

# CS x290/ ECE x100

Pipelining

Spring 2020

Notes copyright © Thomas M. Conte

PIPE-1

1

## New (tiny) ISA for example

- Fixed 32-bit encoding per instruction
  - ◆ Next PC = PC+4 is always the next instruction
- 32, 32-bit registers, plus PC and IR
- Arithmetic (e.g., ADD, AND, OR, SUB, etc):



- ◆ Format: **ADD Rdest, Rsrc1, Rsrc2**
- ◆  $\text{Regs}[\text{dest}] = \text{Regs}[\text{src1}] + \text{Regs}[\text{src2}]$
- Load word:



- ◆ Format: **LW Rdest, Imm(Rsrc1)**
- ◆  $\text{Regs}[\text{dest}] = \text{Memory}[\text{Regs}[\text{src1}] + \text{signextend}(\text{Imm})]$

PIPE-3

3

## New (tiny) ISA (cont.)

### ○ Store word: same encoding as LW

◆ Format: SW Rdest, Imm(Rsrc1)

◆ *NOTE Rdest is a source, etc.!*

◆  $\text{Mem}[\text{Regs}[\text{src1}] + \text{signextend}(\text{Imm})] = \text{Regs}[\text{dest}]$

### ○ Branches



◆ Format: B<cond> Rsrc1, Imm

◆ Example: BEQZ R3, Label

◆ if (Regs[src1] == 0) then

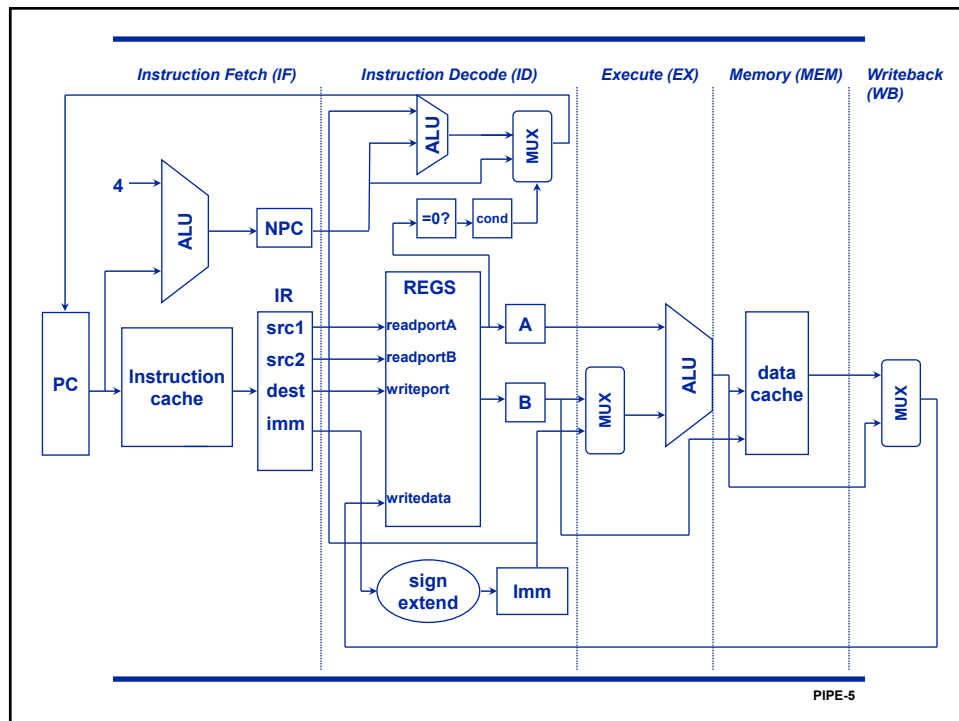
◇  $\text{PC} = \text{NPC} + \text{signextend}(\text{Imm})$

“Imm” here is sometimes called “Pcoffset”

Now let's build the hardware for this...

PIPE-4

4

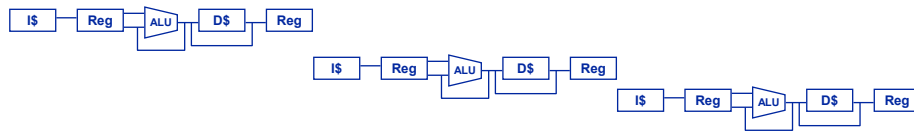


PIPE-5

5

## Analysis

- All instructions (except branch & store):
  - ◆ IF, ID, EX, MEM\*, WB
  - ◆ Each piece takes 1 cycle
- (\*note we don't need MEM for all ops, but we will include it for simplicity and balance)
- CPI = 5
- Graphically (time flowing left to right):

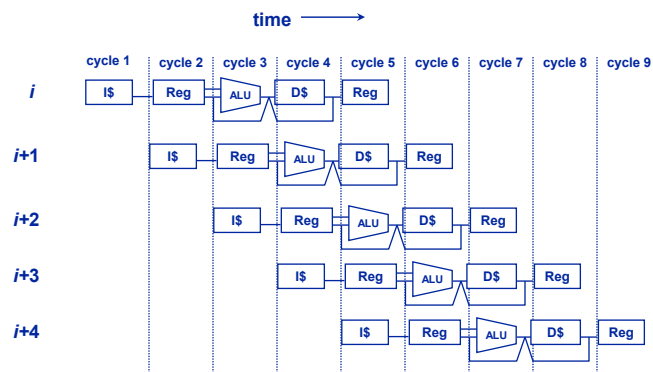


PIPE-6

6

## Pipelining speeds this up!

- Isolate each of IF, ID, EX, MEM, WB with latches
- When instruction  $i$  is in WB,  $i+1$  is in MEM, etc.
- Graphically:



PIPE-7

7

## Pipeline diagrams

Time:	1	2	3	4	5	6	7	8	9
ADD	IF	ID	EX	MEM	WB				
SUB		IF	ID	EX	MEM	WB			
LW			IF	ID	EX	MEM	WB		
ADD				IF	ID	EX	MEM	WB	
SW					IF	ID	EX	MEM	WB

This is a standard way to show pipeline behavior

PIPE-8

8

## Pipeline speedup (no stalls)

- For a pipeline of  $n$  stages for a long program:

$$\begin{aligned}
 \text{speedup} &= \frac{\text{exec. time unpipelined}}{\text{exec time pipelined}} = \frac{IC * CPI_{unpipe} * CT}{IC * CPI_{pipe} * CT} \\
 &= \frac{CPI_{unpipe}}{CPI_{unpipe}/n + T_{latch}} \\
 &= n \quad (\text{ideal case, } T_{latch}=0, \text{ ignore pipeline fill time})
 \end{aligned}$$



(Detail for computer engineers: The latch (level-triggered) could instead be a flip-flop (edge triggered))

PIPE-9

9

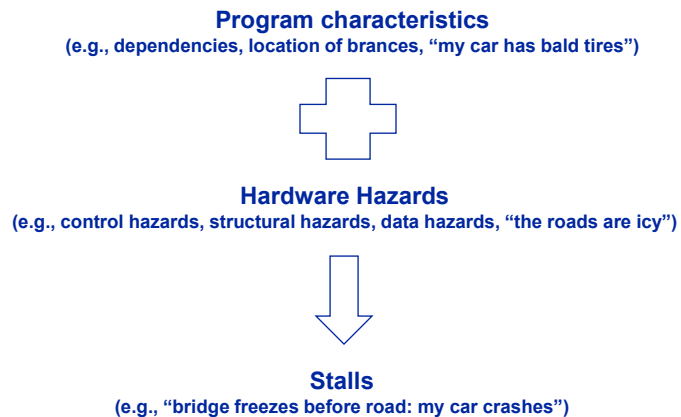
## Pipeline hazards

- A hazard reduces the performance of the pipeline
  - ◆ A problem caused by both hardware design and the program's characteristics
- Three kinds:
  - ◆ Data hazards - Dependencies between instructions prevent their overlapped execution
  - ◆ Structural hazards – Not enough hardware resources for all combinations of instructions
  - ◆ Control hazards - A branch instruction may change the PC, but not until stage 4. What do we fetch before that?

PIPE-10

10

## Program characteristics, Hazards and Stalls



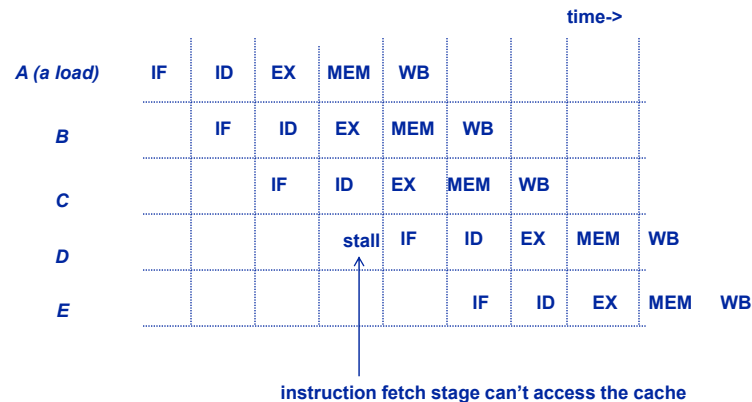
- Not all program characteristics lead to hardware hazards!
- Not all hardware hazards lead to stalls!

PIPE-12

12

## Structural hazards

- Uncommon for our pipeline example
- Consider a pipeline with a unified data+instruction cache:



PIPE-13

13

## Modeling stalls

$$\text{speedup} = \frac{\text{CPI}_{\text{unpipe}}}{\text{CPI}_{\text{pipe}}}$$

$$\text{CPI}_{\text{pipe}} = \text{CPI}_{\text{nostall}} + (\text{ave. stall cycles per instruction})$$

$$= 1 + (\text{ave. stall cycles per instruction})$$

$$= 1 + \text{stalls} \quad \text{where 'stalls' is an abbrev for ave stall cycles per inst.}$$

$$\text{speedup} = \frac{\text{CPI}_{\text{unpipe}}}{\text{CPI}_{\text{pipe}}} = \frac{n}{1 + \text{stalls}}$$

PIPE-14

14



## Delay slots

- Add  $n$  slots to cover  $n$  holes
- ISA is changed to mean “ $n$  instructions after any branch are always executed”
- Problem:
  - ◆ ISA feature that encodes pipeline structure: Difficult to maintain across generations
  - ◆ Typically can fill:
    - ◇ 1 slot 75-80% of time
    - ◇ 2 slots about 30% of time
    - ◇ >2 slots almost never (*fahgettaboutit!*)

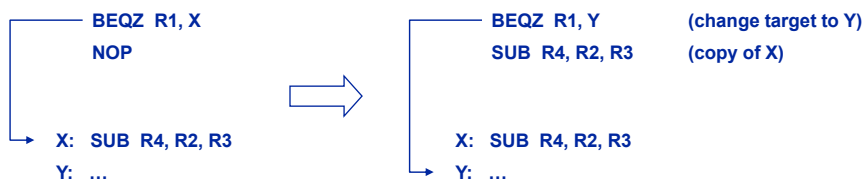


PIPE-18

18

## “Squash” slots

- If you can't fill slot(s) from before the branch, use instructions from either:
  - ◆ Target of branch (if frequently taken)
  - ◆ Fall-through of branch (if frequently not-taken)
- Example: fill from target (branch is frequently taken)



- Disadvantages
  - ◆ Only works if slot instruction is safe to execute (ie, no exceptions occur) when branch goes the opposite (infrequent) direction

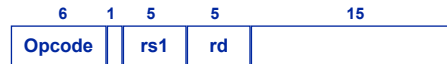
PIPE-19

19



## Squash slots (cont.)

- Each branch instruction now has two opcodes: BEQZ-likely and BEQZ-unlikely



- If compiler thinks branch is frequently taken
  - Compiler sets *likely bit* = 1 (use “BEQZ-likely”)
  - Compiler fills “squash slots” from target
  - Hardware knows to squash slot(s) if branch not-taken
- If compiler thinks branch is frequently not-taken
  - Compiler sets *likely bit* = 0 (use “BEQZ-unlikely”)
  - Compiler fills “squash slots” from fall-through
  - Hardware knows to squash slot(s) if branch taken
- If likely bit set unintelligently, most squash slot instructions must be squashed!

PIPE-20

20

## Methods for setting likely bit

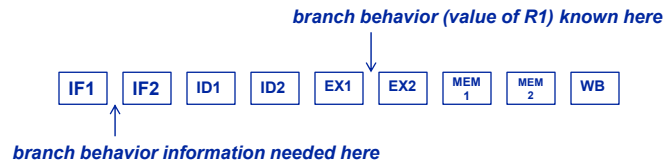
- Likely bit is essentially a *static branch prediction*
  - ◆ Compiler makes a prediction that is fixed for that branch
  - ◆ Likely bit = 1 means predict taken
  - ◆ Likely bit = 0 means predict not-taken
- Static branch prediction methods (compiler branch prediction)
  - ◆ Heuristics
    - ◇ Common one: Backward taken/Forward not taken
      - Could just use the sign bit of the branch offset=> eliminate likely bit from the instruction encoding!
  - ◆ Profiling
    - ◇ Run program once, see what the branches do, recompile and set likely bits

PIPE-21

21

## What about a deeper pipeline?

*hypothetical pipeline of the future: consider BREQZ R1, G*



- Next generation processor wants deeper pipeline
  - ◆ Old “delay slot” code no longer binary compatible!
- Solution:
  - ◆ Do away with delay slots and squashing branches— these go into the dustbin of history!

PIPE-22

22

## Reducing Branch Stalls: The 3 W's

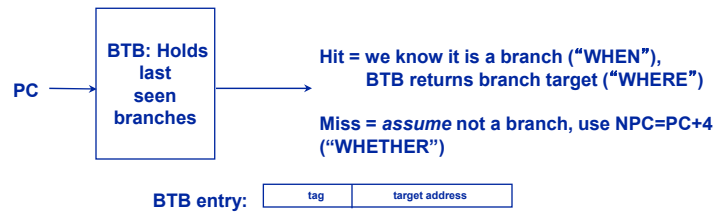
- Issues with Whether to branch:
  - ◆ (For conditional branches) Predict if it will branch or not, *before* execution
- Issues with Where to branch to:
  - ◆ Predicting where the branch will go (if taken)
- Issues with When the branch happens:
  - ◆ When to detect an undecoded instruction is a branch
- Optimal: try to determine all three in IF1 stage
  - ◆ Won't work perfectly (it's a *prediction*), but we can try our best

PIPE-23

23

## Branch Target Buffer

- What does Instruction Fetch stage know about an instruction?
  - ◆ Not much! Only the address of the instruction (PC)
  - ◆ So...Keep table of last known branch targets around
  - ◆ Table is written to by stage that knows the actual branch behavior (e.g., ID2, EX, -- depends on pipeline design)
  - ◆ Traditional name for this is a *Branch Target Buffer (BTB)*

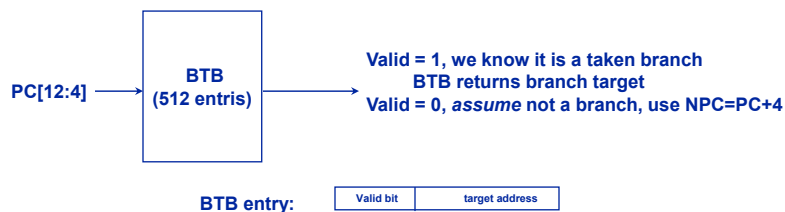


PIPE-24

24

## Branch Target Buffer (cont.)

- Instead of the tag comparison, just hash the PC
- This works because a lookup of the wrong entry just causes a misprediction (ie, a performance hit but the code still works!)
- Example hash: use middle bits from PC, e.g., if table is 512 entries, use bits PC[12:4]



PIPE-25

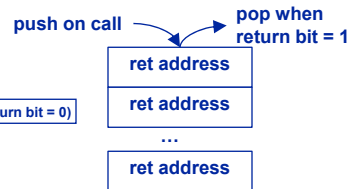
25

## Predicting Where with returns

- Problem: A lot of branches are return instructions – BTB fails on these because target address is dependent on call location – it keeps changing
  - ◆ Holding the last target address is a poor predictor unless the call site is always the same!
- Solution: Keep a hardware “stack” of return addresses
  - ◆ Push return address when a “call” instruction is seen
  - ◆ Pop buffer on returns to get prediction
  - ◆ Empirically need 4 to 8 entries for integer code

Each entry in BTB now contains:

Tag bits or valid bit	return bit	Target address (only when return bit = 0)
-----------------------	------------	---



PIPE-26

26

## About the return address stack

Problem: stack is a fixed size, not  $\infty$

Two program behavior scenarios:

1. Shallow call depth  $\leq N$
2. Very deep call depth  $\gg N$

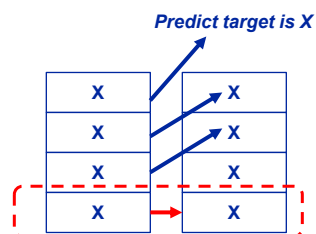
A hardware stack of depth N will do a good job in scenario 1, but it will not in scenario 2.

When does scenario 2 happen?

- It's most common in *recursive code*
- In recursive code, the stack grows very deep
- BUT, the stack contains the same value (e.g., X)
  - stack would be full of many copies of X!

**Solution: Always Copy the bottom of the hardware stack on Pops**

Recursive code:  
call fact(n-1)  
X: next instruction



27

---

## Issues with Whether

- Predicting conditional branches
  - ◆ And sometimes unconditional branches if needed before the instruction's opcode is decoded
  - ◆ Our BTB does the equivalent of “whatever happened last” or a 1-bit scheme

---

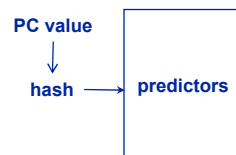
PIPE-28

28

---

## Hardware branch prediction

- OFTEN is a separate table, the branch prediction buffer:



- Only answers the “WHETHER branch is taken or not” question
- Typically it has no tags and conflicts are allowed

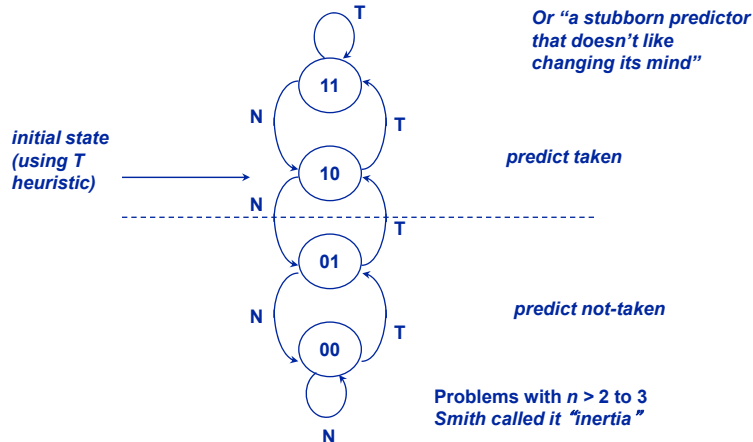
---

PIPE-29

29

## Smith $n$ -bit “bimodal” counter predictor

- Replace prediction *bit* with  $n$ -bit counter:



PIPE-30

30

## Smith counter cont.

- Why it works:
  - ◆ Conditional branches at the end of a loop are taken often but only not taken once
    - ◇ “TTTTTNTTTTN...”
    - ◇ Counter would be 01 11 11 11 11 10 11 11 11 11 10
    - ◇ The infrequent N does not change the prediction
  - ◆ Conditional branches that guard an error condition are rare:
    - ◇ If (error) { handle error }
    - ◇ NNNNNNTNNNN...
    - ◇ The rare T does not change the prediction
  - ◆ Smith called this property “inertia”
- Ravi Nair of IBM Research simulated all possible 2-bit state machines and found variants of the Smith counter were the best predictors
  - ◆ Variation was the initial state choice

PIPE-31

31

## Example of Smith counter

	<u>T</u>	<u>T</u>	<u>T</u>	<u>T</u>	<u>T</u>	<u>T</u>	<u>T</u>	<u>N</u>	<u>T</u>	<u>T</u>	<u>N</u>	<u>N</u>	<u>N</u>	<u>N</u>	<u>N</u>	<u>N</u>	<u>N</u>	<u>N</u>	<u>T</u>	<u>T</u>
previous state	10	11	11	11	11	11	11	10	11	11	10	01	00	00	00	00	00	00	01	01
new state	11	11	11	11	11	11	10	11	11	10	01	00	00	00	00	00	00	01	10	10

5 mispredictions out of 19 branch executions

	<u>N</u>	<u>T</u>	<u>N</u>	<u>T</u>	<u>N</u>	<u>T</u>	<u>N</u>	<u>T</u>	<u>N</u>	<u>T</u>	<u>N</u>	<u>T</u>	<u>N</u>	<u>T</u>	<u>N</u>	<u>T</u>	<u>N</u>	<u>T</u>	
previous state	10	01	10	01	10	01	10	01	10	01	10	01	10	01	10	01	10	01	
new state	01	10	01	10	01	10	01	10	01	10	01	10	01	10	01	10	01	10	

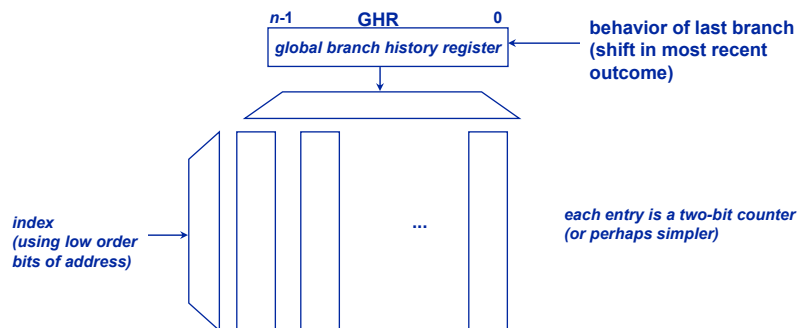
19 mispredictions out of 19: the infamous "toggle branch"

PIPE-32

32

## Improving the Smith counter

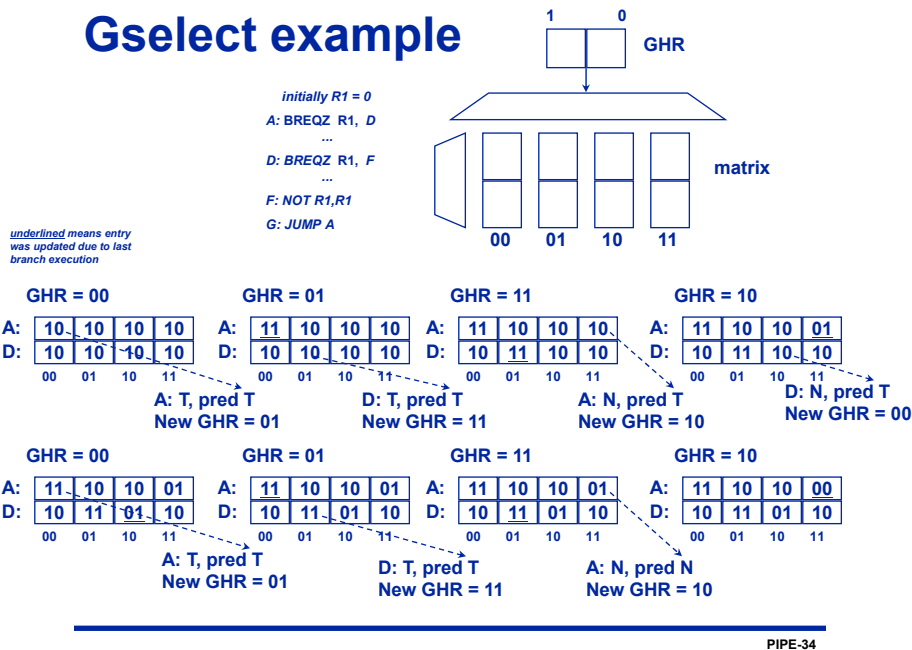
- Idea: Capture correlations between different branches
- Gselect:



PIPE-33

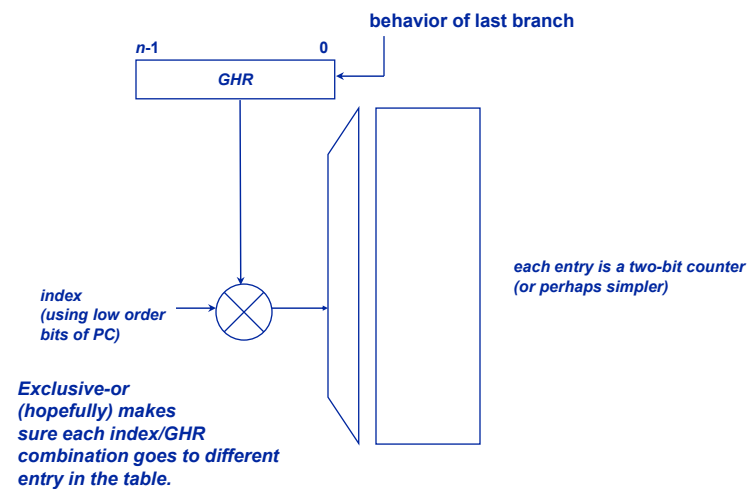
33

## Gselect example



34

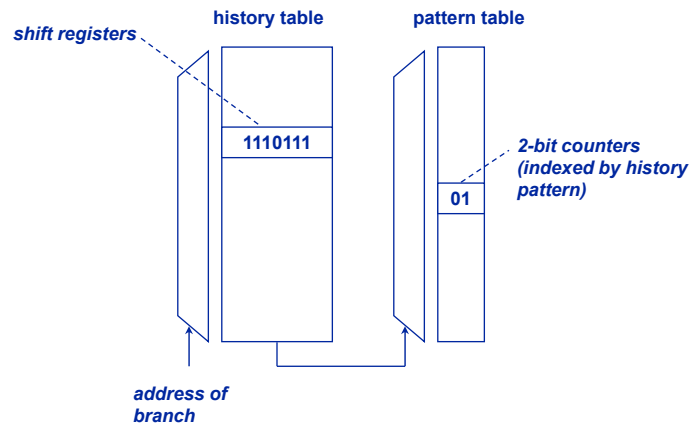
## Another Implementation: Gshare



35



## Yeh/Patt predictor



*Idea: Associate predictions with branch histories, not branch addresses (use different indexing scheme)*

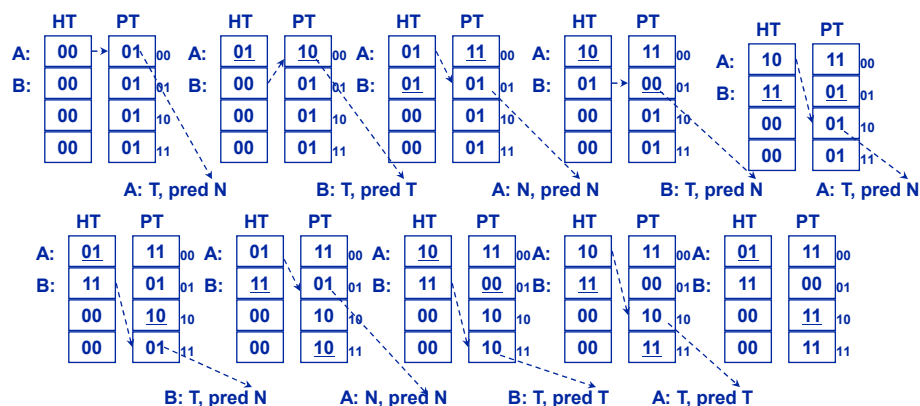
PIPE-36

36

## Yeh/Patt Example: toggle branch

A: TNTNTNTNTNTNTNTNTN  
B: TTTTTTTTTTTTTTTTTTT

Note: We're using Not-Taken heuristic to initialize counters to 01 (weakly not-taken).



4 mispredictions in this case.  
In general: provides 96-98% accuracy for integer code

PT entries 01, 10 are "trained" for A  
and 11 is "trained" for B

PIPE-37

37

## Importance of branch prediction

$CPI = 1 + \text{stalls}$

If a branch is seen 20% of the time, and

The branch predictor is A% accurate, and

The penalty for a misprediction (pipeline depth dependent) is P, then

$CPI = 1 + 0.2 \cdot (1 - A) \cdot P$

Say A = 90% and P = 3 cycles, CPI = 1.06

Now say P = 6, CPI = 1.12

**Remember that CPU Time = IC \* CPI \* CT**

Want low CPI (good branch prediction) and low CT (deep pipelines)

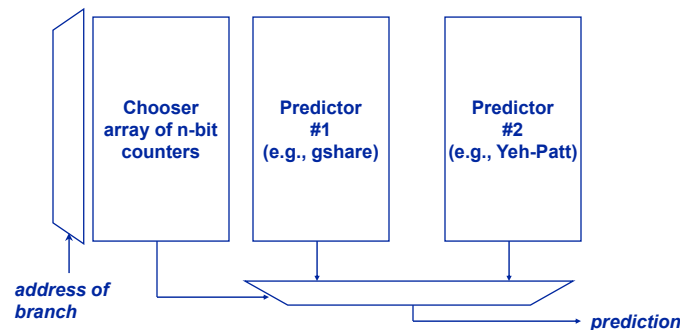
**Branch prediction accuracy allows deeper pipelines, thus faster processors!**

ECE 463/521, Profs Conte/Rotenberg/Sair, Dept. of ECE, NC State University

PIPE-38

38

## Hybrid predictors



- Both predictors supply a prediction-- pipeline uses only one
- Chooser ("confidence predictor") updated based on which predictor was correct
  - ◆ Increment chooser counter if #1 was correct, decrement if #2 was correct

PIPE-39

39

## Program Characteristics: Types of dependencies

### ○ True-dependence (pure-dependence, flow-dependence)

- ◆ ADD R1,R2,R3
- ◆ SUB R4,R5,R1
- ◆ Stall if pipeline has RAW hazards

### ○ Anti-Dependence

- ◆ ADD R3,R2,R1
- ◆ SUB R1,R4,R5
- ◆ Stall if pipeline has WAR hazards
- ◆ Due to reuse: Removed by using another register

### ○ Output-Dependence

- ◆ ADD R1,R2,R3
- ◆ SUB R1,R4,R5
- ◆ Stall if pipeline has WAW hazards
- ◆ Due to reuse: Removed by using another register

PIPE-40

40

## Data Hazards

### ○ Recall that hazards reduce the potential speedup of pipelining

- ◆ Dependencies between instructions prevent their overlapped execution

### ○ Kinds of data hazards

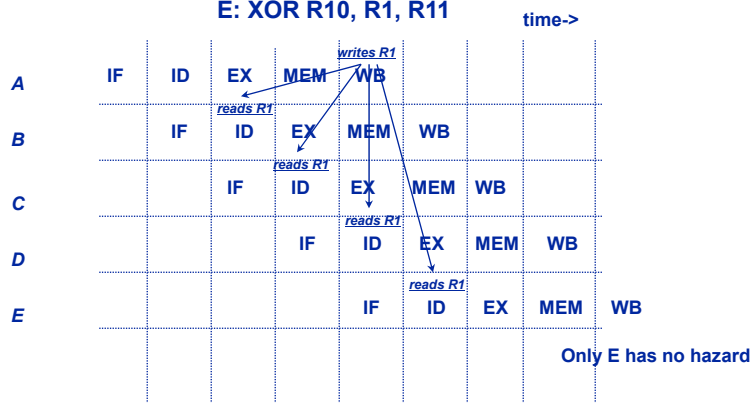
- ◆ Read After Write (RAW)
  - ◇ Caused by pure dependence
- ◆ Write after Read (WAR)
  - ◇ Caused by anti-dependence
- ◆ Write After Write (WAW)
  - ◇ Caused by output-dependence
- ◆ Read after Read (RAR)???
- ◇ Really a structural hazard on read ports

PIPE-41

41

## Data hazards

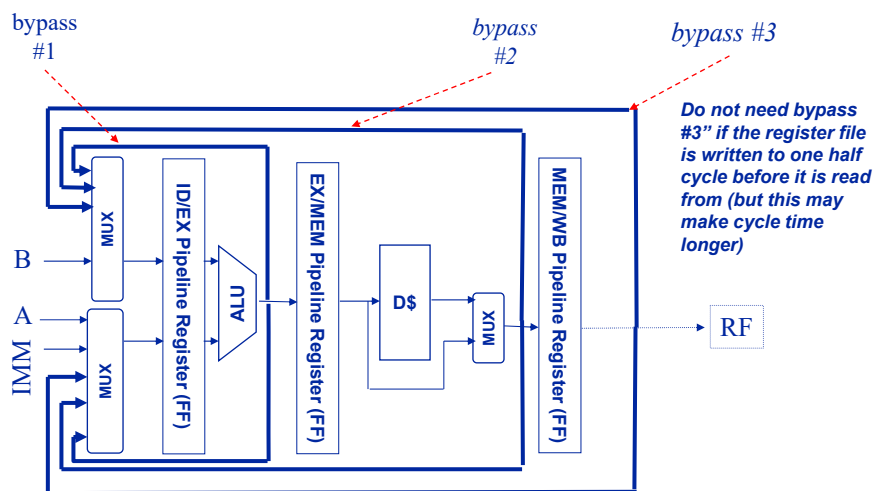
A: ADD R1, R2, R3  
B: SUB R4, R5, R1  
C: AND R6, R1, R7  
D: OR R8, R1, R9  
E: XOR R10, R1, R11



PIPE-43

43

## Data forwarding

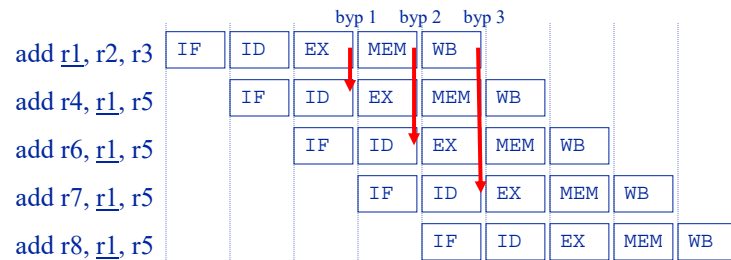


PIPE-44

44

## Data forwarding (cont.)

- Need bypass3 ONLY if we *do not* follow the pipeline model where we *write registers during the first half of a clock cycle and read them in the second*

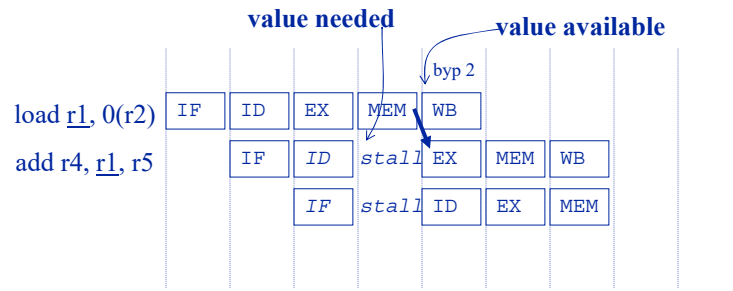


PIPE-45

45

## Stalls due to data hazards

- In our pipeline:
  - ◆ Most RAW hazards => no stall
  - ◆ Loads cause 1-cycle stall (cache hit)
    - ◇ This RAW hazard is so common, it's sometimes given its own name, a "load-use hazard"
  - ◆ Can "feature" the hazard by telling programmer that load values are not available until AFTER the next instruction – called a "Load Delay Slot"



PIPE-46

46

## Other data hazards

### ○ WAR (write-after-read)

- ◆ A ADD r1, r2, r3
- ◆ B ADD r2, r4, r5
- ◆ Hazard if B writes R2 before A reads R2
- ◆ Modify example pipeline so that lcache misses do not stall instruction fetch, and you may get this hazard to occur:



PIPE-47

47

## Other data hazards (cont.)

### ○ WAW (write-after-write)

- ◆ A ADD r1, r2, r3
- ◆ B ADD r1, r4, r5
- ◆ Hazard if B writes R1 before A writes R1
- ◆ Result: later instructions see wrong value in the register
- ◆ Occurs if instructions can write register file out-of-order
- ◆ Using the same pipeline that does not stall fetch on lcache miss:



PIPE-48

48

---

## Other data hazards (cont.)

- Handling WAR/WAW hazards
  - ◆ Stall the later instruction (stalls in WB stage)
  - ◆ Compiler: don't reuse registers (not always possible, because you run out of registers...)
  - ◆ Hardware: register renaming (see *next major topic – ILP*)

---

PIPE-49

49

---

## And finally...

- CPU time = IC x CPI x CT
- $CPI = 1 + \text{stalls}$   
 $= 1 + \text{structural\_stalls} + \text{br\_stalls} + \text{data\_stalls}$
- Best CPI is CPI = 1
- ... Right???

---

PIPE-50

50