# REVIEW OF INSTRUCTION SETS, PIPELINES

12 March 2020

Dr Noor Mahammad Sk

# Computer Architecture Is

□ 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.
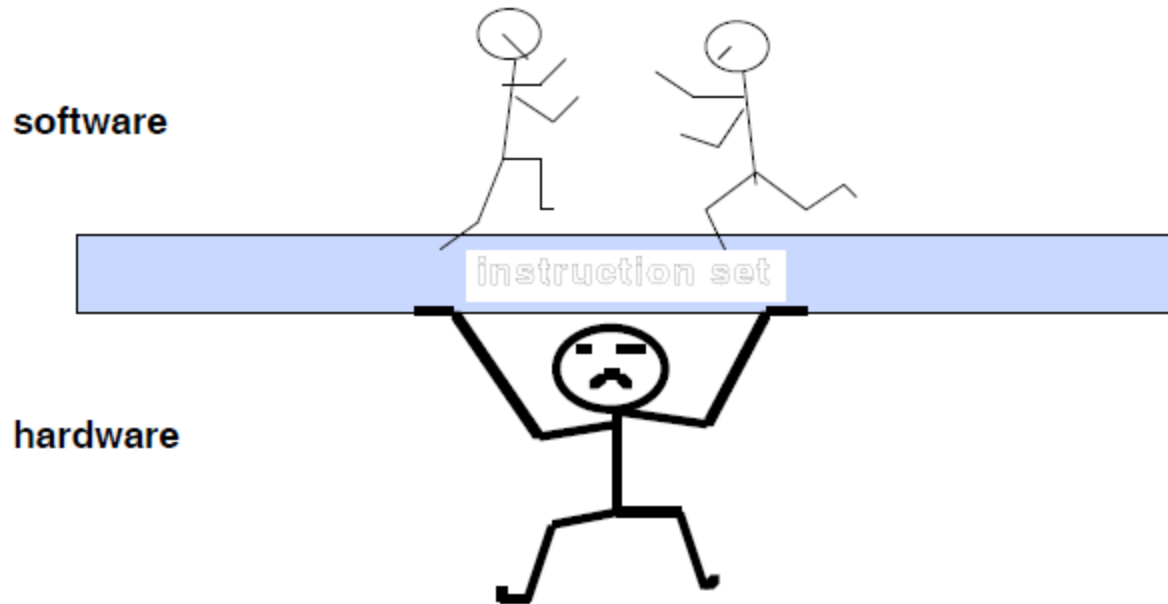
□ Amdahl, Blaaw, and Brooks, 1964

# Computer Architecture's Changing Definition

- 1950s to 1960s:
  - Computer Architecture Course = Computer Arithmetic
- 1970s to mid 1980s:
  - Computer Architecture Course = Instruction Set Design, especially ISA appropriate for compilers
- 1990s to till 2008:
  - Computer Architecture Course = Design of CPU, memory system, I/O system, Multiprocessors
- 2011:
  - Reconfigurable, Application specific computing/Embedded Computing

# Instruction Set Architecture (ISA)
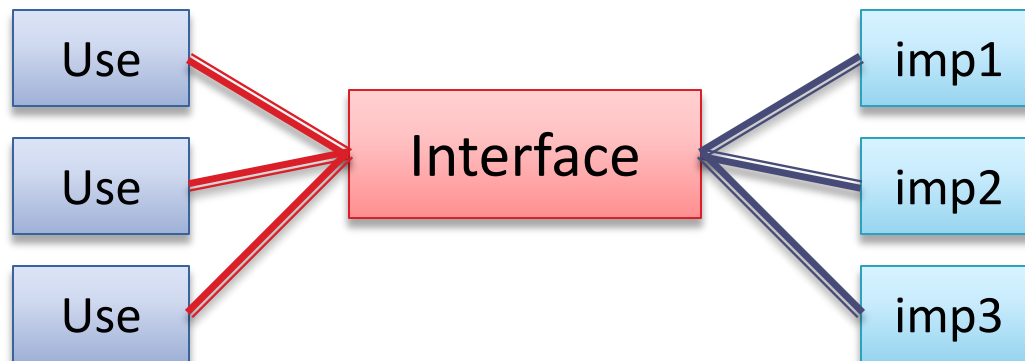
software

instruction set

hardware

# Interface Design
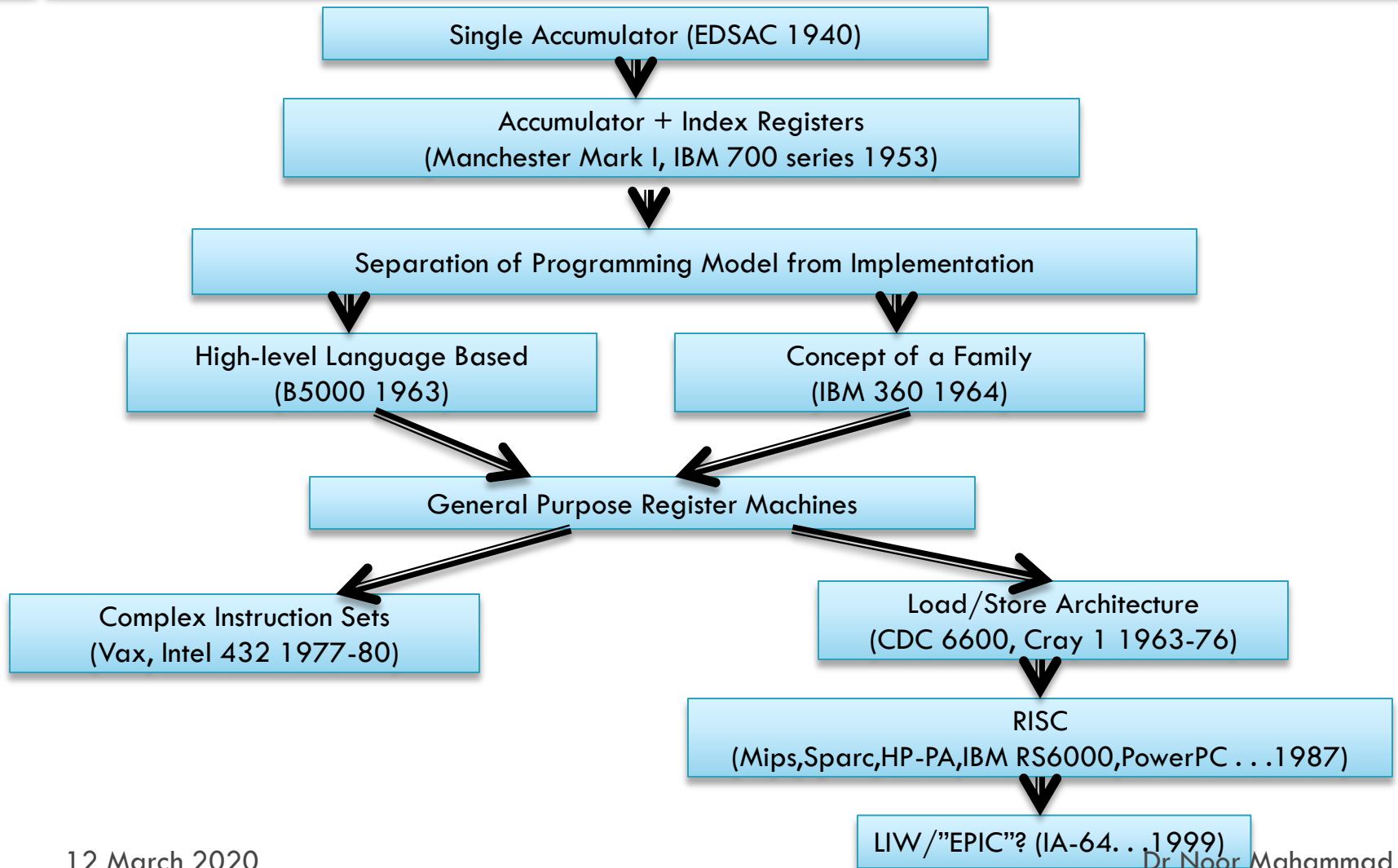
- A good interface:
  - Lasts through many implementations (portability, compatibility)
  - Is used in many different ways (generality)
  - Provides convenient functionality to higher levels
  - Permits an efficient implementation at lower level

```
[Use]  \
[Use]  -->  [Interface]  -->  [imp1]
[Use]  /                      [imp2]
                              [imp3]
```

# Evolution of Instruction Sets

Single Accumulator (EDSAC 1940)

Accumulator + Index Registers
(Manchester Mark I, IBM 700 series 1953)

Separation of Programming Model from Implementation

High-level Language Based
(B5000 1963)

Concept of a Family
(IBM 360 1964)

General Purpose Register Machines

Complex Instruction Sets
(Vax, Intel 432 1977-80)

Load/Store Architecture
(CDC 6600, Cray 1 1963-76)

RISC
(Mips,Sparc,HP-PA,IBM RS6000,PowerPC . . .1987)

LIW/"EPIC"? (IA-64. . .1999)

12 March 2020
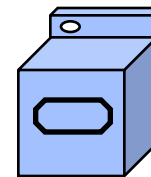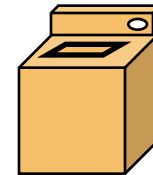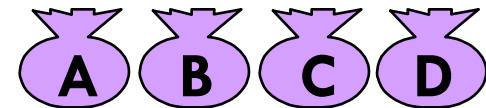
Dr Noor Mahammad Sk

# Evolution of Instruction Sets

- ☐ Major advances in computer architecture are typically associated with landmark instruction set designs
  - ▣ Ex: Stack vs GPR (system 360)
- ☐ Design decisions must take into account:
  - ▣ Technology
  - ▣ Machine organization
  - ▣ Programming languages
  - ▣ Compiler technology
  - ▣ Operating systems

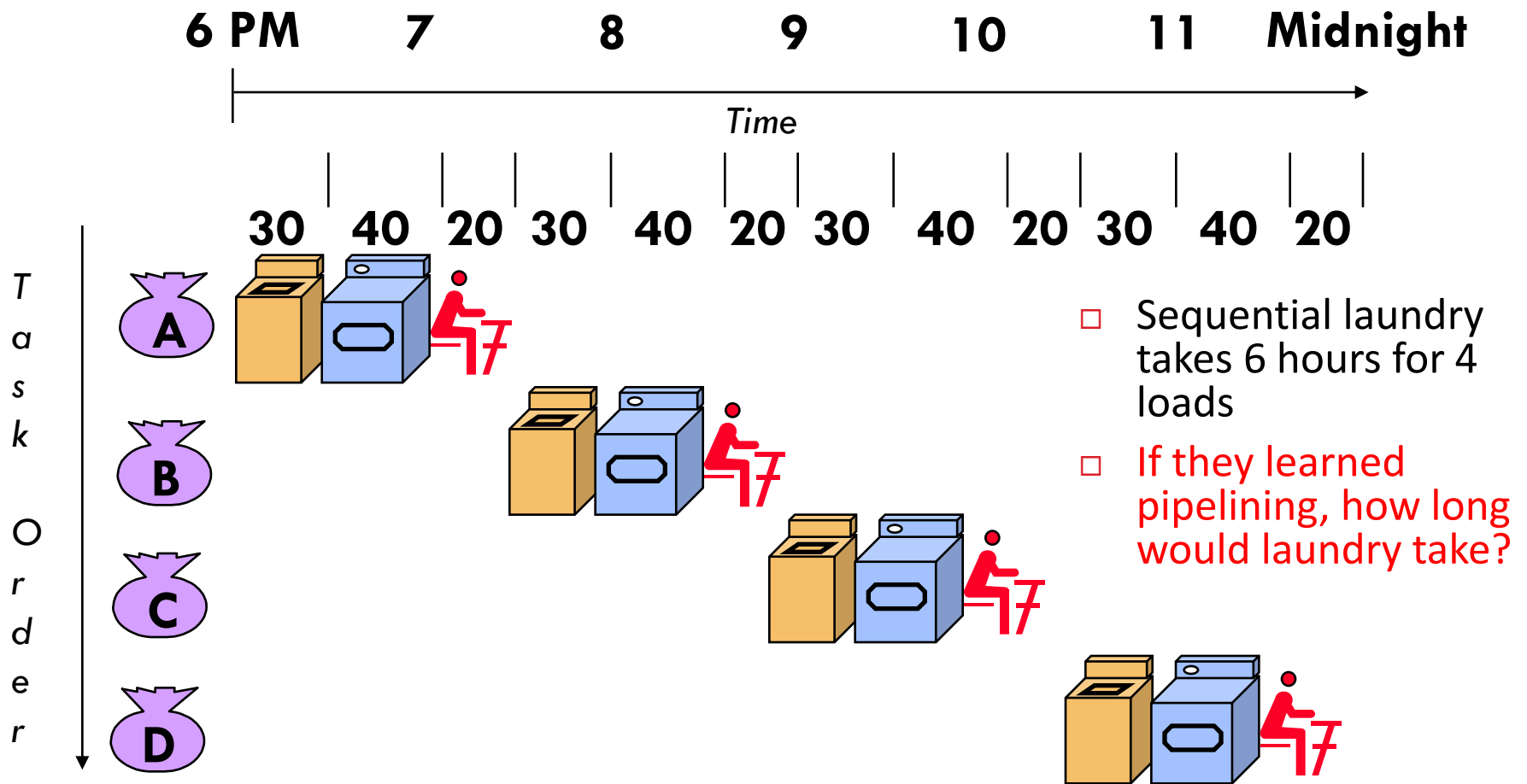　　　　　　　　　　　　　　　　　　　　Dr Noor Mahammad Sk

# 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

# Sequential Laundry

6 PM     7     8     9     10     11     Midnight

*Time*

30  40  20  30  40  20  30  40  20  30  40  20

*Task Order*

A

B

C

D

□ Sequential laundry takes 6 hours for 4 loads

□ If they learned pipelining, how long would laundry take?

Dr Noor Mahammad Sk

# Pipelined Laundry

- Pipelined laundry takes 3.5 hours for 4 loads

# 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

Dr Noor Mahammad Sk

# Computer Pipelines

□ Execute billions of instructions, so *throughput* is what matters

□ What is desirable in instruction sets for pipelining?

◘ Variable length instructions vs.
all instructions same length?

◘ Memory operands part of any operation vs. memory operands only in loads or stores?

◘ Register operand many places in instruction format vs. registers located in same place?

# A "Typical" RISC
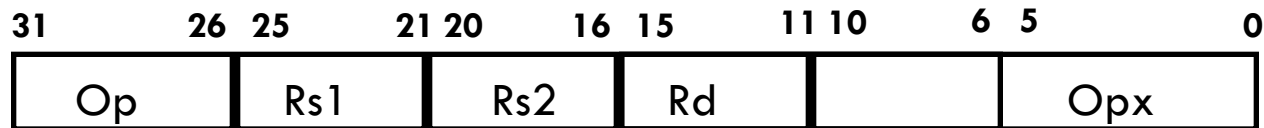
- 32-bit fixed format instruction (3 formats)
- Memory access only via load/store instructions
- 32 32-bit GPR (R0 contains zero)
- 3-address, reg-reg arithmetic instruction; registers in same place
- Single address mode for load/store:
  base + displacement
  - no indirection
- Simple branch conditions
- Delayed branch

Dr Noor Mahammad Sk

# Example: MIPS (Note register location)

**Register-Register**

| 31 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|----|----|----|----|----|----|----|
| Op | Rs1 | Rs2 | Rd | | Opx | |

**Register-Immediate**

| 31 | 26 25 | 21 20 | 16 15 | 0 |
|----|----|----|----|----|
| Op | Rs1 | Rd | immediate | |

**Branch**

| 31 | 26 25 | 21 20 | 16 15 | 0 |
|----|----|----|----|----|
| Op | Rs1 | Rs2/Opx | immediate | |

**Jump / Call**

| 31 | 26 25 | 0 |
|----|----|----|
| Op | target | |

# 5 Steps of MIPS Datapath

| Instruction Fetch | Instr. Decode Reg. Fetch | Execute Addr. Calc | Memory Access | Write Back |
|---|---|---|---|---|

Next PC

Adder

4

Next SEQ PC

Address

Memory

Inst

RS1

RS2

RD

Reg File

Zero?

MUX

MUX

ALU

Sign Extend

Imm

Data Memory

L M D

MUX

MUX

WB Data

Dr Noor Mahammad Sk
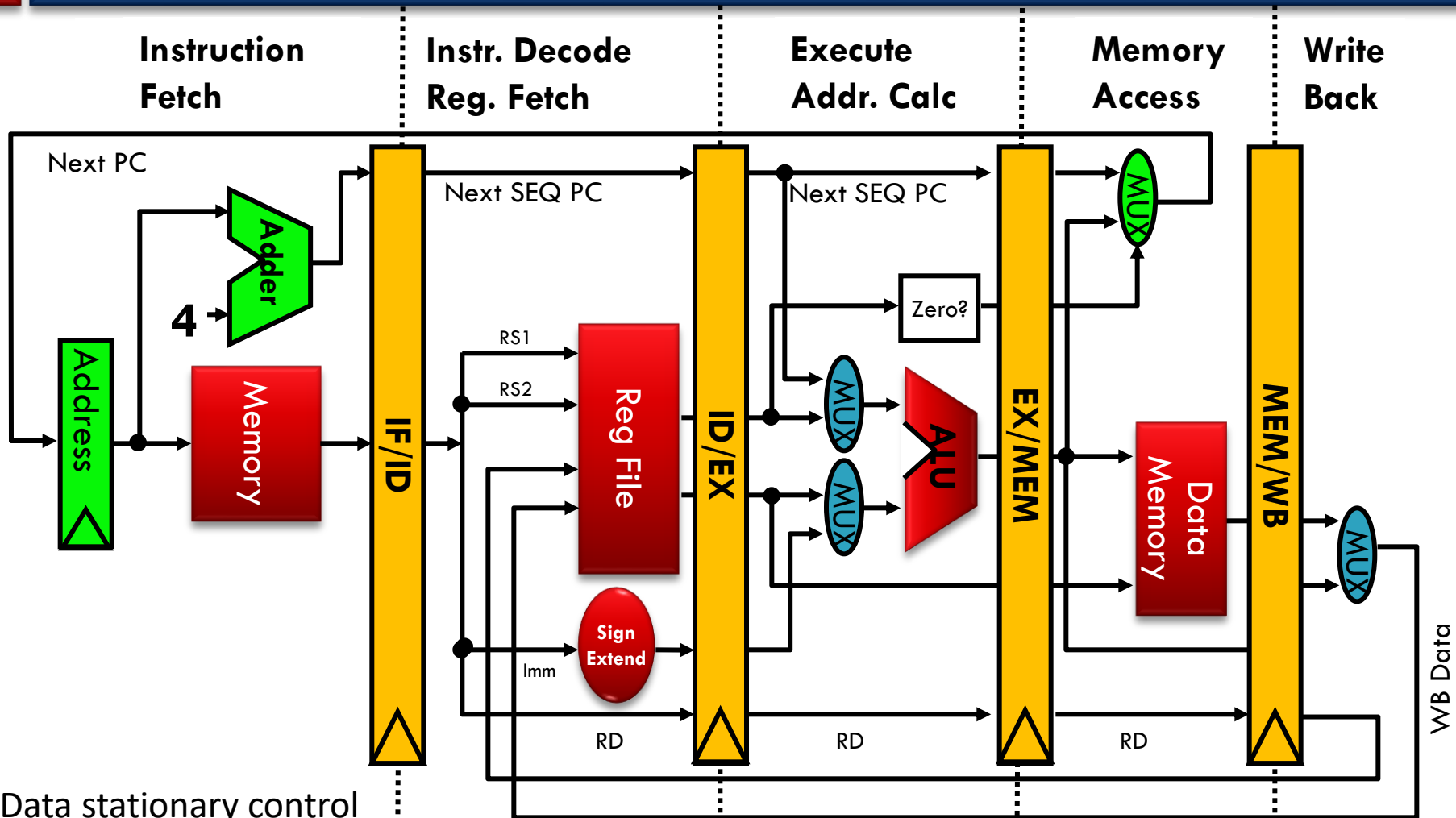
# 5 Steps of MIPS Datapath

□ Data stationary control

  ■ Local decode for each instruction phase/pipeline stage

12 March 2020

Dr Noor Mahammad Sk

# Visualizing Pipelining

# Its Not That Easy for Computers

□ Limits to pipelining: Hazards prevent next instruction from executing during its designated clock cycle

- ◘ <u>Structural hazards</u>: HW cannot support this combination of instructions (single person to fold and put clothes away)

- ◘ <u>Data hazards</u>: Instruction depends on result of prior instruction still in the pipeline (missing sock)

- ◘ <u>Control hazards</u>: Caused by delay between the fetching of instructions and decisions about changes in control flow (branches and jumps).

# One Memory Port/Structural Hazards

19

*Time (clock cycles)*

*Instr. Order*

| | Cycle 1 | Cycle 2 | Cycle 3 | Cycle 4 | Cycle 5 | Cycle 6 | Cycle 7 |
|---|---|---|---|---|---|---|---|
| **Load** | IFetch | Reg | ALU | DMem | Reg | | |
| **Instr 1** | | IFetch | Reg | ALU | DMem | Reg | |
| **Instr 2** | | | IFetch | Reg | ALU | DMem | Reg |
| **Instr 3** | | | | IFetch | Reg | ALU | DMem | Reg |

12 March 2020

Dr Noor Mahammad Sk

# One Memory Port/Structural Hazards

Dr Noor Mahammad Sk

# Data Hazard on R1

**Time (clock cycles)**



*Instr. Order*

add **r1**,r2,r3

sub r4,**r1**,r3

and r6,**r1**,r7

or   r8,**r1**,r9

xor r10,**r1**,r11

Dr Noor Mahammad Sk

# Three Generic Data Hazards

□ Read After Write (RAW)
Instr$_J$ tries to read operand before Instr$_I$ 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.

# Three Generic Data Hazards

□ Write After Read (WAR)
Instr$_J$ writes operand *before* Instr$_I$ 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 MIPS 5 stage pipeline because:

   ◘ All instructions take 5 stages, and

   ◘ Reads are always in stage 2, and

   ◘ Writes are always in stage 5

Dr Noor Mahammad Sk

# Three Generic Data Hazards

- Write After Write (WAW)
  Instr$_J$ writes operand *before* Instr$_I$ 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 MIPS 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

Dr Noor Mahammad Sk

# Forwarding to Avoid Data Hazard

**Time (clock cycles)**

*I n s t r.*

*O r d e r*



add **r1**,r2,r3

sub r4,**r1**,r3

and r6,**r1**,r7

or    r8,**r1**,r9

xor r10,**r1**,r11

Dr Noor Mahammad Sk

# HW Change for Forwarding

Dr Noor Mahammad Sk

# Data Hazard Even with Forwarding

*Time (clock cycles)*

*I n s t r.  O r d e r*

lw r1, 0(r2)

sub r4,r1,r6

and r6,r1,r7

or   r8,r1,r9

Dr Noor Mahammad Sk

# Data Hazard Even with Forwarding

*Time (clock cycles)*

*I
n
s
t
r.

O
r
d
e
r*



lw r1, 0(r2)

sub r4,r1,r6

and r6,r1,r7

or   r8,r1,r9

Dr Noor Mahammad Sk

# 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 |

Dr Noor Mahammad Sk

# Control Hazard on Branches Three Stage Stall

```
10: beq r1,r3,36

14: and r2,r3,r5

18: or  r6,r1,r7

22: add r8,r1,r9

36: xor r10,r1,r11
```

# Example: Branch Stall Impact

- If 30% branch, Stall 3 cycles significant
- Two part solution:
  - Determine branch taken or not sooner, AND
  - Compute taken branch address earlier
- MIPS branch tests if register = 0 or $\neq$ 0
- MIPS 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

Dr Noor Mahammad Sk

# Pipelined MIPS Datapath

**Instruction Fetch**

**Instr. Decode Reg. Fetch**

**Execute Addr. Calc**

**Memory Access**

**Write Back**



Next PC

Next SEQ PC

4

Address

Memory

Adder

Adder

MUX

Zero

RS1

RS2

Reg File

Sign Extend

Imm

IF/ID

ID/EX

MUX

ALU

EX/MEM

Data Memory

MEM/WB

MUX

RD

RD

RD

WB Data

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

12 March 2020

Dr Noor Mahammad Sk

# 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% MIPS branches not taken on average
- PC+4 already calculated, so use it to get next instruction

#3: Predict Branch Taken

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

Dr Noor Mahammad Sk

# Four Branch Hazard Alternatives

## #4: Delayed Branch

- Define branch to take place AFTER a following instruction

```
branch instruction
    sequential successor₁
    sequential successor₂
    ........
    sequential successorₙ
branch target if taken
```

**Branch delay of length $n$**

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

Dr Noor Mahammad Sk

# 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
  - Canceling 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)

Dr Noor Mahammad Sk

# REVIEW OF PERFORMANCE

12 March 2020

# Definitions

- Performance is in units of things per sec
  - Bigger is better
- If we are primarily concerned with response time
- Performance(x) = 1/{excution_time(x)}
- x is n times faster than y means:
- n = Performance(x)/Performance(y)

        = Execution_time(y)/Execution_time(x)

Dr Noor Mahammad Sk

# Aspects of CPU Performance (CPU Law)

| CPU time | = $\dfrac{\text{Seconds}}{\text{Program}}$ | = $\dfrac{\text{Instructions}}{\text{Program}}$ x | $\dfrac{\text{Cycles}}{\text{Instruction}}$ x | $\dfrac{\text{Seconds}}{\text{Cycle}}$ |
|---|---|---|---|---|

|  | **Instruction Count** | **CPI** | **Clock Rate** |
|---|:---:|:---:|:---:|
| Program | X |  |  |
| Compiler | X | (X) |  |
| Instruction Set | X | X |  |
| Organization |  | X | X |
| Technology |  |  | X |

Dr Noor Mahammad Sk

# Cycles Per Instruction (Throughput)

- "Average Cycles per Instruction"
- CPI = (CPU Time * Clock Rate) / Instruction Count
  - = Cycles / Instruction Count

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

- "Instruction Frequency"

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

Dr Noor Mahammad Sk

# Example: Calculating CPI

□ **Base Machine (Reg / Reg)**

| Op | Freq | Cycles | CPI(i) | (% Time) |
|---|---|---|---|---|
| **ALU** | **50%** | **1** | **.5** | **(33%)** |
| **Load** | **20%** | **2** | **.4** | **(27%)** |
| **Store** | **10%** | **2** | **.2** | **(13%)** |
| **Branch** | **20%** | **2** | **.4** | **(27%)** |
| | | | **1.5** | |

**Typical Mix of instruction types in program**

# Example: Branch Stall Impact

□ Assume CPI = 1.0 ignoring branches

□ Assume solution was stalling for 3 cycles

□ If 30% branch, Stall 3 cycles

| □ | Op | Freq | Cycles | CPI(i) | (% Time) |
|---|------|------|--------|--------|----------|
| □ | Other | 70% | 1 | .7 | (37%) |
| □ | Branch | 30% | 4 | 1.2 | (63%) |

□ => new CPI = 1.9, or almost 2 times slower

Dr Noor Mahammad Sk

# Example 2: Speed Up Equation for Pipelining

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

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

**For simple RISC pipeline, CPI = 1:**

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

Dr Noor Mahammad Sk

# Example 3: Evaluating Branch Alternatives (for 1 program)

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

| Scheduling scheme | Branch penalty | CPI | speedup v. stall |
|---|---|---|---|
| Stall pipeline | 3 | 1.42 | 1.0 |
| Predict taken | 1 | 1.14 | 1.26 |
| Predict not taken | 1 | 1.09 | 1.29 |
| Delayed branch | 0.5 | 1.07 | 1.31 |

Conditional & Unconditional = 14%, 65% change PC

# Example 4: 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

$$\text{SpeedUp}_A = \text{Pipeline Depth}/(1 + 0) \times (\text{clock}_{unpipe}/\text{clock}_{pipe})$$
$$= \text{Pipeline Depth}$$

$$\text{SpeedUp}_B = \text{Pipeline Depth}/(1 + 0.4 \times 1) \times (\text{clock}_{unpipe}/(\text{clock}_{unpipe} / 1.05)$$
$$= (\text{Pipeline Depth}/1.4) \times 1.05$$
$$= 0.75 \times \text{Pipeline Depth}$$

$$\text{SpeedUp}_A / \text{SpeedUp}_B = \text{Pipeline Depth}/(0.75 \times \text{Pipeline Depth}) = 1.33$$

- Machine A is 1.33 times faster

Dr Noor Mahammad Sk

# NOW, REVIEW OF MEMORY HIERARCHY

12 March 2020

Dr Noor Mahammad Sk

# Processor-DRAM Memory Gap (latency)

μProc
60%/yr.
(2X/1.5yr)

"Moore's Law"

**Processor-Memory Performance Gap: (grows 50% / year)**

DRAM
9%/yr.
(2X/10 yrs)

**Performance**

1000
100
10
1

1980 1981 1982 1983 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 2000

**Time**

Dr Noor Mahammad Sk

# Levels of the Memory Hierarchy

*Capacity Access Time Cost*

**Staging Xfer Unit**

**Upper Level**

faster

**CPU Registers**
**100s Bytes**
**<1s ns**

Registers

Instr. Operands

prog./compiler
1-8 bytes

*Cache*
**10s-100s K Bytes**
**1-10 ns**
**$10/ MByte**

Cache

Blocks

cache cntl
8-128 bytes

*Main Memory*
**M Bytes**
**100ns- 300ns**
**$1/ MByte**

Memory

Pages

OS
512-4K bytes

*Disk*
**10s G Bytes, 10 ms**
**(10,000,000 ns)**
**$0.0031/ MByte**

Disk

Files

user/operator
Mbytes

Larger

*Tape*
**infinite**
**sec-min**
**$0.0014/ MByte**

Tape

**Lower Level**

12 March 2020

Dr Noor Mahammad Sk

# 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 (hardware) relied on locality for speed

# 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!)



To Processor

From Processor

**Upper Level Memory**

Blk X

**Lower Level Memory**

Blk Y

Dr Noor Mahammad Sk

# 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)

Dr Noor Mahammad Sk

# Simplest Cache: Direct Mapped

**Memory Address**     **Memory**

| | |
|---|---|
| **0** | |
| **1** | |
| **2** | |
| **3** | |
| **4** | |
| **5** | |
| **6** | |
| **7** | |
| **8** | |
| **9** | |
| **A** | |
| **B** | |
| **C** | |
| **D** | |
| **E** | |
| **F** | |

**4 Byte Direct Mapped Cache**

**Cache Index**

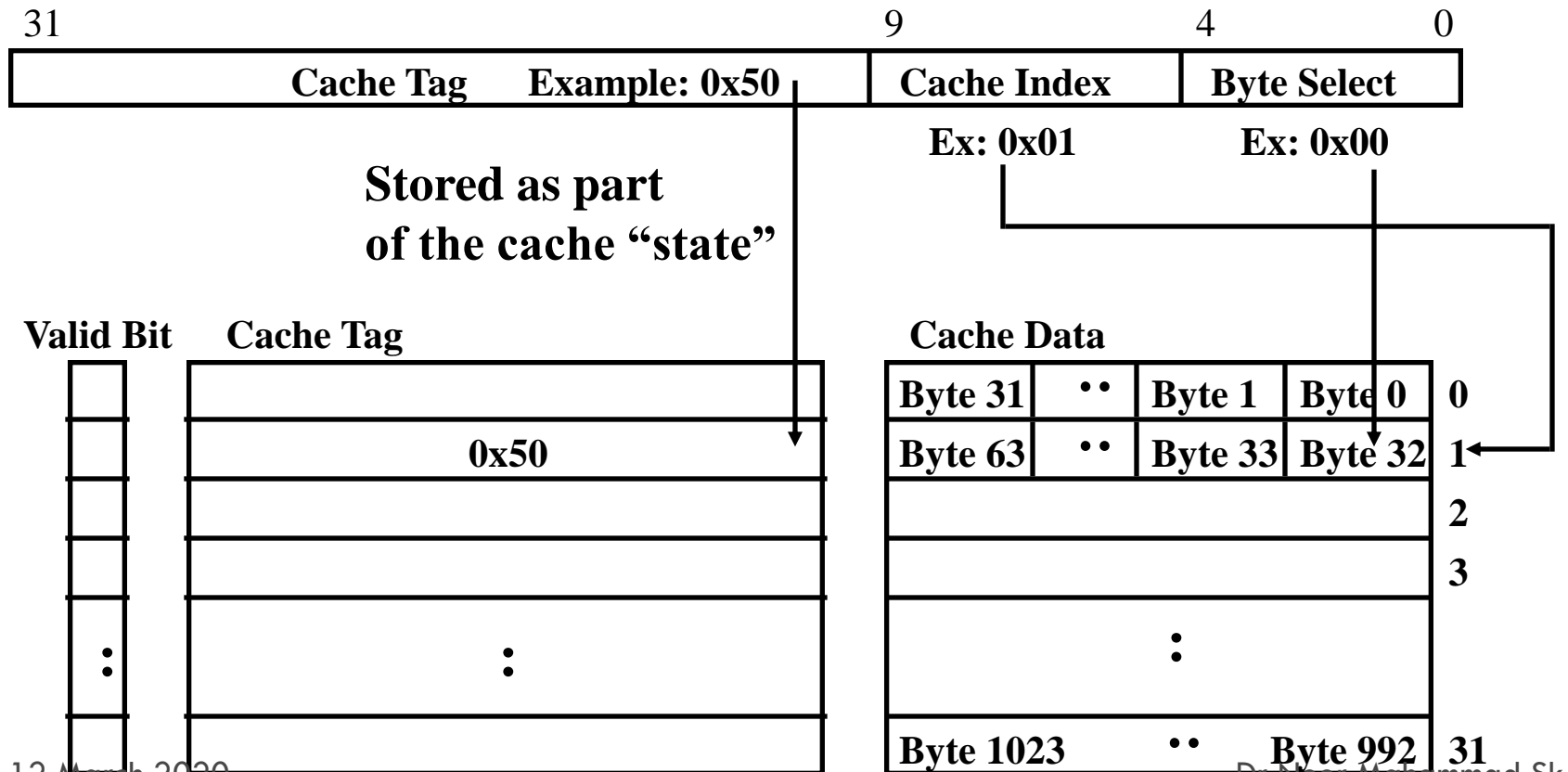| |
|---|
| **0** |
| **1** |
| **2** |
| **3** |

□ 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?

12 March 2020                                     Dr Noor Mahammad Sk

# 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)

| 31 | | 9 | | 4 | | 0 |
|---|---|---|---|---|---|---|
| **Cache Tag** | **Example: 0x50** | | **Cache Index** | | **Byte Select** | |

**Ex: 0x01**          **Ex: 0x00**

**Stored as part
of the cache "state"**

**Valid Bit**          **Cache Tag**          **Cache Data**

| | | Byte 31 | •• | Byte 1 | Byte 0 | **0** |
|---|---|---|---|---|---|---|
| | 0x50 | Byte 63 | •• | Byte 33 | Byte 32 | **1** |
| | | | | | | **2** |
| | | | | | | **3** |
| ⋮ | ⋮ | | ⋮ | | | |
| | | Byte 1023 | •• | | Byte 992 | **31** |

# 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



**Cache Index**

| Valid | Cache Tag | Cache Data | | Cache Data | Cache Tag | Valid |
|---|---|---|---|---|---|---|
| | | Cache Block 0 | | Cache Block 0 | | |

**Adr Tag** **Compare** **Sel1** 1 **Mux** 0 **Sel0** **Compare** **Adr Tag**

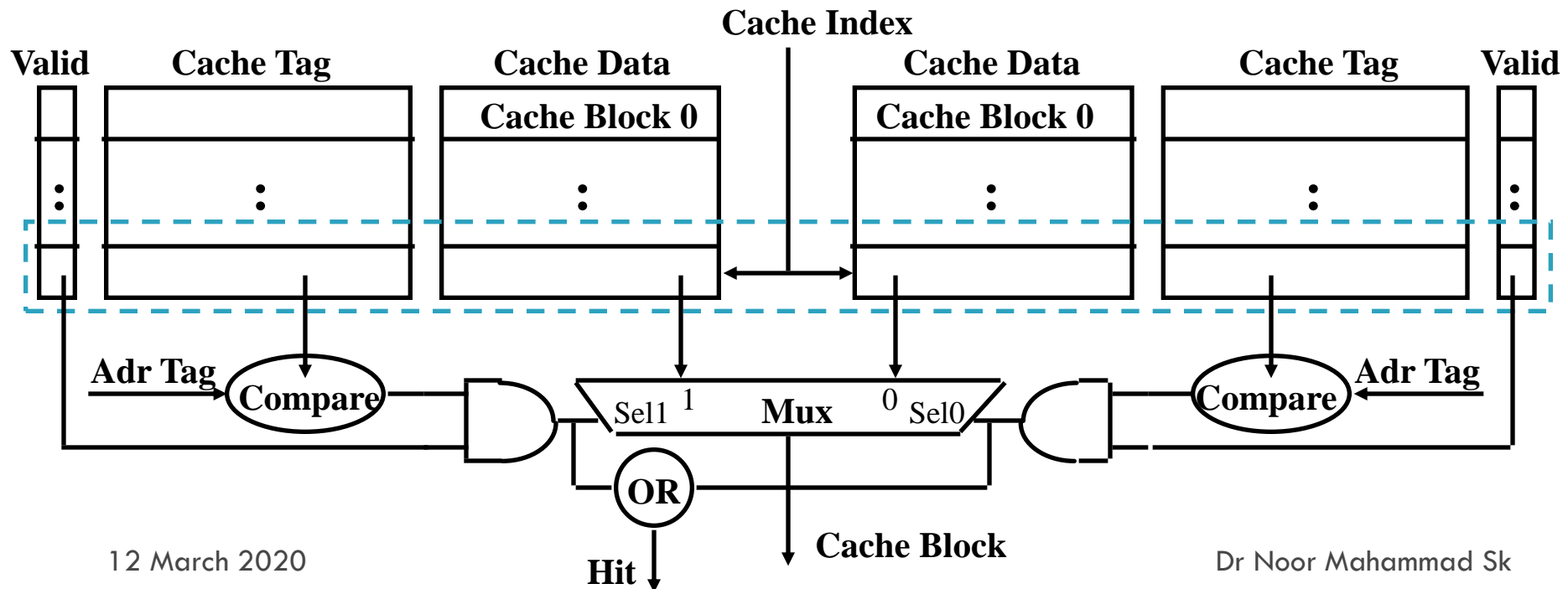**OR**

**Cache Block**

**Hit**

12 March 2020

Dr Noor Mahammad Sk

# 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.



12 March 2020

Dr Noor Mahammad Sk
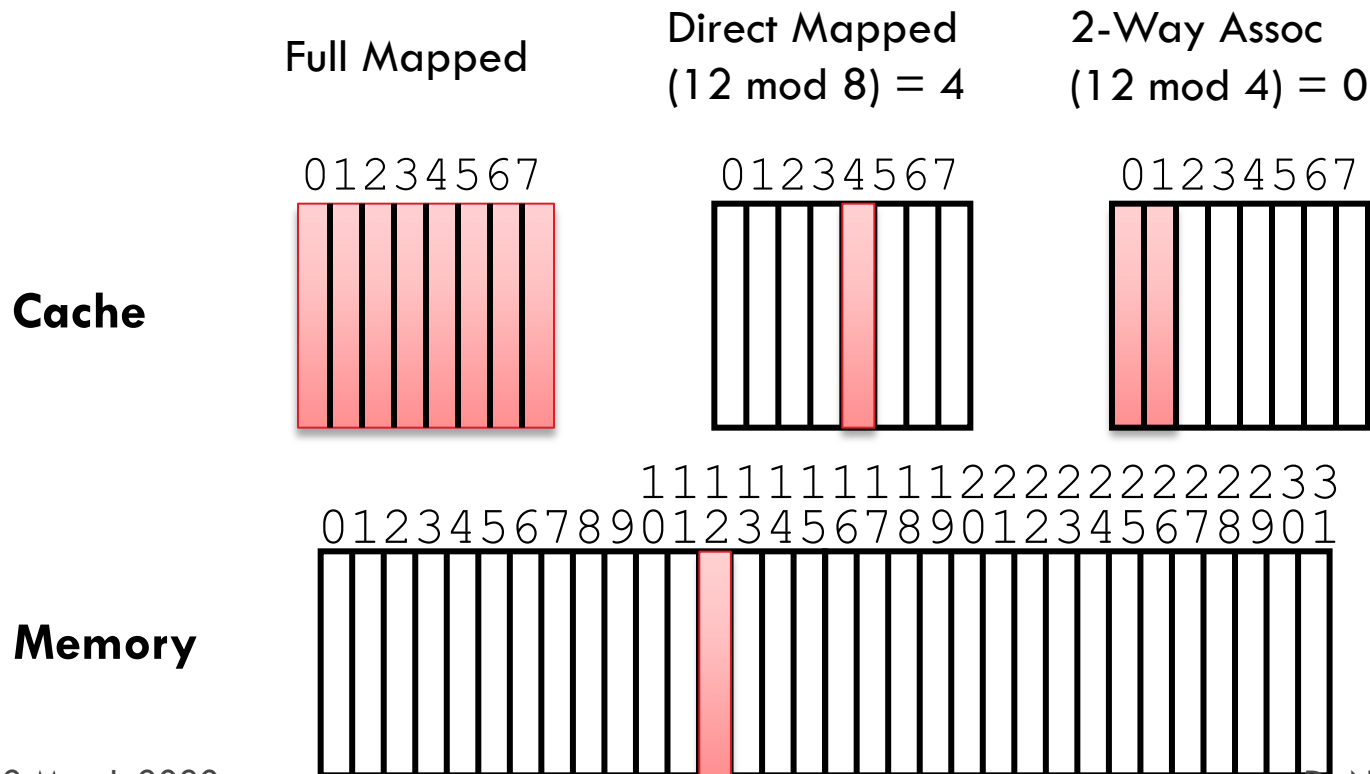
# 4 Questions for Memory Hierarchy

- 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)*

Dr Noor Mahammad Sk

# 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



Full Mapped

Direct Mapped
(12 mod 8) = 4

2-Way Assoc
(12 mod 4) = 0

Cache

Memory

12 March 2020

Dr Noor Mahammad Sk

# 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 | |

# Q3: Which block should be replaced on a miss?

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

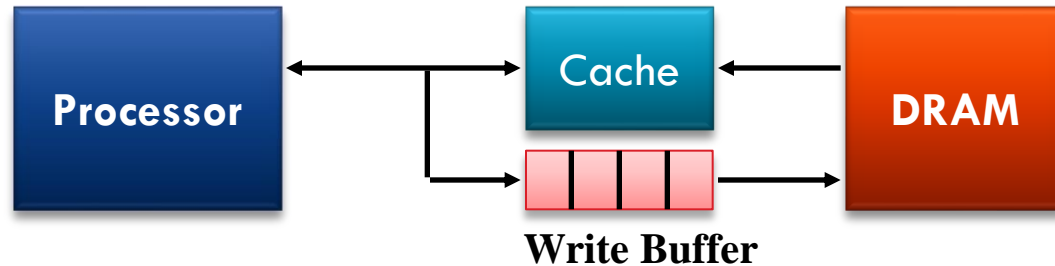| ASSOC | 2 – Way | | 4 – Way | | 8 – Way | |
|-------|-----|--------|-----|--------|-----|--------|
| Size | LRU | Random | LRU | Random | LRU | Random |
| 16 KB | 5.2% | 5.7% | 4.7% | 5.3% | 4.4% | 5.0% |
| 64 KB | 1.9% | 2.0% | 1.5% | 1.7% | 1.4% | 1.5% |
| 256 KB | 1.15% | 1.17% | 1.13% | 1.13% | 1.12% | 1.12% |

Dr Noor Mahammad Sk

# 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
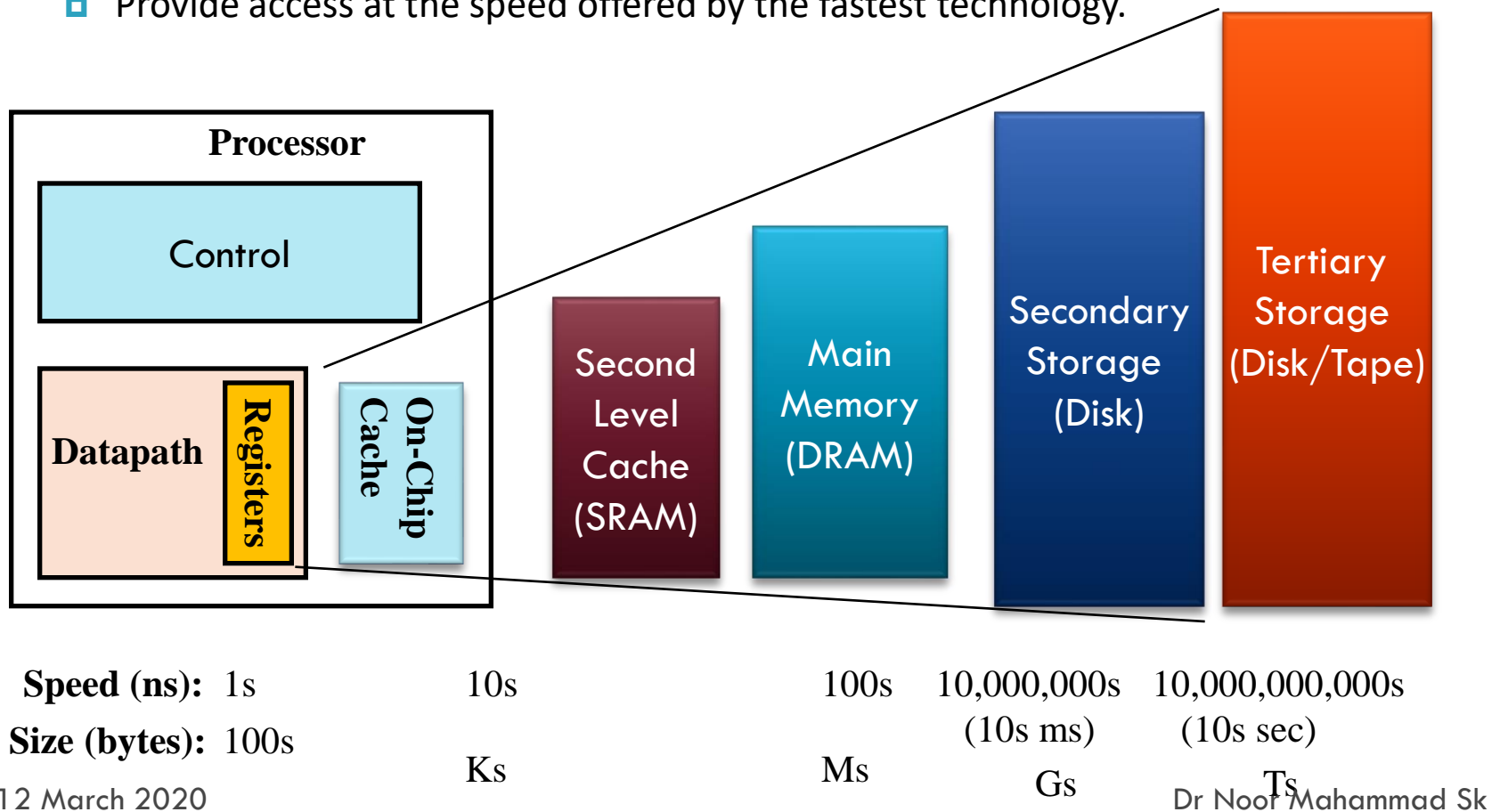
# Write Buffer for Write Through

**Write Buffer**

- 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) << 1 / DRAM write cycle
- Memory system designer's nightmare:
  - Store frequency (w.r.t. time)  >  1 / DRAM write cycle
  - Write buffer saturation

Dr Noor Mahammad Sk

# 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.

**Processor**

Control

Datapath | Registers | On-Chip Cache

Second Level Cache (SRAM)

Main Memory (DRAM)

Secondary Storage (Disk)

Tertiary Storage (Disk/Tape)

| **Speed (ns):** | 1s | 10s | 100s | 10,000,000s (10s ms) | 10,000,000,000s (10s sec) |
|---|---|---|---|---|---|
| **Size (bytes):** | 100s | Ks | Ms | Gs | Ts |

# Summary #1/4: Pipelining & Performance

☐ Just overlap tasks; easy if tasks are independent

☐ Speed Up ≤ 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

| CPU time | = Seconds Program | = Instructions Program | x | Cycles Instruction | x | Seconds Cycle |
|----------|-------------------|------------------------|---|--------------------|---|---------------|

# Summary #2/4: 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.
- ☐ Write Policy:
    - ◰ Write Through: needs a write buffer.
    - ◰ Write Back: control can be complex
- ☐ Today CPU time is a function  of (ops, cache misses) vs. just f(ops): What does this mean to Compilers, Data structures, Algorithms?
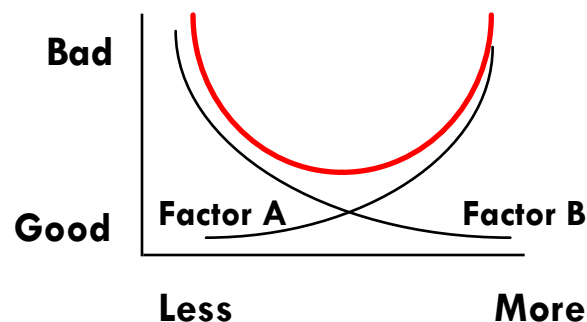
Dr Noor Mahammad Sk

# Summary #3/4:
# The Cache Design Space
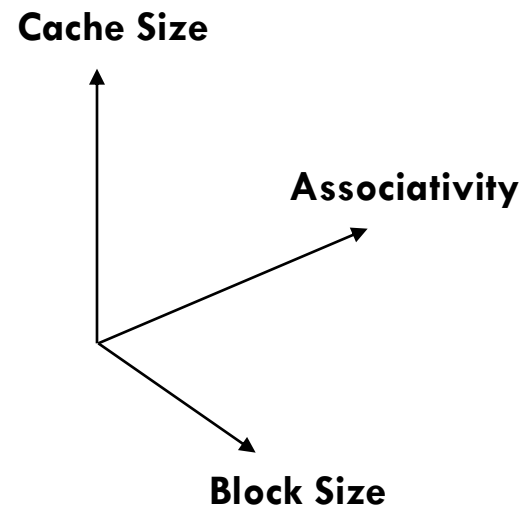
- Several interacting dimensions
  - cache size
  - block size
  - associativity
  - replacement policy
  - write-through vs write-back

**Cache Size**

**Associativity**

**Block Size**

- The optimal choice is a compromise
  - depends on access characteristics
    - workload
    - use (I-cache, D-cache, TLB)
  - depends on technology / cost
- Simplicity often wins

**Bad**

**Good**    Factor A         Factor B

**Less**                    **More**

12 March 2020                                    Dr Noor Mahammad Sk

# Review #4/4: 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 repalced on miss? 4) How are writes handled?

□ Page tables map virtual address to physical address

□ TLBs make virtual memory practical

   ■ Locality in data => locality in addresses of data,
     temporal and spatial

□ TLB misses are significant in processor performance

   ■ funny times, as most systems can't access all of 2nd level cache without TLB misses!

□ 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