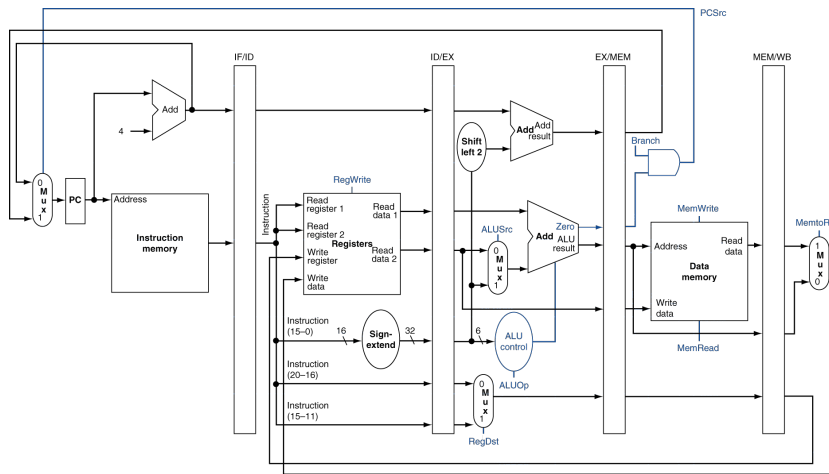


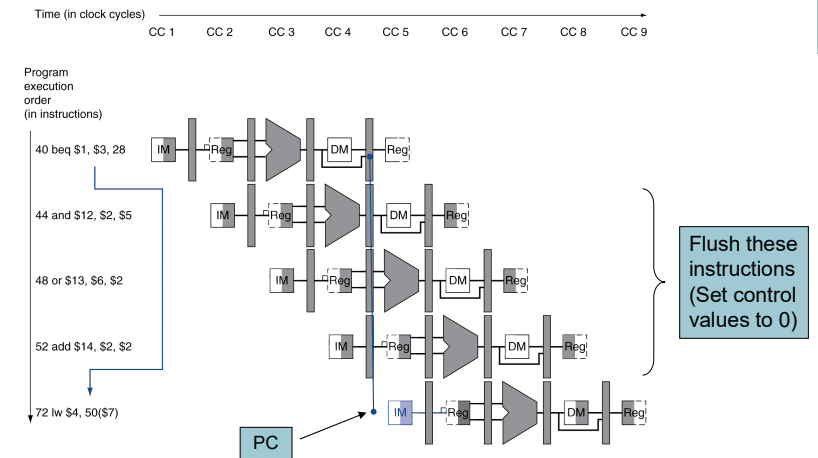
## Pipelined Control (Simplified)



Chapter 4 — The Processor — 2

## Branch Hazards

- If branch outcome determined in MEM



Flush these instructions  
(Set control values to 0)

Chapter 4 — The Processor — 3

## Reducing Branch Delay

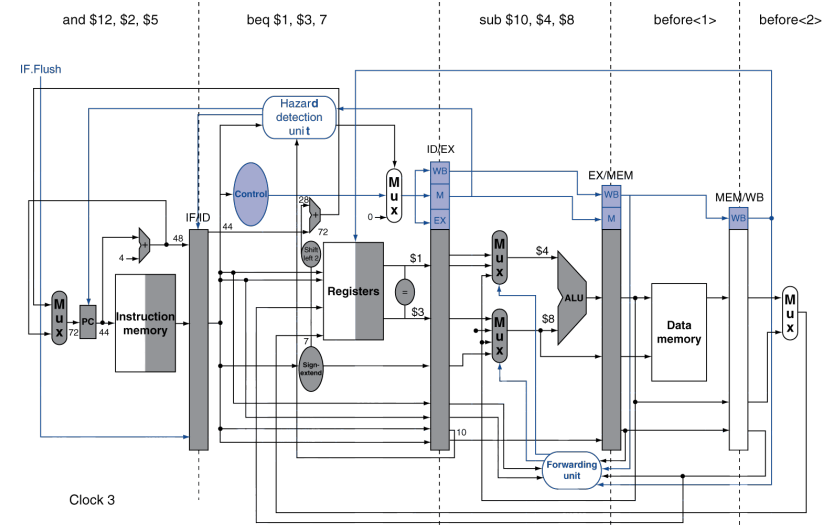
- Move hardware to determine outcome to ID stage
  - Target address adder
  - Register comparator
- Example: branch taken
 

```

36: sub $10, $4, $8
40: beq $1, $3, 7
44: and $12, $2, $5
48: or $13, $2, $6
52: add $14, $4, $2
56: slt $15, $6, $7
...
72: lw $4, 50($7)
      
```

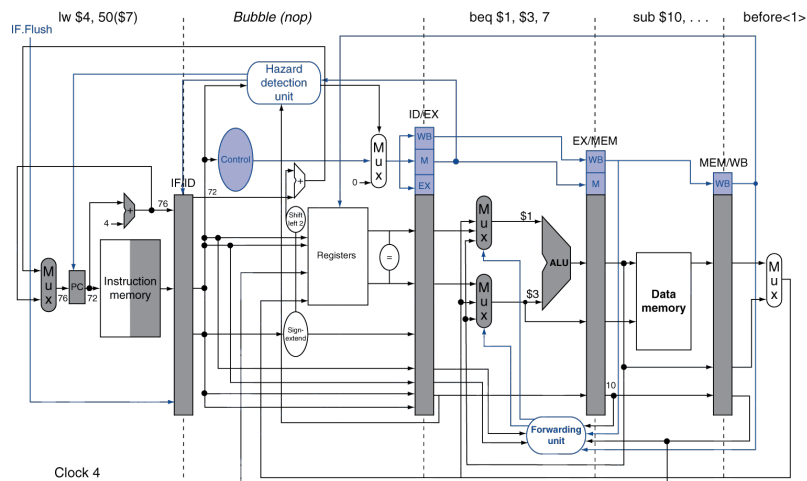
Chapter 4 — The Processor — 4

## Example: Branch Taken



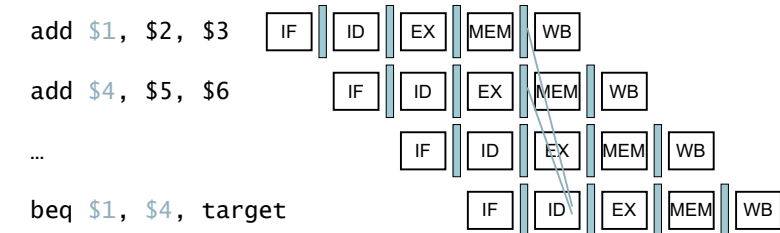
Chapter 4 — The Processor — 5

## Example: Branch Taken



## Data Hazards for Branches

- If a comparison register is a destination of 2<sup>nd</sup> or 3<sup>rd</sup> preceding ALU instruction

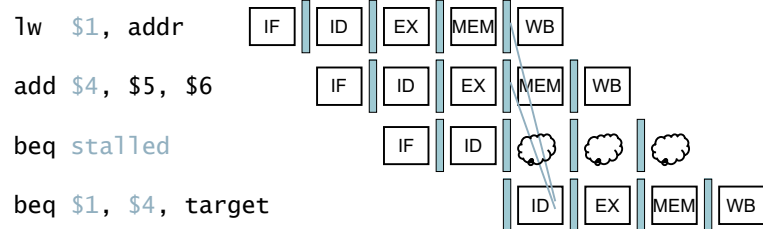


- Can resolve using forwarding



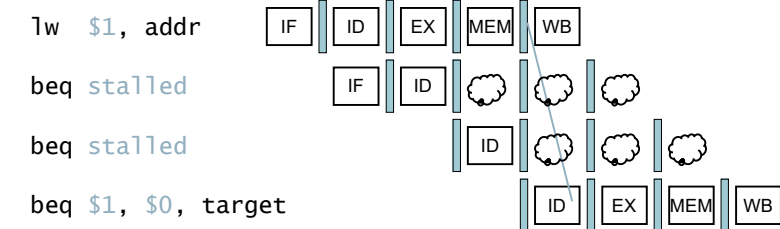
## Data Hazards for Branches

- If a comparison register is a destination of preceding ALU instruction or 2<sup>nd</sup> preceding load instruction
  - Need 1 stall cycle



## Data Hazards for Branches

- If a comparison register is a destination of immediately preceding load instruction
  - Need 2 stall cycles



## Branch (Control) Hazards and Their Impact on the Pipeline Performance

- Control hazards can cause a greater performance loss for a pipelined implementation than data hazards. In first approximation:

$$\text{speedupFromPipelining} = \frac{\text{pipelineDepth}}{1 + [\text{branchFrequency} \times \text{branchPenalty}]}$$

- performance loss can vary between 10% and 30% depending on the branch frequency
- generally, the deeper the pipeline, the worse the branch penalty (measured in clock cycles)
- Various strategies to reduce the **branch penalty**
  - static (compile time) schemes
    - they are fixed for each branch during the entire execution
  - dynamic
    - hardware and software techniques
    - more sophisticated branch prediction schemes

## Reducing Pipeline Branch Penalties - I: Freezing (Flushing) the Pipeline

- Assuming now that the target address calculation and the comparison are performed during the ID stage, then, in case of "branch taken", the PC is changed at the end of the ID stage
  - simplest strategy:** freezing the pipeline as soon as the instruction is *detected as a branch* and regardless of its outcome (**branch penalty is fixed, SW cannot reduce it**)

instruction	1	2	3	4	5	6	7	8	9
BNEZ R5, target	IF	ID							
fall-through inst		IF							
branch target			IF	ID	EX	MEM	WB		
branch target + 1				IF	ID	EX	MEM	WB	
branch target + 2					IF	ID	EX	MEM	WB

## Reducing Pipeline Branch Penalties - II: Predicting "Non-Taken"

- If the branch is taken, all the fetched instruction must be turned into NOPs and it is necessary to restart the fetch at the target address. The state must be unaffected!!

instruction	1	2	3	4	5	6	7	8	9
<b>untaken</b> branch	IF	ID							
Fall-Through Inst.		IF	ID	EX	MEM	WB			
FTI + 1			IF	ID	EX	MEM	WB		
FTI + 2				IF	ID	EX	MEM	WB	
FTI + 3					IF	ID	EX	MEM	WB

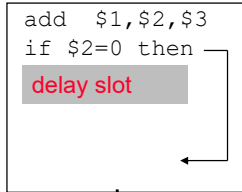
instruction	1	2	3	4	5	6	7	8	9
<b>taken</b> branch	IF	ID							
Fall-Through Inst.		IF							
branch target			IF	ID	EX	MEM	WB		
branch target + 1				IF	ID	EX	MEM	WB	
branch target + 2					IF	ID	EX	MEM	WB

## Delayed Branches

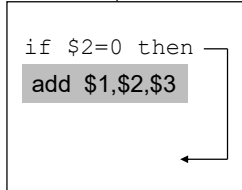
- If the branch hardware has been moved to the ID stage, then we can eliminate all branch stalls with **delayed branches** which are defined as always executing the next sequential instruction after the branch instruction – the branch takes effect **after** that next instruction
  - MIPS compiler moves an instruction to immediately after the branch that is not affected by the branch (a **safe** instruction) thereby **hiding** the branch delay
- With deeper pipelines, the branch delay grows requiring more than one delay slot
  - Delayed branches have lost popularity compared to more expensive but more flexible (dynamic) hardware branch prediction
  - Growth in available transistors has made hardware branch prediction relatively cheaper

## Scheduling Branch Delay Slots

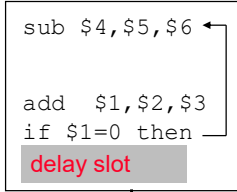
A. From before branch



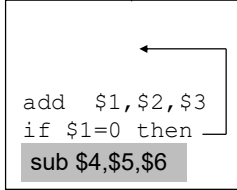
becomes



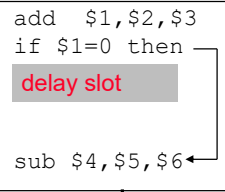
B. From branch target



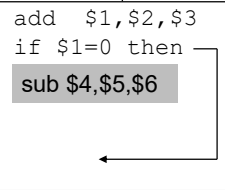
becomes



C. From fall through



becomes



- ❑ A is the best choice, fills delay slot and reduces IC
- ❑ In B and C, the `sub` instruction may need to be copied, increasing IC
- ❑ In B and C, must be okay to execute `sub` when branch fails

CSE431 Chapter 4B.14

Irwin, PSU, 2008

## Basic Branch Prediction: Branch-History Table (or Branch-Prediction Buffer)

### Branch-History Table

- 1-bit prediction scheme
  - remembers last outcome of some branches
- implemented as a small memory indexed by a lower portion (some least-significant bits) of the address of the branch instruction
  - like a direct-map cache where an access is always a hit
- useful only to reduce the branch delay when longer than the time to compute the possible target PC
- We don't even know if prediction is correct. It is just a hint (if wrong the bit is inverted and we restart from the correct PC)
  - e.g., branch instructions at the following addresses share the BHT entry
    - 00F2BC 0101 1100
    - 010A5D 1001 1100

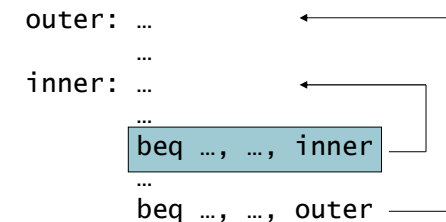
Index	Taken?
0000	1
0001	0
0010	1
0011	0
0100	1
0101	0
0110	1
0111	1
1000	1
1001	0
1010	1
1011	1
1100	0
1101	1
1110	1
1111	1

## Dynamic Branch Prediction

- In deeper and superscalar pipelines, branch penalty is more significant
- Use dynamic prediction
  - Branch prediction buffer (aka branch history table)
  - Indexed by recent branch instruction addresses
  - Stores outcome (taken/not taken)
  - To execute a branch
    - Check table, expect the same outcome
    - Start fetching from fall-through or target
    - If wrong, flush pipeline and flip prediction

## 1-Bit Predictor: Shortcoming

- Inner loop branches mispredicted twice!



- Mispredict as taken on last iteration of inner loop
- Then mispredict as not taken on first iteration of inner loop next time around



## Example: Accuracy of 1-Bit Prediction Scheme in the Presence of "Loop Branching"

```

loop: L.D F0, 0(R1)
      MUL.D F4, F0, F2
      S.D F4, 0(R1)
      DADDIU R1, R1, #-8
      BNEZ R1 loop
  
```

- Assumptions
  - R1 is initialized to #80

Iteration	0	1	2	3	4	5	6	7	8	9	10	0	1	2	3	4	5	6	7	8	9	10
Predicted behavior	N	T	T	T	T	T	T	T	T	T	T	N	T	T	T	T	T	T	T	T	T	T
Actual behavior	T	T	T	T	T	T	T	T	T	T	N	T	T	T	T	T	T	T	T	T	N	T

- Branch taken 90% of the times, but branch prediction accuracy is only 80%
- A 1-bit predictor for "loop branches" mispredicts at twice the rate that the branch is not taken

## Example: 2-Bit Prediction Scheme in Action with "Loop Branching"

```

loop: L.D F0, 0(R1)
      MUL.D F4, F0, F2
      S.D F4, 0(R1)
      DADDIU R1, R1, #-8
      BNEZ R1 loop
  
```

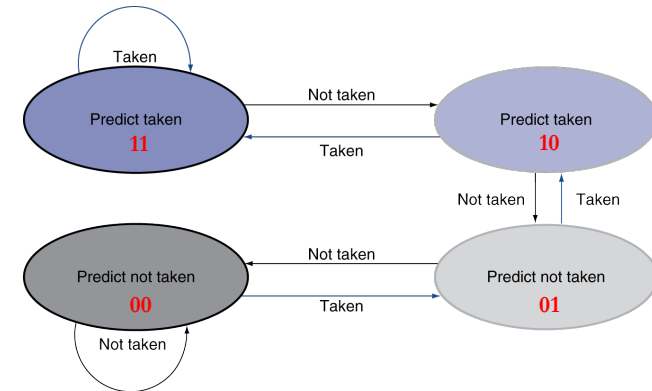
- Assumptions
  - R1 is initialized to #80

Iterations & steps	0	1	2	3	4	5	6	7	8	9	10	0	1	2	3	4	5	6	7	8	9	10
Predicted behavior	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T
Actual behavior	T	T	T	T	T	T	T	T	T	T	N	T	T	T	T	T	T	T	T	T	N	T

- Branch taken 90% of the times, and branch prediction accuracy is now 90%
- The 2-bit predictor mispredicts at the 10<sup>th</sup> step of the 1<sup>st</sup> iteration, but *doesn't change its mind*. It just moves to "weak-predict-taken" for the 1<sup>st</sup> step of the following iteration

## 2-Bit Predictor

- Only change prediction on two successive mispredictions



## Calculating the Branch Target

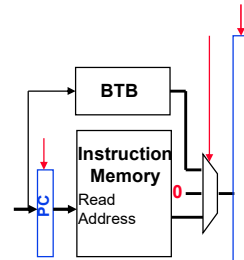
- Even with predictor, still need to calculate the target address
  - 1-cycle penalty for a taken branch
- Branch target buffer
  - Cache of target addresses
  - Indexed by PC when instruction fetched
    - If hit and instruction is branch predicted taken, can fetch target immediately



## Branch Target Buffer

❑ The BHT predicts **when** a branch is taken, but does not tell **where** its taken to!

- A **branch target buffer (BTB)** in the IF stage caches the branch target address, but we also need to fetch the next sequential instruction. The prediction bit in IF/ID selects which “next” instruction will be loaded into IF/ID at the next clock edge
  - Would need a two read port instruction memory



- Or the BTB can cache the branch taken **instruction** while the instruction memory is fetching the next sequential instruction

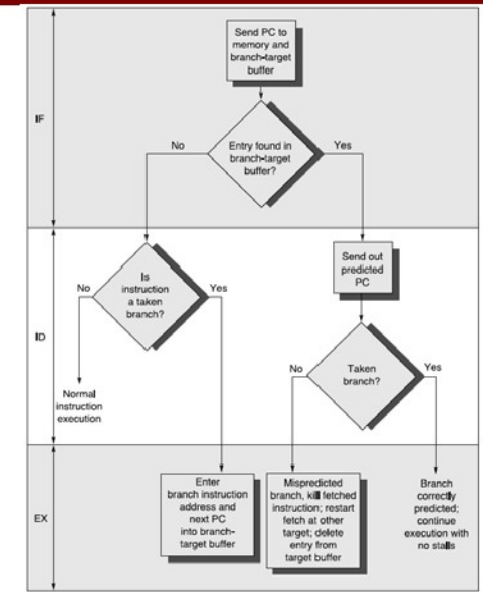
❑ If the prediction is correct, stalls can be avoided no matter which direction they go

CSE431 Chapter 4B.22

Irwin, PSU, 2008

## Predicting the Instruction Target Address: Branch-Target Buffer (or Branch-Target Cache)

- Goal: learn the **predicted instruction address at the end of IF stage**
  - one cycle earlier
    - at best, branch-prediction buffers know the next predicted instruction at the end of ID
- No branch delay
  - if entry found in buffer and prediction is correct
  - if entry not found and branch is not taken
- Two cycle penalty
  - if prediction is incorrect
  - if entry is not found and branch is taken



## Branch Target Buffer: Penalty Table

Instruction in buffer	Prediction	Actual Branch	Penalty Clock Cycles
yes	taken	taken	0
yes	taken	not taken	2
no		taken	2
no		not taken	0

- Assuming
  - 90% buffer hit rate, 85% prediction accuracy, 60% branches taken
- $P(\text{branch in buffer but not taken}) = 0.9 \times 0.15 = 0.135$
- $P(\text{branch not in buffer, but taken}) = 0.1 \times 0.6 = 0.06$
- Total branch penalty =  $(0.135 + 0.06) \times 2 = 0.39$ 
  - the penalty for delayed-branch was about 0.5 clock cycles
  - penalty is lower with better predictors (and bigger branch delays)