

5 October 2016

#### To read more...

This day's papers:

Scott, "Synchronization and Communication in the T3E Multiprocessor" C.mmp—A multi-mini-processor

Supplementary readings: Hennessy and Patterson, section 5.1–2

#### **Homework 1 Post-Mortem**

Almost all students had trouble with: associativity (TLB or cache) TLB size instruction cache size

#### Many not-great results on: latency and throughput

block size

#### HW1: Cache assoc.

#### From Yizhe Zhang's submission:

2



#### HW1: Cache assoc.

Example: 2-way assoc. 4-entry cache pattern: 0/1/2/3/0/1/2/3/...(0/4 misses)

	addresses
set 0 (addr mod 2 == 0)	0/2
set 1 (addr mod 2 $==$ 1)	1/3

pattern: 0/1/2/3/4/0/1/2/3/4/...(3/5 misses)

	addresses		
set 0 (addr mod 2 == 0)	0/2 (after 2)	2/4 (after 4)	4/0 (after 0)
set 1 (addr mod 2 == 1)	1/3		

pattern: 0/1/2/3/4/5/0/1/2/3/4/5/...(6/6 misses)

	addresses		
set 0 (addr mod 2 == 0)	0/2 (after 2)	2/4 (after 4)	4/0 (after 0)
set 1 (addr mod 2 == 1)	1/3 (after 3)	3/5 (after 5)	5/1 (after 1)

#### HW1: Cache assoc. nits

Problem: virtual != physical addresses

Solution 1: Hope addresses are contiguous (often true shortly after boot)

Solution 2: Special large page allocation functions

#### HW1: TLB associativity

Things which seem like they should work (and full credit):

Strategy 1 — same as for cache, but stride by page size

Strategy 1 — stride = TLB reach, how many fit

Strategy 2 — stride = TLB reach / guessed associativity

idea: will get all-misses with # pages = associativity



#### My TLB associativity results (L1)



#### My TLB associativity results (L2)



9

11

# **TLB benchmark: controlling cache behavior**

#### page table:

virtual page number	physical page number
1000	13248
1001	13248
1002	13248
1003	13248
1004	13248
1005	13248

# TLB benchmark: preventing overlapping loads

multiple parallel page table lookups

don't want that for measuring miss time

index = index + stride + array[value];
force dependency

also an issue for many other benchmarks

#### Instruction cache benchmarking

Approx two students successful or mostly successful

Obstacle one: variable length programs?

Obstacle two: aggressive prefetching

#### Variable length programs

Solution 1: Write program to generate source code many functions of different lengths plus timing code multimegabyte source files

Solution 2: Figure out binary for 'jmp to address' allocate memory copy machine code to region add return instruction call as function (cast to function pointer)

#### **Avoiding instruction prefetching**

Lots of jumps (unconditional branches)! Basically requires writing assembly/machine code Might measure branch prediction tables!

#### HW1: Choose two most popular

prefetching stride — see when increasing stride matches random pattern of same size

multicore/thread throughput — run MT code

large pages — straightforward if you can allocate large pages

12

#### **HW1: optimization troubles**

```
int array[1024 * 1024 * 128];
int foo(void) {
    for (int i = 0; i < 1024 * 1024 * 128; ++i) {</pre>
        array[i] = 1;
    }
}
unoptimized loop: gcc -s foo.c
.L3:
        -4(%rbp), %eax // load 'i'
 movl
 cltq
 movl
        $1, array(,%rax,4) // 4-byte store 'array[i]'
        $1, -4(%rbp) // load+add+store 'i'
 addl
        -4(\% rbp), %eax
                          // load 'i'
 movl
        $134217727, %eax
 cmpl
 jbe
         .L3
```

#### HW1: optimization troubles

```
int array[1024 * 1024 * 128];
int foo(void) {
    for (int i = 0; i < 1024 * 1024 * 128; ++i) {</pre>
        array[i] = 1;
    }
}
optimized loop: gcc -S -Ofast -march=native foo.c
.L4:
 addl
         $1, %eax // 'i' in register
 vmovdqa %ymm0, (%rdx) // 32-byte store
         $32, %rdx
 addq
 cmpl
         %ecx, %eax
          .L4
 jb
                                                    16
```

#### HW1: Misc issues

allowing overlap (no dependency/pointer chasing) hard to see cache latencies wrong for measuring latency but right thing for throughput

#### not trying to control physical addresses easy technique — large pages sometimes serious OS limitation

controlling measurement error is it a fluke? how can I tell?

#### Homework 2

checkpoint due Saturday Oct 15

using gem5, a processor simulator

analyzing statistics from 4 benchmark programs

you will need: a 64-bit Linux environment (VM okay, no GUI okay)

... or to build gem5 yourself

# multithreading before: multiple streams of execution within a processor shared almost everything (but extras) shared memory now: on multiple processors duplicated everything except... shared memory (sometimes) 19

# a philosophical question multiprocessor dividing line? dividing of machines

#### **C.mmp worries**

efficient networks

memory access conflicts

OS software complexity

user software complexity

#### **C.mmp worries**

efficient networks

memory access conflicts

OS software complexity

user software complexity

#### topologies for processor networks

crossbar

shared bus

 $\mathsf{mesh}/\mathsf{hypertorus}$ 

fat tree/Clos network

crossbar (approx. C.mmp switch)





# hypermesh/torus









#### trees (alternative)





# fat trees

don't need really thick/fast wires

take bundle of switches



# minimum bisection bandwidth

half of the CPUs communicate with other half what's the worst case:

crossbar: same as best case

```
fat tree, folded Clos: same as best case
or can be built with less (cheaper)
```

```
tree: 1/N of best (everything through root)
```

```
shared bus: 1/N of best (take turns)
```

hypertorus: in between

#### other network considerations

	crossbar	fat tree (full BW)	hypertorus	ł	
bandwidth	N	N	$\frac{d}{\sqrt{N}}$	]	
max hops	1	2k	$d\sqrt[d]{N}$	]	
# switches	1	k	N	(	
switch capacity	N	2k	2d	-	
short cables	no	no	yes	S	
non-asymptotic factors omitted					
N: number of CPUs					
k: switch capacity					
d: number of dimensions in hypertorus					
			3	3	

32

#### metereological simulation



#### barriers

wait for everyone to be done

two messages on each edge of tree



#### **C.mmp worries**

efficient networks

memory access conflicts

OS software complexity

user software complexity

#### memory access conflicts

assumption: memory distributed randomly may need to wait in line for memory bank makes extra processors less effective

### **T3E's solution**

local memories

explicit access to remote memories

programmer/compiler's job to decide

tools to help:

centrifuge — regular distribution across CPUs virtual processor numbers + mapping table

#### hiding memory latencies

C.mmp — don't; CPUs were too slow

T3E local memory — caches

T3E remote memory — many parallel accesses need 100s to hide multi-microsecond latencies

38

#### programming models

T3E — explicit accesses to remote memory programmer must find parallelism in accesses C.mmp — maybe OS chooses memories?

#### caching

C.mmp — read-only data only

T3E — local only remote accesses check cache next to memory



#### **C.mmp** synchronization

C.mmp: locking — exclusive use of OS resource concern about granularity:

too much overhead from locking?

too much overhead from waiting for lock?

### **T3E** synchronization

atomic operations — easy, executed at one location read-modify-write commands to remote memory compare-and-swap-if-equal read-and-add

#### compare-and-swap

}

}

}

```
memory[address] = new-value
```

#### locks with compare-and-swap

```
// compare_and_swap(address,
// expect_old_value,
// new_value) == 1 if changed
while (!compare_and_swap(&lock, 0, 1)) {}
// keep trying until value is 0
```

// use shared resource here

```
lock = 0; // release lock
```

46



#### lock-free/wait-free data structures

use atomic operations 'directly' (no lock)

very tricky to reason about

wait-free: program makes progress even if one thread stops

#### **T3E** alternative: message queues

atomic "add to remote queue" operation only keep retrying if remote queue was full

check own queue — repeatedly read local memory

no extra traffic

#### shared memory synchronization

usually 'just' atomic 'read-modify-write' operations compare-and-swap read-and-add

later: using these to implement locks

later: transactional memory: an alternative to read/modify/write

#### papers for next time

54

56

Goodman, "Using cache memory to reduce processor-memory traffic" seminal paper on cache coherency

Archibald and Baer, "Cache coherence protocols: evaluation using a multiprocessor simulation model several expansions on Goodman's work

#### **Snooping coherency trick**

Listen to other cache's requests to memory

Respond with data if requested even though request was for memory, not you

Allows write-back policy