

RISC Design

Memory System

Virendra Singh

Professor

Computer Architecture and Dependable Systems Lab

Department of Electrical Engineering

Indian Institute of Technology Bombay

<http://www.ee.iitb.ac.in/~viren/>

E-mail: viren@ee.iitb.ac.in

EE-739: Processor Design



Lecture 24 (17 March 2021)

CADSL

Control Hazard

- ① Stall.
- ② Branch delay slot
- ③ Branch Prediction — 98-99% accuracy

Track history → take the same decision

which we took last time. → bit

Branch History Table

PC	BTA	H _B
200	575	1

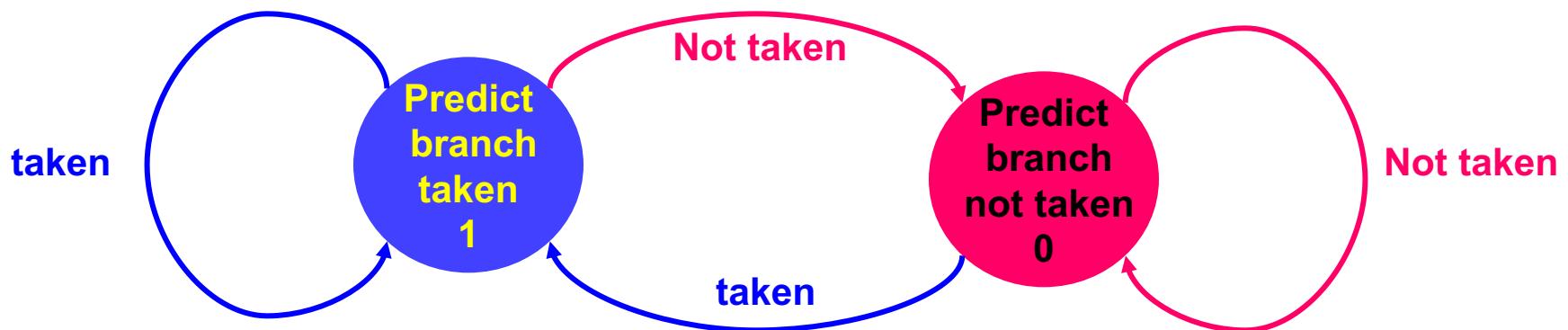
modified
from execution
stage

Fetch ✓

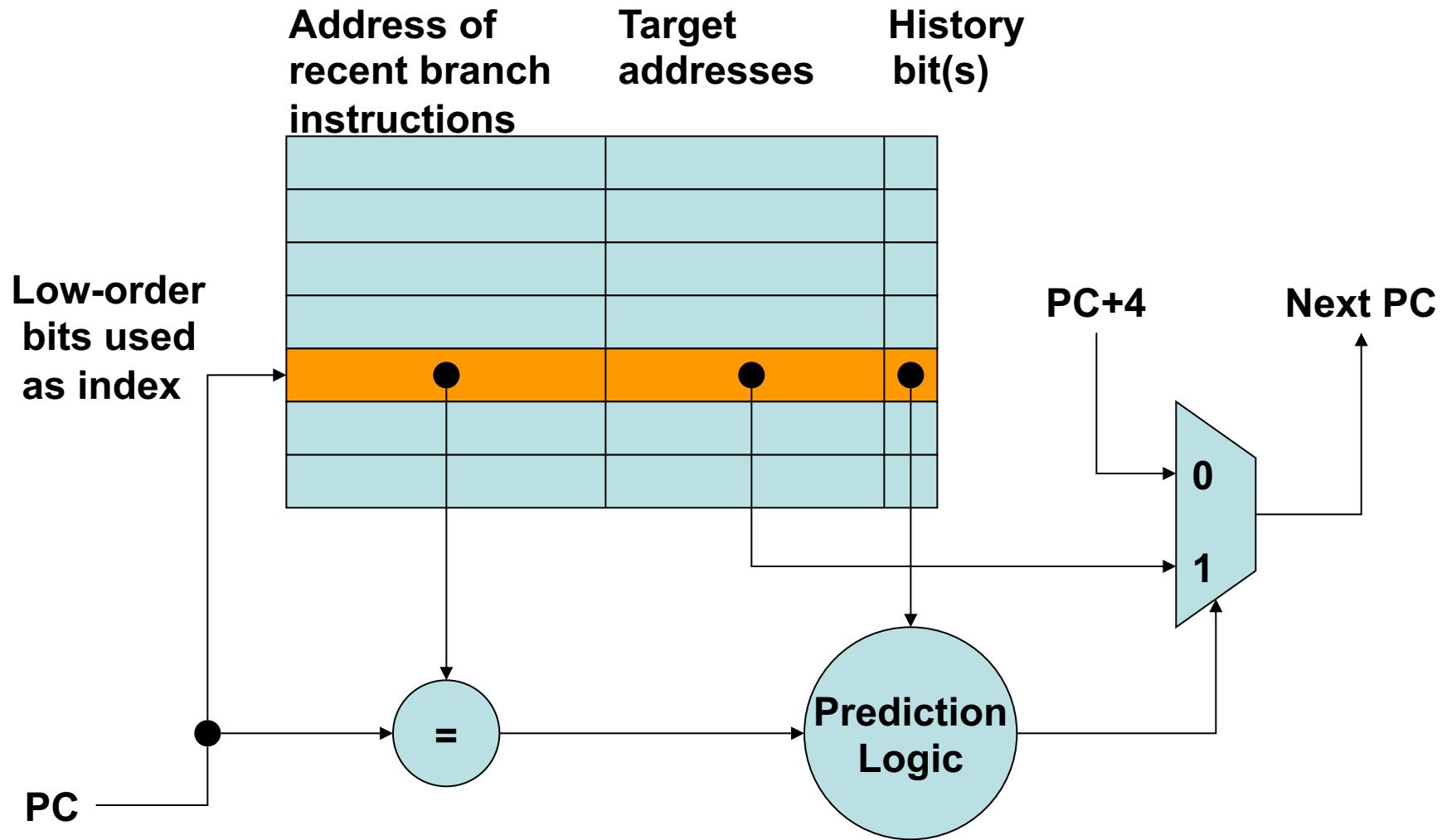


Branch Prediction

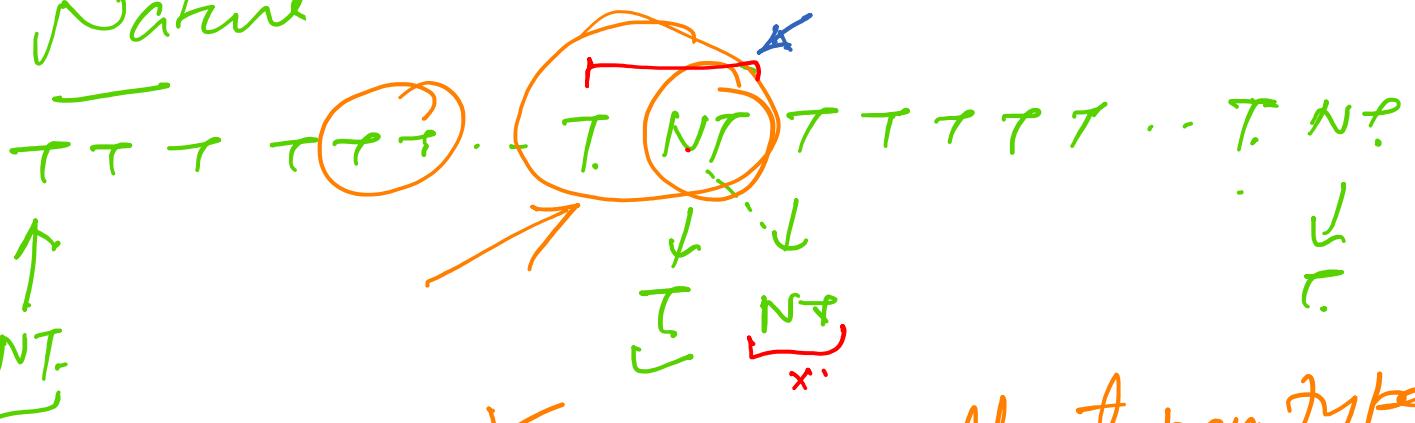
- Useful for program loops.
- A one-bit prediction scheme: a one-bit buffer carries a “history bit” that tells what happened on the last branch instruction
 - History bit = 1, branch was taken
 - History bit = 0, branch was not taken



Branch Prediction

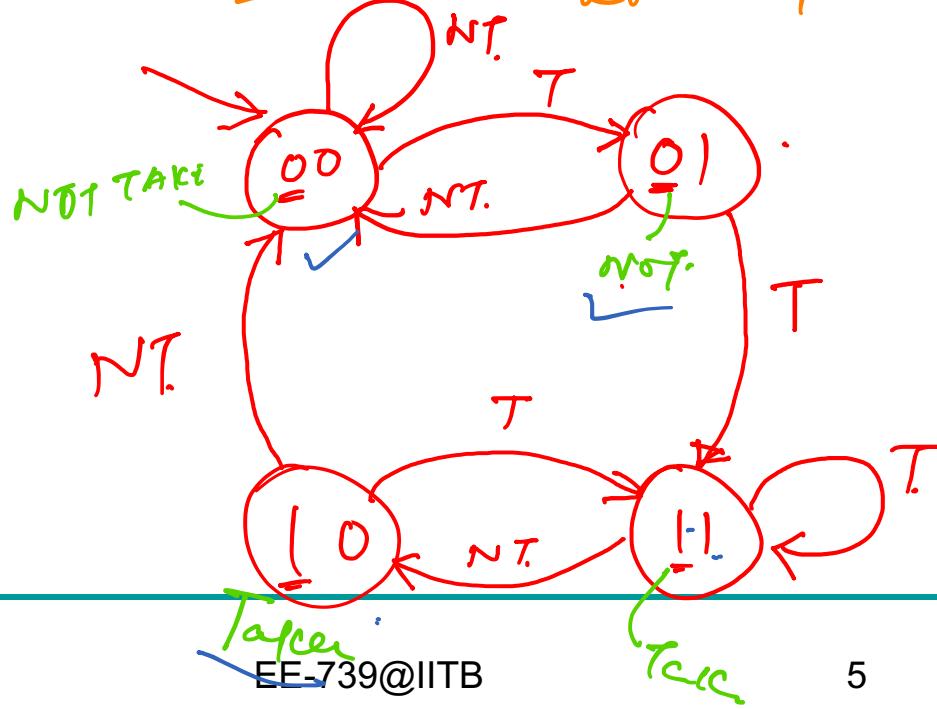


Nature



gbit.

0 0 -
0 1
- 0
—
↓ |
↑ |
decision



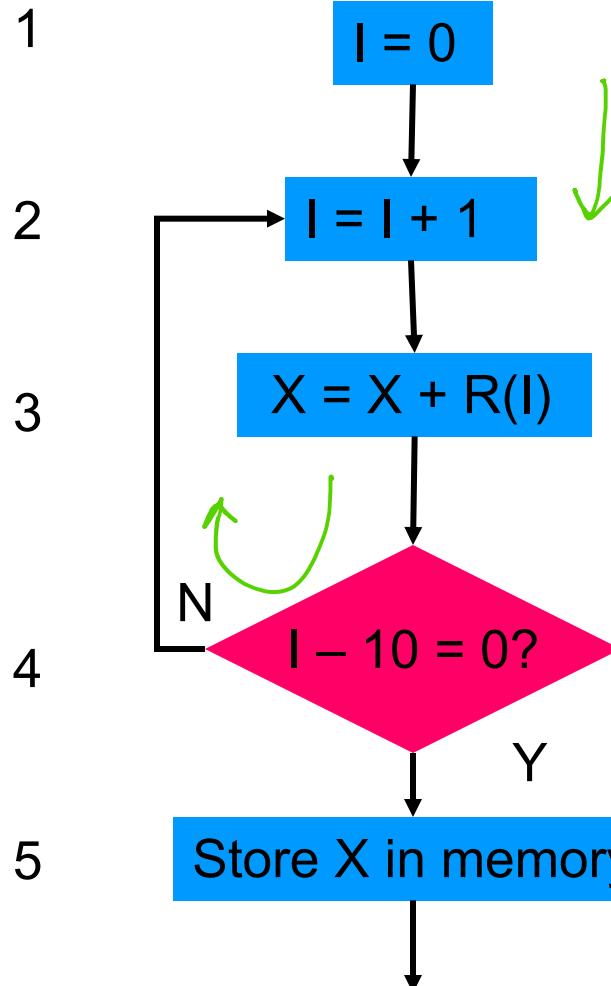
* normally taken type
in deeper history

$$\underline{>90\%} = \underline{\underline{NTNTNT}} \overline{\underline{\underline{T}}} \stackrel{\text{NTNT}}{\overline{\underline{\underline{T}}_2}}$$

8/10 प्र० 100%
2/80%.

Branch Prediction for a Loop

Execution of Instruction 4



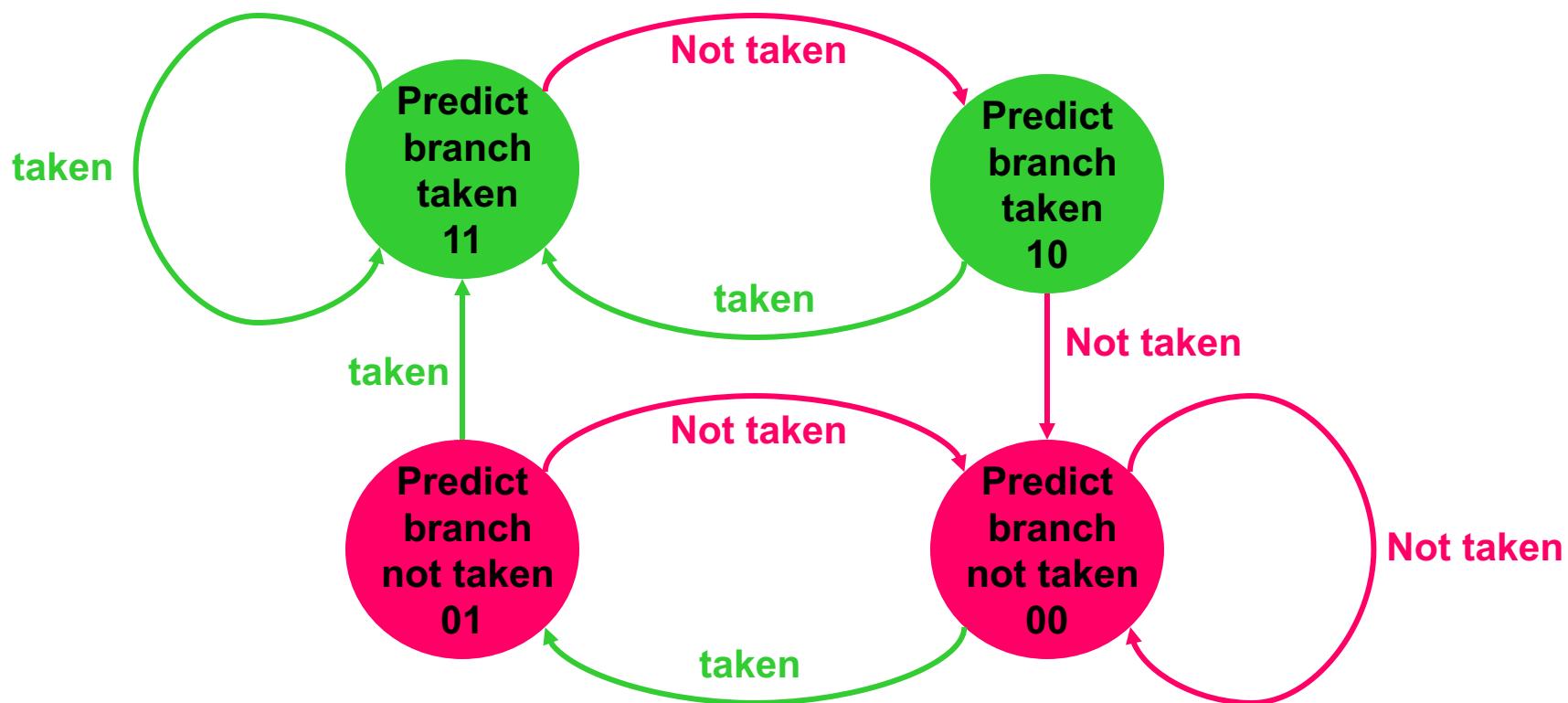
Execution seq.	Old hist. bit	Next instr.			New hist. bit	Prediction
		Pred.	I	Act.		
1	0	5	1	2	1	Bad
2	1	2	2	2	1	Good
3	1	2	3	2	1	Good
4	1	2	4	2	1	Good
5	1	2	5	2	1	Good
6	1	2	6	2	1	Good
7	1	2	7	2	1	Good
8	1	2	8	2	1	Good
9	1	2	9	2	1	Good
10	1	2	10	5	0	Bad

h.bit = 0 branch not taken, h.bit = 1 branch taken.



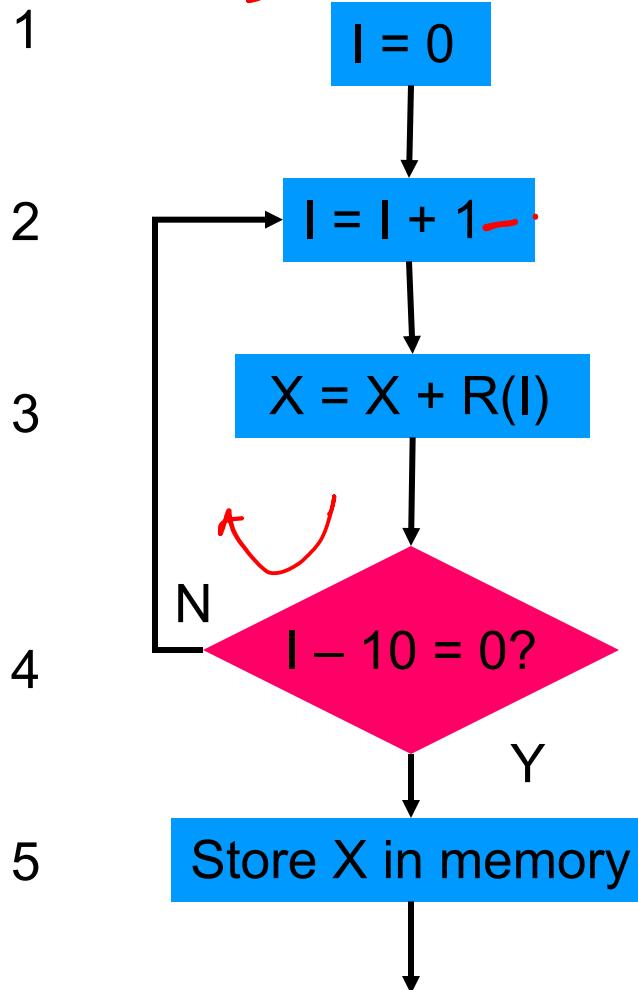
Two-Bit Prediction Buffer

- Can improve correct prediction statistics.



Branch Prediction for a Loop

(D = 90°)



Execution of Instruction 4

Execution seq.	Old Pred. Buf	Next instr.			New pred. Buf	Prediction
		Pred.	I	Act.		
1	10	2	1	2	11	Good
2	11	2	2	2	11	Good
3	11	2	3	2	11	Good
4	11	2	4	2	11	Good
5	11	2	5	2	11	Good
6	11	2	6	2	11	Good
7	11	2	7	2	11	Good
8	11	2	8	2	11	Good
9	11	2	9	2	11	Good
10	11	2	10	5	10	Bad



✓
 \rightarrow 1 bit 85%.
 \rightarrow 2 bit 90%.
 \rightarrow 3 bit \rightarrow 93%.
 \rightarrow 4 bit $\underline{94\%}$.

$\underline{\underline{90-92\%}}$
 accuracy

70% \rightarrow 85% \rightarrow 90%.



17 Mar 2021

EE-739@IITB

9 CADSL

Summary: Hazards

- Structural hazards
 - Cause: resource conflict
 - Remedies: (i) hardware resources, (ii) stall (bubble)
- Data hazards
 - Cause: data unavailability
 - Remedies: (i) forwarding, (ii) stall (bubble), (iii) code reordering
- Control hazards
 - Cause: out-of-sequence execution (branch or jump)
 - Remedies: (i) stall (bubble), (ii) branch prediction/pipeline flush, (iii) delayed branch/pipeline flush



Summary: Hazards

- Structural hazards
 - Cause: resource conflict
 - Remedies: (i) hardware resources, (ii) stall (bubble)
- Data hazards
 - Cause: data unavailability
 - Remedies: (i) forwarding, (ii) stall (bubble), (iii) code reordering
- Control hazards
 - Cause: out-of-sequence execution (branch or jump)
 - Remedies: (i) stall (bubble), (ii) branch prediction/pipeline flush, (iii) delayed branch/pipeline flush



Large memory
Increasing size of
perform & date

Fed - Decod - Enc. - Mem - WD
↓
200
200 ↗

① SRAM
very expensive

② DRAM.
(slow)

Memory

System

Single cycle memory access . →

200
Cycles)

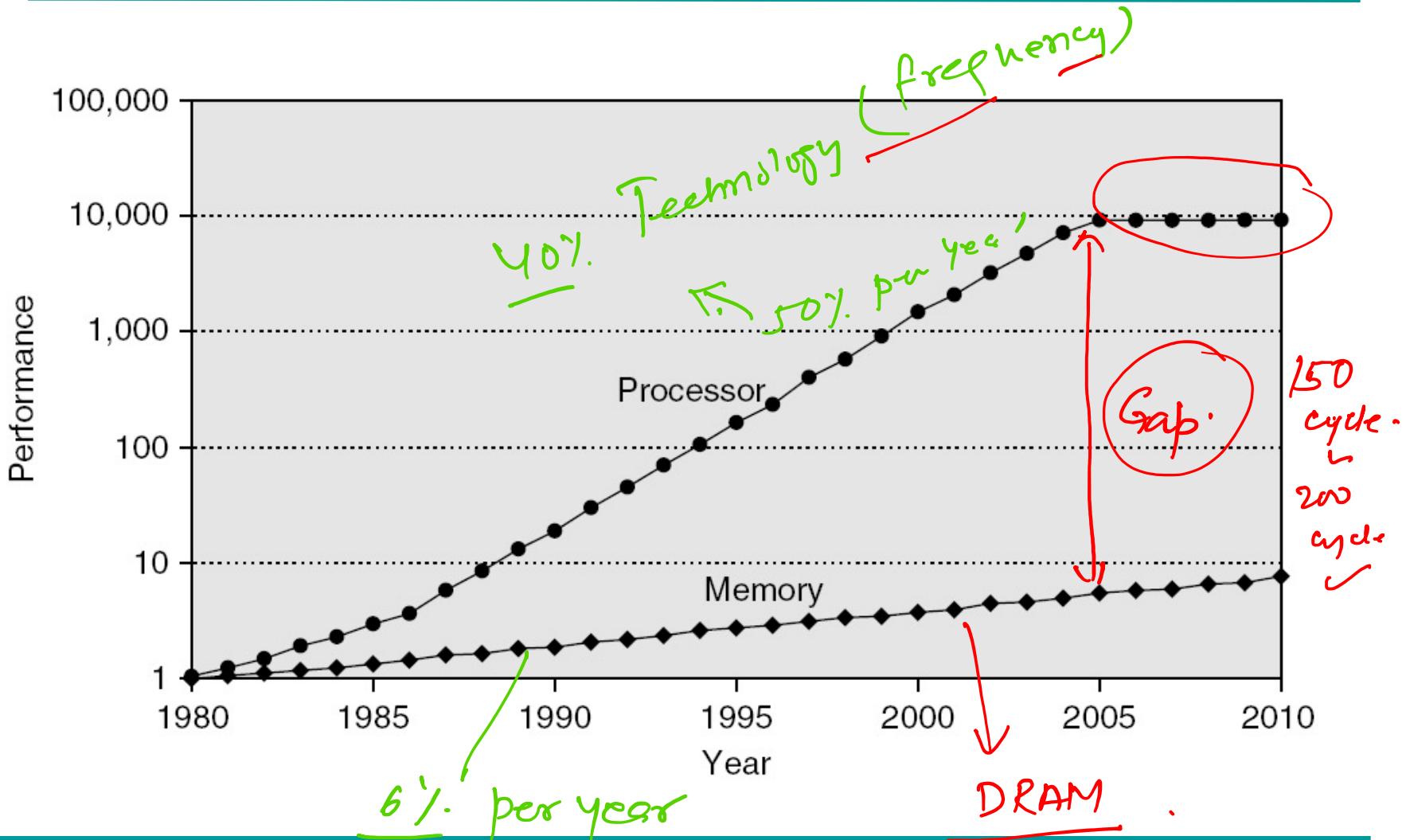


Big Picture

- Memory
 - Just an “ocean of bits”
 - Many technologies are available
- Key issues
 - Technology (how bits are stored)
 - Placement (where bits are stored)
 - Identification (finding the right bits)
 - Replacement (finding space for new bits)
 - Write policy (propagating changes to bits)
- Must answer these regardless of memory type



Memory Performance Gap



Types of Memory

Type	Size	Speed	Cost/bit
Register	< 1KB	< 1ns	\$\$\$\$
On-chip SRAM 	8KB-6MB	< 10ns	\$\$\$
Off-chip SRAM 	1Mb – 16Mb	< 20ns	\$\$
DRAM 	64MB – 1TB	< 100ns	\$
Disk	40GB – 1PB	< 20ms	~0



Memory System

- Programmers want unlimited amounts of memory with low latency
- Fast memory technology is more expensive per bit than slower memory
- **Solution:** organize memory system into a hierarchy
 - Entire addressable memory space available in largest, slowest memory
 - Incrementally smaller and faster memories, each containing a subset of the memory below it, proceed in steps up toward the processor
- Temporal and spatial locality insures that nearly all references can be found in smaller memories
 - Gives the allusion of a large, fast memory being presented to the processor



$\text{for } i=0; i<100; i++)$ temporal locality
 $\underbrace{a(i) = b(i) + c(i)}_{\uparrow} \leftarrow$ $\Rightarrow \begin{cases} e = b + c \\ d = a + e \\ p = s + r \end{cases}$
✓

$\text{for } i=0; i<100)$

instructions

spatial locality

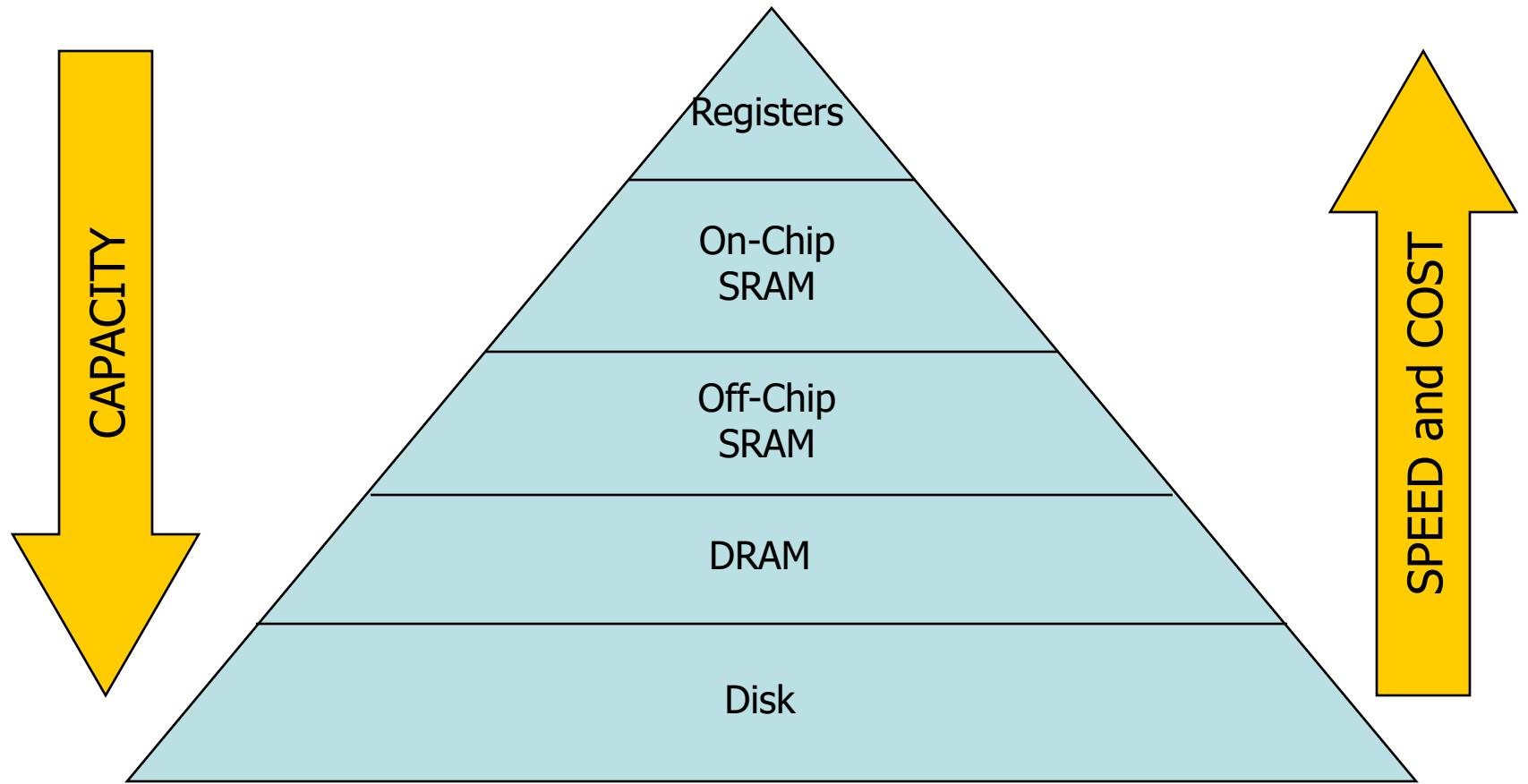
B1i) B1o) D1i)
B1- B12) P13)

$\text{for } i=0; i<100; i++)$
 $\left\{ \begin{array}{l} a(i) = b(i) + c(i) \\ d(i) = a(i) + e(i) \\ p(i) = g(i) * r(i) \end{array} \right.$
}

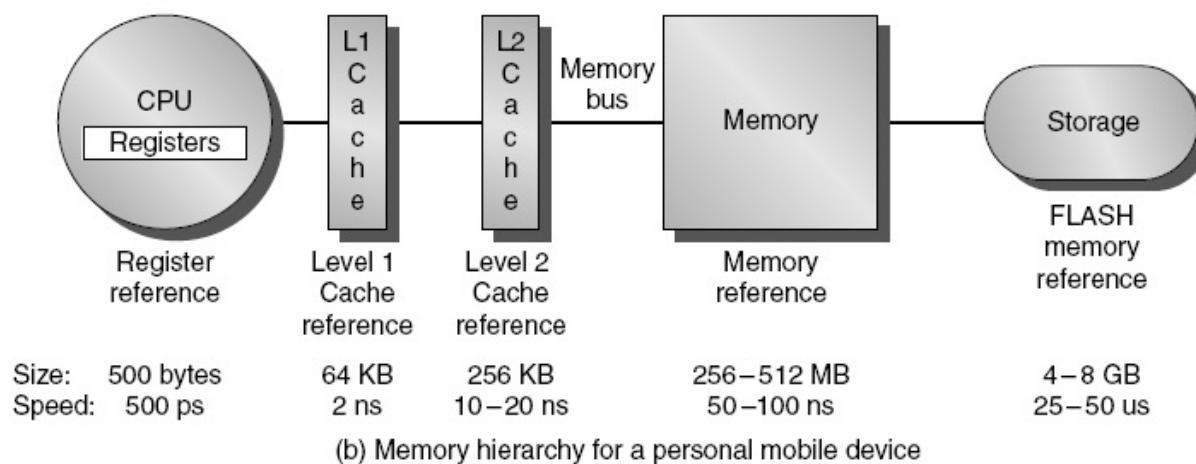
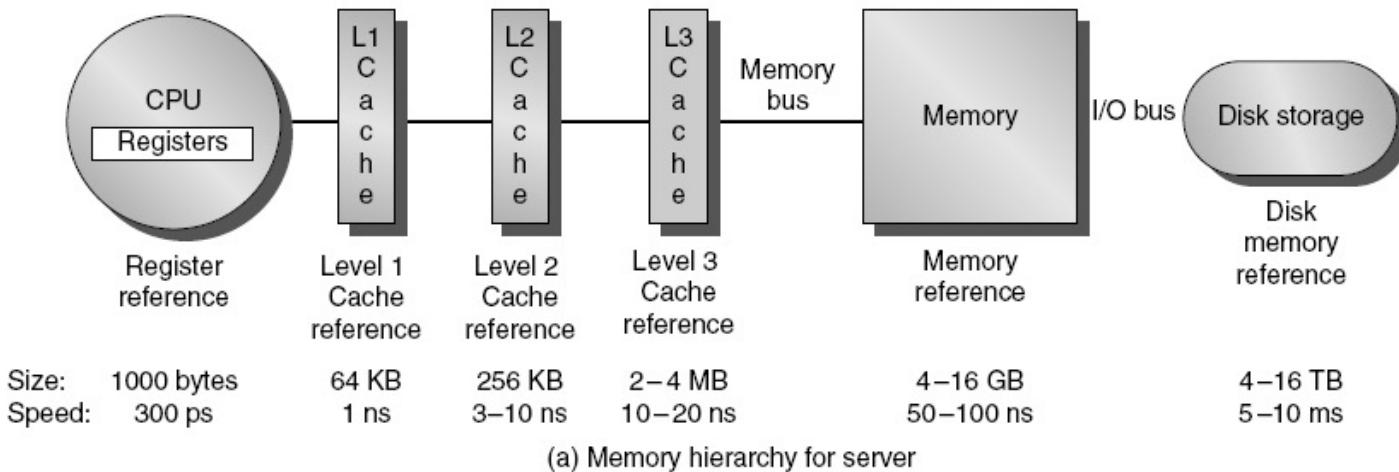
B10 12 12 12 1



Memory Hierarchy



Memory Hierarchy



Why Memory Hierarchy?

- Need lots of bandwidth

$$BW = \frac{1.0inst}{cycle} \times \left[\frac{1Ifetch}{inst} \times \frac{4B}{Ifetch} + \frac{0.2Dref}{inst} \times \frac{4B}{Dref} \right] \times \frac{1Gcycles}{sec}$$
$$= \frac{4.8GB}{sec}$$

$$\frac{\text{data}}{\text{Sec.}} > \frac{\text{Instr.}}{\text{Cycle}} \times \frac{\text{cycle}}{\text{Instr.}} \times \frac{\text{data}}{\text{Instr.}} \times \frac{\text{cycle}}{\text{Sec.}}$$

1 PC ($\frac{1}{CP2}$)

- Need lots of storage
 - 64MB (minimum) to multiple TB
- Must be cheap per bit
 - (TB x anything) is a lot of money!
- These requirements seem incompatible



$$\begin{aligned}
 \text{bw} = \frac{\text{data}}{\text{sec}} : & \quad \frac{\text{data}}{\text{fetch}} \times \frac{\text{fetch}}{\text{instr}} \times \frac{\text{inst}}{\text{cycle}} \times \frac{\text{cycle}}{\text{sec}} \\
 & \quad \boxed{3472} \\
 \text{inst} & \quad [4B + 3 \times 4B] \times 1 \times \textcircled{1} \times 3 \times 10^9 \\
 & \quad \text{data only fr } \underline{1\text{ sec / instr}} \\
 & \quad \times 30 \\
 & = \cancel{4} \times 4 \times 1.3 \times 3 \times 10^9 = 12 \times 1.3 \times 10^9 \\
 & = \underline{15.6 \times 10^9} \\
 & = \underline{15.6 \text{ GBps}} \quad \checkmark
 \end{aligned}$$



Memory Hierarchy Design

- Memory hierarchy design becomes more crucial with recent multi-core processors:
 - Aggregate peak bandwidth grows with # cores:
 - Intel Core i7 can generate two references per core per clock
 - Four cores and 3.2 GHz clock
 - 25.6 billion 64-bit data references/second +
 - 12.8 billion 128-bit instruction references
 - = 409.6 GB/s! ↛
 - DRAM bandwidth is only 6% of this (25 GB/s)
 - Requires:
 - Multi-port, pipelined caches
 - Two levels of cache per core
 - Shared third-level cache on chip

30-40
GB/s



Why Locality?

- Analogy:
 - Library (Disk) ✓
 - Bookshelf (Main memory)
 - Stack of books on desk (off-chip cache)
 - Opened book on desk (on-chip cache)
- Likelihood of:
 - Referring to same book or chapter again?
 - Probability decays over time
 - Book moves to bottom of stack, then bookshelf, then library
 - Referring to chapter n+1 if looking at chapter n?



Thank You

