

# AI Hardware

---

Severaj Chetput

---

---

---

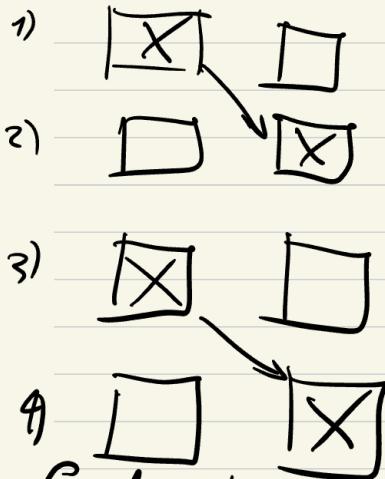


# Review Computer Structures 6.C04 MIT ocw

## 7 - Performance Measurements Latency & Throughput

Laundry ex:

washer 30m    dryer 60m



Combinational

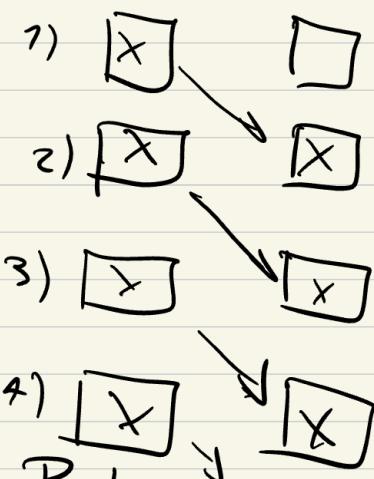
$$T_{PD} = N (\text{washer}_{PD} + \text{dryer}_{PD})$$

$= N \cdot 90$

Loads

Comb.

washer 30m    dryer 60m



2-stage  
Pipeline

Pipeline

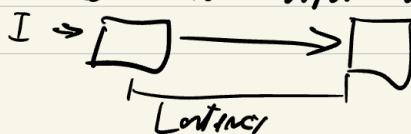
$$T_{PD} = N \max(\text{washer}_{PD}, \text{dryer}_{PD})$$

$= N \cdot 60$

> don't account for setup time

Performance Measures:  
Latency:

Delay from when an input is established until its associated output w/ that input becomes valid



$$\text{Comb. latency} = 90 \text{ min}$$

$$\text{Pipeline latency} = 120 \text{ min}$$

## Throughput

Rate at which I/O's are processed

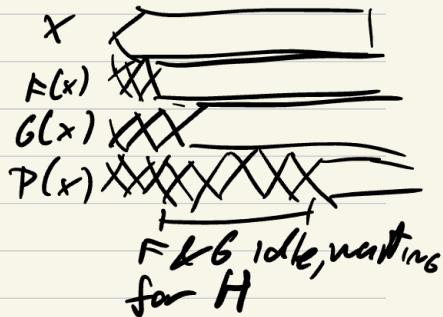
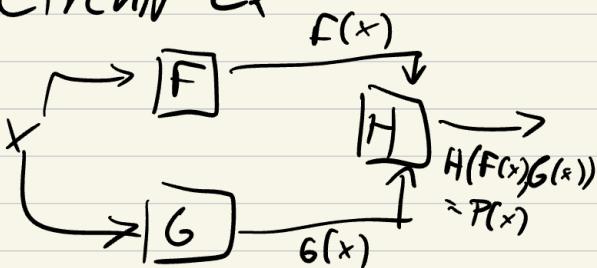
$$\text{Comb throughput} = \frac{1}{90} \text{ min}$$

$$\text{Pipeline throughput} = \frac{1}{60} \text{ min}$$

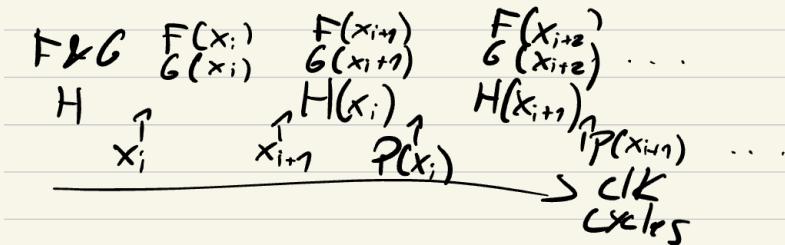
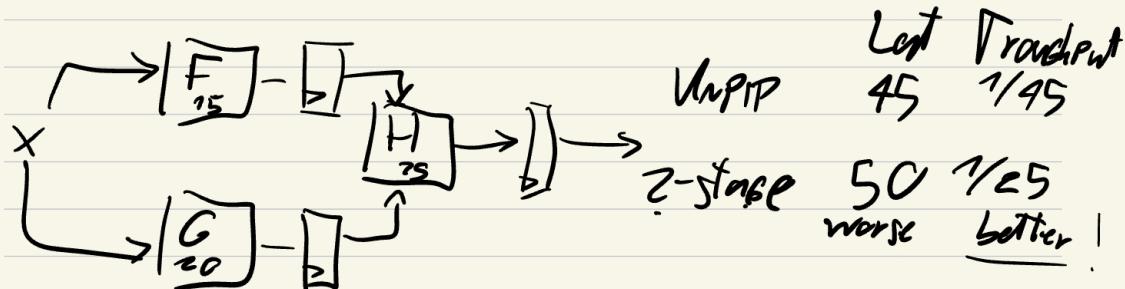
Comb does 1 job in less time

Pipeline produces > output / time

Circuit ex:



Bad throughput!



# Pipeline Conventions

Def:

Well-formed K-stage Pipeline ("K Pipeline")  
acyclic circuit w/ exactly K reg & path  $I \rightarrow O$   
Comb. circuit = 0-stage Pipeline

## Composition Convention

( $\wedge$  pipeline stage  $\Rightarrow$  K-stage Pipeline) has a reg on its O

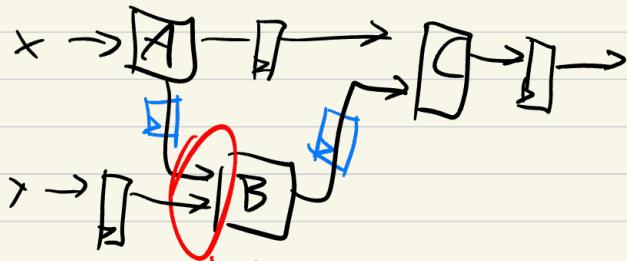
Always:

CLK common to all regs must have T sufficient to  
cover  $t_{\text{comb}} + t_{\text{PD}} + t_{\text{setup}}$  ( $= T_{\text{clk}}$ )

Latency of K-Pipeline =  $K \cdot T_{\text{clk}}$

Throughput of K-Pipeline =  $f_{\text{clk}}$

## III-formed Pipelines



Successive inputs mixed!

fix

## Pipeline Methodology

- 1) Draw line that crosses TO & mark endpoints as terminal points
- 2) Continue drawing lines b/w terminal points across various stages ensuring to make crosses to line in same direction

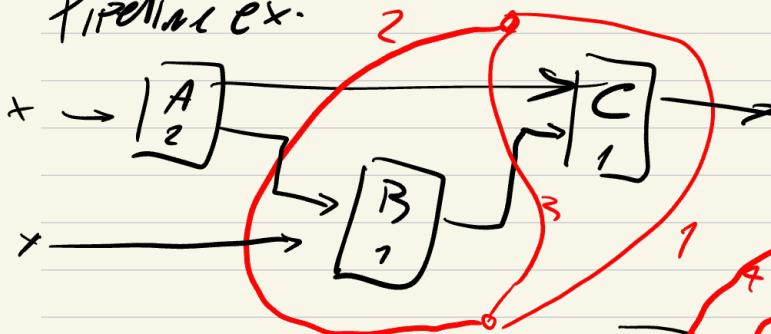
Lines demarcate pipeline stages

Adding pipeline reg b/w int where separating line crosses a connec.

will always generate a valid pipeline

Strategy: focus on tent or on placing reg around bottlenecks

Pipelined ex:



L	T
0	4
1	4
2	4
3	6

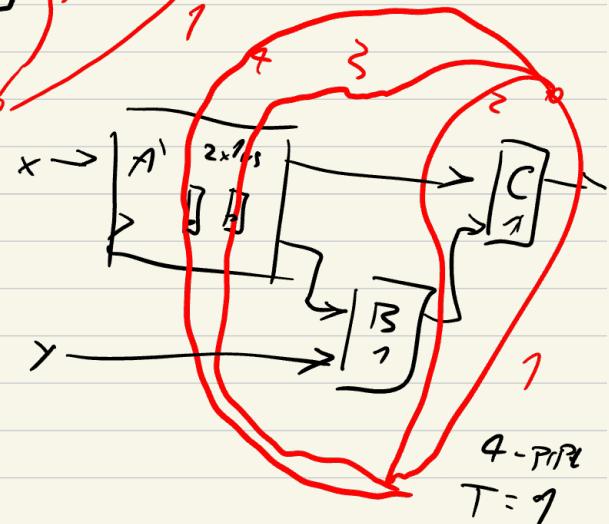
Pipelined Comp.

Pipelined sys's can  
be hierarchical

Replacing a slow comp  
comp. w/ K-pipe version

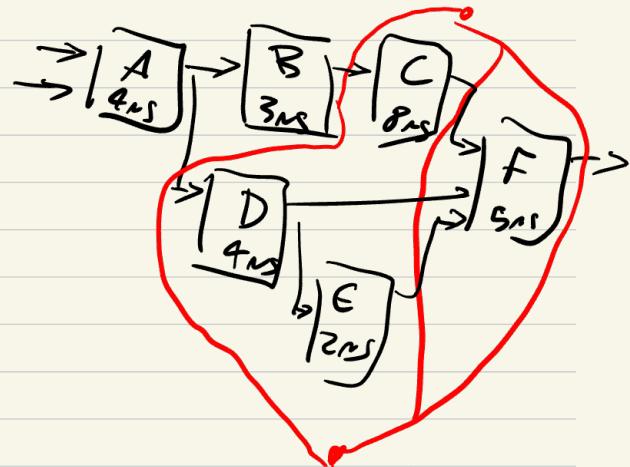
may  $\downarrow$   $T_{CIK}$

Must account for new  
pipeline stages in new plan



4-pipe  
 $T = 9$

Ex: C slowest comp.  
limits  $T_{CIK} = 8\text{ns}$   
 $\Rightarrow T = 7/8\text{ns}$

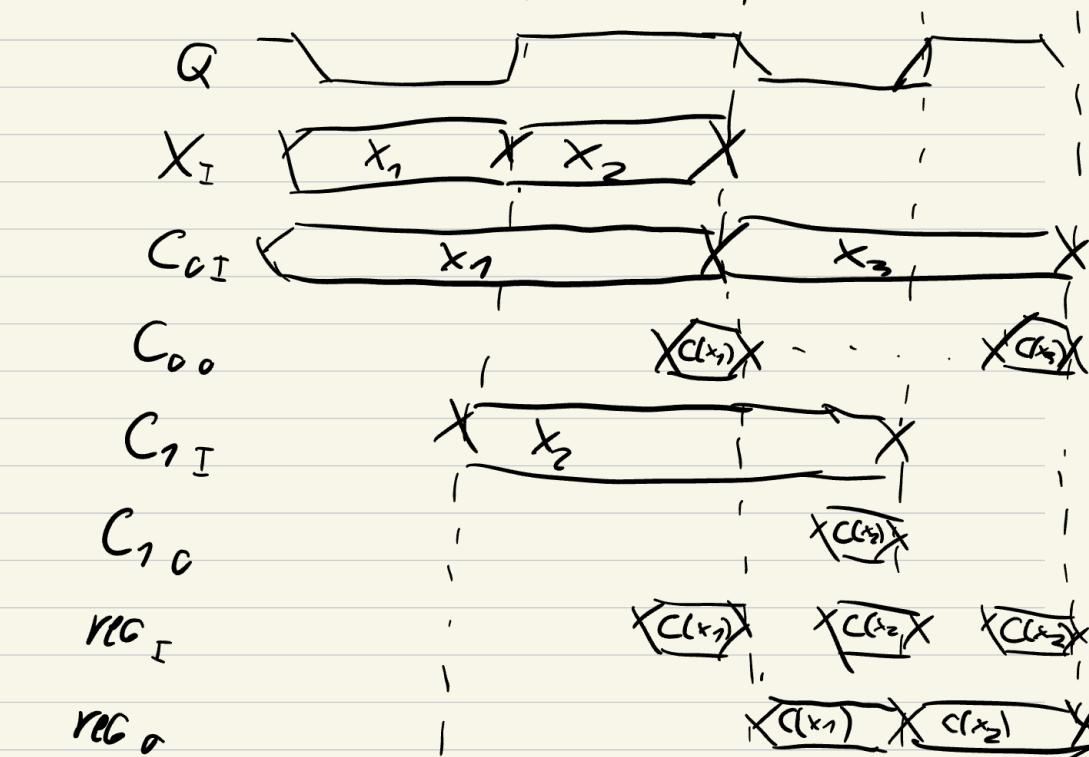
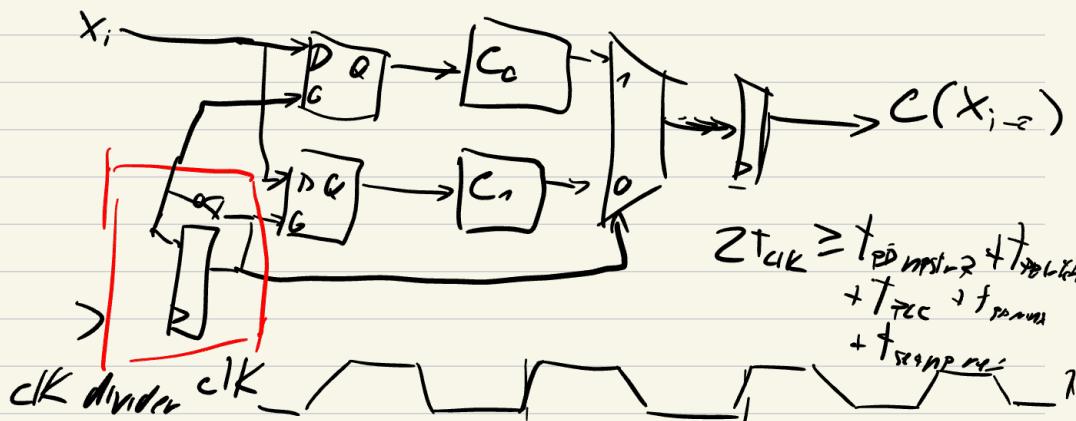


Improve by:

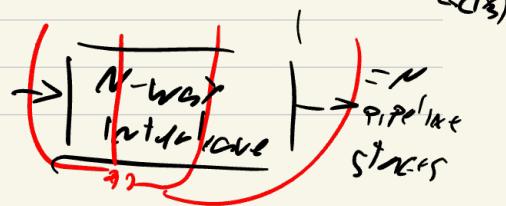
a) Pipelined C

b) Interleaving  $\Rightarrow ?$   
→ copies of C

# Circuit interleaving



2-clock martinxing  
"in by  $t_i$ , out by  $T_{i+2}$   
 $T = 1/clk$   
 $L = 2clk$

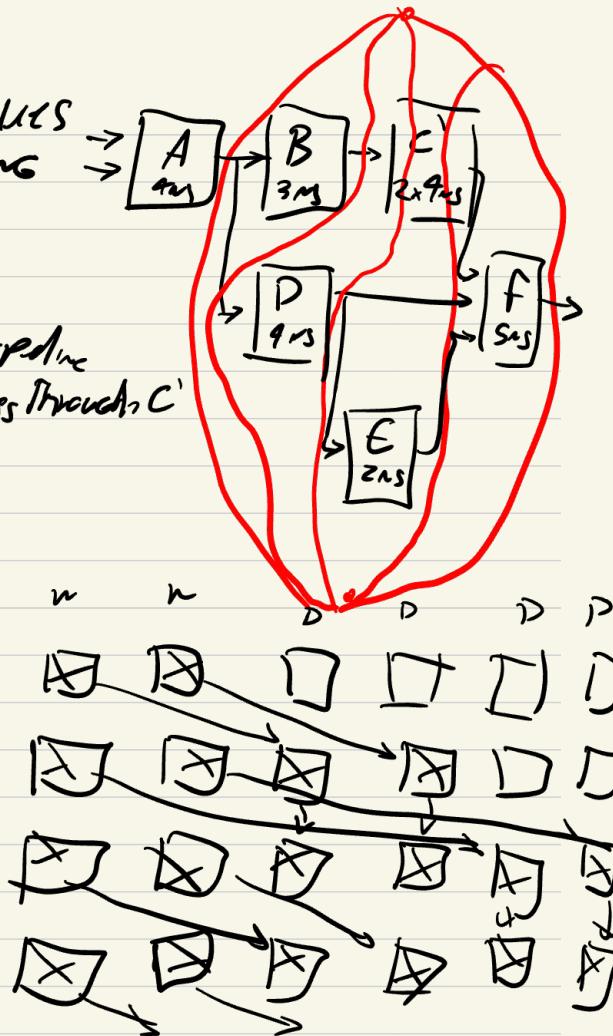


Combine Techniques  
 Interleaving + Pipelining  
 $C'$  interleaves  $Z \times C$   
 Effective  $T_{clk} = 4\text{ ns}$ ,  $L = 8\text{ ns}$

$C'$  behaves as 2-stage pipeline  
 2 pipelines contours must pass through  $C'$

Bottleneck:  $C \rightarrow F$

Parallelism  
 Interleaving  
 + Pipelining  
 + Parallelism

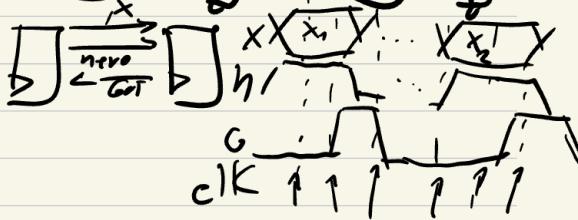


## Control Structure Alternatives

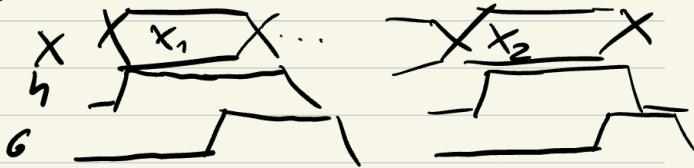
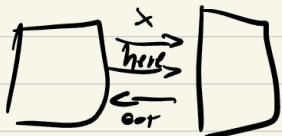
Sync, globally-timed:



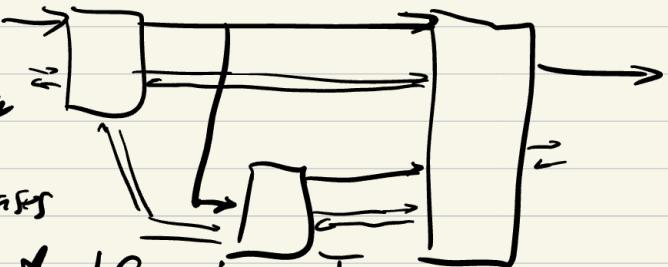
Sync, locally-timed:  
 local FSMs; control flow  
 at data w/ handshake



Asyc, locally-timed:  
using Transition  
signalling



Soft-timed ex:  
each comp. spec. its own  
time constraints  
local adapt to spec. cases  
asyc or sync



### Control Structure Taxonomy

easy to design  
fixed-sized interval  
(can be wasteful)

glob.

timed

centralized c/KFSM

Generates & ctrl signals

Ssync

locally  
timed

Start/Finish signals gen.  
by major subsyst. sync.  
w/ global c/K

best way to  
break up large sys  
w/ indep. timed  
comps

Central ctrl unit

Tailors ctrl time slice  
to ctrl tasks

Asyc

→ inner sys/  
very comp. imp.  
gen. req!

each subcomp. take  
asyc start, gen asyc  
finish (perhaps w/ local c/K)

↓ lot of design work load  
extra work is worth it in  
specia'l cases

### Summary

Latency:  $t$  for  $I \rightarrow O$   
Throughput:  $\otimes O's/t$   
Comb:  $L = t_{PD} + T = K/T$

K-Pipeline

Vreg on output

K reads & push  $I \rightarrow O$

I avail. after c/K:

O avail after dK/k

$t_{c/K} = t_{PD, reg} + t_{PD,bus} + t_{ctrl,p}$

?  $T \Rightarrow$  split slow PIP.  
stages  
req. width, if no  
further split

$L = K + c/K = K/T$

$L_{IP} \gg L_{comb}$

# Caches & Memory Hierarchy

Single-Cycle CPU

Don't do Inst. Mem. or the slowest piece of HW in system  
↑ \* MFTG's  
↑ propagation delay

(also 1 stage)

Execution starts by fetching instr. from main mem

All data processed by CPU loaded/stored from main mem

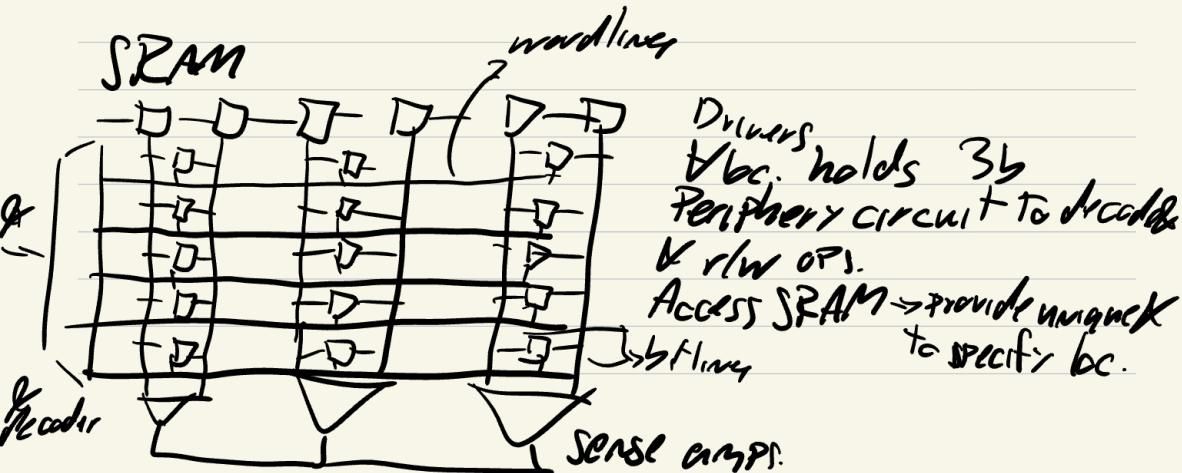
Most prog. manipulate much more data than can be accom. in reg file

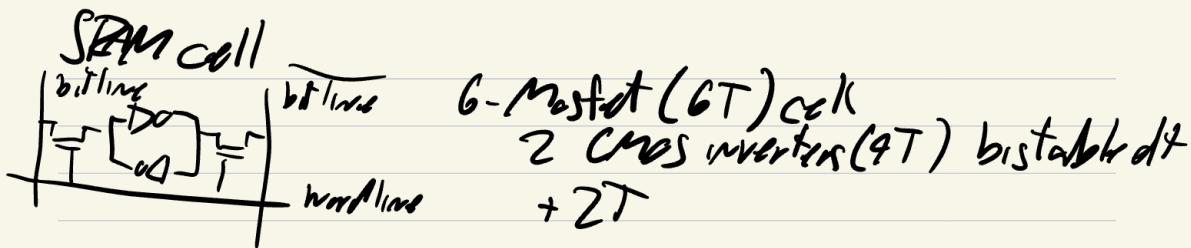
Most CPUs limited by b/w (B/s) CPU  $\rightarrow$  main mem  
(Memory bottleneck)

## Memory Technologies

	cap	latency	Cost / GB	
ROM	1KB	20ps	\$888	— processor data path
SRAM	10KB-1GB	1-10ns	~\$1K	mem hierarchy
DRAM	~10GB	80ns	~\$10	
Flash*	~100GB	100ms	~\$1	I/O
HD *	~1TB	1Gms	~\$0.20	subsys

\* non volatile





During R, drivers recharge & bitlines to VDD & disconnect bitlines floating @?

& Decoder sets 1 of wordlines high, connecting word bitcells to their bitline

& cell in active. word pulls 1 of the bitline to 0 Gnd  
 $\Rightarrow$  slow, since bitline  $\leq$   $C_{total}$ ,  $\tau = t_{on} + t_{off}$   
 keep cell as small as possible

Senseamps sense change in bitline voltages, produce Odata

During W, drivers set & hold bitline to desired vals  
 & decoders activates 1 holdline

Each cell in word is overpowered by drivers, stores val

### Multiplexed SRAMs

so far can do either  $r/w$

Multi  $r/w$  by  $\sim$  1 set of wordlines & bitlines per port  
 Cost/bit?  $\sim$  ports:

$N$  wordlines

$2N$  bitlines

$2N$  accessFETs

wires often

dominate area

$\rightarrow O(N^2)$  area!

### Summary:

Any of  $K^b$  cells ( $K$  words,  $b$  cells/word)

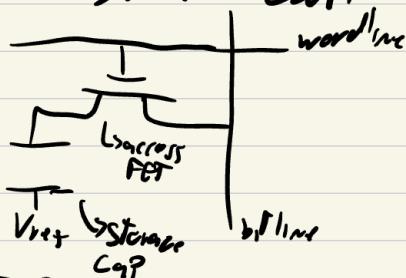
Cell bistable off + access T

Analog circuit to allow  $r/w$

Read: precharge bitlines, activate wordline, sense  
 write: drive bitlines, activate wordline, overdrive cells

# DRAM

1T DRAM cell



20x smaller area  
dense & cheap!  
C tracks q. must be refreshed periodically (~ms)

DRAM wr

w: drive bit line, activate word line, discharge C

r: Precharge bit line to  $V_{DD}/2$

activate word line

C & bit line share q

If C discharged,  $V_{bit\ line} \downarrow$

If C charged,  $V_{bit\ line} \uparrow$

Sense bit line to determine C/q

Reads are destructive!  
data must be rewritten  
to cell & r.

Summary: DRAM

1T DRAM cell: T + C

↓ size vs SRAM vs destructive reads & Clocks

DRAM only needs access to

n word after Vr

refresh every word periodically

is SRAM:

20x denser

~2-10x slower

Non-Volatile Storage: Flash

e- @ float gate & strength from charge

⇒ no inversion

⇒ NFT off even when word is high

Very dense: mult bits/T, r/w in blocks

Slow: 10-100ns

Limited: floating gate (n) PV damage T

Non-Volatile Storage: HD

Extremely slow ( $\sim 10\text{ms}$ )

Cheap

Memory Hierarchy

Want large, fast & cheap memory, but

Large = slow

Fast = expensive

Idea: Hierarchical system of memories w/ different tradeoffs to emulate large, fast, cheap memory?



Interface

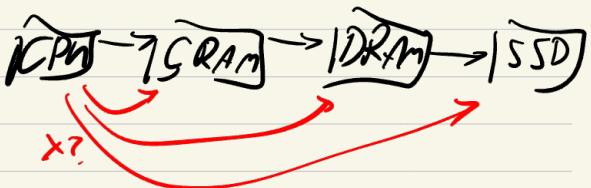
Approach 1: Expose hierarchy | CPU |  
regs, SRAM, DRAM, flash, HD  
each avail as storage option.

Tell programmers "use clusters"

Approach 2: hide hierarchy

Program. model: single mem,  
single space &

Machine transp. stores data  
in fast/slow mem depend.  
on user & wait times



Locality Principle

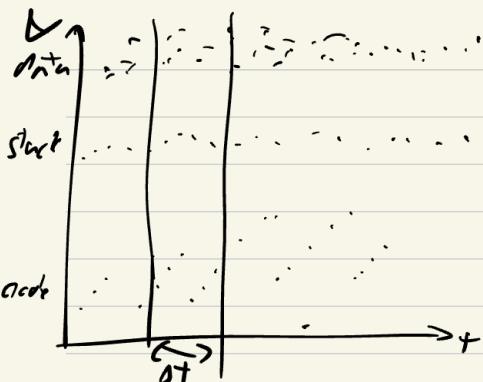
Keep most often-used data in small, fast SRAM (after local to CPU)

Refer to my notes for rem. data

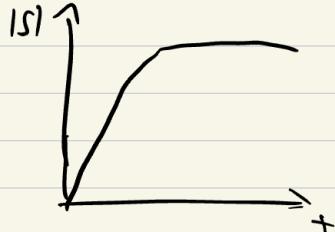
Why it works: LOCALITY

Locality of reference: access to  $X$  at  $t \Rightarrow$  access to  $X + \Delta X$  at  $t + \Delta t$  more probable as  $\Delta X, \Delta t \rightarrow 0$

# Memory Reference Patterns



$S$  set of locations accessed during  $\Delta t$   
Working set  $S$  changes slowly w.r.t. according  $\Delta t$



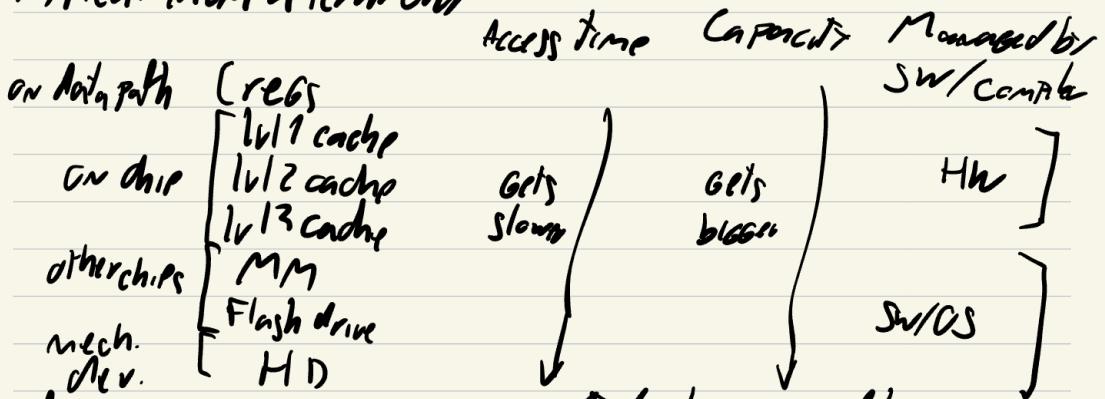
## Caches

Small, interim component that transp. retains (caches) data from recently accessed locations

- ?? fast access if data is cached other access slower, larger cache/mem
- exploits locality principle

Comp. systs often use multiple levels of cache  
widely applied beyond HW (e.g. web caches)

## Typical mem hierarchy



Cache access  
process must deal w/ variable mem access

CPU  $\leftrightarrow$  Cache

1) Cache hit

2) Cache miss

return data from memory



# Cache Metrics

$$\text{Hit Ratio} : HR = \frac{\text{hit}}{\text{hits + misses}} = 1 - MR$$

$$\text{Miss Ratio} : MR = \frac{\text{misses}}{\text{hits + misses}} = 1 - HR$$

$$\text{Avg. Mem. Access Time (AMAT)} = \text{HitTime} + MR \cdot \text{MissPenalty}$$

Goal of caching:  $\downarrow$  AMAT

Formula can be applied recursively w/ multi-level hierarchies

Ex: How high of an HR?

HR needed to break even  
(MM only: AMAT = 100)

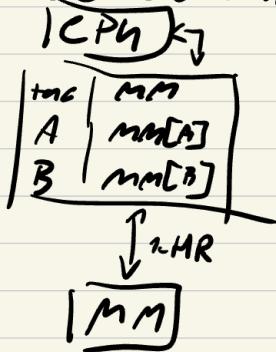
$$100 = 4 + (1 - HR)100 \quad HR = 9\%$$

HR to achieve AMAT = 5?

$$5 = 4 + (1 - HR)100 \quad HR = 99\%$$



## Basic Cache Algorithm



on ref. to  $mm[x]$  look for  $X$  among tags  
HIT:  $X = TAC(i)$   
Read: return  $DATA(i)$   
Write: change  $DATA(i)$ , start w/  $mm[x]$

MISS:  $X \notin TAC$

Repl. Sel: sel same k to hold  $mm[x]$   
Read: Read  $mm[x]$ , set  $TAC(k) = X$ ,  $DATA(k) = mm[x]$   
Write: start w/  $mm[x]$   
set  $TAC(k) = X$ ,  $DATA(k) = new\ mm[x]$

Q: How to "search" Cache?

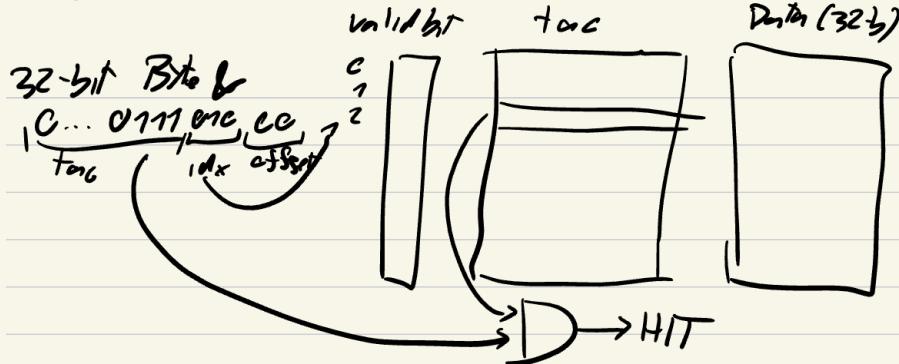
## Direct Mapped Caches

Each word in memory maps into a single cache line

Access (for cache + 2<sup>W</sup> lines)

Index into cache w/ W bits (idx bits)  
read out valid bit, tag to match  
if (valid bit) & tag matches w/ idx bits, HIT

Ex: 8-location PM cache ( $w=3$ )



Block Size

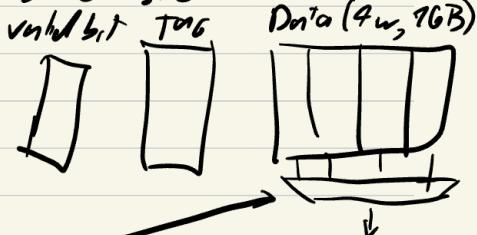
✓ Take adv. of locality; increase block size

✓ Reduces size of tag mem

✗ Fewer blocks in cache

32b Block A

1 tag 26b 14 block offset  
25 4b



Block Size tradeoffs

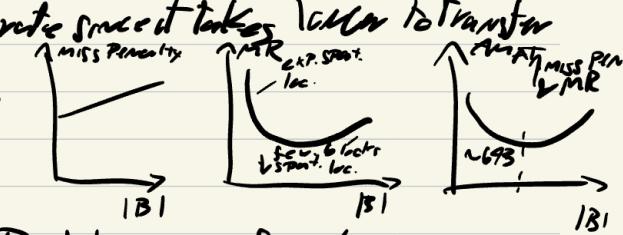
Larger Block sizes

Advantage of spatial locality

Incur larger miss penalty rate since it takes longer to transfer block into cache

Can increase avg HT & MR

$$AMAT = HT + Miss\text{ Penalty} \cdot MR$$



Direct-Mapped Cache Problem: Conflict / Misses

Assume 1024-line DM cache,  $|B| = 1 w$

Consider looping code, stand, static. Word  $\neq$  (not Blc)

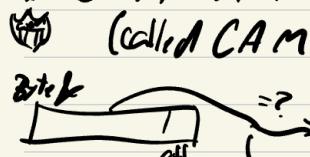
Inflexible mapping ( $W$  can be in only 1 cache loc.)  $\Rightarrow$  conflict misses!

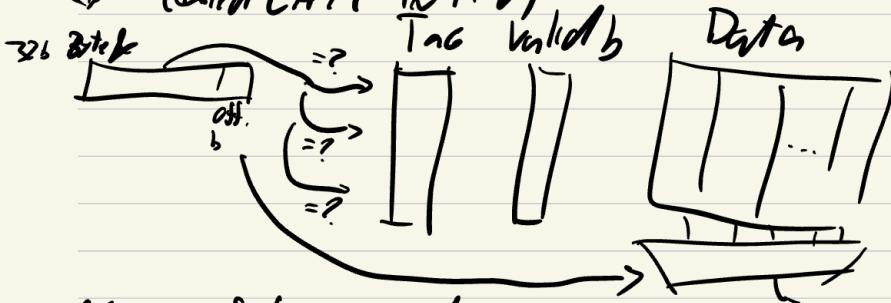
# Fully-Associative Cache

Opposite extreme:  $\forall \text{ cache can be } \forall \text{ loc.}$

No cache idx

Flexible (no conflict misses)

- \* Extensive: Must comp. tags w/ entry in parallel to find match  




## $N$ -way Set-associative

Compromise b/w direct + mapped to fully-associative

$$\# \text{ rows} = \# \text{ sets}$$

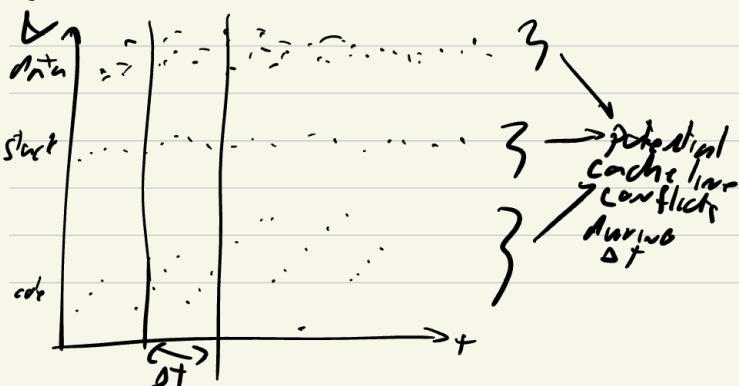
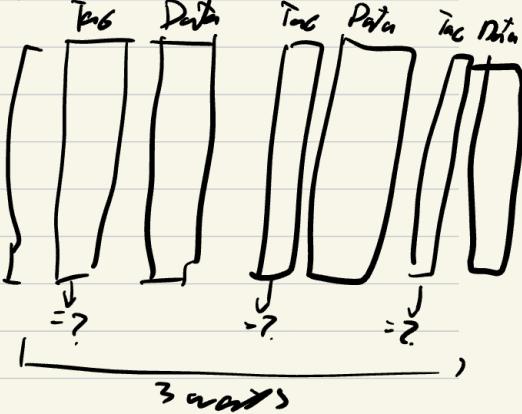
$$\# \text{ cols} = \# \text{ ways}$$

SetSize = "set associativity"       $\xrightarrow{\text{sets}}$   
 (e.g. 4-way  $\Rightarrow 4 \text{ entries/SET}$ )

Compare  $\forall \text{ tag } \forall \text{ ways in Parallel}$

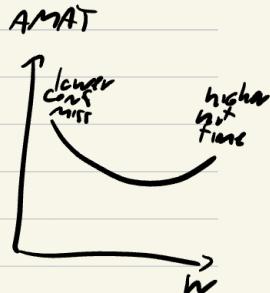
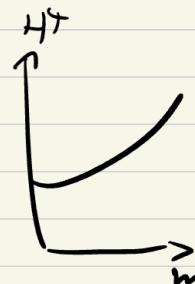
$N$ -way cache =  $N$  direct-mapped  
caches in parallel

Direct-mapped & fully-associative  
special cases of  $N$ -way set-associative



# Associativity Tradeoffs

More ways  
↓ Conflict misses  
↑ Hit Time



## Associativity Implies Choices

Direct-Mapped

$N$ -way set-associ.

Fully assoc

$\Delta \stackrel{?}{=} 1 \text{ Tag}$   
Loc A can be stored  
in exact 1 cache line

$\Delta \stackrel{?}{=} N \text{ Tags}$   
Loc A can be stored  
in exact 1 set, but in  
any of the  $N$  cache  
lines belonging to  
that set

$\Delta \stackrel{?}{=} \sqrt{N} \text{ tags}$   
Loc A can be stored  
in any cache line

## Replacement Policies

Optimal Policy (Belady's Mru): replace the block that is accessed the furthest in the future  $\rightarrow$  req. Knowing the future  
Idea: If a block has not been used recently, it's often less likely  
to be accessed in the near future

Least Recently Used(LRU): replace block that was accessed  
furthest in the past

- works well in practice

- Needs to keep ordered list of  $N$  items  $\rightarrow N!$  orderings  
 $\rightarrow O(\log N!) = O(N \log N)$  "LRU" bits + CMPT logic

Caches often implement cheaper approx of LRU

Other Policies: FIFO

Random  $\rightarrow$  not very good, but doesn't have adversarial access patterns

## Write Policy

Write-through: CPU writes cached, but also written to MM  
imm. (stall CPU until write completed). MM always holds cache.  
Simple, slow, wastes b/w

Write-behind: CPU writes Cached, not to MM may be buffered.  
CPU keeps exec. while w comp. in background  
Faster, still uses a lot of b/w

Write-back: CPU writes Cached, not written to MM until block replaced.  
MM contents can be "stale"  
Faster, low b/w, more cpx  
Commonly implemented in curr sys

Write-back: on ref to  $mem[x]$  look for X in TAC  
HIT:  $X = TAC(:)$

Read: return DATA(:)

Writer: change DATA(:), ~~start w to  $mem[x]$~~

MISS:  $X \notin TAC$

Repl. Sel: sel same K to hold  $mem[x]$

Read: Read  $mem[x]$ , set  $TAC(K)=X$ ,  $DATA(K)=mem[x]$

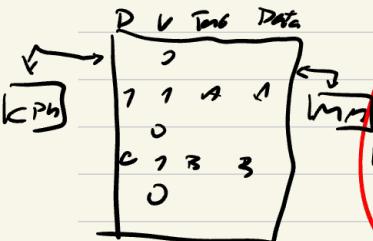
Writer: ~~start w to  $mem[x]$~~

set  $TAC(K)=X$ ,  $DATA(K)=new\ mem[x]$

write  $DATA[K]$  to  
 $mem[&TAC[K]]$

Write-back  
w/ "dirty" bits

Add 1 bit to record whether or not block has been written to. Early write-back via  $\times$  blocks



on  $\text{refToMem}[x]$  look for  $X$  among tags  
 $\text{HIT}: X = \text{TAC}(i)$

Read: return  $\text{DATA}(i)$

$D[i] = 1$

Write: change  $\text{DATA}(i)$ , ~~start w to  $\text{Mem}[x]$~~

MISS:  $X \notin \text{TAC}$

Repl. Sel: sel same  $k$  to hold  $\text{MM}[x]$

Read: Read  $\text{MM}[x]$ , set  $\text{TAC}(K) = X$ ,  $\text{DATA}(K) = \text{MM}[x]$   
 $\text{call}_\text{loc}$   $D[K] = 1$

Write: ~~start w to  $\text{MM}[x]$~~   $D[K] = 1$ ,  
set  $\text{TAC}(K) = X$ ,  $\text{DATA}(K) = \text{new MM}[x]$

$\rightarrow \text{IS}(D[K])$  w  $\text{Data}[K]$  to  $\text{MM}[K]$  &  $\text{TAC}[K]$

## Summary: Cache Tradeoffs

$$AT = \text{Hit Time} + MR \cdot M_{i,j}, \text{Penalty}$$

• Cache size:  $\downarrow MR$ ,  $\uparrow$  Hit Time

•  $|IB|$ : To start vs Temp. locality,  $\uparrow$  Miss Penalty

• Assoc. (ways):  $\downarrow MR$ ,  $\uparrow$  Hit Time

•写入替换:  $\downarrow MR$ ,  $\uparrow$  Cost

• Write Policy:  $\downarrow$  b/w,  $\uparrow$  Comp.

How to navigate all these dimensions? Simulate different cache organizations on real programs.

# Pipelining CPU

32-bit, single-cycle CPU

1 instruction / clk cycle

Circuit Loads new PC → sent to MM as & instr.

→ 32-bit word instr. fetched from MM

→ Opcode decoded by ctrl logic → ctrl signals

→ Operands read from register → passed to ALU

→ Mem ops: ALU → Reg / Data[1..n]

PC + 4 & ALU values can be written to regfile

$t_{CLK}$  determined by delay through & comp. involved in instr. exec.  
How to make it faster?

## Single-Cycle Beta Performance

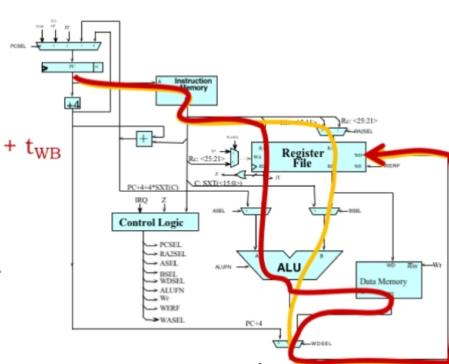
$$\frac{\text{Time}}{\text{Program}} = \frac{\text{Instructions}}{\text{Program}} \cdot \frac{\text{Cycles}}{\text{Instruction}} \cdot \frac{\text{Time}}{\text{Cycle}}$$

CPI       $t_{CLK}$

- CPI = 1
- $t_{CLK}$  = Longest path for any instruction

$$t_{CLK} \approx t_{IFETCH} + t_{RF} + t_{ALU} + t_{MEM} + t_{WB}$$

- Slow
- Inflexible: Instructions with smaller critical path cannot execute faster



## Pipelined Implementation

Divide datapath in Pipeline stages to fit clk

1 instr. exec. over N cycles

Consec. instruc. overlapped

To keep CPI  $\approx 1.0$

I: - Register File: monitors PC,  
↓ fetches instr., passes to

RF: regfile: reads src from regfile, passes to

ALU: performs op in ALU, passes result to

Mem: if LD, use ALU res. as L, pass mem data to

WB: Write-back: writes res. back to reg file

CPU has state: PC, regfile, Mem

Independent branches

- compute next PC

- write res into regfile

## Pipeline Hazards

Pipeline: tries to overlap instrs.  $\rightarrow$  may depend on results of earlier instr.

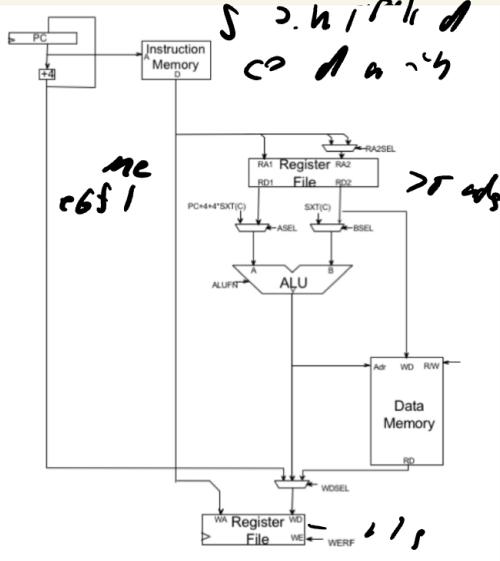
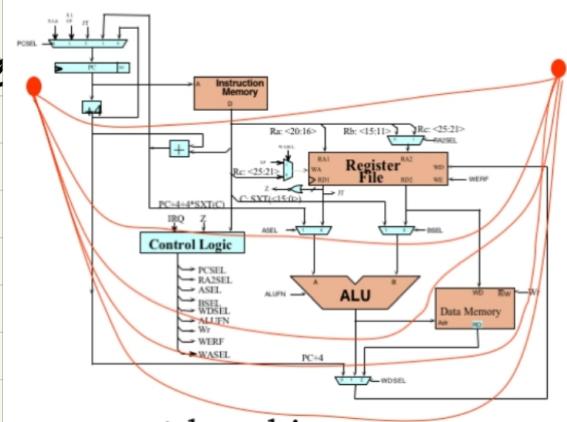
Data reading  $\rightarrow$  data hazard

PC  $\rightarrow$  ctrl hazard (branches, jumps, exceptions)

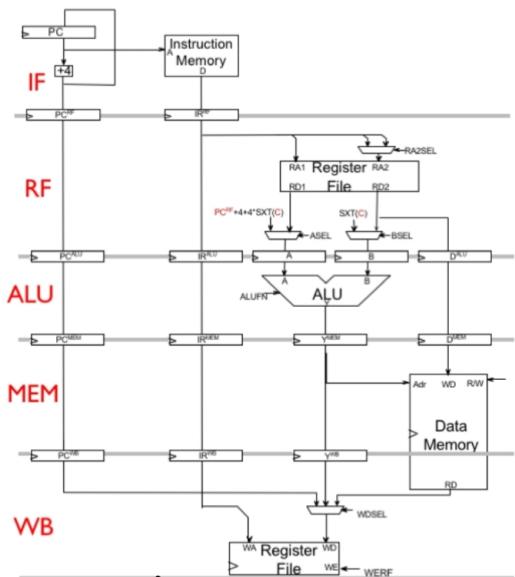
- 1) 5-stage Pipeline

- 2) Handle data hazards

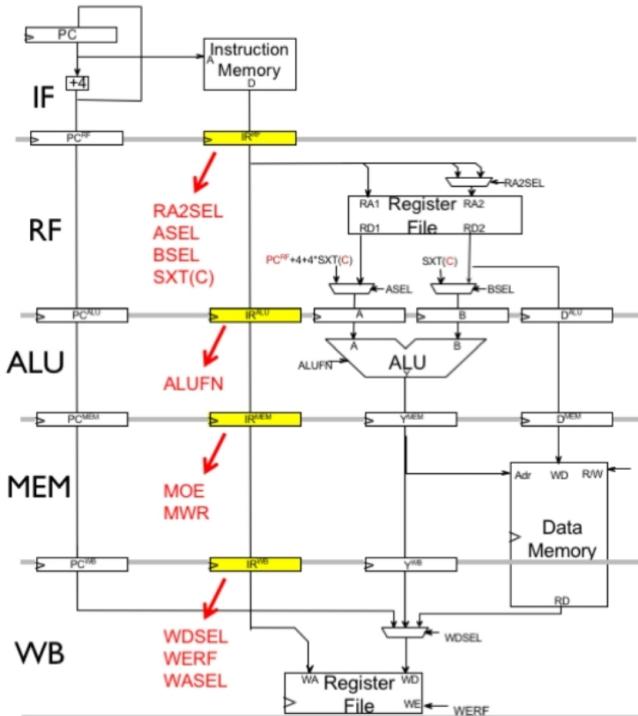
- 3) handle ctrl hazards



# 5-stage Pipeline Datapath



I turned control



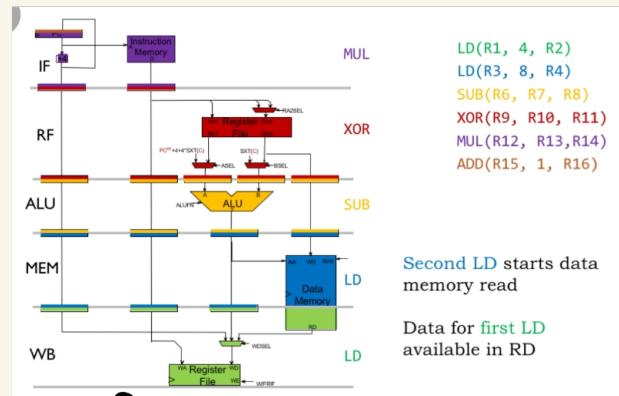
- 5 stages services 1 instr/cycle
- Data mem reads are pipelined, not combinational
- Data reads appear in RD the next cycle

Instr. contents propagated through pipeline in Instruction Registers (IR's)

Ctrl signals 1 stage generated from corresponding IR

Pipeline hazards will require new ctrl signals

# Example of Pipelined Instructions



Second LD starts data memory read

Data for first LD available in RD

## Timeline Diagrams When do r/w happen?

	Cycles →					
Stages	I	2	3	4	5	6
IF	LD	LD	SUB	XOR	MUL	ADD
RF		LD	LD	SUB	XOR	MUL
ALU			LD	LD	SUB	XOR
MEM				LD	LD	SUB
WB					LD	LD

reads @ RF

writes @ WB

## Data Hazards

- Consider this instruction sequence:

ADD(C(R1, 1, R2))  
SUBC(R2, 4, R3)  
MUL(R6, R7, R8)  
XOR(R9, R10, R11)

	I	2	3	4	...	5	6
IF	ADDC	SUBC	MUL	XOR			
RF		ADDC	SUBC	MUL	XOR		
ALU			ADDC	SUBC	MUL	XOR	
MEM				ADDC	SUBC	MUL	
WB					ADDC	SUBC	

- SUBC reads R2 on cycle 3, but ADDC does not update it until end of cycle 5 → R2 is stale!
- Pipeline must maintain correct behavior...

## Resolving Stalls:

- 1) Stall. Wait for res to be avail. by freezing earlier stages
- 2) Bypass/Forward. Route data to earlier stage once it's calculated

## Speculate

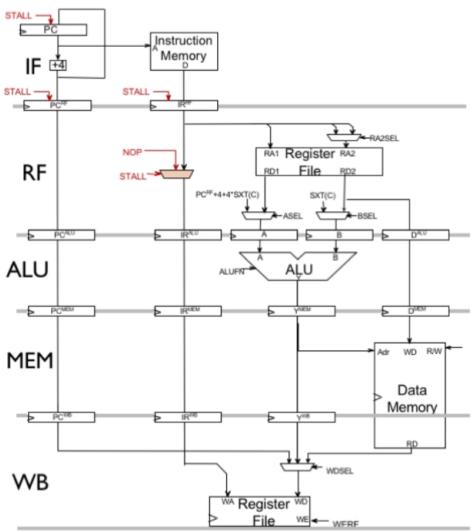
- Guess on value & exec.
- acting/val over.:
  - a) ✓ → continue
  - b) X → Kill & restart

## 1) Stall

	1	2	3	4	5	6	7	8
IF	ADDC	SUBC	MUL	MUL	MUL	MUL	XOR	
RF		ADDC	SUBC	SUBC	SUBC	SUBC	MUL	XOR
ALU			ADDC	NOP	NOP	NOP	SUBC	MUL
MEM				ADDC	NOP	NOP	NOP	SUBC
WB					ADDC	NOP	NOP	NOP

Stalls increase CPI!

## Stall logic



$\text{addc } r1, r1, r2$   
 $\text{subc } r2, r4, r3$   
 $\text{mul } r6, r7, r8$   
 $\text{xor } x9, x7C, x11$

## New Stall Ctrl Signal

$\text{STALL} = 1$

- Disables PC & RF pipe. recs
- injects NOP into ALU

Ctrl logic sets  $\text{STALL} = 1$  if  
src of instr. in RF match w/  
dest on ALU, MEM, WB

$\text{addc } r1, r1, r2$   
 $\text{subc } r2, r4, r3$   
 $\text{mul } r6, r7, r8$   
 $\text{xor } x9, x7C, x11$

## 2) By Pass

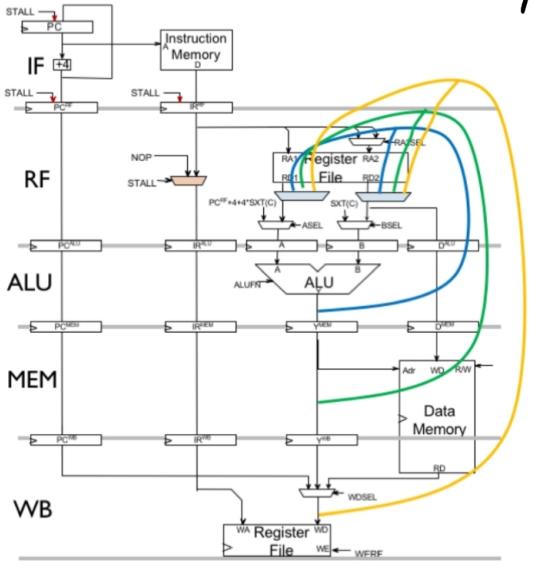
	1	2	3	4	5
IF	ADDC	SUBC	MUL	XOR	
RF		ADDC	SUBC	MUL	XOR
ALU			ADDC	SUBC	MUL
MEM				ADDC	SUBC
WB					ADDC

ADDC result computed ↑

R2 updated

$\text{addc w/ r2 by end of}$   
 $\text{cycles}$   
 $\text{res. onward at end of ALU}$   
 $\text{stage!}$

# Bypass Logic



Promiss marks to RF outputs

route ALU, Mem, WB output to mux inputs

1. Pass Val. if dest at instr. in ALU, Mem, WB  
= Src instr. RF

most if multiple matches?  
Select most recent instr.  
ALU > Mem > WB

## Fully Bypassed Pipeline

Same instr. write PC+R  
Route PC ALU & PC Mem as add.  
Bypass max inputs

BYPASSING EXPENSIVE!  
Wiring & large muxes

Not needed! We can stall  
e.g. just bypass from ALU

w/ FBP, do we need stall?

# Load-to-use Stall

Bypassing cannot elim. load delays  
→ data avail. until WB!

	1	2	3	4	5	6	7
IF	LD	SUBC	MUL	MUL	MUL	XOR	
RF		LD	SUBC	SUBC	SUBC	MUL	XOR
ALU			LD	NOP	NOP	SUBC	MUL
MEM				LD	NOP	NOP	SUBC
WB					LD	NOP	NOP

LD data available ↑      R2 updated ↘

Subc r2, r3

MUL...

XOR...

Summary: Pipelining w/ Data Hazards

1) Stall

Simple, wastes cycles, ↑ CPI

2) Bypass

Expensive, ↓ CPI

Still needs stalls when res. produced after ALU

- can use ↓ bypass, ↑ stall

↑ Pipelining → ↑ Data Hazards  
↓ t<sub>CIK</sub>, ↑ CPI

Compilers can help!

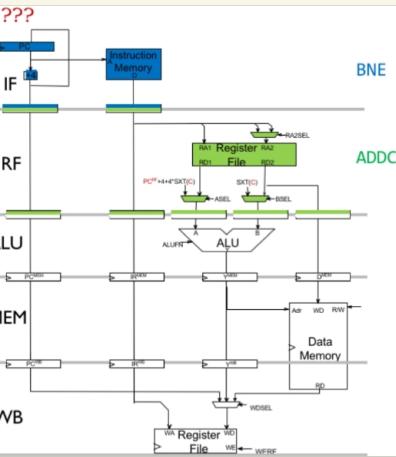
rearrange code to put dependent instr. further apart  
Only works well when compiler can find indep. instrs.  
to move around

Change the ISA? So reg updated w/ 3 instr delay?

ISAs outline implementations.

Bypassing from WB still  
saves a cycle

# Control Hazards



loop: ADDC(R1, 4, R2)  
BNE(R3, loop)  
SUB(R6, R7, R8)

...

How do we set NextPC?

What do we need for NextPC?  
Branch:  $\{ \begin{cases} \text{reg} == ? \\ PC = PC + \text{offset} \end{cases} \}$   
else  $PC += 4$

Jump:  $reg = PC + R$   
 $PC = reg$

Need reg  
unknown until RF

## Resolving Ctrl Hazards

- 1) Stall
- 2) Speculate

Stall

addc r1, -1, r3

⋮

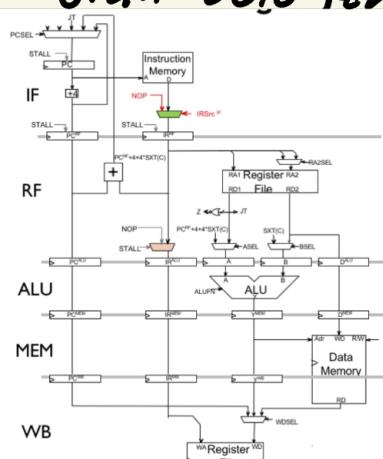
bne r3, loop

If branch/jump in IF, stall HF 1 cycle

(Assume BNE taken in ex code)

	1	2	3	4	5	6	7	8	9
IF	ADDC	MUL	BNE	<b>NOP</b>	ADDC	MUL	BNE	<b>NOP</b>	ADDC
RF		ADDC	MUL	BNE	<b>NOP</b>	ADDC	MUL	BNE	<b>NOP</b>
ALU			ADDC	<b>MUL</b>	BNE	<b>NOP</b>	ADDC	MUL	BNE
MEM				<b>ADDC</b>	MUL	BNE	<b>NOP</b>	ADDC	MUL
WB					ADDC	MUL	BNE	<b>NOP</b>	ADDC

Stall logic for Control hazards



IR Src<sup>IF</sup> control signal

If opcode<sup>RF</sup> == J, B

IR Src<sup>IF</sup> = 1, inject NOP

Set PC sel to J/B

# ISA Issue, Simple vs Comp Branches

Other ISAs  $\geq$  comp branches resolved in ALU

	1	2	3	4	5	6	7	8
IF	ADDC	MUL	BNE	SUB	NOP	ADDC	MUL	BNE
RF		ADDC	MUL	BNE	NOP	NOP	ADDC	MUL
ALU			ADDC	MUL	BNE	NOP	NOP	ADDC
MEM				ADDC	MUL	BNE	NOP	NOP
WB					ADDC	MUL	BNE	NOP

More annulments (but sometimes fewer instructions)

## 2) Speculatc

	1	2	3	4	5	6	7	8	9
IF	ADDC	MUL	BNE	SUB	XOR				
RF		ADDC	MUL	BNE	SUB	XOR			
ALU			ADDC	MUL	BNE	SUB	XOR		
MEM				ADDC	MUL	BNE	SUB	XOR	
WB					ADDC	MUL	BNE	SUB	XOR

Good Guess for NextPC? = PC+4  
Assume BNE not taken

Start fetching at PC+4 (SUB) but ↑  
BNE not resolved yet...

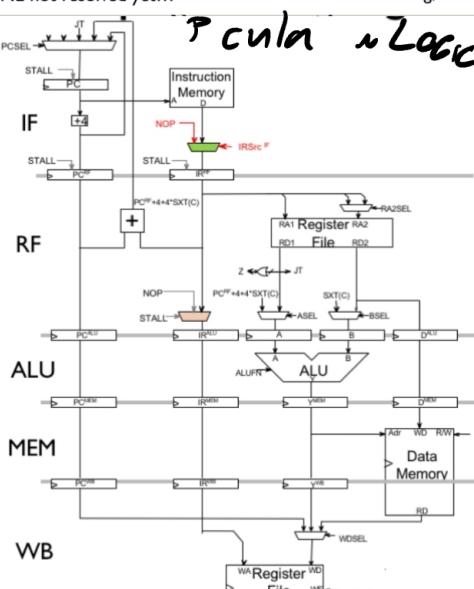
Guessed right, keep going

	1	2	3	4	5	6	7	8	9
IF	ADDC	MUL	BNE	SUB	ADDC	MUL	BNE	SUB	ADDC
RF		ADDC	MUL	BNE	NOP	ADDC	MUL	BNE	NOP
ALU			ADDC	MUL	BNE	NOP	ADDC	MUL	BNE
MEM				ADDC	MUL	BNE	NOP	ADDC	MUL
WB					ADDC	MUL	BNE	NOP	ADDC

Start fetching at PC+4 (SUB) but ↑  
BNE not resolved yet...

Guessed wrong, annul SUB

Assume BNE taken



Familiar: same Path Circuity  
Instead of setting IRSRC<sup>IF</sup>=1 if branch,  
set =1 when branch taken

If op code<sup>RF</sup> = J/B taken

- IRSrc<sup>IF</sup>=1, inject NOP to annul fetched instr. ("branch annulment")
- Set PCsel to J/B

# Branch Prediction

Always Guessing PC<sup>+4</sup> wastes a cycle on taken J/B  
~ 10% CPI

w/ deeper pipelines, taken B wastes many more cycles

Modern CPUs dynamically predict outcome of control-flow instr.

- Predict both B condition / target
- Works well bc B have repeated behavior
  - e.g. B loops taken
  - term/limit/error not taken

## Branch Delay Slots

Change ISA so instr. following J/B always exec.

loop: addc r7, -7, +3

bne r3, loop

mul r4, r5, r6

(→ executes regardless  
of B outcome)

	1	2	3	4	5	6	7	8
IF	ADDC	BNE	MUL	ADDC	BNE	MUL	ADDC	BNE
RF		ADDC	BNE	MUL	ADDC	BNE	MUL	ADDC
ALU			ADDC	BNE	MUL	ADDC	BNE	MUL
MEM				ADDC	BNE	MUL	ADDC	BNE
WB					ADDC	BNE	MUL	ADDC

✓ If compiler can fill in w/ useful instr., no penalty

✗ Can't fill slot ~50% of time → NCP must, larger code

Longer pipeline → more delay slots?

B pred. works better in practice

## Exceptions:

Need to: Save curr PC+4 in XP

Load PC w/ excep. vector (11(CP, Xaddr))

Cause control flow hazards!

Implicit B

Want precise exceptions

✓ Prec. instr. must have completed

instr. causing excep. & future instr. must not have exec

No updates to reg / MM

Simple in single-cycle, CMpx in pipeline

## When Can Exceptions Happen?

Memory fault, illegal instruction, or arithmetic exception  
instr. following exceptions may already be in pipeline  
→ none other written to reg/mem!

## Resolving Exceptions

If instr. except at stage i

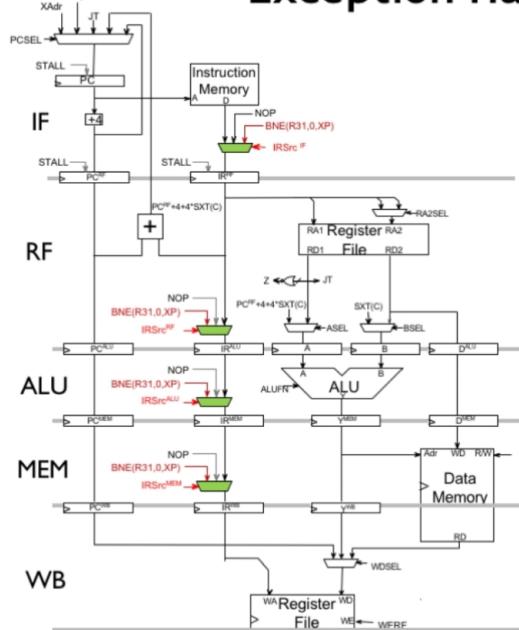
- Turning tr. into B(r31, C, Xp) to save PC+P
- All instr. in stages i-1, ..., 1 (flush PIPE, inc)
- Set PC ← I1Op / Xaddr

Ex: LD w/ mem fault

	1	2	3	4	5	6
IF	LD	ST	MUL	SUB	ADDC	ST
RF		LD	ST	MUL	NOP	ADDC
ALU			LD	ST	NOP	NOP
MEM				LD	NOP	NOP
WB					BNE	NOP

## Exception Handling Logic

### exception handling



IRSrc1IF, RF, ALU, MEM, makes to inject NOP/BNE

- NOP if preceding instr. except.
- BNE if instr. in cur stage has an exception

# Multiple Exceptions?

Id → mem fault  
 ??  
 min  
 sub

	1	2	3	4	5	6
IF	LD	???	MUL	XORC	ADDC	ST
RF	LD	???	NOP	NOP	NOP	ADDC
ALU		LD	BNE	NOP	NOP	
MEM			LD	NOP	NOP	
WB				BNE	NOP	

Invalid opcode detected

Memory fault detected

## Asynchronous Interrupts

Interrupts are easier:

```
// Interrupted code:
...
LD(...)
ADD(...)
SUB(...) ← Taken
HERE
...
// Interrupt handler:
XAddr: OR(...)

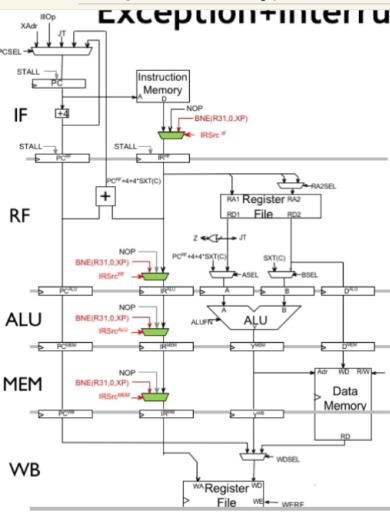
...
SUBC(xp,4,xp)
JMP(xp)
```

	1	2	3	4
IF	ADD	<b>BNE</b>	OR	...
RF	LD	ADD	<b>BNE</b>	OR
ALU		LD	ADD	<b>BNE</b>
MEM			LD	ADD
WB				LD

- Suppose interrupt req. while sub in IF (cycle 2)  
**Handling:**
  - replace sub w/ BNE(..., XP)
  - select Xaddr as next PC
- code handler to return SUB instr.  
 ADD & earlier insts. unaffected

## Exception + Interrupt Handling Logic

### exception+interrupt



Same as before

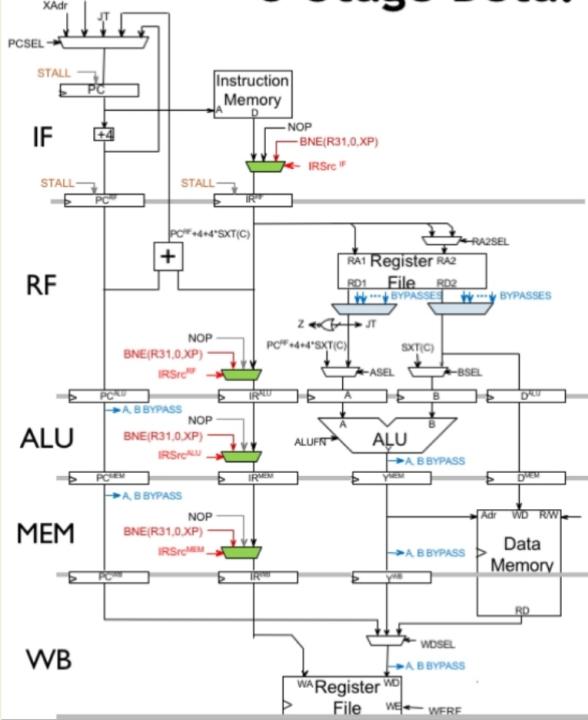
IRSrc<sup>IF, RF, ALU, MEM</sup> muxes to inject NOP/BNE

- NOP if exec. instr. has an except.

- BNE if instruc. in curr stage has an except

Use IRSrc<sup>IF</sup> mux to inject BNE on an interrupt  
 (same as an except. in IF)

# 5-stage CPU: Final Version



## Data hazards

- Stall IF, RF (STALL=1)
  - + IR Src<sup>IF</sup> = NOP

## Control hazards:

### Speculative PC+R &

- JMP or Taken B in RF
- IR Src<sup>IF</sup> = NOP
- $PC Sel \rightarrow JT / B + \text{target}$

### If except. at stage X

$$IR Src^X = BNG$$

$$\text{Prev } IR Src^X = NOP$$

$$PC Sel = X_{addr} \text{ or } 1110_2$$

### 1} interrupt

$$IR Src^I = B$$

$$PC Sel = X_{addr}$$