

Lecture 13: Dense Linear Algebra II

David Bindel

8 Mar 2010

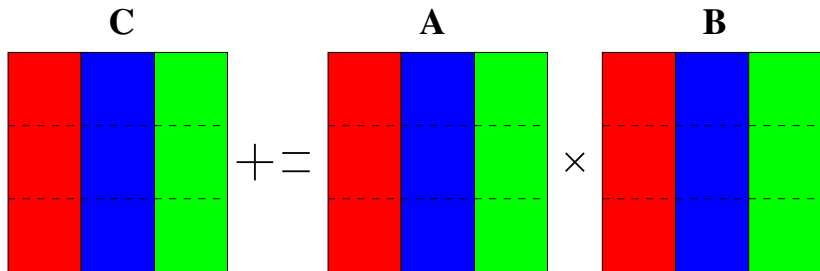
Logistics

- ▶ Tell me your project idea today (if you haven't already)!
- ▶ HW 2 extension to Friday
 - ▶ Meant to provide more flexibility, not more work!
 - ▶ See comments at start of last time about expectation
- ▶ HW 2 common issues
 - ▶ Segfault in binning probably means particle out of range
 - ▶ Particles too close together means either an interaction skipped or a time step too short

Review: Parallel matmul

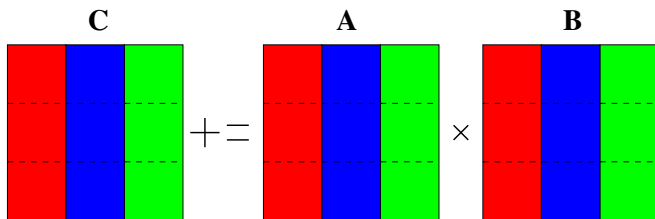
- ▶ Basic operation: $C = C + AB$
- ▶ Computation: $2n^3$ flops
- ▶ Goal: $2n^3/p$ flops per processor, minimal communication

1D layout



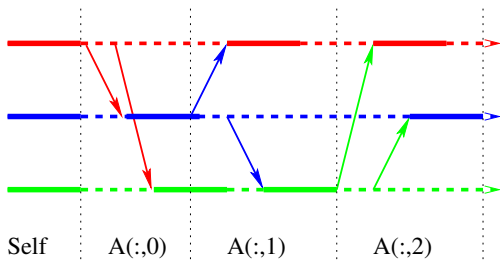
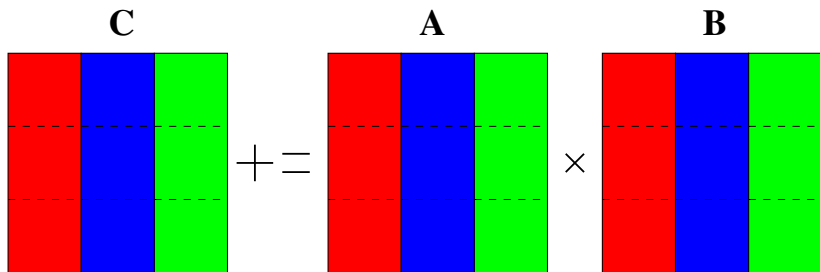
- ▶ Block MATLAB notation: $A(:,j)$ means j th block column
- ▶ Processor j owns $A(:,j)$, $B(:,j)$, $C(:,j)$
- ▶ $C(:,j)$ depends on *all* of A , but only $B(:,j)$
- ▶ How do we communicate pieces of A ?

1D layout on bus (no broadcast)



- ▶ Everyone computes local contributions first
- ▶ **P0** sends $A(:, 0)$ to each processor j in turn; processor j receives, computes $A(:, 0)B(0, j)$
- ▶ **P1** sends $A(:, 1)$ to each processor j in turn; processor j receives, computes $A(:, 1)B(1, j)$
- ▶ **P2** sends $A(:, 2)$ to each processor j in turn; processor j receives, computes $A(:, 2)B(2, j)$

1D layout on bus (no broadcast)



1D layout on bus (no broadcast)

```
C(:,myproc) += A(:,myproc)*B(myproc,myproc)
for i = 0:p-1
  for j = 0:p-1
    if (i == j)      continue;
    if (myproc == i) i
      send A(:,i) to processor j
    if (myproc == j)
      receive A(:,i) from i
      C(:,myproc) += A(:,i)*B(i,myproc)
    end
  end
end
end
```

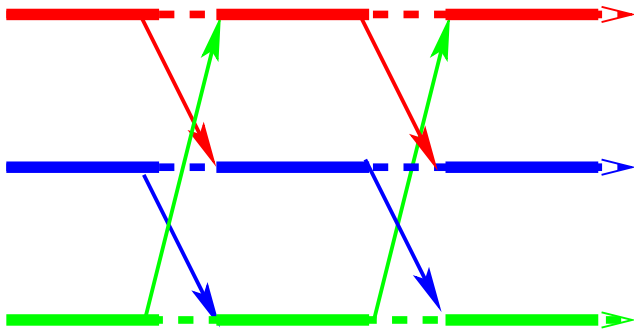
Performance model?

1D layout on bus (no broadcast)

No overlapping communications, so in a simple $\alpha - \beta$ model:

- ▶ $p(p - 1)$ messages
- ▶ Each message involves n^2/p data
- ▶ Communication cost: $p(p - 1)\alpha + (p - 1)n^2\beta$

1D layout on ring



- ▶ Every process j can send data to $j + 1$ simultaneously
- ▶ Pass slices of A around the ring until everyone sees the whole matrix ($p - 1$ phases).

1D layout on ring

```
tmp = A(myproc)
C(myproc) += tmp*B(myproc,myproc)
for j = 1 to p-1
    sendrecv tmp to myproc+1 mod p,
              from myproc-1 mod p
    C(myproc) += tmp*B(myproc-j mod p, myproc)
```

Performance model?

1D layout on ring

In a simple $\alpha - \beta$ model, at each processor:

- ▶ $p - 1$ message sends (and simultaneous receives)
- ▶ Each message involves n^2/p data
- ▶ Communication cost: $(p - 1)\alpha + (1 - 1/p)n^2\beta$

Outer product algorithm

Serial: Recall outer product organization:

```
for k = 0:s-1
    C += A(:,k)*B(k,:);
end
```

Parallel: Assume $p = s^2$ processors, block $s \times s$ matrices.
For a 2×2 example:

$$\begin{bmatrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{bmatrix} = \begin{bmatrix} A_{00}B_{00} & A_{00}B_{01} \\ A_{10}B_{00} & A_{10}B_{01} \end{bmatrix} + \begin{bmatrix} A_{01}B_{10} & A_{01}B_{11} \\ A_{11}B_{10} & A_{11}B_{11} \end{bmatrix}$$

- ▶ Processor for each $(i, j) \implies$ parallel work for each $k!$
- ▶ Note everyone in row i uses $A(i, k)$ at once, and everyone in row j uses $B(k, j)$ at once.

Parallel outer product (SUMMA)

```
for k = 0:s-1
  for each i in parallel
    broadcast A(i,k) to row
  for each j in parallel
    broadcast A(k,j) to col
  On processor (i,j), C(i,j) += A(i,k)*B(k,j);
end
```

If we have tree along each row/column, then

- ▶ $\log(s)$ messages per broadcast
- ▶ $\alpha + \beta n^2/s^2$ per message
- ▶ $2 \log(s)(\alpha s + \beta n^2/s)$ total communication
- ▶ Compare to 1D ring: $(p-1)\alpha + (1-1/p)n^2\beta$

Note: Same ideas work with block size $b < n/s$

Cannon's algorithm

$$\begin{bmatrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{bmatrix} = \begin{bmatrix} A_{00}B_{00} & A_{01}B_{11} \\ A_{11}B_{10} & A_{10}B_{01} \end{bmatrix} + \begin{bmatrix} A_{01}B_{10} & A_{00}B_{01} \\ A_{10}B_{00} & A_{11}B_{11} \end{bmatrix}$$

Idea: Reindex products in block matrix multiply

$$\begin{aligned} C(i, j) &= \sum_{k=0}^{p-1} A(i, k)B(k, j) \\ &= \sum_{k=0}^{p-1} A(i, k + i + j \bmod p) B(k + i + j \bmod p, j) \end{aligned}$$

For a fixed k , a given block of A (or B) is needed for contribution to *exactly one* $C(i, j)$.

Cannon's algorithm

```
% Move A(i, j) to A(i, i+j)
for i = 0 to s-1
    cycle A(i, :) left by i

% Move B(i, j) to B(i+j, j)
for j = 0 to s-1
    cycle B(:, j) up by j

for k = 0 to s-1
    in parallel;
        C(i, j) = C(i, j) + A(i, j)*B(i, j);
    cycle A(:, i) left by 1
    cycle B(:, j) up by 1
```

Cost of Cannon

- ▶ Assume 2D torus topology
- ▶ Initial cyclic shifts: $\leq s$ messages each ($\leq 2s$ total)
- ▶ For each phase: 2 messages each ($2s$ total)
- ▶ Each message is size n^2/s^2
- ▶ Communication cost: $4s(\alpha + \beta n^2/s^2) = 4(\alpha s + \beta n^2/s)$
- ▶ This communication cost is optimal!
... but SUMMA is simpler, more flexible, almost as good

Speedup and efficiency

Recall

$$\text{Speedup} := t_{\text{serial}}/t_{\text{parallel}}$$

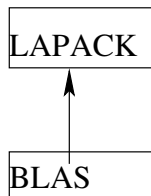
$$\text{Efficiency} := \text{Speedup}/p$$

Assuming no overlap of communication and computation, efficiencies are

$$\begin{array}{ll} \text{1D layout} & \left(1 + O\left(\frac{p}{n}\right)\right)^{-1} \\ \text{SUMMA} & \left(1 + O\left(\frac{\sqrt{p} \log p}{n}\right)\right)^{-1} \\ \text{Cannon} & \left(1 + O\left(\frac{\sqrt{p}}{n}\right)\right)^{-1} \end{array}$$

Reminder: Why matrix multiply?

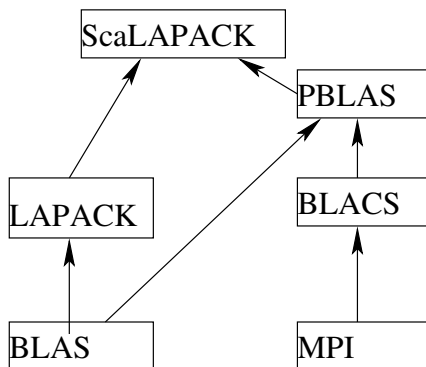
LAPACK structure



Build fast serial linear algebra (LAPACK) on top of BLAS 3.

Reminder: Why matrix multiply?

ScaLAPACK structure

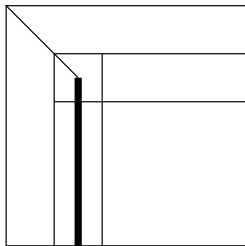


ScaLAPACK builds additional layers on same idea.

Reminder: Evolution of LU

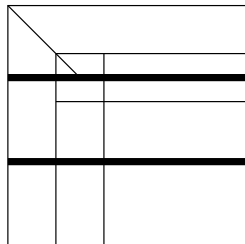
On board...

Blocked GEPP



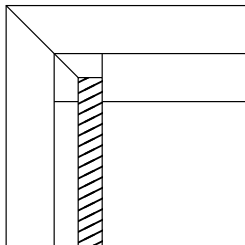
Find pivot

Blocked GEPP



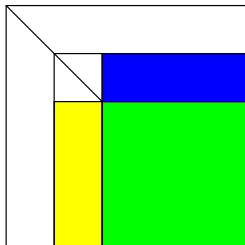
Swap pivot row

Blocked GEPP



Update within block

Blocked GEPP



Delayed update (at end of block)

Big idea

- ▶ *Delayed update* strategy lets us do LU fast
 - ▶ Could have also delayed application of pivots
- ▶ Same idea with other one-sided factorizations (QR)
- ▶ Can get decent multi-core speedup with parallel BLAS!
... assuming n sufficiently large.

There are still some issues left over (block size? pivoting?)...

Explicit parallelization of GE

What to do:

- ▶ *Decompose* into work chunks
- ▶ *Assign* work to threads in a balanced way
- ▶ *Orchestrate* the communication and synchronization
- ▶ *Map* which processors execute which threads

Possible matrix layouts

1D column blocked: bad load balance

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \end{bmatrix}$$

Possible matrix layouts

1D column cyclic: hard to use BLAS2/3

$$\begin{bmatrix} 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \\ 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \\ 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \\ 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \\ 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \\ 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \\ 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \\ 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \\ 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \end{bmatrix}$$

Possible matrix layouts

1D column block cyclic: block column factorization a bottleneck

$$\begin{bmatrix} 0 & 0 & 1 & 1 & 2 & 2 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 2 & 2 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 2 & 2 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 2 & 2 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 2 & 2 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 2 & 2 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 2 & 2 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 2 & 2 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 2 & 2 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 2 & 2 & 0 & 0 & 1 & 1 \end{bmatrix}$$

Possible matrix layouts

Block skewed: indexing gets messy

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 2 & 2 & 2 & 0 & 0 & 0 & 1 & 1 & 1 \\ 2 & 2 & 2 & 0 & 0 & 0 & 1 & 1 & 1 \\ 2 & 2 & 2 & 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 1 & 1 & 2 & 2 & 2 & 0 & 0 & 0 \\ 1 & 1 & 1 & 2 & 2 & 2 & 0 & 0 & 0 \\ 1 & 1 & 1 & 2 & 2 & 2 & 0 & 0 & 0 \end{bmatrix}$$

Possible matrix layouts

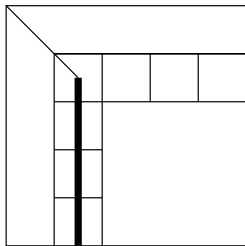
2D block cyclic:

$$\begin{bmatrix} 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 2 & 2 & 3 & 3 & 2 & 2 & 3 & 3 \\ 2 & 2 & 3 & 3 & 2 & 2 & 3 & 3 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 2 & 2 & 3 & 3 & 2 & 2 & 3 & 3 \\ 2 & 2 & 3 & 3 & 2 & 2 & 3 & 3 \end{bmatrix}$$

Possible matrix layouts

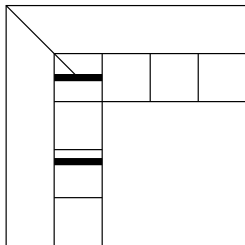
- ▶ 1D column blocked: bad load balance
- ▶ 1D column cyclic: hard to use BLAS2/3
- ▶ 1D column block cyclic: factoring column is a bottleneck
- ▶ Block skewed (a la Cannon): just complicated
- ▶ 2D row/column block: bad load balance
- ▶ 2D row/column block cyclic: win!

Distributed GEPP



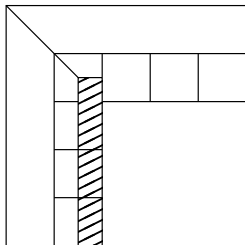
Find pivot (column broadcast)

Distributed GEPP



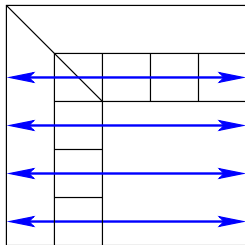
Swap pivot row within block column + broadcast pivot

Distributed GEPP



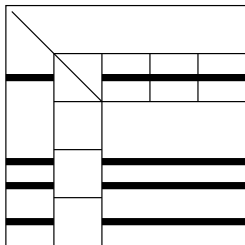
Update within block column

Distributed GEPP



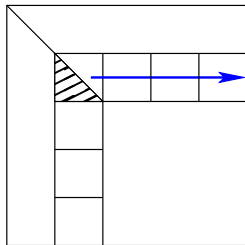
At end of block, broadcast swap info along rows

Distributed GEPP



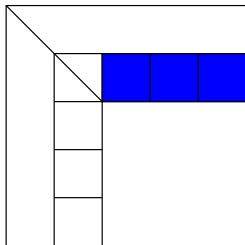
Apply all row swaps to other columns

Distributed GEPP



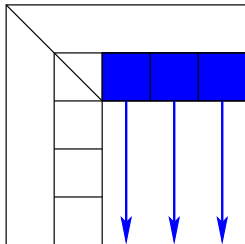
Broadcast block $L_{//}$ right

Distributed GEPP



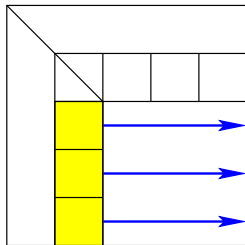
Update remainder of block row

Distributed GEPP



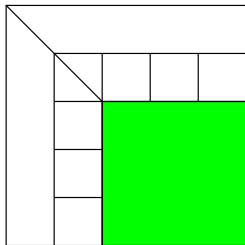
Broadcast rest of block row down

Distributed GEPP



Broadcast rest of block col right

Distributed GEPP



Update of trailing submatrix

Cost of ScaLAPACK GEPP

Communication costs:

- ▶ Lower bound: $O(n^2/\sqrt{P})$ words, $O(\sqrt{P})$ messages
- ▶ ScaLAPACK:
 - ▶ $O(n^2 \log P/\sqrt{P})$ words sent
 - ▶ $O(n \log p)$ messages
 - ▶ Problem: reduction to find pivot in each column
- ▶ Recent research on stable variants without partial pivoting

Onward!

Next up: Sparse linear algebra and iterative solvers!