Thread-Level Parallelism— **Parallel Programming**

Hung-Wei





Recap: What about "linked list"

Static instructions

LOOP: ld X10, 8(X10) addi X7, X7, 1 bne X10, X0, LOOP

Dynamic instructions



3

5

7

9

eu

6

8



Simultaneous multithreading





ld 1 2 add 3 4 bne ld 5 6 add \bigcirc • (8) bne ld 9 10 (11)add (12)

X1, 0(X10)addi X10, X10, 8 X20, X20, X1 X10, X2, LOOP X1, 0(X10) addi X10, X10, 8 X20, X20, X1 X10, X2, LOOP X1, 0(X10)addi X10, X10, 8 X20, X20, X1 bne X10, X2, LOOP





Concept of CMP



What software thinks about "multiprogramming" hardware



Coherency & Consistency

- Coherency Guarantees all processors see the same value for a variable/memory address in the system when the processors need the value at the same time
 - What value should be seen
- Consistency All threads see the change of data in the same order
 - When the memory operation should be done



Snooping Protocol



read miss/hit

What happens when we write in coherent caches?



Cache coherency

• Assuming that we are running the following code on a CMP with a cache coherency protocol, how many of the following outputs are possible? (a is initialized to 0 as assume we will output more than 10 numbers)

thread 1			
while(1) printf("%d ",a);	while(1) a++;		
① 0123456789			
② 1259368101213			
③ 111111164100			
④ 11111111100			
A. 0			
B. 1			
C. 2			
D. 3			
E. 4			

thread 2



- Performance/correctness in multiprogramming
- Dark Silicon and modern architectures

Performance/correctness issues in multiprogramming environments

False Sharing





```
Version L
void *threaded_vadd(void *thread_id)
                                               void *threaded_vadd(void *thread_id)
{
                                                {
 int tid = *(int *)thread_id;
                                                 int tid = *(int *)thread_id;
 int i;
                                                  int i;
                                                  for(i=tid*(ARRAY_SIZE/NUM_OF_THREADS);i<(tid+1)*(ARRAY_SIZE/NUM_OF_THREADS);i++)</pre>
 for(i=tid;i<ARRAY_SIZE;i+=NUM_OF_THREADS)</pre>
                                                  {
        c[i] = a[i] + b[i];
                                                      c[i] = a[i] + b[i];
                                                  }
  }
                                                  return NULL;
 return NULL;
                                                }
}
```





Version R



4Cs of cache misses

- 3Cs:
 - Compulsory, Conflict, Capacity
- Coherency miss:
 - A "block" invalidated because of the sharing among processors.



False sharing

- True sharing
 - Processor A modifies X, processor B also want to access X.
- False sharing
 - Processor A modifies X, processor B also want to access Y. However, Y is invalidated because X and Y are in the same block!

Performance comparison

 Comparing implementations of thread_vadd — L and R, please identify which one will be performing better and why Version L

```
void *threaded_vadd(void *thread_id)
                                                void *threaded vadd(void *thread id)
 int tid = *(int *)thread_id;
                                                  int tid = *(int *)thread_id;
 int i;
                                                  int i;
                                                  for(i=tid*(ARRAY_SIZE/NUM_OF_THREADS);i<(tid+1)*(ARRAY_SIZE/NUM_OF_THREADS);i++)</pre>
  for(i=tid;i<ARRAY_SIZE;i+=NUM_OF_THREADS)</pre>
        c[i] = a[i] + b[i];
                                                      c[i] = a[i] + b[i];
                                                  }
 return NULL;
                                                  return NULL;
```

A. L is better, because the cache miss rate is lower

- B. R is better, because the cache miss rate is lower
- C. L is better, because the instruction count is lower
- D. R is better, because the instruction count is lower
- E. Both are about the same

```
tids[i] = i;
```



Version R

Main thread for(i = 0 ; i < NUM_OF_THREADS ; i++)</pre>

pthread_create(&thread[i], NULL, threaded_vadd, &tids

for(i = 0 ; i < NUM_OF_THREADS ; i++)</pre> pthread_join(thread[i], NULL);

Possible scenarios



Why (0,0)?

- Processor/compiler may reorder your memory operations/ instructions
 - Coherence protocol can only guarantee the update of the same memory address
 - Processor can serve memory requests without cache miss first
 - Compiler may store values in registers and perform memory operations later
- Each processor core may not run at the same speed (cache misses, branch mis-prediction, I/O, voltage scaling and etc..)
- Threads may not be executed/scheduled right after it's spawned

Again — how many values are possible?

 Consider the given program. You can safely assume the caches are coherent. How many of the following outputs will you see?

#include <stdio.h>

(0, 0)② (0,1) ③ (1,0) **(1, 1)** A. 0 B. 1 C. 2 D. 3 E. 4

```
#include <stdlib.h>
#include <pthread.h>
#include <unistd.h>
volatile int a,b;
volatile int x,y;
volatile int f;
void* modifya(void *z) {
  a=1;
  x=b;
  return NULL;
void* modifyb(void *z) {
  b=1;
  y=a;
  return NULL;
```

int main() { int i; pthread_t thread[2]; pthread_join(thread[0], NULL); pthread_join(thread[1], NULL); fprintf(stderr,"(%d, %d)\n",x,y); return 0;

```
pthread_create(&thread[0], NULL, modifya, NULL);
pthread_create(&thread[1], NULL, modifyb, NULL);
```

fence instructions

- x86 provides an "mfence" instruction to prevent reordering across the fence instruction
- x86 only supports this kind of "relaxed consistency" model. You still have to be careful enough to make sure that your code behaves as you expected



thread 2

mfenceb=1 must occur/update before mfence

Take-aways of parallel programming

- Processor behaviors are non-deterministic
 - You cannot predict which processor is going faster
 - You cannot predict when OS is going to schedule your thread
- Cache coherency only guarantees that everyone would eventually have a coherent view of data, but not when
- Cache consistency is hard to support



Power/Energy Consumption & Dark Silicon

Hung-Wei Tseng



Power v.s. Energy

- Power is the direct contributor of "heat"
 - Packaging of the chip
 - Heat dissipation cost
- Energy = P * ET
 - The electricity bill and battery life is related to energy!
 - Lower power does not necessary means better battery life if the processor slow down the application too much

ergy! Dattery life if the

Power & Energy

- Regarding power and energy, how many of the following statements are correct?
 - ① Lowering the power consumption helps extending the battery life
 - ② Lowering the power consumption helps reducing the heat generation
 - ③ Lowering the energy consumption helps reducing the electricity bill
 - ④ A CPU with 10% utilization can still consume 33% of the peak power
 - A. 0
 - **B**. 1
 - C. 2



Power

Dynamic/Active Power

- The power consumption due to the switching of transistor states
- Dynamic power per transistor $P_{dynamic} \sim \alpha \times C \times V^2 \times f \times N$
 - α : average switches per cycle
 - C: capacitance
 - V: voltage
 - f: frequency, usually linear with V
 - N: the number of transistors



Static/Leakage Power

- The power consumption due to leakage transistors do not turn all the way off during no operation
- Becomes the dominant factor in the most advanced process technologies. 1000

$$P_{leakage} \sim N \times V \times e^{-V_t}$$

- N: number of transistors
- V: voltage
- V_t : threshold voltage where transistor conducts (begins to switch)



Figure 1: Leakage power becomes a growing problem as demands for more performance and functionality drive chipmakers to nanometer-scale process nodes (Source: IBS).



Dennardian Scaling

Given a scaling factor S

Parameter	Relation	Classical Scali
Power Budget		1
Chip Size		1
Vdd (Supply Voltage)		1/S
Vt (Threshold Voltage)	1/S	1/S
tex (oxide thickness)		1/S
W, L (transistor dimensions)		1/S
Cgate (gate capacitance)	WL/tox	1/S
Isat (saturation current)	WVdd/tox	1/S
F (device frequency)	lsat/(CgateVdd)	S
D (Device/Area)	1/(WL)	S ²
p (device power)	IsatVdd	1/S ²
P (chip power)	Dp	1
U (utilization)	1/P	1



Dennardian Broken

Given a scaling factor S

Parameter	Relation	Classical Scaling	Leakage Limited
Power Budget		1	1
Chip Size		1	1
Vdd (Supply Voltage)		1/S	1
Vt (Threshold Voltage)	1/S	1/S	1
tex (oxide thickness)		1/S	1/S
W, L (transistor		1/S	1/S
Cgate (gate capacitance)	WL/tox	1/S	1/S
Isat (saturation current)	WVdd/tox	1/S	1
F (device frequency)	Isat/(CgateVdd)	S	S
D (Device/Area)	1/(WL)	S ²	S ²
p (device power)	IsatVdd	1/S ²	1
P (chip power)	Dp	1	S ²
U (utilization)	1/P	1	1/S ²

Power consumption to light on all transistors

Chip								
1	1	1	1	1	1	1		
1	1	1	1	1	1	1		
1	1	1	1	1	1	1		
1	1	1	1	1	1	1		
1	1	1	1	1	1	1		
1	1	1	1	1	1	1		
1	1	1	1	1	1	1		

Dennardian Scaling

Chip

0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5

=50W

=49W

Dennardian Broken



=100W!

Dark Silicon and the End of Multicore Scaling

H. Esmaeilzadeh, E. Blem, R. St. Amant, K. Sankaralingam and D. Burger University of Washington, University of Wisconsin—Madison, University of Texas at Austin, Microsoft Research

More cores per chip, slower per core

Products	Solutions Support		intel			
		X Intel® Xeon® Processor E7-8890 v4	Intel® Xeon® Processor E7-8893 v4			
	Status	Launched	Launched			
	Launch Date 🧕	Q2''16	Q2'16			
	Lithography 🟮	14 nm	14 nm			
	Performance					
	# of Cores 1	24	4			
	# of Threads 🧿	48	8			
	Processor Base Frequency 🜖	2.20 GHz	3.20 GHz			
	Max Turbo Frequency 🟮	3.40 GHz	3.50 GHz			
	Cache 🚯	60 MB	60 MB			
	Bus Speed 🚯	9.6 GT/s	9.6 GT/s			
	# of QPI Links 🟮	3	3			
	TDP 🟮	165 W	140 W			

.

×	Intel® Xeon® Processor E7-8880 v4	×	
	Launched		
	Q2'16		
	14 nm		
	22		
	44		
	2.20 GHz		
	3.30 GHz		
	55 MB		

9.6 GT/s 3

150 W

What happens if power doesn't scale with process technologies?

- If we are able to cram more transistors within the same chip area (Moore's law continues), but the power consumption per transistor remains the same. Right now, if we power the chip with the same power consumption but put more transistors in the same area because the technology allows us to. How many of the following statements are true?
 - ① The power consumption per chip will increase
 - ² The power density of the chip will increase
 - Given the same power budget, we may not able to power on all chip area if we maintain the (3) same clock rate
 - ④ Given the same power budget, we may have to lower the clock rate of circuits to power on all chip area

52

- A. 0
- B. 1
- C. 2
- D. 3

E. 4

Announcement

- Final Review tonight 7pm-8:20pm @ WCH 143
- Homework #4 due 12/4
- iEval submission attach your "confirmation" screen, you get an extra/bonus homework
- Project due on 12/2
- Office hour for Hung-Wei this week MWF 1p-2p