## CS 110 Computer Architecture Review Midterm II

http://shtech.org/courses/ca/

School of Information Science and Technology SIST

ShanghaiTech University

Slides based on UC Berkley's CS61C

# Midterm II

- Date: Thursday, May 17
- Time: 10:15 12:15 (similar to last time)
- Venue: Teaching Center 301 + 302
- Closed book:
  - You can bring <u>two</u> A4 pages with notes ((both sides; English preferred; Chinese is OK): Write your Chinese and **Pinyin** name on the top!
  - This time **HANDWRITTEN** only!
  - You will be provided with the MIPS "green sheet"
  - No other material allowed!

# Midterm II

- Switch cell phones off! (not silent mode off!)
   Put them in your bags.
- Bags to the front. Nothing except paper, pen, 1 drink, 1 snack on the table!
- No other electronic devices are allowed!
   No ear plugs, music, ...
- Anybody touching any electronic device will FAIL the course!
- Anybody found cheating (copy your neighbors answers, additional material, ...) will FAIL the course!















#### Do not look inside the midterm until we tell you to (until the time starts) !!

# Midterm II

- Ask questions today!
- This review session does not/ can not cover all possible topics!

• Topics: SDS till Data-Level-Parallelism

# Synchronous Digital Systems

Hardware of a processor, such as the MIPS, is an example of a Synchronous Digital System

Synchronous:

- All operations coordinated by a central clock
  - "Heartbeat" of the system!

Digital:

- Represent all values by discrete values
- Two binary digits: 1 and 0
- Electrical signals are treated as 1's and 0's
  - 1 and 0 are complements of each other
- High /low voltage for true / false, 1 / 0

## **Combinational Logic Symbols**

 Common combinational logic systems have standard symbols called logic gates



## **Boolean Algebra**

• Use plus "+" for OR

- "logical sum" 1+0 = 0+1 = 1 (True); 1+1=2 (True); 0+0 = 0 (False)

- Use product for AND (a•b or implied via ab)
  - "logical product"
    0\*0 = 0\*1 = 1\*0 = 0 (False); 1\*1 = 1 (True)
- "Hat" to mean complement (NOT)
- Thus

 $ab + a + \overline{c}$ 

- $= a \cdot b + a + \overline{c}$
- = (a AND b) OR a OR (NOT c)







Exhaustive list of the output value generated for each combination of inputs

How many logic functions can be defined with N inputs?

a	b	c	d	У
0	0	0	0	F(0,0,0,0)
0	0	0	1	F(0,0,0,1)
0	0	1	0	F(0,0,1,0)
0	0	1	1	F(0,0,1,1)
0	1	0	0	F(0,1,0,0)
0	1	0	1	F(0,1,0,1)
0	1	1	0	F(0,1,1,0)
0	1	1	1	F(0,1,1,1)
1	0	0	0	F(1,0,0,0)
1	0	0	1	F(1,0,0,1)
1	0	1	0	F(1,0,1,0)
1	0	1	1	F(1,0,1,1)
1	1	0	0	F(1,1,0,0)
1	1	0	1	F(1,1,0,1)
1	1	1	0	F(1,1,1,0)
1	1	1	1	F(1,1,1,1)

#### Boolean Algebra: Circuit & Algebraic Simplification



original circuit

equation derived from original circuit

algebraic simplification

simplified circuit

# Representations of Combinational Logic (groups of logic gates)



#### Question

• Simplify  $Z = A + BC + \overline{A(BC)}$ 

- A: Z = 0
- B:  $Z = \overline{A(1 + BC)}$
- C: Z = (A + BC)
- D: Z = BC
- E: Z = 1



# Type of Circuits

- *Synchronous Digital Systems* consist of two basic types of circuits:
  - Combinational Logic (CL) circuits
    - Output is a function of the inputs only, not the history of its execution
    - E.g., circuits to add A, B (ALUs)
  - Sequential Logic (SL)
    - Circuits that "remember" or store information
    - aka "State Elements"
    - E.g., memories and registers (Registers)

#### **Uses for State Elements**

- Place to store values for later re-use:
  - Register files (like \$1-\$31 in MIPS)
  - Memory (caches and main memory)
- Help control flow of information between combinational logic blocks
  - State elements hold up the movement of information at input to combinational logic blocks to allow for orderly passage

#### Recap of Timing Terms

- Clock (CLK) steady square wave that synchronizes system
- Setup Time when the input must be stable <u>before</u> the rising edge of the CLK
- Hold Time when the input must be stable <u>after</u> the rising edge of the CLK
- "CLK-to-Q" Delay how long it takes the output to change, measured from the rising edge of the CLK
- Flip-flop one bit of state that samples every rising edge of the CLK (positive edge-triggered)
- Register several bits of state that samples on rising edge of CLK or on LOAD (positive edge-triggered)

#### Maximum Clock Frequency

• What is the maximum frequency of this circuit?



Max Delay = Setup Time + CLK-to-Q Delay + CL Delay

#### FSM Example: 3 ones... FSM to detect the occurrence of 3 consecutive 1's in the input. INPUT a 1 1 6 1 1 1 0 1 1 1 0 1 1 1 1 0 OUTPUT Input/output 1/1 Draw the FSM... 0/0 10 1/0 $\zeta \sigma$

Assume state transitions are controlled by the clock: on each clock cycle the machine checks the inputs and moves to a new state and produces a new output...

## Datapath: Five Stages of Instruction Execution

- Stage 1: Instruction Fetch
- Stage 2: Instruction Decode
- Stage 3: ALU (Arithmetic-Logic Unit)
- Stage 4: Memory Access
- Stage 5: Register Write

#### Stages of Execution on Datapath



#### **Datapath Control Signals**

• MemWr: 1 => write memory "zero", "sign" ExtOp: • MemtoReg: 0 => ALU; 1 => Mem ALUsrc:  $0 \Rightarrow regB;$ • RegDst: 0 => "rt"; 1 => "rd" 1 => immed"ADD", "SUB", "OR" RegWr: 1 => write register ALUctr: nPC\_sel: 1 => branchALUctr **MemtoReg** Rd Rt RegDst MemWr 0 Inst Address Rs Rt RegWr nPC sel & Equal 5 5/ Adder 32 busA Ra Rb Rw busW 32 00 ALU RegFile busB 32 Mux PC 32 clk | Adder Extender 32 WrEn Adr imm16 Data In PC clk Data 32 16 Ext Memory clk **ALUSrc** ExtOp

R	ΓL:	The	Add	d Ins <sup>.</sup>	truct	ion
31	26	21	16	11	6	0
op	)	rs	rt	rd	shamt	funct
6 k	oits	5 bits	5 bits	5 bits	5 bits	6 bits

#### add rd, rs, rt

- MEM[PC] Fetch the instruction from memory
- R[rd] = R[rs] + R[rt] The actual operation
- PC = PC + 4 Calculate the next instruction's address

Instruction Fetch Unit at the Beginning of Add

 Fetch the instruction from Instruction memory: Instruction = MEM[PC]



# Single Cycle Datapath during Add312621161160oprsrtrdshamtfunct

R[rd] = R[rs] + R[rt]



#### Summary of the Control Signals (1/2)

- <u>inst</u> <u>Register Transfer</u>
- add R[rd] ← R[rs] + R[rt]; PC ← PC + 4 ALUSrc=RegB, ALUCtr="ADD", RegDst=rd, RegWr, nPC\_sel="+4"
- sub R[rd] ← R[rs] R[rt]; PC ← PC + 4
  ALUsrc=RegB, ALUctr="SUB", RegDst=rd, RegWr, nPC\_sel="+4"
- ori R[rt] ← R[rs] + zero\_ext(Imm16); PC ← PC + 4
  ALUsrc=Im, Extop="Z", ALUctr="OR", RegDst=rt,RegWr, nPC\_sel="+4"
- lw R[rt] ← MEM[ R[rs] + sign\_ext(Imm16)]; PC ← PC + 4
  ALUsrc=Im, Extop="sn", ALUctr="ADD", MemtoReg, RegDst=rt, RegWr,
  nPC\_sel = "+4"
- sw MEM[ R[rs] + sign\_ext(Imm16)] ← R[rs]; PC ← PC + 4
  ALUsrc=Im, Extop="sn", ALUctr = "ADD", MemWr, nPC\_sel = "+4"
- beq if (R[rs] == R[rt]) then PC  $\leftarrow$  PC + sign\_ext(Imm16)] || 00 else PC  $\leftarrow$  PC + 4

nPC\_sel = "br", ALUctr = "SUB"

#### Summary of the Control Signals (2/2)

See ───────────────────────────────────		10 0000	10 0010	We Don't Care :-)							
Appendix	Ά		ор	00 0000	00 0000	00 110	1 10 0011	10 1011	00 0100	00 0010	
[				add	sub	ori	lw	sw	beq	jump	
[	Re	egDst		1	1	0	0	х	х	х	
[	Α	LUSrc		0	0	1	1	1	0	х	
[	Μ	lemtoReg		0	0	0	1	х	х	х	
	R	egWrite		1	1	1	1	0	0	0	
	Μ	lemWrite		0	0	0	0	1	0	0	
	nl	PCsel		0	0	0	0	0	1	х	
	Ju	Imp		0	0	0	0	0	0	1	
	E>	кtОр		х	х	0	1	1	х	х	
[	Α	LUctr<2:0>	>	Add	Subtract	Or	Add	Add	Subtract	х	
31 26		2	1	16	11	6		0			
R-typ	e	ор		rs	rt		rd	shamt	fund	ct ado	l, sub
I-type op		rs	rt	immediate			ori,	lw, sw, beq			
J-typ	J-type op target address						jum	າp			

### Pipelining: Single Cycle Datapath



## **Pipeline registers**



Need registers between stages

To hold information produced in previous cycle

### **Graphical Pipeline Representation**

 RegFile: left half is write, right half is read Time (clock cycles)



# Pipelining Performance (3/3)


## **Pipelining Hazards**

A *hazard* is a situation that prevents starting the next instruction in the next clock cycle

#### 1) Structural hazard

 A required resource is busy (e.g. needed in multiple stages)

#### 2) Data hazard

- Data dependency between instructions
- Need to wait for previous instruction to complete its data read/write

#### 3) Control hazard

- Flow of execution depends on previous instruction

#### Structural Hazard #1: Single Memory



#### Structural Hazard #2: Registers (1/2)



## 2. Data Hazards (2/2)

• Data-flow *backwards* in time are hazards

Time (clock cycles)



#### **Data Hazard Solution: Forwarding**

Forward result as soon as it is available
 OK that it's not stored in RegFile yet



#### Data Hazard: Loads (1/4)

• **Recall:** Dataflow backwards in time are hazards



- Can't solve all cases with forwarding
  - Must *stall* instruction dependent on load, then forward (more hardware)

## Data Hazard: Loads (2/4)



#### Data Hazard: Loads (3/4)

• Stalled instruction converted to "bubble", acts like nop



### 3. Control Hazards

- Branch determines flow of control
  - Fetching next instruction depends on branch outcome
  - Pipeline can't always fetch correct instruction
    - Still working on ID stage of branch
- BEQ, BNE in MIPS pipeline
- Simple solution Option 1: *Stall* on every branch until branch condition resolved
  - Would add 2 bubbles/clock cycles for every Branch! (~ 20% of instructions executed)

## Caches



#### Adding Cache to Computer



Total Cash Capacity =Associativity \* # of sets \* block\_sizeBytes = blocks/set \* sets \* Bytes/blockC = N \* S \* BTagIndexIndex

address\_size = tag\_size + index\_size + offset\_size = tag\_size + log2(S) + log2(B)

Clicker Question: C remains constant, S and/or B can change such that C = 2N \* (SB)' => (SB)' = SB/2 Tag\_size = address\_size - (log2(S) + log2(B)) = address\_size - log2(SB) = address\_size - (log2(SB) - 1)

- Principle of Locality for Libraries /Computer Memory
- Hierarchy of Memories (speed/size/cost per bit) to Exploit Locality
- Cache copy of data lower level in memory hierarchy
- Direct Mapped to find block in cache using Tag field and Valid bit for Hit
- Cache design choice:
  - Write-Through vs. Write-Back

#### **Cache Organizations**

- "Fully Associative": Block can go anywhere
  - First design in lecture
  - Note: No Index field, but 1 comparator/block
- "Direct Mapped": Block goes one place
  - Note: Only 1 comparator
  - Number of sets = number blocks
- "N-way Set Associative": N places for a block
  - Number of sets = number of blocks / N
  - N comparators
  - Fully Associative: N = number of blocks
  - Direct Mapped: N = 1

#### Processor Address Fields used by Cache Controller

- Block Offset: Byte address within block
- Set Index: Selects which set
- Tag: Remaining portion of processor address

Processor Address (32-bits total)

Tag	Set Index	Block offset
-----	-----------	--------------

- Size of Index = log2 (number of sets)
- Size of Tag = Address size Size of Index – log2 (number of bytes/block)

#### Write Policy Choices

- Cache hit:
  - write through: writes both cache & memory on every access
    - Generally higher memory traffic but simpler pipeline & cache design
  - write back: writes cache only, memory `written only when dirty entry evicted
    - A dirty bit per line reduces write-back traffic
    - Must handle 0, 1, or 2 accesses to memory for each load/store
- Cache miss:
  - **no write allocate**: only write to main memory
  - write allocate (aka fetch on write): fetch into cache
- Common combinations:
  - write through and no write allocate
  - write back with write allocate

#### **Direct-Mapped Cache Review**

• One word blocks, cache size = 1K words (or 4KB)



## Sources of Cache Misses (3 C's)

- *Compulsory* (cold start, first reference):
  - 1<sup>st</sup> access to a block, "cold" fact of life, not a lot you can do about it.
    - If running billions of instructions, compulsory misses are insignificant
- Capacity:
  - Cache cannot contain all blocks accessed by the program
    - Misses that would not occur with infinite cache
- *Conflict* (collision):
  - Multiple memory locations mapped to same cache set
    - Misses that would not occur with ideal fully associative cache

#### Impact of Cache Parameters on Performance

• AMAT = Hit Time + Miss Rate \* Miss Penalty

– Note, we assume always first search cache, so must charge hit time for both hits and misses!

• For misses, characterize by 3Cs

#### Local vs. Global Miss Rates

- Local miss rate the fraction of references to one level of a cache that miss
- Local Miss rate L2\$ = \$L2 Misses / L1\$ Misses
- Global miss rate the fraction of references that miss in all levels of a multilevel cache
  - L2\$ local miss rate >> than the global miss rate
- Global Miss rate = L2\$ Misses / Total Accesses
   = (L2\$ Misses / L1\$ Misses) × (L1\$ Misses / Total Accesses)
   = Local Miss rate L2\$ × Local Miss rate L1\$
- AMAT = Time for a hit + Miss rate × Miss penalty
- AMAT = Time for a L1\$ hit + (local) Miss rate L1\$ × (Time for a L2\$ hit + (local) Miss rate L2\$ × L2\$ Miss penalty)

#### In Conclusion, Cache Design Space

- Several interacting dimensions
  - Cache size
  - Block size
  - Associativity
  - Replacement policy
  - Write-through vs. write-back
  - Write-allocation
- Optimal choice is a compromise
  - Depends on access characteristics
    - Workload
    - Use (I-cache, D-cache)
  - Depends on technology / cost
- Simplicity often wins



#### Iron Law of Performance

- A program executes instructions
- CPU Time/Program
  - = Clock Cycles/Program x Clock Cycle Time
  - = Instructions/Program
     x Average Clock Cycles/Instruction
     x Clock Cycle Time
- 1<sup>st</sup> term called *Instruction Count*
- 2<sup>nd</sup> term abbreviated CPI for average
   Clock Cycles Per Instruction
- 3rd term is 1 / Clock rate

#### IEEE 754 Floating-Point Standard (1/3)

Single Precision (Double Precision similar):



- Sign bit: 1 means negative 0 means positive
- Significand in *sign-magnitude* format (not 2's complement)
  - To pack more bits, leading 1 implicit for normalized numbers
  - 1 + 23 bits single, 1 + 52 bits double
  - always true: 0 < Significand < 1 (for</li>
- (for normalized numbers)
- Note: 0 has no leading 1, so reserve exponent value 0 just for number 0

#### IEEE 754 Floating Point Standard (2/3)

- IEEE 754 uses "biased exponent" representation
  - Designers wanted FP numbers to be used even if no FP hardware; e.g., sort records with FP numbers using integer compares
  - Wanted bigger (integer) exponent field to represent bigger numbers
  - 2's complement poses a problem (because negative numbers look bigger)
    - Use just magnitude and offset by half the range

#### IEEE 754 Floating Point Standard (3/3)

- Called <u>Biased Notation</u>, where bias is number subtracted to get final number
  - IEEE 754 uses bias of 127 for single prec.
  - Subtract 127 from Exponent field to get actual value for exponent

<ul> <li>Summary (single precision):</li> </ul>			
<u>31 30 23</u>	22 0		
S Exponent	Significand		
1 bit 8 bits	23 bits		
•(-1) <sup>S</sup> x (1 + Significand) x 2 <sup>(Exponent-127)</sup>			

 Double precision identical, except with exponent bias of 1023 (half, quad similar)

### Flynn\* Taxonomy, 1966

		Data Streams		
		Single	Multiple	
Instruction Streams	Single	SISD: Intel Pentium 4	SIMD: SSE instructions of x86	
	Multiple	MISD: No examples today	MIMD: Intel Xeon e5345 (Clovertown)	

- Since about 2013, SIMD and MIMD most common parallelism in architectures – usually both in same system!
- Most common parallel processing programming style: Single Program Multiple Data ("SPMD")
  - Single program that runs on all processors of a MIMD
  - Cross-processor execution coordination using synchronization primitives
- SIMD (aka hw-level *data parallelism*): specialized function units, for handling lock-step calculations involving arrays
  - Scientific computing, signal processing, multimedia (audio/video processing)

\*Prof. Michael Flynn, Stanford

#### Big Idea: Amdahl's Law



Example: the execution time of half of the program can be accelerated by a factor of 2. What is the program speed-up overall?

$$\frac{1}{\frac{0.5+0.5}{2}} = \frac{1}{\frac{0.5+0.25}{2}} = 1.33$$

### Strong and Weak Scaling

- To get good speedup on a parallel processor while keeping the problem size fixed is harder than getting good speedup by increasing the size of the problem.
  - Strong scaling: when speedup can be achieved on a parallel processor without increasing the size of the problem
  - Weak scaling: when speedup is achieved on a parallel processor by increasing the size of the problem proportionally to the increase in the number of processors
- Load balancing is another important factor: every processor doing same amount of work
  - Just one unit with twice the load of others cuts speedup almost in half

#### Data Level Parallelism

- Loop Unrolling
- Intel SSE SIMD Instructions
  - Exploit data-level parallelism in loops
  - One instruction fetch that operates on multiple operands simultaneously
  - 128-bit XMM registers
- SSE Instructions in C using Intrinsics
  - Embed the SSE machine instructions directly into C programs through use of intrinsics
  - Achieve efficiency beyond that of optimizing compiler

#### CS 110 Computer Architecture

#### Thread-Level Parallelism (TLP) and OpenMP Intro

Instructor: Sören Schwertfeger

http://shtech.org/courses/ca/

School of Information Science and Technology SIST

ShanghaiTech University

Slides based on UC Berkley's CS61C

#### Review

- Sequential software is slow software
   SIMD and MIMD are paths to higher performance
- MIMD thru: multithreading processor cores (increases utilization), Multicore processors (more cores per chip)
- Synchronization coordination among threads
  - MIPS: atomic read-modify-write using loadlinked/store-conditional
- OpenMP as simple parallel extension to C
  - Pragmas for forking multiple Threads
  - ≈ C: small so easy to learn, but not very high level and it's easy to get into trouble

#### **OpenMP Programming Model - Review**

• Fork - Join Model:



- OpenMP programs begin as single process (*master thread*) and executes sequentially until the first parallel region construct is encountered
  - FORK: Master thread then creates a team of parallel threads
  - Statements in program that are enclosed by the parallel region construct are executed in parallel among the various threads
  - JOIN: When the team threads complete the statements in the parallel region construct, they synchronize and terminate, leaving only the master thread

# parallel Pragma and Scope - Review

• Basic OpenMP construct for parallelization:

```
#pragma omp parallel
{
    /* code goes here */
}
```

- *Each* thread runs a copy of code within the block
- Thread scheduling is *non-deterministic*
- OpenMP default is *shared* variables

   To make private, need to declare with pragma:
   #pragma omp parallel private (x)

### **OpenMP Directives (Work-Sharing)**

These are defined within a parallel section



#### Parallel Statement Shorthand



#pragma omp parallel for
 for(i=0; i<len; i++) { ... }</pre>

• Also works for sections
## Building Block: for loop

for (i=0; i<max; i++) zero[i] = 0;</pre>

- Breaks for loop into chunks, and allocate each to a separate thread
  - e.g. if max = 100 with 2 threads: assign 0-49 to thread 0, and 50-99 to thread 1
- Must have relatively simple "shape" for an OpenMPaware compiler to be able to parallelize it
  - Necessary for the run-time system to be able to determine how many of the loop iterations to assign to each thread
- No premature exits from the loop allowed <
   <ul>
   In general, don't jump
   i.e. No break, return, exit, goto statements
   outside of any pragma block

#### Parallel for pragma

#### #pragma omp parallel for for (i=0; i<max; i++) zero[i] = 0;</pre>

- Master thread creates additional threads, each with a separate execution context
- All variables declared outside for loop are shared by default, except for loop index which is *private* per thread (Why?)
- Implicit "barrier" synchronization at end of for loop
- Divide index regions sequentially per thread
  - Thread 0 gets 0, 1, ..., (max/n)-1;
  - Thread 1 gets max/n, max/n+1, ..., 2\*(max/n)-1
  - Why?

master

master

FORK

DO / for loop

.101

# **OpenMP** Timing

- Elapsed wall clock time:
  - double omp\_get\_wtime(void);
  - Returns elapsed wall clock time in seconds
  - Time is measured per thread, no guarantee can be made that two distinct threads measure the same time
  - Time is measured from "some time in the past," so subtract results of two calls to omp\_get\_wtime to get elapsed time

#### Matrix Multiply in OpenMP



### Notes on Matrix Multiply Example

- More performance optimizations available:
  - Higher compiler optimization (-O2, -O3) to reduce number of instructions executed
  - Cache blocking to improve memory performance
  - Using SIMD SSE instructions to raise floating point computation rate (*DLP*)

#### Simple Multi-core Processor



### **Multiprocessor Caches**

- Memory is a performance bottleneck even with one processor
- Use caches to reduce bandwidth demands on main memory
- Each core has a local private cache holding data it has accessed recently
- Only cache misses have to access the shared common memory



#### **Shared Memory and Caches**

- What if?
  - Processors 1 and 2 read Memory[1000] (value 20)



## **Shared Memory and Caches**

- Now:
  - Processor 0 writes Memory[1000] with 40



# **Keeping Multiple Caches Coherent**

- Architect's job: shared memory
   => keep cache values coherent
- Idea: When any processor has cache miss or writes, notify other processors via interconnection network
  - If only reading, many processors can have copies
  - If a processor writes, invalidate any other copies
- Write transactions from one processor, other caches "snoop" the common interconnect checking for tags they hold
  - Invalidate any copies of same address modified in other cache

#### **Shared Memory and Caches**

- Example, now with cache coherence
  - Processors 1 and 2 read Memory[1000]
  - Processor 0 writes Memory[1000] with 40



#### Question: Which statement(s) are true?

- A: Using write-through caches removes the need for cache coherence
- B: Every processor store instruction must check contents of other caches
- C: Most processor load and store accesses only need to check in local private cache
- D: Only one processor can cache any memory location at one time

## Cache Coherency Tracked by Block



- Suppose block size is 32 bytes
- Suppose Processor 0 reading and writing variable X, Processor 1 reading and writing variable Y
- Suppose in X location 4000, Y in 4012
- What will happen?

# **Coherency Tracked by Cache Block**

- Block ping-pongs between two caches even though processors are accessing disjoint variables
- Effect called *false sharing*
- How can you prevent it?

### Shared Memory and Caches

- Use valid bit to "unload" cache lines (in Processors 1 and 2)
- Dirty bit tells me: "I am the only one using this cache line"! => no need to announce on Network!



#### Review: Understanding Cache Misses: The 3Cs

- Compulsory (cold start or process migration, 1<sup>st</sup> reference):
  - First access to block, impossible to avoid; small effect for long-running programs
  - Solution: increase block size (increases miss penalty; very large blocks could increase miss rate)
- Capacity (not compulsory and...)
  - Cache cannot contain all blocks accessed by the program *even with perfect replacement policy in fully associative cache*
  - Solution: increase cache size (may increase access time)
- **Conflict** (not compulsory or capacity and...):
  - Multiple memory locations map to the same cache location
  - Solution 1: increase cache size
  - Solution 2: increase associativity (may increase access time)
  - Solution 3: improve replacement policy, e.g.. LRU

#### Fourth "C" of Cache Misses: Coherence Misses

- Misses caused by coherence traffic with other processor
- Also known as *communication* misses because represents data moving between processors working together on a parallel program
- For some parallel programs, coherence misses can dominate total misses

#### Example: Calculating $\pi$

#### Numerical Integration



Mathematically, we know that:

$$\int_{0}^{1} \frac{4.0}{(1+x^2)} \, dx = \pi$$

We can approximate the integral as a sum of rectangles:

 $\sum_{i=0}^{N} F(\mathbf{x}_{i}) \Delta \mathbf{x} \approx \pi$ 

Where each rectangle has width  $\Delta x$  and height F(x<sub>i</sub>) at the middle of interval i.

#### Sequential Calculation of $\pi$ in C

```
#include <stdio.h> /* Serial Code */
static long num steps = 100000;
double step;
void main () {
    int i;
    double x, pi, sum = 0.0;
    step = 1.0/(double)num steps;
    for (i = 1; i <= num steps; i++) {</pre>
      x = (i - 0.5) * step;
      sum = sum + 4.0 / (1.0 + x*x);
                                              Ν
     }
                                                F(x_i)\Delta x \approx \pi
    pi = sum * step;
    printf ("pi = %6.12f\n", pi);
```

#### Parallel OpenMP Version (1)

```
#include <omp.h>
#define NUM THREADS 4
static long num steps = 100000; double step;
void main () {
  int i; double x, pi, sum[NUM THREADS];
  step = 1.0/(double) num steps;
  #pragma omp parallel private ( i, x )
    int id = omp get thread num();
    for (i=id, sum[id]=0.0; i< num steps; i=i+NUM THREADS)</pre>
    ł
      x = (i+0.5) * step;
      sum[id] += 4.0/(1.0+x*x);
    }
  for(i=1; i<NUM THREADS; i++) sum[0] += sum[i];</pre>
  pi = sum[0] * step;
 printf ("pi = %6.12f\n", pi);
```

#### **OpenMP Reduction**

```
double avg, sum=0.0, A[MAX]; int i;
#pragma omp parallel for private ( sum )
for (i = 0; i <= MAX ; i++)
    sum += A[i];
avg = sum/MAX; // bug</pre>
```

- Problem is that we really want sum over all threads!
- Reduction: specifies that 1 or more variables that are private to each thread are subject of reduction operation at end of parallel region:

#### reduction(operation:var) where

- Operation: operator to perform on the variables (var) at the end of the parallel region: +, \*, -, &, ^, |, &&, or ||.
- *Var*: One or more variables on which to perform scalar reduction.

```
double avg, sum=0.0, A[MAX]; int i;
#pragma omp for reduction(+ : sum)
for (i = 0; i <= MAX ; i++)
    sum += A[i];
avg = sum/MAX;</pre>
```

#### Version 2: parallel for, reduction

```
#include <omp.h>
#include <stdio.h>
/static long num steps = 100000;
double step;
void main () {
    int i; double x, pi, sum = 0.0;
    step = 1.0 / (double)num steps;
#pragma omp parallel for private(x) reduction(+:sum)
    for (i=1; i<= num steps; i++) {</pre>
       x = (i - 0.5) * step;
       sum = sum + 4.0 / (1.0+x*x);
    }
    pi = sum * step;
    printf ("pi = %6.8f\n", pi);
```

}

# And in Conclusion, ...

- Multiprocessor/Multicore uses Shared Memory
  - Cache coherency implements shared memory even with multiple copies in multiple caches
  - False sharing a concern; watch block size!
- OpenMP as simple parallel extension to C
   Threads, Parallel for, private, reductions ...
  - ≈ C: small so easy to learn, but not very high level and it's easy to get into trouble
  - Much we didn't cover including other synchronization mechanisms (locks, etc.)