

# Course Learning Outcomes

- Students shall demonstrate an understanding of the internal organization of a computer system through assembly language.
- Students shall design and simulate the data path and the control unit of a simple computer based on an instruction set.
- Students shall demonstrate an understanding of pipelining including instruction sequencing, register value forwarding, data interlocking.
- **Students shall demonstrate an understanding of the basic concepts of multiprocessor and multi-core designs.**
- Students shall demonstrate an understanding of the history **and possible future of the field necessary for staying at the forefront of computing systems development (life-long learning).**



# Course Learning Outcomes

- Students shall demonstrate an understanding of the internal organization of a computer system through assembly language.
- Students shall design and simulate the data path and the control unit of a simple computer based on an instruction set.
- Students shall demonstrate an understanding of pipelining including instruction sequencing, register value forwarding, data interlocking.
- Students shall demonstrate an understanding of the basic concepts of multiprocessor and multi-core designs.
- Students shall demonstrate an understanding of the history and possible future of the field necessary for staying at the forefront of computing systems development (life-long learning).



# Chapter 7

## Multicores, Multiprocessors, and Clusters

### Introduction

§9.1 Introduction

- Goal: connecting multiple computers to get higher performance
  - Multiprocessors
  - Scalability, availability, power efficiency
- Job-level (process-level) parallelism
  - High throughput for independent jobs
- Parallel processing program
  - Single program run on multiple processors
- Multicore microprocessors
  - Chips with multiple processors (cores)

# Hardware and Software

- Hardware
  - Serial: e.g., Pentium 4
  - Parallel: e.g., quad-core Xeon e5345
- Software
  - Sequential: e.g., matrix multiplication
  - Concurrent: e.g., operating system
- Sequential/concurrent software can run on serial/parallel hardware
  - Challenge: making effective use of parallel hardware



# Parallel Programming

- Parallel software is the problem
- Need to get significant performance improvement
  - Otherwise, just use a faster uniprocessor, since it's easier!
- Difficulties
  - Partitioning
  - Coordination
  - Communications overhead



# Amdahl's Law

- Sequential part can limit speedup
- Example: 100 processors, 90× speedup?
  - $T_{\text{new}} = T_{\text{parallelizable}}/100 + T_{\text{sequential}}$
  - $\text{Speedup} = \frac{1}{(1 - F_{\text{parallelizable}}) + F_{\text{parallelizable}}/100} = 90$
  - Solving:  $F_{\text{parallelizable}} = 0.999$
- Need sequential part to be 0.1% of original time



# Scaling Example

- Workload: sum of 10 scalars, and 10 × 10 matrix sum
  - Speed up from 10 to 100 processors
- Single processor: Time = (10 + 100) ×  $t_{\text{add}}$
- 10 processors
  - Time = 10 ×  $t_{\text{add}}$  + 100/10 ×  $t_{\text{add}}$  = 20 ×  $t_{\text{add}}$
  - Speedup = 110/20 = 5.5 (55% of potential)
- 100 processors
  - Time = 10 ×  $t_{\text{add}}$  + 100/100 ×  $t_{\text{add}}$  = 11 ×  $t_{\text{add}}$
  - Speedup = 110/11 = 10 (10% of potential)
- Assumes load can be balanced across processors



## Scaling Example (cont)

- What if matrix size is  $100 \times 100$ ?
- Single processor: Time =  $(10 + 10000) \times t_{\text{add}}$
- 10 processors
  - Time =  $10 \times t_{\text{add}} + 10000/10 \times t_{\text{add}} = 1010 \times t_{\text{add}}$
  - Speedup =  $10010/1010 = 9.9$  (99% of potential)
- 100 processors
  - Time =  $10 \times t_{\text{add}} + 10000/100 \times t_{\text{add}} = 110 \times t_{\text{add}}$
  - Speedup =  $10010/110 = 91$  (91% of potential)
- Assuming load balanced



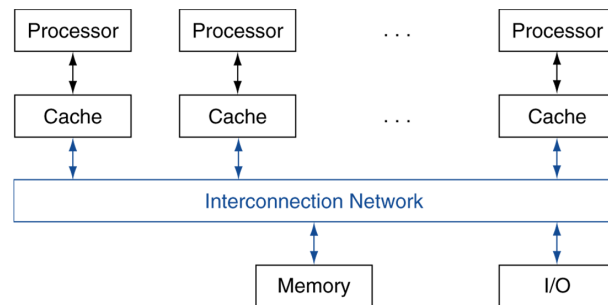
## Strong vs Weak Scaling

- Strong scaling: problem size fixed
  - As in example
- Weak scaling: problem size proportional to number of processors
  - 10 processors,  $10 \times 10$  matrix
    - Time =  $20 \times t_{\text{add}}$
  - 100 processors,  $32 \times 32$  matrix
    - Time =  $10 \times t_{\text{add}} + 1000/100 \times t_{\text{add}} = 20 \times t_{\text{add}}$
  - Constant performance in this example



# Shared Memory

- SMP: shared memory multiprocessor
  - Hardware provides single physical address space for all processors
  - Synchronize shared variables using locks
  - Memory access time
    - UMA (uniform) vs. NUMA (nonuniform)



## Example: Sum Reduction

- Sum 100,000 numbers on 100 processor UMA
  - Each processor has ID:  $0 \leq P_n \leq 99$
  - Partition 1000 numbers per processor
  - Initial summation on each processor
 

```
sum[Pn] = 0;
for (i = 1000*Pn;
    i < 1000*(Pn+1); i = i + 1)
  sum[Pn] = sum[Pn] + A[i];
```
- Now need to add these partial sums
  - Reduction: divide and conquer
  - Half the processors add pairs, then quarter, ...
  - Need to synchronize between reduction steps



# Example: Sum Reduction

```
half = 100;  
repeat
```

```
  synch();
```

```
  if (half%2 != 0 && Pn == 0)
```

```
    sum[0] = sum[0] + sum[half-1];
```

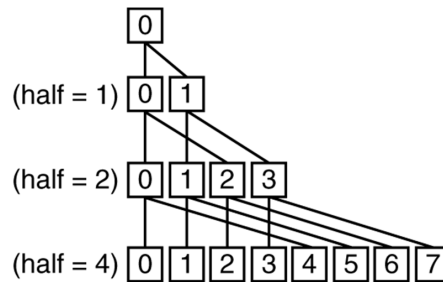
```
    /* Conditional sum needed when half is odd;
```

```
       Processor0 gets missing element */
```

```
  half = half/2; /* dividing line on who sums */
```

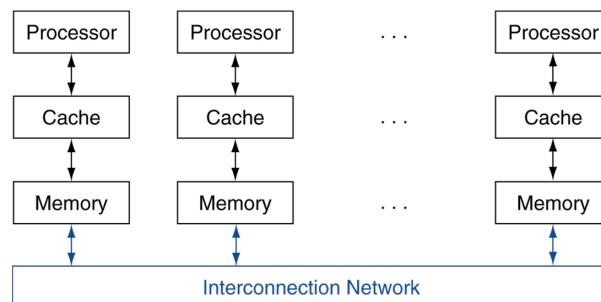
```
  if (Pn < half) sum[Pn] = sum[Pn] + sum[Pn+half];
```

```
until (half == 1);
```



# Message Passing

- Each processor has private physical address space
- Hardware sends/receives messages between processors



# Loosely Coupled Clusters

- Network of independent computers
  - Each has private memory and OS
  - Connected using I/O system
    - E.g., Ethernet/switch, Internet
- Suitable for applications with independent tasks
  - Web servers, databases, simulations, ...
- High availability, scalable, affordable
- Problems
  - Administration cost (prefer virtual machines)
  - Low interconnect bandwidth
    - c.f. processor/memory bandwidth on an SMP



# Grid Computing

- Separate computers interconnected by long-haul networks
  - E.g., Internet connections
  - Work units farmed out, results sent back
- Can make use of idle time on PCs
  - E.g., SETI@home, World Community Grid





# Multithreading

- Performing multiple threads of execution in parallel
  - Replicate registers, PC, etc.
  - Fast switching between threads
- Fine-grain multithreading
  - Switch threads after each cycle
  - Interleave instruction execution
  - If one thread stalls, others are executed
- Coarse-grain multithreading
  - Only switch on long stall (e.g., L2-cache miss)
  - Simplifies hardware, but doesn't hide short stalls (eg, data hazards)

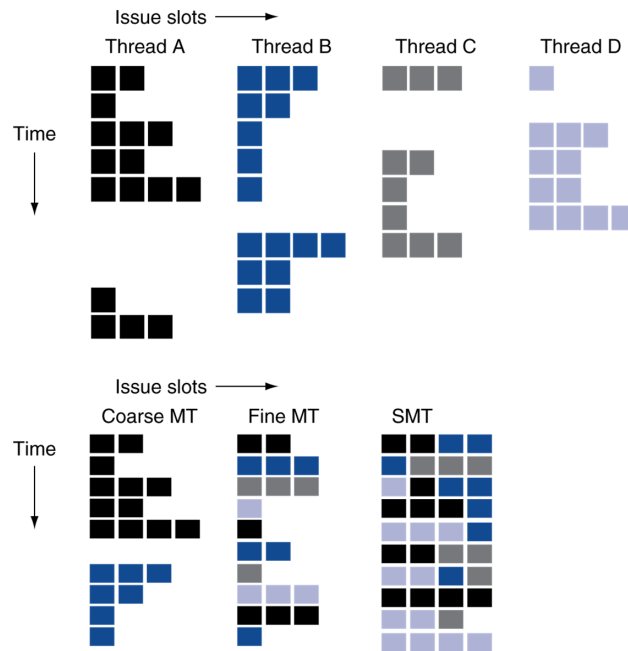


# Simultaneous Multithreading

- In multiple-issue dynamically scheduled processor
  - Schedule instructions from multiple threads
  - Instructions from independent threads execute when function units are available
  - Within threads, dependencies handled by scheduling and register renaming



# Multithreading Example



# Instruction and Data Streams

- An alternate classification

		Data Streams	
		Single	Multiple
Instruction Streams	Single	<b>SISD:</b> Intel Pentium 4	<b>SIMD:</b> SSE instructions of x86
	Multiple	<b>MISD:</b> No examples today	<b>MIMD:</b> Intel Xeon e5345

- SPMD: Single Program Multiple Data
  - A parallel program on a MIMD computer
  - Conditional code for different processors

§7.6 SISD, MIMD, SIMD, SPMD, and Vector



# History of GPUs

- Early video cards
  - Frame buffer memory with address generation for video output
- 3D graphics processing
  - Originally high-end computers (e.g., SGI)
  - Moore's Law  $\Rightarrow$  lower cost, higher density
  - 3D graphics cards for PCs and game consoles
- Graphics Processing Units
  - Processors oriented to 3D graphics tasks
  - Vertex/pixel processing, shading, texture mapping, rasterization

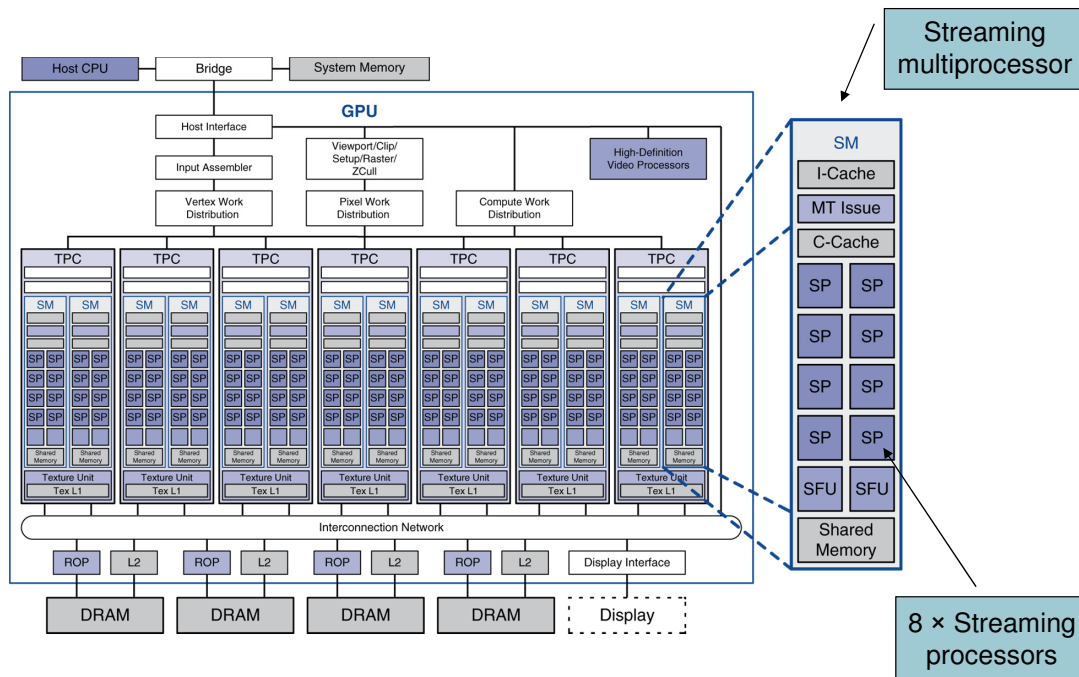


# GPU Architectures

- Processing is highly data-parallel
  - GPUs are highly multithreaded
  - Use thread switching to hide memory latency
    - Less reliance on multi-level caches
  - Graphics memory is wide and high-bandwidth
- Trend toward general purpose GPUs
  - Heterogeneous CPU/GPU systems
  - CPU for sequential code, GPU for parallel code
- Programming languages/APIs
  - DirectX, OpenGL
  - C for Graphics (Cg), High Level Shader Language (HLSL)
  - Compute Unified Device Architecture (CUDA)

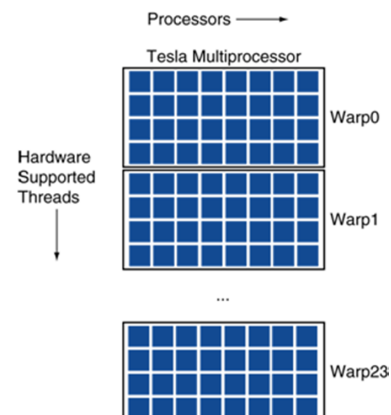


# Example: NVIDIA Tesla



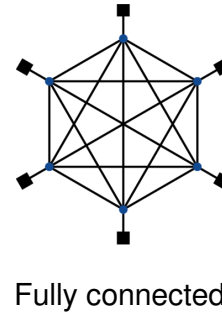
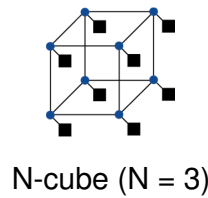
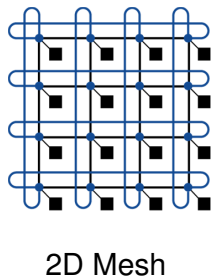
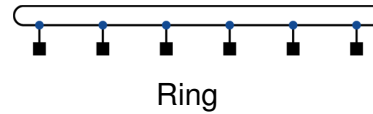
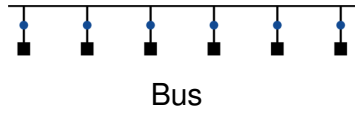
# Example: NVIDIA Tesla

- Streaming Processors
  - Single-precision FP and integer units
  - Each SP is fine-grained multithreaded
- Warp: group of 32 threads
  - Executed in parallel,
  - SIMD
    - 8 SPs
    - × 4 clock cycles

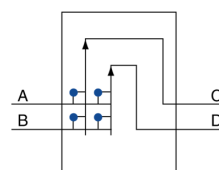
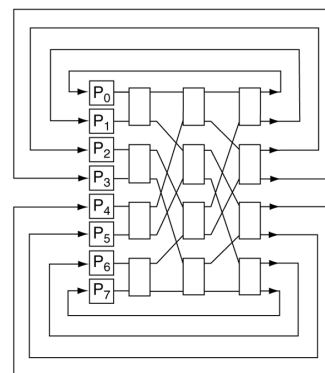
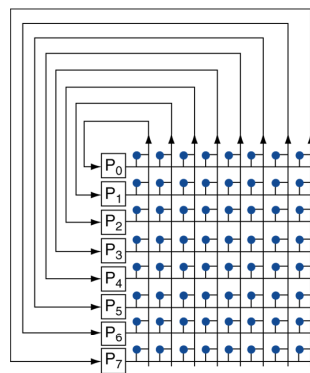


# Interconnection Networks

- Network topologies
  - Arrangements of processors, switches, and links



# Multistage Networks



# Network Characteristics

- Performance
  - Latency per message (unloaded network)
  - Throughput
  - Congestion delays (depending on traffic)
- Cost
- Power
- Routability in silicon



# Concluding Remarks

- Goal: higher performance by using multiple processors
- Difficulties
  - Developing parallel software
  - Devising appropriate architectures
- Many reasons for optimism
  - Changing software and application environment
  - Chip-level multiprocessors with lower latency, higher bandwidth interconnect
- An ongoing challenge for computer architects!

