaining the lower bound**
- Other Parallel Algorithms (next lecture)

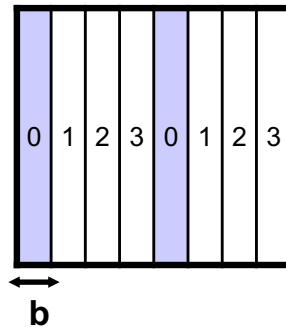
Different Parallel Data Layouts for Matrices (not all!)



1) 1D Column Blocked Layout

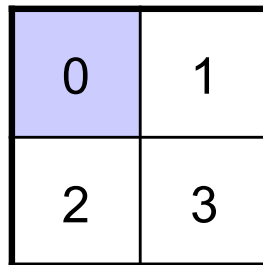


2) 1D Column Cyclic Layout

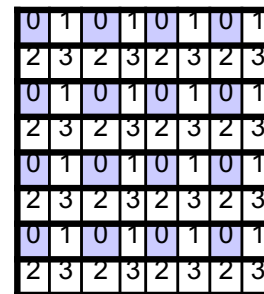


3) 1D Column Block Cyclic Layout

4) Row versions of the previous layouts



5) 2D Row and Column Blocked Layout



6) 2D Row and Column Block Cyclic Layout

Generalizes others

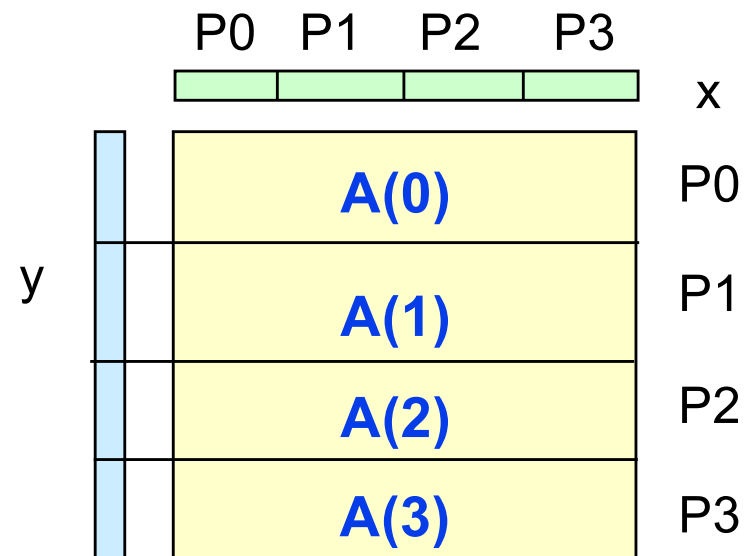
Parallel Matrix-Vector Product

- Compute $y = y + A*x$, where A is a dense matrix
- Layout:

- **1D row blocked**

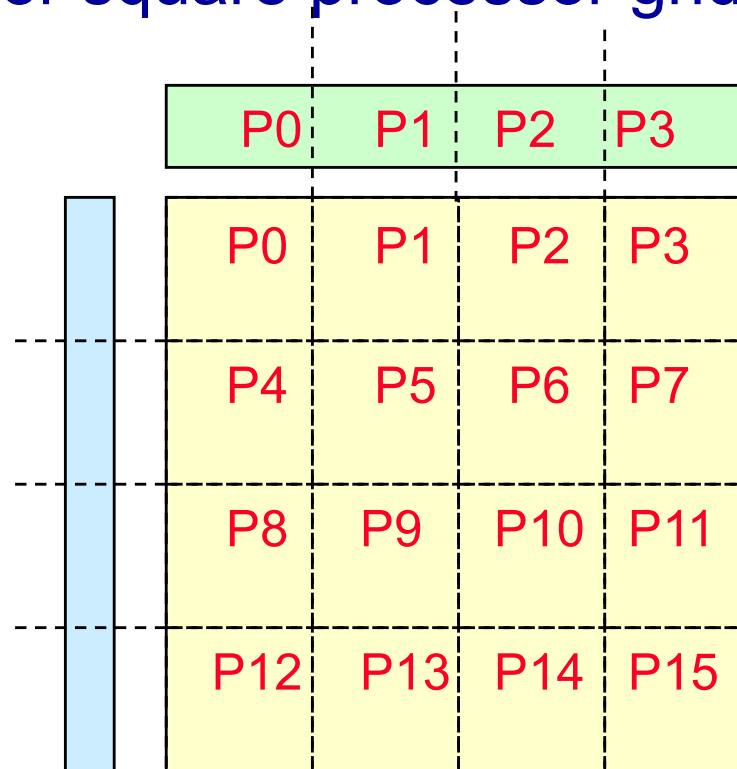
- $A(i)$ refers to the n by n/p block row that processor i owns,
- $x(i)$ and $y(i)$ similarly refer to segments of x, y owned by i
- **Algorithm:**
 - **Foreach processor i**
 - **Broadcast $x(i)$**
 - **Compute $y(i) = A(i)*x$**
- Algorithm uses the formula

$$y(i) = y(i) + A(i)*x = y(i) + \sum_j A(i,j)*x(j)$$



Matrix-Vector Product $y = y + A*x$

- A **column layout** of the matrix eliminates the broadcast of x
 - But adds a reduction to update the destination y
- A **2D blocked layout** uses a broadcast and reduction, both on a subset of processors
 - \sqrt{p} for square processor grid

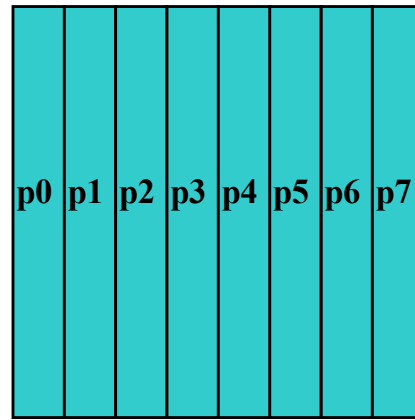


Parallel Matrix Multiply

- Computing $C = C + A * B$
- Using basic algorithm: $2 * n^3$ Flops
- Variables are:
 - Data layout: 1D? 2D? Other?
 - Topology of machine: Ring? Torus?
 - Scheduling communication
- Use of performance models for algorithm design
 - **Message Time = “latency” + #words * time-per-word**
 $= \alpha + n * \beta$
- Efficiency (in any model):
 - serial time / (p * parallel time)
 - perfect (linear) speedup \leftrightarrow efficiency = 1

Matrix Multiply with 1D Column Layout

- Assume matrices are $n \times n$ and n is divisible by p



May be a reasonable
assumption for analysis,
not for code

- $A(i)$ refers to the n by n/p block column that processor i owns (similarly for $B(i)$ and $C(i)$)
- $B(i,j)$ is the n/p by n/p subblock of $B(i)$
 - in rows $j*n/p$ through $(j+1)*n/p - 1$
- Algorithm uses the formula

$$C(i) = C(i) + A*B(i) = C(i) + \sum_j A(j)*B(j,i)$$

Matrix Multiply: 1D Layout on Bus or Ring

- Algorithm uses the formula

$$C(i) = C(i) + A * B(i) = C(i) + \sum_j A(j) * B(j,i)$$

- First consider a bus-connected machine without broadcast: only one pair of processors can communicate at a time (ethernet)
- Second consider a machine with processors on a ring: all processors may communicate with nearest neighbors simultaneously

MatMul: 1D layout on Bus without Broadcast

Naïve algorithm:

$C(\text{myproc}) = C(\text{myproc}) + A(\text{myproc}) * B(\text{myproc}, \text{myproc})$

for $i = 0$ to $p-1$

 for $j = 0$ to $p-1$ except i

 if ($\text{myproc} == i$) send $A(i)$ to processor j

 if ($\text{myproc} == j$)

 receive $A(i)$ from processor i

$C(\text{myproc}) = C(\text{myproc}) + A(i) * B(i, \text{myproc})$

 barrier

Cost of inner loop:

computation: $2 * n * (n/p)^2 = 2 * n^3 / p^2$

communication: $\alpha + \beta * n^2 / p$

Naïve MatMul (continued)

Cost of inner loop:

computation: $2*n*(n/p)^2 = 2*n^3/p^2$

communication: $\alpha + \beta*n^2/p$

Only 1 pair of processors (i and j) are active on any iteration,
and of those, only i is doing computation

=> the algorithm is almost entirely serial

Running time:

$= (p*(p-1) + 1)*\text{computation} + p*(p-1)*\text{communication}$

$\approx 2*n^3 + p^2*\alpha + p*n^2*\beta$

This is worse than the serial time and grows with p.

Matmul for 1D layout on a Processor Ring

- Pairs of adjacent processors can communicate simultaneously

Copy $A(\text{myproc})$ into Tmp

$C(\text{myproc}) = C(\text{myproc}) + \text{Tmp} * B(\text{myproc}, \text{myproc})$

for $j = 1$ to $p-1$

 Send Tmp to processor $\text{myproc}+1 \bmod p$

 Receive Tmp from processor $\text{myproc}-1 \bmod p$

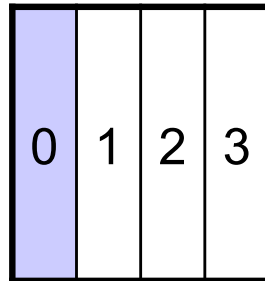
$C(\text{myproc}) = C(\text{myproc}) + \text{Tmp} * B(\text{myproc}-j \bmod p, \text{myproc})$

- Same idea as for gravity in simple sharks and fish algorithm
 - May want double buffering in practice for overlap
 - Ignoring deadlock details in code
- Time of inner loop = $2 * (\alpha + \beta * n^2 / p) + 2 * n * (n/p)^2$

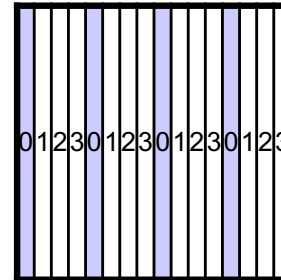
Matmul for 1D layout on a Processor Ring

- Time of inner loop = $2*(\alpha + \beta*n^2/p) + 2*n*(n/p)^2$
- Total Time = $2*n*(n/p)^2 + (p-1) * \text{Time of inner loop}$
- $\approx 2*n^3/p + 2*p*\alpha + 2*\beta*n^2$
- (Nearly) Optimal for 1D layout on Ring or Bus, even with Broadcast:
 - Perfect speedup for arithmetic
 - $A(\text{myproc})$ must move to each other processor, costs at least $(p-1)*\text{cost of sending } n*(n/p) \text{ words}$
- Parallel Efficiency = $2*n^3 / (p * \text{Total Time})$
$$= 1/(1 + \alpha * p^2/(2*n^3) + \beta * p/(2*n))$$
$$= 1/ (1 + O(p/n))$$
- Grows to 1 as n/p increases (or α and β shrink)
- But far from communication lower bound

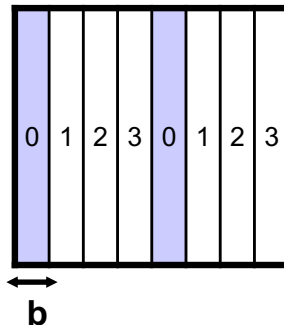
Need to try 2D Matrix layout



1) 1D Column Blocked Layout

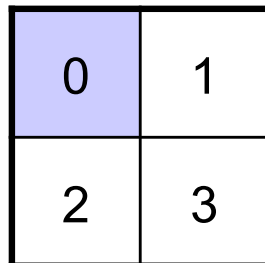


2) 1D Column Cyclic Layout

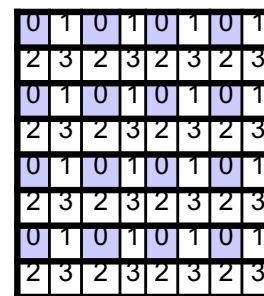


3) 1D Column Block Cyclic Layout

4) Row versions of the previous layouts



5) 2D Row and Column Blocked Layout



6) 2D Row and Column Block Cyclic Layout

Generalizes others

Summary of Parallel Matrix Multiply

- SUMMA
 - Scalable Universal Matrix Multiply Algorithm
 - Attains communication lower bounds (within $\log p$)
- Cannon
 - Historically first, attains lower bounds
 - More assumptions
 - A and B square
 - P a perfect square
- 2.5D SUMMA
 - Uses more memory to communicate even less
- Parallel Strassen
 - Attains different, even lower bounds

SUMMA Algorithm

- SUMMA = Scalable Universal Matrix Multiply
- Presentation from van de Geijn and Watts
 - www.netlib.org/lapack/lawns/lawn96.ps
 - Similar ideas appeared many times
- Used in practice in PBLAS = Parallel BLAS
 - www.netlib.org/lapack/lawns/lawn100.ps

SUMMA uses Outer Product form of MatMul

- $C = A*B$ means $C(i,j) = \sum_k A(i,k)*B(k,j)$

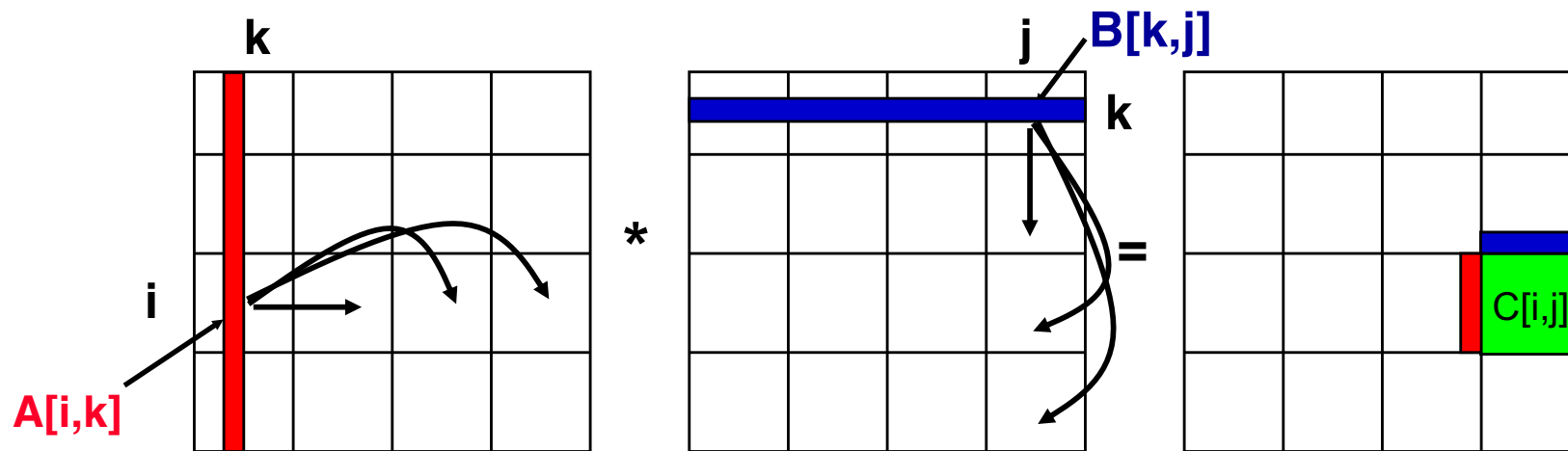
- Column-wise outer product:

$$\begin{aligned} C &= A*B \\ &= \sum_k A(:,k)*B(k,:) \\ &= \sum_k (\text{k-th col of } A) * (\text{k-th row of } B) \end{aligned}$$

- Block column-wise outer product
(block size = 4 for illustration)

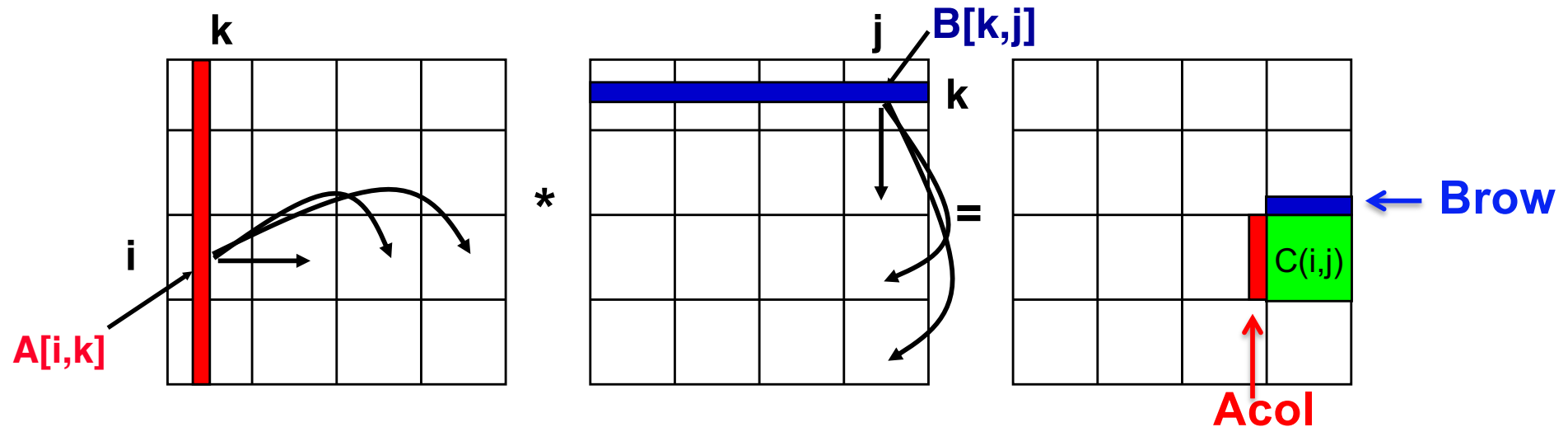
$$\begin{aligned} C &= A*B \\ &= A(:,1:4)*B(1:4,:) + A(:,5:8)*B(5:8,:) + \dots \\ &= \sum_k (\text{k-th block of 4 cols of } A) * \\ &\quad (\text{k-th block of 4 rows of } B) \end{aligned}$$

SUMMA – $n \times n$ matmul on $P^{1/2} \times P^{1/2}$ grid



- $C[i, j]$ is $n/P^{1/2} \times n/P^{1/2}$ submatrix of C on processor P_{ij}
- $A[i, k]$ is $n/P^{1/2} \times b$ submatrix of A
- $B[k, j]$ is $b \times n/P^{1/2}$ submatrix of B
- $C[i, j] = C[i, j] + \sum_k A[i, k] * B[k, j]$
 - summation over submatrices
- Need not be square processor grid

SUMMA– $n \times n$ matmul on $P^{1/2} \times P^{1/2}$ grid



For $k=0$ to $n/b-1$

for all $i = 1$ to $P^{1/2}$

owner of $A[i,k]$ broadcasts it to whole processor row (using binary tree)

for all $j = 1$ to $P^{1/2}$

owner of $B[k,j]$ broadcasts it to whole processor column (using bin. tree)

Receive $A[i,k]$ into $Acol$

Receive $B[k,j]$ into $Brow$

$C_{myproc} = C_{myproc} + Acol * Brow$

SUMMA Costs

For $k=0$ to $n/b-1$

for all $i = 1$ to $P^{1/2}$

owner of $A[i,k]$ broadcasts it to whole processor row (using binary tree)

... #words = $\log P^{1/2} * b * n / P^{1/2}$, #messages = $\log P^{1/2}$

for all $j = 1$ to $P^{1/2}$

owner of $B[k,j]$ broadcasts it to whole processor column (using bin. tree)

... same #words and #messages

Receive $A[i,k]$ into A_{col}

Receive $B[k,j]$ into B_{row}

$C_{myproc} = C_{myproc} + A_{col} * B_{row}$... #flops = $2n^2 * b / P$

- Total #words = $\log P * n^2 / P^{1/2}$
 - Within factor of $\log P$ of lower bound
 - (more complicated implementation removes $\log P$ factor)
- Total #messages = $\log P * n/b$
 - Choose b close to maximum, $n/P^{1/2}$, to approach lower bound $P^{1/2}$
- Total #flops = $2n^3/P$

Performance of PBLAS

**PDGEMM = PBLAS routine
for matrix multiply**

Observations:

For fixed N, as P increases
Mflops increases, but
less than 100% efficiency
For fixed P, as N increases,
Mflops (efficiency) rises

Speed in Mflops of PDGEMM					
Machine	Procs	Block Size	N		
			2000	4000	10000
Cray T3E	4=2x2	32	1055	1070	0
	16=4x4		3630	4005	4292
	64=8x8		13456	14287	16755
IBM SP2	4	50	755	0	0
	16		2514	2850	0
	64		6205	8709	10774
Intel XP/S MP Paragon	4	32	330	0	0
	16		1233	1281	0
	64		4496	4864	5257
Berkeley NOW	4	32	463	470	0
	32=4x8		2490	2822	3450
	64		4130	5457	6647

**DGEMM = BLAS routine
for matrix multiply**

**Maximum speed for PDGEMM
= # Procs * speed of DGEMM**

Observations (same as above):
Efficiency always at least 48%
For fixed N, as P increases,
efficiency drops
For fixed P, as N increases,
efficiency increases

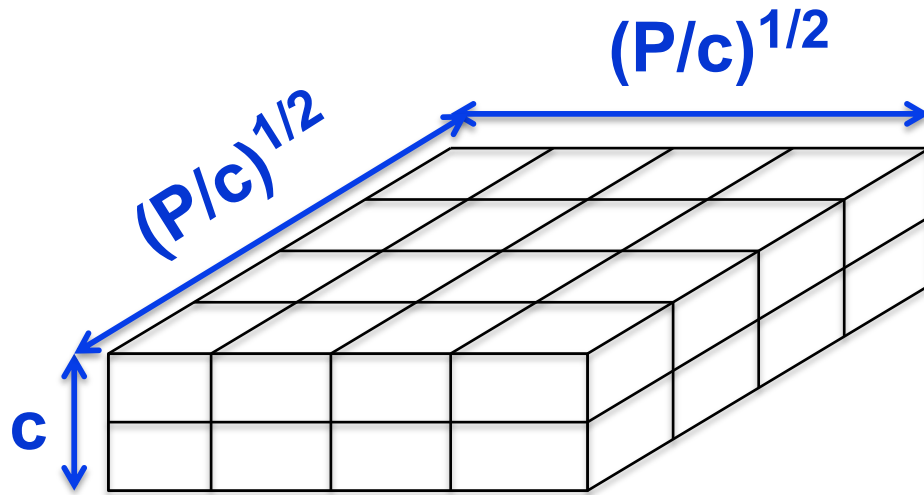
Efficiency = MFlops(PDGEMM)/(Procs*MFlops(DGEMM))						
Machine	Peak/ proc	DGEMM Mflops	Procs	N		
				2000	4000	10000
Cray T3E	600	360	4	.73	.74	
			16	.63	.70	.75
			64	.58	.62	.73
IBM SP2	266	200	4	.94		
			16	.79	.89	
			64	.48	.68	.84
Intel XP/S MP Paragon	100	90	4	.92		
			16	.86	.89	
			64	.78	.84	.91
Berkeley NOW	334	129	4	.90	.91	
			32	.60	.68	.84
			64	.50	.66	.81

Can we do matrix multiply better?

- Lower bound assumed 1 copy of data: $M = O(n^2/P)$ per proc.
- What if matrix small enough to fit $c > 1$ copies, so $M = cn^2/P$?
 - $\#words_moved = \Omega(\#flops / M^{1/2}) = \Omega(n^2 / (c^{1/2} P^{1/2}))$
 - $\#messages = \Omega(\#flops / M^{3/2}) = \Omega(P^{1/2} / c^{3/2})$
- Can we attain new lower bound?
 - Special case: “3D Matmul”: $c = P^{1/3}$
 - Bernstein 89, Agarwal, Chandra, Snir 90, Aggarwal 95
 - Processors arranged in $P^{1/3} \times P^{1/3} \times P^{1/3}$ grid
 - Processor (i,j,k) performs $C(i,j) = C(i,j) + A(i,k)*B(k,j)$, where each submatrix is $n/P^{1/3} \times n/P^{1/3}$
 - Not always that much memory available...

2.5D Matrix Multiplication

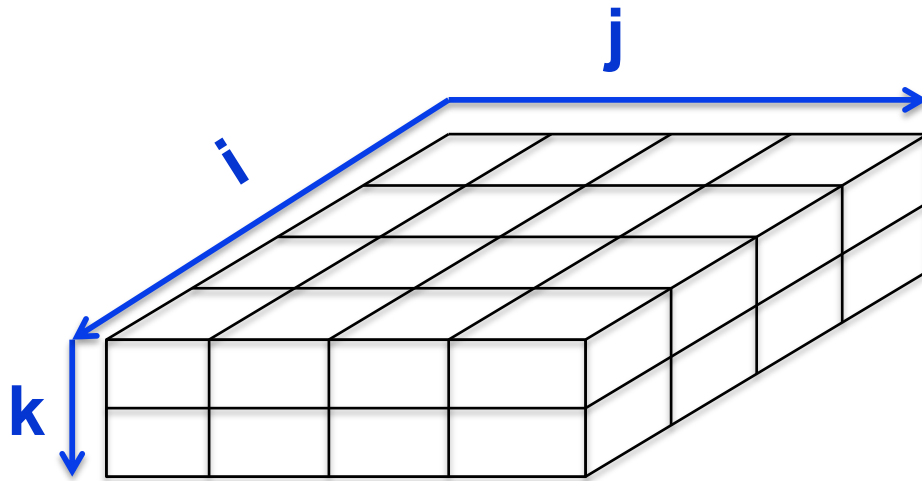
- Assume can fit cn^2/P data per processor, $c > 1$
- Processors form $(P/c)^{1/2} \times (P/c)^{1/2} \times c$ grid



Example: $P = 32$, $c = 2$

2.5D Matrix Multiplication

- Assume can fit cn^2/P data per processor, $c > 1$
- Processors form $(P/c)^{1/2} \times (P/c)^{1/2} \times c$ grid



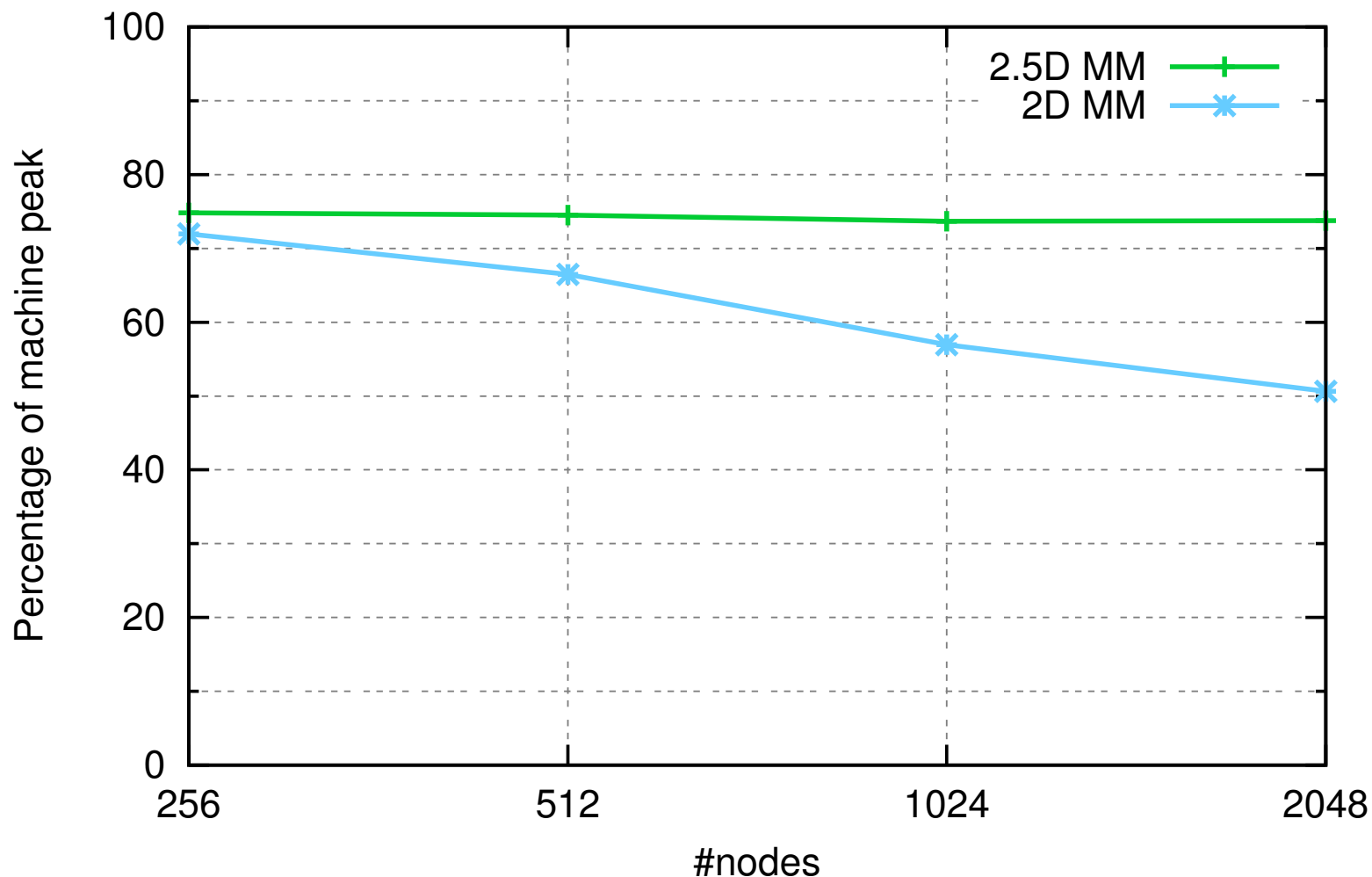
Initially $P(i,j,0)$ owns $A(i,j)$ and $B(i,j)$
each of size $n(c/P)^{1/2} \times n(c/P)^{1/2}$

- (1) $P(i,j,0)$ broadcasts $A(i,j)$ and $B(i,j)$ to $P(i,j,k)$
- (2) Processors at level k perform $1/c$ -th of SUMMA, i.e. $1/c$ -th of $\sum_m A(i,m) * B(m,j)$
- (3) Sum-reduce partial sums $\sum_m A(i,m) * B(m,j)$ along k -axis so $P(i,j,0)$ owns $C(i,j)$

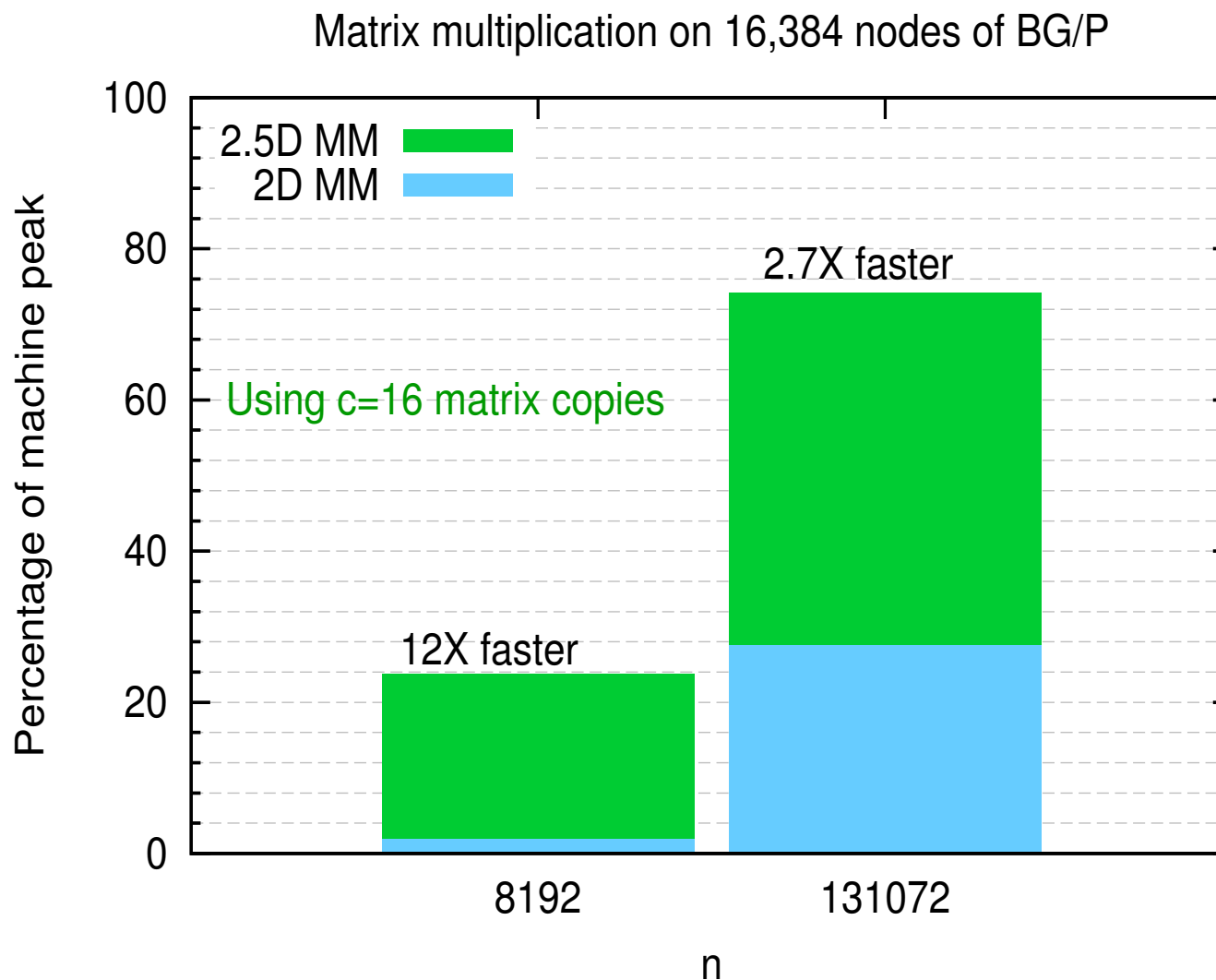
2.5D Matmul on IBM BG/P, n=64K

- As P increases, available memory grows \rightarrow c increases proportionally to P
 - #flops, #words_moved, #messages per proc all decrease proportionally to P
 - $\#words_moved = \Omega(\#flops / M^{1/2}) = \Omega(n^2 / (c^{1/2} P^{1/2}))$
 - $\#messages = \Omega(\#flops / M^{3/2}) = \Omega(P^{1/2} / c^{3/2})$
- Perfect strong scaling! But only up to $c = P^{1/3}$

Matrix multiplication on BG/P (n=65,536)



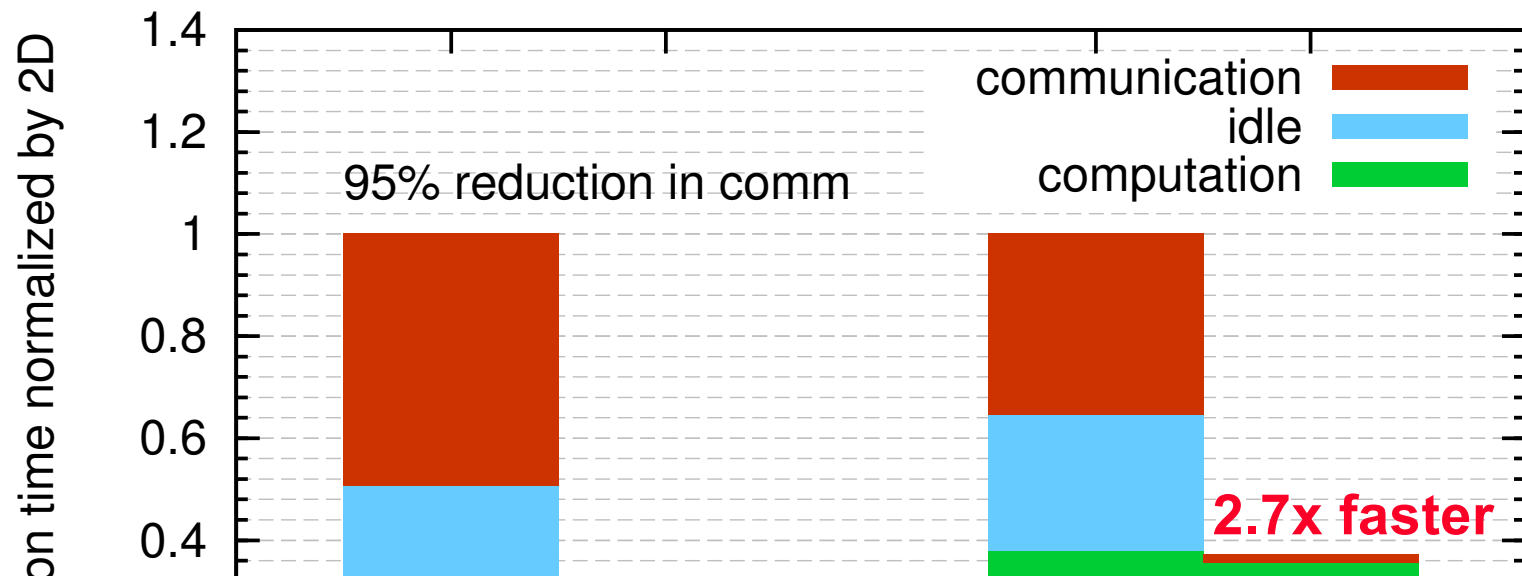
2.5D Matmul on IBM BG/P, 16K nodes / 64K cores



2.5D Matmul on BG/P, 16K nodes / 64K cores

c = 16 copies

Matrix multiplication on 16,384 nodes of BG/P



**Ideas adopted by Nervana, “deep learning” startup,
acquired by Intel in August 2016**

n=8192, 2D

n=8192, 2.5D

n=131072, 2D

n=131072, 2.5D

**Distinguished Paper Award, EuroPar’11 (Solomonik, D.)
SC’11 paper by Solomonik, Bhatele, D.**

Perfect Strong Scaling – in Time and Energy

- Every time you add a processor, you should use its memory M too
- Start with minimal number of procs: $PM = 3n^2$
- Increase P by a factor of $c \rightarrow$ total memory increases by a factor of c
- Notation for timing model:
 - $\gamma_T, \beta_T, \alpha_T$ = secs per flop, per word_moved, per message of size m
- $T(cP) = n^3/(cP) [\gamma_T + \beta_T/M^{1/2} + \alpha_T/(mM^{1/2})]$
 $= T(P)/c$
- Notation for energy model:
 - $\gamma_E, \beta_E, \alpha_E$ = joules for same operations
 - δ_E = joules per word of memory used per sec
 - ϵ_E = joules per sec for leakage, etc.
- $E(cP) = cP \{ n^3/(cP) [\gamma_E + \beta_E/M^{1/2} + \alpha_E/(mM^{1/2})] + \delta_E MT(cP) + \epsilon_E T(cP) \}$
 $= E(P)$
- c cannot increase forever: $c \leq P^{1/3}$ (3D algorithm)
 - Corresponds to lower bound on #messages hitting 1
- Perfect scaling extends to Strassen's matmul, direct N-body, ...
 - "Perfect Strong Scaling Using No Additional Energy"
 - "Strong Scaling of Matmul and Memory-Indep. Comm. Lower Bounds"
 - Both at bebop.cs.berkeley.edu

Strassen's Matrix Multiply

- The traditional algorithm (with or without tiling) has $O(n^3)$ flops
- Strassen discovered an algorithm with asymptotically lower flops
 - $O(n^{2.81})$
- Consider a 2x2 matrix multiply, normally takes 8 multiplies, 4 adds
 - Strassen does it with 7 multiplies and 18 adds

$$\text{Let } M = \begin{pmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix}$$

$$\text{Let } p1 = (a_{12} - a_{22}) * (b_{21} + b_{22})$$

$$p5 = a_{11} * (b_{12} - b_{22})$$

$$p2 = (a_{11} + a_{22}) * (b_{11} + b_{22})$$

$$p6 = a_{22} * (b_{21} - b_{11})$$

$$p3 = (a_{11} - a_{21}) * (b_{11} + b_{12})$$

$$p7 = (a_{21} + a_{22}) * b_{11}$$

$$p4 = (a_{11} + a_{12}) * b_{22}$$

$$\text{Then } m_{11} = p1 + p2 - p4 + p6$$

$$m_{12} = p4 + p5$$

$$m_{21} = p6 + p7$$

$$m_{22} = p2 - p3 + p5 - p7$$

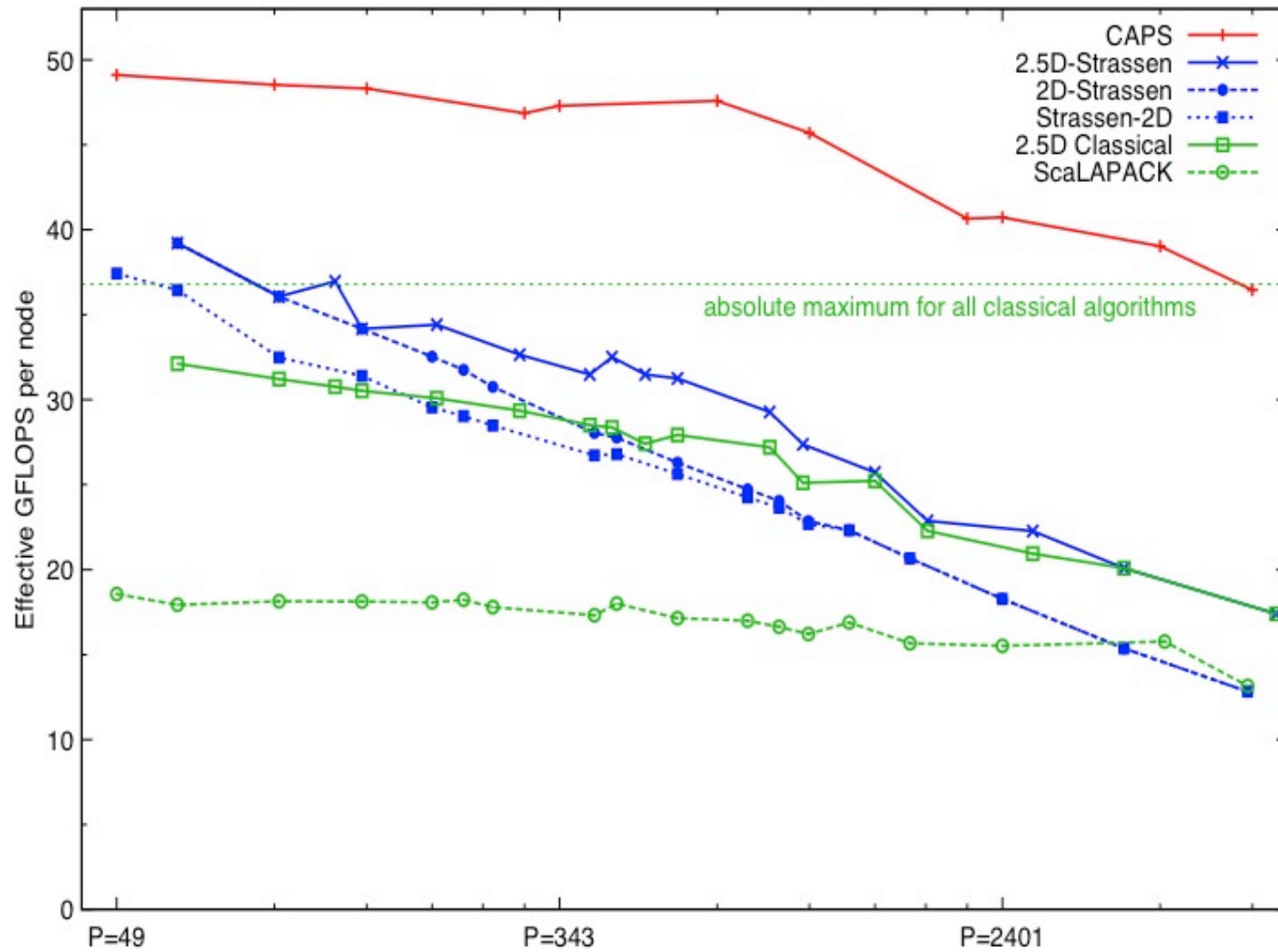
Extends to nxn by divide&conquer
 $F(n) = 7 * F(n/2) + O(n^2) = O(n^{\log_2 7})$

Classical Matmul

- Complexity of classical Matmul
- Flops: $O(n^3/p)$
- Communication lower bound on #words:
$$\Omega((n^3/p)/M^{1/2}) = \Omega(M(n/M^{1/2})^3/p)$$
- Communication lower bound on #messages:
$$\Omega((n^3/p)/M^{3/2}) = \Omega((n/M^{1/2})^3/p)$$
- All attainable as M increases past $O(n^2/p)$, up to a limit:
can increase M by factor up to $p^{1/3}$
#words as low as $\Omega(n^2/p^{2/3})$

Strong scaling of Matmul on Hopper (n=94080)

G. Ballard, D., O. Holtz, B. Lipshitz, O. Schwartz



03/03/2020

“Communication-Avoiding Parallel Strassen”
bebop.cs.berkeley.edu, Supercomputing’12

ScaLAPACK Parallel Library

ScaLAPACK SOFTWARE HIERARCHY

