

Lecture on Multicores

Darius Sidlauskas Post-doc



Danmarks Grundforskningsfond Danish National Research Foundation

Outline

- Part 1
 - Background
 - Current multicore CPUs
- Part 2
 - To share or not to share
- Part 3
 - Demo
 - War story



Darius Sidlauskas, 12/3-2013

mapalgo



Outline

Part 1

- Background
- Current multicore CPUs
- Part 2
 - To share or not to share
- Part 3
 - Demo
 - War story



Darius Sidlauskas, 12/3-2013

mapalgo



Software crisis

"The major cause of the software crisis is that the machines have become several orders of magnitude more powerful! To put it quite bluntly: as long as there were no machines, programming was no problem at all; when we had a few weak computers, programming became a mild problem, and now we have gigantic computers, programming has become an equally gigantic problem."

-- E. Dijkstra, 1972 Turing Award Lecture



Danmarks Grundforskningsfond Danish National Research Foundation

Darius Sidlauskas, 12/3-2013

Before..

- The 1st Software Crisis
 - When: around '60 and 70'
 - Problem: large programs written in assembly
 - Solution: abstraction and portability via high-level languages like C and FORTRAN
- The 2nd Software Crisis
 - When: around '80 and '90
 - Problem: building and maintaining large programs written by hundreds of programmers
 - Solution: software as a process (OOP, testing, code reviews, design patterns)
 - Also better tools: IDEs, version control, component libraries, etc.



Darius Sidlauskas, 12/3-2013

mapalgo -- - --



Recently..

- Processor-oblivious programmers
 - A Java program written on PC works on your phone
 - A C program written in '70 still works today and is faster
 - Moore's law takes care of good speedups



Darius Sidlauskas, 12/3-2013



Currently..

- Software crisis again?
 - When: 2005 and ...
 - Problem: sequential performance is stuck
 - Required solution: continuous and reasonable performance improvements
 - To process large datasets (BIG Data!)
 - To support new features
 - Without loosing portability and maintainability



Moore's law





Uniprocessor performance



Uniprocessor performance (cont.)



Uniprocessor performance (cont.)



MHZ

Uniprocessor performance (cont.)



MHZ

Why

- Power considerations
 - Consumption
 - Cooling
 - Efficiency
- DRAM access latency
 - Memory wall
- Wire delays
 - Range of wire in one clock cycle
- Diminishing returns of more instruction-level parallelism
 - Out-of-order execution, branch prediction, etc.



Darius Sidlauskas, 12/3-2013

mapalco -- - - - -

Overclocking [2]

- Air-water cooling: 4.8 GHz
 - Possible at home
- Phase change cooling: 5.89 GHz
- Liquid helium: 8.429 GHz
 - Made into Guinness Book of World Records!



Darius Sidlauskas, 12/3-2013



Shift to multicores

- Instead of going faster --> go more parallel!
 - Transistors are used now for multiple cores





Darius Sidlauskas, 12/3-2013

mapalgo ---

16/59

Multi-socket architecture



Current commercial multicore CPUs

Intel

- i7-3960X: 6-core (12 threads), 15M Cache, max 3.9 GHz
- Xeon E7-8870: 10-core (20 threads), 30M Cache, 2.8 GHz
- Xeon Phi: 60 cores (240 threads), 1.053 GHz, 320 GB/s memory bandwidth
- AMD
 - FX-8350: 8-core, 8M Cache, 4.0 GHz
 - Opteron 6274: 16-core, 16M Cache, 2.2 GHz
- Oracle
 - SPARC T4: 8-core (64 threads), 4M Cache, 3.0 GHz
- Cell



Darius Sidlauskas, 12/3-2013



Concurrency vs. Parallelism

Parallelism

- A condition that arises when at least two threads are executing simultaneously
- A specific case of concurrency
- Concurrency:
 - A condition that exists when at least two threads are making progress.
 - A more generalized form of parallelism
 - E.g., concurrent execution via time-slicing in uniprocessors (virtual parallelism)
- Distribution:
 - As above but running simultaneously on different

anaico

19/59

Danmarks Grundforskningsford machines (e.g., Conducts a a a started st

Amdhal's law

- Potential program speedup is defined by the fraction of code that can be parallelized
- Serial components rapidly become performance limiters as thread count increases
 - p fraction of work that can parallelized
 - n the number of processors



Amdhal's law

Danmarks





Cache coherence

Ensures the consistency between all the caches.





MESIF protocol

- Modified (M): present only in the current cache and dirty. A write-back to main memory will make it (E).
- Exclusive (E): present only in the current cache and clean. A read request will make it (S), a write-request will make it (M).
- Shared (S): maybe stored in other caches and *clean*.
 Maybe changed to (I) at any time.
- Invalid (I): unusable
- Forward (F): a specialized form of the S state



You've seen this..

L1 and L2 Cache Sizes



Read latencies on multicores [3]



Latency in nsec on 2-socket Intel Nehalem [3]



Read latencies on multicores [3]

	Exclusive cache lines			Modified cache lines			Shared cache lines			
Source	L1	L2	L3	L1	L2	L3	L1	L2	L3	RAM
Local	1.3 (4)	3.4 (10)	13.0 (38)	1.3 (4)	3.4 (10)	13.0 (38)	1.3 (4)	3.4 (10)	13.0 (38)	65.1
Core1 (on die)	22.2 (65)			28.3 (83)	25.5 (75)	15.0 (50)	13.0 (38)			05.1
Core4 (QPI)	63.4 (186)			102 - 109			58.0 (170)			106.0

Latency in nsec (cycles) on 2-socket Intel Nehalem [3]



Does it have effect in practice?

Processing 1600M tuples on 32-core machine [4]





mapalgo -----



Commandments [4]

- C1: Thou shalt not write thy neighbor's memory randomly – chunk the data, redistribute, and then sort/work on your data locally.
- C2: Thou shalt read thy neighbor's memory only sequentially – let the prefetcher hide the remote access latency.
- C3: Thou shalt not wait for thy neighbors don't use fine grained latching or locking and avoid synchronization points of parallel threads.



Outline

- Part 1
 - Background
 - Current multicore CPUs
- Part 2
 - To share or not to share?
- Part 3
 - Demo
 - War story



Darius Sidlauskas, 12/3-2013

mapalgo



Automatic contention detection and amelioration for data-intensive operations

- A generic framework (similar to Google's MapReduce) that
 - Efficiently parallelizes generic tasks
 - Automatically detects contention
 - Scales on multicore CPUs
 - Makes programmer's life easier :-)
- Based on
 - J. Cieslewicz, K. A. Ross, K. Satsumi, and Y. Ye.
 "Automatic contention detection and amelioration for data-intensive operations." In SIGMOD 2010.
 - Y. Ye, K. A. Ross, and N. Vesdapunt. Scalable aggregation on multicore processors. In *DaMoN* 2011



To Share or not to share

- Independent computation
 - Shared-nothing (disjoint processing)
 - No coordination (synchronization) overhead
 - No contention
 - Each thread use only 1/N of CPU resources
 - Merge step required
- Shared computation
 - Common data structures
 - Coordination (synchronization) overhead
 - Potential contention
 - All threads enjoy all CPU resources
 - No merge step required Darius Sidlauskas, 12/3-2013

Grundforskningsfo Danish National Research Foundati





Thread level parallelism

- On-chip coherency enables fine-grain parallelism
 - that was previously unprofitable (e.g., on SMPs)
- However, *beware*:
 - Correct parallel code does not mean no contention bottlenecks (hotspots)
 - Naive implementation can lead to huge performance pitfalls
 - Serialization due to shared access
 - E.g., many threads attempt to modify the same hash cell



Aggregate computation

Parallelizing simple DB operation:

```
SELECT R.G, count(*), sum(R.V)
FROM R
GROUP BY R.G
```

- What happens when values in R.G are highly skew?
- What happens when number of cores is much higher than |G|?
- Remember a key question: to share or not to share?
 Darius Sidlauskas, 12/3-2013



mapalgo -- -- --



Atomic CAS instruction

- Notation: CAS(&L, A, B)
- The meaning:
 - Compare the old value in location L with the expected old value A. If they are the same, then exchange the new value B with the value in location L.
 - Otherwise do not modify the value at location L because some other thread has changed the value at location L (since last time A was read). Return the current value of location L in B.
- After a CAS operation, one can determine whether the location L was successfully updated by comparing the contents of A and B.



Darius Sidlauskas, 12/3-2013

59

Atomic operations via CAS

```
atomic inc 64( &target ) {
    do {
       cur val = Load(&target);
       new val = curr + 1;
       ret val = CAS(&target, cur val, new val);
    } while (ret val != cur val);
  }
atomic dec 64( &target );
 atomic add 64( &target, value);
 atomic mul 64( &target, value);
. . .
```

Darius Sidlauskas, 12/3-2013

Datco



What is contention then?

Number of CAS retries



Darius Sidlauskas, 12/3-2013



mapalgo -- -- -.

Measuring contention (pseudo-code)

- my_atomic_inc_64(&target, &cas_counter) {
- do { cur val = Load(&target); new val = curr + 1; ret val = CAS(&target, cur val, new val); cas counter++; } while (ret val != cur val); } my atomic dec 64(&target, &cas counter); my atomic add 64(&target, value, &cas counter); my atomic mul 64(&target, value, &cas counter);



Darius Sidlauskas, 12/3-2013

paigo



Measuring contention (assembly code)

- .inline my_atomic_add_64,0! %o1 contains update value
- ldx [%00], %04 ! load current sum into %04; ld [%02], %05 ! load update-counter into %05 1: inc 1, %05 ! increment update-counter add %04, %01, %03 ! add value to current sum; put in %03 casx [%00], %04, %03 ! compare-and-swap %o3 into memory ! location of sum; ! %04 contains the value seen cmp %04, %03 ! check if compare-and-swap succeeded ! i.e., if %04 is equal to %03 bne,a,pn %xcc, 1b ! if not, retry loop starting at 1: mov %03, %04 ! statement executed even when branch ! taken; %04 now has a more recent value ! of the current sum and we have to add
 - ! %ol over again

! store the update-counter

- st %o5, [%o2]
- .end



Darius Sidlauskas, 12/3-2013

mapaico



Contention management

- Applies only to commutative operations
 - I.e., changing the order of the operands does not change the result
 - E.g., aggregation and partitioning
- General idea:
 - Perform operation on X and measure contention
 - Create extra version of X when contented
 - Spread the subsequent accesses among the two copies of X
 - Combine the results at the end





Framework

- Requires 4 user-defined template functions
 - create-clone: how a new version is created (x = 0)
 - combine: how multiple versions are merged (x + x1)
 - simple-update: how the new value of a data item is obtained from the current value and an update (x += v)
 - atomic-update: user defined function (next slide)
- Framework takes care
 - When to clone
 - Which clone is accessed by which thread



Example of atomic-update

 bool AggregatorAtomicUpdate(Aggregator *agg, const uint64_t value) {

```
int32_t cas_counter = 0;
```

```
my_atomic_inc_64(&agg->count, &cas_counter);
my_atomic_add_64(&agg->sum, value, &cas_counter);
return (3 < cas_counter);
```

```
}
```

```
Recall:
SELECT R.G, count(*), sum(R.V)
FROM R
GROUP BY R.G
```



Darius Sidlauskas, 12/3-2013



Techniques for managing contention

- Main concerns:
 - What information to maintain about the current number of clones?
 - How to map threads to clones in a balanced fashion?
- Two broad approaches for managing clones:
 - Global
 - Local





Managing clones globally

- New clones are created in shared address space
- Clone allocation happens in response to a single contention event (no threshold counters)
- The number of clones is always doubled
 - E.g., we can get to 64 clones of a heavy-hitter element after 6 contention steps
 - With few very popular items, each thread might end up having its own clone (no atomic operations needed afterwards!)



Managing clones locally

- Each thread creates clones in a local table used by that thread alone
- Table size is kept small
 - e.g., smaller than the thread's share of the L1 data cache
- When the table is full, new insertions are accomplished by spilling an existing value into the global data element



Darius Sidlauskas, 12/3-2013



Managing clones locally (cont.)

Danmarks

Danish National

Research Foundation



mapalgo



Experimental platforms

Platform	Sun T2	Intel Nehalem Xeon E5620		
Operating System	Solaris 10	Ubuntu Linux 2.6.32.25-server		
Processors	1	2		
Cores/processor (Threads/core)	8 (8)	4 (2)		
RAM	32GB	48 GB		
L1 Data Cache	8KB per core	32 KB per core		
L1 Inst. Cache	16KB per core	32 KB per core		
L2 Cache	4MB, 12-way Shared by 8 cores	256 KB per core		
L3 Cache	N.A.	12MB, 16-way Shared by 4 cores		

Danmarks Grundforskningsfond Danish National Research Foundation

Darius Sidlauskas, 12/3-2013



madalgo -- ---

Input data

- Refers to the characteristics of the group-by key in the input relation
- Synthetically generated distributions (N = 2²⁴):
 - Uniform
 - Sorted (1 1 1 2 3 3 4 5 ... N)
 - Heavy hitter (50%)
 - Repeated-run (1 2 3 ... N 1 2 3 ... N 1 2 ...)
 - Zipf (exponent of 0.5)
 - Self-similar (80-20 proportion)
 - Moving-cluster (locality window)
- During input generation a targeted group-by cardinality is specified

Danmarks Grundforskningsfond Danish National Research Foundation

Darius Sidlauskas, 12/3-2013

madalgo **-- -**- -



Cache and memory issues



Number of group by values where contention has been detected and at least one clone constructed





Results





Effects of the local table size





Conclusions

- Automatic contention detection
- Effective contention amelioration
- Both proposed schemes (global and local) mitigate contention
 - Global slightly faster
 - Local uses less memory
- However
 - Works just for commutative operations
 - Different architectures favor different approaches





Outline

- Part 1
 - Background
 - Current multicore CPUs
- Part 2
 - To share or not to share
- Part 3
 - Demo
 - War story



Darius Sidlauskas, 12/3-2013

mapalgo



Demo: false sharing

- Threads operate on different variables
- But variables reside on the same cache line







Darius Sidlauskas, 12/3-2013

madalgo ---

War story



Looking for a master thesis topic?

- ACM SIGMOD 2013 Programming Contest
- ACM SIGSPATIAL GIS 2013



Darius Sidlauskas, 12/3-2013

anaigo



References

- [1] Samuel H. Fuller and Lynette I. Millett, "The Future of Computing Performance: Game Over or Next Level?" The National Academies Press, 2010. [link]
- [2] AMD Showcases World's Fastest CPU [link]
- [3] D. Molka, D. Hackenberg, R. Schone, and M. S. Muller. "Memory performance and cache coherency effects on an intel nehalem multiprocessor system." In *PACT* 2009.
- [4] Albutiu, M. C., Kemper, A., & Neumann, T. "Massively parallel sortmerge joins in main memory multi-core database systems." *VLDB* 2012.
- [5] J. Cieslewicz, K. A. Ross, K. Satsumi, and Y. Ye. "Automatic contention detection and amelioration for data-intensive operations." In *SIGMOD* 2010.
- [6] Y. Ye, K. A. Ross, and N. Vesdapunt. "Scalable aggregation on multicore processors." In *DaMoN* 2011





Thank you

Darius Sidlauskas Post-doc

Contact: dariuss@madalgo.au.dk



Danmarks Grundforskningsfond Danish National Research Foundation

All in one [1]

