# Multithreaded Architectures and Programming on Multithreaded Architectures (II)

Hung-Wei Tseng

# 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

ene

6

8





## **Recap: Simultaneous multithreading**



ld 1 2 add 3 4 bne ld 5 6 add  $\overline{\mathcal{O}}$ 8 bne ld 9 10 (11)add (12)

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

## **Recap: SMT**

- Improve the throughput of execution
  - May increase the latency of a single thread
- Less branch penalty per thread
- Increase hardware utilization
- Simple hardware design: Only need to duplicate PC/Register **Files**
- Real Case:
  - Intel HyperThreading (supports up to two threads per core)
    - Intel Pentium 4, Intel Atom, Intel Core i7
  - AMD RyZen (Zen microarchitecture)

### Wide-issue SS processor v.s. multiple narrower-issue SS processors



### **Recap: Concept of CMP**





## **Recap: SMT v.s. CMP**

- An SMT processor is basically a SuperScalar processor with multiple instruction front-end. Assume within the same chip area, we can build an SMT processor supporting 4 threads, with 6-issue pipeline, 64KB cache or — a CMP with 4x 2-issue pipeline & 16KB cache in each core. Please identify how many of the following statements are/is correct when running programs on these processors.
  - If we are just running one program in the system, the program will perform better on an SMT processor
     If we are running 4 applications simultaneously, the cache miss rates will be higher in the SMT processor

  - If we are running 4 applications simultaneously, the branch mis-prediction will be higher in the SMT processor (3)
  - ④ If we are running one program with 4 parallel threads, the cache miss rates will be higher in the SMT processor — it depends!
  - ⑤ If we are running one program with 4 parallel threads simultaneously, the branch mis-prediction will be longer in the SMT processor — it depends!
  - A. 1
  - B. 2
  - C. 3 The only thing we know for sure — if we don't parallel the program, it won't get any faster on CMP
  - D. 4
  - E. 5



— it depends!

# Architectural Support for Parallel Programming

## **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



## Simple cache coherency protocol

- Snooping protocol
  - Each processor broadcasts / listens to cache misses
- State associate with each block (cacheline)
  - Invalid
    - The data in the current block is invalid
  - Shared
    - The processor can read the data
    - The data may also exist on other processors
  - Exclusive
    - The processor has full permission on the data
    - The processor is the only one that has up-to-date data



### What happens when we write in coherent caches?



### **Observer**

```
thread 1
                                                                    thread 2
int loop;
                                                void* modifyloop(void *x)
                                                {
                                                  sleep(1);
int main()
                                                  printf("Please input a number:\n");
{
  pthread_t thread;
                                                  scanf("%d",&loop);
  loop = 1;
                                                  return NULL;
                                                }
  pthread_create(&thread, NULL, modifyloop,
NULL);
  while(loop == 1)
  {
    continue;
  }
  pthread_join(thread, NULL);
  fprintf(stderr,"User input: %d\n", loop);
  return ⊘;
}
```

### **Observer**

### prevents the compiler from putting the variable "loop" in the "register"

```
thread 1
                                                                     thread 2
volatile int loop;
                                                void* modifyloop(void *x)
                                                {
int main()
                                                  sleep(1);
                                                  printf("Please input a number:\n");
{
                                                  scanf("%d",&loop);
  pthread_t thread;
  loop = 1;
                                                  return NULL;
                                                }
  pthread_create(&thread, NULL, modifyloop,
NULL);
  while(loop == 1)
  {
    continue;
  }
  pthread_join(thread, NULL);
  fprintf(stderr,"User input: %d\n", loop);
  return ⊘;
}
```



- Parallel programming
- GPU

### **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		
	16	

### https://www.pollev.com/hungweitseng close in 1:30

thread 2

Coherency



## **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

### **Cache coherency**



### https://www.pollev.com/hungweitseng close in 1:30 **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)
 int tid = *(int *)thread_id;
 int i;
  for(i=tid;i<ARRAY_SIZE;i+=NUM_OF_THREADS)</pre>
        c[i] = a[i] + b[i];
  return NULL;
```

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

```
ł
 tids[i] = i;
```

### Version R

Ε

D



for(i = 0 ; i < NUM\_OF\_THREADS ; i++)</pre> pthread\_join(thread[i], NULL);

В

С

### Lv.s.R

```
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)
 int tid = *(int *)thread_id;
 int i;
  for(i=tid;i<ARRAY_SIZE;i+=NUM_OF_THREADS)</pre>
        c[i] = a[i] + b[i];
 return NULL;
```

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

## 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?

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

```
#include <stdio.h>
#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;

## https://www.pollev.com/hungweitseng close in 1:30

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

Consistency

А

С

D

Е

### **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



# **Alternative Parallel Architectures**

### **GPU (Graphics Processing Unit)**



These shaders need to be "programmable" to apply different rendering effects/algorithms (Phong shading, Gouraud shading, and etc...)



Raster **Operations** / Output merger



vec3 diffuseColor = texture(u\_diffuseTexture, o\_texcoords).rgb;

// combination of all components and diffuse color of the object

### https://www.pollev.com/hungweitseng close in 1:30 What do you want from a GPU?

- Given the basic idea of shading algorithms, how many of the following statements would fit the agenda of designing a GPU?
  - Many ALUs to process multiple pixels simultaneously
  - ② Low latency memory bus to supply pixels, vectors and textures
  - ③ High performance branch predictors
  - ④ Powerful ALUs to process many different kinds of operators



GPU

А

Such the presentation to see live content. For screen share software, share the entire screen. Get help at polycocorajep

С

D

Ε

## What do you want from a GPU?

- Given the basic idea of shading algorithms, how many of the following statements would fit the agenda of designing a GPU?
  - ① Many ALUs to process multiple pixels simultaneously 1920\*1080 pixels!
  - Low latency memory bus to supply pixels, vectors and textures Actually, high bandwidth since each pixel requires different L, N, R, V and we need to feed thousands of pixels simultaneously High performance branch predictors not really, the behavior is uniform across all pixels Powerful ALUs to process many different kinds of operators not really, we only need vector add, vector mul, vector div. Low frequency is OK since we have many threads
  - A. 0 In terms of latency, even for 120 frames, you still have 8ms latency to get everything done!
  - B. 1
  - C. 2
  - D. 3
  - E. 4



### **Nvidia GPU architecture Connect to PCIe system interconnect**



### **Inside each SMX**

		_			_			ins	uuuu	un ca	LITE	_		_	_	_		_	
Warp Scheduler			Warp Scheduler				Warp Scheduler					Warp Scheduler							
Di	spatc	h	Dispat	tch	D	ispatc ₽	h I	Dispat J	tch	Di	spatc	h	Dispat	ch	D	ispato	:h	Dispat	cl
				Regi	ster	File (	65,536	x 32-	bit G	K110	)   (1:	31,07	2 x32-k	oit Gl	<b>{210</b> )	)			
	<b>.</b>	+	+	+	+	÷	+	+	÷	+	÷	+	+	+	÷	+	+	+	
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	S
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	9
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	\$
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	5
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	
Core	Core	Core	Anit	to	ta	bre	Dinit		SFU	12	Cole	Core	9	2e	CC		es	LD/ST	5
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	5
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	5
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	I D/ST	SELL	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	
Coro	Coro	Coro	DR Unit	Coro	Coro	Coro	DB Unit	LD/ST	eeu	Coro	Coro	Core	DR Unit	Coro	Corro	Coro	DP Heit	LDICT	
core	Core	Core		Core	Core	Core	DP Unit	LD/ST	oru	Core	Core	Core	DP Offic	Core	Core	Core		LD/31	
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	~
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	-
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	
								Inter	conne	ect Net	work						01/04	<u></u>	
	(	04 K	B Shar	ea Me	emor	y / L1		PBa		(128	KB S	onare	a wem	ory /	LIC	ache	GK21	0)	
	T		-			<b>T</b>	40 K	D Re	au-O			acine				τ.		-	
	Tex		lex	<b>`</b>		Tex		lex	S		Tex		lex	\$		lex		lex	

Each of these performs the same operation, but each of these is also a "thread"

### **Nvidia GPU architecture**





### **AMD GPU Architecture**



50





### **A CU in an AMD GPU**





### https://www.pollev.com/hungweitseng close in 1:30

## **CPU v.s. GPU**

- Comparing the performance of solving the following set of problems using modern CPU and GPU architectures, how many can GPUs outperform CPUs?
  - ① Matrix multiplications
  - Minimum Spanning Trees (2)
  - ③ Shortest Path Problems
  - **Gaussian Elimination** (4)
  - A. 0
  - B. 1
  - C. 2
  - D. 3

### E. 4

CPUvsGPU

А



D

С

### **CPU v.s. GPU**

- Comparing the performance of solving the following set of problems using modern CPU and GPU architectures, how many can GPUs outperform CPUs?
  - ① Matrix multiplications
- Minimum Spanning Trees
  - Shortest Path Problems
  - ④ Gaussian Elimination
  - A. 0
  - B. 1
  - C. 2
  - D. 3
  - E. 4

### How things are connected









### New overhead/bottleneck emerges



Second Storage Devices







### **APU (Accelerated Processing Unit)**



### It's now very common in intel and AMD lineups





# https://www.pollev.com/hungweitseng close in 1:30

- Regarding the pros and cons of an APU, how many of the followings are correct
  - ① APU eliminates the need of moving data from DRAM to GPU device memory
  - The memory bandwidth of an APU is generally better than that of GPU device (2) memory
  - The total number of ALUs that an APU can provide must be fewer than a discrete (3) GPU given the same power budget and chip area
  - ④ A memory intensive GPU kernel can slowdown the performance of another CPU program

A. C	
B. 1	
C. 2	
D. 3	
E. 4	

APU

Such the presentation to see live content. For screen share software, share the entire screen. Get help at polycocorajep

В

А

С

D

Ε

### ΑΡυ

- Regarding the pros and cons of an APU, how many of the followings are correct
  - APU eliminates the need of moving data from DRAM to GPU device memory CPUs/GPUs in an APU are all connected to DRAM
     The memory bandwidth of an APU is generally better than that of GPU device
  - Not true, APU uses CPU memory bus that is optimized for "latencies" memory
  - The total number of ALUs that an APU can provide must be fewer than a discrete 3
  - GPU given the same power budget and chip area
     CPU cores are really big due to the dynamic scheduling logic
     A memory intensive GPU kernel can slowdown the performance of another CPU Because they have to compete with the DRAM bandwidth program
  - A. 0
  - B. 1
  - C. 2

### D. 3

E. 4

### Announcement

- Project due next Monday
- Last reading quiz due next Monday
- Assignment #4 due next Wednesday
- iEVAL, until 12/3
  - Please fill the survey to let us know your opinion!
  - Don't forget to take a screenshot of your submission and submit through iLearn it counts as a **full credit assignment**
  - We will drop your lowest 2 assignment grades
- Final Exam
  - Starting from 12/6 to 12/10 12:00pm, any consecutive 180 minutes you pick
  - Similar to the midterm, but more time and about 1.5x longer
  - Two of the questions will be comprehensive exam questions
  - Will release a sample final at the end of the last lecture

## Computer Science & Engineering



66