## Chapter 2: Machine Models

Elements of Parallel Computing

Eric Aubanel

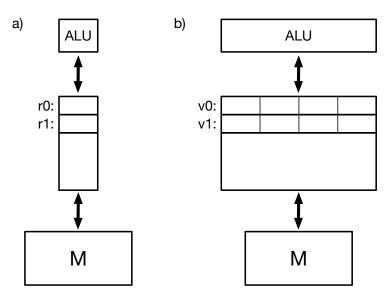
#### Outline

#### Three machine models:

- SIMD
- Shared memory
  - Multicore: latency-oriented
  - Manycore: Throughput-oriented (Understanding Throughput-Oriented Architectures, Garland and Kirk, CACM, Nov. 2010)
- Distributed memory

More than one model can be present in one computer

### **SIMD**



#### Scalar ALU: $c \leftarrow a + b$

1: 
$$r1 \leftarrow \text{load } a[0]$$
  
2:  $r2 \leftarrow \text{load } b[0]$   
3:  $r2 \leftarrow \text{add } r1, r2$   
4:  $c[0] \leftarrow \text{store } r2$   
5:  $r1 \leftarrow \text{load } a[1]$   
6:  $r2 \leftarrow \text{load } b[1]$   
7:  $r2 \leftarrow \text{add } r1, r2$   
8:  $c[1] \leftarrow \text{store } r2$ 

9: 
$$r1 \leftarrow \text{load } a[2]$$
  
10:  $r2 \leftarrow \text{load } b[2]$   
11:  $r2 \leftarrow \text{add } r1, r2$   
12:  $c[2] \leftarrow \text{store } r2$   
13:  $r1 \leftarrow \text{load } a[3]$   
14:  $r2 \leftarrow \text{load } b[3]$   
15:  $r2 \leftarrow \text{add } r1, r2$   
16:  $c[3] \leftarrow \text{store } r2$ 

#### SIMD ALU: $c \leftarrow a + b$

- 1:  $v1 \leftarrow vload a$
- 2:  $v2 \leftarrow vload b$
- 3:  $v2 \leftarrow \text{vadd } v1, v2$
- 4:  $c \leftarrow \text{vstore } v2$

#### SIM**T** $c \leftarrow a + b$

Nvidia GPUs: Single Instruction Multiple **Threads**: group (*warp*) of threads executes in SIMD fashion

```
t0: r00 \leftarrow \text{load } a[0] r10 \leftarrow \text{load } a[1] \cdots

t1: r01 \leftarrow \text{load } b[0] r11 \leftarrow \text{load } b[1] \cdots

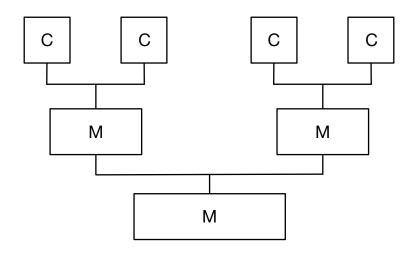
t2: r01 \leftarrow \text{add } r00, r01 \ r11 \leftarrow \text{add } r10, r11 \cdots

t3: c[0] \leftarrow \text{store } r01 c[1] \leftarrow \text{store } r11 \cdots
```

#### SIMT Model

- Programmer writes code for single thread
  - SPMD programming model (see later)
- Warps of threads execute in SIMD fashion
  - Not the same type of threads as in multicore processsor

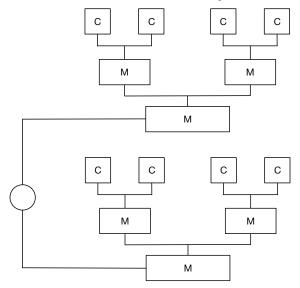
#### Multicore



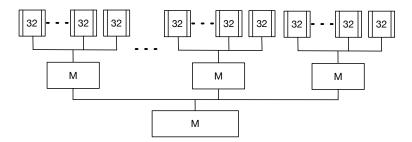
#### Multiprocessor or Multicomputer

- ▶ Multiprocessor: single address space
- Multicomputer: private address space per node

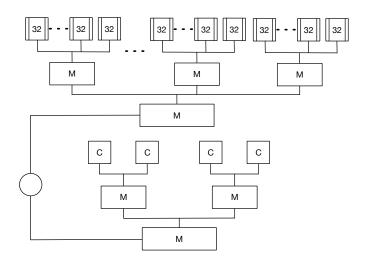
#### Multiprocessor or Multicomputer



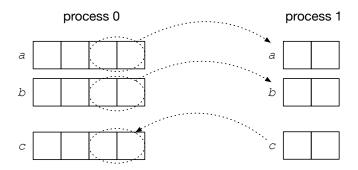
### Manycore



# Manycore/Multicore Hybrid



## Distributed Memory Execution



## Distributed Memory Execution

```
// process 0 // process 1

1 send a[2...3] to process 1 1 receive a[0...1] from process 0

2 send b[2...3] to process 1 2 receive b[0...1] from process 0

3 c[0] \leftarrow a[0] + b[0] 3 c[0] \leftarrow a[0] + b[0]

4 c[1] \leftarrow a[1] + b[1] 4 c[1] \leftarrow a[1] + b[1]

5 receive c[2...3] from process a send c[0...1] to process 0
```

## Shared Memory Execution

```
// thread 0 // thread 1 // thread r00 \leftarrow \text{load } a[0] r10 \leftarrow \text{load } a[1] r01 \leftarrow \text{load } b[0] r11 \leftarrow \text{load } b[1] \cdots r01 \leftarrow \text{add } r00, r01 r11 \leftarrow \text{add } r10, r11 \cdots c[0] \leftarrow \text{store } r01 c[1] \leftarrow \text{store } r11 \cdots
```

#### Not in lockstep!

## Pitfalls of Shared Memory Execution

- Data races
- Memory consistency
- Cache coherence

#### Watch out!

```
// thread 0 // thread 1 // thread r00 \leftarrow load \ a[0] \quad r10 \leftarrow load \ a[1] r01 \leftarrow load \ b[0] \quad r11 \leftarrow load \ b[1] \quad \cdots r01 \leftarrow add \ r00, r01 \ r11 \leftarrow add \ r10, r11 \quad \cdots c[0] \leftarrow store \ r01 \quad c[1] \leftarrow store \ r11 \quad \cdots r00 \leftarrow load \ c[3] x \leftarrow store \ r00
```

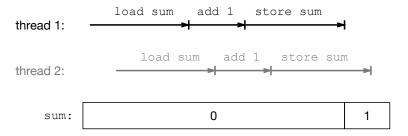
#### Data Race

One thread stores a value to a memory location while other threads load or store to the same location

#### Classic Data Race Example

```
// thread 0
                          // thread 1
sum \leftarrow sum + 1 sum \leftarrow sum + 1
At machine level:
// thread 0
                          // thread 1
r00 \leftarrow load sum
                          r10 \leftarrow load sum
r01 \leftarrow \text{add } r00.1 \qquad r11 \leftarrow \text{add } r10.1
sum \leftarrow store \ r01
                          sum \leftarrow store r11
```

#### Data Race Illustration



## Don't Roll Your Own Synchronization

```
// thread 0 // thread 3
L: r00 \leftarrow load \ flag \ c[3] \leftarrow store \ r31
If flag \neq 1 goto L flag \leftarrow store 1
r00 \leftarrow load \ c[3]
x \leftarrow store \ r00
Bad! Why?
```

### Possible $c \leftarrow a + b$ Interleaving

1: 
$$r00 \leftarrow load\ a[0]$$
 $r20, r21$ 

 2:  $r10 \leftarrow load\ a[1]$ 
 10:  $r01 \leftarrow add\ r00, r01$ 

 3:  $r20 \leftarrow load\ a[2]$ 
 $r30, r31$ 

 4:  $r30 \leftarrow load\ a[3]$ 
 12:  $r11 \leftarrow add$ 

 5:  $r01 \leftarrow load\ b[0]$ 
 13:  $c[0] \leftarrow store\ r01$ 

 6:  $r11 \leftarrow load\ b[1]$ 
 14:  $c[1] \leftarrow store\ r11$ 

 7:  $r21 \leftarrow load\ b[2]$ 
 15:  $c[2] \leftarrow store\ r21$ 

 8:  $r31 \leftarrow load\ b[3]$ 
 16:  $c[3] \leftarrow store\ r31$ 

 $r21 \leftarrow add$ 

#### Sequential Consistency

Each thread executes its instructions in program order. The order of the interleaving is arbitrary.

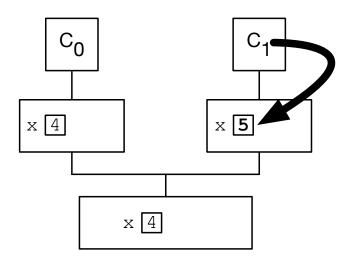
### Synchronization Race

Non-sequential (relaxed) consistency:

```
thread 3: flag \leftarrow \text{store } 1
thread 0: L: r00 \leftarrow \text{load}
flag
thread 0: If flag \neq 1 goto L
thread 0: r00 \leftarrow \text{load} c[3]
thread 0: x \leftarrow \text{store } r00
thread 3: c[3] \leftarrow \text{store } r31
```

# If no data (or synchronization) race, sequential consistency applies

#### Cache Coherence Problem



#### Cache Coherence

At any given time either any number of cores can read or only one core can write to a *cache block*, and in addition once a write occurs it is immediately visible to a subsequent read