

# LECTURE 24

06/03

## Computer Architecture So Far...



### ECE243 Simple Processor

- Four registers  $r_0 \rightarrow r_3$ , word size 8 bits
- 256 bytes of memory, has a program counter
- Only 10 instructions  $\rightarrow$  each takes two arguments
- NEW  $\rightarrow$  Z-bit = 1 if result of most recent computation is zero  
 $\hookrightarrow$  N-bit = 1 if result of most recent computation is negative

$\hookrightarrow$  Used in conditional branch instruction

- Consider add instruction  $\rightarrow$  add  $r_1, r_2 \quad r_3 \leftarrow r_1 + r_2$

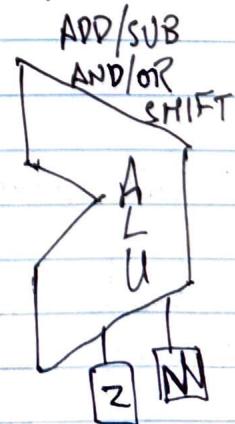
① Sketch encoded instruction from mem

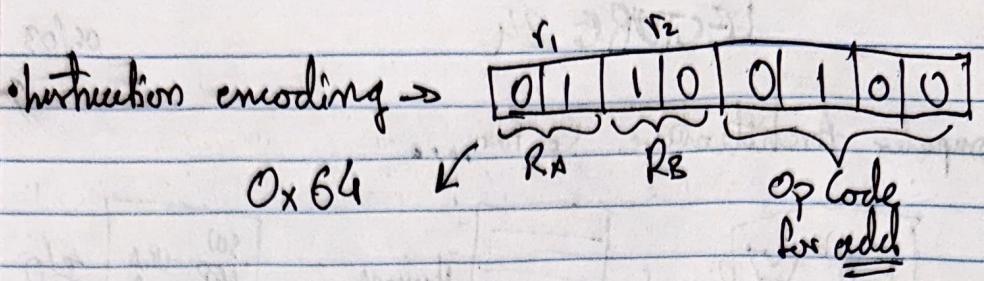
② Figure out what instruction is

③ Compute  $R_3 \leftarrow R_1 + R_2$

④ Set N & Z

⑤ Increment pc





## LECTURE 25

10/03

- OR-immediate  $\rightarrow$  immediate is 5 bits  $\rightarrow$  no room

$\hookrightarrow$  So we always store into R<sub>1</sub> to specify a register and zero-extend the immediate

We reserve two 4-bit opcodes  $\leftarrow$  (0111) & (1111) to ORI

				$1 \ 1 \ 1 \ 1$			
				imm		opcode	

## Conditional instructions

- In SimPloc we use result of prev operation to compare to zero

$\hookrightarrow$  bnez imm4  $\rightarrow$  execution  $\rightarrow$  if ( $z == \text{imm}^4$ )  $PC \leftarrow SE(\text{imm}^4)$   
else  $PC + 1$

- Relative branching gives us movable code

## LECTURE 26

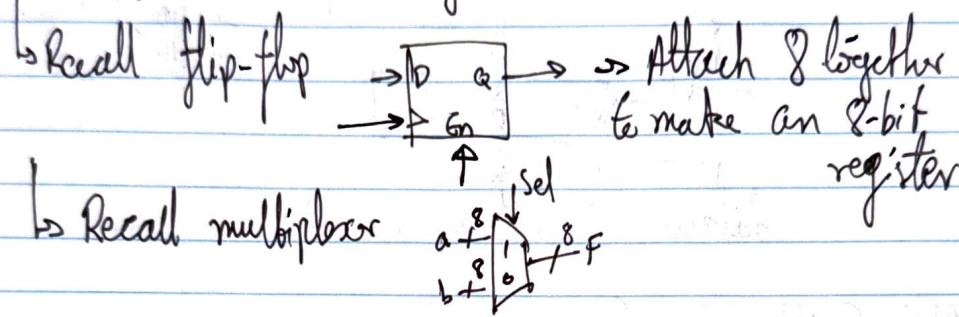
12/03

# Let's write a small assembly program. (we are assembly)

$\hookrightarrow$  subtracts 2 numbers and branches back to the beginning

Address	Instruction	Encoding	Hex
0	loop: load R <sub>2</sub> , (R <sub>0</sub> ) #R <sub>2</sub> ← mem(R <sub>0</sub> )	10000000	0x80
1	load R <sub>3</sub> , (R <sub>1</sub> ) #R <sub>3</sub> ← mem(R <sub>1</sub> )	11010000	0xA0
2	sub R <sub>2</sub> , R <sub>3</sub> #R <sub>2</sub> ← R <sub>2</sub> - R <sub>3</sub>	10110110	0xB6
3	brnz loop	11011001	0xD9

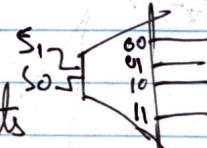
- Recall from EC6241 a few things



## LECTURE 27

13/03

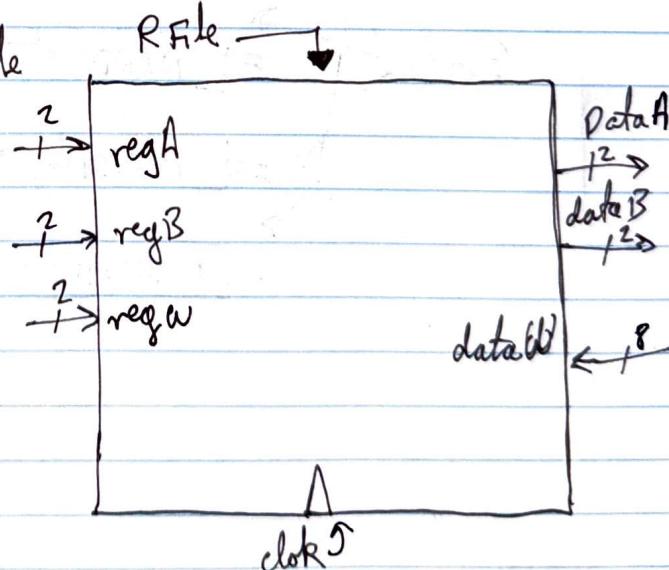
- Recall also from 241 the decoder →

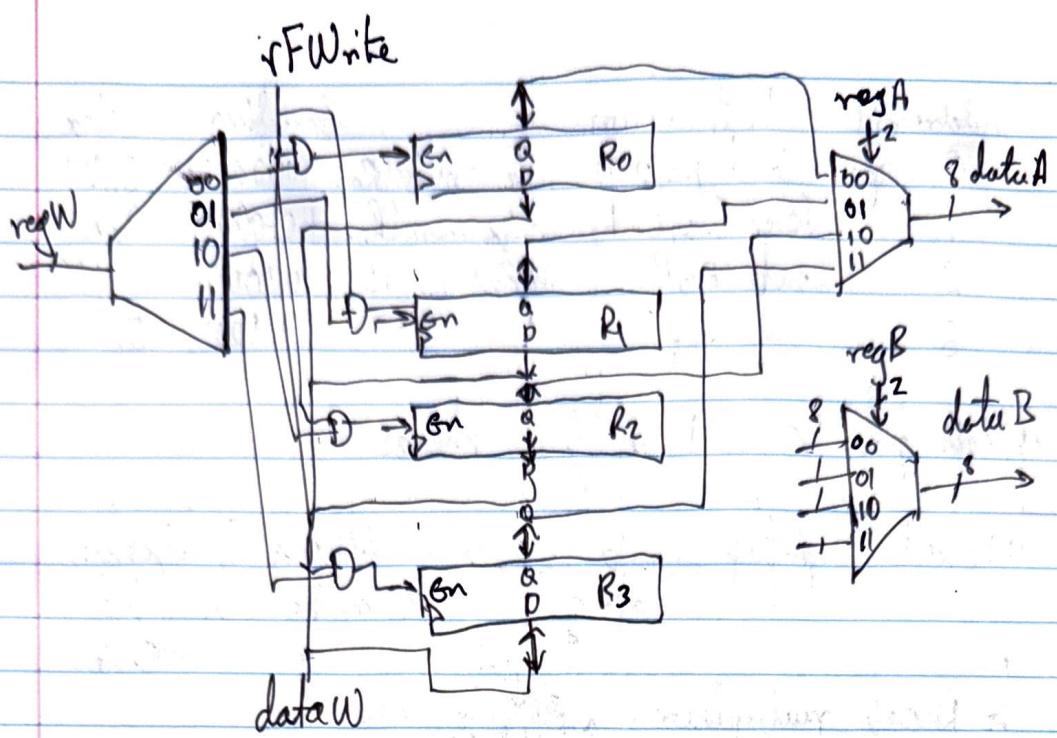


- Building more detail for SimProc elements

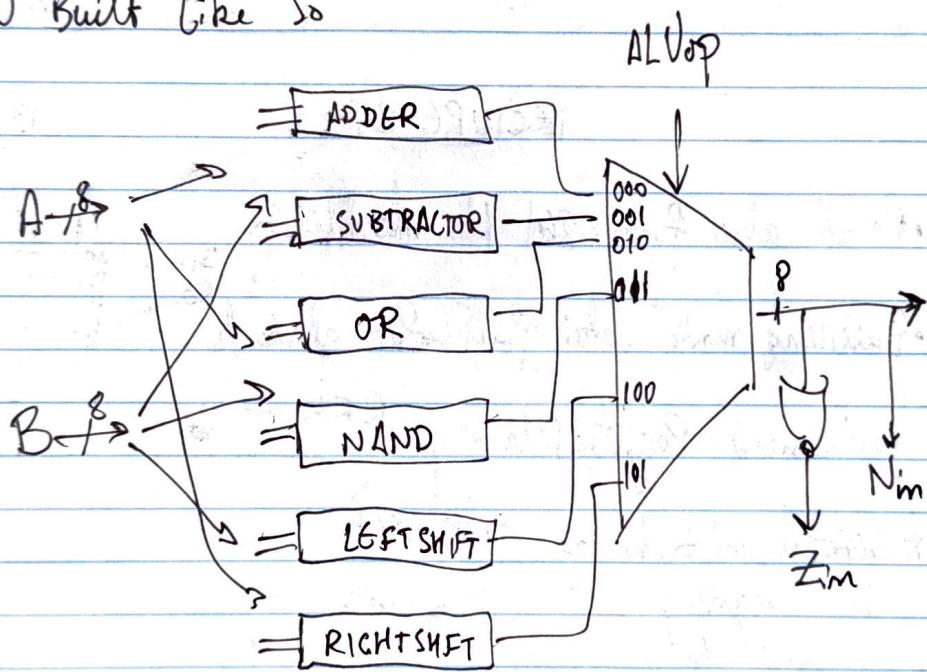
- ↳ We need Register File

\* Careful not to change  
the registers  
accidentally

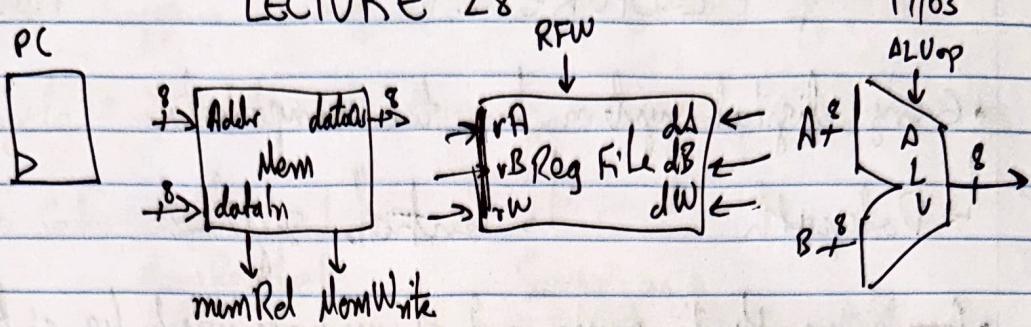




ALU Built Like So

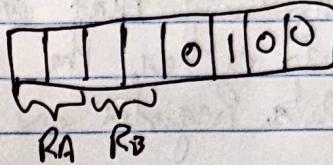


## LECTURE 28



17/03

- Here's what happens in full execution



- ↳ ① Fetch op code from memory
- ↳ Put this into the Instruction Register
- ↳ ② Decode the instruction → Stored & decoded in IR
- ↳ ③ Read RA & RB from register
- ↳ ④ Compute ALUout  $\leftarrow R_A + R_B$
- ↳ ⑤ Write ALUout back into RA
- ↳ ⑥ Set Zm & Nm based on the operation's result
- ↳ ⑦ Compute PC = PC + 1

\* Special case happening w/ ORI instruction

↳ Look at last 3-5 mins of lecture to review

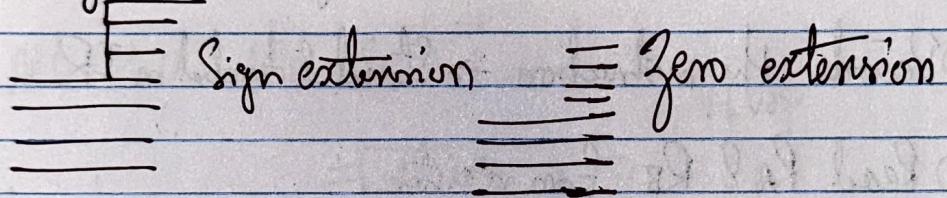
## LECTURE 29

19/03

- Every digital system has two components
  - ↳ Datapath
  - ↳ Control system

- Anything that comes out of memory must be sorted into either the Instruction Register or Memory Data Register

- Dereferencing data/storing into memory done by interfacing w/ Address & the memsel mux.



- How does IR know if opcode is 3 or 6 bits?

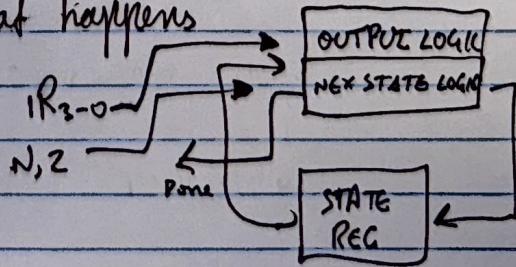
- ↳ SimProz automatically responds to X011 as a shift regardless of what X is

## LECTURE 30

20/03

### Finite State Control of SimProz

- FSM controls which signals are activated & when (on which clock signal) that happens



• What happens during each cycle for ADD/SUB/NAND

① Cycle 1 → Fetch instruction from memory:  $IR \leftarrow Mem(pc)$

↳ Set AddrSel = 1

MemRead = 1

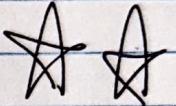
IRLoad = 1

Saves as a

cycle of compute!

↳ We can also get  $PC \leftarrow PC + 1$  in the same cycle

↳ Set ALU\_A = 0 ( $PC \rightarrow ALU$ ) } clock-to-Q time  
ALU\_B = 001 ( $ALU_B \rightarrow 1$ ) } ensures we don't  
ALU\_op = 000 (add) } screw up here  
PCwrite = 1 } ↳ setup & hold time



Shows up in many  
interviews

② Cycle 2 → Decode the signal in IR (needs 1 full cycle)

↳ We guess we'll probably need to read R<sub>A</sub>/R<sub>B</sub>

↳ Set RASel = 0 (selects R<sub>A</sub> to regA)

ABLD = 1

③ Cycle 3 → Add/Sub/NAND specific

↳ Set ALU\_A = 1 (select A to ALU)

Set ALU\_B = 000 (select B to ALU)

Set ALU\_op = one of add/sub/nand

ALU\_outLD = 1

flagWrite = 1

④ Cycle 4  $\rightarrow$  Get ALU.out to Ra

$\hookrightarrow$  Set RegIn = 0 (connect ALU.out  $\rightarrow$  dataW)

RFWrite = 1

Done = 1  $\rightarrow$  This is for next-state logic

Now let's do "Load Ra, (RB)"

$\hookrightarrow$  Only need cycle 3 onwards

③ Cycle 3  $\rightarrow$  Addr Sel = 0

MDRload = 1

MemRead = 1

② Cycle 4  $\rightarrow$  RegIn = 1 (MDR load  $\rightarrow$  dataW)

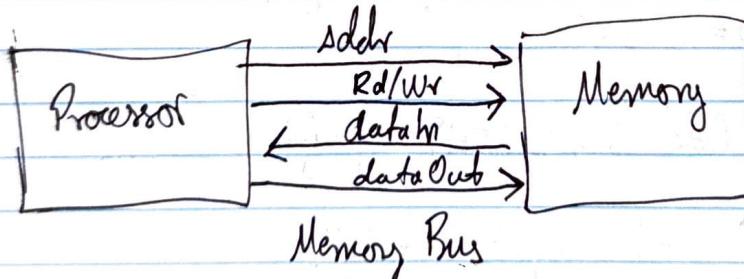
RFWrite = 1

Done = 1

## LECTURE 31

24/03

### • Introduction to Cache Memory



Processor has two memory access types

$\hookrightarrow$  Read (memory or instructions)  $\rightarrow$  both look the same to memory

$\hookrightarrow$  Write (store)

- Some physics for context:
  - ↳ Larger memories are slower  $\rightarrow$  smaller memory faster
    - ↳ Longer wires have more resistance & capacitance, making signal propagation slower
    - ↳ Decoder time will get slower in larger memories
  - Interesting point is that as chips get physically bigger they don't necessarily get slower
  - 16 GB memory system takes ~300 cycles to respond to memory-access request

How to solve the problem of a fast processor but a big, slow memory?

- ↳ Put a small fast memory in between the processor and main memory & keep the much smaller active memory in between.
- ↳ Called the cache memory



- Sets of instructions and certain data structures benefit from having close addresses
- ↳ Locality of reference
- ↳ Some sets of instructions called multiple times

- Two types of Locality
  - ↳ Spatial locality  $\rightarrow$  physically closer together
  - ↳ Temporal Locality  $\rightarrow$  Addresses same over time

### Operation of a Memory Read Using a Cache

- ① Processor requests to read data from address A.
- ② Cache (somehow) looks inside to see if it has a copy of data at address A.
  - ↳ If no, we call a cache MISS. Cache requests data from A (and some chunk of block).
  - ↳ Memory responds slowly, cache copies to itself and sends to the processor
  - ↳ If yes, then we call it a cache HIT. Quickly send that data to the processor in just a few cycles
- ③ If processor then requests A+1 it will come quickly after A b/c we fetched block surrounding A.

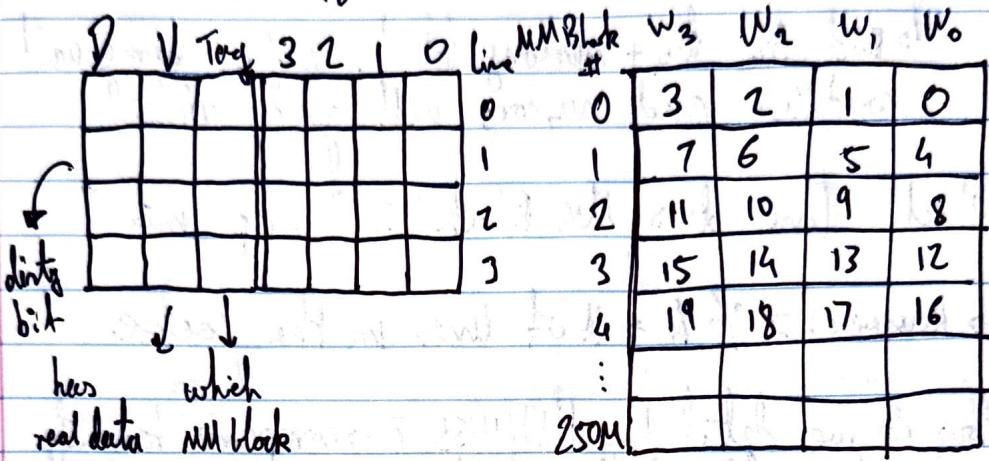
### LECTURE 32

26/03

### Logistics of the Cache Memory

- ① Because cache asks for 31 mem location when MISS, we consider main memory as organized into blocks (say 6 words for e.g.).

- We assume every word has an address



## The Direct-Mapped Cache

- Cache has rows & columns like MM, each row called cache line

↳ line size equals main memory block size b/c we take whole blocks into cache lines

- Tag is group of bits that tell you which MM block is in the cache line

- Valid bit tells us if there's real data in the cache or not

- Dirty bit tells us if we write something new to cache, needs to send information back to MM.

- Now consider when an address is requested by the processor → where can the data go?

↳ 1 line → "Direct-Mapped" Cache

↳ good b/c we only have to look at 1 tag/line in cache for HIT/MISS

↳ bad b/c two + memory blocks using same will conflict and memory will get enriched

Which line does MM Block "J" map into?

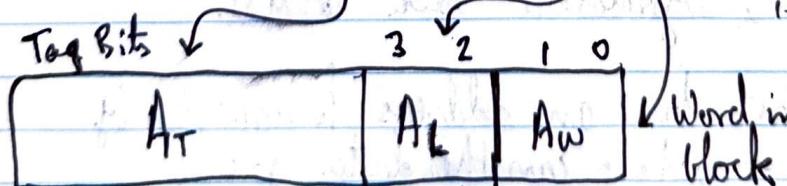
↳ Answer:  $J \% 4 \rightarrow \#$  of lines in the cache

How do we detect the HIT/MISS?  $\rightarrow$  compare tag of cache line (if valid) w/ the tag of requested memory block.

Consider these possible memory access request addresses

Block # Address Bits

	31 30.....5 4 3 2 1 0	
0	0 0 0 0 0 XX	
1	0 0 0 0 1 XX	
2	0 0 0 1 0 XX	
3	0 0 0 0 1 1 XX	
4	0 0 0 1 0 0 XX	
5	0 0 0 1 0 1 XX	Tag (A2)
6	0 0 0 1 1 0 XX	
7	0 0 0 1 1 1 XX	If matches, and
8	0 0 1 0 0 0 XX	is valid then it's a HIT



Line the  
addr gets  
mapped to

## LECTURE 33

27/03

Consider what happens when the processor makes a read request at MM address  $A_T$

↳ Cache looks inside itself to see if  $\text{tag}(A_L) = A_T$

↳ If yes, and  $V(A_L) = 1$ , then cache has a HIT

↳ Then the cache sends that word quickly to the processor

↳ Otherwise it's a MISS

↳ The line in  $S_L$  is valid, it must be evicted  
The direct-mapped cache only allows an MM block in one place

↳ The new block is slowly retrieved from MM and copied into cache line  $A_L$

↳ Set  $\text{tag}(A_L) = A_T$ ,  $V(A_L) = 1$ ,  $D(A_L) = 1$

Eviction of a cache line

• Method 1 → Write-Back Cache

↳ If there was a write/store to any word in the line, the Dirty bit will be set

↳ When  $D(A_L) = 1$ , upon eviction, the line must be written-back so MM is correct

★ Multiple writes to a line results in only 1 write to MM

- Method 2  $\rightarrow$  Write-Through Cache

- ↳ Has no Dirty bit

- ↳ All writes to a line are also automatically sent to Main Memory

- ↳ Generally undesirable as we may have multiple writes to Main Memory in very short span.

---

## LECTURE 34

31/03

### New Cache Type: Associative Caches

- Every Main Memory block can map into every cache line

- ↳ Good: Solves the conflict problem in best way possible

- ↳ Bad: Needs a lot of digital circuitry (comparators) to simultaneously check all tags of all lines

- There is a compromise b/w the 2 (direct-mapped & fully associative).

- ↳ Allows each Main Memory block to be placed in a smaller set of lines

- ↳ Say "M" different lines  $\rightarrow$  called M-Way Set-Associative Cache

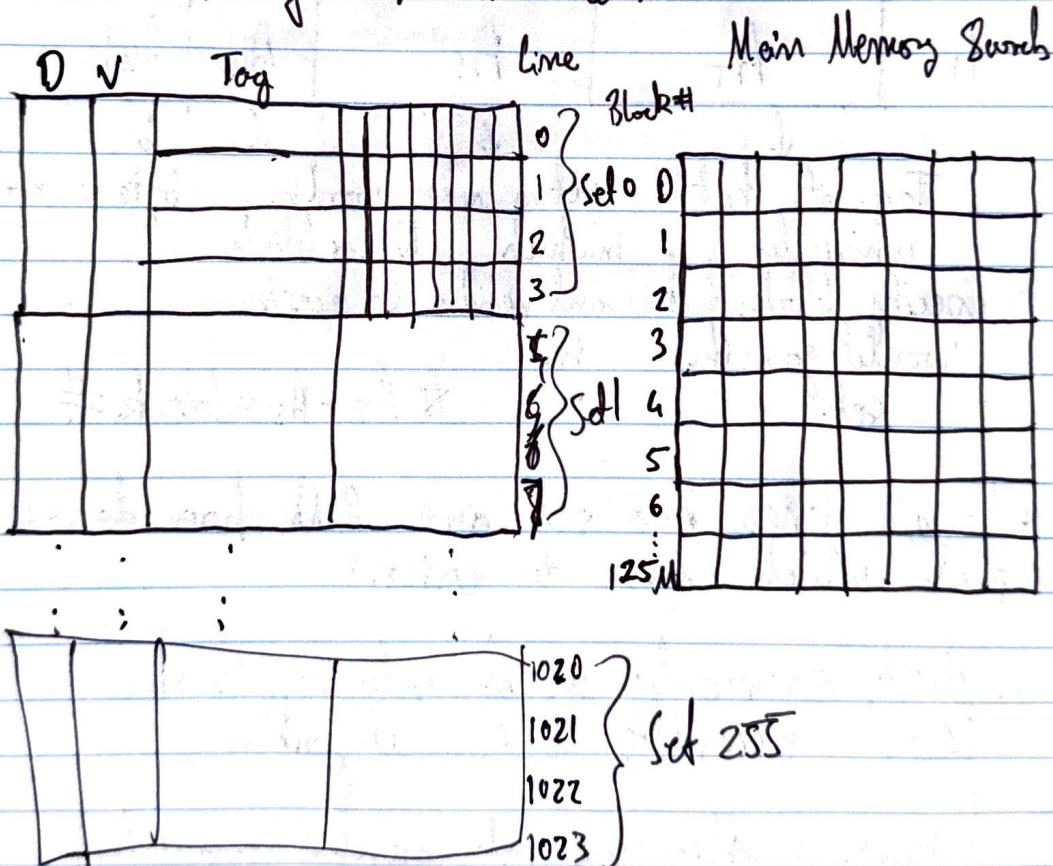
E.g. 1  $\rightarrow$  A 2-Way Set-Associative Cache allows an MM block to be placed in 2 different cache lines

$\hookrightarrow$  Those 2 lines are called a set.

E.g. 2  $\rightarrow$  For NIOS-V (addressability w/ bytes), cache size is 32 Kbytes

$\hookrightarrow$  Cache has 32 bytes  $\star$  per line,  $32K / 64K$  always means nearest power of 2  $\star$

$\hookrightarrow$  It is a 4-Way Set-Associative Cache



\* For an N-Way Set-Associative Cache, each memory block can be placed into N-different cache lines

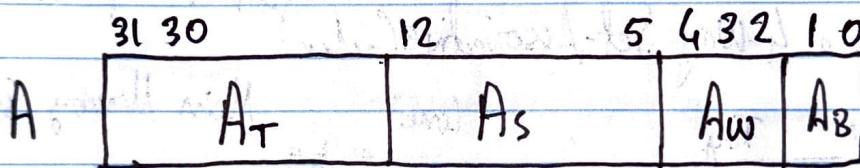
To determine the set-associative cache mapping / HIT/MISS, we need to compute:

$$\hookrightarrow \# \text{ of lines} = \frac{\# \text{ bytes in cache}}{\# \text{ bytes/line}} = \frac{32K}{32} = 1K = \underline{1024}$$

$$\hookrightarrow \# \text{ of sets} = \frac{\# \text{ lines}}{N} = \frac{1024}{8} = 256$$

$\hookrightarrow$  Need  $\log_2(256) = 8$  bits to select a set

$\hookrightarrow 8 \text{ words/line} \rightarrow \log_2(8) = 3$  bits to determine which word in a line/block to address



↓  
Tag bits that uniquely specify blocks that would go in the set.  
set.      ↓  
set m cache word in which line or block resides to access  
byte or word

$$\star AT + As = \text{block \#}$$

- If all lines in a set are full, how do we pick which one to evict.

Choosing randomly in general about same as LRU

$\hookrightarrow$  Most commonly used solution is the Least Recently Used algorithm

$\hookrightarrow$  Have another two bit counter associated w/ each line in the set

$\hookrightarrow$  Set to 0x00 if line is accessed and increment by 1 each time not used

## LECTURE 35

02/04

- Spent some time on cache simulators → check notes

### Analyzing a Program's Cache Behaviour

- Consider some code examples w/ cache behaviour

Eg 1 → Conflicting access (assume DM w/ BS=8)

```
int A[100], B[100], sum;  
sum = 0;  
for (int i=0; i<100; i++) sum += A[i] + B[i];
```

↳ Suppose every element  $A[i]$  maps to same line as  $B[i]$

↳ Accessing  $B[i]$  after accessing  $A[i]$  will kick  
8 integers out of cache

↳ Happens each iteration of for loop, destroys cache

↳ Soln: maybe split into two loops

Eg 2 → Accessing 2D arrays e.g. frame buffer from lab 7

```
void bad_clear(void) {  
    for (int c=0; c<240; c++) {  
        for (int r=0; r<240; r++) {  
            buf[r][c] = 0x0;  
        }  
    }  
}
```

↳ Horrendous because of row-major array interferes  
w/ column-major clear function

### Eg 3 → Accounting for cache size

```
int A[N], sum, prod;  
for (int i=0; i<N; i++) sum = sum + A[i];  
for (int i=0; i<N; i++) prod = prod * A[i];
```

- ↳ If  $N \gg$  cache size, we have to clear cache and redo whole cache retrieval for the second loop
- ↳ Soln: do both sum & prod calculations in same for loop

## LECTURE 36

03/04

### Performance Analysis of Caches Running a Program

- Recall how memory calls work
  - ↳ Processor emits a request to access memory
  - ↳ If memory address in the cache, then we have HIT and responds quickly, say in "c" cycles
    - ↳ These days  $c \approx 4$  cycles
  - ↳ If cache miss, takes M cycles for memory to find that address
    - ↳ M is called the miss penalty,  $\approx 300$  cycles

• Define hit rate = num cache HITs / total num access

↳ always  $\Leftrightarrow 0 \leq h < 1$

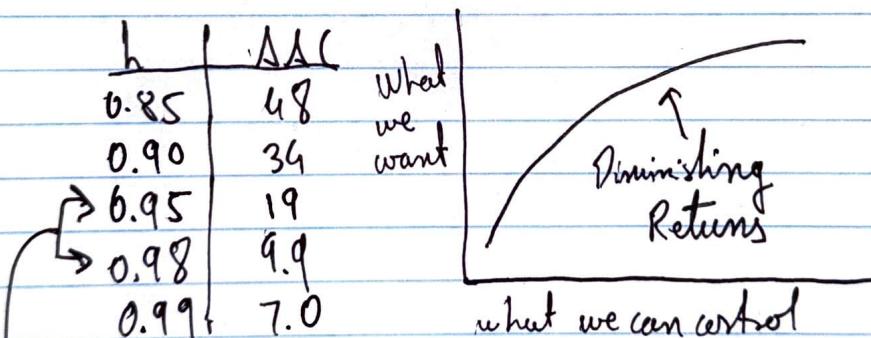
↳ We can formulate by using  $h, C, M$  compute average number of cycles to complete a memory access request.

Eq. How many cycles would 1000 accesses take?

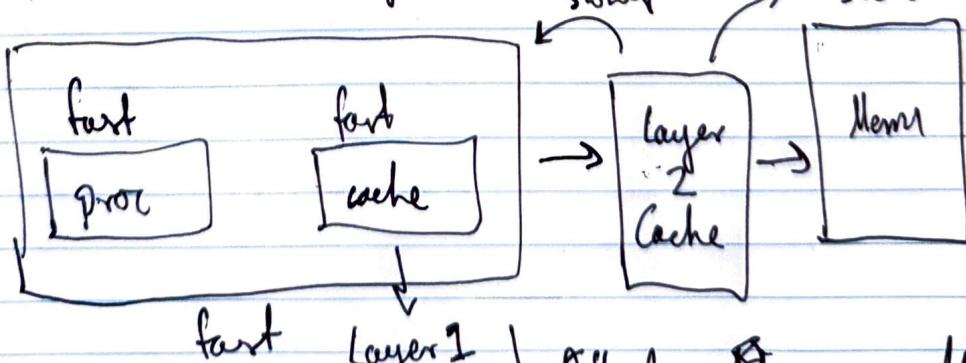
$$\frac{1000 \cdot h \cdot C + 1000 \cdot (1-h) \cdot M}{1000}$$

• Have quantity Average Access Cycles (AAC)

↳  $h \cdot C + (1-h) \cdot M$



↳ 3% diff  $\rightarrow$  2x performance



↳ Effective miss penalties seen by L1 cache is lower

- Because instruction memory accesses have different patterns than data accessed, the L1 cache is split into
  - ↳ L1 instruction cache
  - ↳ L1 data cache
- Discourse on Intel's reign and arrival of competition
- Discourse on China-Taiwan-TSMC relations