

EECS 252 Graduate Computer Architecture

Lec 11 – Mid Term Review

David Culler
Electrical Engineering and Computer Sciences
University of California, Berkeley

<http://www.eecs.berkeley.edu/~culler>
<http://www-inst.eecs.berkeley.edu/~cs252>

Review Exercise

- The 1+X “accumulator” based ISA never seems to go away because of its “minimal” processor state – witness the longevity of the 8051
- You are given the task of designing a “high performance 8051”. Having learned about the separation of architected state and microarchitecture, you are ready to attack the problem. A simple analysis suggests that 8051 code has very strong sequential dependences. You will need to use serious instruction lookahead, branch prediction, and register renaming to get at the ILP.
- Assume a MIPS 10K-like data path with multiple function units, lots of physical registers. You need to design the instruction issue and register mapping logic to get ILP out of this beast.
- When is a physical register available for reuse?

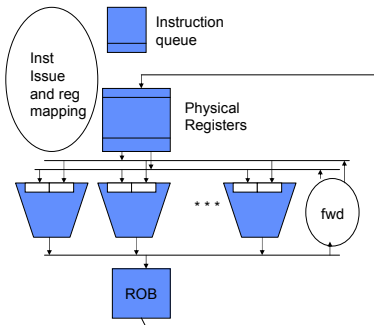
2/22/2005

CS252 L11-review

2

Solution Framework

- ISA?
- Typical sequence
- Dependences
- Names?
- Mapping
- Free



2/22/2005

CS252 L11-review

3

Review of Memory Hierarchy that we skipped

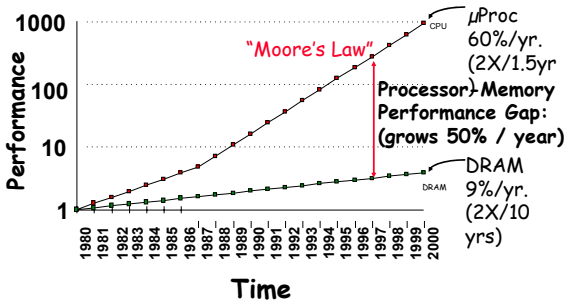
2/22/2005

CS252 L11-review

4

Recap: Who Cares About the Memory Hierarchy?

Processor-DRAM Memory Gap (latency)

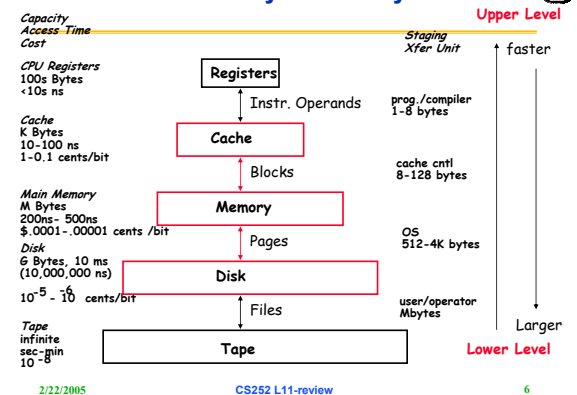


2/22/2005

CS252 L11-review

5

Levels of the Memory Hierarchy



2/22/2005

CS252 L11-review

6

The Principle of Locality

- **The Principle of Locality:**
 - Program access a relatively small portion of the address space at any instant of time.
- **Two Different Types of Locality:**
 - **Temporal Locality** (Locality in Time): If an item is referenced, it will tend to be referenced again soon (e.g., loops, reuse)
 - **Spatial Locality** (Locality in Space): If an item is referenced, items whose addresses are close by tend to be referenced soon (e.g., straightline code, array access)
- Last 15 years, HW relied on locality for speed

It is a property of programs which is exploited in machine design.

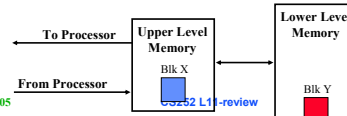
2/22/2005

CS252 L11-review

7

Memory Hierarchy: Terminology

- **Hit:** data appears in some block in the upper level (example: Block X)
 - **Hit Rate:** the fraction of memory access found in the upper level
 - **Hit Time:** Time to access the upper level which consists of RAM access time + Time to determine hit/miss
- **Miss:** data needs to be retrieve from a block in the lower level (Block Y)
 - **Miss Rate** = 1 - (Hit Rate)
 - **Miss Penalty:** Time to replace a block in the upper level + Time to deliver the block the processor
- Hit Time << Miss Penalty (500 instructions on 21264!)



2/22/2005

CS252 L11-review

8

Cache Measures

- **Hit rate:** fraction found in that level
 - So high that usually talk about **Miss rate**
 - Miss rate fallacy: as MIPS to CPU performance, miss rate to average memory access time in memory
- **Average memory-access time**

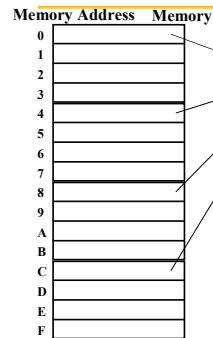
$$= \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$
 (ns or clocks)
- **Miss penalty:** time to replace a block from lower level, including time to replace in CPU
 - **access time:** time to lower level = f(latency to lower level)
 - **transfer time:** time to transfer block = f(BW between upper & lower levels)

2/22/2005

CS252 L11-review

9

Simplest Cache: Direct Mapped



- **Location 0 can be occupied by data from:**
 - Memory location 0, 4, 8, ... etc.
 - In general: any memory location whose 2 LSBs of the address are 0s
 - Address <1:0> => cache index
- Which one should we place in the cache?
- How can we tell which one is in the cache?

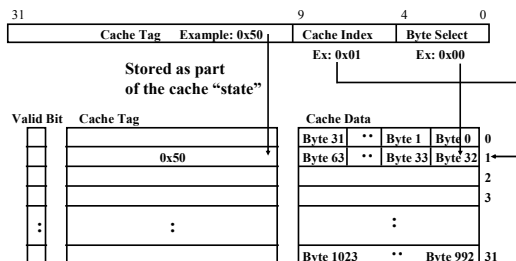
2/22/2005

CS252 L11-review

10

1 KB Direct Mapped Cache, 32B blocks

- For a $2^{**} N$ byte cache:
 - The uppermost (32 - N) bits are always the Cache Tag
 - The lowest M bits are the Byte Select (Block Size = $2^{**} M$)



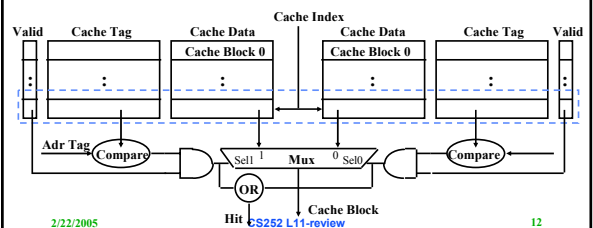
2/22/2005

CS252 L11-review

11

Two-way Set Associative Cache

- **N-way set associative:** N entries for each Cache Index
 - N direct mapped caches operates in parallel (N typically 2 to 4)
- **Example: Two-way set associative cache**
 - Cache Index selects a "set" from the cache
 - The two tags in the set are compared in parallel
 - Data is selected based on the tag result



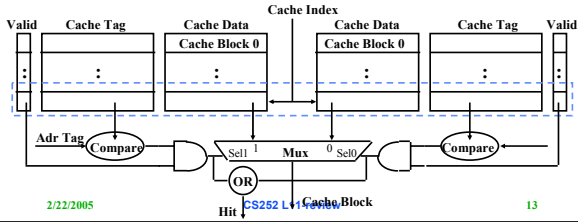
2/22/2005

CS252 L11-review

12

Disadvantage of Set Associative Cache

- **N-way Set Associative Cache v. Direct Mapped Cache:**
 - N comparators vs. 1
 - Extra MUX delay for the data
 - Data comes AFTER Hit/Miss
- In a direct mapped cache, Cache Block is available BEFORE Hit/Miss:
 - Possible to assume a hit and continue. Recover later if miss.



2/22/2005

13

4 Questions for Memory Hierarchy

- Q1: Where can a block be placed in the upper level?
(Block placement)
- Q2: How is a block found if it is in the upper level?
(Block identification)
- Q3: Which block should be replaced on a miss?
(Block replacement)
- Q4: What happens on a write?
(Write strategy)

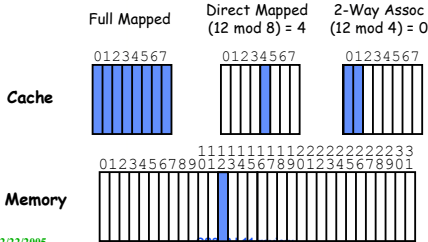
2/22/2005

CS252 L11-review

14

Q1: Where can a block be placed in the upper level?

- Block 12 placed in 8 block cache:
 - Fully associative, direct mapped, 2-way set associative
 - S.A. Mapping = Block Number Modulo Number Sets

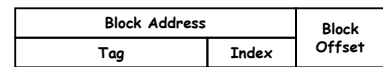


2/22/2005

15

Q2: How is a block found if it is in the upper level?

- Tag on each block
 - No need to check index or block offset
- Increasing associativity shrinks index, expands tag



2/22/2005

CS252 L11-review

16

Q3: Which block should be replaced on a miss?

- Easy for Direct Mapped
- Set Associative or Fully Associative:
 - Random
 - LRU (Least Recently Used)

| Assoc: | 2-way | | 4-way | | 8-way | |
|--------|-------|-------|-------|-------|-------|-------|
| Size | LRU | Ran | LRU | Ran | LRU | Ran |
| 16 KB | 5.2% | 5.7% | 4.7% | 5.3% | 4.4% | 5.0% |
| 64 KB | 1.9% | 2.0% | 1.5% | 1.7% | 1.4% | 1.5% |
| 256 KB | 1.15% | 1.17% | 1.13% | 1.13% | 1.12% | 1.12% |

2/22/2005

CS252 L11-review

17

Q4: What happens on a write?

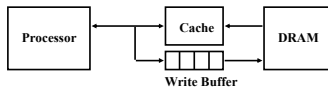
- **Write through**—The information is written to both the block in the cache and to the block in the lower-level memory.
- **Write back**—The information is written only to the block in the cache. The modified cache block is written to main memory only when it is replaced.
 - is block clean or dirty?
- Pros and Cons of each?
 - WT: read misses cannot result in writes
 - WB: no repeated writes to same location
- WT always combined with write buffers so that don't wait for lower level memory

2/22/2005

CS252 L11-review

18

Write Buffer for Write Through



- **A Write Buffer is needed between the Cache and Memory**
 - Processor: writes data into the cache and the write buffer
 - Memory controller: write contents of the buffer to memory
- **Write buffer is just a FIFO:**
 - Typical number of entries: 4
 - Works fine if: Store frequency (w.r.t. time) $\ll 1 / \text{DRAM write cycle}$
- **Memory system designer's nightmare:**
 - Store frequency (w.r.t. time) $\rightarrow 1 / \text{DRAM write cycle}$
 - Write buffer saturation

2/22/2005

CS252 L11-review

19

Impact of Memory Hierarchy on Algorithms

- Today CPU time is a function of (ops, cache misses) vs. just $f(\text{ops})$: What does this mean to Compilers, Data structures, Algorithms?
- "The Influence of Caches on the Performance of Sorting" by A. LaMarca and R.E. Ladner. *Proceedings of the Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, January, 1997, 370-379.
- Quicksort: fastest comparison based sorting algorithm when all keys fit in memory
- Radix sort: also called "linear time" sort because for keys of fixed length and fixed radix a constant number of passes over the data is sufficient independent of the number of keys
- For AlphaStation 250, 32 byte blocks, direct mapped L2 2MB cache, 8 byte keys, from 4000 to 4000000

2/22/2005

CS252 L11-review

20

Key topics firehose...

2/22/2005

CS252 L11-review

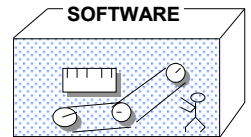
21

Instruction Set Architecture

... the attributes of a [computing] system as seen by the programmer, i.e. the conceptual structure and functional behavior, as distinct from the organization of the data flows and controls the logic design, and the physical implementation.

– Amdahl, Blaaw, and Brooks, 1964

- Organization of Programmable Storage
- Data Types & Data Structures: Encodings & Representations
- Instruction Formats
- Instruction (or Operation Code) Set
- Modes of Addressing and Accessing Data Items and Instructions
- Exceptional Conditions

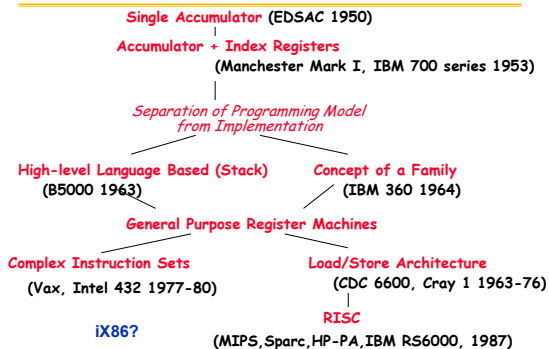


2/22/2005

CS252 L11-review

22

Evolution of Instruction Sets



2/22/2005

CS252 L11-review

23

Components of Performance

$$\text{CPU time} = \frac{\text{Seconds}}{\text{Program}} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Cycle}}$$

inst count
CPI
Cycle time

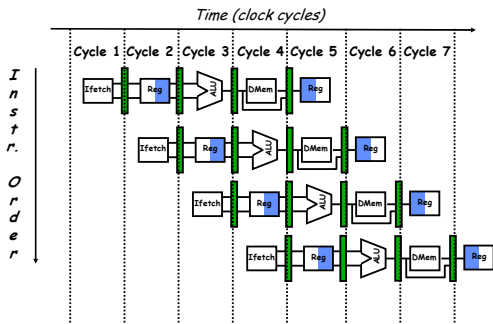
| | Inst Count | CPI | Clock Rate |
|--------------|------------|-----|------------|
| Program | X | | |
| Compiler | X | (X) | |
| Inst. Set. | X | X | |
| Organization | | X | X |
| Technology | | | X |

2/22/2005

CS252 L11-review

24

Pipelined Instruction Execution



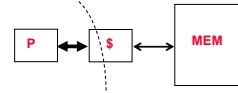
2/22/2005

CS252 L11-review

25

The Principle of Locality

- **The Principle of Locality:**
 - Program access a relatively small portion of the address space at any instant of time.
- **Two Different Types of Locality:**
 - **Temporal Locality** (Locality in Time): If an item is referenced, it will tend to be referenced again soon (e.g., loops, reuse)
 - **Spatial Locality** (Locality in Space): If an item is referenced, items whose addresses are close by tend to be referenced soon (e.g., straightline code, array access)
- Last 30 years, HW relied on locality for speed

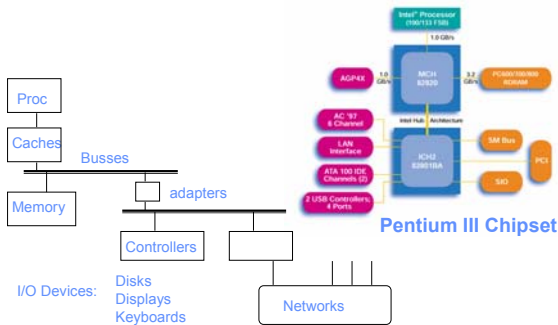


2/22/2005

CS252 L11-review

26

System Organization: It's all about communication



2/22/2005

CS252 L11-review

27

Amdahl's Law

$$\text{ExTime}_{\text{new}} = \text{ExTime}_{\text{old}} \times \left[(1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}} \right]$$

$$\text{Speedup}_{\text{overall}} = \frac{\text{ExTime}_{\text{old}}}{\text{ExTime}_{\text{new}}} = \frac{1}{(1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}}}$$

Best you could ever hope to do:

$$\text{Speedup}_{\text{maximum}} = \frac{1}{(1 - \text{Fraction}_{\text{enhanced}})}$$



2/22/2005

CS252 L11-review

28

Cycles Per Instruction (Throughput)

"Average Cycles per Instruction"

$$\text{CPI} = \frac{\text{CPU Time} \times \text{Clock Rate}}{\text{Instruction Count}} = \frac{\text{Cycles}}{\text{Instruction Count}}$$

$$\text{CPU time} = \text{Cycle Time} \times \sum_{j=1}^n \text{CPI}_j \times I_j$$

$$\text{CPI} = \sum_{j=1}^n \text{CPI}_j \times F_j \quad \text{where } F_j = \frac{I_j}{\text{Instruction Count}}$$

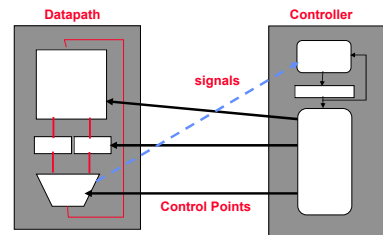
"Instruction Frequency"

2/22/2005

CS252 L11-review

29

Datapath vs Control



- **Datapath:** Storage, FU, interconnect sufficient to perform the desired functions
 - Inputs are Control Points
 - Outputs are signals
- **Controller:** State machine to orchestrate operation on the datapath
 - Based on desired function and signals

2/22/2005

CS252 L11-review

30

Pipelining is not quite that easy!

- Limits to pipelining: **Hazards** prevent next instruction from executing during its designated clock cycle
 - Structural hazards:** HW cannot support this combination of instructions (single person to fold and put clothes away)
 - Data hazards:** Instruction depends on result of prior instruction still in the pipeline (missing sock)
 - Control hazards:** Caused by delay between the fetching of instructions and decisions about changes in control flow (branches and jumps).

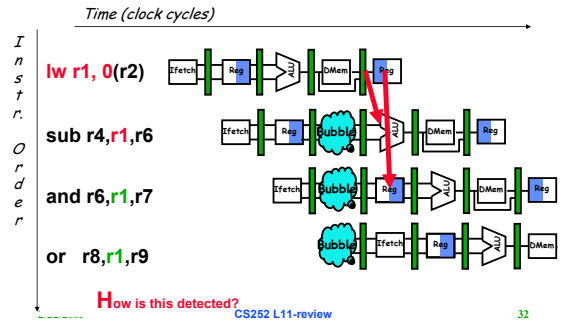
2/22/2005

CS252 L11-review

31

Data Hazard Even with Forwarding

Figure 3.13, Page 154



CS252 L11-review

32

Speed Up Equation for Pipelining

$$CPI_{\text{pipelined}} = \text{Ideal CPI} + \text{Average Stall cycles per Inst}$$

$$\text{Speedup} = \frac{\text{Ideal CPI} \times \text{Pipeline depth}}{\text{Ideal CPI} + \text{Pipeline stall CPI}} \times \frac{\text{Cycle Time}_{\text{unpipelined}}}{\text{Cycle Time}_{\text{pipelined}}}$$

For simple RISC pipeline, $CPI = 1$:

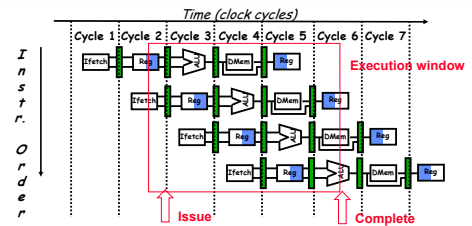
$$\text{Speedup} = \frac{\text{Pipeline depth}}{1 + \text{Pipeline stall CPI}} \times \frac{\text{Cycle Time}_{\text{unpipelined}}}{\text{Cycle Time}_{\text{pipelined}}}$$

2/22/2005

CS252 L11-review

33

Ordering Properties of basic inst. pipeline



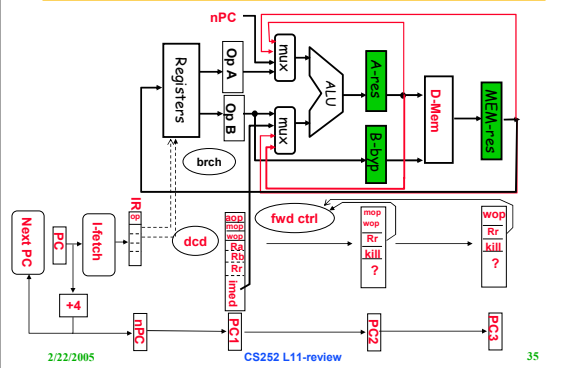
- Instructions issued in order
- Operand fetch is stage 2 => operand fetched in order
- Write back in stage 5 => no WAW, no WAR hazards
- Common pipeline flow => operands complete in order
- Stage changes only at "end of instruction"

2/22/2005

CS252 L11-review

34

Control Pipeline



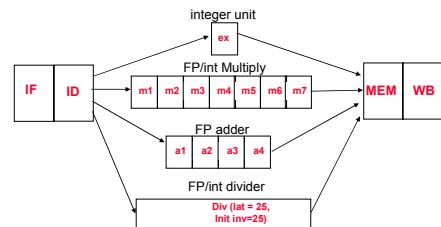
2/22/2005

CS252 L11-review

35

Typical "simple" Pipeline

- Example: MIPS R4000



2/22/2005

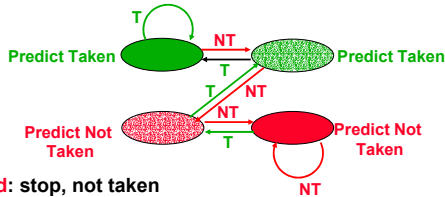
CS252 L11-review

36

2-bit Dynamic Branch Prediction

(J. Smith, 1981)

- 2-bit scheme where change prediction only if get misprediction **twice**:



- Red:** stop, not taken
- Green:** go, taken
- Adds **hysteresis** to decision making process
- Generalize to n-bit saturating counter

2/22/2005

CS252 L11-review

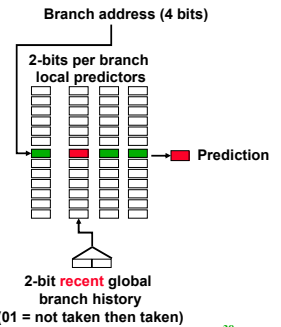
37

Correlating Branches

Idea: taken/not taken of recently executed branches is related to behavior of next branch (as well as the history of that branch behavior)

- Then behavior of recent branches selects between, say, 4 predictions of next branch, updating just that prediction

- (2,2) predictor: 2-bit global, 2-bit local



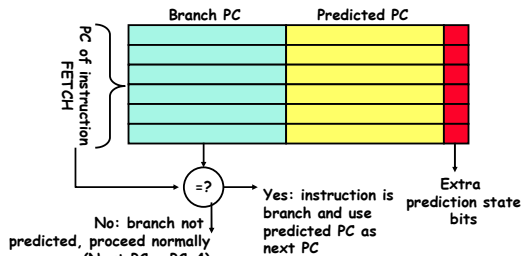
2/22/2005

CS252 L11-review

38

Need Address at Same Time as Prediction

- Branch Target Buffer (BTB): Address of branch index to get prediction AND branch address (if taken)
 - Note: must check for branch match now, since can't use wrong branch address (Figure 3.19, 3.20)



2/22/2005

CS252 L11-review

39

Pipelining with Reg. Reservations

- Assumptions

- Multiple pipelined function units of different latency
 - able to accept operations at issue rate
 - may be exceptions (e.g., divide)
- Issue instructions in order
- Operand fetch in order
- Completion out of order
 - short ops may bypass long ones
- Some shared resources (e.g., reg write port)

- Implications

- WAR hazard still resolved by pipeline flow (2 & 3)
- RAW, WAW, and structural still present

- Design philosophy (ala Cray)

- Resolve hazards as instruction is issued into pipeline
- Pipeline is non-blocking

2/22/2005

CS252 L11-review

40

Hazard Resolution

- Structural

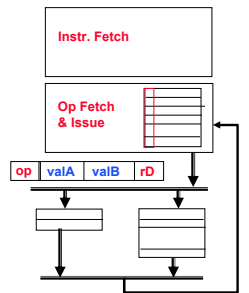
- Op code => resource usage
- Check resource resv
- Set on issue

- Data

- Add reservation bit one each register
- Check RegRsv for source and destination registers
- Hold issue till clear
- Set bit on destination register
- Clear bit on dest reg. Write

- Questions:

- Forwarding?



2/22/2005

CS252 L11-review

41

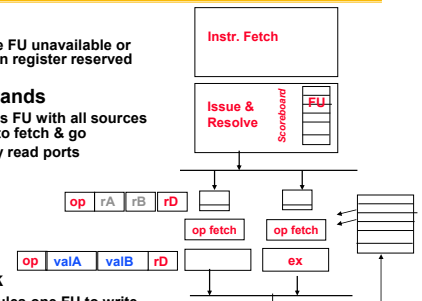
Scoreboard Operation

- Issue

- Hold while FU unavailable or destination register reserved (by FU f)

- Read operands

- SB informs FU with all sources available to fetch & go
- Limited by read ports



- Write back

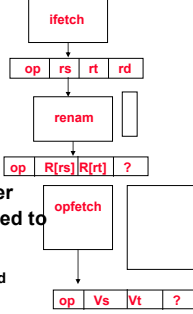
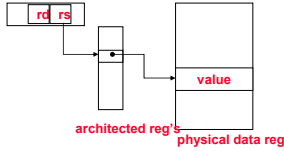
- SB schedules one FU to write
- Waits no FU waiting to fetch (old version) of reg

2/22/2005

CS252 L11-review

42

Register Renaming (less Conceptual)



- Separate the functions of the register
- Reg identifier in instruction is mapped to "physical register" id for current instance of the register
 - Physical reg set may be larger than allocated
- What are the rules for allocating / deallocating physical registers?

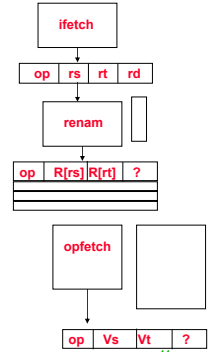
2/22/2005

CS252 L11-review

43

Reg renaming

- **Source Reg s:**
 - physical reg $P=R[s]$
- **Destination reg d:**
 - Old physical register $R[d]$ "terminates"
 - $R[d] := \text{get_free}$
- **Free physical register when**
 - No longer referenced by any architected register (terminated)
 - No incomplete instructions waiting to read it
 - » Easy with in-order
 - » Out of order?

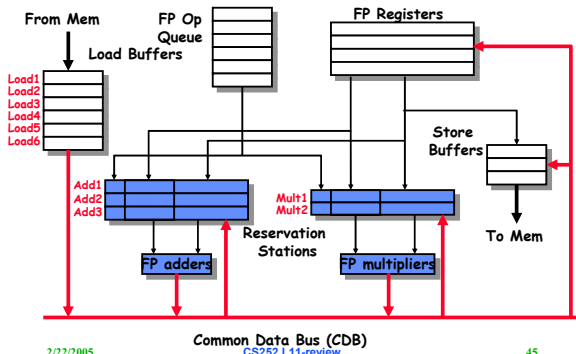


2/22/2005

CS252 L11-review

44

Tomasulo Organization



2/22/2005

CS252 L11-review

45

Three Stages of Tomasulo Algorithm

- 1. Issue**—get instruction from FP Op Queue
 - If reservation station free (no structural hazard), control issues instr & sends operands (renames registers).
 - 2. Execution**—operate on operands (EX)
 - When both operands ready then execute;
 - if not ready, watch Common Data Bus for result
 - 3. Write result**—finish execution (WB)
 - Write on Common Data Bus to all awaiting units;
 - mark reservation station available
- Normal data bus: data + destination ("go to" bus)
 - **Common data bus:** data + **source** ("come from" bus)
 - 64 bits of data + 4 bits of Functional Unit **source** address
 - Write if matches expected Functional Unit (produces result)
 - Does the broadcast

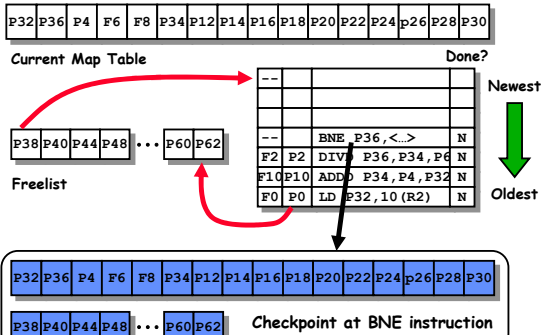
2/22/2005

CS252 L11-review

46

Explicit register renaming:

R10000 Freelist Management

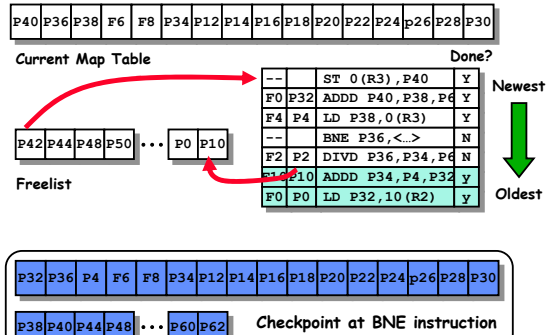


CS252 L11-review

47

Explicit register renaming:

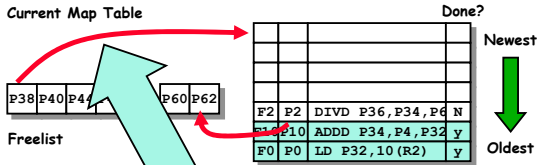
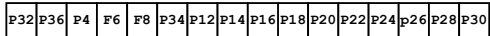
R10000 Freelist Management



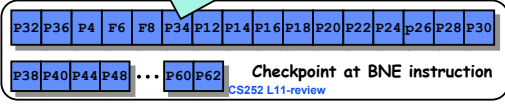
CS252 L11-review

48

Explicit register renaming: R10000 Freelist Management



Speculation error fixed by restoring map table and freelist



CS252 L11-review

49

Problems with scalar approach to ILP extraction

Limits to conventional exploitation of ILP:

- **pipelined clock rate:** at some point, each increase in clock rate has corresponding CPI increase (branches, other hazards)
- **branch prediction:** branches get in the way of wide issue. They are too unpredictable.
- **instruction fetch and decode:** at some point, its hard to fetch and decode more instructions per clock cycle
- **register renaming:** Rename logic gets really complicate for many instructions
- **cache hit rate:** some long-running (scientific) programs have very large data sets accessed with poor locality; others have continuous data streams (multimedia) and hence poor locality

2/22/2005

CS252 L11-review

50

Exception classifications

- **Traps:** relevant to the current process
 - Faults, arithmetic traps, and system "calls"
 - Invoke software on behalf of the currently executing process
- **Interrupts:** caused by asynchronous, outside events
 - I/O devices requiring service (DISK, network)
 - Clock interrupts (real time scheduling)
- **Machine Checks:** caused by serious hardware failure
 - Not always restartable
 - Indicate that bad things have happened.
 - » Non-recoverable ECC error
 - » Machine room fire
 - » Power outage

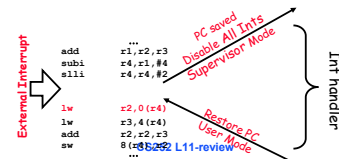
2/22/2005

CS252 L11-review

51

Precise Interrupts/Exceptions

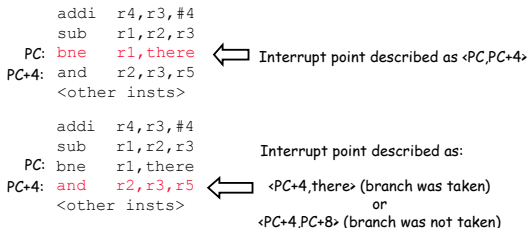
- An interrupt or exception is considered *precise* if there is a single instruction (or interrupt point) for which:
 - All instructions before that have committed their state
 - No following instructions (including the interrupting instruction) have modified any state.
- This means, that you can restart execution at the interrupt point and "get the right answer"
 - Implicit in our previous example of a device interrupt:
 - » Interrupt point is at first `lw` instruction



2/22/2005

52

Precise interrupt point may require multiple PCs



- On SPARC, interrupt hardware produces "pc" and "npc" (next pc)
- On MIPS, only "pc" – must fix point in software

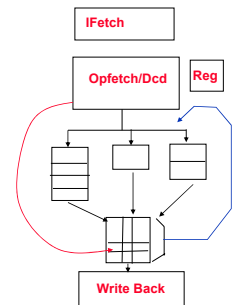
2/22/2005

CS252 L11-review

53

Reorder Buffer + Forwarding + Speculation

- **Idea:**
 - Issue branch into ROB
 - Mark with prediction
 - Fetch and issue predicted instructions speculatively
 - Branch must resolve before leaving ROB
 - Resolve correct
 - » Commit following instr
 - Resolve incorrect
 - » Mark following instr in ROB as invalid
 - » Let them clear



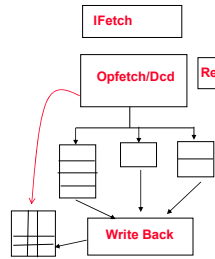
2/22/2005

CS252 L11-review

54

History File

- Maintain issue order, like ROB
- Each entry records dest reg and old value of dest. Register
 - What if old value not available when instruction issues?
- FUs write results into register file
 - Forward into correct entry in history file
- When exception reaches head
 - Restore architected registers from tail to head



2/22/2005

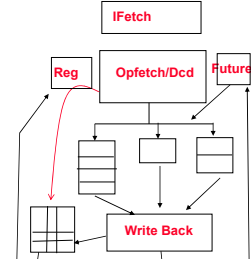
CS252 L11-review

55

Future file

Idea

- Arch registers reflect state at commit point
- Future register reflect whatever instructions have completed
- On WB update future
- On commit update arch
- On exception
 - » Discard future
 - » Replace with arch
 - Dest w/ ROB

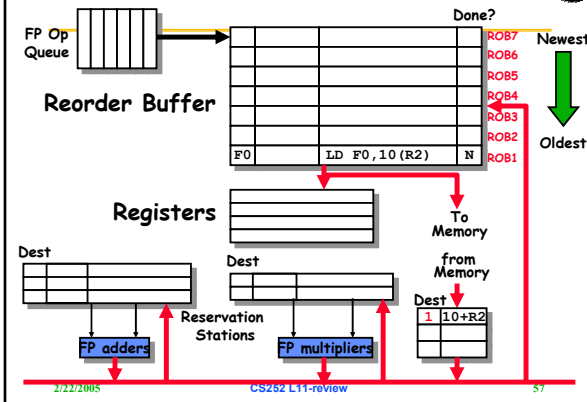


2/22/2005

CS252 L11-review

56

Tomasulo With Reorder buffer:



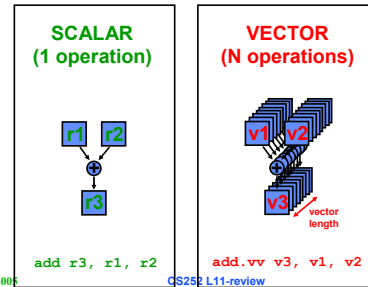
2/22/2005

CS252 L11-review

57

Alternative Model: Vector Processing

- Vector processors have high-level operations that work on linear arrays of numbers: "vectors"



2/22/2005

CS252 L11-review

58

25

What needs to be specified in a Vector Instruction Set Architecture?

- ISA in general
 - Operations, Data types, Format, Accessible Storage, Addressing Modes, Exceptional Conditions
- Vectors
 - Operations
 - Data types (Float, int, V op V, S op V)
 - Format
 - Source and Destination Operands
 - » Memory?, register?
 - Length
 - Successor (consecutive, stride, indexed, gather/scatter, ...)
 - Conditional operations
 - Exceptions

2/22/2005

CS252 L11-review

59

"DLX" Vector Instructions

| Instr. | Operands | Operation | Comment |
|----------------------|----------------|-------------------------|-------------------|
| • ADD _V | V1, V2, V3 | V1 = V2 + V3 | vector + vector |
| • ADD _S V | V1, F0, V2 | V1 = F0 + V2 | scalar + vector |
| • MULTV | V1, V2, V3 | V1 = V2 x V3 | vector x vector |
| • MULSV | V1, F0, V2 | V1 = F0 x V2 | scalar x vector |
| • LV | V1, R1 | V1 = M[R1..R1+63] | load, stride=1 |
| • LV _{WS} | V1, R1, R2 | V1 = M[R1..R1+63*R2] | load, stride=R2 |
| • LV _I | V1, R1, V2 | V1 = M[R1+V2i, i=0..63] | indir. ("gather") |
| • CeqV | VM, V1, V2 | VMASKi = (V1i = V2i)? | comp. setmask |
| • MOV | <u>V</u> L, R1 | Vec. Len. Reg. = R1 | set vector length |
| • MOV | <u>V</u> M, R1 | Vec. Mask = R1 | set vector mask |

2/22/2005

CS252 L11-review

60

Vector Execution Time

- Time = $f(\text{vector length, data dependencies, struct. hazards})$
- **Initiation rate**: rate that FU consumes vector elements (= number of lanes; usually 1 or 2 on Cray T-90)
- **Convoy**: set of vector instructions that can begin execution in same clock (no struct. or data hazards)
- **Chime**: approx. time for a vector operation
- **m convoys take m chimes**; if each vector length is n , then they take approx. $m \times n$ clock cycles (ignores overhead; good approximation for long vectors)

```

1: LV V1,Rx ;load vector X
2: MULV V2,F0,V1 ;vector-scalar mult.
   LV V3,Ry ;load vector Y
3: ADDV V4,V2,V3 ;add
4: SV Ry,V4 ;store the result
    
```

4 convoys, 1 lane, VL=64
 $\Rightarrow 4 \times 64 = 256$ clocks
 (or 4 clocks per result)

61

Strip Mining

- Suppose Vector Length > Max. Vector Length (MVL)?
- **Strip mining**: generation of code such that each vector operation is done for a size \leq to the MVL
- 1st loop do short piece ($n \bmod \text{MVL}$), rest VL = MVL


```

low = 1
VL = (n mod MVL) /*find the odd size piece*/
do 1 j = 0,(n / MVL) /*outer loop*/
    do 10 i = low,low+VL-1 /*runs for length VL*/
        Y(i) = a*X(i) + Y(i) /*main operation*/
    10 continue
    low = low+VL /*start of next vector*/
    VL = MVL /*reset the length to max*/
1 continue
            
```

2/22/2005

CS252 L11-review

62

Vector Opt #1: Chaining

- Suppose:
 MULV V1,V2,V3
 ADDV V4,V1,V5 ; separate convoy?
- **chaining**: vector register (V1) is not as a single entity but as a group of individual registers, then **pipeline forwarding can work on individual elements of a vector**
- **Flexible chaining**: allow vector to chain to any other active vector operation \Rightarrow more read/write ports
- As long as enough HW, increases convoy size

Unchained

Total=141

Chained

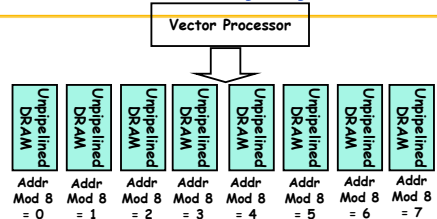
Total=77

2/22/2005

CS252 L11-review

63

Interleaved Memory Layout



- **Great for unit stride**:
 - Contiguous elements in different DRAMs
 - Startup time for vector operation is latency of single read
- **What about non-unit stride?**
 - Above good for strides that are relatively prime to 8
 - Bad for: 2, 4
 - Better: prime number of banks...

2/22/2005

CS252 L11-review

64