

Lecture 3: Intro to parallel machines and models

David Bindel

1 Sep 2011

Logistics

Remember:

<http://www.cs.cornell.edu/~bindel/class/cs5220-f11/>
<http://www.piazza.com/cornell/cs5220>

- ▶ Note: the entire class will *not* be as low-level as lecture 2!
- ▶ Crocus cluster setup is in progress.
- ▶ If you drop/add, tell me so I can update CMS.
- ▶ Lecture slides are posted (in advance) on class web page.

A little perspective

“We should forget about small efficiencies, say about 97% of the time: premature optimization is the root of all evil.”

– C.A.R. Hoare (quoted by Donald Knuth)

- ▶ Best case: good algorithm, efficient design, *obvious code*
- ▶ Speed vs readability, debuggability, maintainability?
- ▶ A sense of balance:
 - ▶ Only optimize when needed
 - ▶ Measure before optimizing
 - ▶ Low-hanging fruit: data layouts, libraries, compiler flags
 - ▶ Concentrate on the bottleneck
 - ▶ Concentrate on inner loops
 - ▶ Get correctness (and a test framework) first

Matrix multiply

Consider naive square matrix multiplication:

```
#define A(i, j) AA[j*n+i]
#define B(i, j) BB[j*n+i]
#define C(i, j) CC[j*n+i]

for (i = 0; i < n; ++i) {
    for (j = 0; j < n; ++j) {
        C(i, j) = 0;
        for (k = 0; k < n; ++k)
            C(i, j) += A(i, k) * B(k, j);
    }
}
```

How fast can this run?

Note on storage

Two standard matrix layouts:

- ▶ Column-major (Fortran): $A(i,j)$ at $A+j*n+i$
- ▶ Row-major (C): $A(i,j)$ at $A+i*n+j$

I default to column major.

Also note: C doesn't really support matrix storage.

1000-by-1000 matrix multiply on my laptop

- ▶ Theoretical peak: 10 Gflop/s using both cores
- ▶ Naive code: 330 MFlops (3.3% peak)
- ▶ Vendor library: 7 Gflop/s (70% peak)

Tuned code is $20\times$ faster than naive!

Can we understand naive performance in terms of membench?

1000-by-1000 matrix multiply on my laptop

- ▶ Matrix sizes: about 8 MB.
- ▶ Repeatedly scans B in memory order (column major)
- ▶ 2 flops/element read from B
- ▶ 3 ns/flop = 6 ns/element read from B
- ▶ Check menchmark — gives right order of magnitude!

Simple model

Consider two types of memory (fast and slow) over which we have complete control.

- ▶ m = words read from slow memory
- ▶ t_m = slow memory op time
- ▶ f = number of flops
- ▶ t_f = time per flop
- ▶ $q = f/m$ = average flops / slow memory access

Time:

$$ft_f + mt_m = ft_f \left(1 + \frac{t_m/t_f}{q} \right)$$

Larger q means better time.

How big can q be?

1. Dot product: n data, $2n$ flops
2. Matrix-vector multiply: n^2 data, $2n^2$ flops
3. Matrix-matrix multiply: $2n^2$ data, $2n^3$ flops

These are examples of level 1, 2, and 3 routines in *Basic Linear Algebra Subroutines* (BLAS). We like building things on level 3 BLAS routines.

q for naive matrix multiply

$q \approx 2$ (on board)

Better locality through blocking

Basic idea: rearrange for smaller working set.

```
for (I = 0; I < n; I += bs) {
    for (J = 0; J < n; J += bs) {
        block_clear(&(C(I, J)), bs, n);
        for (K = 0; K < n; K += bs)
            block_mul(&(C(I, J)), &(A(I, K)), &(B(K, J)),
                      bs, n);
    }
}
```

Q: What do we do with “fringe” blocks?

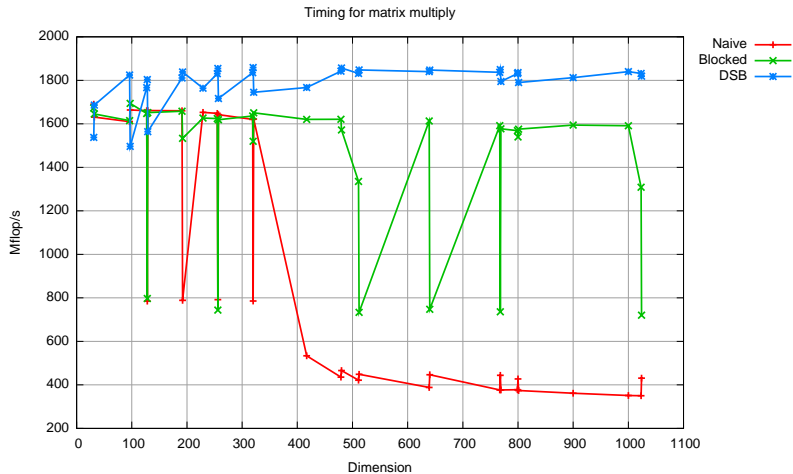
q for naive matrix multiply

$q \approx b$ (on board). If M_f words of fast memory, $b \approx \sqrt{M_f/3}$.

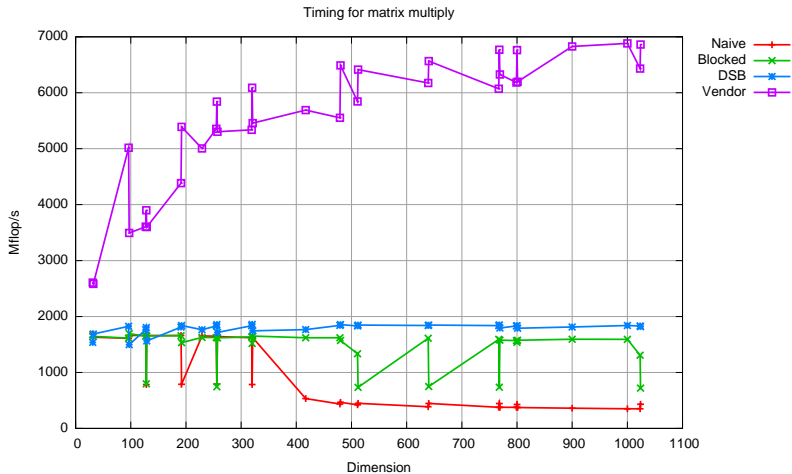
Th: (Hong/Kung 1984, Irony/Tishkin/Toledo 2004): Any reorganization of this algorithm that uses only associativity and commutativity of addition is limited to $q = O(\sqrt{M_f})$

Note: Strassen uses distributivity...

Better locality through blocking



Truth in advertising



Coming attractions

HW 1: You will optimize matrix multiply yourself!

Some predictions:

- ▶ You will make no progress without addressing memory.
- ▶ It will take you longer than you think.
- ▶ Your code will be rather complicated.
- ▶ Few will get anywhere close to the vendor.
- ▶ Some of you will be sold anew on using libraries!

Not all assignments will be this low-level.

Class cluster basics

`crocus.csuglab.cornell.edu` is a Linux Rocks cluster

- ▶ Six nodes (one head node, five compute nodes)
- ▶ Head node is virtual — *do not overload!*
- ▶ Compute nodes are dedicated — *be polite!*
- ▶ Batch submissions using Sun Grid Engine
- ▶ Read docs on assignments page

Class cluster basics

- ▶ Compute nodes are dual quad-core Intel Xeon E5504
- ▶ Nominal peak per core:
 - 2 SSE instruction/cycle ×
 - 2 flops/instruction ×
 - 2 GHz = 8 GFlop/s per core
- ▶ Caches:
 1. L1 is 32 KB, 4-way
 2. L2 is 256 KB (unshared) per core, 8-way
 3. L3 is 4 MB (shared), 16-way associativeL1 is relatively slow, L2 is relatively fast.
- ▶ Inter-node communication is switched gigabit Ethernet
- ▶ 16 GB memory per node

Cluster structure

Consider:

- ▶ Each core has vector parallelism
- ▶ Each chip has four cores, shares memory with others
- ▶ Each box has two chips, shares memory
- ▶ Cluster has five compute nodes, communicate via Ethernet

How did we get here? Why this type of structure? And how does the programming model match the hardware?

Parallel computer hardware

Physical machine has *processors, memory, interconnect*.

- ▶ Where is memory physically?
- ▶ Is it attached to processors?
- ▶ What is the network connectivity?

Parallel programming model

Programming *model* through languages, libraries.

- ▶ Control
 - ▶ How is parallelism created?
 - ▶ What ordering is there between operations?
- ▶ Data
 - ▶ What data is private or shared?
 - ▶ How is data logically shared or communicated?
- ▶ Synchronization
 - ▶ What operations are used to coordinate?
 - ▶ What operations are atomic?
- ▶ Cost: how do we reason about each of above?

Simple example

Consider dot product of x and y .

- ▶ Where do arrays x and y live? One CPU? Partitioned?
- ▶ Who does what work?
- ▶ How do we combine to get a single final result?

Shared memory programming model

Program consists of *threads* of control.

- ▶ Can be created dynamically
- ▶ Each has private variables (e.g. local)
- ▶ Each has shared variables (e.g. heap)
- ▶ Communication through shared variables
- ▶ Coordinate by synchronizing on variables
- ▶ Examples: OpenMP, pthreads

Shared memory dot product

Dot product of two n vectors on $p \ll n$ processors:

1. Each CPU evaluates partial sum (n/p elements, local)
2. Everyone tallies partial sums

Can we go home now?

Race condition

A race condition:

- ▶ Two threads access same variable, at least one write.
- ▶ Access are concurrent – no ordering guarantees
 - ▶ Could happen simultaneously!

Need synchronization via lock or barrier.

Race to the dot

Consider `S += partial_sum` on 2 CPU:

- ▶ P1: Load `S`
- ▶ P1: Add `partial_sum`
- ▶ P2: Load `S`
- ▶ P1: Store new `S`
- ▶ P2: Add `partial_sum`
- ▶ P2: Store new `S`

Shared memory dot with locks

Solution: consider `S += partial_sum` a *critical section*

- ▶ Only one CPU at a time allowed in critical section
- ▶ Can violate invariants locally
- ▶ Enforce via a lock or mutex (mutual exclusion variable)

Dot product with mutex:

1. Create global mutex `l`
2. Compute `partial_sum`
3. Lock `l`
4. `S += partial_sum`
5. Unlock `l`

Shared memory with barriers

- ▶ Lots of scientific codes have distinct phases (e.g. time steps)
- ▶ Communication only needed at end of phases
- ▶ Idea: synchronize on end of phase with *barrier*
 - ▶ More restrictive (less efficient?) than small locks
 - ▶ But much easier to think through! (e.g. less chance of deadlocks)
- ▶ Sometimes called *bulk synchronous programming*

Shared memory machine model

- ▶ Processors and memories talk through a bus
- ▶ Symmetric Multiprocessor (SMP)
- ▶ Hard to scale to lots of processors (think ≤ 32)
 - ▶ Bus becomes bottleneck
 - ▶ *Cache coherence* is a pain
- ▶ Example: Quad-core chips on cluster

Multithreaded processor machine

- ▶ May have more threads than processors! Switch threads on long latency ops.
- ▶ Called *hyperthreading* by Intel
- ▶ Cray MTA was one example

Distributed shared memory

- ▶ Non-Uniform Memory Access (NUMA)
- ▶ Can *logically* share memory while *physically* distributing
- ▶ Any processor can access any address
- ▶ Cache coherence is still a pain
- ▶ Example: SGI Origin (or multiprocessor nodes on cluster)

Message-passing programming model

- ▶ Collection of named processes
- ▶ Data is *partitioned*
- ▶ Communication by send/receive of explicit message
- ▶ Lingua franca: MPI (Message Passing Interface)

Message passing dot product: v1

Processor 1:

1. Partial sum s_1
2. Send s_1 to P2
3. Receive s_2 from P2
4. $s = s_1 + s_2$

Processor 2:

1. Partial sum s_2
2. Send s_2 to P1
3. Receive s_1 from P1
4. $s = s_1 + s_2$

What could go wrong? Think of phones vs letters...

Message passing dot product: v1

Processor 1:

1. Partial sum s_1
2. Send s_1 to P2
3. Receive s_2 from P2
4. $s = s_1 + s_2$

Processor 2:

1. Partial sum s_2
2. Receive s_1 from P1
3. Send s_2 to P1
4. $s = s_1 + s_2$

Better, but what if more than two processors?

MPI: the de facto standard

- ▶ Pro: *Portability*
- ▶ Con: least-common-denominator for mid 80s

The “assembly language” (or C?) of parallelism...
but, alas, assembly language can be high performance.

Distributed memory machines

- ▶ Each node has local memory
 - ▶ ... and no direct access to memory on other nodes
- ▶ Nodes communicate via network interface
- ▶ Example: our cluster!
- ▶ Other examples: IBM SP, Cray T3E

Why clusters?

- ▶ Clusters of SMPs are everywhere
 - ▶ Commodity hardware – economics! Even supercomputers now use commodity CPUs (though specialized interconnects).
 - ▶ Relatively simple to set up and administer (?)
- ▶ But still costs room, power, ...
- ▶ Will grid/cloud computing take over next?