

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

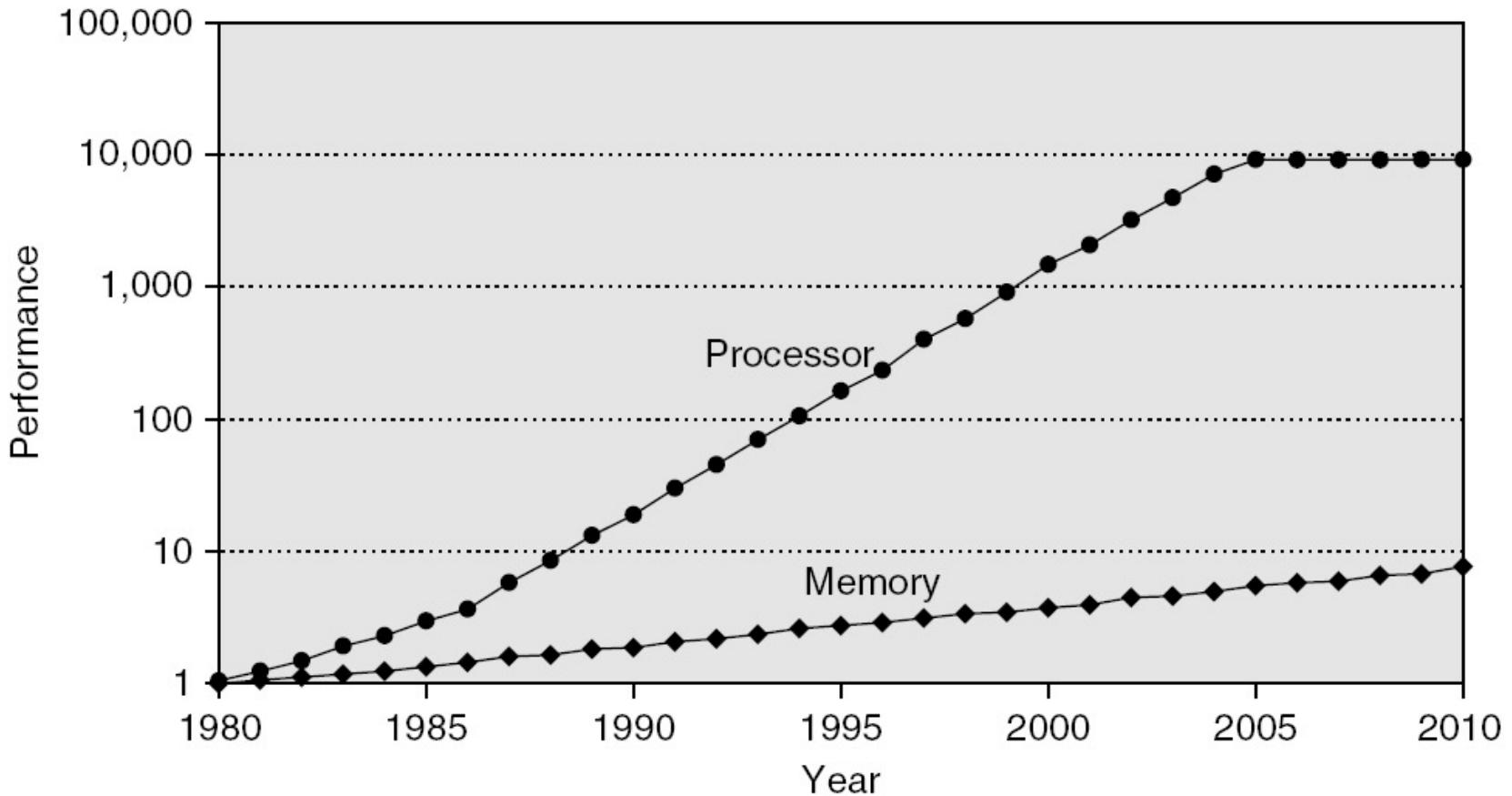
EE-739: Processor Design



Lecture 26 (30 March 2021)

CADSL

Memory Performance Gap



Why Memory Hierarchy?

- Need lots of bandwidth

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

- 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



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



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?



Why Does a Hierarchy Work?

- Locality of reference
 - Temporal locality
 - Reference same memory location repeatedly
 - Spatial locality
 - Reference near neighbors around the same time
- Empirically observed
 - Significant!
 - Even small local storage (8KB) often satisfies >90% of references to multi-MB data set



Four Burning Questions

- These are:
 - Placement
 - Where can a block of memory go?
 - Identification
 - How do I find a block of memory?
 - Replacement
 - How do I make space for new blocks?
 - Write Policy
 - How do I propagate changes?
- Consider these for caches
 - Usually SRAM
- Will consider main memory, disks later

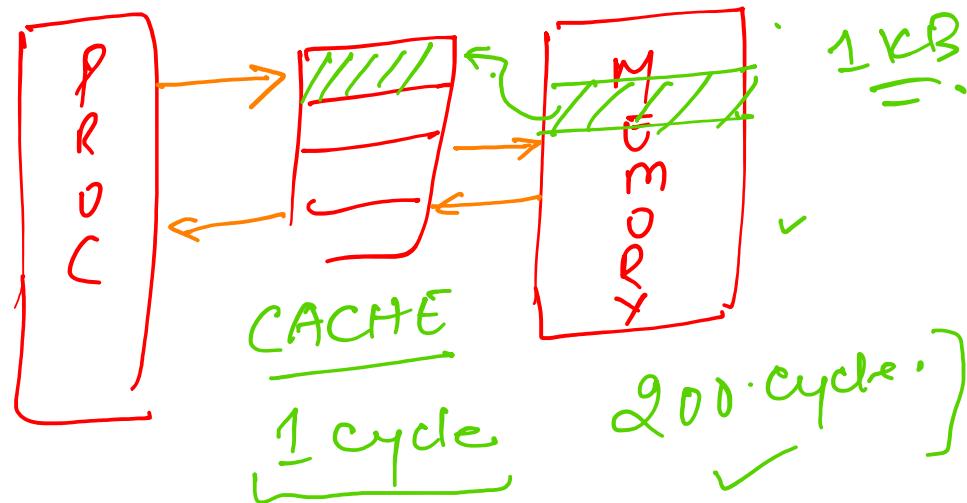


Placement

Memory Type	Placement	Comments
Registers	Anywhere; Int, FP, SPR	Compiler/programmer manages
Cache (SRAM)	Fixed in H/W	<i>Direct-mapped, set-associative, fully-associative</i>
DRAM	Anywhere	O/S manages
Disk	Anywhere	O/S manages

HUH?





- ① Where to place? }
 - ② How to identify }
 - ③ whom to evict
 - ④ When to write

Reading - same image
in RAM (Memory) ?

Write —

Chumic - 16 bytes

454-

Block

↓
Block no.

#Rim

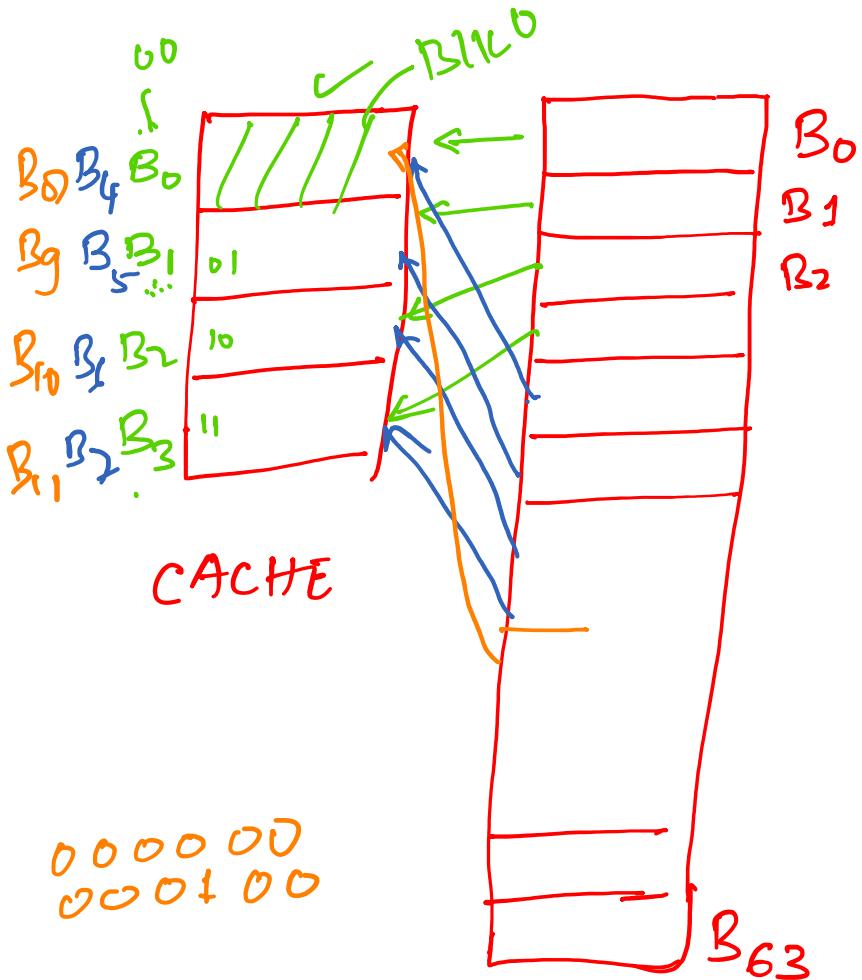
offset within layer

lock

$$\# \text{R} \text{racks} = \underline{64} \quad B_{CJ} - B_0$$



PLACEMENT POLICY



Map Sequentially

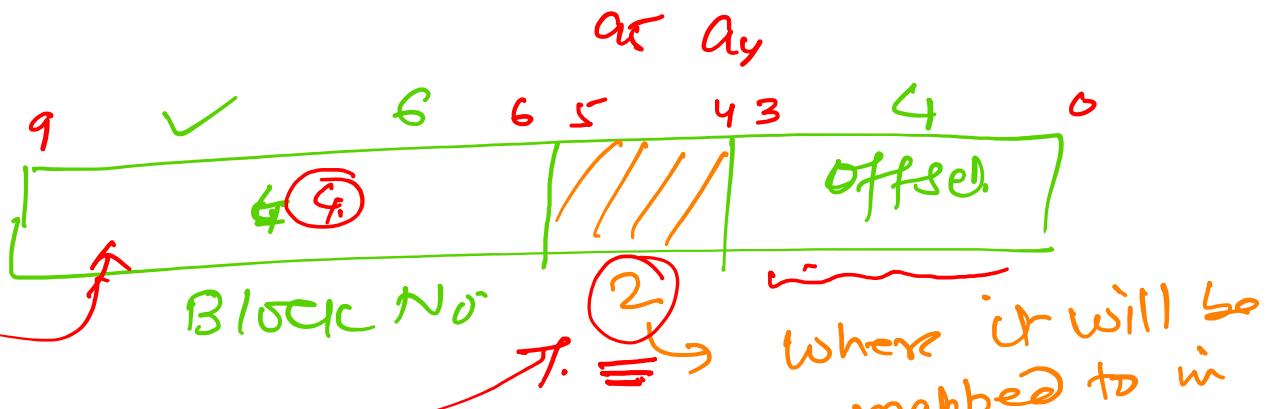
$\underbrace{B_0 - B_{63}}$ to $\underbrace{B_0 - B_4}$.

$\boxed{B_{11K0}} \rightarrow \underbrace{B_0, B_4, B_8,}_{\downarrow} \underbrace{B_{12}, B_{16}, B_{20}, B_{24}, \dots}$

every fourth block gets mapped to $\underline{B_{11K0}}$,

Block No % 4.

$$\underline{\text{BLK No}} = \underline{\text{Block No \% 4}}$$



Identify

BLK# = B0, B4, D8, B12, B16...

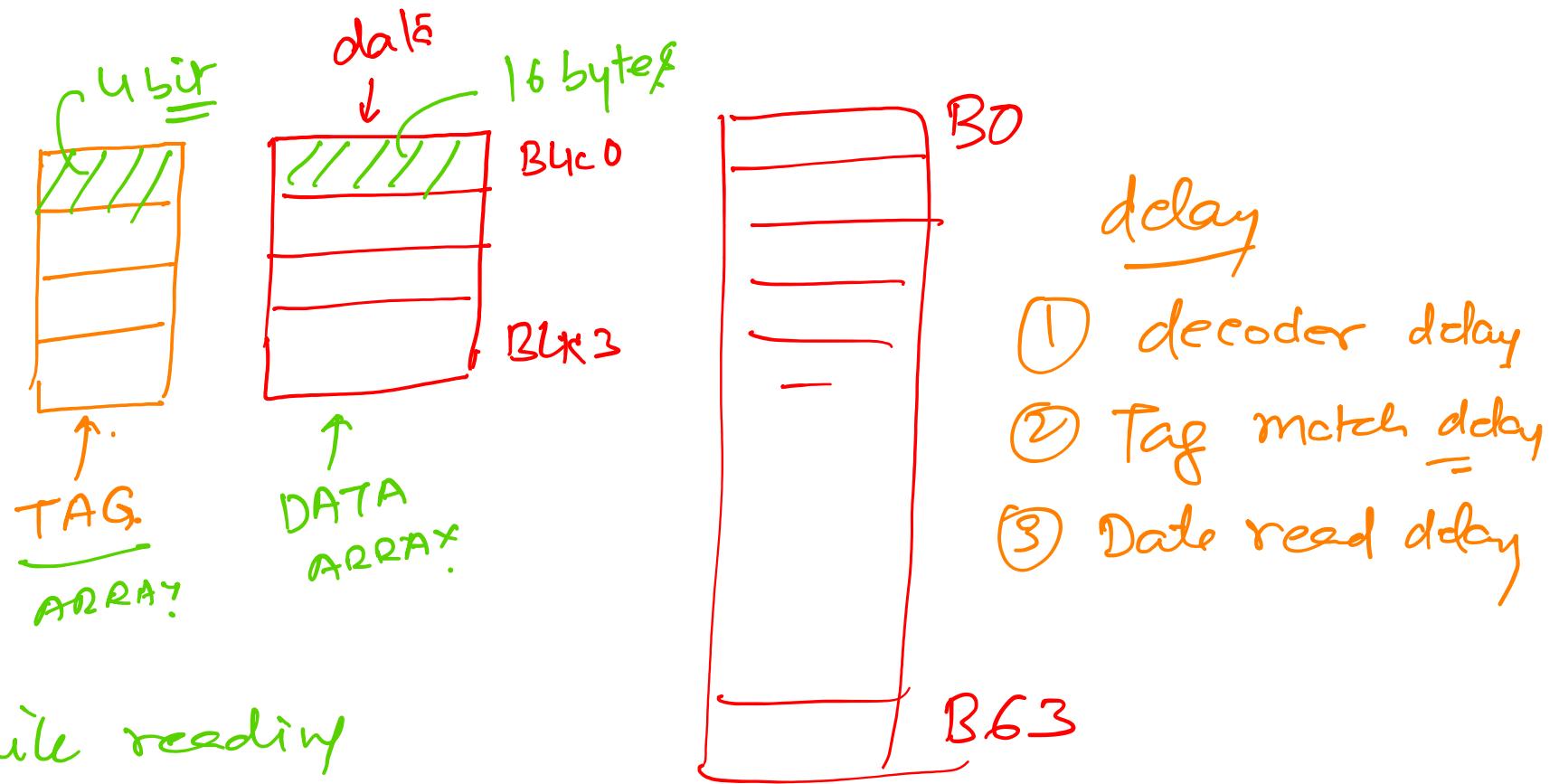
To identify → store the rest of the bits from block address. $6 - 2 = 4$.

TAG

Store the tag.

10001
0010x





while reading

- ① Identify cache block (last 6 bits of Block No.)
- ② match the tag (mes 6 bits of Block no.:)
- ③ If match then read the from data array



BLK0 - B0, B4, B8, B12 - - -

for ($i=0; i<100; i++$)
if ($i \mod 2$)

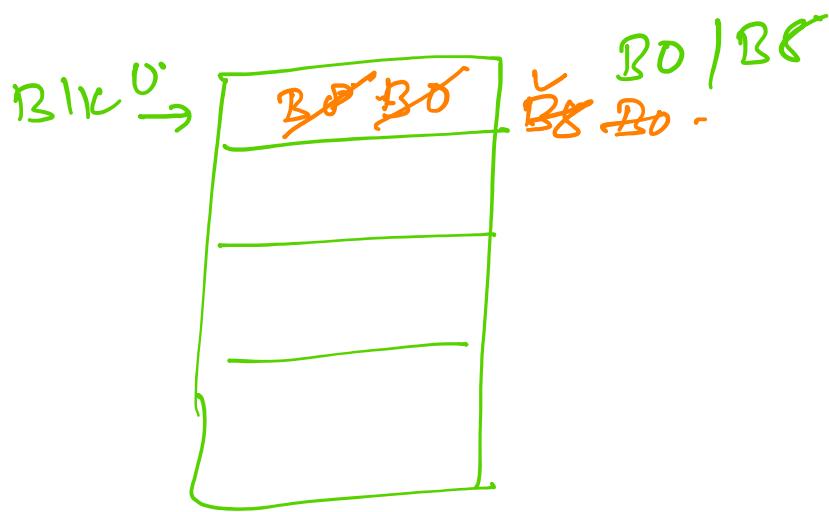
{
 $j = j + 1$ } $\leftarrow B8$

else
 $j = j - 1$ } $\leftarrow B0$

}

Simple.

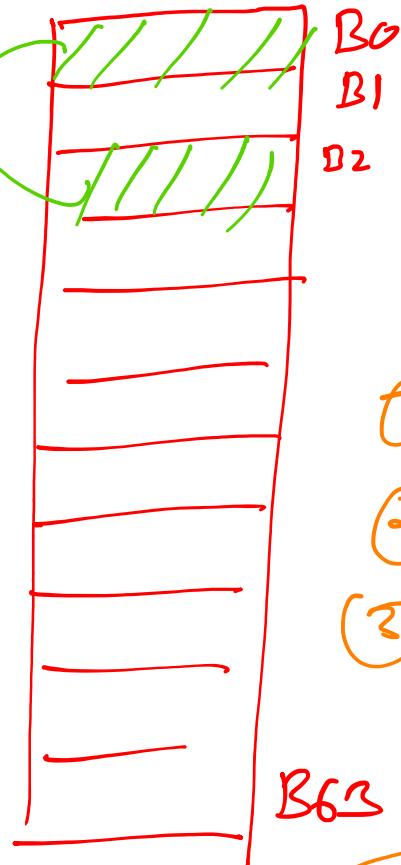
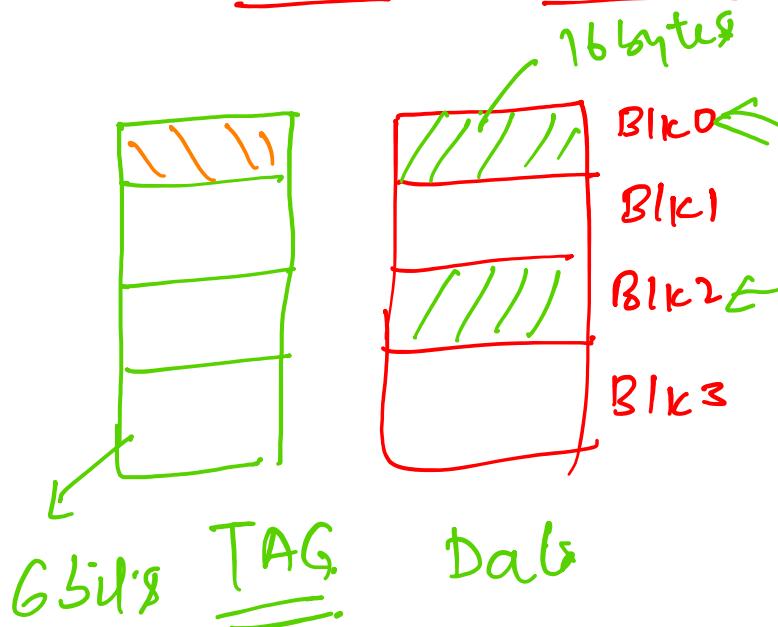
- not very efficient ✓



{
 strict mapping.
 DIRECT MAPPINGS.



Place wherever you want



6	4
Block No	offset

TAG.

delay

- ① Tag comparison.
- ② Delay to decoder.
- ③ Data read. delay

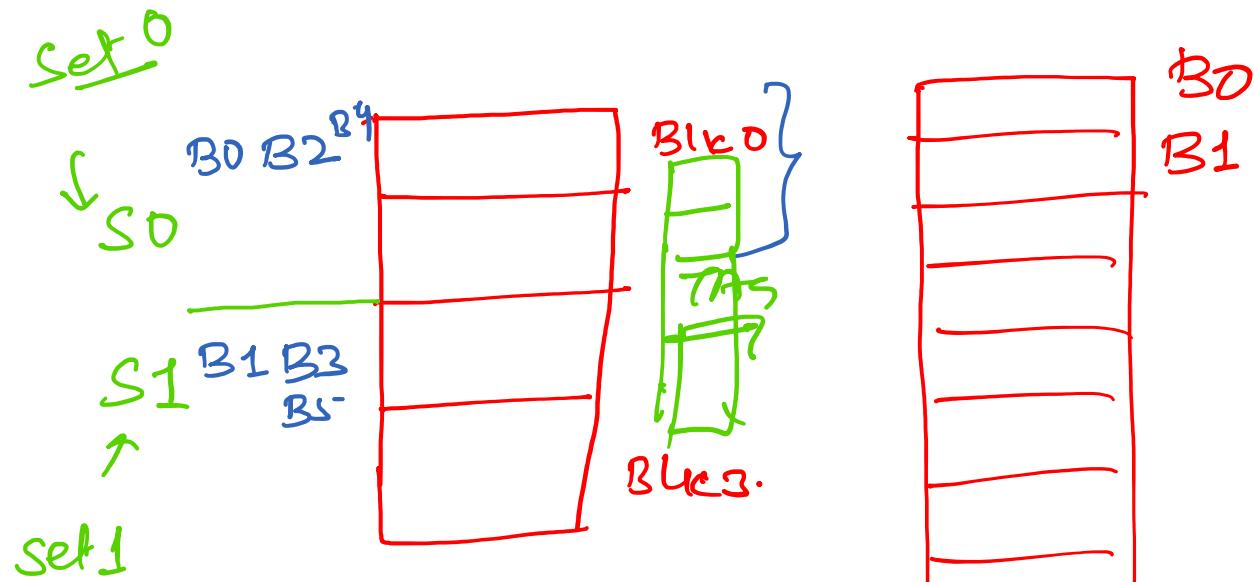
Fully Associative ✓

deep tree A comparisons →

slow.

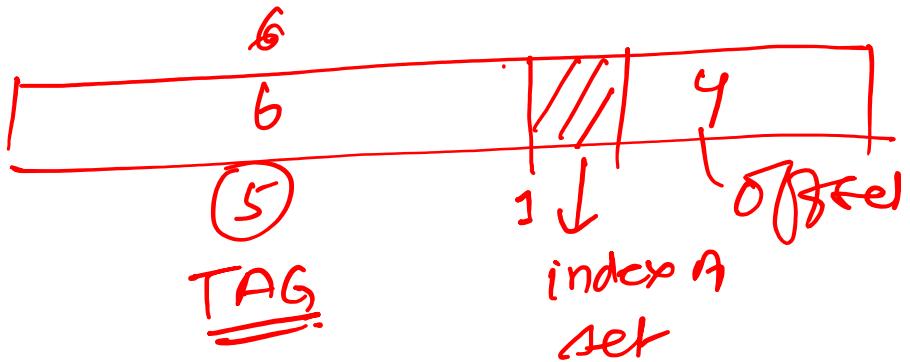


Set Associative Mapping



- Mapping to set is fixed (DM)
- Within set it can be placed anywhere (like FA)





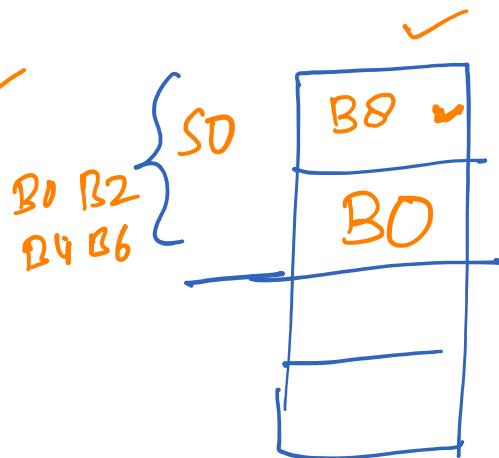
`for (i=0; i<100; i++)`

`if (i%2)`
then
`{ j=j+1 } B8 ✓`

`else`
`{ j=j-1 } B0 ✓`

`3`

only 2 tape ('in this case')
has to be matched in parallel.

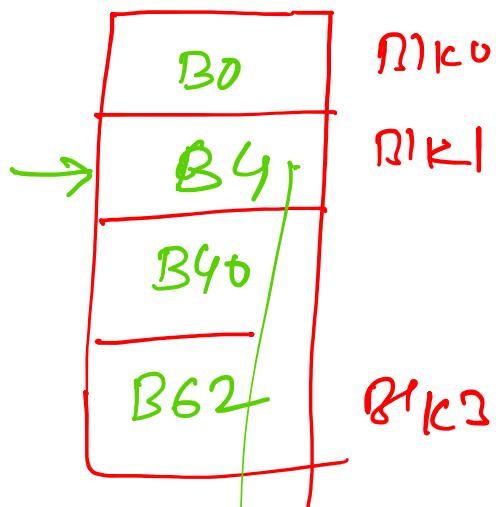
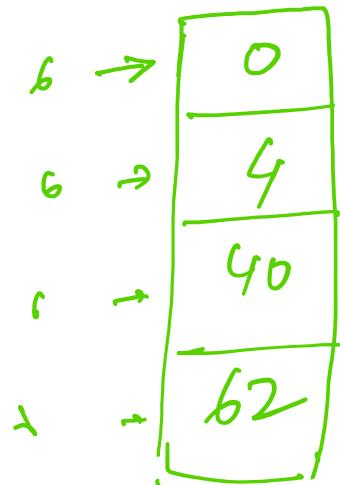


How to place & identify

- ① Direct mapped — fast/instruction
- ② Fully Associative — slow/efficient.
- ③ Set Associative



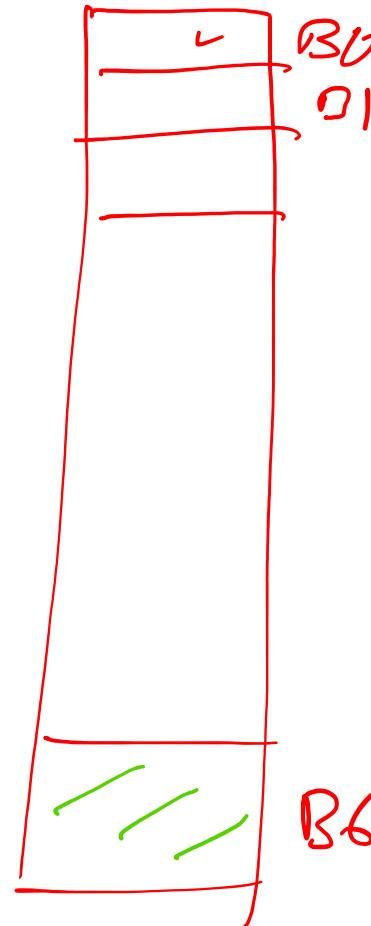
FULLY ASSOCIATIVE



Tag:



4 cycles
sequentially X
Parallel -



Aim → Access in single cycle.

Thank You



30 Mar 2021

EE-739@IITB

19

CADSL