



# Chapter 7

## Multicores, Multiprocessors, and Clusters

# Introduction

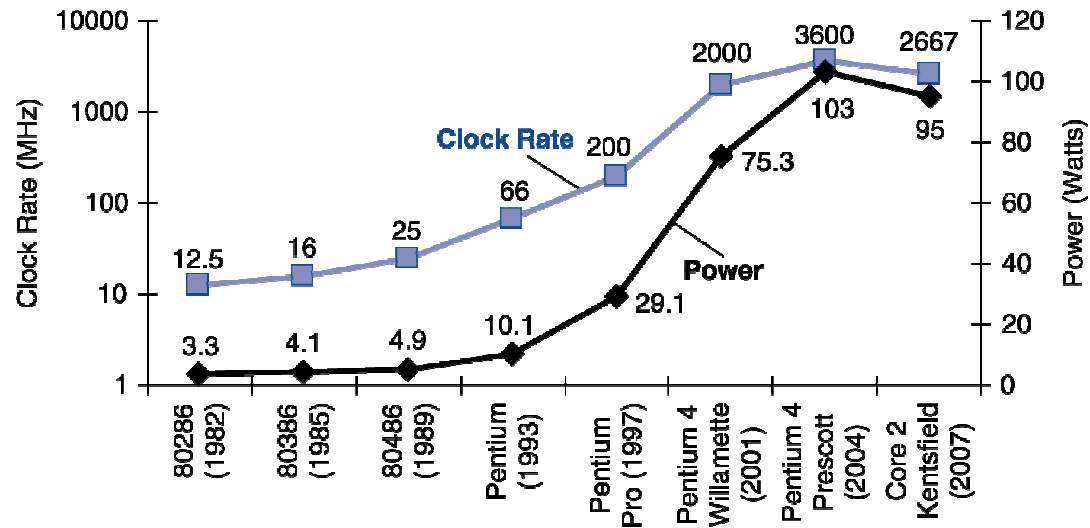
- Goal: connecting multiple computers to get higher performance
  - Multiprocessors
  - Scalability, availability, power efficiency
- Job-level (process-level) parallelism
  - High throughput for independent jobs
- Parallel processing program
  - Single program run on multiple processors
- Multicore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



# Recap

- Why processor designers moved towards multiprocessors / multicores?

# Power Trends



- In CMOS IC technology

$$\text{Power} = \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency}$$

×30

5V → 1V

×1000



# Reducing Power

- The power wall
  - We can't reduce voltage further
  - We can't remove more heat
- How else can we improve performance?

=>

- Multicore microprocessors
  - More than one processor per chip

# What We've Already Covered

- §2.11: Parallelism and Instructions
  - Synchronization
- §3.6: Parallelism and Computer Arithmetic
  - Associativity
- §4.10: Parallelism and Advanced Instruction-Level Parallelism
- §5.8: Parallelism and Memory Hierarchies
  - Cache Coherence [Lec.22]
- §6.9: Parallelism and I/O:
  - Redundant Arrays of Inexpensive Disks [L.27]



# Synchronization (from §2.11)

- Two processors sharing an area of memory
  - P1 writes, then P2 reads
  - Data race if P1 and P2 don't synchronize
    - Result depends of order of accesses
- Hardware support required
  - Atomic read/write memory operation
  - No other access to the location allowed between the read and write
- Could be a single instruction
  - E.g., atomic swap of register  $\leftrightarrow$  memory
  - Or an atomic pair of instructions



# Instruction-Level Parallelism (ILP) (§4.10)

- Pipelining: executing multiple instructions in parallel
- To increase ILP
  - Deeper pipeline
    - Less work per stage  $\Rightarrow$  shorter clock cycle
  - Multiple issue
    - Replicate pipeline stages  $\Rightarrow$  multiple pipelines
    - Start multiple instructions per clock cycle
    - $CPI < 1$ , so use Instructions Per Cycle (IPC)
    - E.g., 4GHz 4-way multiple-issue
      - 16 BIPS, peak  $CPI = 0.25$ , peak  $IPC = 4$
    - But dependencies reduce this in practice





# Multiple Issue

- Static multiple issue
  - Compiler groups instructions to be issued together
  - Packages them into “issue slots”
  - Compiler detects and avoids hazards
- Dynamic multiple issue
  - CPU examines instruction stream and chooses instructions to issue each cycle
  - Compiler can help by reordering instructions
  - CPU resolves hazards using advanced techniques at runtime

# Static Multiple Issue

- Compiler groups instructions into “issue packets”
  - Group of instructions that can be issued on a single cycle
  - Determined by pipeline resources required
- Think of an issue packet as a very long instruction
  - Specifies multiple concurrent operations
  - $\Rightarrow$  Very Long Instruction Word (VLIW)

# Dynamic Multiple Issue

- “Superscalar” processors
- CPU decides whether to issue 0, 1, 2, ... each cycle
  - Avoiding structural and data hazards
- Avoids the need for compiler scheduling
  - Though it may still help
  - Code semantics ensured by the CPU

# Why Do Dynamic Scheduling?

- Why not just let the compiler schedule code?
- Not all stalls are predictable
  - e.g., cache misses
- Can't always schedule around branches
  - Branch outcome is dynamically determined
- Different implementations of an ISA have different latencies and hazards

# Does Multiple Issue Work?

## The BIG Picture

- Yes, but not as much as we'd like
- Programs have real dependencies that limit ILP
- Some dependencies are hard to eliminate
  - e.g., pointer aliasing
- Some parallelism is hard to expose
  - Limited window size during instruction issue
- Memory delays and limited bandwidth
  - Hard to keep pipelines full
- Speculation can help if done well

# Cache Coherence Problem (§5.8)

- Suppose two CPU cores share a physical address space
  - Write-through caches

| Time step | Event               | CPU A's cache | CPU B's cache | Memory |
|-----------|---------------------|---------------|---------------|--------|
| 0         |                     |               |               | 0      |
| 1         | CPU A reads X       | 0             |               | 0      |
| 2         | CPU B reads X       | 0             | 0             | 0      |
| 3         | CPU A writes 1 to X | 1             | 0             | 1      |

# Coherence Defined

- Informally: Reads return most recently written value
- Formally:
  - $P$  writes  $X$ ;  $P$  reads  $X$  (no intervening writes)  
 $\Rightarrow$  read returns written value
  - $P_1$  writes  $X$ ;  $P_2$  reads  $X$  (sufficiently later)  
 $\Rightarrow$  read returns written value
    - c.f. CPU B reading  $X$  after step 3 in example
  - $P_1$  writes  $X$ ,  $P_2$  writes  $X$   
 $\Rightarrow$  all processors see writes in the same order
    - End up with the same final value for  $X$

# Cache Coherence Protocols

- Operations performed by caches in multiprocessors to ensure coherence
  - Migration of data to local caches
    - Reduces bandwidth for shared memory
  - Replication of read-shared data
    - Reduces contention for access
- Snooping protocols
  - Each cache monitors bus reads/writes
- Directory-based protocols
  - Caches and memory record sharing status of blocks in a directory



# Invalidating Snooping Protocols

- Cache gets exclusive access to a block when it is to be written
  - Broadcasts an invalidate message on the bus
  - Subsequent read in another cache misses
    - Owning cache supplies updated value

| CPU activity        | Bus activity     | CPU A's cache | CPU B's cache | Memory |
|---------------------|------------------|---------------|---------------|--------|
|                     |                  |               |               | 0      |
| CPU A reads X       | Cache miss for X | 0             |               | 0      |
| CPU B reads X       | Cache miss for X | 0             | 0             | 0      |
| CPU A writes 1 to X | Invalidate for X | 1             |               | 0      |
| CPU B read X        | Cache miss for X | 1             | 1             | 1      |



# Memory Consistency

- When are writes seen by other processors
  - “Seen” means a read returns the written value
  - Can’t be instantaneously
- Assumptions
  - A write completes only when all processors have seen it
  - A processor does not reorder writes with other accesses
- Consequence
  - P writes X then writes Y  
⇒ all processors that see new Y also see new X
  - Processors can reorder reads, but not writes

# 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

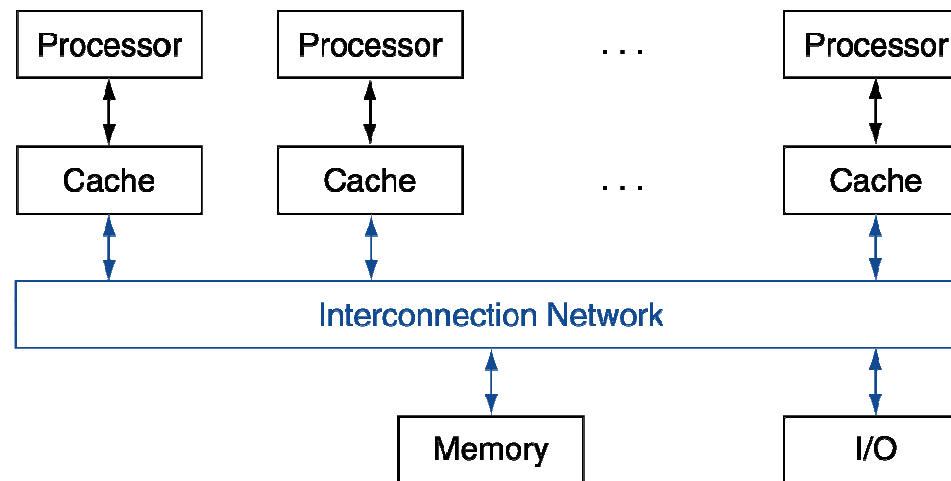
# Parallel Programming

- 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



# 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];
```

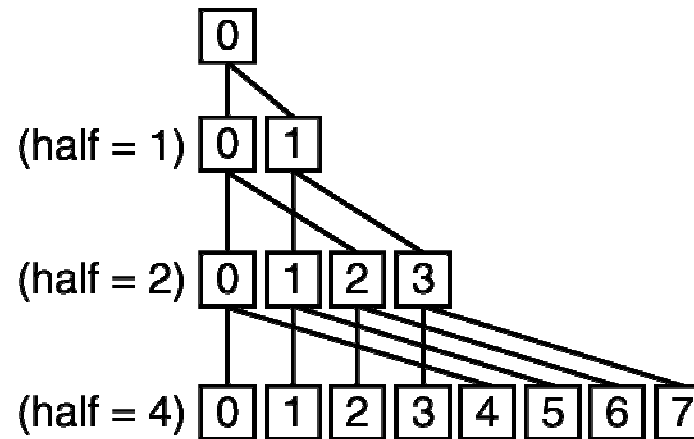
```
        /* Conditional sum needed when half is odd;
```

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

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

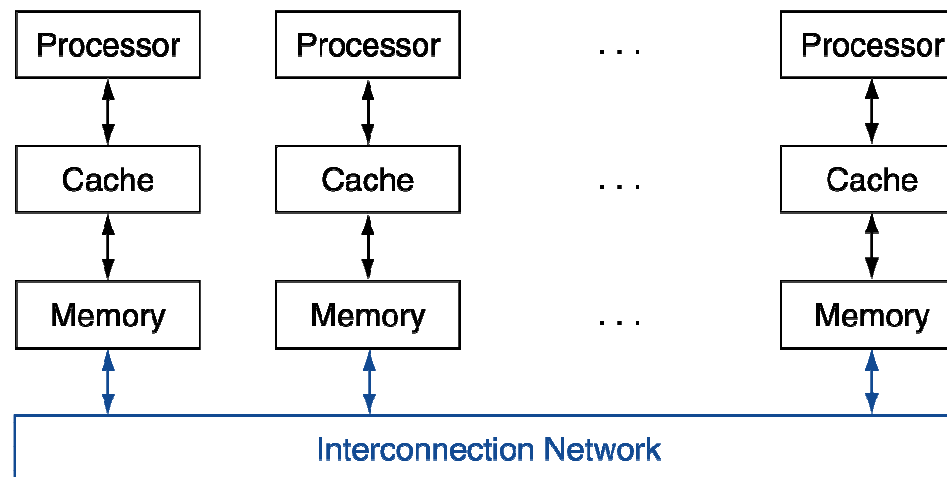
```
    if (Pn < half) sum[Pn] = sum[Pn] + sum[Pn+half];
```

```
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

# Sum Reduction (Again)

- Sum 100,000 on 100 processors
- First distribute 100 numbers to each
  - They do partial sums

```
sum = 0;  
for (i = 0; i < 1000; i = i + 1)  
    sum = sum + AN[i];
```
- Reduction
  - Half the processors send, other half receive and add
  - The quarter send, quarter receive and add, ...

# Sum Reduction (Again)

- Given send() and receive() operations

```
limit = 100; half = 100; /* 100 processors */
repeat
    half = (half+1)/2; /* send vs. receive
                        dividing line */
    if (Pn >= half && Pn < limit)
        send(Pn - half, sum);
    if (Pn < (limit/2))
        sum = sum + receive();
    limit = half; /* upper limit of senders */
until (half == 1); /* exit with final sum */
```

- Send/receive also provide synchronization
- Assumes send/receive take similar time to addition



# 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

# Processes vs. Threads

- Processes are independent,
  - threads are subsets of a process
- Processes carry considerably more state information than threads,
  - multiple threads within a process share process state as well as memory and other resources
- Processes have separate address spaces,
  - threads share their address space

# Processes vs. Threads (2)

- Processes interact only through system-provided inter-process communication (IPC) mechanisms
- Context switching between threads in the same process is typically faster than context switching between processes

# 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) – pipeline start-up costs

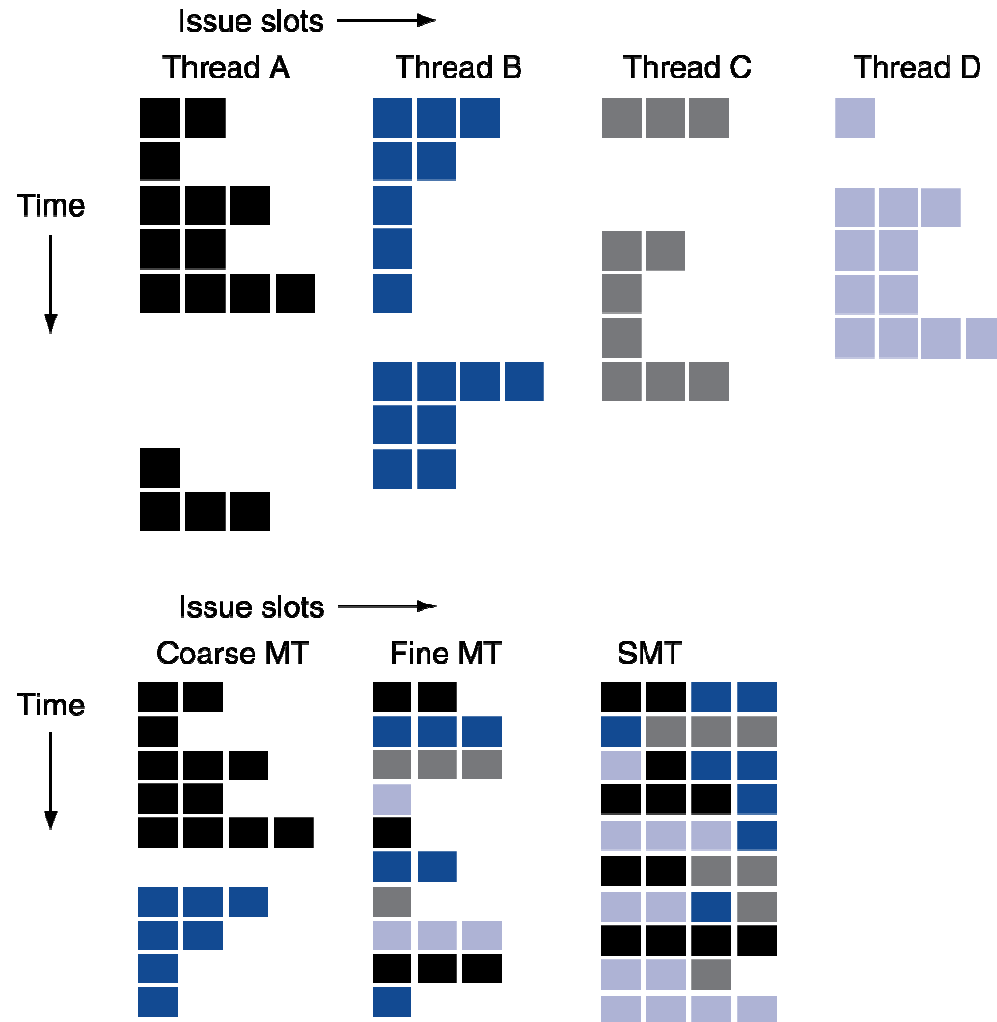


# Simultaneous Multithreading

- 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 etc
- Example: Intel Pentium-4 HT
  - Two threads: duplicated registers, shared function units and caches



# Multithreading Example



# Hyper-threading

- Hyper-threading works by
  - duplicating certain sections of the processor—(some registers), but not duplicating the main execution resources
  - appear as 2 "logical" processors to the host OS
  - allows the OS to schedule two threads or processes simultaneously
  - When execution resources are not used by the current task, they can be used to execute another scheduled task
    - The processor may stall due to a cache miss, branch mis-prediction, or data dependency.)

# Future of Multithreading

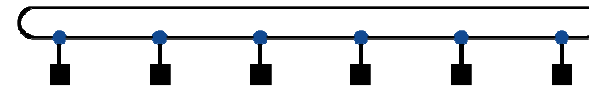
- Will it survive? In what form?
- Power considerations  $\Rightarrow$  simplified microarchitectures
  - Simpler forms of multithreading
- Tolerating cache-miss latency
  - Thread switch may be most effective
- Multiple simple cores might share resources more effectively

# Interconnection Networks

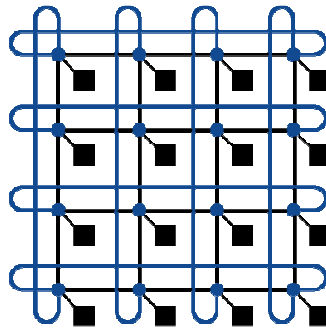
- Network topologies
  - Arrangements of processors, switches, and links



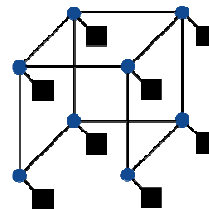
Bus



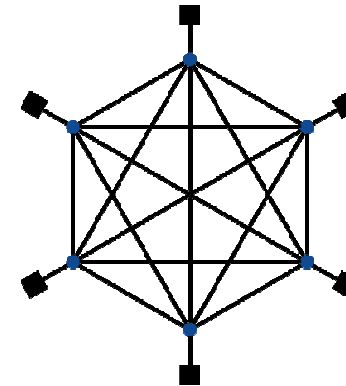
Ring



2D Mesh

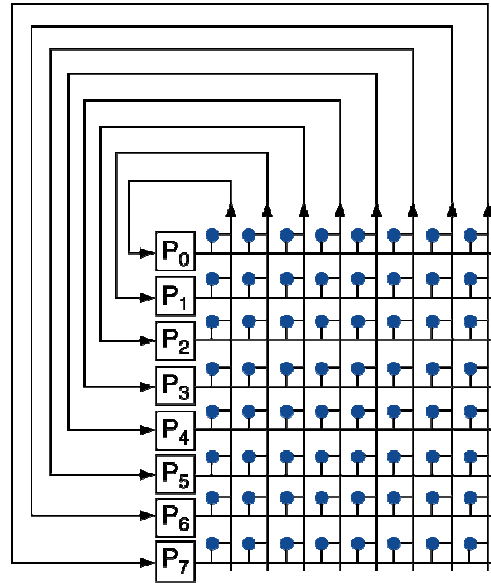


N-cube ( $N = 3$ )

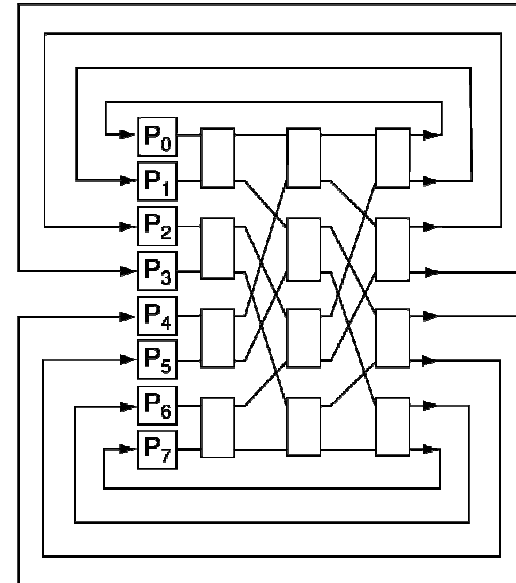


Fully connected

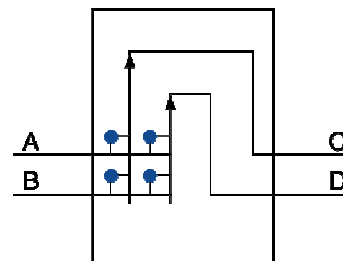
# Multistage Networks



a. Crossbar



b. Omega network



c. Omega network switch box

# Network Characteristics

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

# Instruction and Data Streams

## ■ An alternate classification

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

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

# SIMD

- Operate elementwise on vectors of data
  - E.g., MMX and SSE instructions in x86
    - Multiple data elements in 128-bit wide registers
- All processors execute the same instruction at the same time
  - Each with different data address, etc.
- Simplifies synchronization
- Reduced instruction control hardware
- Works best for highly data-parallel applications



# Vector Processors

- Highly pipelined function units
- Stream data from/to vector registers to units
  - Data collected from memory into registers
  - Results stored from registers to memory
- Example: Vector extension to MIPS
  - $32 \times 64$ -element registers (64-bit elements)
  - Vector instructions
    - `lv, sv`: load/store vector
    - `addv.d`: add vectors of double
    - `addvs.d`: add scalar to each element of vector of double
- Significantly reduces instruction-fetch bandwidth

# Vector vs. Scalar

- Vector architectures and compilers
  - Simplify data-parallel programming
  - Explicit statement of absence of loop-carried dependences
    - Reduced checking in hardware
  - Regular access patterns benefit from interleaved and burst memory
  - Avoid control hazards by avoiding loops
- More general than ad-hoc media extensions (such as MMX, SSE)
  - Better match with compiler technology

# 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!

