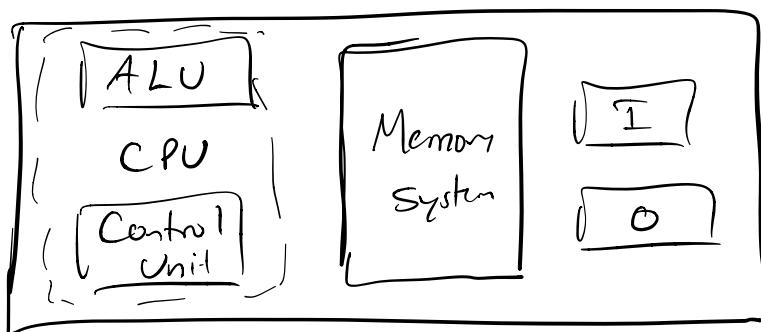


BASICS OF COMPUTER ARCHITECTURE

Program → Assembly / Instruction Set (IS) → Machine Language

↓

C++
C
Java

- Processor-specific language
- x86, ARM
- Seq. of bytes that tell processor to do things
- 1011... . . .

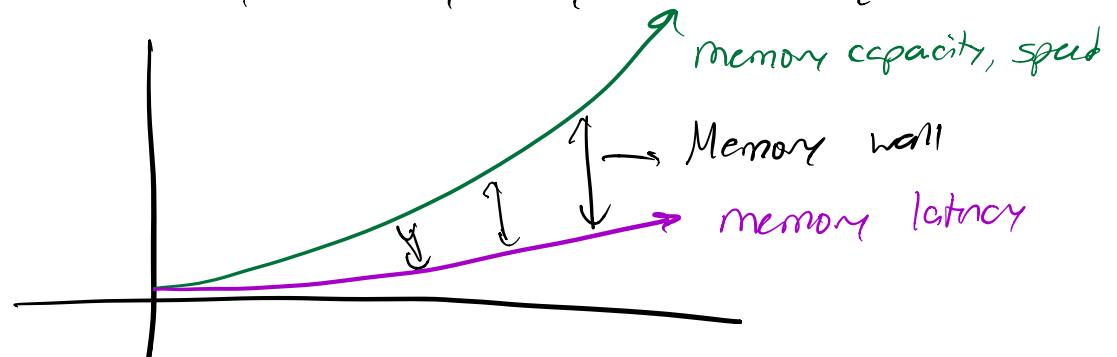
program → OS → compiled → IS → microarch → register → circuits

- Difference between processor arch. + microarch:

- Processor arch: interface of hardware to software
- Microarchitecture: implementation of hardware

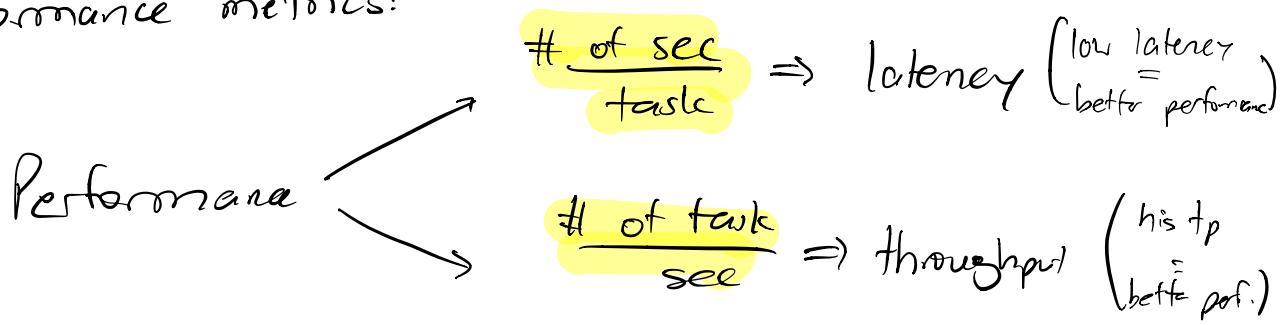
METRICS

- Moore's Law: every 1.5-2 yrs, 2x transistors on chip area.
- Processing speed ↑, energy consumption ↓, memory capacity ↑
- Memory wall: memory latency only increases by 1.1x



▫ More and more caches.

- Performance metrics:



- Comparison of Ferrari w/ bus.
- Comparison: speedup

$$\text{Speedup of } A \text{ over } B = N = \frac{\text{throughput}_A}{\text{throughput}_B}$$

Or

$$N = \frac{\text{latency}_B}{\text{latency}_A} \Rightarrow N \propto \frac{1}{\text{latency}}$$

- Iron Law of Performance: Whole programs

$$\text{CPU time} = \frac{\text{instructions}}{\text{Program}} \times \frac{\text{cycles}}{\text{instruction}} \times \frac{\text{Seconds}}{\text{cycle}}$$

$\text{instructions} \xrightarrow{\text{Compiler}}$
 $\text{cycles} \xrightarrow{\text{IS}}$
 $\text{Seconds} \xrightarrow{\text{Clock}}$

$\text{Program} \xrightarrow{\text{Algorithm}}$
 $\text{instruction} \xrightarrow{\text{Microarch}}$
 $\text{cycle} \xrightarrow{\text{Circuit}}$

$\text{cycles} \xrightarrow{\text{Microarch}}$
 $\text{Seconds} \xrightarrow{\text{Transistor}}$
 $\text{cycle} \xrightarrow{\text{Microarch}}$

$$\circ \text{ Latency: } \frac{\text{seconds}}{\text{instruction}} = \frac{\text{cycles}}{\text{instruction}} \times \frac{\text{seconds}}{\text{cycle}}$$

$$\circ \text{ Throughput: } \frac{\text{instruction}}{\text{second}} = \frac{\text{instn.}}{\text{cycle}} \times \frac{\text{cycle}}{\text{seconds}}$$

↳ MIPS: millions of instr./second.

Ex:// CPI: 2 . Clock = 500 MHz ($\frac{\text{cycles}}{\text{sec}}$) , Find MIPS

$$\begin{aligned}\text{MIPS} &= \frac{\text{Instr}}{\text{sec}} \times 10^6 \Rightarrow \frac{\text{Instr}}{2\text{cycles}} \times \frac{500 \times 10^6 \text{cycles}}{\text{second}} \\ &= \frac{250 \text{ instr}}{\text{sec}} \times 10^6 \\ &= 250 \text{ MIPS}\end{aligned}$$

- Amdahl's Law: improving part of the program \Rightarrow overall speedup?

$$\text{Overall speedup} = \frac{1}{(1 - \text{Frac}_{EHN}) + \frac{\text{Frac}_{EHN}}{\text{Speedup}_{EHN}}}$$

- Frac_{EHN} = proportion of program improved (% of execution time)
- Speedup_{EHN} = how fast did proportion speedup.

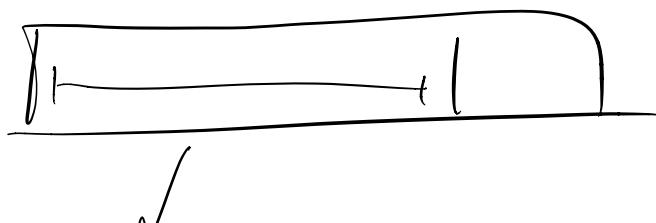
Ex:// 2 improvements for program

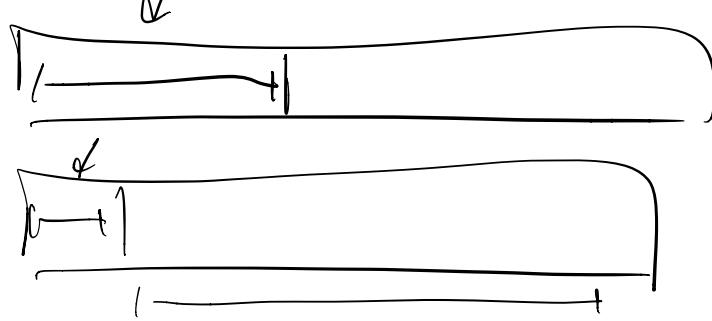
- 1) Speedup 20x on 10% of total time
- 2) Speedup 1.6x on 80% of time

$$1) \text{ Speedup}_{\text{overall}} = \frac{1}{0.9 + \frac{0.1}{20}} = 1.05$$

$$2) \text{ Speedup}_{\text{overall}} = \frac{1}{0.2 + \frac{0.8}{1.6}} = 1.43$$

- Common improvement \rightarrow high improvement on small program
- Law of Diminishing Returns: improve multiple aspects of program.





- Power improvements :

- o Dynamic power: activity consumes energy.

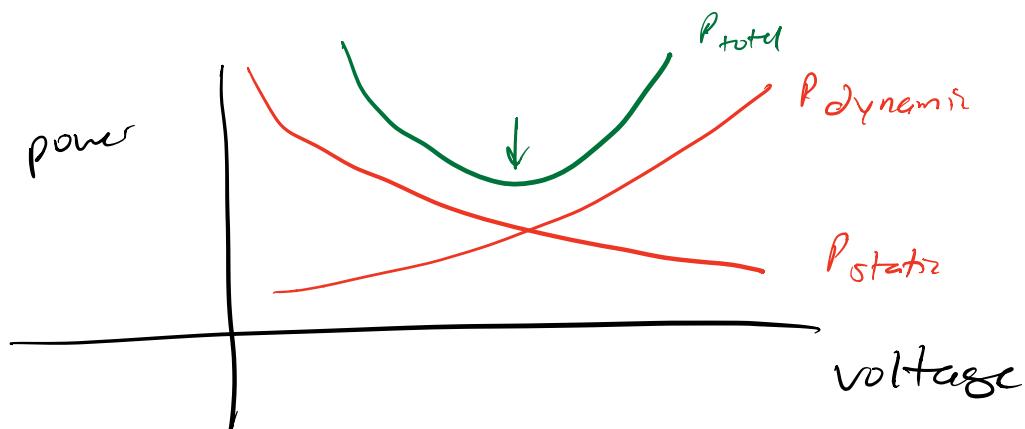
$$\text{Dynamic power} = \frac{1}{2} C V^2 f \alpha$$

\uparrow \downarrow \hookrightarrow

exp. freq. % of active
volt. transistor

- o Static power: idle consumption.

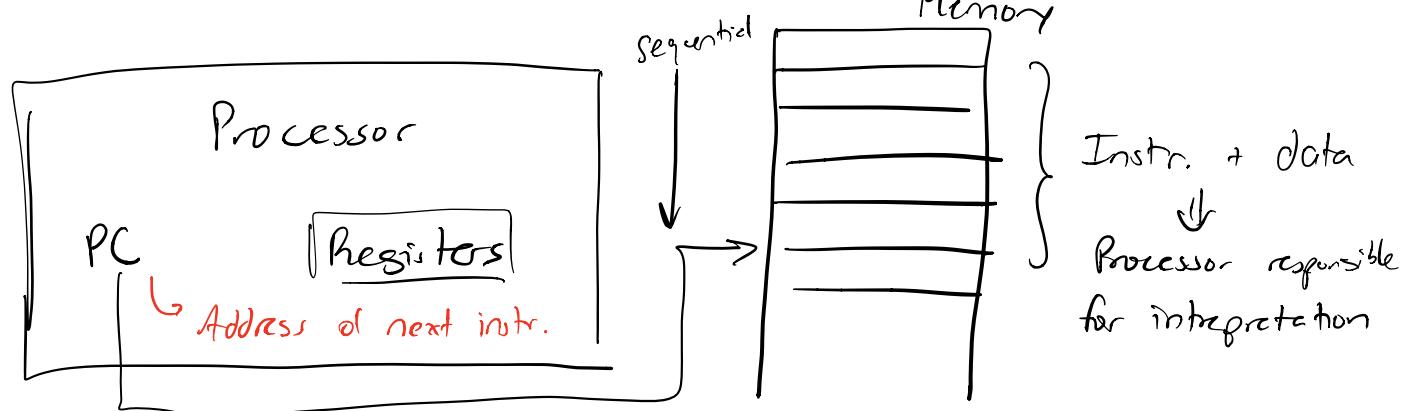
voltage $\downarrow \rightarrow$ leakage \uparrow (power wall)



INSTRUCTION SET ARCHITECTURES

- Defn: contract between software + hardware

- Von Neumann: model of all computers. Stored program model



Design of ISAs

RISC (ARM)

- Simple instructions

- $\frac{\text{instr}}{\text{program}} \uparrow, \frac{\text{cycles}}{\text{instr}} \downarrow$

CISC (x86)

Code

- Complex instruction \Rightarrow dense

- $\frac{\text{instr}}{\text{program}} \downarrow, \frac{\text{cycles}}{\text{instr}} \uparrow$

- Instruction format: length + encoding of instruction

Length:

- Fixed: easy computation but difficult to program
- Variable: opposite pros + cons

Encoding

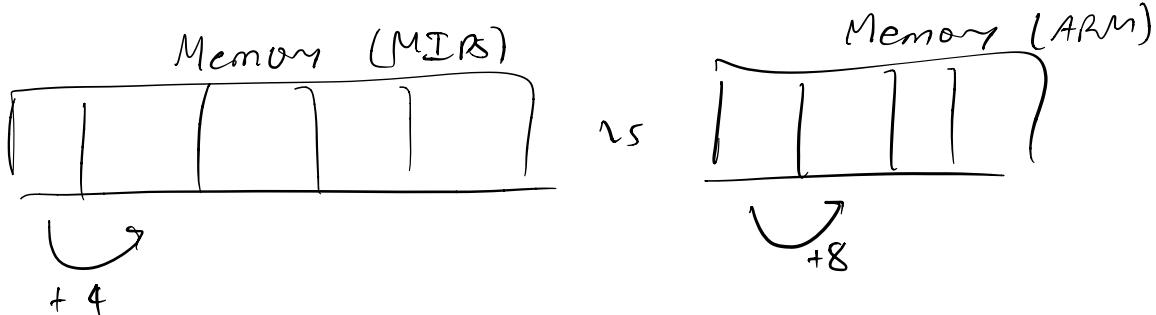
- # of ways to encode op.
- Simple + low # of encodings

- Opcode: tells us which operation to perform

- Operand model: # of operands, sequence, types of operands (memory / register operands)

ARM ISA OVERVIEW

- 64-bit architecture: every memory address is 64 bits



- Registers: X0 - X30 (64 bit)

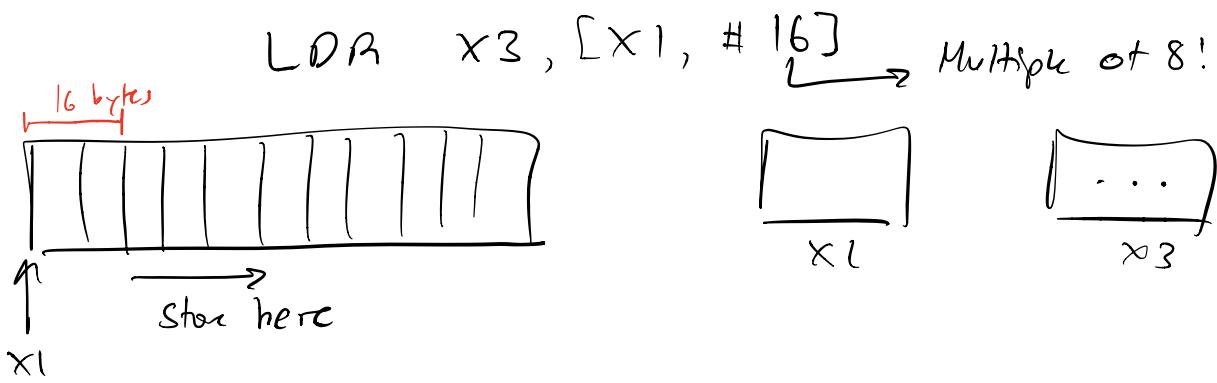
- Special: XZR (zero register), SP (stack ptr)

- 32-bit registers: W0 - W30, WZR, WSP

- Labels: refers to specific lines \Rightarrow "alpha:"

- Associated LI line #

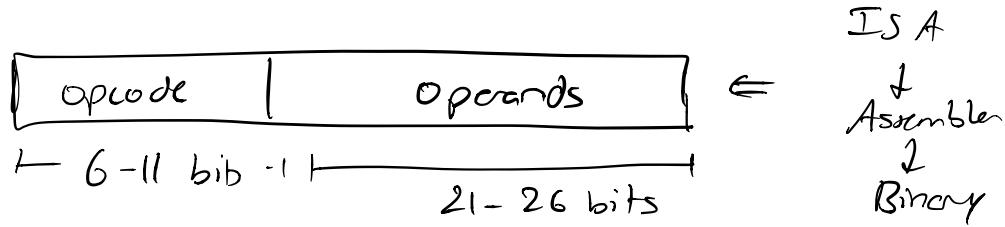
- Load + store from memory: LDR / SDR destReg, [addrReg, offset]



- Arithmetic: operation destRes, Res1, Res2 \Rightarrow destRes = Res1 op Res2
 - Branching: PC \rightarrow place in order to execute operation

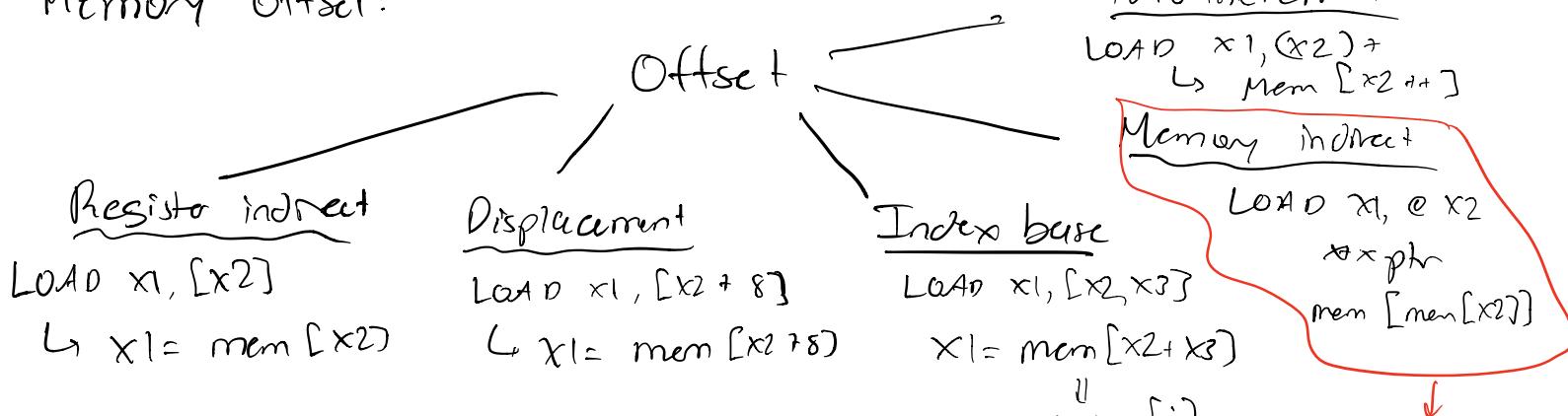
B ← label, offset

- Machine code counting:



- Variations: R, I, D, B, CB, IW

- ### - Memory offset:

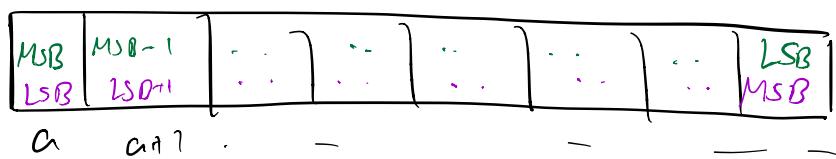


- ## o Autoincrement in ARM:

- Preindex: LDR x1, [x2, #8] ! \Rightarrow increment x2 after access
 - Post index: LDR x1, [x2], #8 \Rightarrow access x2 \rightarrow increment

- ## - Endianess:

Big Endianess, Small endianess



- Load-store architecture: register - register operations
 - Opposite x86: register - memory

ARM ASSEMBLY

- R format:

opcode	Rm	shamt	Rn	Rd	→ Sbit <u>dest. reg</u>
S bit, <u>2^w op.</u>	Logical shift, <u>shift</u>	S bit <u>L^S opnd</u>			

- Arithmetic or logical instructions

- Logical: LSL, LSR, AND/ANDI, ORR, XOR

□ Ex:// $LSL \quad X1, X2, 4 \Rightarrow X1 = X2 \ll 4$
 $= X2 * 2^4$

$AND \quad X1, X2, X3 \Rightarrow X1 = X2 \& X3$

□ Tip: Use AND w/ 1 in order to isolate certain bits

$$\begin{array}{r} 00110 \\ \& 00100 \\ \hline 00100 \end{array} \Rightarrow \text{MASK}$$

- D format:

opcode	address	op2	Rn	Rt
	9 bit			

- Memory access (Load/Store)

- Address is scaled by data size

- Aligned:

□ 32 bit \Rightarrow addr $\% 4 == 0$, 64 bit \Rightarrow addr $\% 8 == 0$

□ ARM supports aligned + unaligned

□ Aligned: LDR/SDR \Rightarrow offset is multiple of 8

□ Unaligned: LDUR / SDUR \Rightarrow $\neg 11$

- I format: immediate instructions (ADDI, SUBF)

opcode	immediate	Rn	Rd
	12 bit		

- CB format: branching

opcode	address	Rt
--------	---------	----

19 bit, PC-relative

- o ARM instructions: CBZ x1, label (go to label if $x1 == 0$), CBNZ x1, Label
- o Use inequality:

N (negative)	Z (zero)	V (overflow)	C (carry)
-----------------	-------------	-----------------	--------------

- ① Set flags in operation

ADDS, SUBS,

- ② Branch on conditions

B.EQ, B.GE, B.GT, B.NZ . . . Label

- o Ex:// ADDI W1, WZR, #3

ADDI W2, WZR, #2

ADDI W3, WZR, #0

SUBS W4, W1, W2

B.GT Label \Rightarrow IF $W1 > W2 \Rightarrow \text{label}$

- B format:

opcode	address
26 bit	Offset to PC

- o Unconditional branches to label (B address)

- o Address is in terms of words (32 bits)

▫ Represent 2^{26} words \Rightarrow 256 MB (+/- 128 MB)

▫ Farther than 128 MB \Rightarrow BR address (not an offset)

- IW format:

opcode	shamt	immediate	Rt
2-bit	16 bit	5 bits	

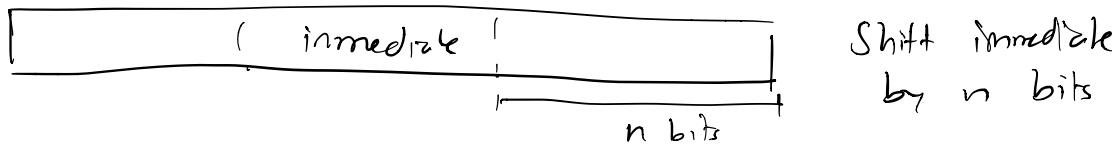
\hookrightarrow Bigger than I format

- o Used for MOV instructions

MOVZ x1, 2SS, LSL, 0

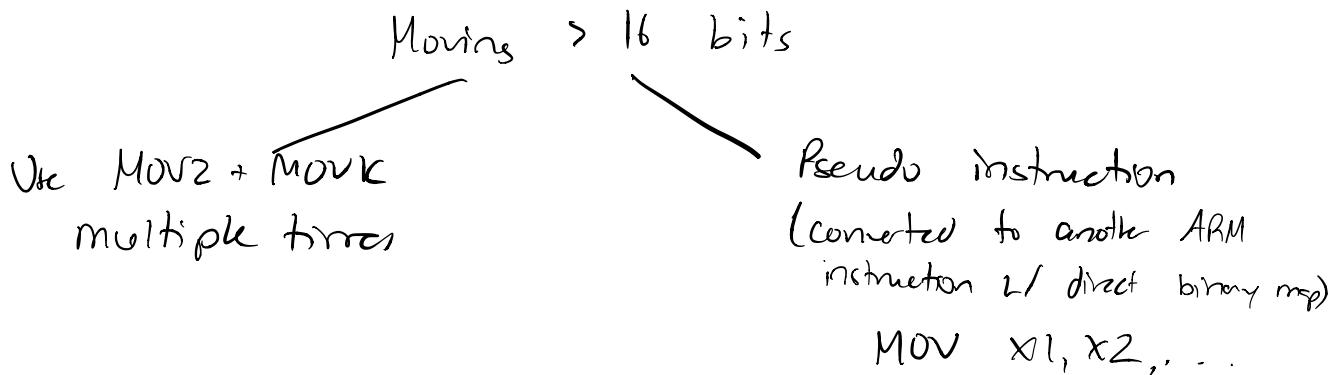
Translation: 2SS \rightarrow X1, zero out rest of the bits

- Shift amount: used to move immediate to particular location



Aligned: 0, 16, 32, 48

- MOVK: does not zero out rest of bits
- Moving more than 16 bits:



PROGRAMMING PARADIGMS

- Array access:

C code:

$$a = A[8] + b \Rightarrow$$

ARM: (assume addr of A is X3)

LDR X4, [X3, #64] \Rightarrow A[8] w
ADD X1, X2, X4 \Rightarrow add offset in bytes

C code:

$$A[12] = b + A[8] \Rightarrow$$

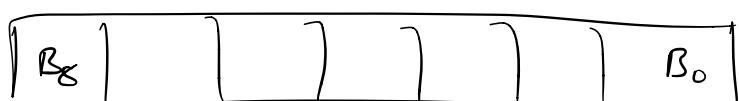
ARM:

LDR X4, [X3, #64]
ADD X1, X2, X4
STR X1, [X3, #96] \Rightarrow A[12]

- Assumption: each item in array was 8 bytes. Known beforehand

- Move/store bytes:

LDRB / STRB targetRes [addrRes, offset]



Operated on first

- Sign extend: dealing w/ bytes loses sign information

LDR SB, STR SB (doesn't matter if unsigned)

- If-then-else:

C code:

```
if (a < b) a++;
else b++;
```

ARM ($x_1 = a, x_2 = b$)

SUBS x_3, x_1, x_2

B. GE else

↳ Condition

ADDI $x_1, x_1, 1 \quad \left\{ \begin{array}{l} a \leq b \\ B \text{ join } \leftarrow \text{skip eln!} \end{array} \right.$

else: ADDI $x_2, x_2, 1 \Rightarrow a > b$

join: . . .

- For loops:

C code

```
long int A[100], sum, n, i;
for (i = 0; i < 10; i++) {
    sum += A[i]
```

ARM ($x_1 = i, x_2 = N, x_3 = \text{sum}[]$)

$x_4 = \text{sum}$

MOVZ $x_1, 0, LSL, 0$

loop: SUBS $x_5, x_1, x_2 \quad \left\{ \begin{array}{l} \text{cond} \\ \text{check} \end{array} \right.$

B. GE exit

i should access $\Rightarrow LSL x_6, x_1, 3$

array \Rightarrow byte access $\Rightarrow i * 8 \quad \left\{ \begin{array}{l} LDR x_7, [x_3, x_6] \\ ADD x_4, x_4, x_7 \end{array} \right.$

ADDI $x_1, x_1, 1 \quad \left\{ \begin{array}{l} \text{sum} = .. \\ \text{sum} = .. \end{array} \right.$

B loop

exit: . . .

- Cuse statement:

- Using if-else ARM paradigm is really inefficient

- Use branch table:

br-table: label-case-0

:
:
label-case-n

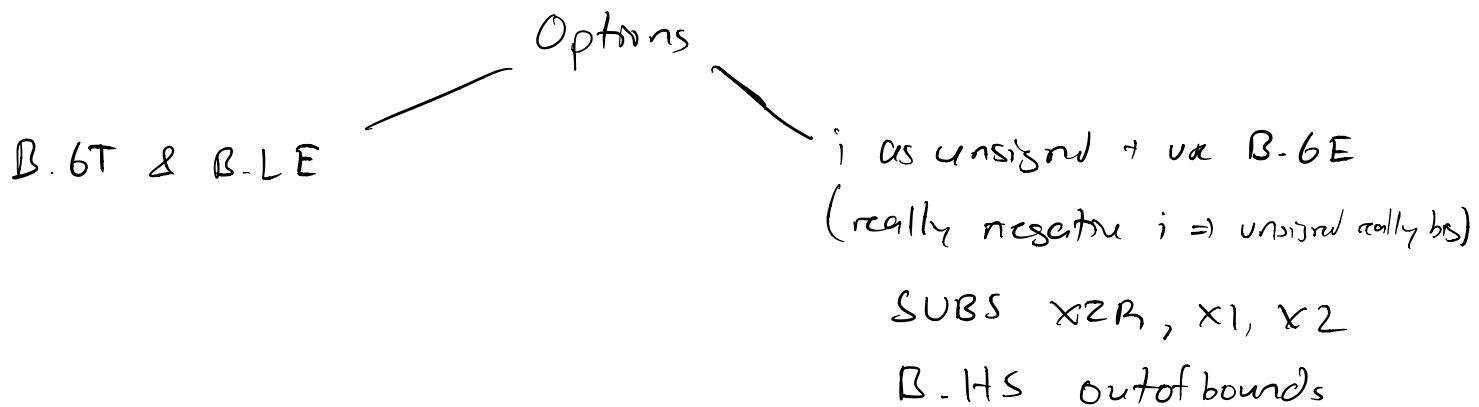
} Array of double words stored in memory

- Implementation:

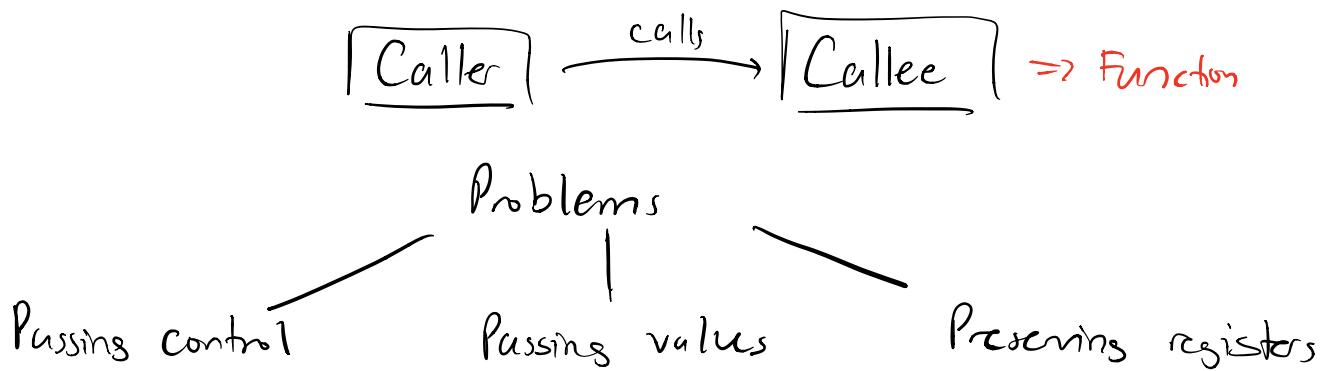
ADR x_0 , br-table \Rightarrow br-table absolut addr $\Rightarrow x_0$
LSL $x_1, x_1, 3 \Rightarrow$ Offset (set x_1 to be numeric
case, S, B...)
LDR $x_0, [x_0, x_1] \Rightarrow$ collecting label
BR $x_0 \Rightarrow$ Branching on label

- Bound checking:

- Out-of-bound: lower than min or higher than max



PROCEDURES



- Pussing control:

Main: - . ,
- - - ,

B1 fog

1. PC set to foo
 2. Previous PC stored in LR

foo: . . .

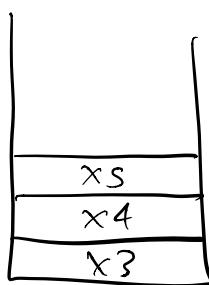
BR LR → (x30) tells program to go to stored PC

- Passing values:

- Soln: designate static special registers

- o Pass by value (src values) / pass by reference (src addr)

- Preserving registers:

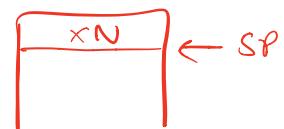


Stack: saving registers so we can restore later

- o Formula:

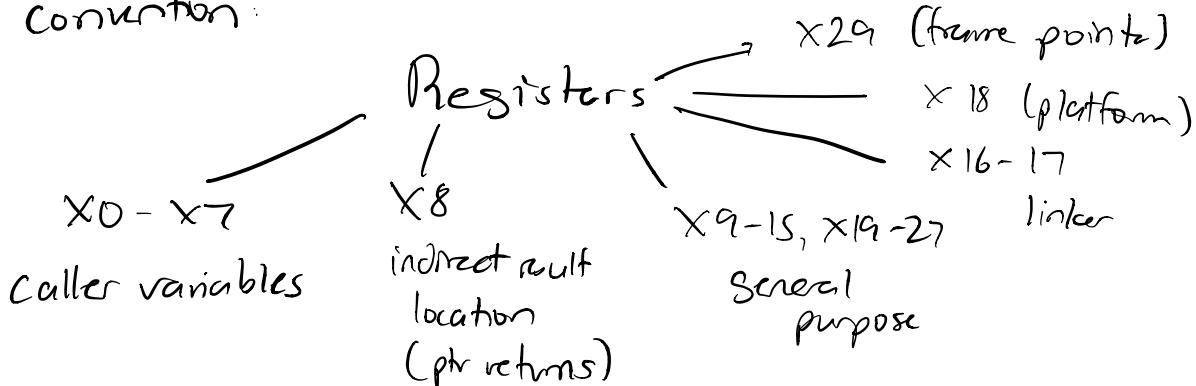
foo: $\text{SUBI SP, SP, } \#(8 \times n)$ } Decreasing stack ptr
 $\text{STR } x_N, [SP, \#0]$ } + storing

⋮
function work
⋮



$\text{LDR } x_N, [SP, \#0]$ } Loading + Increasing
 $\text{ADDI SP, SP, } \#(8 \times n)$ } stack ptr
 BR LR

- Register convention:



- Procedure of procedures:

1. Pass args to procedures: $x_0 - x_7$ registers / stack
2. Save caller-save registers + adjust SP
3. BL
4. Establish stack frame: $\text{SUBI, SP, SP, } \#(8 \times n)$
5. Save callee-save registers: variables that callee uses (e.g. LR)
6. Work

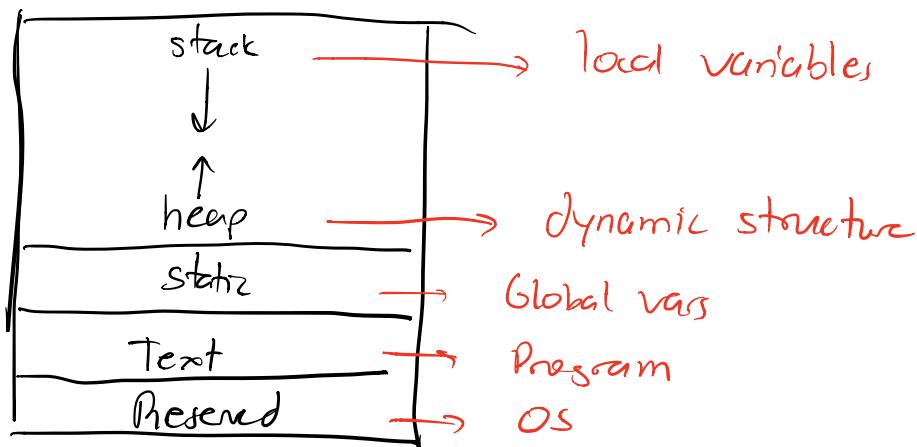
7. Put return values into $x0 - x7$

8. Restore callee-saved registers

9. Pop stack

10. Return to caller

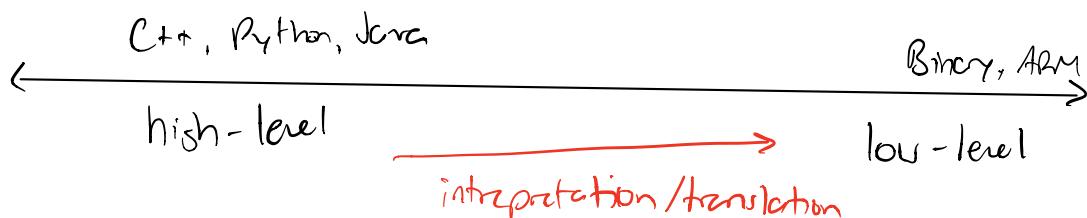
- Structure of memory



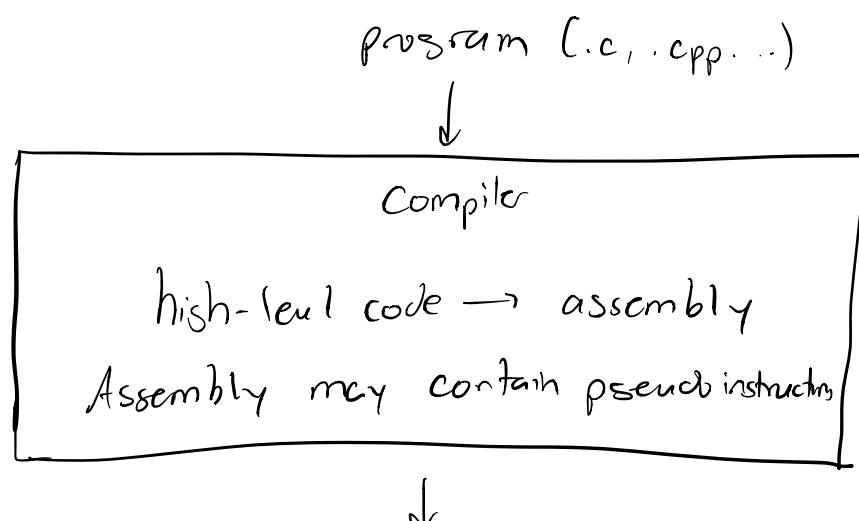
- Character data: use appropriate loading/storing based on bit encodings

° Ex:// 8 bit encoding = LDRB/STRB

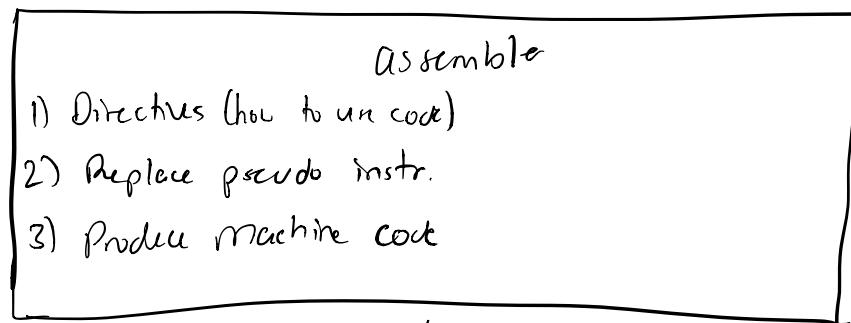
COMPILING, ASSEMBLING, LINKING, LOADING



- Interpretation is quite slow (es. Python), so we use compilers/translators



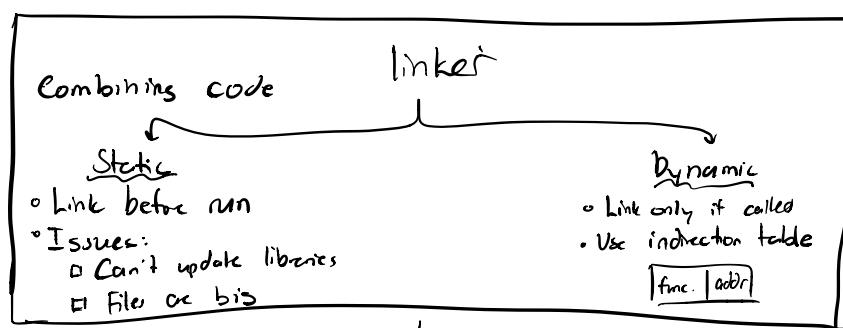
assembly program



⇒ 2 passes to build 2 tables
1. Symbol table (label)
2. Relocation table (unaccessible file)

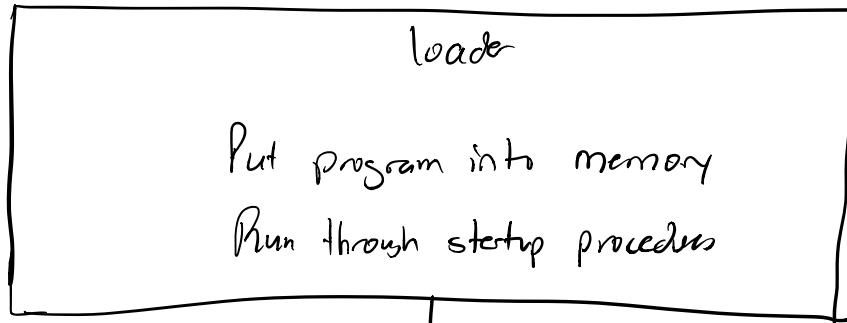
Object program (ELF)

↳ Header, text segment (machine code),
data segment (static data in binary), tables, debugging



← library object files

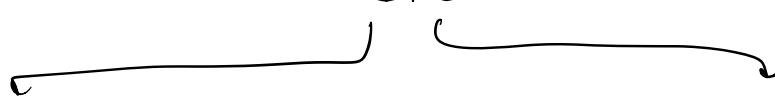
executable



memory

PROCESSOR DESIGN: DATAPATH I

CPU



Datapath

- Hardware to perform ops.

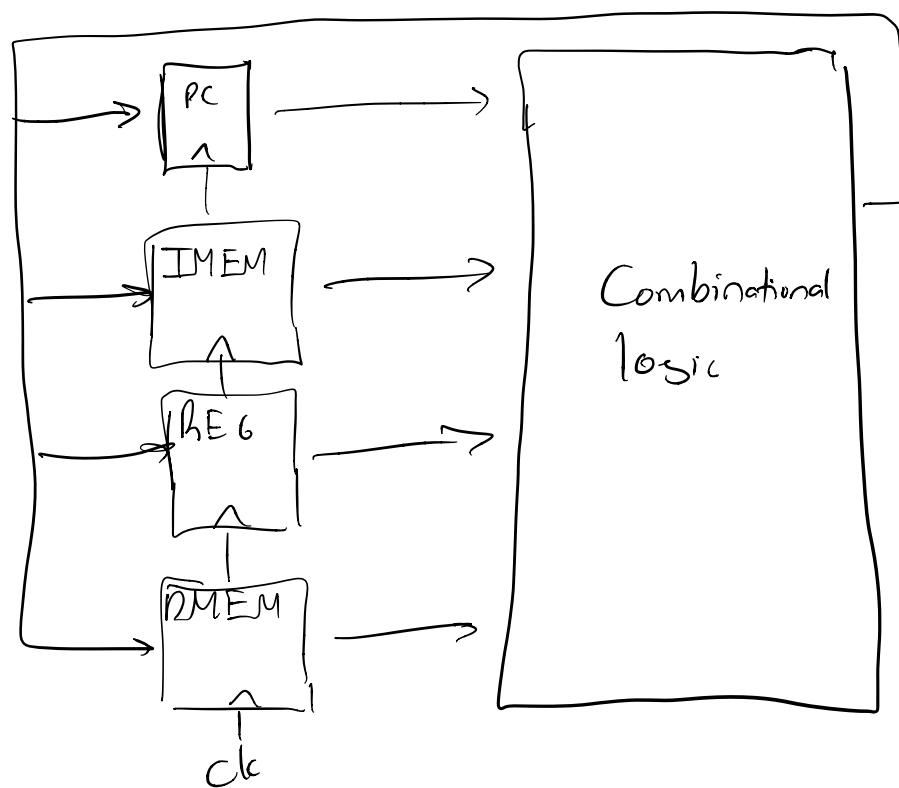
Control

- How datapath should work

- Designing processors:

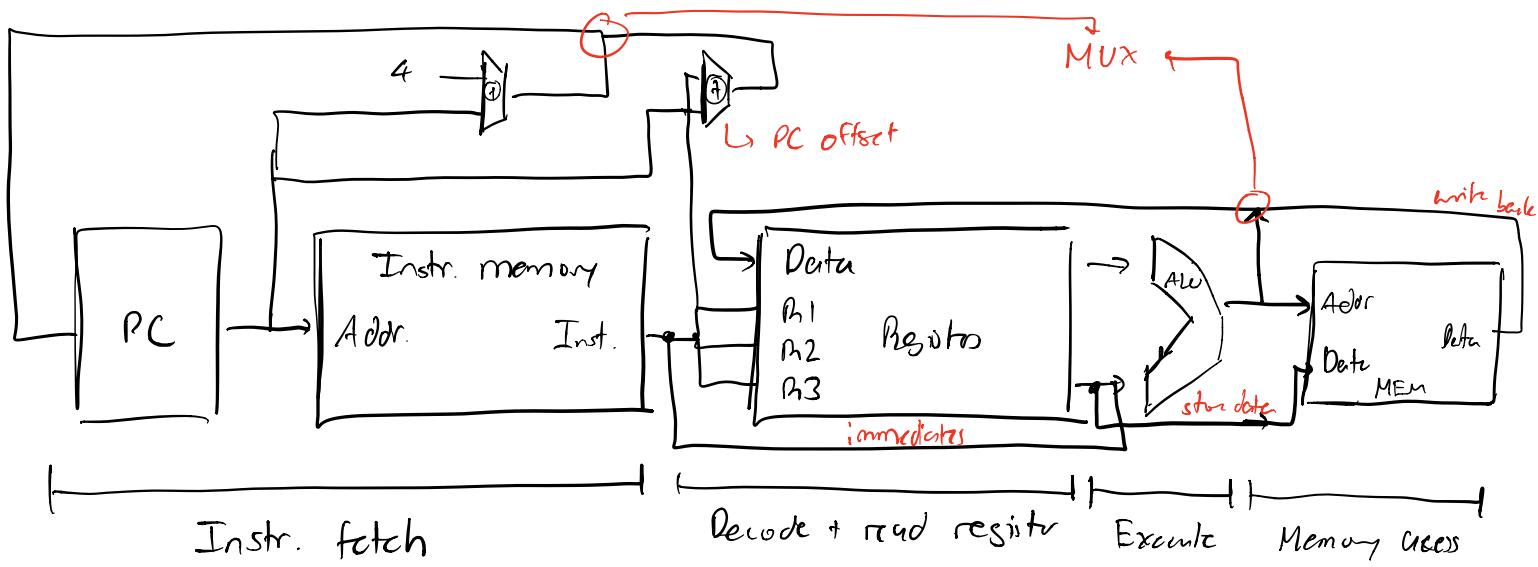
1. Create data path acc. bnd on ISA
 - o ISA model → register transfer model (RTL)
 2. Select components + choose clocking patterns
 3. Assemble components to meet RTL model
 4. Analyze ISA implementations + determine control points
 5. Assemble circuit logic for control points
- Processor design impacts CPI & clock cycle speed
- Single-cycle processor: 1 cycle/instruction \Rightarrow very low clock speed

Iterate until good



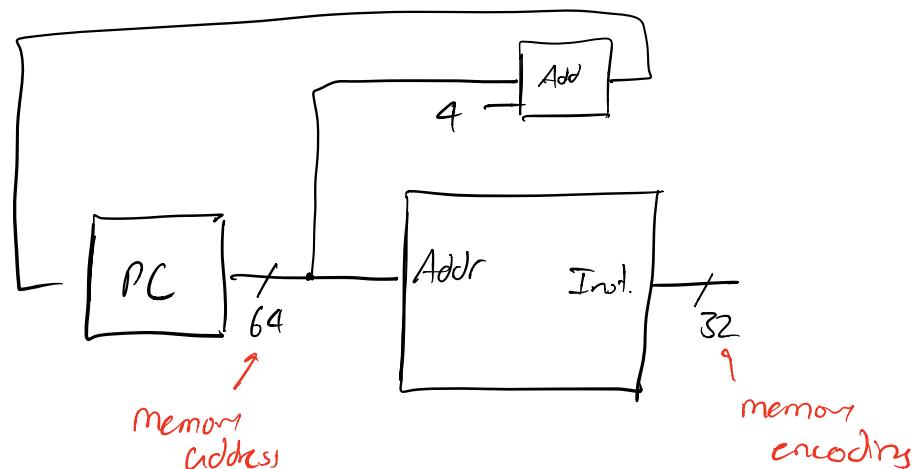
- Datapath stages:

1. Fetching instructions : puts instr from mem to reg., $PC += 8\dots$
2. Decoding || : read opcode \rightarrow extracted data \rightarrow reg.
3. Execute : logical/arithmetic \Rightarrow ALU. Load/stores: offset arithmetic
4. Memory access : load/store
5. Write to reg: results \rightarrow dest. reg (not done by load/stores)

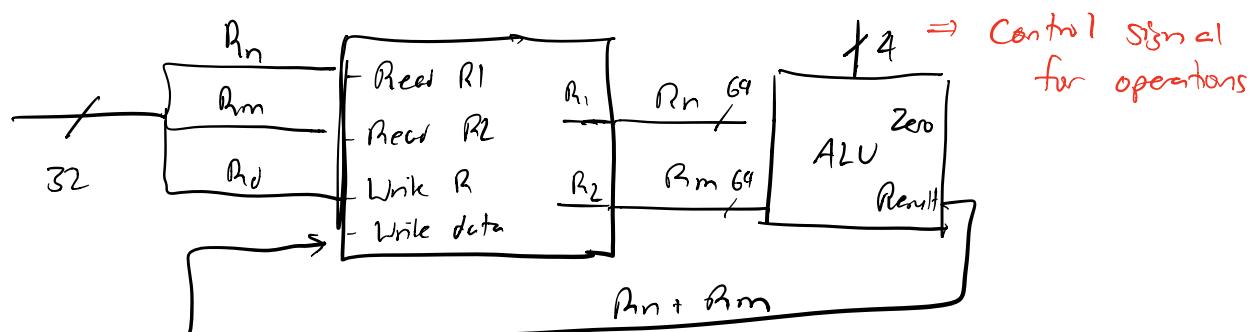


PROCESSOR DESIGN : DATAPATH-II

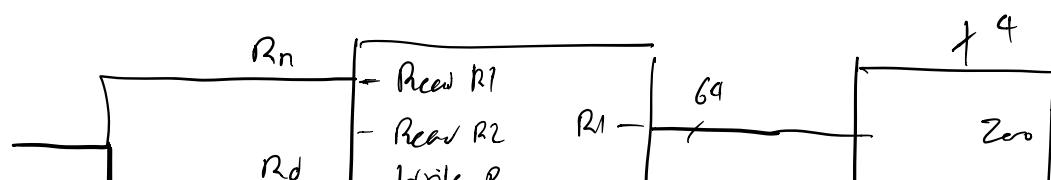
- Need processor to deal w/ 6 different inst. formats
- Fetch:

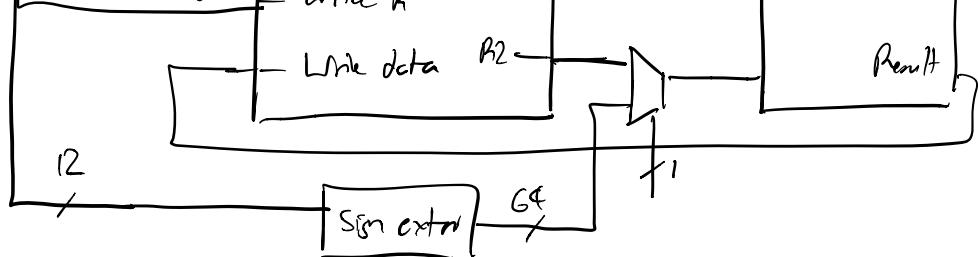


- ADD: R format instructions



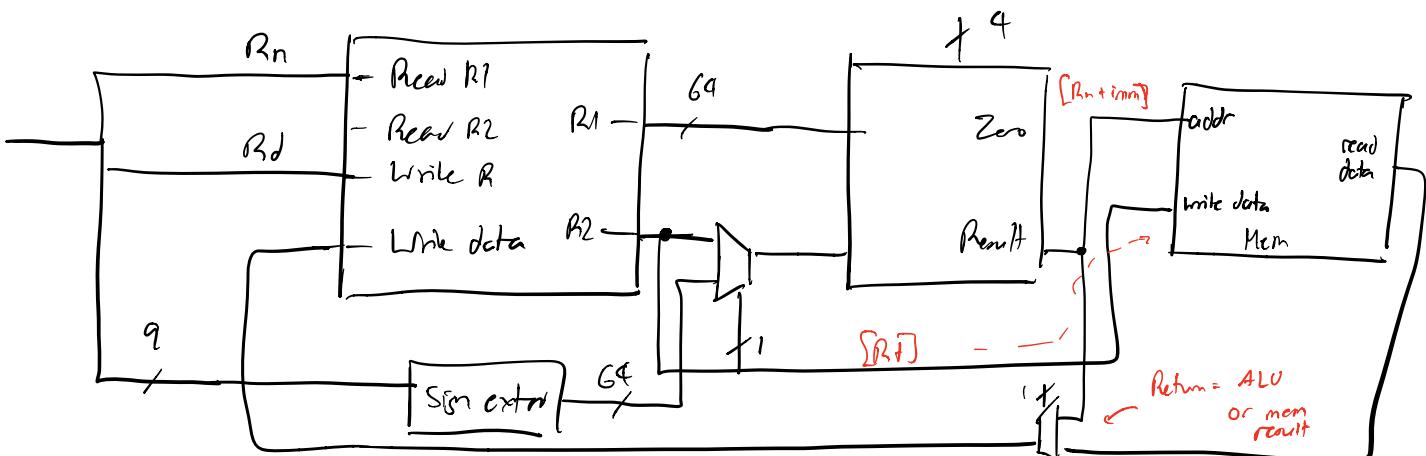
- ADDI : I format



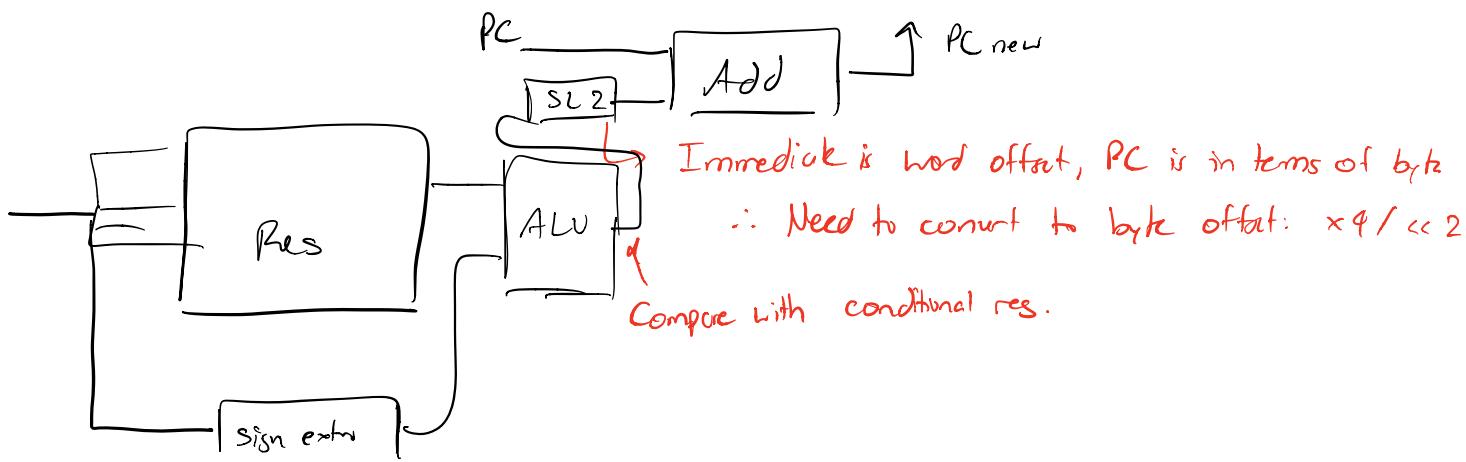


↳ We need to make immediate 12 bits → sign ext.

- LDR/STR: D format



- CBZ: B format



- Conditional flags: set in conditional registers ($R[31:26]$)

- If asked to do datapath for other instructions:

① Break down what hardware needs to do via datapath stages

② Connect on diagram

- Reasoning behind 5 stages:

1. Longest instruction can still function

2. Minimize amount of hazard

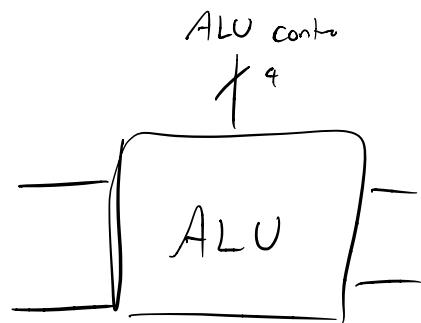
PROCESSOR DESIGN: CONTROL

Opcode $\xrightarrow{\text{comb. logic}}$ control signals

- Control signals:

- ALU control:

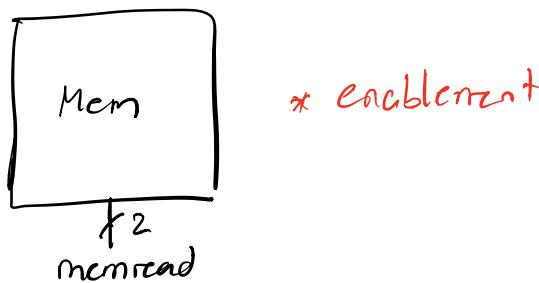
ALU control	function	ARM
0000	AND	AND
0001	OR	ORR
0010	add	ADD/ADDI
0110	sub	SUB/SUBS
0111	psub	1 reg. inputs (CBZ)
1100	NOP	ORN
:	:	:



▫ To generate ALU control: multi-level control gen.

ALU Opcode + inst. Opcode \Rightarrow ALU control \Rightarrow efficient

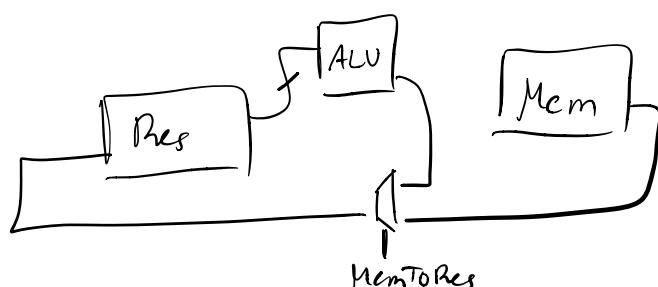
- Mem read: do we need to read from memory



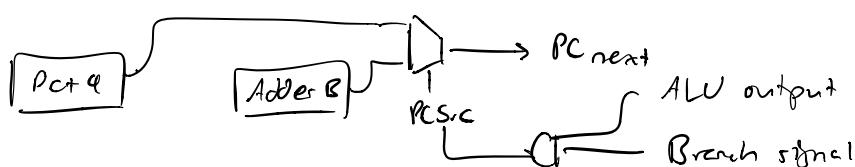
- Mem write: do we need to write date to memory



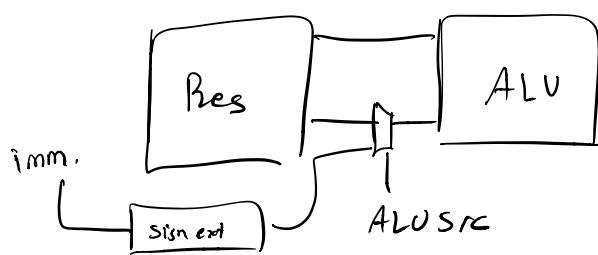
- MemToReg: do we need to write ALU (0) / Mem (1) date to reg?



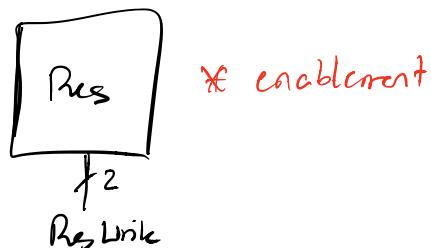
- PC Src: should we select normal, PC+4 (0) or branch target (1)



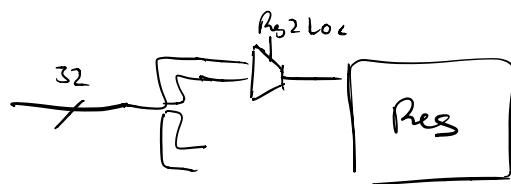
- o ALU Src: add w/ register (0) or w/ immediate (1)



- o Reg Write: should we allow write back to reg.



- o Reg2Loc: Select reg from R format [16-21] (0), or select D/CR [9-0], (1)



↓
x for LDH because
only reading 1 register!

- Strategy to determine control signals for operations

- ① Write down all signals (inc. branch signal)
- ② Think back to defn. of signal \Rightarrow whether needed or not
- ③ If signal is not needed \Rightarrow X
- ④ If it is \Rightarrow 0, 1, ...

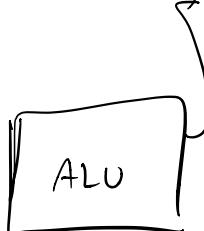
- o Ex:// CBZ format. Find control signals

Control Signals	Values
Reg2Loc	1
Branch	1
MemRead	0
Mem Write	0
Mem to Reg	X
ALUOp	0111
ALUSrc	0
Reg Write	0

- Dealing w/ unconditional branches:

- o No ALU input necessary



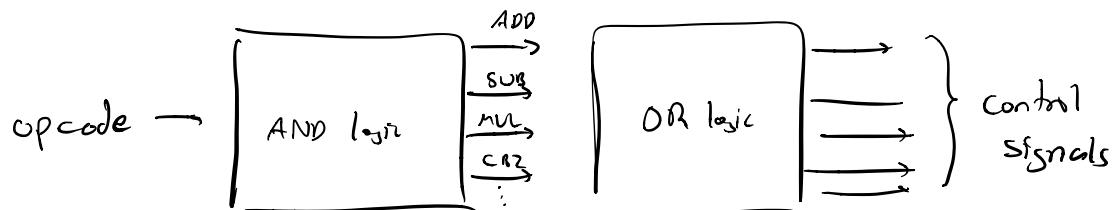


- Implementing control signals

1. Use ROM to store all control signals for all instruction possibilities

↳ Problem: unrealistic for large ISAs

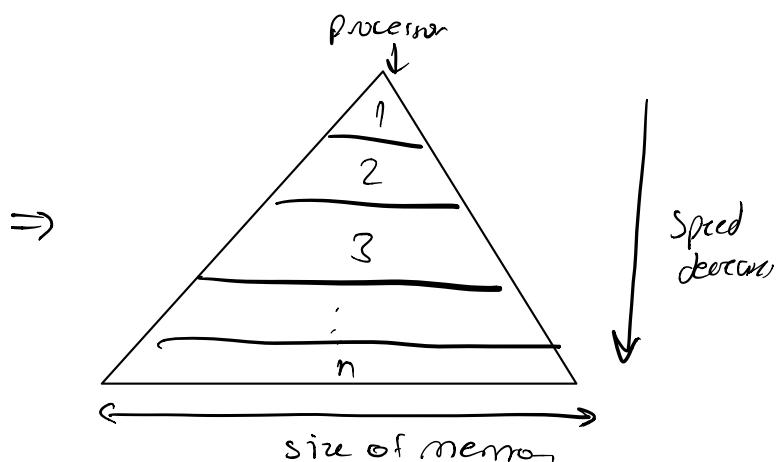
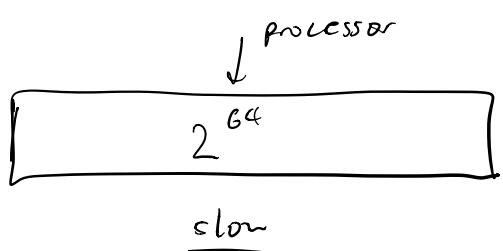
2. Combinational logic



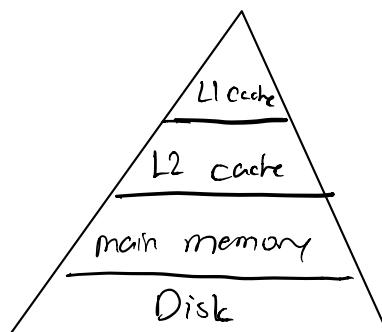
▫ Hidden efficiency: group similar inst \Rightarrow similar AND logic

- Overall problem: single-cycle is really slow!! Cycle time = time for longest instruction

MEMORY HIERARCHY



- Implementation:



- Reasons why memory hierarchy works:

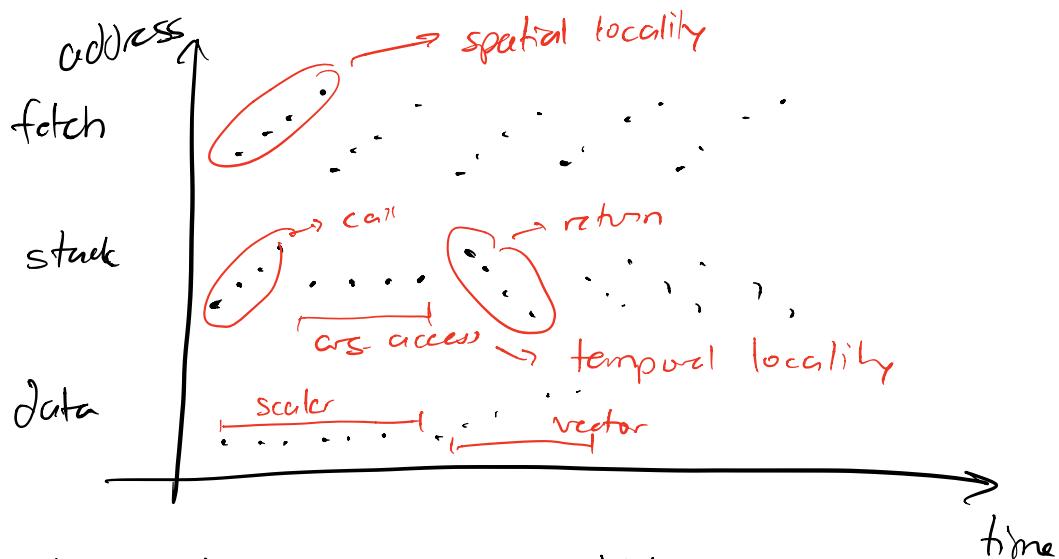
1. Temporal locality: recent exec. inst. more likely to be exec. again

▫ Can bring inst. into cache \Rightarrow use it again + again

2. Spatial locality: data near exec. inst. likely to be exec. as well

▫ Can bring set of local data/inst. into cache

- Patterns in program:



- Reason why spatial + temporal exist: 10/90 rule

- Instructions: 10% of static inst. account for 90% of dynamic inst.
- Data: 10% of variables account for 90% of accesses
 - Variables \Rightarrow global.

- Hit: access mem. level + find data

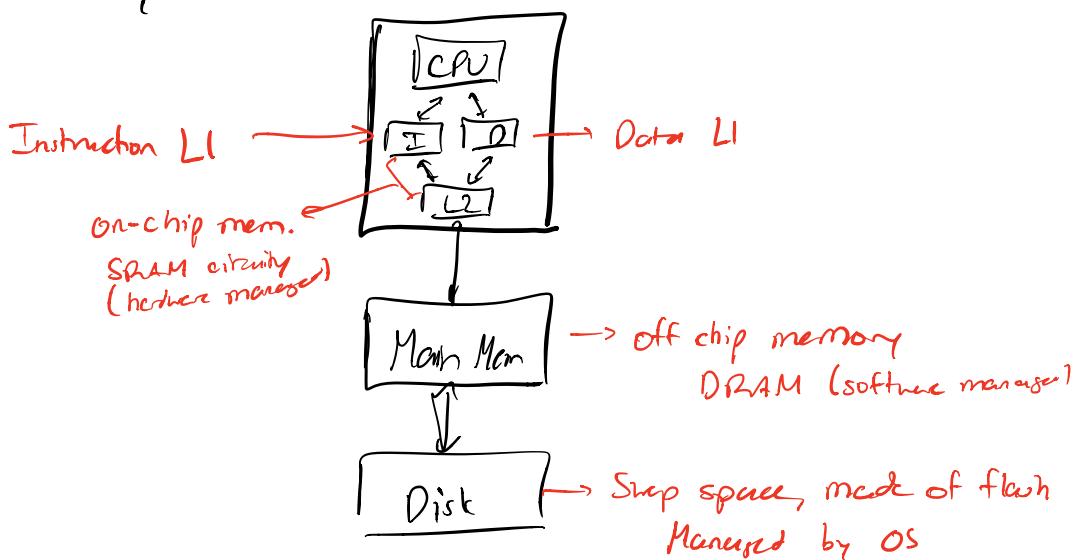
- Miss: access mem. level + don't find data \Rightarrow move to next mem. level

◦ Expensive: requires you to bring additional data back to top level

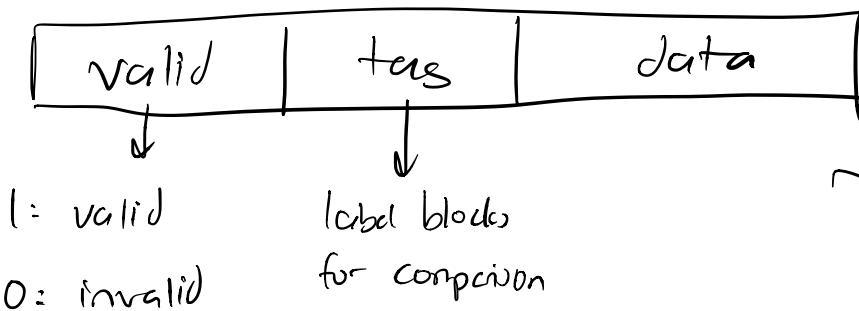
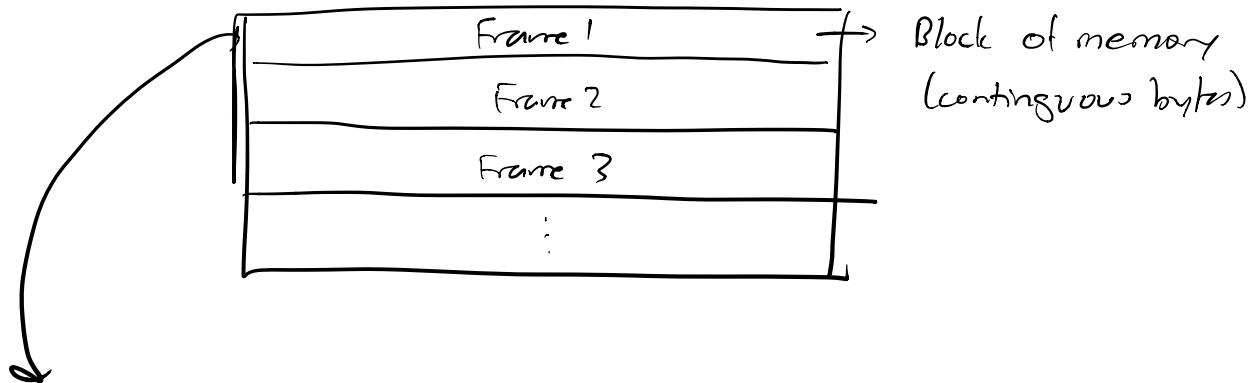
- Goal of memory arch:

1. High hit rate / minimize miss rate

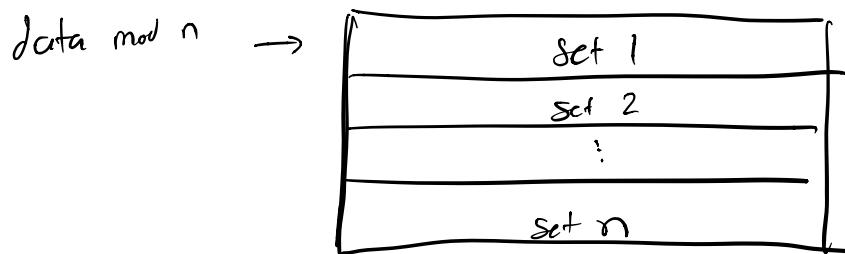
2. Low latency



CACHE ORGANIZATION

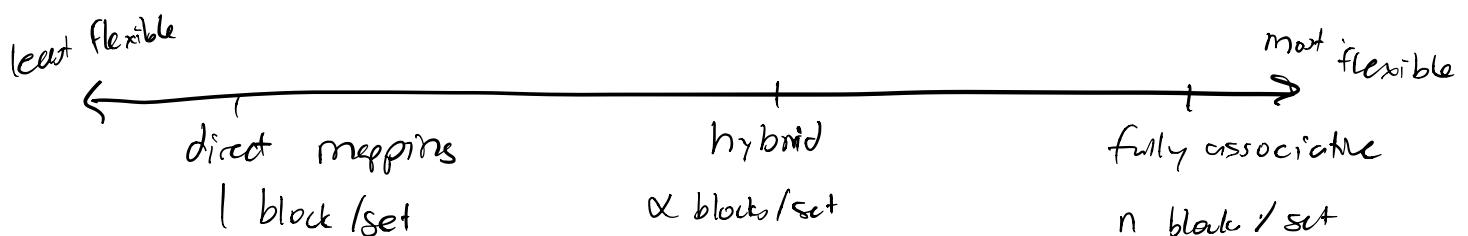
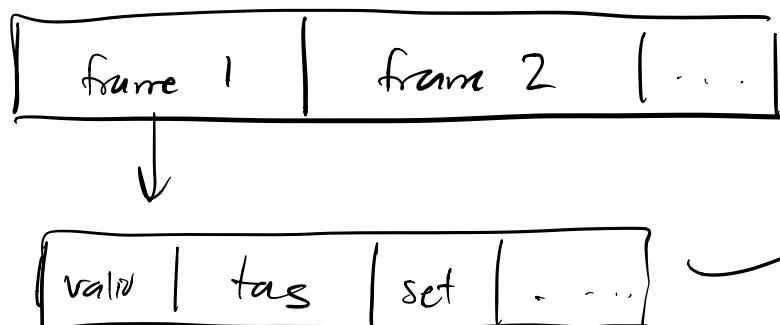


- Like to group frames into sets



Bits of the block act as the valid, tag, set tags

- Each set has a number of frames



- If data maps to already filled set \Rightarrow replace existing data
- Ex:// 1. Computer turns on

valid	tags	block
0	[0-31]	junk
:	:	:
0	:	:

2. CPU asks for word at address [32-38]

1. Which block does this belong to?

Because aligned \Rightarrow [32-63] block has data

2. Check cache: miss

3. Pull from lower level of memory

valid	tags	block
1	[32-63]	...
:	:	:
0	:	:

3. CPU asks for word at address [32-38] again

1. Block? [32-63]

2. Check cache: hit

4. CPU asks for word at address [36-39]

1. Block? [32-63]

2. Check cache: hit
(8B block)

- Ex:// Direct mapped, 8 set cache. Stream of 1 byte requests:

2, 11, S, 50, 67, S1, 3

Sequence of hits \Rightarrow misses.

- 1. 2 → block [0-7] → miss (set 0)
- 3. 11 → block [8-15] → miss (set 1)
- 4. 8 → block [0-7] → hit (set 0)
- 5. 50 → block [48-55] → miss (set 6)
- 6. 67 → block [64-71] → miss (set 0) \Rightarrow 67 is not in set 0 !! set 0 repl.
- 7. 51 → block [48-55] → hit (set 6)
- 8. 3 → block [0-7] → miss (set 0)

◦ To solve: n-way associative set (n frames / set)

- Mappings address to sets:

- Block size = N bytes. Need $\log_2 N$ bits to identify bytes
- Set = S sets. Need $\log_2 S$ bits to identify set
- Rest of bits: tags



32 - index - set

◦ Ex:// Address 10:

0...01010	set index	}
tag	offset	

dependent on architecture

↳ 2nd byte access

- If mapping to set + miss, which frame to replace?

- LRU (least recently used): find block not used recently + replace
 - Expensive to implement
- NMRU (not most recently used): any block not used replaced (LRU for 2-set associativity)
- Belady's algo: optimal, but need to know future access

CACHE ORGANIZATION II

Cache design

Associativity
(# of frames / set)

Block Size
amount of data
brought into cache

Capacity
amount of data
that cache can hold

Rewards for cache misses

Compulsory (cold) miss

- Cache has never seen data

Capacity miss

- Seen by cache but kicked \Leftarrow cache too small
- If N blocks stored by cache:
$$\text{access}, \dots, \dots, \dots, \text{access}$$
$$\underbrace{\quad\quad\quad}_{N \text{ distinct blocks}}$$

Conflict miss

- Not compulsory/capacity
- Cache associativity is low

- Problem solving strat: classifying a miss

1. Draw out cache : size of cache, block size

2. Fill in w/ stream of data

3. Use properties to figure out miss

- Ex:// 8B blocks. Stream: 2, 11, 5, 50, 67, 125, 256, 512, 1024, 2.
Last access a conflict / capacity miss?

2: [0 - 87]

11: [9 - 167]

5: [0 - 87] * kicked it out $\xrightarrow{\text{multiple frames/ set}}$

50: [49 - 557] $\xrightarrow{\text{1 frame/ set}}$

Only 7
distinct blocks
 \Downarrow
Conflict miss.

- Minor optimizations:

- Associativity \uparrow : \downarrow conflict misses, hit latency \uparrow
- Block size \uparrow : compulsory misses \downarrow , capacity misses \downarrow , conflict misses \uparrow
 - Reasoning behind conflict misses \uparrow : cache size is sum \Rightarrow set # \downarrow
- Capacity \uparrow : capacity misses \downarrow , hit latency \uparrow

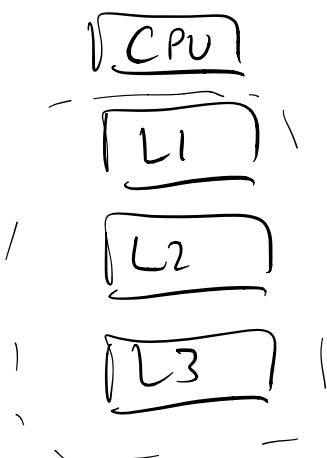
- Performance metric:

- Average memory access time:

$$AMAT = t_{hit} + \frac{\# \text{ of misses}}{\# \text{ of accesses}} \times t_{miss}$$

\hookrightarrow time to replace block

- For multiple caches:



Abstract:

$$t_{miss-L_n} = AMAT_{L_{n+1}}$$

$$AMAT = t_{hit} + \%_{miss} \times t_{miss}$$

$$= t_{hit-L_1} + \%_{miss-L_1} \times t_{miss-L_1}$$

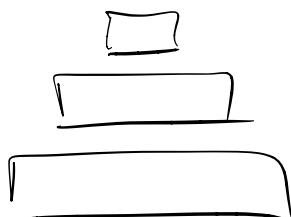
recursion!

$$= t_{hit-L_1} + \%_{miss-L_1} \times (t_{hit-L_2} + \%_{miss-L_2})$$

Hierarchy

Inclusive:

L_{n+1} is superset of L_n



(assume this is true)

Exclusive:

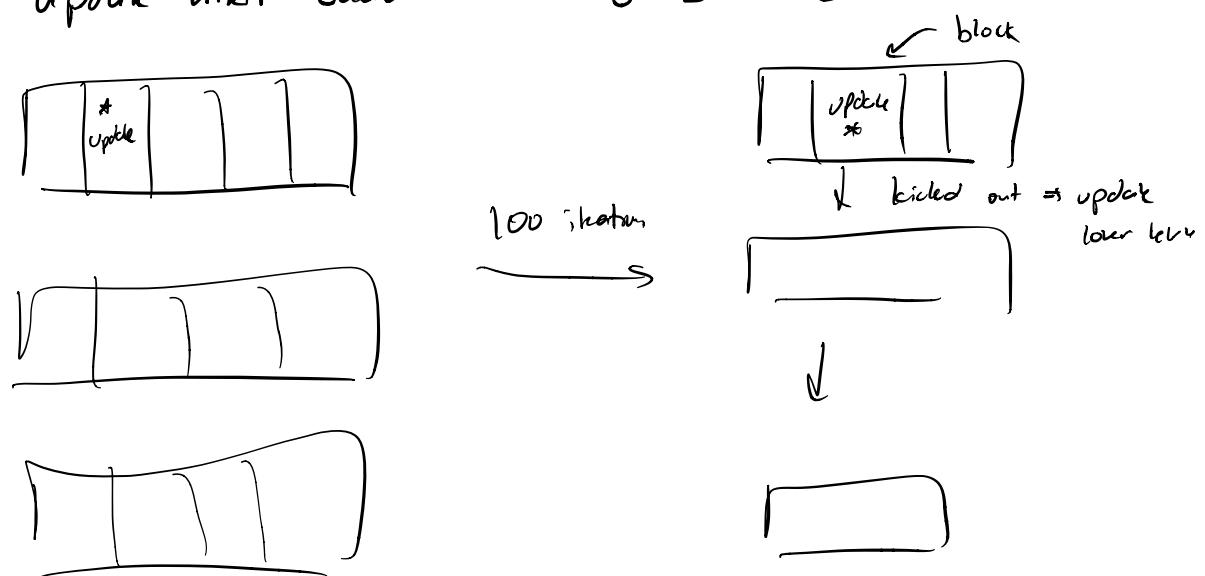
$L_{n+1} \cap L_n = \emptyset$



- How to change data in our cache (wrt hit)?
 - o Method #1: write-through
 - Immediately update lower level if new value
 - Pros: simple, miss latency is same, Cache coherent (cache have same data)
 - Con: time!

- o Method #2: write-back

- Only update when cache block is going to be kicked out



- Issue: cache needs to know if block replaced is a modified block
 - Dirty bit: tells cache if modified or not
 - Pros: minimum time
 - Con: non-uniform miss latency (completely dependent on dirty bit!)

- Write miss policies?

- o Method #1: Write-allocate

- Updates block in memory → update cache
- Works well with write-back

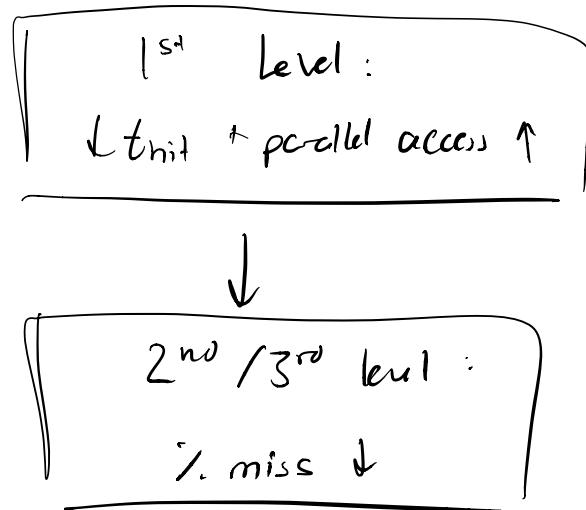
- o Method #1: no write-allocate

- Update block in memory → no cache update
- Works well with write-through

- Optimization for writing: write buffer

- If miss ⇒ write to buffer. Will execute later

- Optimization for diff. levels of hierarchy



- Performance metrics in general:

- Ex:// Stream of data: 20% stores, 80% loads.

L1: $t_{hit} = 1\text{ns}$, $\% \text{ miss} = 5\%$ (write through \rightarrow write buffer)

L2: $t_{init} = 10\text{ns}$, $\% \text{ miss} = 20\%$ (write back, 50% dirty)

Main: $t_{hit} = 50\text{ns}$, $\% \text{ miss} = 0\%$

1) What is t_{avg} w/out L2:

$$\begin{aligned}
 \text{AMAT} &= t_{avg} = t_{hit} + \% \text{ miss} \cdot t_{miss} \\
 &= t_{hit-L1} + (\% \text{ miss}_{load} + \% \text{ miss}_{store}) \cdot t_{miss} \\
 &= t_{init-L1} + (\% \text{ load} \cdot \% \text{ miss}) \cdot t_{miss} \\
 &= 1\text{ns} + 0.8 \cdot 0.05 \cdot t_{hit-main} \\
 &= 1\text{ns} + 0.8 \cdot 0.05 \cdot 50\text{ns}
 \end{aligned}$$

$0\% \text{ miss}$
 \rightarrow write buffer

2) What is t_{avg} w/ L2:

$$\begin{aligned}
 \text{AMAT} &= t_{avg} = t_{init} + \% \text{ miss} \cdot t_{miss} \\
 &= t_{init-L1} + (\% \text{ load} \cdot \% \text{ miss}) \cdot t_{AMAT-L2} \\
 &= t_{init-L1} + (\% \text{ load} \cdot \% \text{ miss}) \cdot (t_{hit-L2} + \% \text{ miss}_{L2} \cdot t_{miss-L2}) \\
 &= t_{init-L1} + (\% \text{ load} \cdot \% \text{ miss}) \cdot (t_{hit-L2} + \underbrace{(2 \cdot (\% \text{ miss} \cdot \% \text{ dirty} + \% \text{ miss}))}_{\text{dirty miss}} \cdot t_{hit-men}) \\
 &\quad + (\% \text{ load} \cdot \% \text{ miss} \cdot t_{hit-men})
 \end{aligned}$$

- Tag storage cost:

- Considered overhead (not included in capacity of cache)
- Ex:// 4 kB direct-mapped cache. 1024 4B frames.

4B frames: 2 bit offset (block)

$$\frac{4\text{KB}}{4\text{B}} = 1024 \text{ frames} \Rightarrow 1024 \text{ frames: } 10 \text{ bit offset (frame offset)}$$

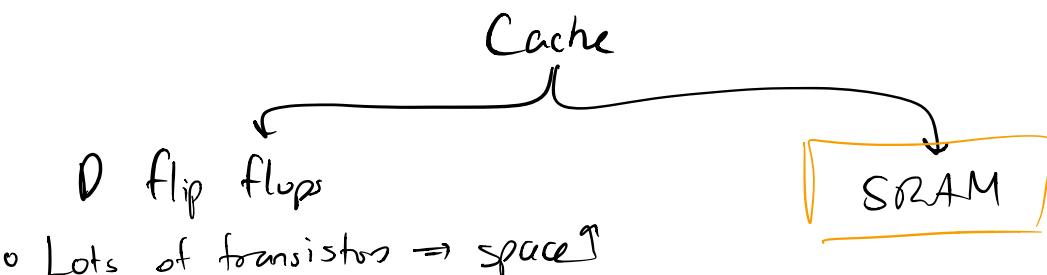
Assume 32-bit addr: 20 bits leftover \rightarrow tag

$$\text{For all frames: } 20 \text{ bits} \times 1024 \text{ frames} = 2.5 \text{ kB}$$

- 2 more optimizations:

- Victim buffer: used to reduce conflict misses
 - Fully associative cache, but much smaller
 - Block kicked out \rightarrow VIB. Miss \rightarrow check VIB
- Prefetching: reduce capacity + compulsory miss
 - Anticipate future misses via hardware logic / compiler
 - Next block prefetch: if miss on block X \rightarrow take X + (X + block-size)
 - Problems: more time + inaccurate (don't cover useful cache blocks).

CACHE IMPLEMENTATION



- Lots of transistors \Rightarrow space↑

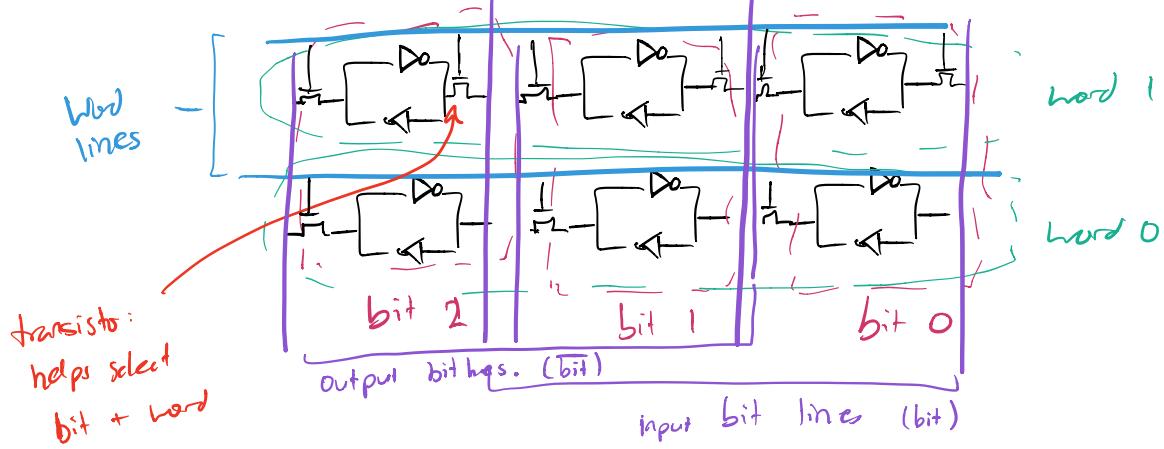
- SRAM: cross coupled inverter. 4-6 transistors / cell



- Static: bit maintains value, no refreshing

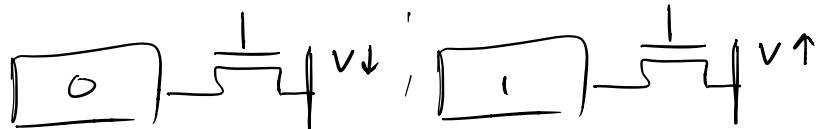
- Volatile: power gone \rightarrow memory gone

- Organization:



- Reading to SRAM:

1. Set interested bitlines to 1
2. Select interest word to 1
3. Cell will pull bitline up / down.
 - o Transistors use voltage differences.



4. Use sense amplifier to detect difference in voltage (too small to measure) between output \rightarrow input bitlines.

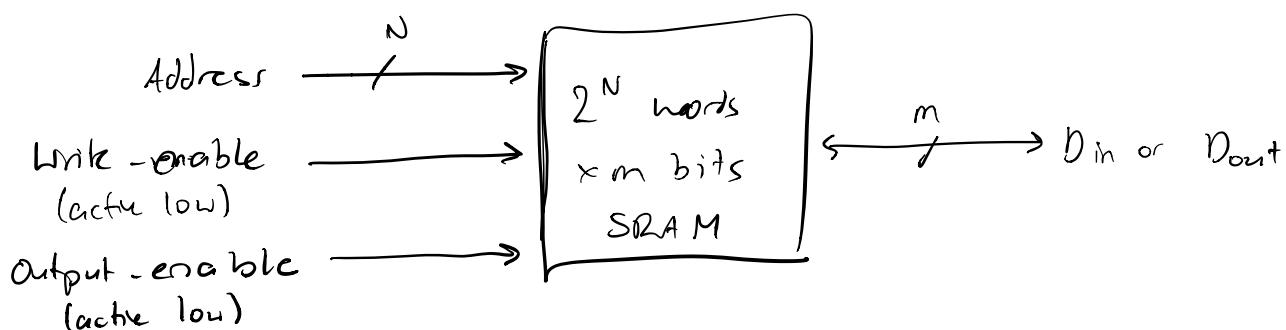
- Write to SRAM:

1. Drive bit lines
 - o Store 1: $\text{bit} = 1, \overline{\text{bit}} = 0 \}$
 - o Store 0: $\text{bit} = 0, \overline{\text{bit}} = 1 \}$

2. Select word via word line

- Length of bitlines \gg length of wordlines \leftarrow # of words \gg # of bits

- Logical diagram of SRAM:

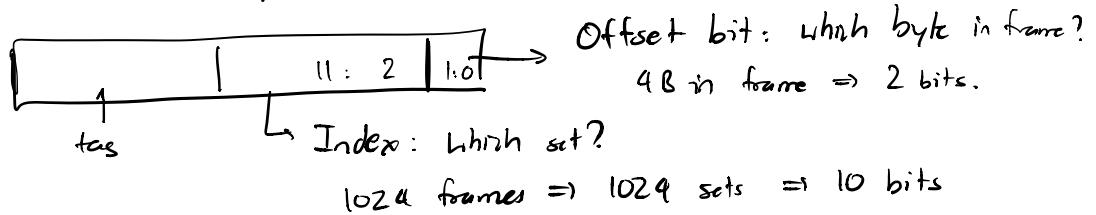


- WE \downarrow & OE $\uparrow \Rightarrow D = D_{in}$

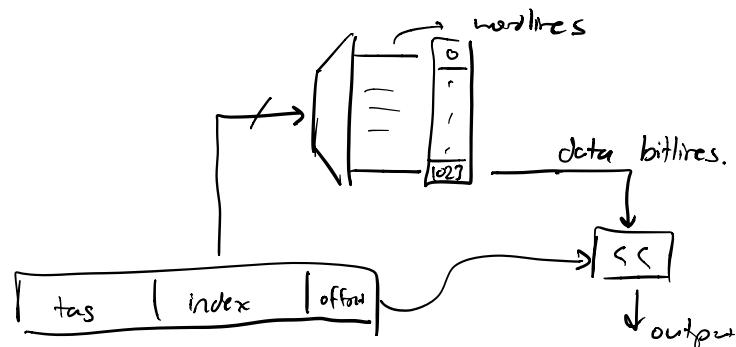
- WE \uparrow & OE $\downarrow \Rightarrow D = D_{out}$

- Ex:// 4KB cache w/ 1024 4B frames & direct mapped.

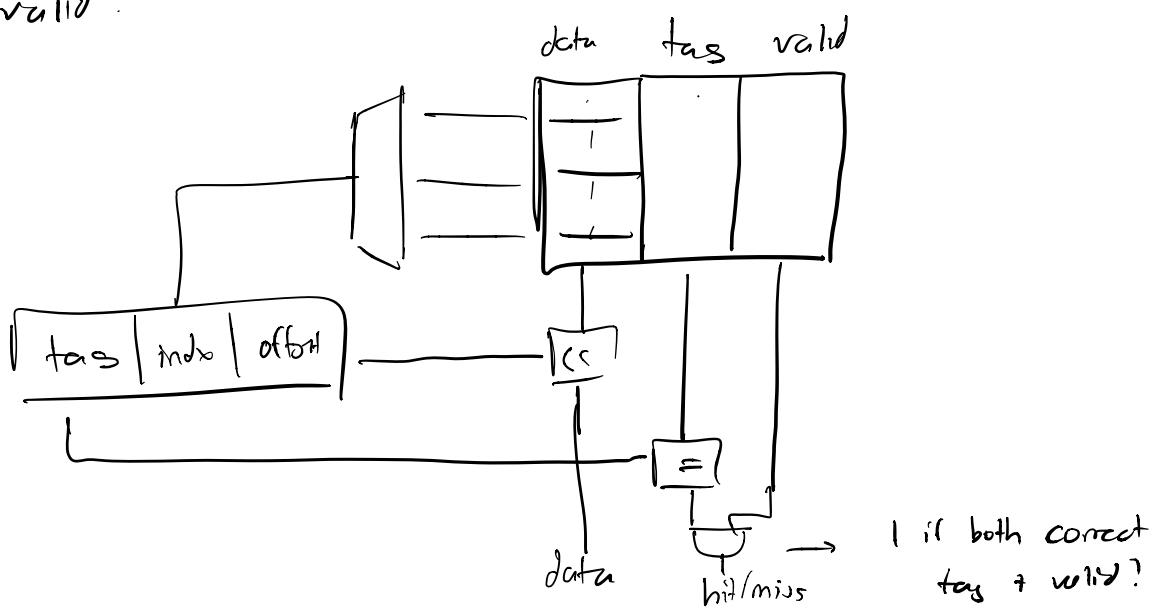
① Dissect input address given problem + 32-bit arch.



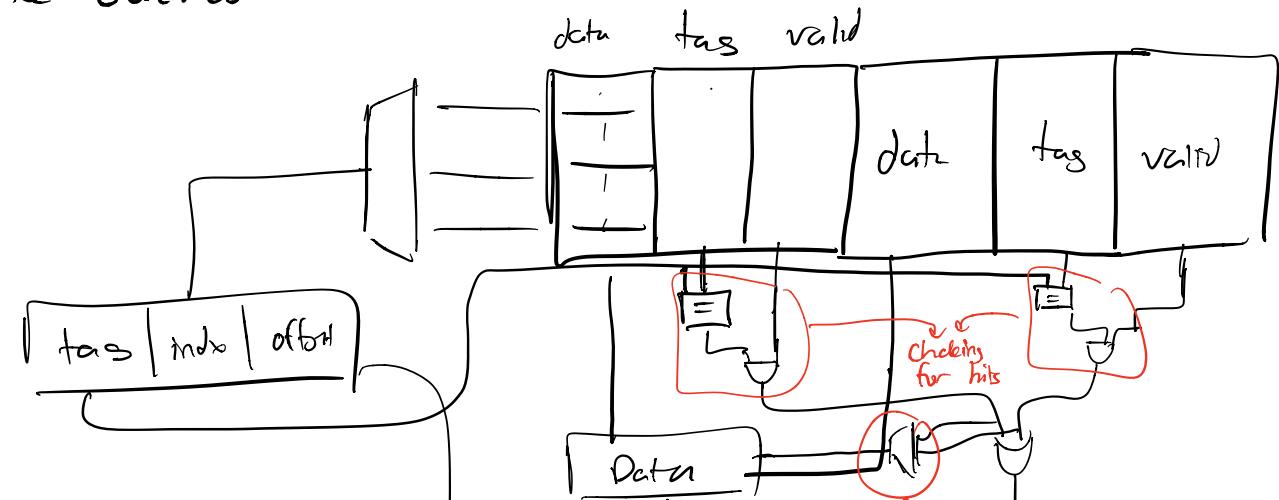
② How to take data from cache?



- Tag and valid?

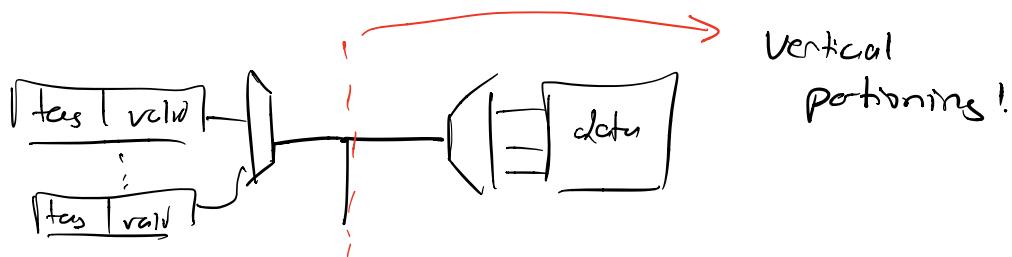


- Associative caches?

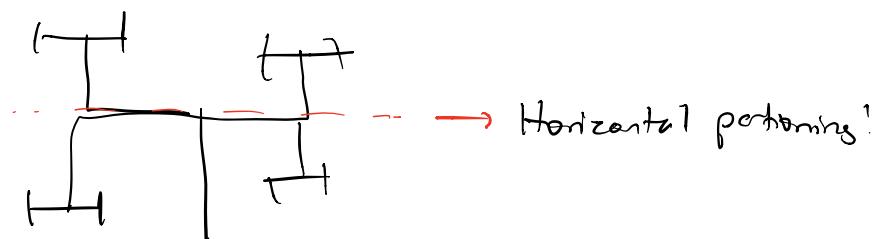




Physical implementation (conserve material)



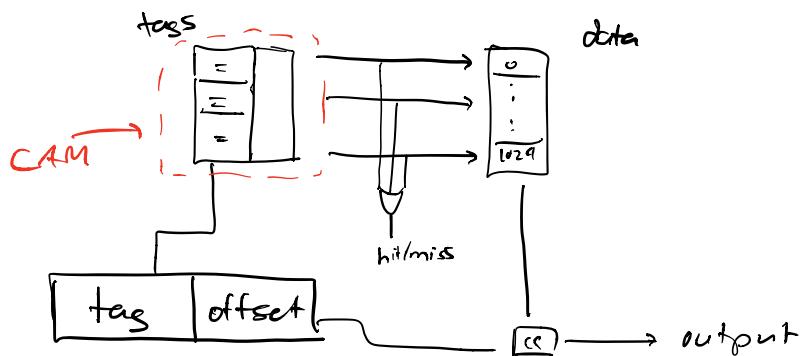
High level:



Fully associative cache:

- Problem: prov. solution \Rightarrow too many gates + muxes.
- CAM: content addressable memory. Expensive

Block \rightarrow CAM \rightarrow onehot encoding



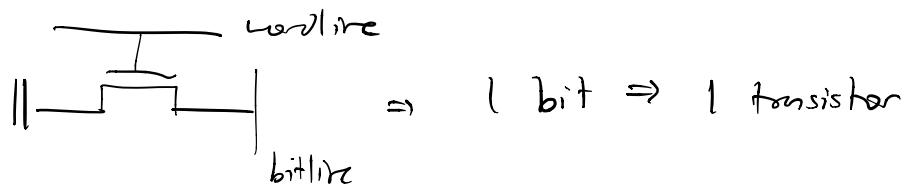
Stores:

- Match tag
- Write to matching set

} Parallel \rightarrow corruption of data

MAIN MEMORY IMPLEMENTATION

Optimizing for density \rightarrow DRAM. Sane architecture



- Decay over time \Rightarrow DRAM must be refreshed
- Access time \propto parts # $\times \sqrt{\# \text{ of bits}}$ (longer than SRAM)

- Read:

1. Precharge bitline to $V_{DD}/2$

2. Select word via wordlines

- 3.

4. Use sense amplifier to detect difference

- Write:

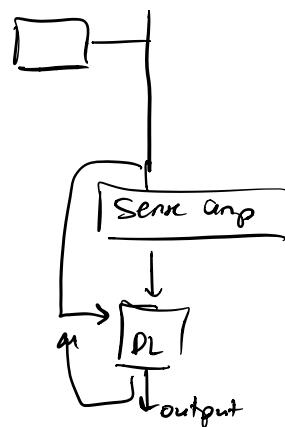
1. Select wordline + bit line

2. Pass in data \Rightarrow charging/discharging capacitor

- Problem: destructive read (every time we read \rightarrow data will change).

- Soln: D latch as buffer.

- Read:



1. Read

2. Write to buffer

3. Buffer will write back to cell
(borrowing)

- Write:

1. Read word to buffer

2. Write data into buffer

3. Write buffer into cell

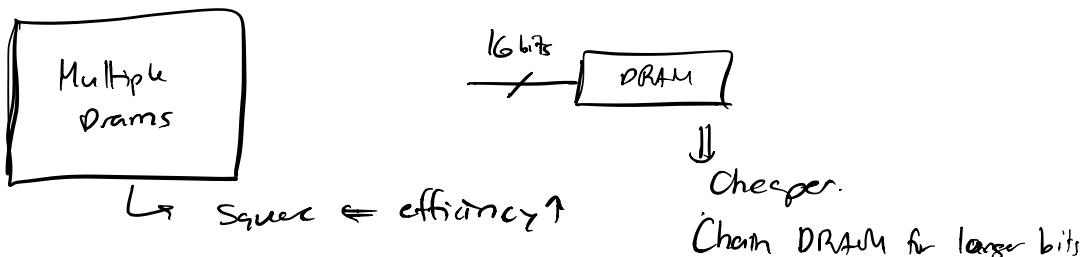
- Problem: leakage from capacitor

- Soln: refresh DRAM

for all words \rightarrow put into buffer \rightarrow write back to DRAM

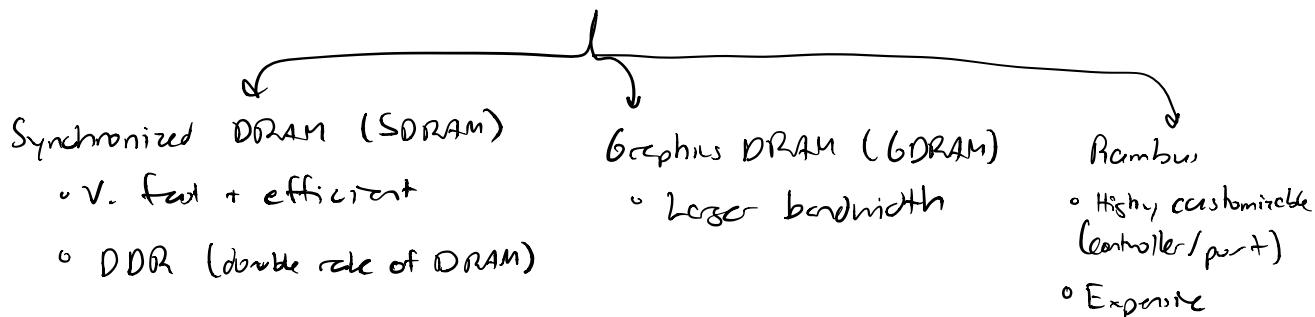
- Refresh interval \leftarrow capacity, tech., temp. (\uparrow temp \rightarrow more refresh)
- 1-2% of DRAM time spent on refreshing

- Design choices for DRAM:

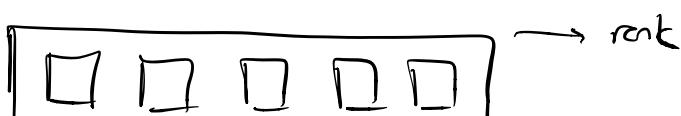


- Cycle time is $2 \times$ SRAM (refresh)
- Runs on diff. clock than processor \Rightarrow memory bus clock
 - If freq. \uparrow \Rightarrow memory bandwidth cannot handle

DRAM types

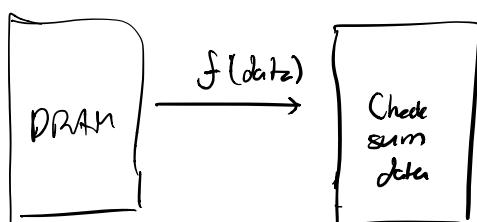


- Dual inline memory module: DDR



- Problem: errors

- Soln: checksum-style redundancy

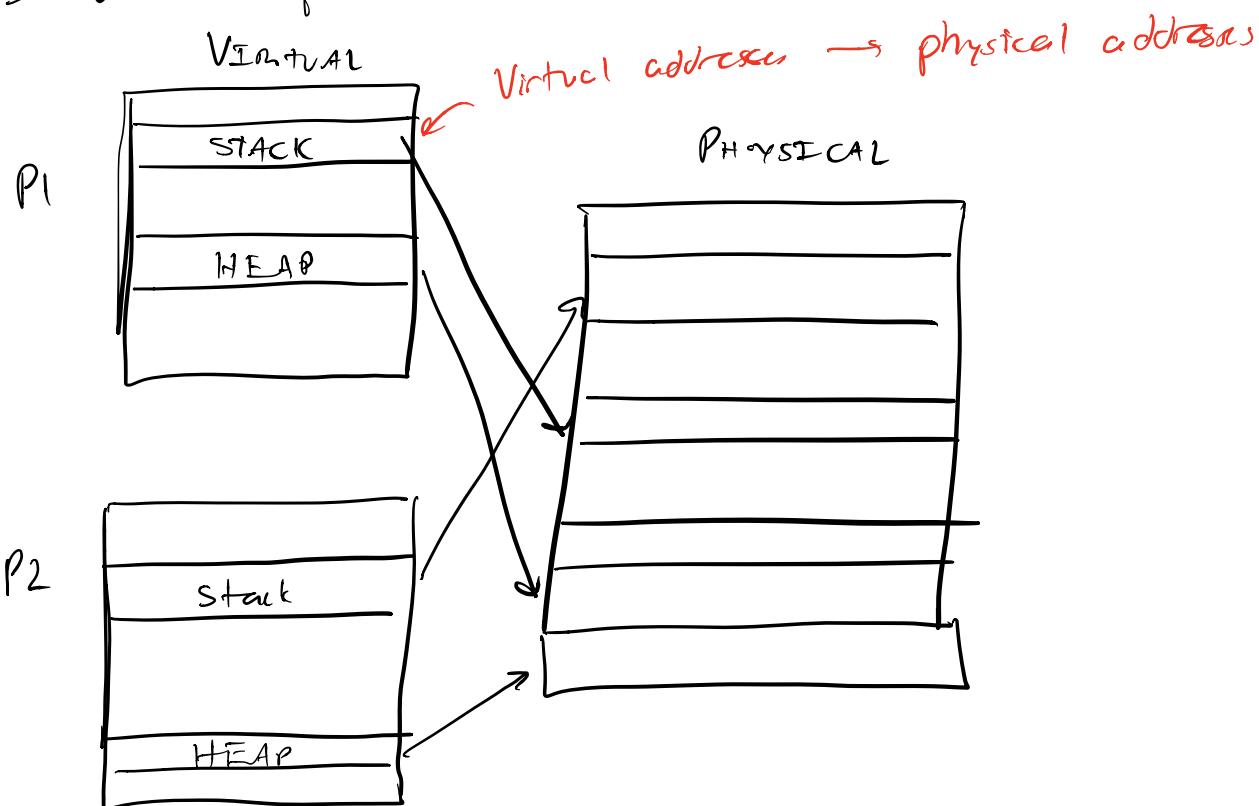


(On R/W \Rightarrow check if $f(data) = stored\ f(data)$)

$f(\dots)$: parity (XOR bits). Works if 1 bit flip, not for even # of flips

VIRTUAL MEMORY

- Memory is treated like cache \Rightarrow dynamic subset of total mem.
- Program thinks it has access to full memory \Rightarrow only has access to a couple of areas.



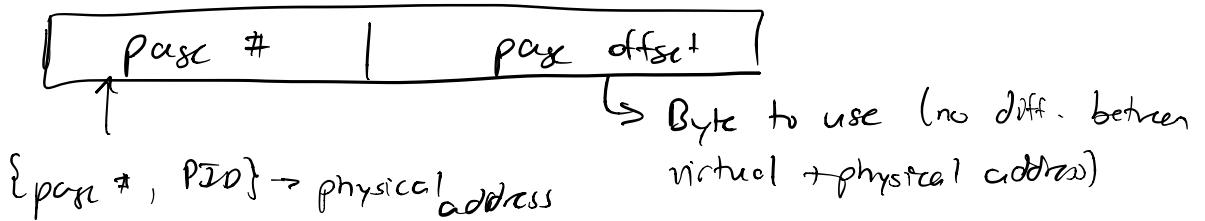
- Process:

1. Assign virtual memory + set data

2. OS + hardware: $\{ \text{Process ID, virtual address} \} \rightarrow \text{physical address}$

\hookrightarrow Virtual address can be same \Rightarrow need PID to distinguish

3. Translation:



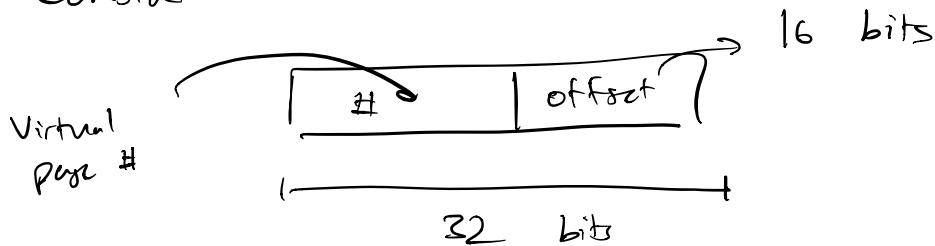
- Miss: page fault

1. Go to disk & get the data
2. Replaces page via LRU policy done by OS

- Finding # of page # bits:

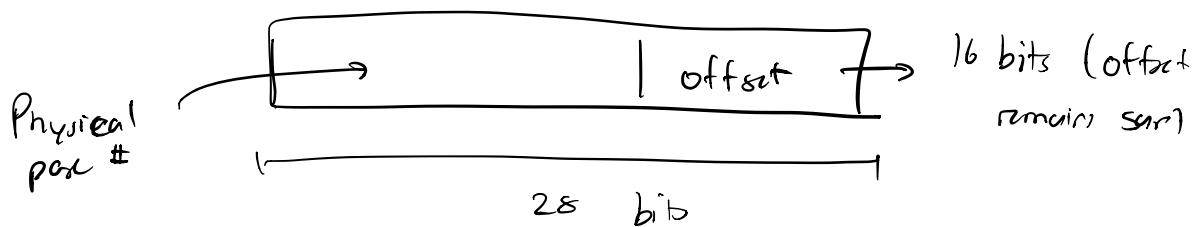
- Ex:// 64 KB page $\rightarrow 2^{16}$ bytes \rightarrow 16 bit offset for byte ID

Consider 32-bit arch:



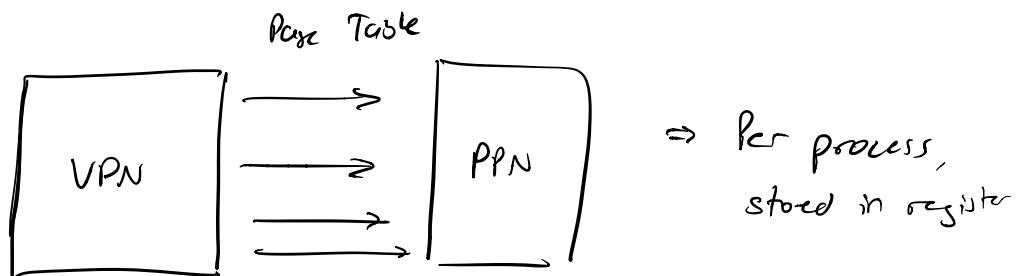
Consider 256 MB memory

To reference: 2^{28} bytes \Rightarrow 28 bits for physical address



Page table: 16 bit of VPN \rightarrow 12 bit of PPN

- How do we map between VPN + PPN?



- Ex:// 4B page table entries. 32-bit machine. 4KB pages. Page table size.

Q: How many entries in page table?

① Figure out entire memory capacity

32-bit VA \rightarrow 4GB memory (virtual)

② How many pages?

$$\frac{4\text{GB}}{4\text{KB}} = 1\text{M pages}$$

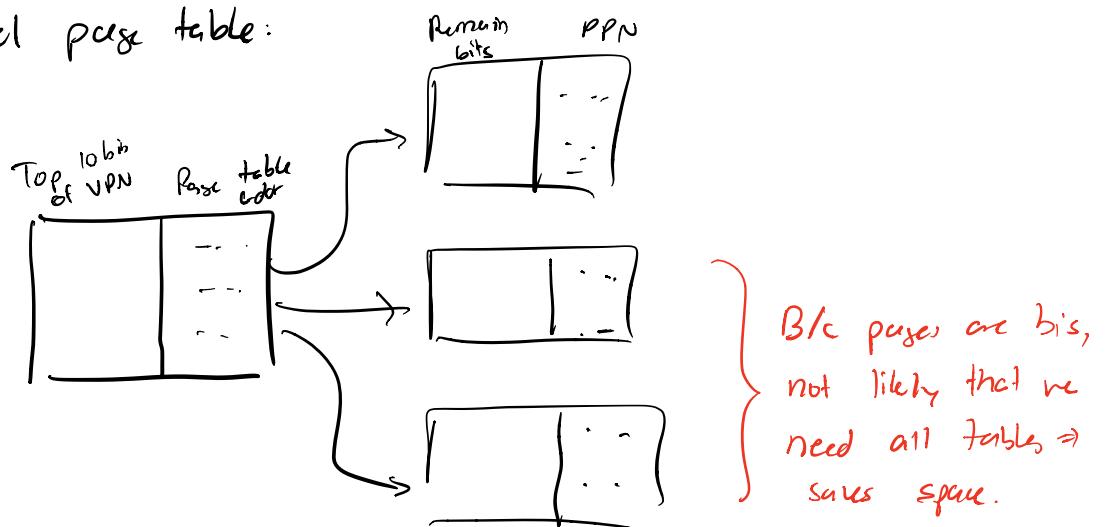
(?) How many bytes per page entry?

$$1 \text{ mil pages} \times \frac{4B}{page} = 4MB \text{ page table}$$

- To reduce page table size:

1. Inexact page size: unused memory \rightarrow more misses, longer time

2. Multilevel page table:



- Who does translation?

1. Process translates own address

- Page table would be in process virtual mem.
- Problem: process owner has to manage access \Rightarrow messy

2. OS translates:

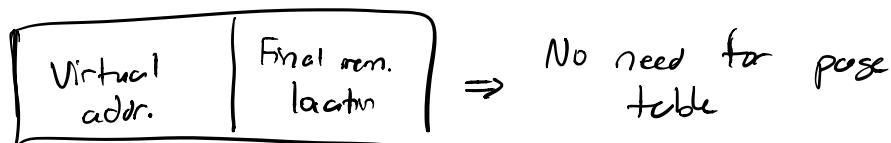
- Page table in OS virtual mem.
- OS can determine access + when to use (L2 miss)

- L2 miss \Rightarrow interrupt

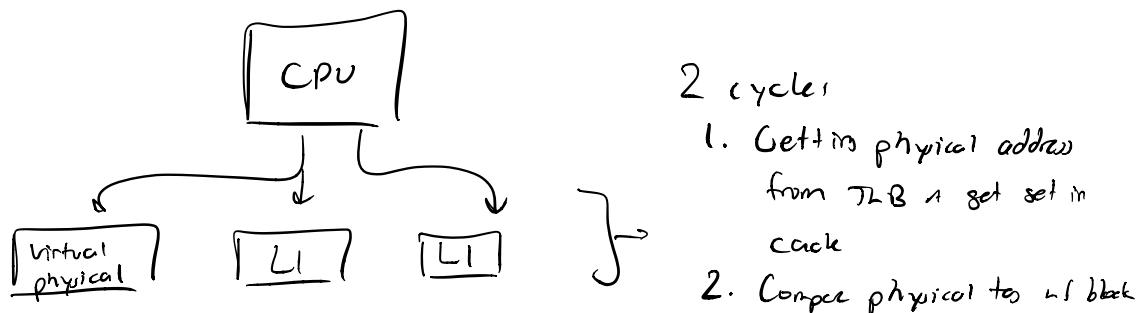
- OS takes control \Rightarrow VPN: PPN \Rightarrow return mem.

- Problem: takes time to get main mem. data

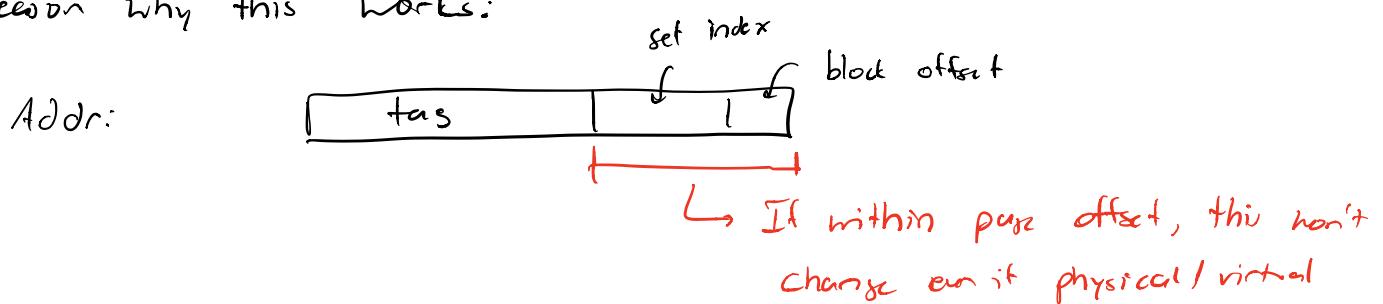
- Soln: translation lookaside buffer (TLB)



- o Miss: use page table to get mem. location + update TLB
- How is page table implemented?
 - 1) OS: flexible organization, slow
 - 2) Hardware: not flexible, fast
- Ex:// 32 KB cache. 64 B block size. 4 KB page size. # of entries if we want similar miss rates in cache + TLB?
 Similar miss rates \Rightarrow covering same amount of memory.
 Low-end: $32 \text{ KB} \div 4 \text{ KB page size} \Rightarrow 8 \text{ pages}$
 High end: $32 \text{ KB} \div 64 \text{ B blocks} = 512 \text{ blocks (512 pages)}$
 $\therefore \text{TLB entries} \in \{8, 512\}$
- Problem: Blk virtual addresses used to access cache \Rightarrow referring physical address/block in cache!
- o Soln: convert virtual address to physical at start (slow, but saves pain)



- o Reason why this works:



Constraint \Rightarrow more associative sets (less set bits)
 smaller blocks (less block offset)
 bigger pages (increased page offset)

PIPELINING

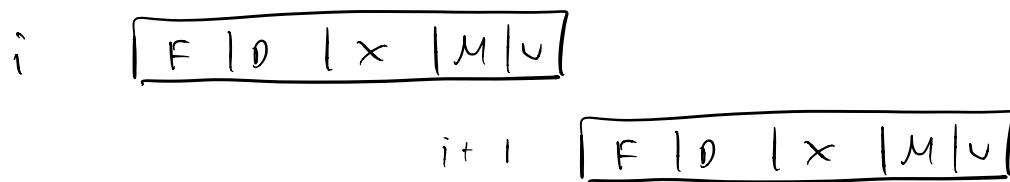
- Problem:

1. single instr \Rightarrow single cycle \Rightarrow cycle period is long!

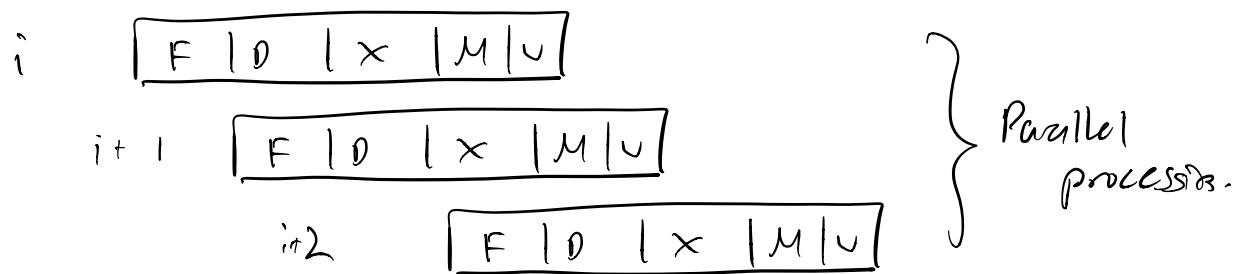
2. Sequential processing slow

- Pipelining: throughput \uparrow

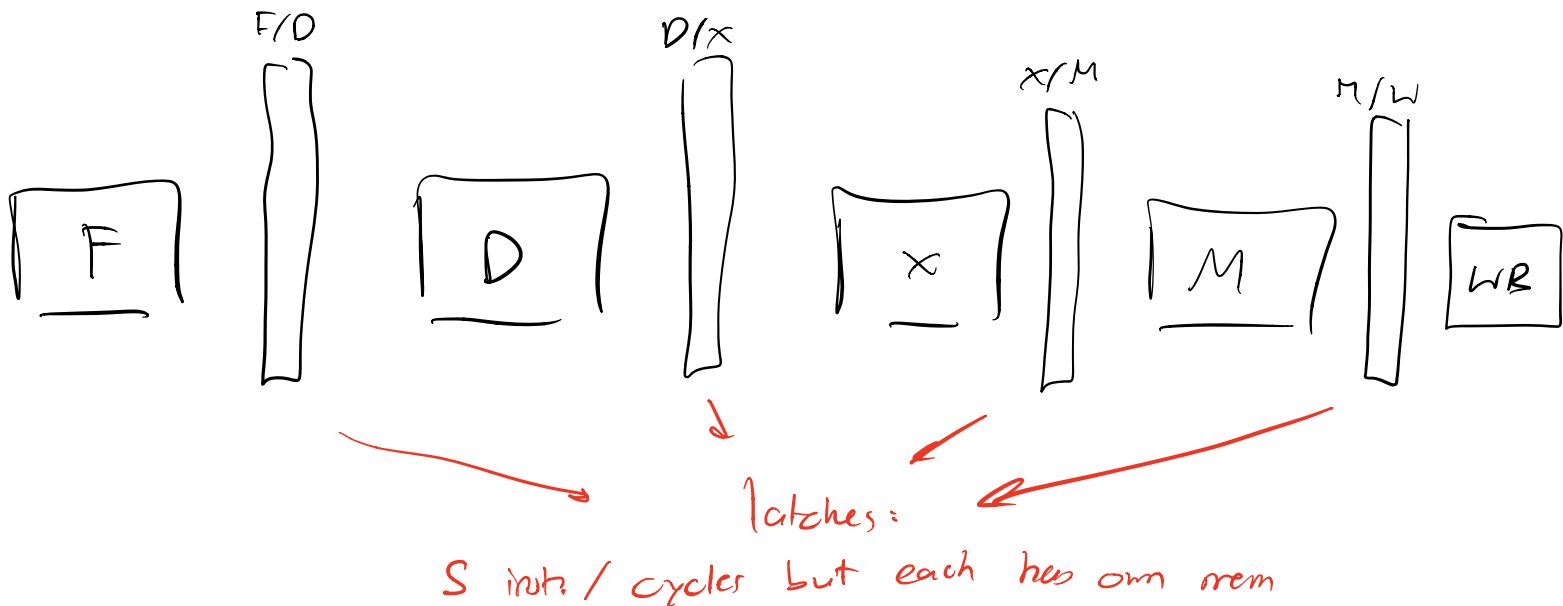
◦ Non-pipelined:



◦ Pipelined:



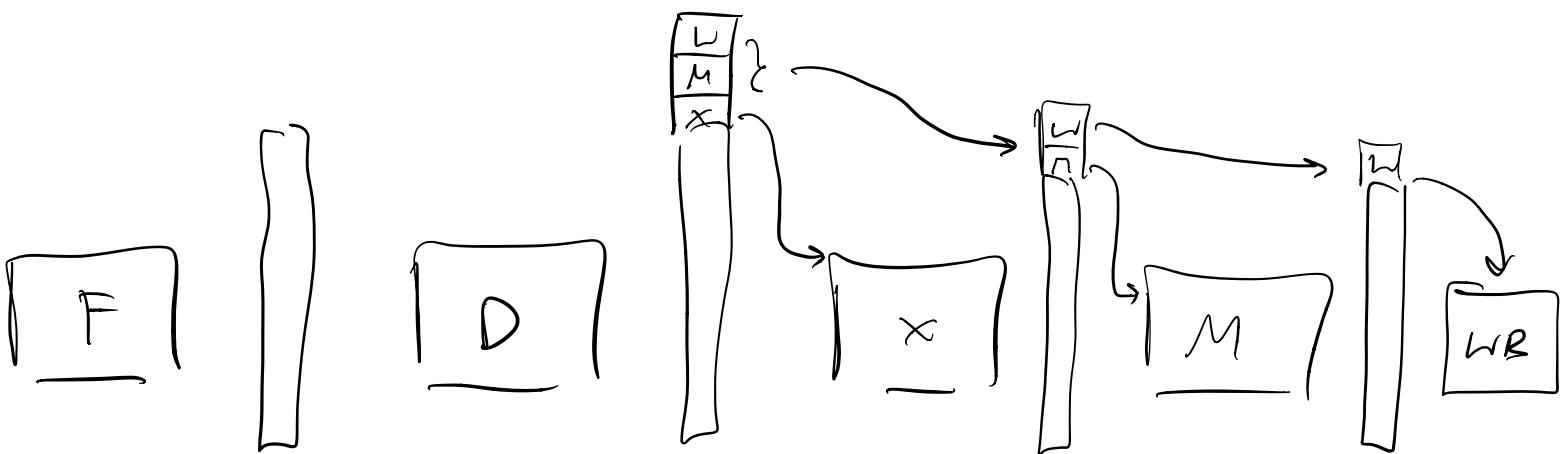
- Hardware:



- Pipeline diagram:

Ins:	Cycles					
	1	2	3	4	5	..
X	F	D	x	m	w	..
Y		F	D	x	m	w
Z			F	D	x	..

- Pipelining controls for each instruction so that controls not mixed for different instructions.



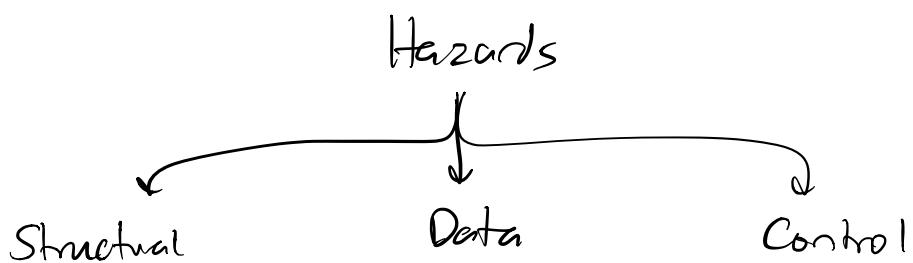
- Performance of pipelined instructions:

- o Single cycle processor: performance = clock period

- o Pipelined processor: 1 clock / stage \Rightarrow clock period \downarrow

$$\text{Clock period} > \frac{\text{single cycle period}}{s} \leftarrow \underline{\text{latches}}$$

- Hazards:



- Structural hazards: Using same resource @ diff. stages

- Ex:// 1 instr. @ decode stage, another at write back stage
Both using reg. file \Rightarrow var precedence

- Soln:

1. Properly design ISA

a) Each instr. uses a structure exactly once
(ports, mux, hardware component)

\square In ex, how diff ports in reg. file

b) Each structure uses struct exactly once

c) Instr. in only 1 stage

- Data hazard

- Dependencies: instr. depends on data from another instr.

Read after write

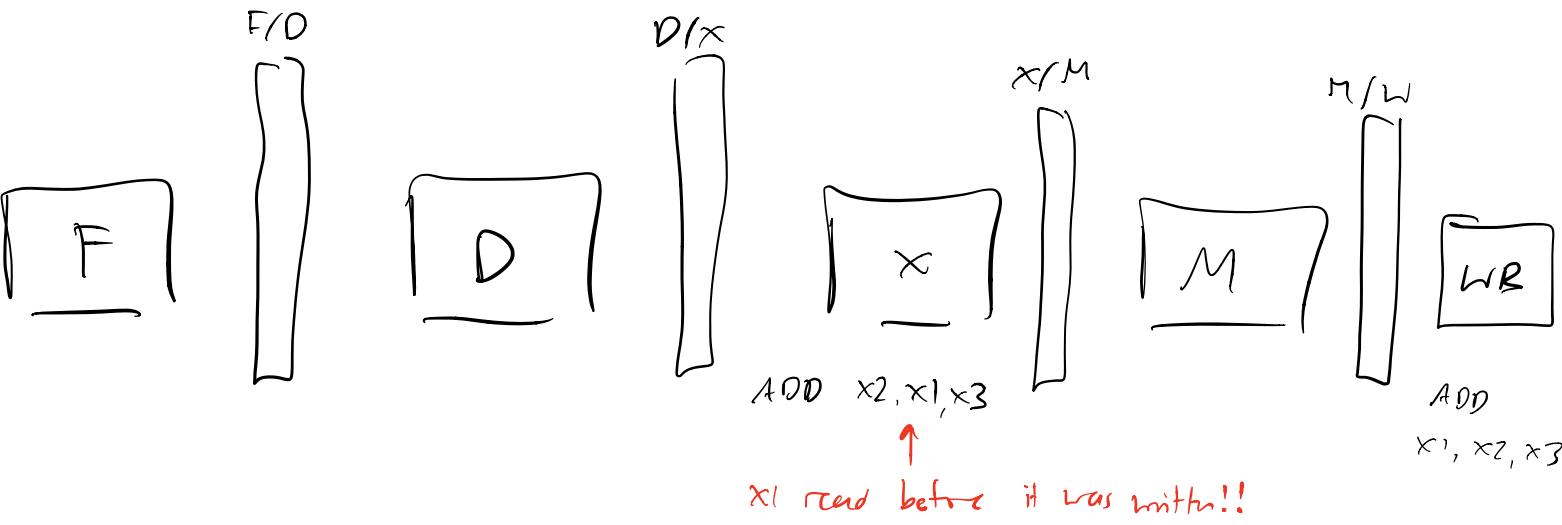
\rightarrow ADD x_1, x_2, x_3
LDR $x_4, [x_1, \#4]$

Write after write

ADD x_1, x_2, x_3
ADD x_2, x_1, x_3

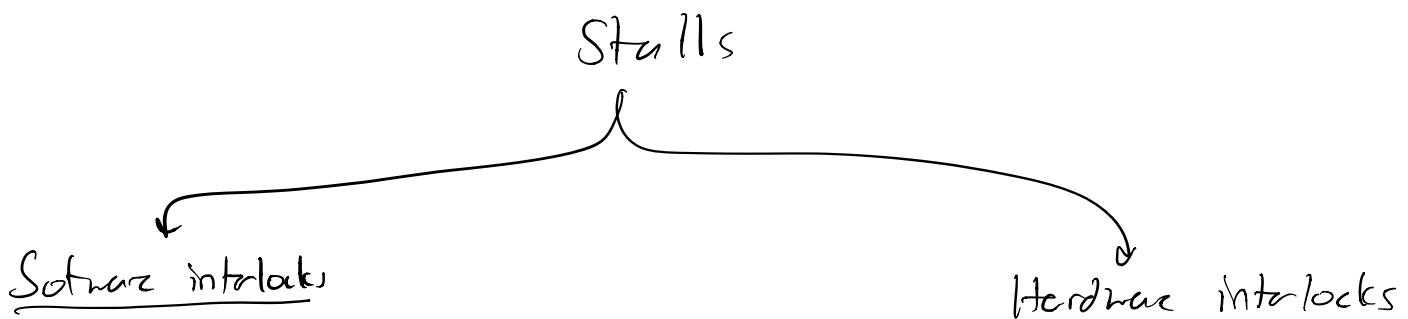
Write after read

LDR $x_4, [x_1, \#4]$
ADD x_1, x_2, x_3



- If problem where: no solutions to core problem + they give sequence of instructions. Will it work?

- ① Find data dependencies in sequence.
 - ② Use pipeline table as you move cycle \Rightarrow update regts.
- Tip: if you read < 3 cycles after write instr. \Rightarrow data hazard
 - Soln: stalls s.t. get correct answers/hypassing



- Compiler puts msp. instructions/ Nops to prevent hazard

◦ Ex:// ADD x3, x2, x1 ①
 LDR x4, [x3] ②
 STR x7, [x3] ③
 ADD x6, x2, x8 ④

①	F	D	x	M	w
④	F	D	x	M	w
②	F	D	x	M	w
③					

This still leaves us w/ hazard,
 so add NOP \Rightarrow delays for later

- Performance: not great
 20% instrs need NOP, 5% need 2 NOPs

$$CPI = 1 + \frac{0.2 \times 1}{1 \text{ NOP}} + \frac{2 \times 0.05}{2 \text{ NOP}}$$

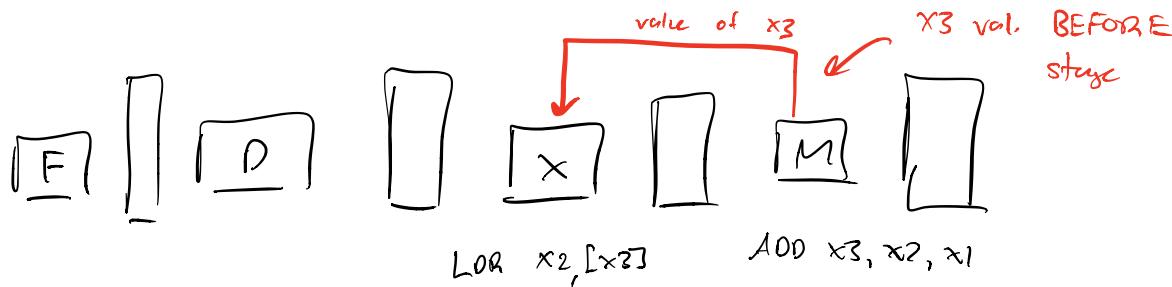
- More compatible \Rightarrow indep. of # of stages
- Detect hazards:
 1. Pass instruction reg. w/ latch
 2. Compare names w/ comb. circuit
 3. Delay if same name is used!
- Fix hazards:
 1. add NOP into D/X latch
 2. Clear controls
 3. Disable F/D latch + PC write (no future instr. should be affected).

◦ Ex://

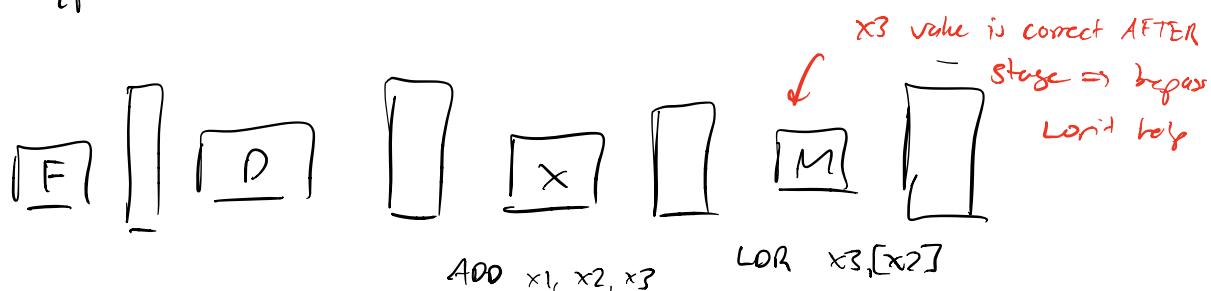
ADD x3, x2, x1	1	2	3	4	5		
LDR x4, [x3]	F	D	x	M	w		
STR x6, [x7]							

Delay propagates!

▫ Bypassing: reduces stall delay by linking stages



- Add bypass logic: select between normal input / bypass input
- Bypass first, then stall if cannot bypass

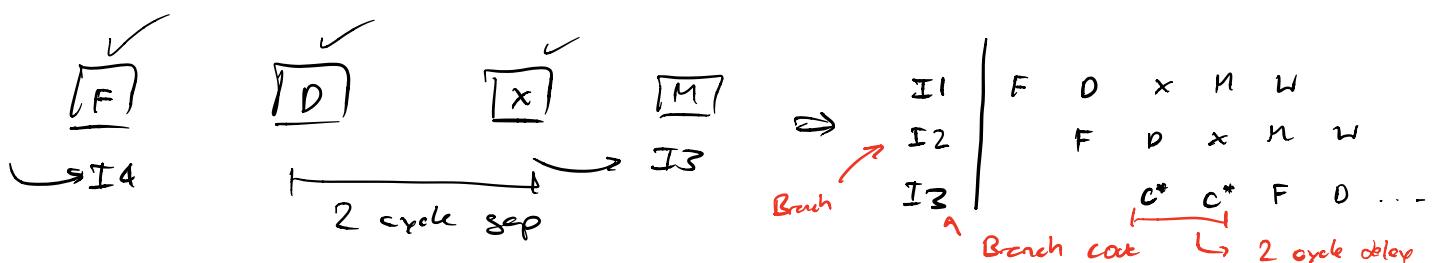
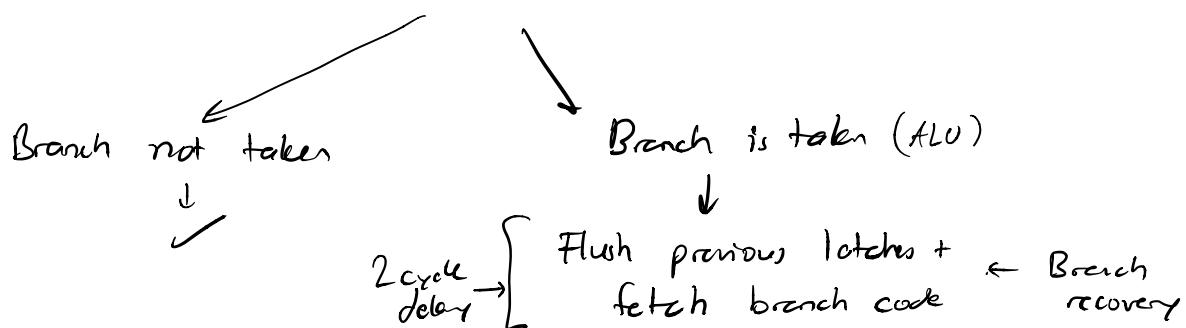


- If value after latch is NOT correct for bypass \Rightarrow stall

- Control hazards: need to get next instr. w/out knowing outcome of branch

- Default solution:

Assume branch not taken & fetch next line of code



- Performance: since so many taken branches \rightarrow CPI \uparrow

- Optimizations:

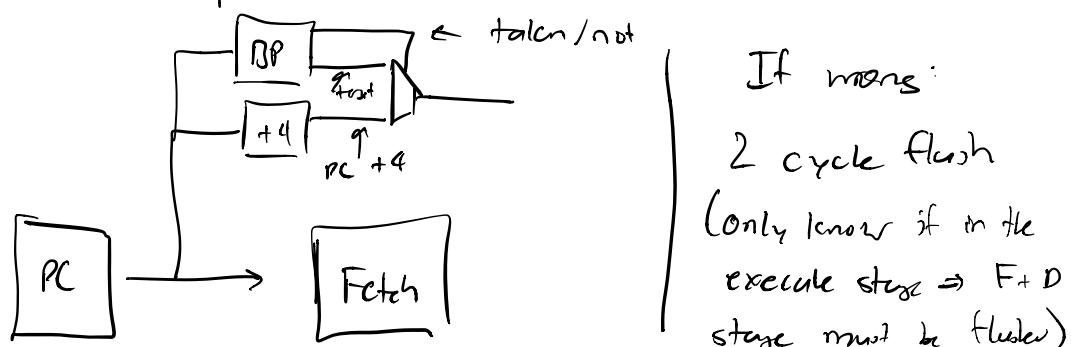
1) Fast branch: move up branch calc. in decode stage

- Pros: flush out only 1 instruction \Rightarrow delay ↓
- Cons: Need to support bypassing into decode
- Performance^{CPI}: $1 + (\% \text{ of branches}) (\% \text{ taken}) (1 \text{ cycle})$

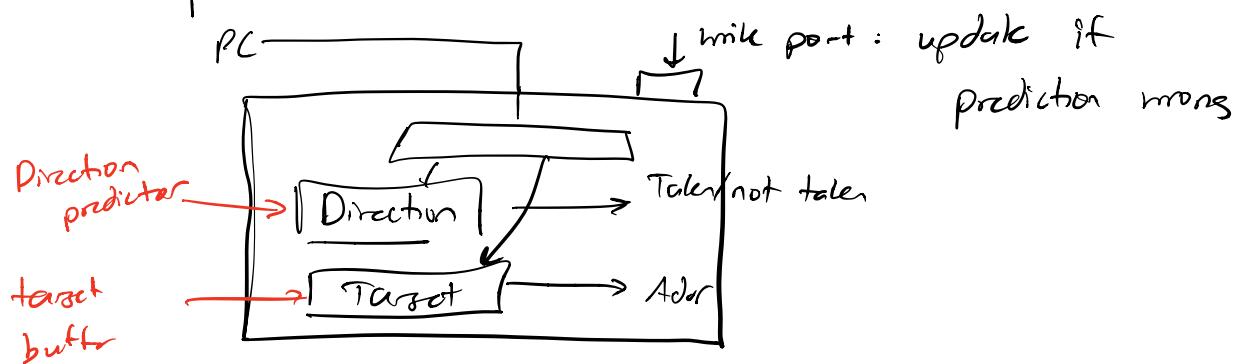
2) Delayed branch: compiler puts indep. insns after branch \Rightarrow no flush

- Pros: delay ↓
- Cons: ISA modification + compiler
- Performance(CPI): $1 + (\% \text{ bran}) (\% \text{ taken}) (1 \text{ cycle}) + (\% \text{ bran})(\% \text{ delay or NOPs}) (1 \text{ cycle})$

3) Dynamic branch prediction: predict next instruction



Branch predictor:

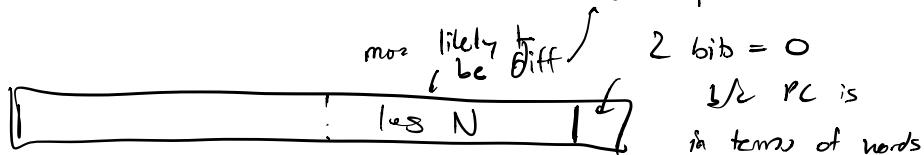


1. What part of PC are we using for prediction?

a) Number of data entries in buffer: N

$\therefore \log N$ bits to index better predictions

b)



2. What happens if 2 PCs have same target/destination
b/c $\log N$ bits same?

↪ Aliasins \Rightarrow not a problem b/c just prediction! :)

3. Updating predictions: link port for incorrect pred.

Performance inc:

$$\text{LPI: } 1 + (\gamma_{\text{branch}}) (\gamma_{\text{mispredicted}}) (2 \text{ cycle})$$

v. low

- If # of stages in pipeline \uparrow :

- Clock cycle \downarrow , but delays longer (more stages to connect)

Parallelism

Instruction

Thread

Data

- Pipelining
- Superscalar (multiple instr.
/ stage)
- Out of order execution