

## Lecture 2: Review of Instruction Sets, Pipelines, and Caches

Prof. John Kubiatowicz  
Computer Science 252  
Fall 1998

JDK.F98  
Slide 1

## Review, #1

- Technology is changing rapidly:
 

	<u>Capacity</u>	<u>Speed</u>
Logic	2x in 3 years	2x in 3 years
DRAM	4x in 3 years	2x in 10 years
Disk	4x in 3 years	2x in 10 years
Processor	( n.a.)	2x in 1.5 years
- What was true five years ago is not necessarily true now.
- Execution time is the REAL measure of computer performance!
  - Not clock rate, not CPI
- "X is n times faster than Y" means:

$$\frac{\text{ExTime}(y)}{\text{ExTime}(X)} = \frac{\text{Performance}(X)}{\text{Performance}(Y)}$$

JDK.F98  
Slide 2

## Review, #2

- Amdahl's Law: (or Law of Diminishing Returns)

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

- CPI Law:

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

- The "End to End Argument" is what RISC was ultimately about -- it is the performance of the complete system that matters, not individual components!

JDK.F98  
Slide 3

Today:  
Quick review of everything you  
should have learned

∞<sub>0</sub>

( A countably-infinite set of  
computer architecture concepts )

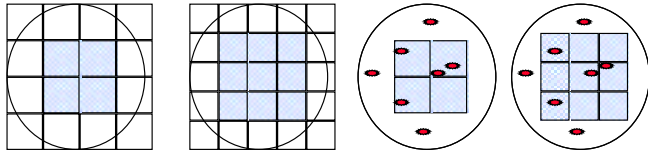
JDK.F98  
Slide 4

# Integrated Circuits Costs

$$IC \text{ cost} = \frac{\text{Die cost} + \text{Testing cost} + \text{Packaging cost}}{\text{Final test yield}}$$

$$\text{Die cost} = \frac{\text{Wafer cost}}{\text{Dies per Wafer} \times \text{Die yield}}$$

$$\text{Dies per wafer} = \frac{\pi (\text{Wafer\_diam}/2)^2}{\text{Die\_Area}} - \frac{\pi \times \text{Wafer\_diam}}{\sqrt{2} \cdot \text{Die\_Area}} - \text{Test\_Die}$$



$$\text{Die Yield} = \text{Wafer\_yield} \times \left\{ 1 + \left( \frac{\text{Defect\_Density} \times \text{Die\_area}}{\alpha} \right)^\alpha \right\}$$

Die Cost goes roughly with die area<sup>4</sup>

JDK.F98  
Slide 5

# Real World Examples

Chip	Metal layers	Line width	Wafer cost	Defect /cm <sup>2</sup>	Area mm <sup>2</sup>	Dies/ wafer	Yield	Die Cost
386DX	2	0.90	\$900	1.0	43	360	71%	\$4
486DX2	3	0.80	\$1200	1.0	81	181	54%	\$12
PowerPC 601	4	0.80	\$1700	1.3	121	115	28%	\$53
HP PA 7100	3	0.80	\$1300	1.0	196	66	27%	\$73
DEC Alpha	3	0.70	\$1500	1.2	234	53	19%	\$149
SuperSPARC	3	0.70	\$1700	1.6	256	48	13%	\$272
Pentium	3	0.80	\$1500	1.5	296	40	9%	\$417

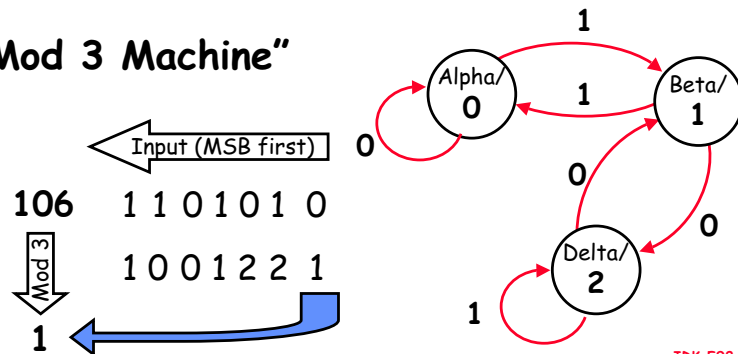
– From "Estimating IC Manufacturing Costs," by Linley Gwennap, *Microprocessor Report*, August 2, 1993, p. 15

JDK.F98  
Slide 6

# Finite State Machines:

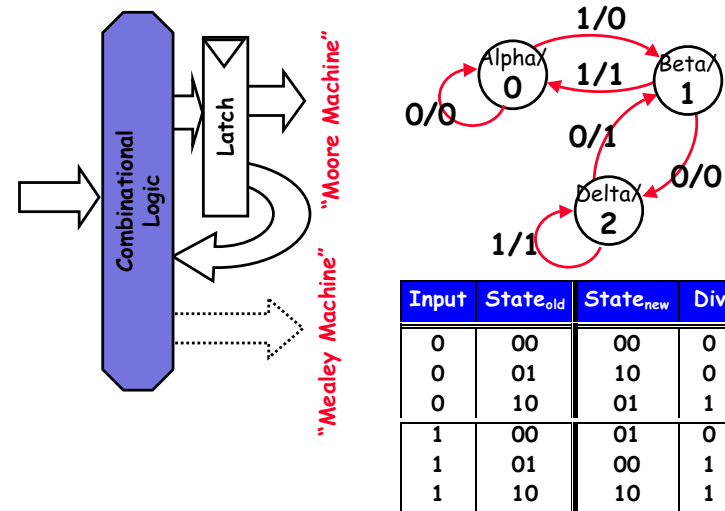
- System state is explicit in representation
- Transitions between states represented as arrows with inputs on arcs.
- Output may be either part of state or on arcs

## "Mod 3 Machine"



JDK.F98  
Slide 7

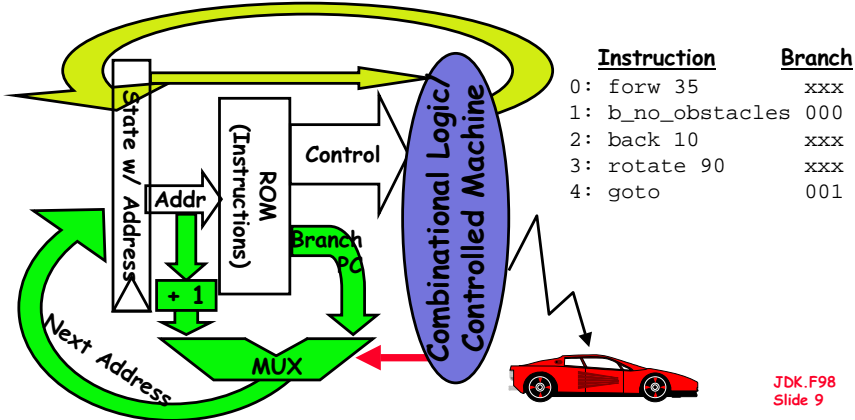
# Implementation as Combinational logic + Latch



JDK.F98  
Slide 8

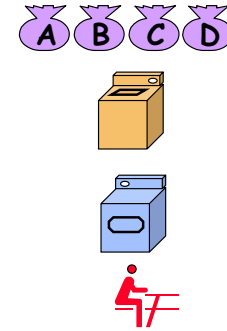
## Microprogrammed Controllers

- State machine in which part of state is a "micro-pc".
  - Explicit circuitry for incrementing or changing PC
- Includes a ROM with "microinstructions".
  - Controlled logic implements at least branches and jumps



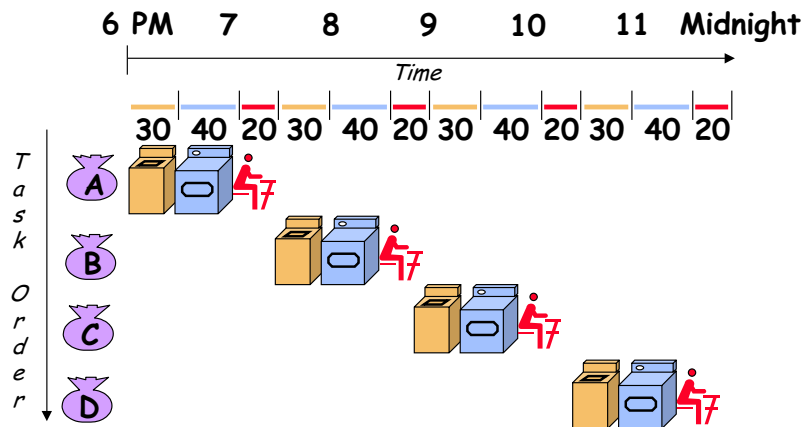
## Pipelining: Its Natural!

- Laundry Example
- Ann, Brian, Cathy, Dave each have one load of clothes to wash, dry, and fold
- Washer takes 30 minutes
- Dryer takes 40 minutes
- "Folder" takes 20 minutes



JDK.F98  
Slide 10

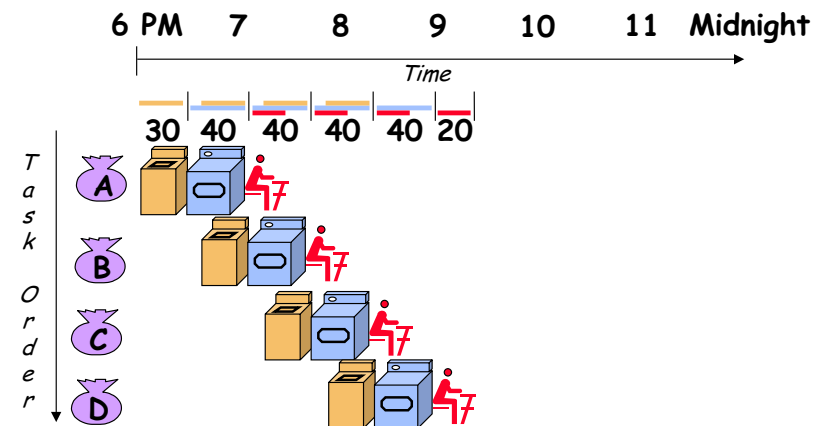
## Sequential Laundry



- Sequential laundry takes 6 hours for 4 loads
- If they learned pipelining, how long would laundry take?

JDK.F98  
Slide 11

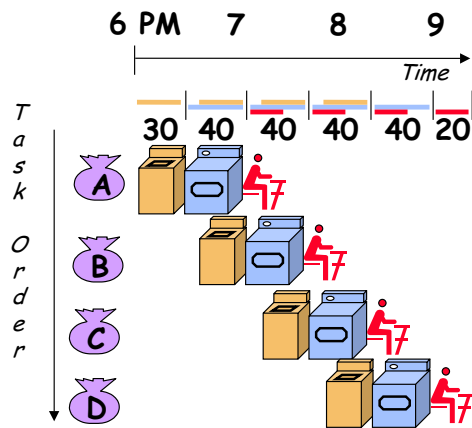
## Pipelined Laundry Start work ASAP



- Pipelined laundry takes 3.5 hours for 4 loads

JDK.F98  
Slide 12

## Pipelining Lessons



- Pipelining doesn't help **latency** of single task, it helps **throughput** of entire workload
- Pipeline rate limited by **slowest** pipeline stage
- **Multiple** tasks operating simultaneously
- Potential speedup = **Number pipe stages**
- Unbalanced lengths of pipe stages reduces speedup
- Time to "fill" pipeline and time to "drain" it reduces speedup

JDK.F98  
Slide 13

## Computer Pipelines

- Execute billions of instructions, so **throughput** is what matters
- DLX desirable features: all instructions same length, registers located in same place in instruction format, memory operands only in loads or stores

JDK.F98  
Slide 14

## A "Typical" RISC

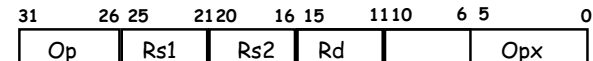
- 32-bit fixed format instruction (3 formats)
- 32 32-bit GPR (R0 contains zero, DP take pair)
- 3-address, reg-reg arithmetic instruction
- Single address mode for load/store: base + displacement
  - no indirection
- Simple branch conditions
- Delayed branch

see: SPARC, MIPS, HP PA-Risc, DEC Alpha, IBM PowerPC, CDC 6600, CDC 7600, Cray-1, Cray-2, Cray-3

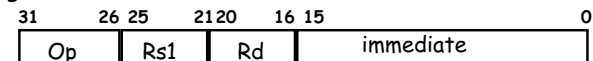
JDK.F98  
Slide 15

## Example: MIPS (- DLX)

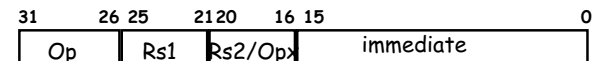
### Register-Register



### Register-Immediate



### Branch



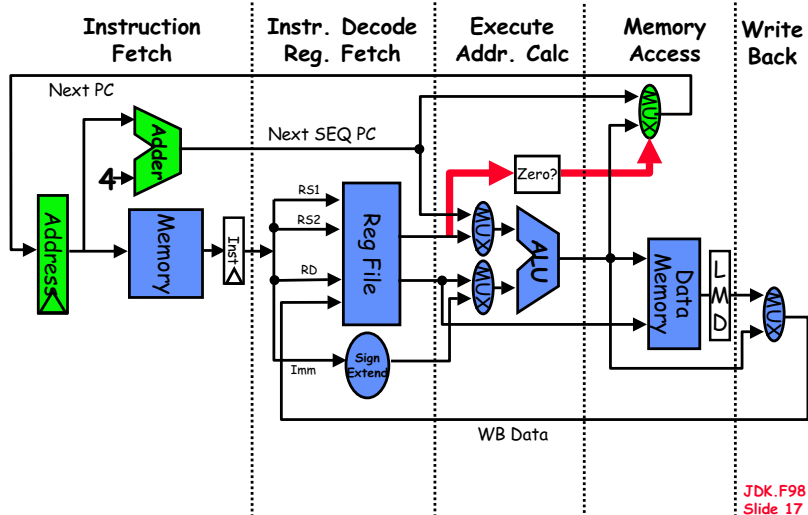
### Jump / Call



JDK.F98  
Slide 16

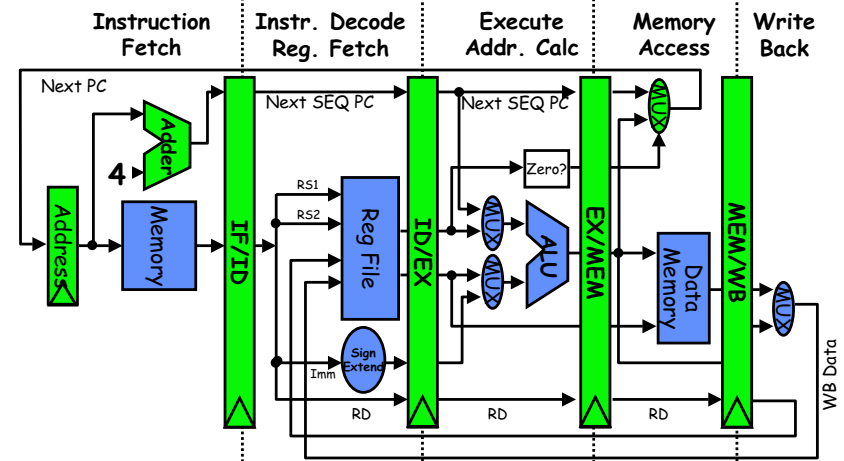
## 5 Steps of DLX Datapath

Figure 3.1, Page 130



## 5 Steps of DLX Datapath

Figure 3.4, Page 137

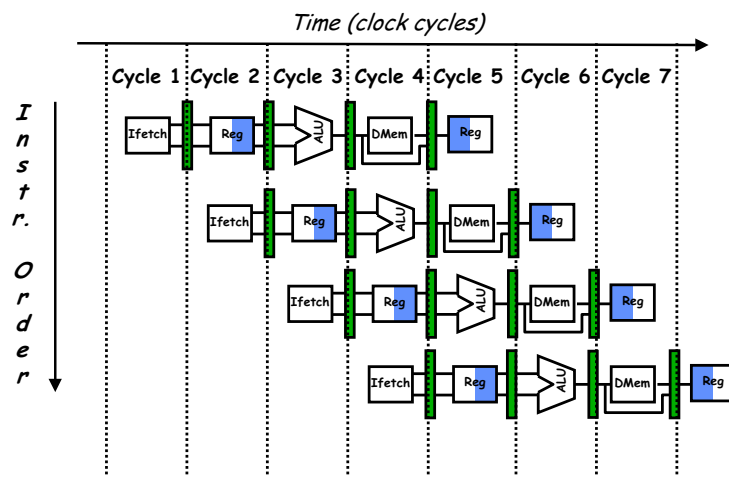


- Data stationary control
  - local decode for each instruction phase / pipeline stage

JDK.F98  
Slide 18

## Visualizing Pipelining

Figure 3.3, Page 133



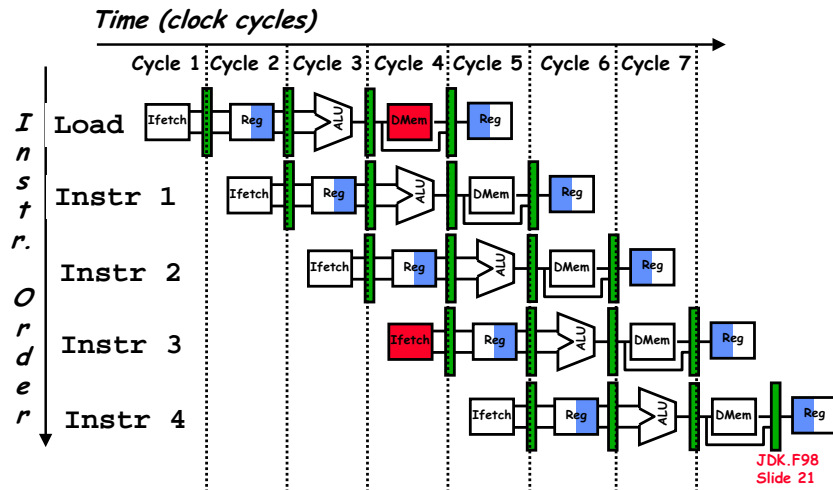
## Its Not That Easy for Computers

- 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).

JDK.F98  
Slide 20

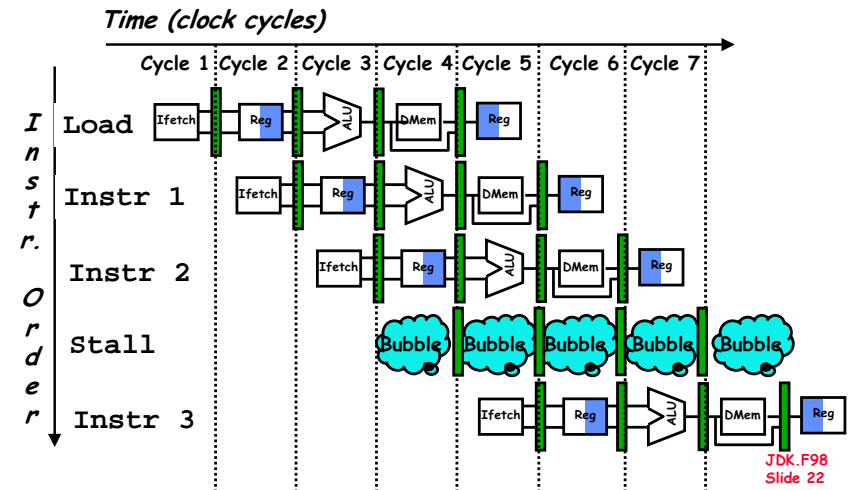
## One Memory Port/Structural Hazards

Figure 3.6, Page 142



## One Memory Port/Structural Hazards

Figure 3.7, Page 143



## 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,  $CPI = 1$ :

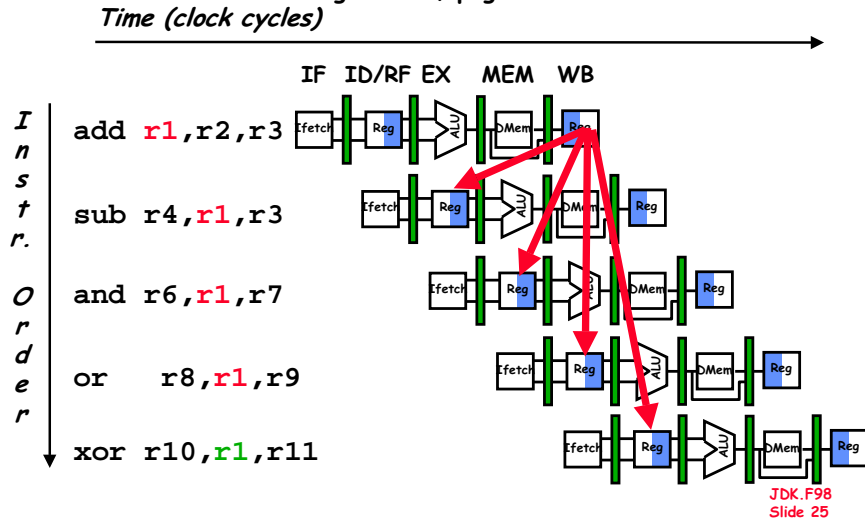
$$\text{Speedup} = \frac{\text{Pipeline depth}}{1 + \text{Pipeline stall CPI}} \times \frac{\text{Cycle Time}_{\text{unpipelined}}}{\text{Cycle Time}_{\text{pipelined}}}$$

## Example: Dual-port vs. Single-port

- Machine A: Dual ported memory ("Harvard Architecture")
- Machine B: Single ported memory, but its pipelined implementation has a 1.05 times faster clock rate
- Ideal CPI = 1 for both
- Loads are 40% of instructions executed
  - $\text{SpeedUp}_A = \text{Pipeline Depth} / (1 + 0) \times (\text{clock}_{\text{unpipe}} / \text{clock}_{\text{pipe}})$
  - $= \text{Pipeline Depth}$
  - $\text{SpeedUp}_B = \text{Pipeline Depth} / (1 + 0.4 \times 1) \times (\text{clock}_{\text{unpipe}} / (\text{clock}_{\text{unpipe}} / 1.05))$
  - $= (\text{Pipeline Depth} / 1.4) \times 1.05$
  - $= 0.75 \times \text{Pipeline Depth}$
  - $\text{SpeedUp}_A / \text{SpeedUp}_B = \text{Pipeline Depth} / (0.75 \times \text{Pipeline Depth}) = 1.33$
- Machine A is 1.33 times faster

## Data Hazard on R1

Figure 3.9, page 147



## Three Generic Data Hazards

- **Read After Write (RAW)**  
Instr<sub>J</sub> tries to read operand before Instr<sub>I</sub> writes it

↻ I: add r1,r2,r3  
 ↻ J: sub r4,r1,r3

- Caused by a "**Dependence**" (in compiler nomenclature). This hazard results from an actual need for communication.

JDK.F98  
Slide 26

## Three Generic Data Hazards

- **Write After Read (WAR)**  
Instr<sub>J</sub> writes operand before Instr<sub>I</sub> reads it

↻ I: sub r4,r1,r3  
 ↻ J: add r1,r2,r3  
 K: mul r6,r1,r7

- Called an "**anti-dependence**" by compiler writers. This results from reuse of the name "r1".
- Can't happen in DLX 5 stage pipeline because:
  - All instructions take 5 stages, and
  - Reads are always in stage 2, and
  - Writes are always in stage 5

JDK.F98  
Slide 27

## Three Generic Data Hazards

- **Write After Write (WAW)**  
Instr<sub>J</sub> writes operand before Instr<sub>I</sub> writes it.

↻ I: sub r1,r4,r3  
 ↻ J: add r1,r2,r3  
 K: mul r6,r1,r7

- Called an "**output dependence**" by compiler writers. This also results from the reuse of name "r1".
- Can't happen in DLX 5 stage pipeline because:
  - All instructions take 5 stages, and
  - Writes are always in stage 5
- Will see WAR and WAW in later more complicated pipes

JDK.F98  
Slide 28

## CS 252 Administrivia

- Sign up today! Web site is:  
<http://www.cs.berkeley.edu/~kubitron/cs252-F98>
- **In class exam on Wednesday Sept 2nd**
  - Improve 252 experience if recapture common background
  - **Bring 1 sheet of paper with notes on both sides**
  - Doesn't affect grade, only admission into class
  - 2 grades: Admitted or audit/take CS 152 1st (before class Friday)
- Review: Chapters 1- 3, CS 152 home page, maybe "Computer Organization and Design (COD)2/e"
  - If did take a class, be sure COD Chapters 2, 5, 6, 7 are familiar
  - Copies in Bechtel Library on 2-hour reserve
- Aaron Brown will be holding a review session this weekend: Sunday, Aug 30. 1:00pm. 310 Soda

JDK.F98  
Slide 29

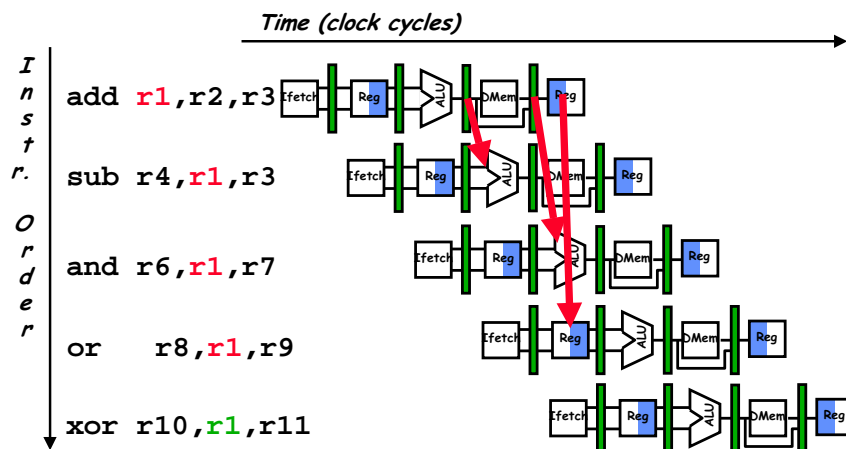
## CS 252 Administrivia

- Resources for course on web site:
  - Check out the ISCA (International Symposium on Computer Architecture) **25th year retrospective** on web site. Look for "Additional reading" below text-book description
  - Pointers to previous CS152 exams and resources
  - Lots of old CS252 material
  - Interesting pointers at bottom. Check out the: **WWW Computer Architecture Home Page**
- To give proper attention to projects (as well as homeworks and quizzes), I can handle up to 36 students
  - First priority is students taking ARCH prelims in next year
  - Second priority is students taking this for breadth
  - Third priority is EECS students
  - Fourth priority College of Engineering grad students

JDK.F98  
Slide 30

## Forwarding to Avoid Data Hazard

Figure 3.10, Page 149



JDK.F98  
Slide 31

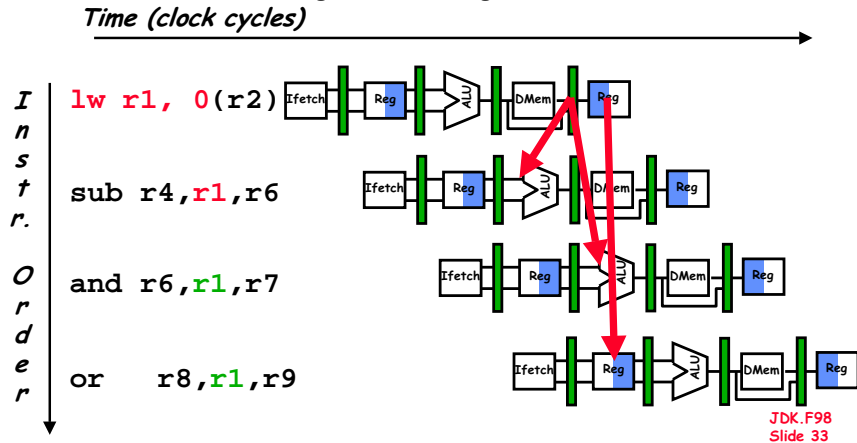
## HW Change for Forwarding

Figure 3.20, Page 161

JDK.F98  
Slide 32

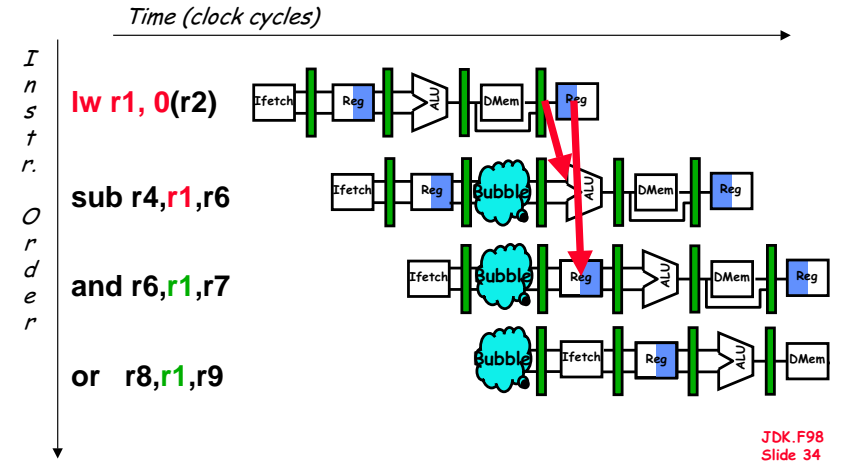
## Data Hazard Even with Forwarding

Figure 3.12, Page 153



## Data Hazard Even with Forwarding

Figure 3.13, Page 154



## Software Scheduling to Avoid Load Hazards

Try producing fast code for

a = b + c;

d = e - f;

assuming a, b, c, d, e, and f in memory.

Slow code:

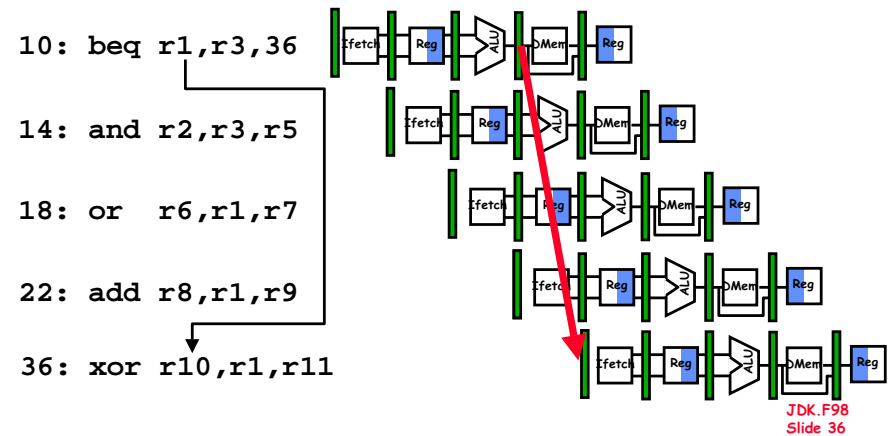
LW Rb,b  
LW Rc,c  
ADD Ra,Rb,Rc  
SW a,Ra  
LW Re,e  
LW Rf,f  
SUB Rd,Re,Rf  
SW d,Rd

Fast code:

LW Rb,b  
LW Rc,c  
LW Re,e  
ADD Ra,Rb,Rc  
LW Rf,f  
SW a,Ra  
SUB Rd,Re,Rf  
SW d,Rd

JDK.F98 Slide 35

## Control Hazard on Branches Three Stage Stall



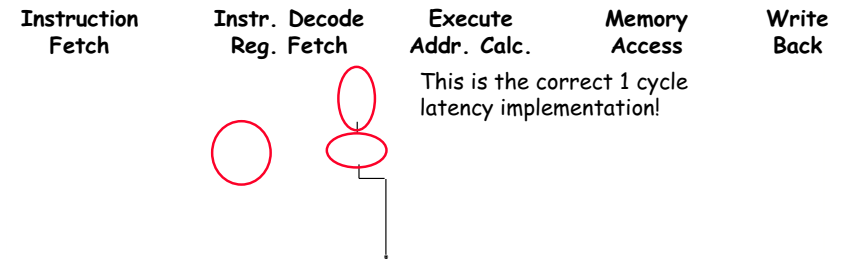
## Branch Stall Impact

- If  $CPI = 1$ , 30% branch, Stall 3 cycles => new  $CPI = 1.9!$
- Two part solution:
  - Determine branch taken or not sooner, AND
  - Compute taken branch address earlier
- DLX branch tests if register = 0 or  $\neq$
- 0
- DLX Solution:
  - Move Zero test to ID/RF stage
  - Adder to calculate new PC in ID/RF stage
  - 1 clock cycle penalty for branch versus 3

JDK.F98  
Slide 37

## Pipelined DLX Datapath

Figure 3.22, page 163



JDK.F98  
Slide 38

## Four Branch Hazard Alternatives

**#1: Stall until branch direction is clear**

**#2: Predict Branch Not Taken**

- Execute successor instructions in sequence
- "Squash" instructions in pipeline if branch actually taken
- Advantage of late pipeline state update
- 47% DLX branches not taken on average
- PC+4 already calculated, so use it to get next instruction

**#3: Predict Branch Taken**

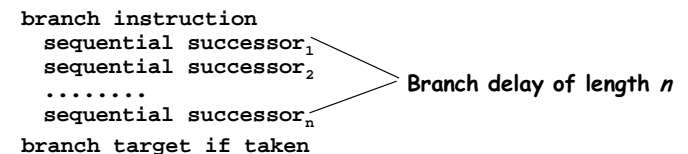
- 53% DLX branches taken on average
- But haven't calculated branch target address in DLX
  - » DLX still incurs 1 cycle branch penalty
  - » Other machines: branch target known before outcome

JDK.F98  
Slide 39

## Four Branch Hazard Alternatives

**#4: Delayed Branch**

- Define branch to take place **AFTER** a following instruction



- 1 slot delay allows proper decision and branch target address in 5 stage pipeline
- DLX uses this

JDK.F98  
Slide 40

## Delayed Branch

- Where to get instructions to fill branch delay slot?
  - Before branch instruction
  - From the target address: only valuable when branch taken
  - From fall through: only valuable when branch not taken
  - Cancelling branches allow more slots to be filled
- Compiler effectiveness for single branch delay slot:
  - Fills about 60% of branch delay slots
  - About 80% of instructions executed in branch delay slots useful in computation
  - About 50% (60% x 80%) of slots usefully filled
- Delayed Branch downside: 7-8 stage pipelines, multiple instructions issued per clock (superscalar)

JDK.F98  
Slide 41

## Evaluating Branch Alternatives

$$\text{Pipeline speedup} = \frac{\text{Pipeline depth}}{1 + \text{Branch frequency} \times \text{Branch penalty}}$$

Scheduling scheme	Branch penalty	CPI	speedup v. unpipelined	speedup v. stall
Stall pipeline	3	1.42	3.5	1.0
Predict taken	1	1.14	4.4	1.26
Predict not taken	1	1.09	4.5	1.29
Delayed branch	0.5	1.07	4.6	<b>1.31</b>

Conditional & Unconditional = 14%, 65% change PC

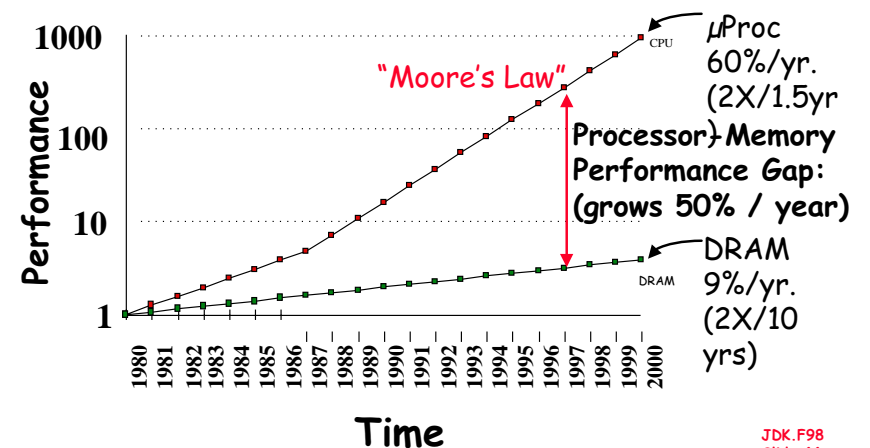
JDK.F98  
Slide 42

## Now, Review of Memory Hierarchy

JDK.F98  
Slide 43

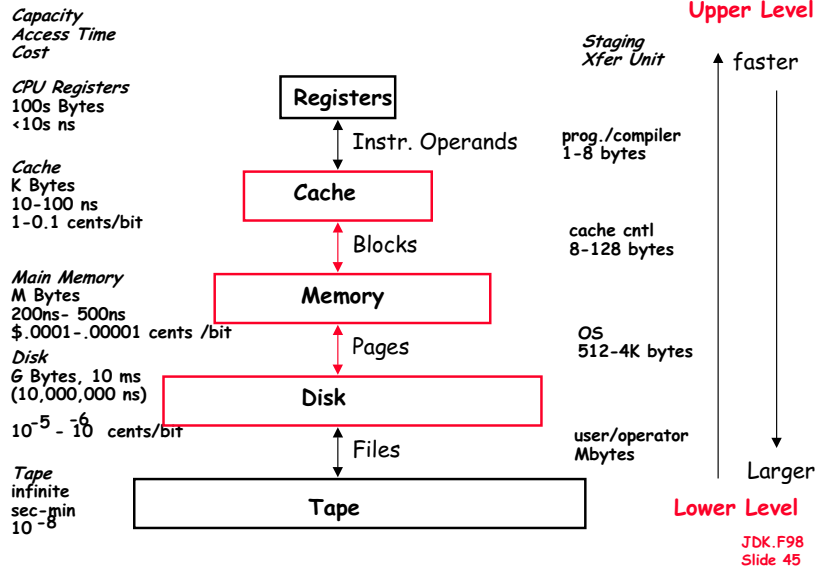
## Recap: Who Cares About the Memory Hierarchy?

Processor-DRAM Memory Gap (latency)



JDK.F98  
Slide 44

# Levels of the Memory Hierarchy



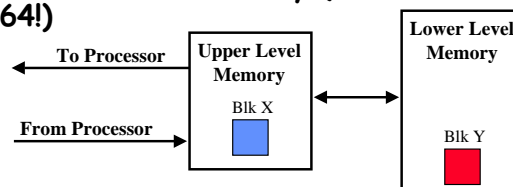
# The Principle of Locality

- 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

JDK.F98 Slide 46

# 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!)



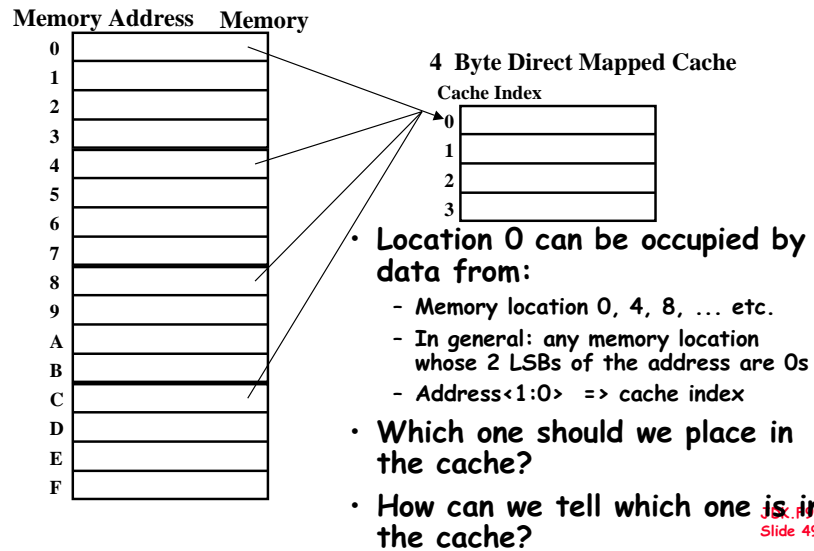
JDK.F98 Slide 47

# Cache Measures

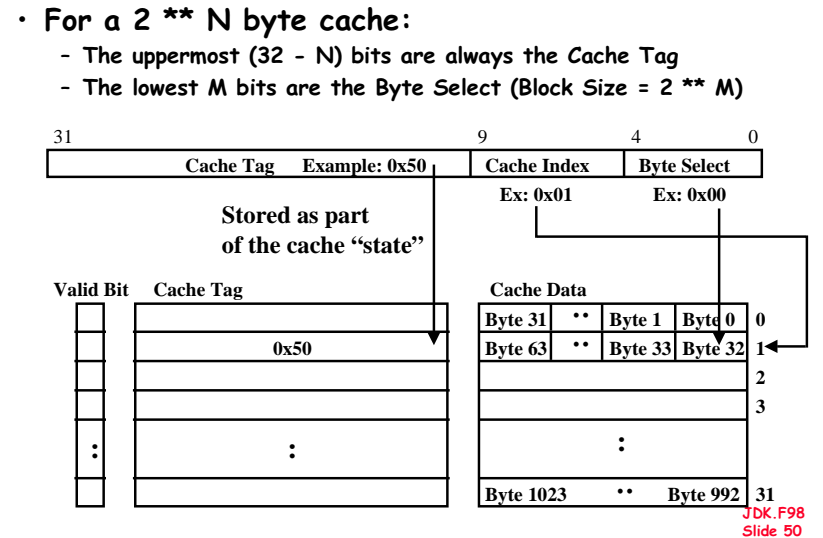
- **Hit rate**: fraction found in that level
  - So high that usually talk about **Miss rate**
  - Miss rate fallacy: as MIPS to CPU performance, miss rate to average memory access time in memory
- Average memory-access time  
 = Hit time + Miss rate x Miss penalty (ns or clocks)
- **Miss penalty**: 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)

JDK.F98 Slide 48

# Simplest Cache: Direct Mapped

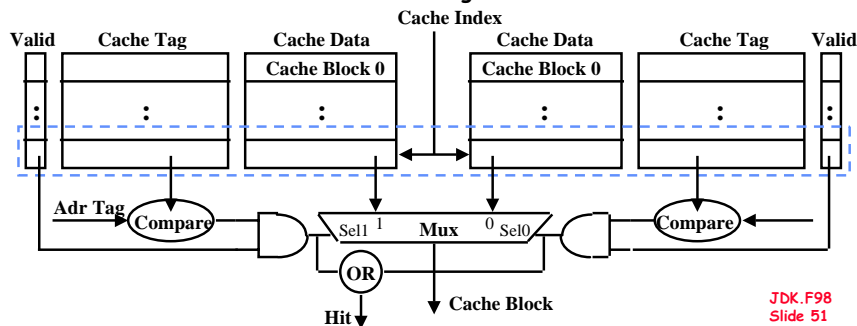


# 1 KB Direct Mapped Cache, 32B blocks



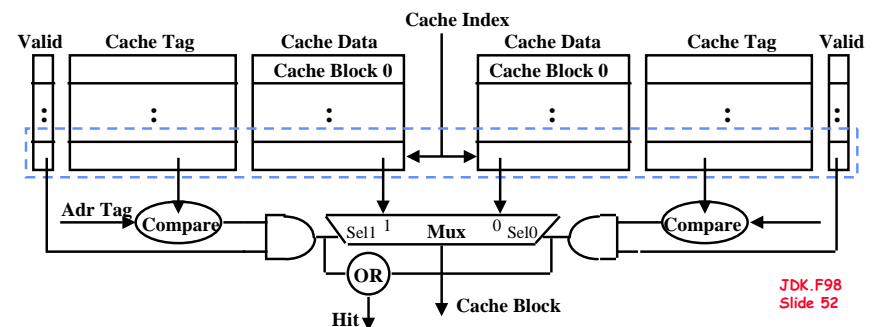
# Two-way Set Associative Cache

- 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



# 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.



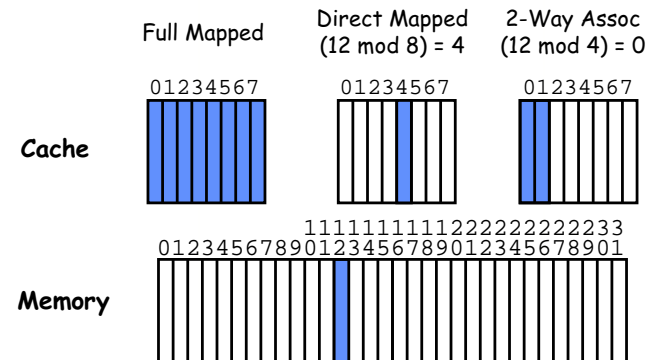
## 4 Questions for Memory Hierarchy

- 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)*

JDK.F98  
Slide 53

## 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



JDK.F98  
Slide 54

## 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 →

Block Address		Block Offset
Tag	Index	

JDK.F98  
Slide 55

## 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%

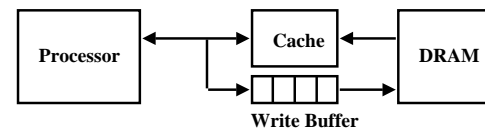
JDK.F98  
Slide 56

## 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

JDK.F98  
Slide 57

## 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

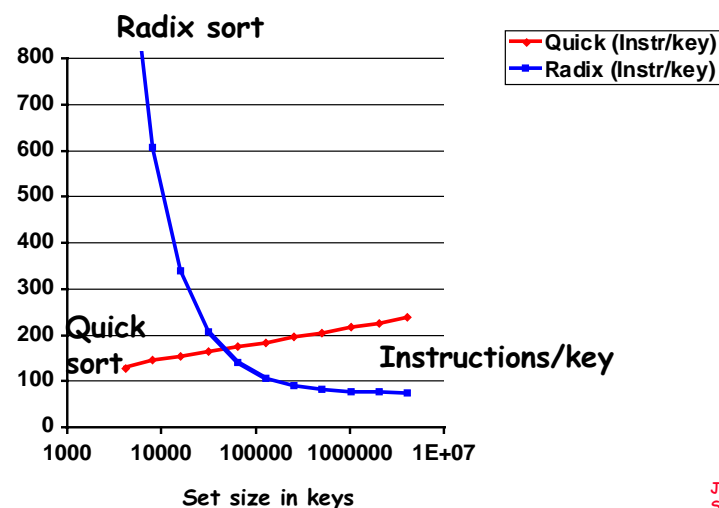
JDK.F98  
Slide 58

## 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

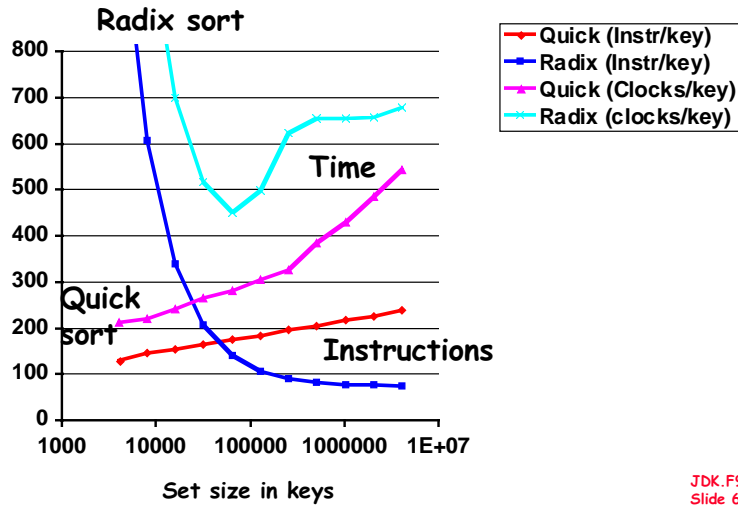
JDK.F98  
Slide 59

## Quicksort vs. Radix as vary number keys: Instructions



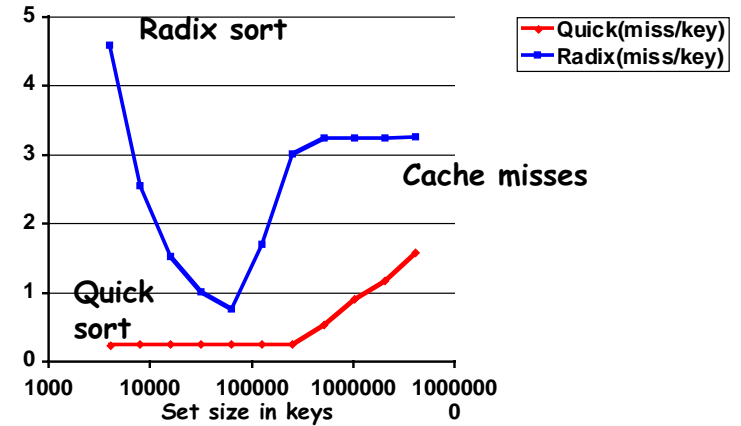
JDK.F98  
Slide 60

## Quicksort vs. Radix as vary number keys: Instrs & Time



JDK.F98  
Slide 61

## Quicksort vs. Radix as vary number keys: Cache misses



What is proper approach to fast algorithms?

JDK.F98  
Slide 62

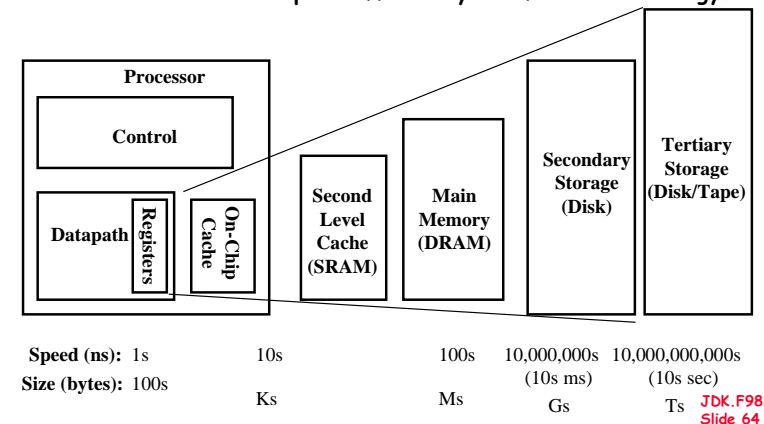
## 5 minute Class Break

- 80 minutes straight is too long for me to lecture (12:40:00 - 2:00:00):
  - 1 minute: review last time & motivate this lecture
  - 20 minute lecture
  - 3 minutes: **discuss class management**
  - 25 minutes: lecture
  - 5 minutes: **break**
  - 25 minutes: lecture
  - 1 minute: summary of today's important topics

JDK.F98  
Slide 63

## A Modern Memory Hierarchy

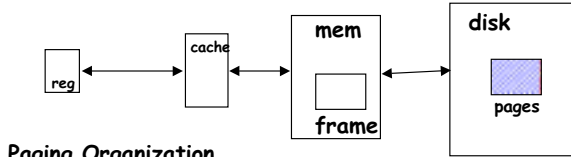
- By taking advantage of the principle of locality:
  - Present the user with as much memory as is available in the cheapest technology.
  - Provide access at the speed offered by the fastest technology.



JDK.F98  
Slide 64

# Basic Issues in VM System Design

- size of information blocks that are transferred from secondary to main storage (M)
- block of information brought into M, and M is full, then some region of M must be released to make room for the new block --> *replacement policy*
- which region of M is to hold the new block --> *placement policy*
- missing item fetched from secondary memory only on the occurrence of a fault --> *demand load policy*



## Paging Organization

virtual and physical address space partitioned into blocks of equal size  
 pages  
 page frames

JDK.F98  
Slide 65

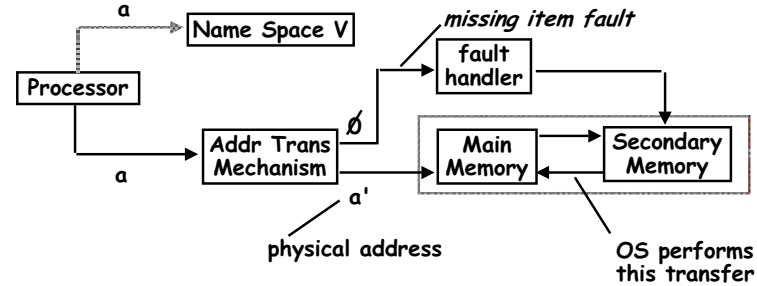
# Address Map

$V = \{0, 1, \dots, n - 1\}$  virtual address space  $n > m$   
 $M = \{0, 1, \dots, m - 1\}$  physical address space

MAP:  $V \rightarrow M \cup \{0\}$  address mapping function

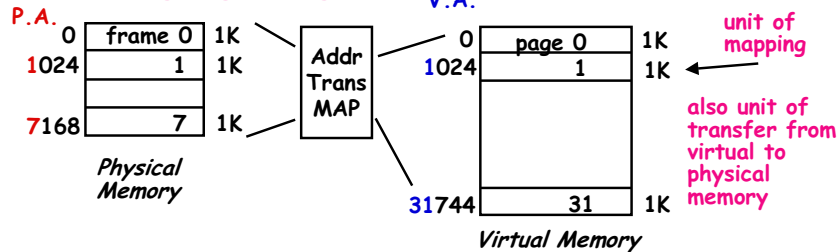
MAP(a) = a' if data at virtual address a is present in physical address a' and a' in M

/ = 0 if data at virtual address a is not present in M

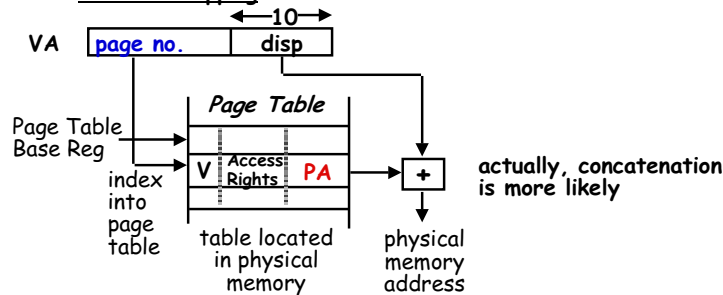


JDK.F98  
Slide 66

# Paging Organization

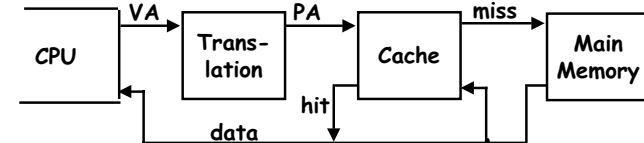


## Address Mapping



JDK.F98  
Slide 67

## Virtual Address and a Cache



It takes an extra memory access to translate VA to PA

This makes cache access very expensive, and this is the "innermost loop" that you want to go as fast as possible

ASIDE: Why access cache with PA at all? VA caches have a problem!  
*synonym / alias problem*: two different virtual addresses map to same physical address => two different cache entries holding data the same physical address!

for update: must update all cache entries with same physical address or memory becomes inconsistent

determining this requires significant hardware, essentially an associative lookup on the physical address tags to see if you have multiple hits; or

software enforced *alias boundary*: same lsb of VA & PA > cache size

JDK.F98  
Slide 68

## TLBs

A way to speed up translation is to use a special cache of recently used page table entries -- this has many names, but the most frequently used is *Translation Lookaside Buffer* or *TLB*

Virtual Address	Physical Address	Dirty	Ref	Valid	Access

Really just a cache on the page table mappings

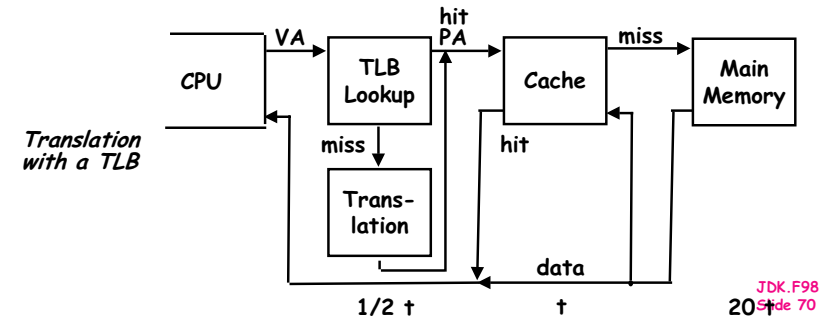
TLB access time comparable to cache access time  
(much less than main memory access time)

JDK.F98  
Slide 69

## Translation Look-Aside Buffers

Just like any other cache, the TLB can be organized as fully associative, set associative, or direct mapped

TLBs are usually small, typically not more than 128 - 256 entries even on high end machines. This permits fully associative lookup on these machines. Most mid-range machines use small n-way set associative organizations.



## Reducing Translation Time

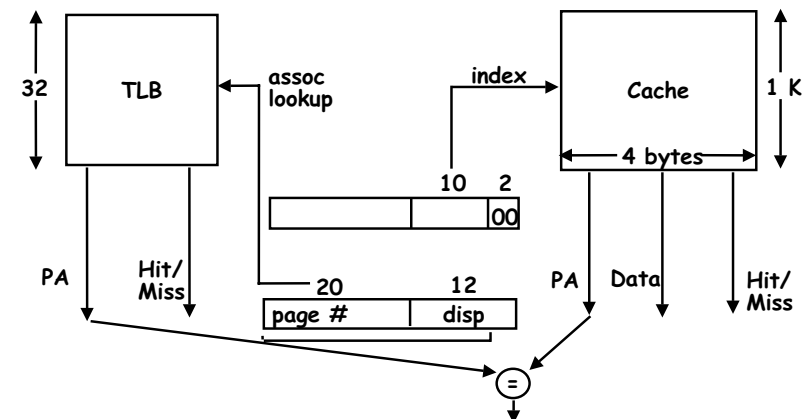
Machines with TLBs go one step further to reduce # cycles/cache access

They overlap the cache access with the TLB access:

high order bits of the VA are used to look in the TLB while low order bits are used as index into cache

JDK.F98  
Slide 71

## Overlapped Cache & TLB Access



IF cache hit AND (cache tag = PA) then deliver data to CPU  
ELSE IF [cache miss OR (cache tag = PA)] and TLB hit THEN  
access memory with the PA from the TLB  
ELSE do standard VA translation

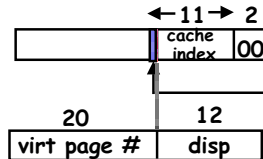
JDK.F98  
Slide 72

## Problems With Overlapped TLB Access

Overlapped access only works as long as the address bits used to index into the cache *do not change* as the result of VA translation

This usually limits things to small caches, large page sizes, or high n-way set associative caches if you want a large cache

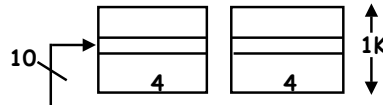
Example: suppose everything the same except that the cache is increased to 8 K bytes instead of 4 K:



This bit is changed by VA translation, but is needed for cache lookup

Solutions:

go to 8K byte page sizes;  
go to 2 way set associative cache; or  
SW guarantee VA[13]=PA[13]



2 way set assoc cache

JDK.F98  
Slide 73

## Summary #1/5: Control and Pipelining

- Control VIA **State Machines** and **Microprogramming**
- Just overlap tasks; easy if tasks are independent
- Speed Up  $\leq$  Pipeline Depth; if ideal CPI is 1, then:

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

- Hazards limit performance on computers:

- Structural: need more HW resources
- Data (RAW,WAR,WAW): need forwarding, compiler scheduling
- Control: delayed branch, prediction

JDK.F98  
Slide 74

## Summary #2/5: Caches

- The Principle of Locality:

- Program access a relatively small portion of the address space at any instant of time.
  - » Temporal Locality: Locality in Time
  - » Spatial Locality: Locality in Space

- Three Major Categories of Cache Misses:

- **Compulsory Misses**: sad facts of life. Example: cold start misses.
- **Capacity Misses**: increase cache size
- **Conflict Misses**: increase cache size and/or associativity.  
Nightmare Scenario: ping pong effect!

- Write Policy:

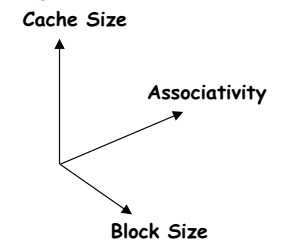
- **Write Through**: needs a **write buffer**. Nightmare: WB saturation
- **Write Back**: control can be complex

JDK.F98  
Slide 75

## Summary #3/5: The Cache Design Space

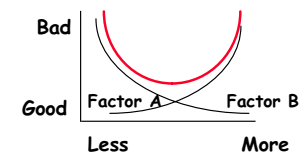
- Several interacting dimensions

- cache size
- block size
- associativity
- replacement policy
- write-through vs write-back
- write allocation



- The optimal choice is a compromise

- depends on access characteristics
  - » workload
  - » use (I-cache, D-cache, TLB)
- depends on technology / cost



- Simplicity often wins

JDK.F98  
Slide 76

## Summary #4/5: TLB, Virtual Memory

- Caches, TLBs, Virtual Memory all understood by examining how they deal with 4 questions: 1) Where can block be placed? 2) How is block found? 3) What block is replaced on miss? 4) How are writes handled?
- Page tables map virtual address to physical address
- TLBs are important for fast translation
- TLB misses are significant in processor performance
  - funny times, as most systems can't access all of 2nd level cache without TLB misses!

JDK.F98  
Slide 77

## Summary #5/5: Memory Hierarchy

- Virtual memory was controversial at the time: can SW automatically manage 64KB across many programs?
  - 1000X DRAM growth removed the controversy
- Today VM allows many processes to share single memory without having to swap all processes to disk; today VM protection is more important than memory hierarchy
- Today CPU time is a function of (ops, cache misses) vs. just  $f(\text{ops})$ :  
What does this mean to Compilers, Data structures, Algorithms?

JDK.F98  
Slide 78