

# 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](mailto:viren@ee.iitb.ac.in)

*EE-739: Processor Design*

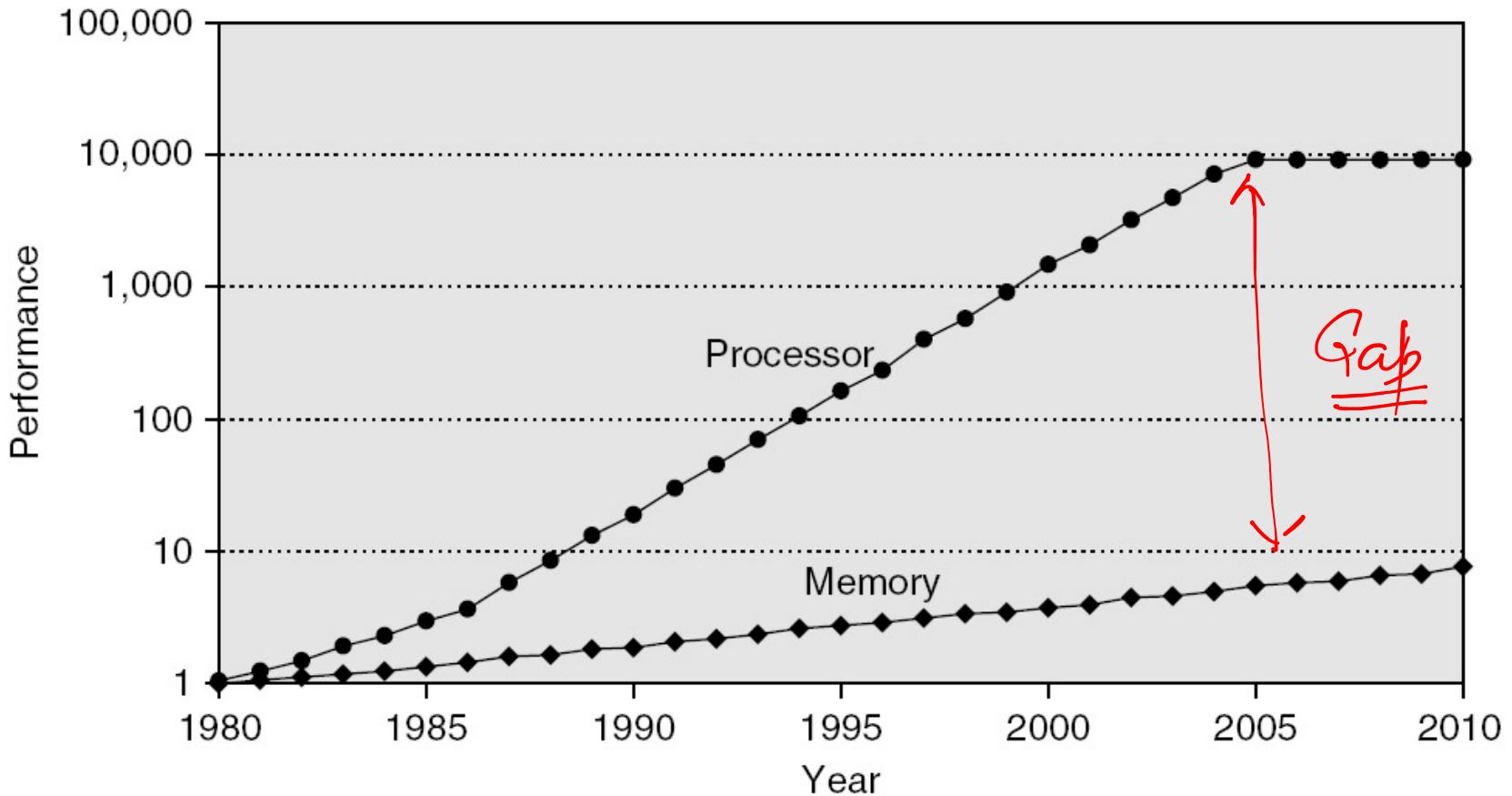
---



Lecture 25 (18 March 2021)

**CADSL**

# Memory Performance Gap



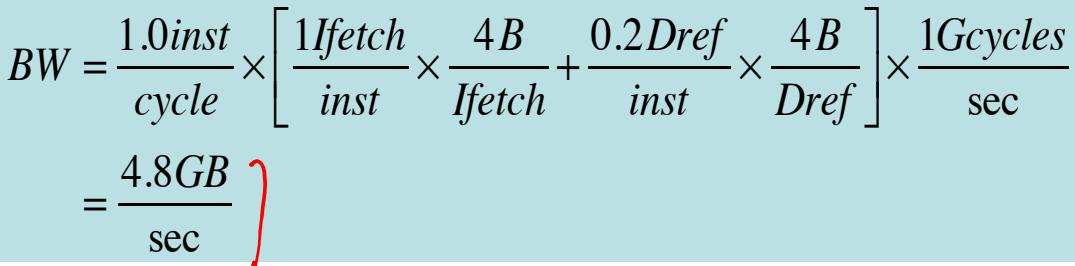
# Why Memory Hierarchy?

---

- Need lots of bandwidth

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

$= \frac{4.8GB}{sec}$



- 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?

{ temporal locality  
spatial 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?

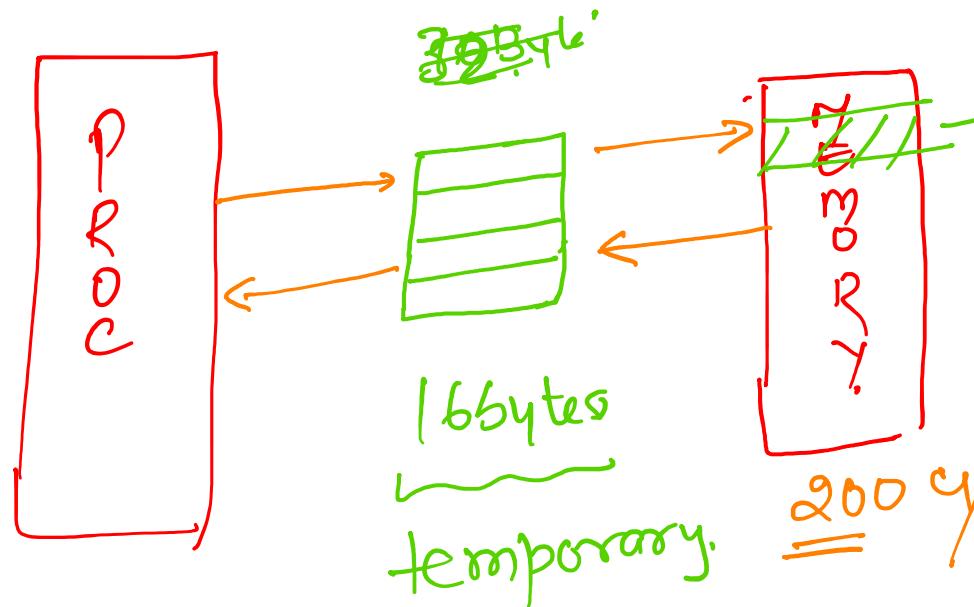


# Why Does a Hierarchy Work?

---

- 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





Instruction

chunk

Bring a chunk from memory & place it in temporary

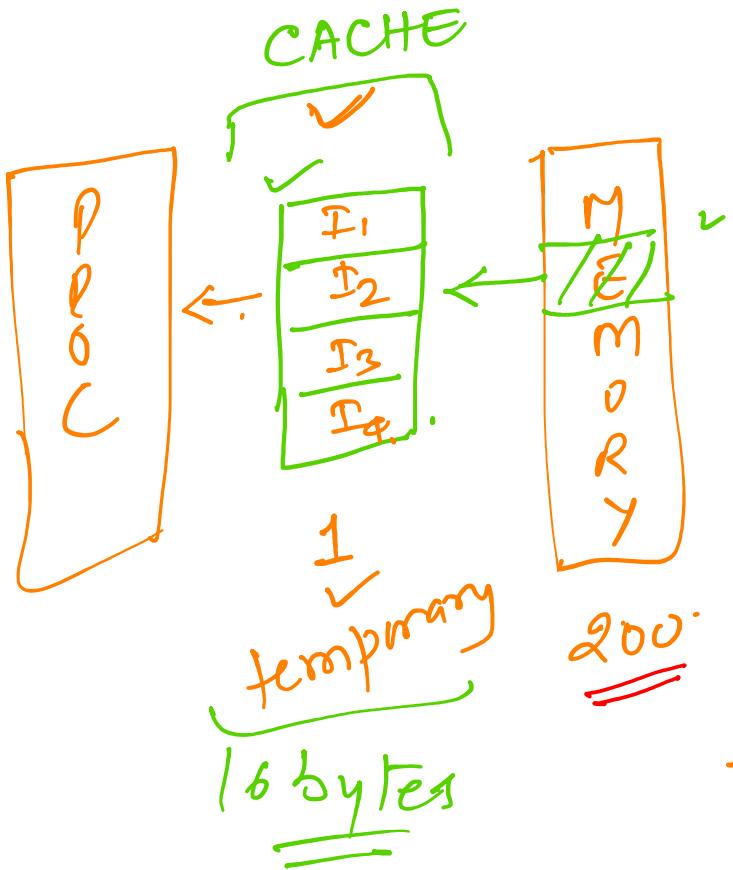
for  $i=0 ; i<100 ; i++$       F - D - E - M A - W

$T_1: a = b + c$   
 $T_2: d = a + e$   
 $T_3: p = \sum r$   
 $T_4: s = p + t$

time:  $(840) \times 100$  cycles

$\underbrace{200+10+200+10}_{T_1}$     $\underbrace{200+10+200+10}_{T_2}$     $\underbrace{200+10+200+10}_{T_3}$     $\underbrace{200+10+200+10}_{T_4}$





average time. fetch instr  
 $= 200$  (without temp memory)

$$= \frac{200 + 100 + 4 \times 99}{400}$$

Spatial locality

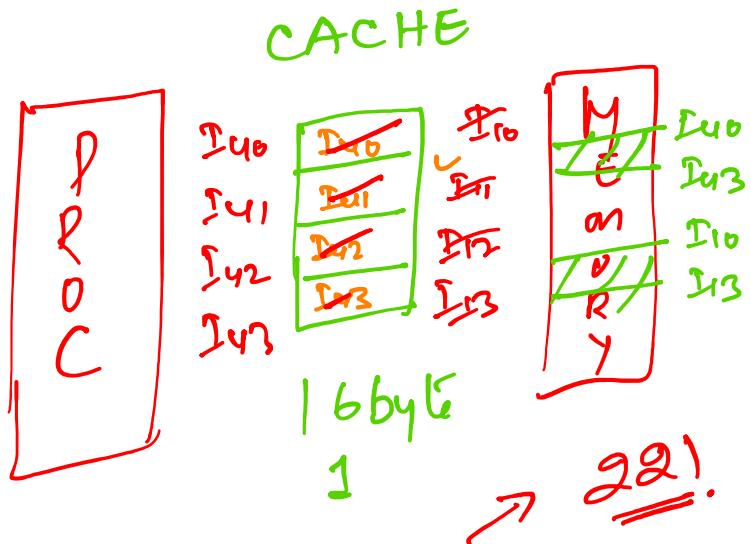
$$\begin{aligned} & 200 + 10 + \underbrace{1 + 10 + 1 + 10 + 1 + 10}_{243} \\ & + 1 + 10 + 1 + 10 + 1 + 10 \\ & = 243 + \underbrace{44 \times 99}_{\text{temporal.}} \end{aligned}$$

$$\text{Speedup: } \frac{840 \times 100}{243 + 44 \times 99}$$

with temp  
memory

$$\begin{aligned} & 1 - 203 + \\ & 1 - 49 \approx 1.5 \end{aligned}$$





$$\frac{200+10}{f} + 1 + 10 +$$

$$200+10+1+10+$$

$$200+10+1+10+$$

$$200+10+1+10$$

+ for  $i = 0 ; i < 100 ; i++$   
 196: if ( $i \% 2$ ) then  
 $\{$   
 $a = b + c$        $I_{40}$   
 $d = a + e$        $I_{41}$   
 $\}$   
 else  
 $\{$   
 $a = b - c$        $I_{70}$   
 $d = a - e$        $I_{41}$   
 $\}$  ]  $\frac{100}{104}$

$$\begin{cases} a = b + c & I_{70} \\ d = a + e & I_{41} \end{cases} \begin{cases} 200 \\ 204. \end{cases}$$

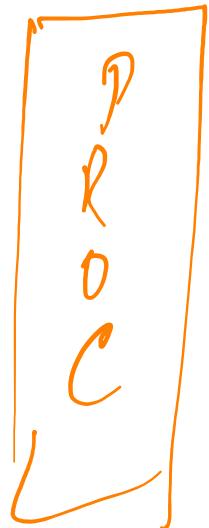
total time  $\rightarrow \underline{221 \times 100}$

average time  $= \frac{200+10}{2} \approx 101$

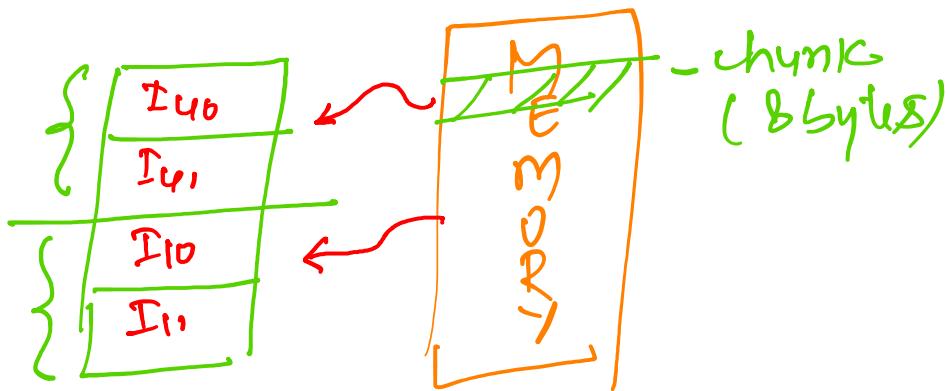
memory.

$$\text{Speedup} = \frac{\frac{420}{241 \times 100}}{\frac{221 \times 100}{221 \times 100}} \approx 2$$





CACHE



16 bytes  
1

2 chunks

$$\underbrace{200+10+1+10}_{1tr1} + \underbrace{200+10+1+10}_{1tr2} + \underbrace{1+10+1+10}_{2tr3} + \underbrace{1+10+1+10}_{2tr4}$$

$$\Rightarrow 221 * 5 221 + 22 \times 98 = 442 + 2156 = \underline{\underline{2598}}$$

$$\text{Speedup} = \frac{420 \times 100}{2598} = \frac{420 \times 100}{2600} = \frac{420}{26}$$

```

for (i=0; i<100; i++)
{
    if (a[i] > b[i])
        c = b+c - I40
        d = c+d - I41
    else
        c = b-c - I10
        d = a-c - I11
}

```



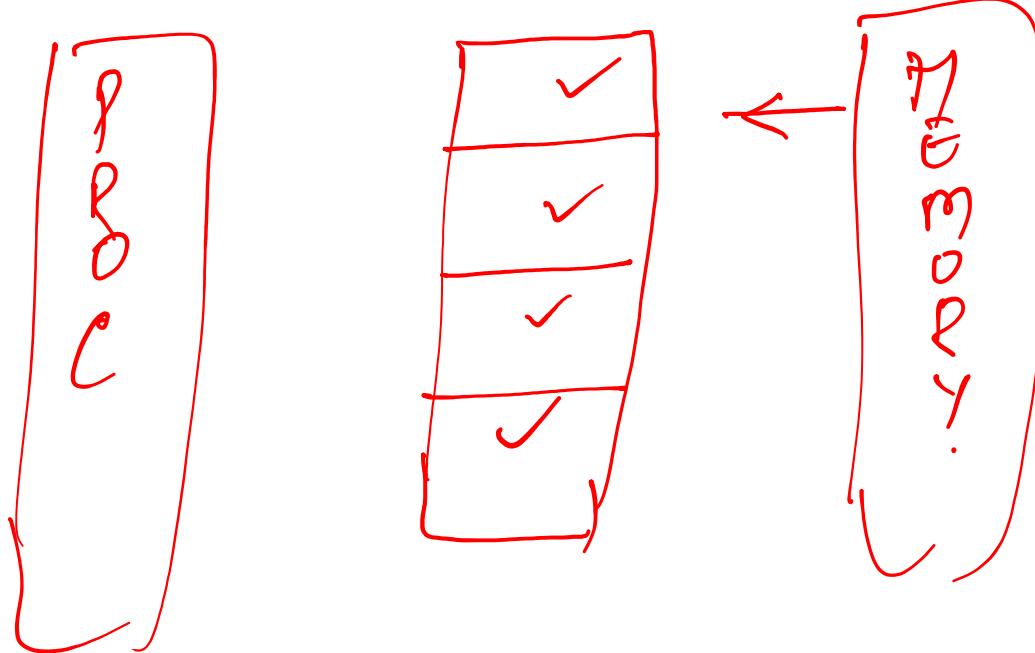
Average memory access time

$$= \frac{200H + 200T + (\underline{1+1}) \times 98}{200}$$

?



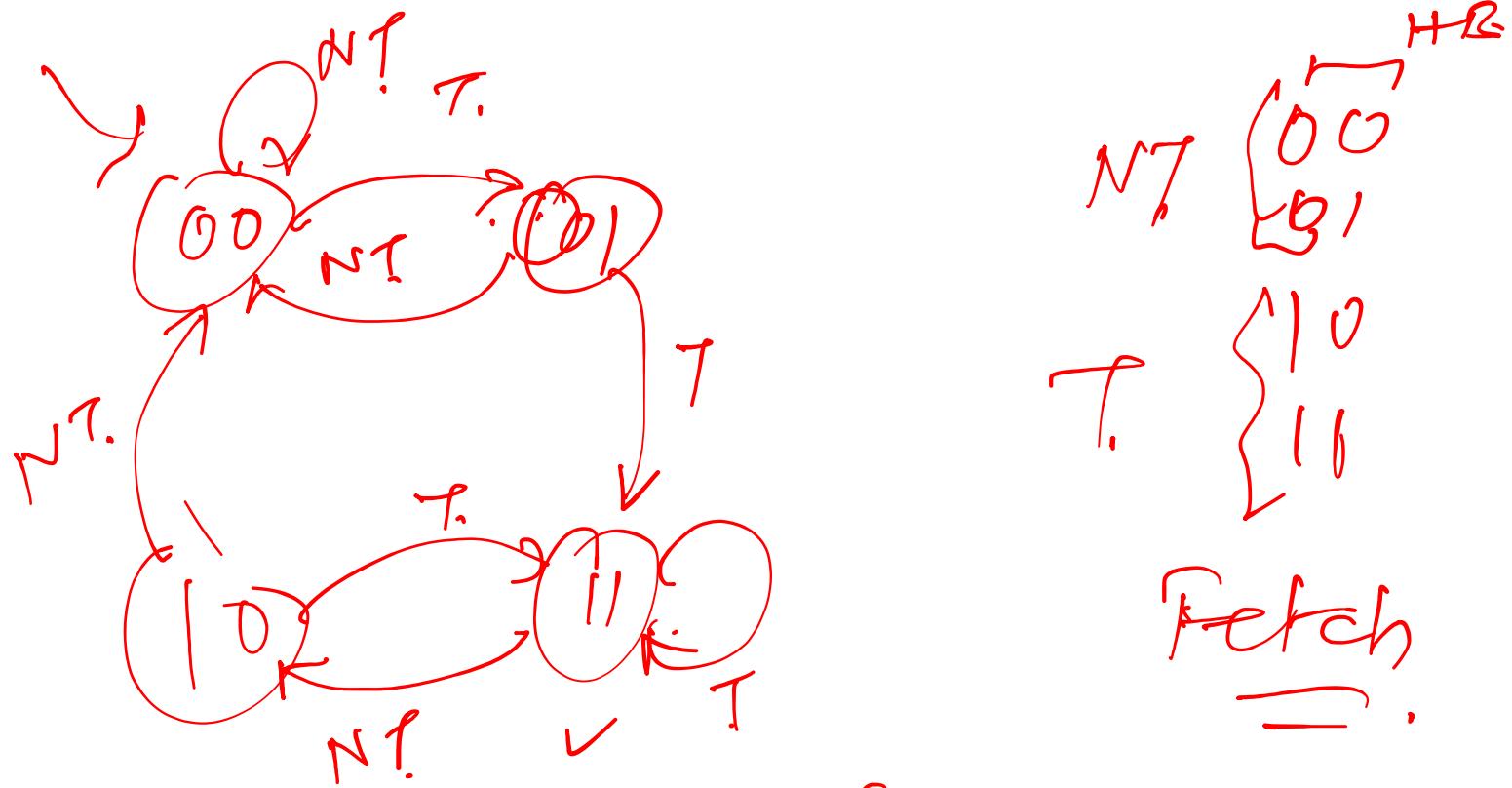
# CACHE



- ① Where to place.
- ② How to identify
- ③ How whom to evict

32 byte





Fetch =

Predict → fetch ← NT.

Validate → Execution ← Pass → fail → flush instruction  
fetched after this  
(IF & ID stage)

~~for~~  $i = 0$ ;  $i < 100 \geq i++$  for ( $i = 0$ ;  $i < 100$ ,  
 $i++$ )  
for ( $j = 0$ ;  $j < 100$ ;

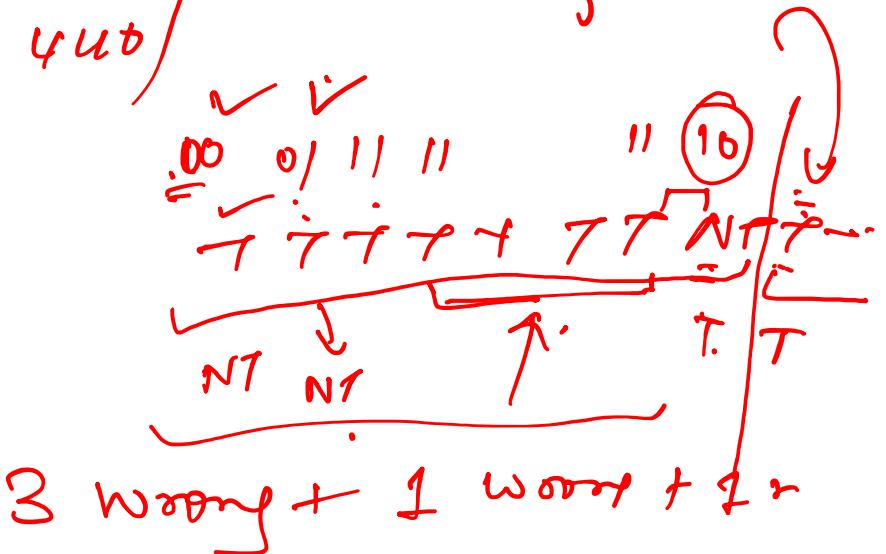
for (j=0; i<100;  
j++)

$\stackrel{10^0}{\Rightarrow}$  it { 2nd }  $\leftarrow$  5 )  
= 3 400 )  
+ 3 400 )  
-----  
0 ]

10

1

6



# 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



# Thank You



18 Mar 2021

EE-739@IITB

16

**CADSL**