

Parallel Programming

Kai Shen

1/26/2011 CSC 258/458 - Spring 2011 1

A Sidetrack

- Matrix multiplication

- Efficient use of cache
- Block-by-block multiplication

1/26/2011 CSC 258/458 - Spring 2011 2

Gaussian Elimination: Dependencies

- Simplification - ignore pivoting and final solving step
- Reduce an equation matrix into an equivalent upper-diagonal

for $i=1$ to N
 for $j=i+1$ to N
 zero out $A_{i,j}$ by adding $p*A_{i,i}$ to $A_{i,j}$

- Dependencies:**
 - Outer loop instances (e.g., $i=2$ depends on $i=1$)?
 - Inner loop instances (e.g., $j=3$ depends on $j=2$)?
 - Within one inner loop instance?

1/26/2011 CSC 258/458 - Spring 2011 3




Gaussian Elimination: Task Decomposition

for $i=1$ to N
 for $j=i+1$ to N
 zero out $A_{i,j}$ by adding $p*A_{i,i}$ to $A_{i,j}$

- Parallelism:**
 - Inner loop instances (row-wise)
 - Within one inner loop instance (column-wise)
- Task decomposition:**
 - Row, column, 2-dimensional

1/26/2011 CSC 258/458 - Spring 2011 4

Gaussian Elimination: Partitioning Granularity

- Row partitioning at different granularities
 - Row by row, cyclically
 
 - Every proc gets a contiguous chunk
 
 - Every proc gets a block of rows
 
- A good block size to efficiently utilize processor cache

1/26/2011 CSC 258/458 - Spring 2011 5

Comparison of Alternative Approaches

- Row partitioning, column partitioning, 2-dimensional partitioning
- Blocking at different sizes
- Static vs. dynamic task assignment
-

1/26/2011 CSC 258/458 - Spring 2011 6

Suggestions on Testing

- Incremental testing:
 - Try an artificial, small, dense input matrix first
 - Try to parallelize without partial pivoting first

1/26/2011 CSC 258/458 - Spring 2011 7

Irregular Parallelism

- Real problems contain large, sparse matrices
- Solve them as dense matrices waste time on zero-element operations
- Sparse matrix computation
 - Load imbalance
 - Managing nonzero fillins

1/26/2011 CSC 258/458 - Spring 2011 8

Example Speedup Results

- The speed ratio over the best sequential run

| Number of processors | Speedup |
|----------------------|---------|
| 1 | 1.0 |
| 2 | 2.0 |
| 4 | 3.5 |
| 8 | 5.5 |
| 12 | 6.5 |
| 16 | 7.0 |

| Number of processors | Speedup |
|----------------------|---------|
| 1 | 1.0 |
| 2 | 1.5 |
| 4 | 3.0 |
| 8 | 5.0 |
| 12 | 7.0 |
| 16 | 9.0 |

1/26/2011 CSC 258/458 - Spring 2011 9

Alternative Methods for Solving Linear Equations

- Faster solutions if we tolerate less precision
 - Less-than-precise pivoting
 - Iterative method

1/26/2011 CSC 258/458 - Spring 2011 10

Instructional-Level Parallelism

- Instructional-level parallelism
 - A single CPU core, but multiple functional units (arithmetic computation, floating-point computation, ...)
 - Fetch/decode multiple instructions at a time and execute them in parallel
 - possibly retire multiple instructions in a cycle
- Manage control dependencies
 - Stop before branch
 - Speculative execution on branch but rollback if mis-speculate
- Manage data dependencies
 - Recognize and observe dependencies
 - Resolve artificial dependencies through renaming

1/26/2011 CSC 258/458 - Spring 2011 11

Cache Coherence Problem

- In shared memory multiprocessors (except multithreading), each processor has a local cache

```

    graph TD
      P1[processor  
cache] --- Bus
      P2[processor  
cache] --- Bus
      Bus --- Mem[memory/more cache]
    
```

- For each data item in memory, additional copies may exist in processor local caches
 - after one processor updates the data, another processor's local copy may be incoherent
 - in cases of write-through and write-back caches

1/26/2011 CSC 258/458 - Spring 2011 12



Cache Coherence

- Coherence means the system semantics is the same as that of a system without multiple processor-local caches
- Multiprocessor cache coherent if there exists a hypothetical sequential order of all operations for each data location:
 - returned value in the read operation is that written by last write in the sequential order
 - the sequential order matches the order of operations from each processor

1/26/2011

CSC 258/458 - Spring 2011

13



Sequential Memory Consistency

- It means the memory semantics is the same as that of a uni-processor system
- Sequential consistent if there exists a hypothetical sequential order of all operations on all locations:
 - returned value in the read operation is that written by last write in the sequential order
 - the sequential order matches the order of operations from each processor

1/26/2011

CSC 258/458 - Spring 2011

14