

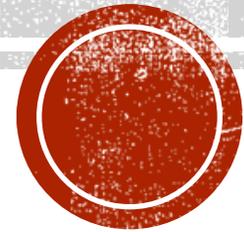
ARCHITECTURE OF COMPUTER SYSTEMS

LECTURE 5 - PIPELINING

II

(BRANCHES,

EXCEPTIONS)



LAST TIME IN LECTURE 4

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

Increases because of pipeline bubbles

 Reduces because fewer logic gates on critical paths between flip-flops

- Pipelining increases clock frequency, while growing CPI more slowly, hence giving greater performance
- Pipelining of instructions is complicated by HAZARDS:
 - Structural hazards (two instructions want same hardware resource)
 - Data hazards (earlier instruction produces value needed by later instruction)
 - Control hazards (instruction changes control flow, e.g., branches or exceptions)
- Techniques to handle hazards:
 - 1) Interlock (hold newer instruction until older instructions drain out of pipeline and write back results)
 - 2) Bypass (transfer value from older instruction to newer instruction as soon as available somewhere in machine)
 - 3) Speculate (guess effect of earlier instruction)

CONTROL HAZARDS

What do we need to calculate next PC?

- For Jumps
 - Opcode, PC and offset
- For Jump Register
 - Opcode, Register value, and PC
- For Conditional Branches
 - Opcode, Register (for condition), PC and offset
- For all other instructions
 - Opcode and PC (and have to know it's not one of above)

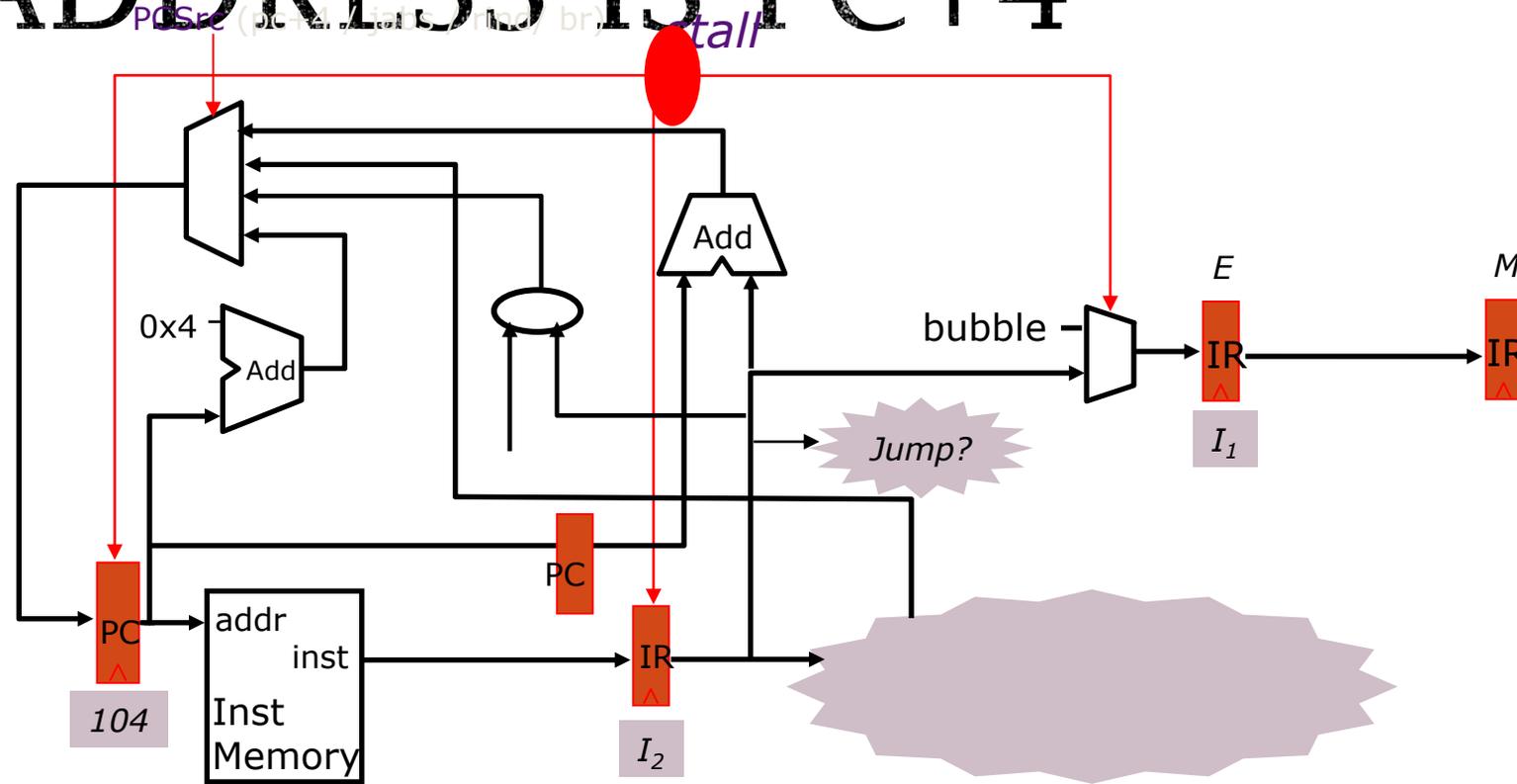
PC CALCULATION BUBBLES

	<i>time</i>										
	t0	t1	t2	t3	t4	t5	t6	t7		
(I ₁) $x1 \leftarrow x0 + 10$	IF ₁	ID ₁	EX ₁	MA ₁	WB ₁						
(I ₂) $x3 \leftarrow x2 + 17$		IF ₂	IF ₂	ID ₂	EX ₂	MA ₂	WB ₂				
(I ₃)				IF ₃	IF ₃	ID ₃	EX ₃	MA ₃	WB ₃		
(I ₄)						IF ₄	IF ₄	ID ₄	EX ₄	MA ₄	WB ₄

	<i>time</i>										
	t0	t1	t2	t3	t4	t5	t6	t7		
IF	I ₁	-	I ₂	-	I ₃	-	I ₄				
ID		I ₁	-	I ₂	-	I ₃	-	I ₄			
EX			I ₁	-	I ₂	-	I ₃	-	I ₄		
MA				I ₁	-	I ₂	-	I ₃	-	I ₄	
WB					I ₁	-	I ₂	-	I ₃	-	I ₄

- ⇒ *pipeline bubble*

SPECULATE NEXT ADDRESS IS PC+4

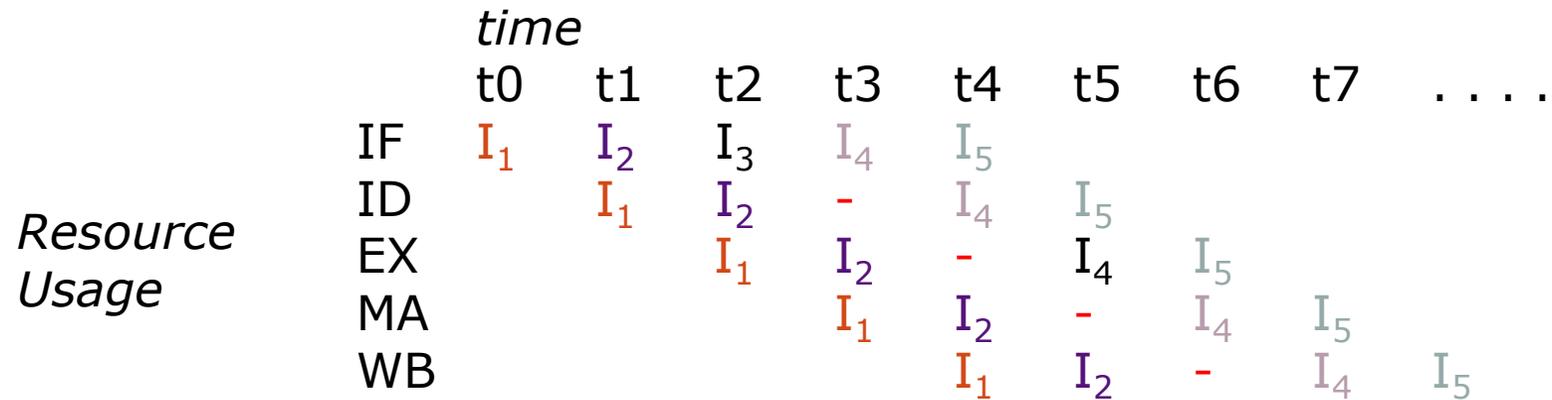
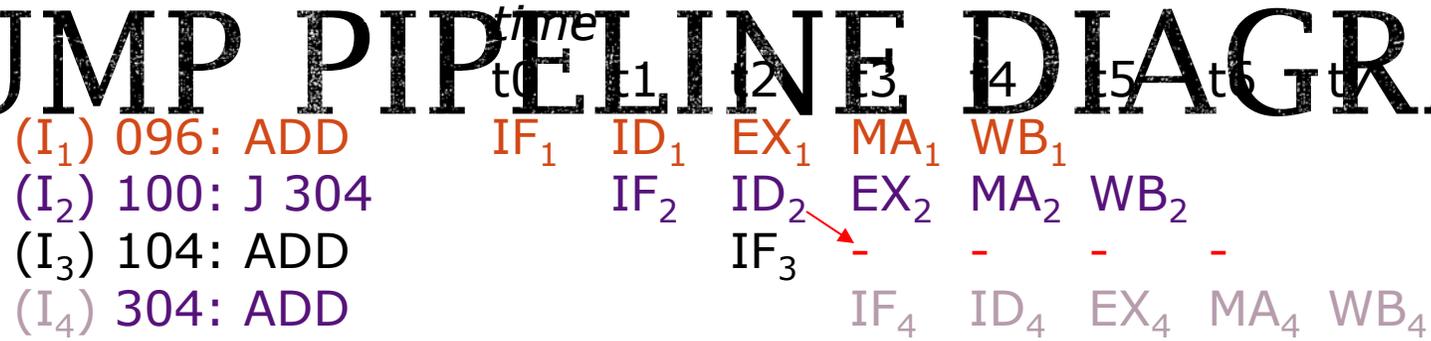


- I₁ 096ADD
- I₂ 100J 304
- I₃ 104ADD — *kill*
- I₄ 304ADD

A jump instruction kills (not stalls) the following instruction

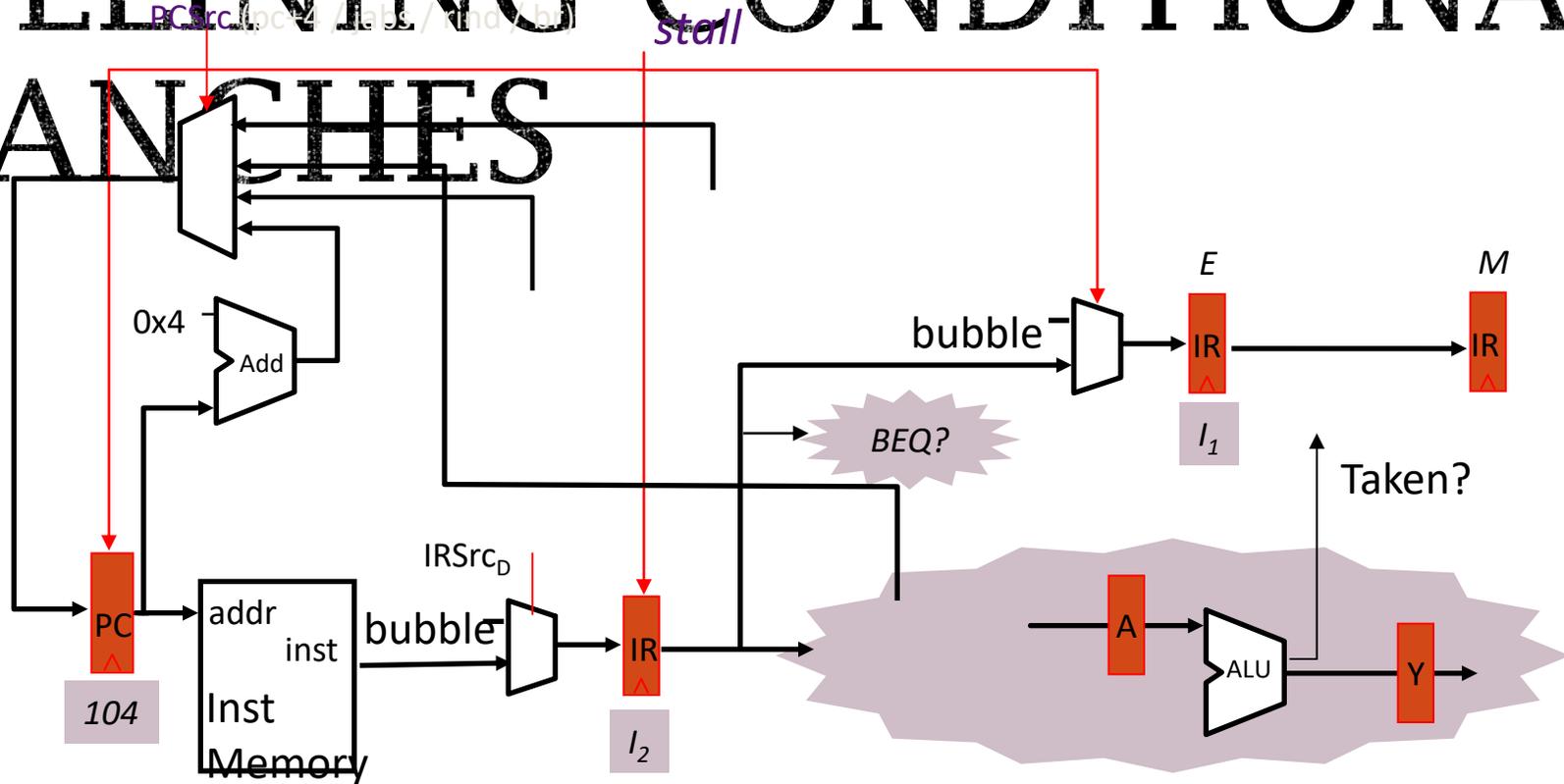
How?

JUMP PIPELINE DIAGRAMS



- ⇒ pipeline bubble

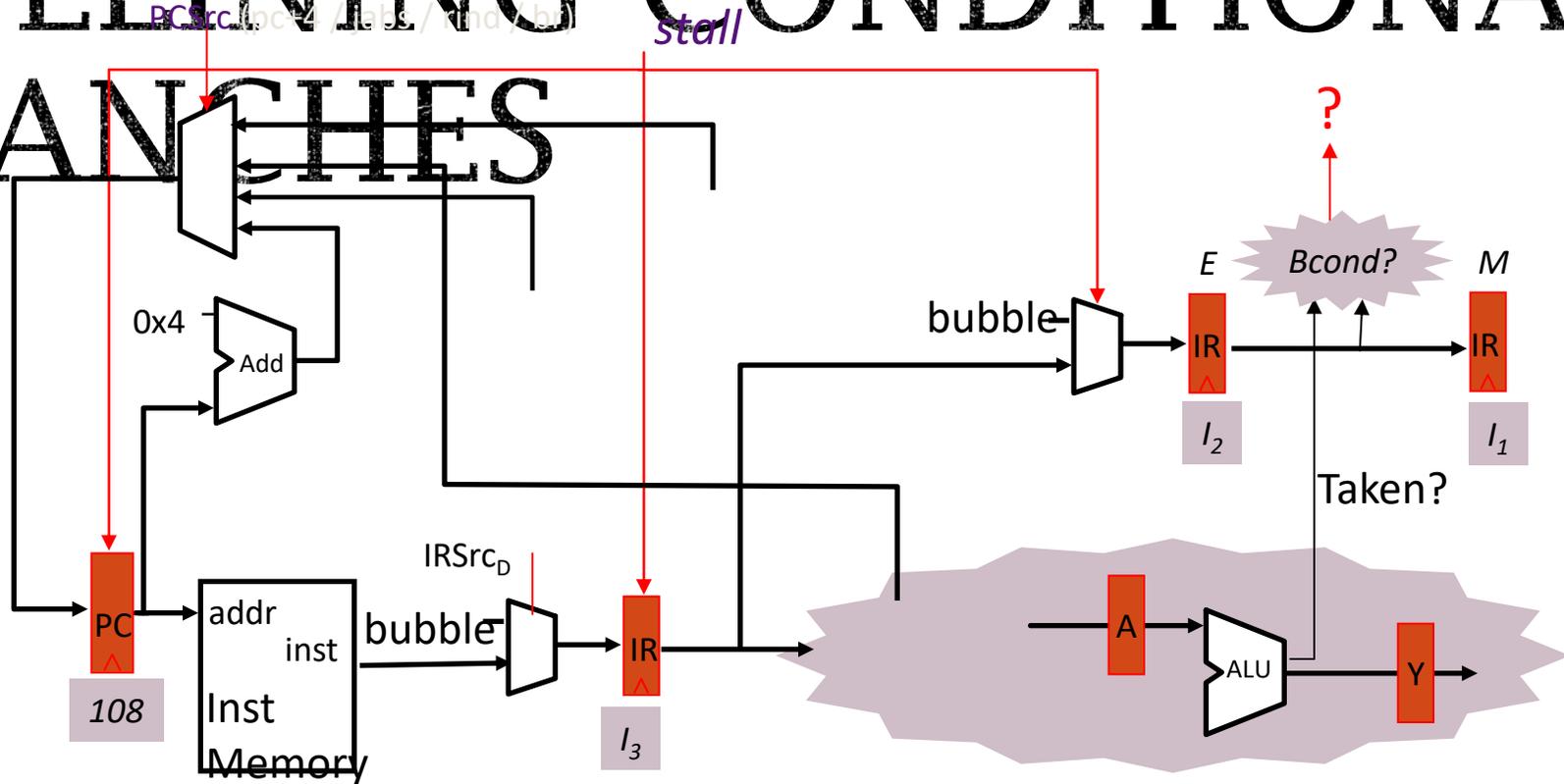
PIPELINING CONDITIONAL BRANCHES



- I₁ 096 ADD
- I₂ 100 BEQ x1,x2 +200
- I₃ 104 ADD
- I₄ 300 ADD

Branch condition is not known until the execute stage
what action should be taken in the decode stage ?

PIPELINING CONDITIONAL BRANCHES

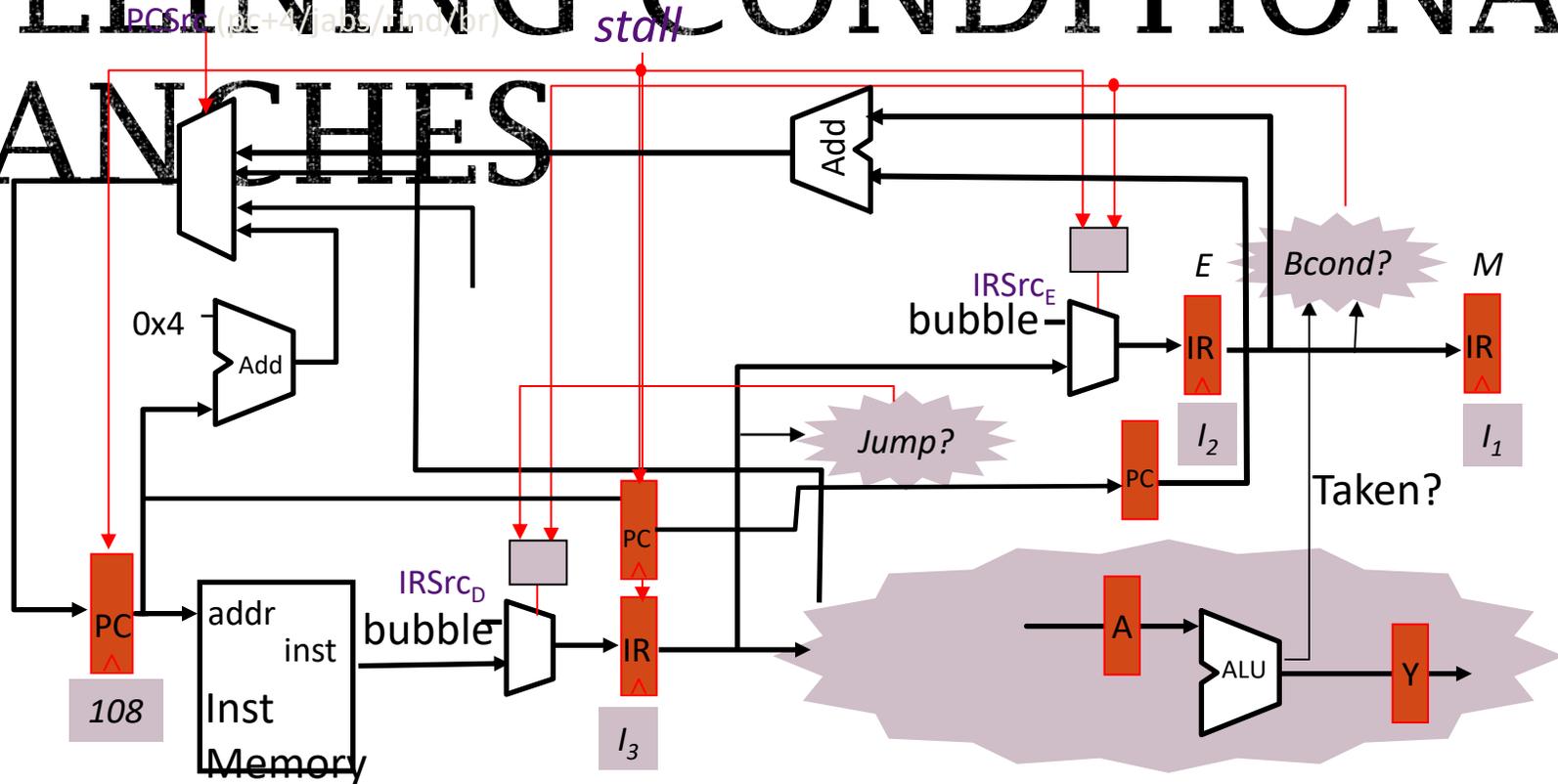


- I₁ 096 ADD
- I₂ 100 BEQ x1,x2 +200
- I₃ 104 ADD
- I₄ 300 ADD

If the branch is taken

- kill the two following instructions
- the instruction at the decode stage is not valid ⇒ *stall signal is not valid*

PIPELINING CONDITIONAL BRANCHES

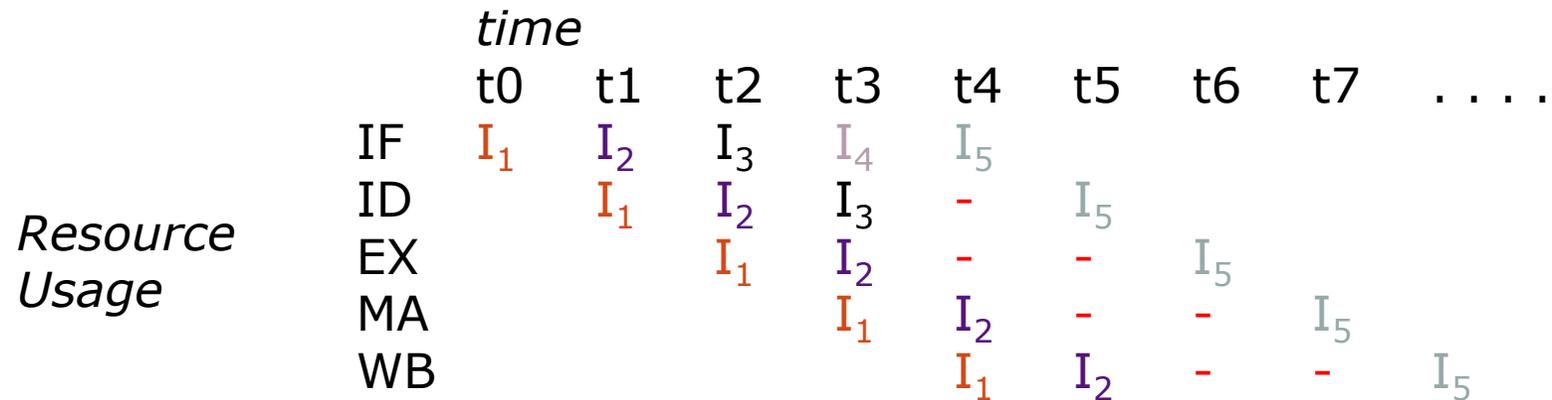
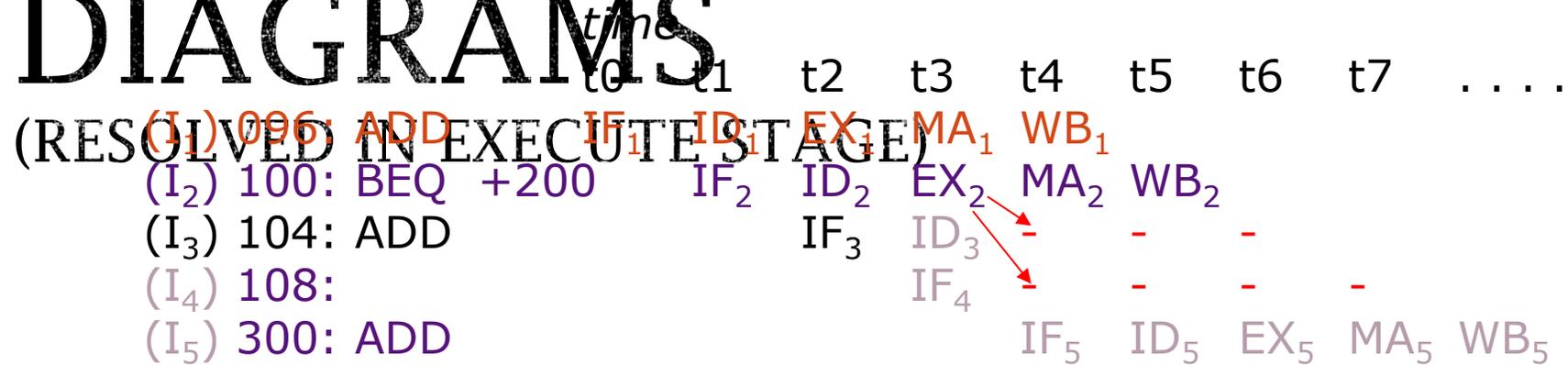


- I₁: 096 ADD
- I₂: 100 BEQ x1,x2 +200
- I₃: 104 ADD
- I₄: 300 ADD

If the branch is taken

- kill the two following instructions
- the instruction at the decode stage is not valid \Rightarrow *stall signal is not valid*

BRANCH PIPELINE DIAGRAMS



- ⇒ pipeline bubble

(EIGHT) OVER-COMPAK ONE REG AGAINST ZERO) WITH COMPARE IN DECODE

STAGE	<i>time</i>								
	t0	t1	t2	t3	t4	t5	t6	t7
(I ₁) 096: ADD	IF ₁	ID ₁	EX ₁	MA ₁	WB ₁				
(I ₂) 100: BEQZ +200		IF ₂	ID ₂	EX ₂	MA ₂	WB ₂			
(I ₃) 104: ADD			IF ₃	-	-	-	-		
(I ₄) 300: ADD				IF ₄	ID ₄	EX ₄	MA ₄	WB ₄	

<i>Resource Usage</i>	<i>time</i>								
	t0	t1	t2	t3	t4	t5	t6	t7
IF	I ₁	I ₂	I ₃	I ₄	I ₅				
ID		I ₁	I ₂	-	I ₄	I ₅			
EX			I ₁	I ₂	-	I ₄	I ₅		
MA				I ₁	I ₂	-	I ₄	I ₅	
WB					I ₁	I ₂	-	I ₄	I ₅

- ⇒ *pipeline bubble*

BRANCH DELAY SLOTS (EXPOSE CONTROL HAZARD TO SOFTWARE)

Change the ISA semantics so that the instruction that follows a jump or branch is always executed

- gives compiler the flexibility to put in a useful instruction where normally a pipeline bubble would have resulted.

I₁ 096 ADD

I₂ 100 BEQZ r1, +200

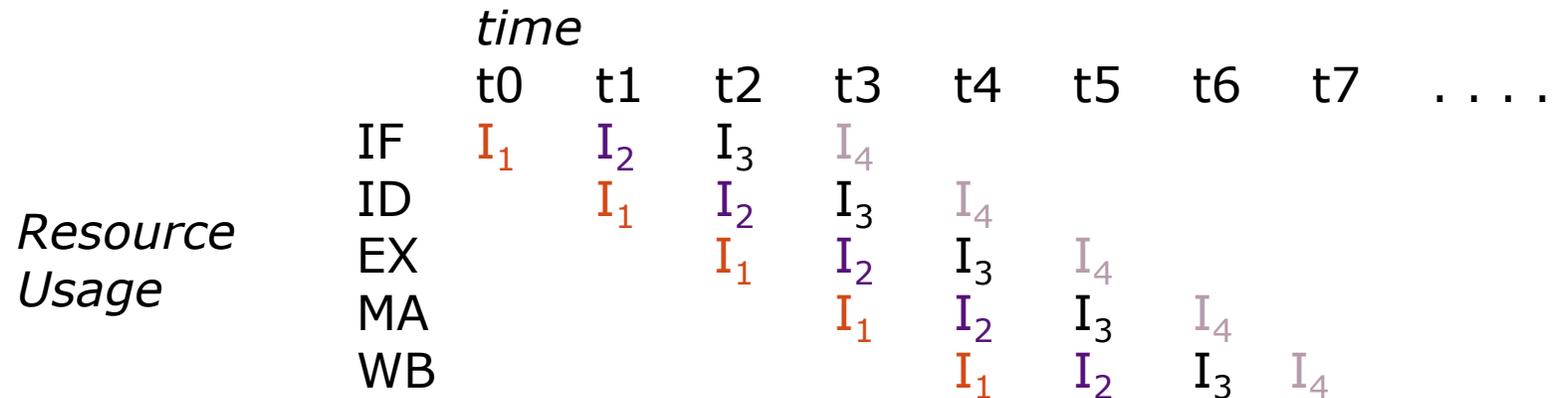
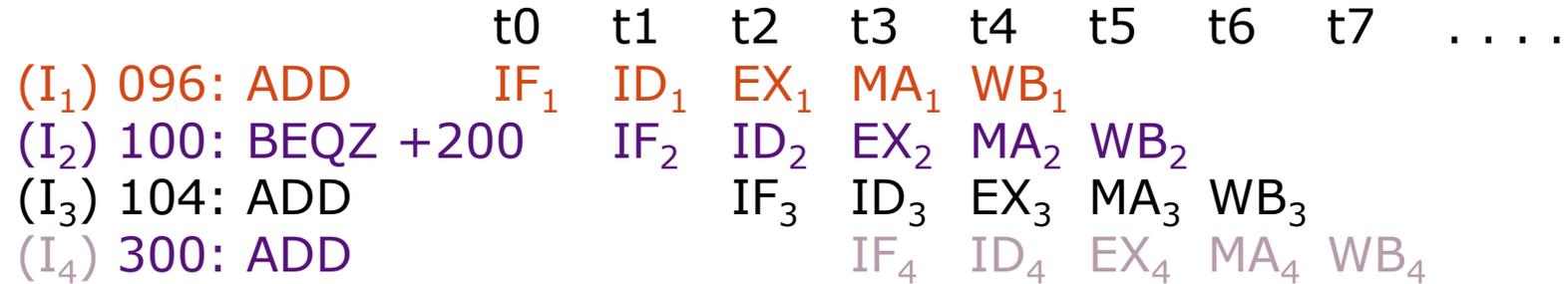
I₃ 104 ADD

I₄ 300 ADD

← *Delay slot instruction
executed regardless of
branch outcome*

BRANCH PIPELINE DIAGRAMS

(BRANCH DELAY SLOT)
time



POST-1990 RISC ISAS DON'T HAVE DELAY SLOTS

- Encodes microarchitectural detail into ISA
 - C.f. IBM 650 drum layout
- Performance issues
 - E.g., I-cache miss on delay slot causes machine to wait, even if delay slot is a NOP
- Complicates more advanced microarchitectures
 - 30-stage pipeline with four-instruction-per-cycle issue
- Better branch prediction reduced need

INSTRUCTION MAY NOT BE DISPATCHED

- Full bypassing may be too expensive to implement
 - typically all frequently used paths are provided (CPI > 1)
 - some infrequently used bypass paths may increase cycle time and counteract the benefit of reducing CPI
- Loads have two-cycle latency
 - Instruction after load cannot use load result
 - MIPS-I ISA defined *load delay slots*, a software-visible pipeline hazard (compiler schedules independent instruction or inserts NOP to avoid hazard). Removed in MIPS-II (pipeline interlocks added in hardware)
 - MIPS: “Microprocessor without Interlocked Pipeline Stages”
- Conditional branches may cause bubbles
 - kill following instruction(s) if no delay slots

Machines with software-visible delay slots may execute significant number of NOP instructions inserted by the compiler. NOPs increase instructions/program!

RISC-V BRANCHES AND JUMPS

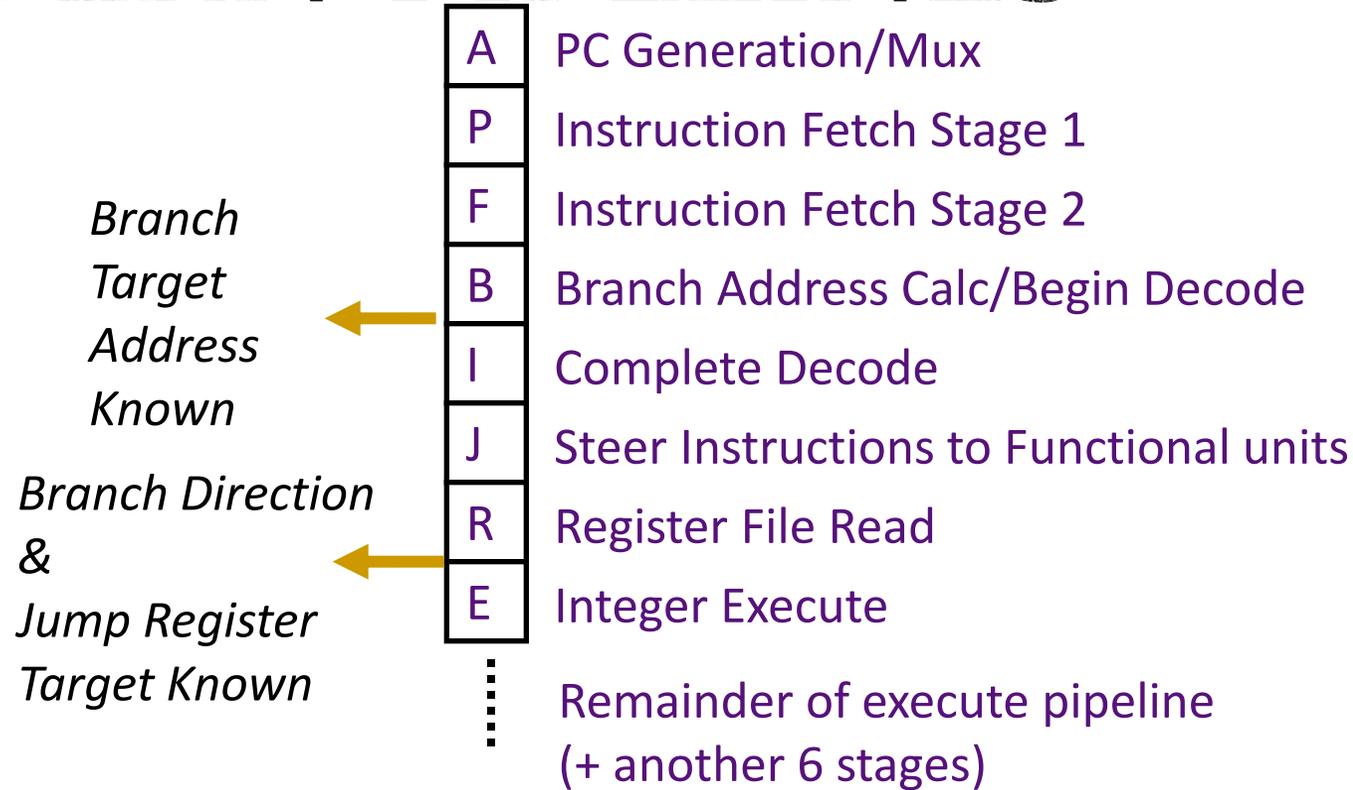
Each instruction fetch depends on one or two pieces of information from the preceding instruction:

- 1) Is the preceding instruction a taken branch?
- 2) If so, what is the target address?

<i>Instruction</i>	<i>Taken known?</i>	<i>Target known?</i>
J	After Inst. Decode	After Inst. Decode
JR	After Inst. Decode	After Reg. Fetch
B<cond.>	After Execute	After Inst. Decode

BRANCH PENALTIES IN MODERN PIPELINES

UltraSPARC-III instruction-fetch-pipeline stages
(in-order issue, 4-way superscalar, 750MHz, 2000)



REDUCING CONTROL FLOW PENALTY

- Software solutions
 - Eliminate branches - loop unrolling
 - Increases the run length
 - Reduce resolution time - instruction scheduling
 - Compute the branch condition as early as possible (of limited value because branches often in critical path through code)
- Hardware solutions
 - Find something else to do - delay slots
 - Replaces pipeline bubbles with useful work (requires software cooperation)
 - Speculate - branch prediction
 - Speculative execution of instructions beyond the branch

Motivation:

BRANCH PREDICTION

Branch penalties limit performance of deeply pipelined processors

Modern branch predictors have high accuracy (>95%) and can reduce branch penalties significantly

Required hardware support:

Prediction structures:

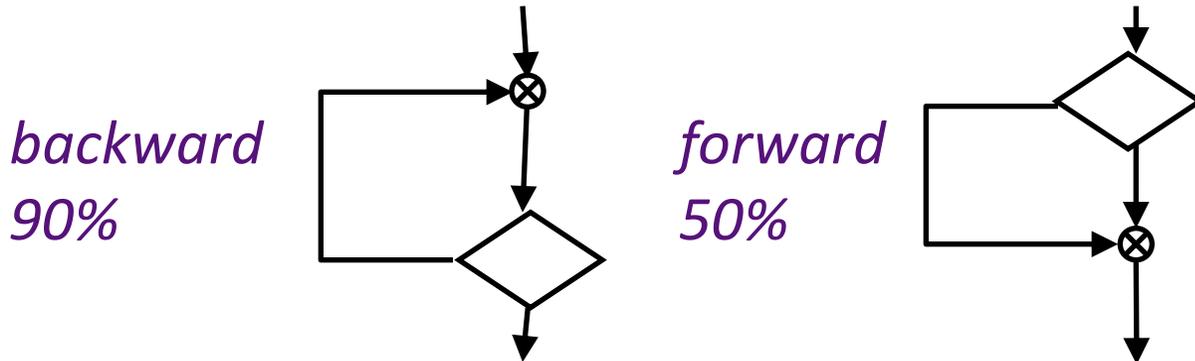
- Branch history tables, branch target buffers, etc.

Mispredict recovery mechanisms:

- *Keep result computation separate from commit*
- Kill instructions following branch in pipeline
- Restore state to that following branch

STATIC BRANCH PREDICTION

Overall probability a branch is taken is ~60-70% but:



ISA can attach preferred direction semantics to branches, e.g.,
Motorola MC88110

bne0 (preferred taken) *beq0 (not taken)*

ISA can allow arbitrary choice of statically predicted direction,
e.g., HP PA-RISC, Intel IA-64

typically reported as ~80% accurate

DYNAMIC BRANCH PREDICTION LEARNING BASED ON PAST BEHAVIOR

■ Temporal correlation

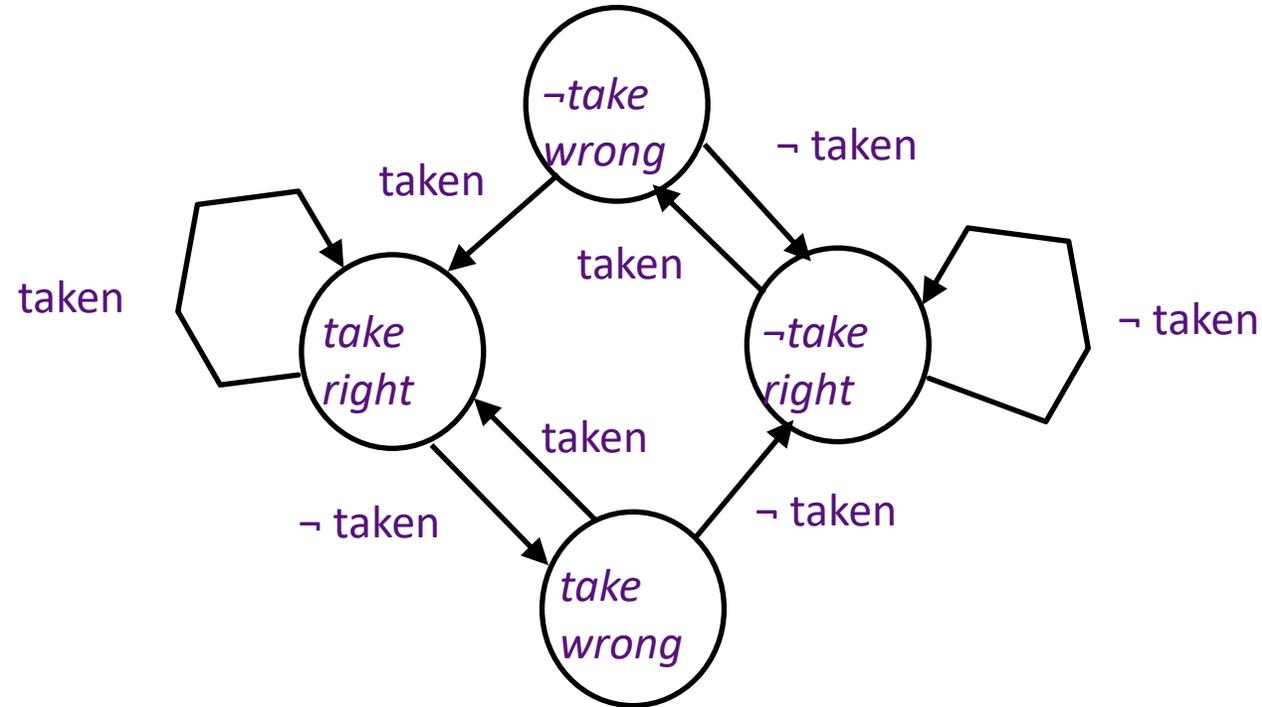
- The way a branch resolves may be a good predictor of the way it will resolve at the next execution

■ Spatial correlation

- Several branches may resolve in a highly correlated manner (a preferred path of execution)

BRANCH PREDICTION BITS

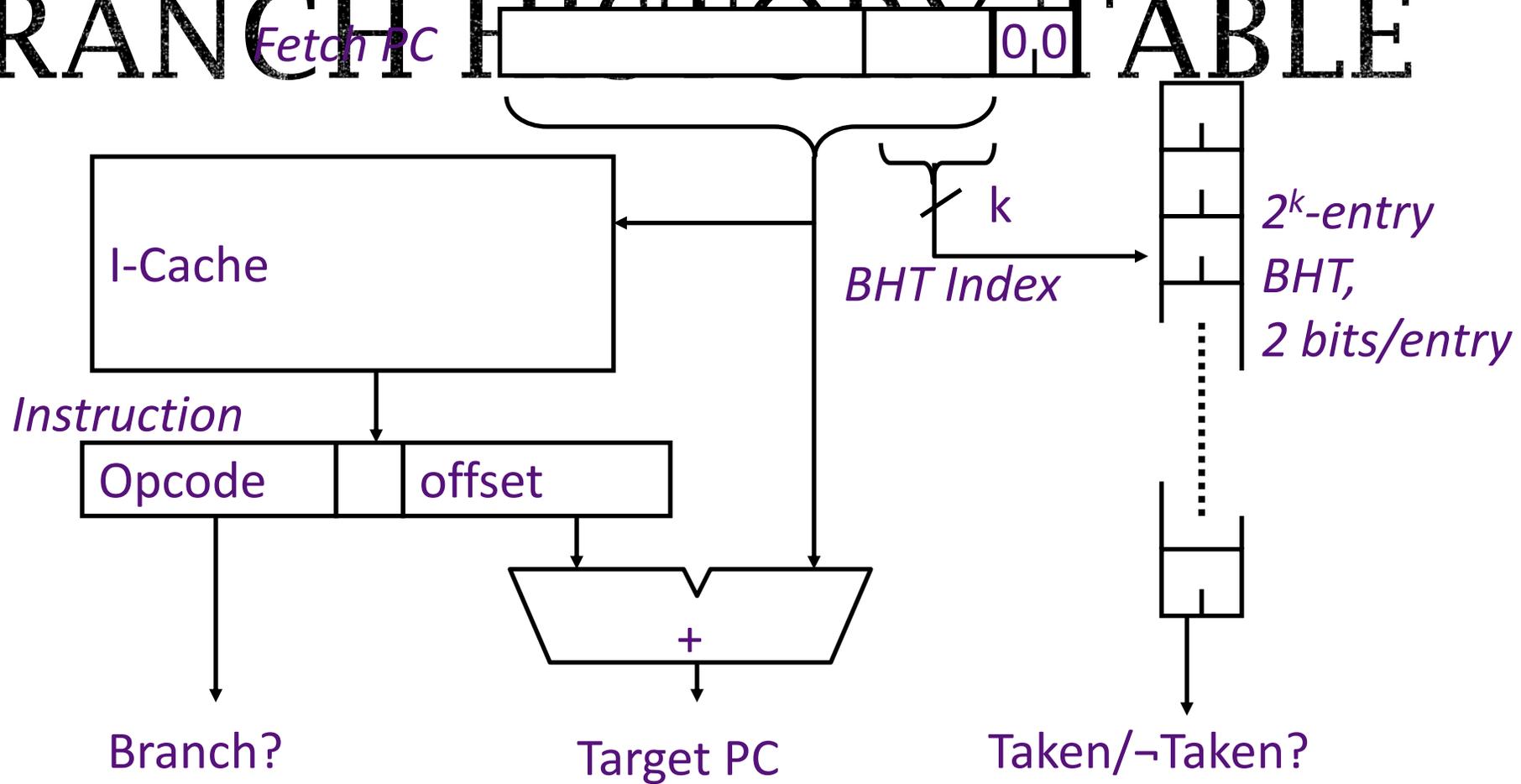
- Assume 2 BP bits per instruction
- Change the prediction after two consecutive mistakes!



BP state:

(predict take/-take) x (last prediction right/wrong)

BRANCH HISTORY TABLE



4K-entry BHT, 2 bits/entry, ~80-90% correct predictions

EXPLOITING SPATIAL CORRELATION

YEH AND PATT, 1992

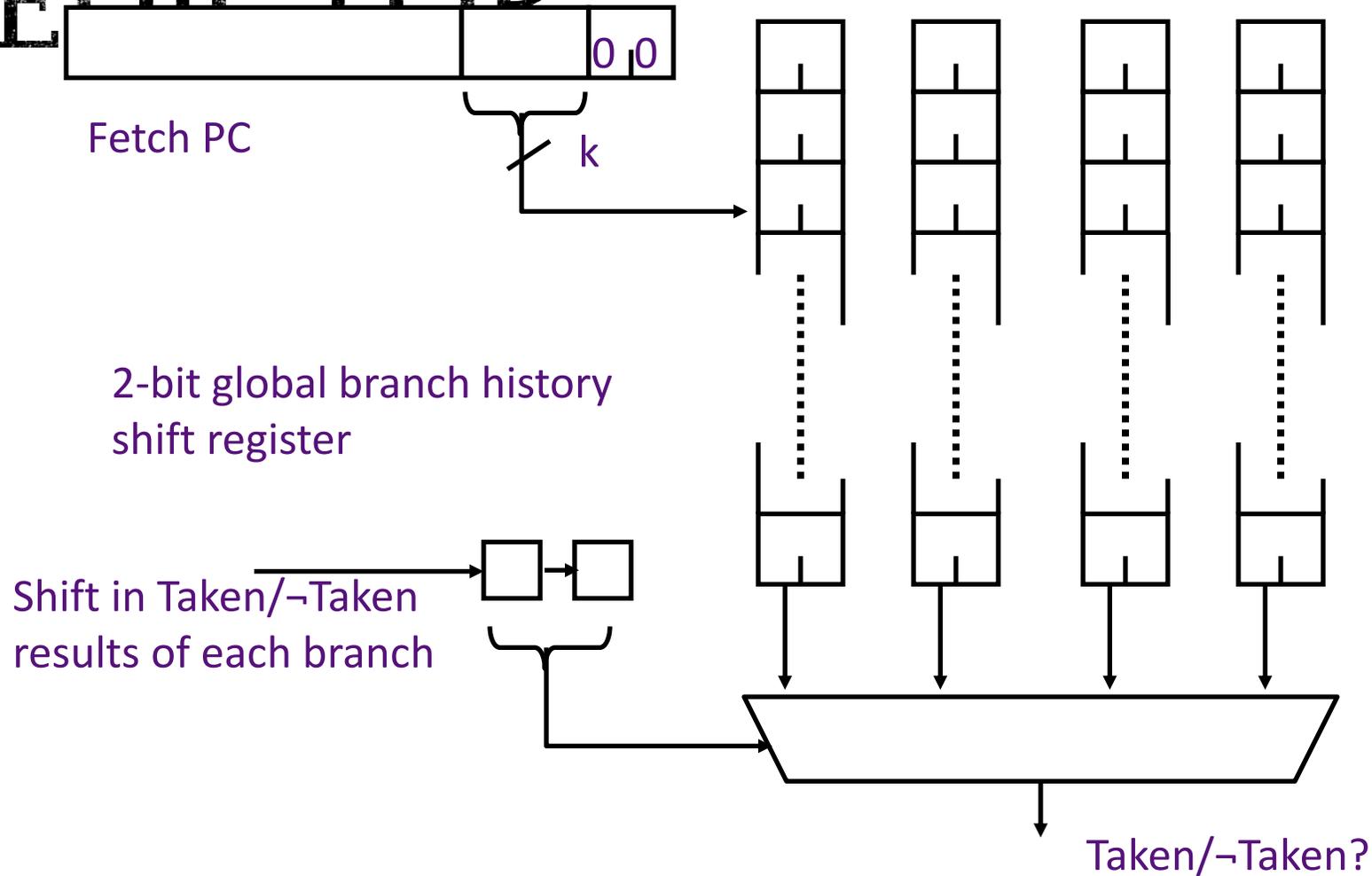
```
if (x[i] < 7) then  
  y += 1;  
if (x[i] < 5) then  
  c -= 4;
```

If first condition false, second condition also false

History register, H, records the direction of the last N branches executed by the processor

TWO-LEVEL BRANCH PREDICTOR

Pentium Pro uses the result from the last two branches to select one of the four sets of BHT bits (~95% correct)



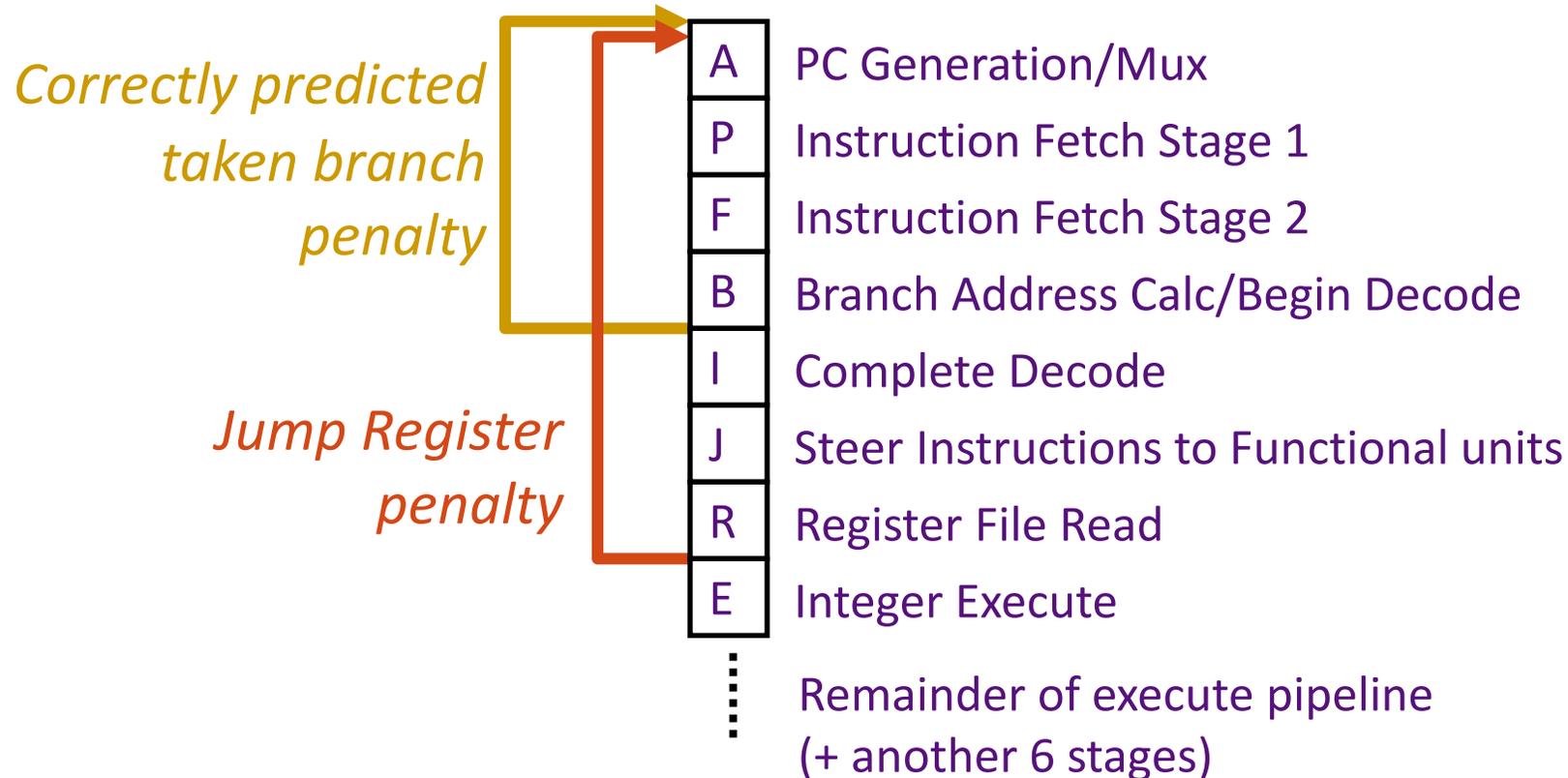
SPECULATING BOTH

DIRECTIONS

- An alternative to branch prediction is to execute both directions of a branch speculatively
 - resource requirement is proportional to the number of concurrent speculative executions
 - only half the resources engage in useful work when both directions of a branch are executed speculatively
 - branch prediction takes less resources than speculative execution of both paths
- With accurate branch prediction, it is more cost effective to dedicate all resources to the predicted direction!

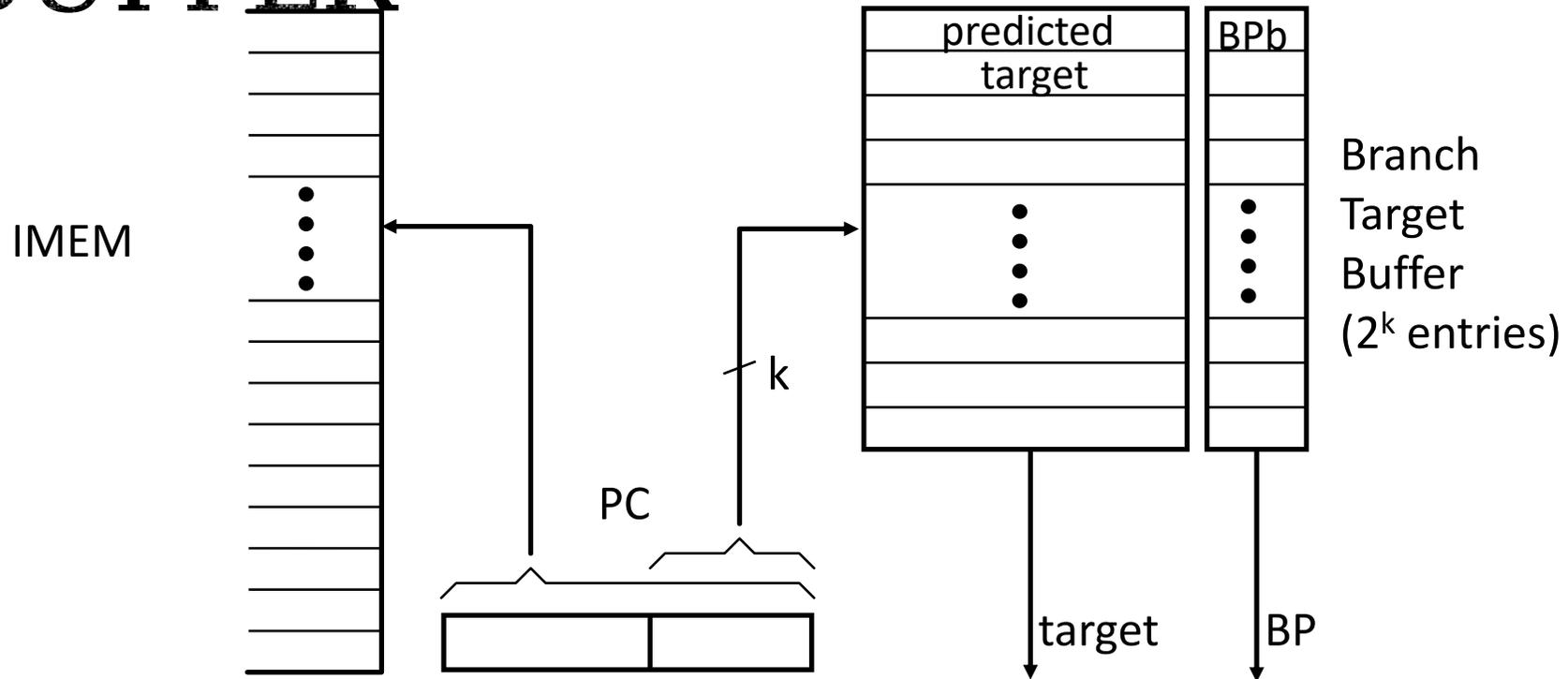
LIMITATIONS OF BRITS

Only predicts branch direction. Therefore, cannot redirect fetch stream until after branch target is determined.



UltraSPARC-III fetch pipeline

BRANCH TARGET BUFFER



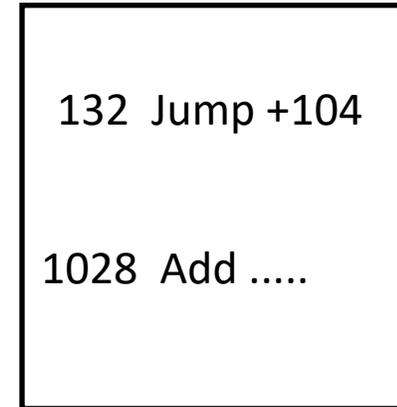
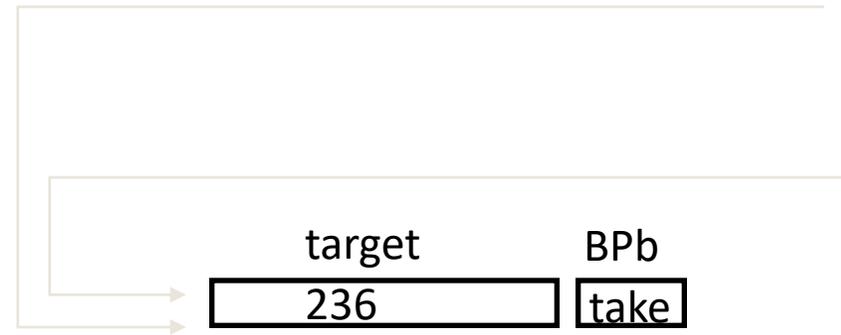
BP bits are stored with the predicted target address.

IF stage: *If (BP=taken) then nPC=target else nPC=PC+4*

Later: *check prediction, if wrong then kill the instruction and update BTB & BPb else update BPb*

ADDRESS COLLISIONS

Assume a
128-entry
BTB



Instruction
Memory

What will be fetched after the instruction at 1028?

BTB prediction = 236
Correct target = 1032

?? *kill* PC=236 and *fetch* PC=1032

*Is this a common occurrence?
Can we avoid these bubbles?*

BTB IS ~~ONLY FOR CONTROL INSTRUCTIONS~~

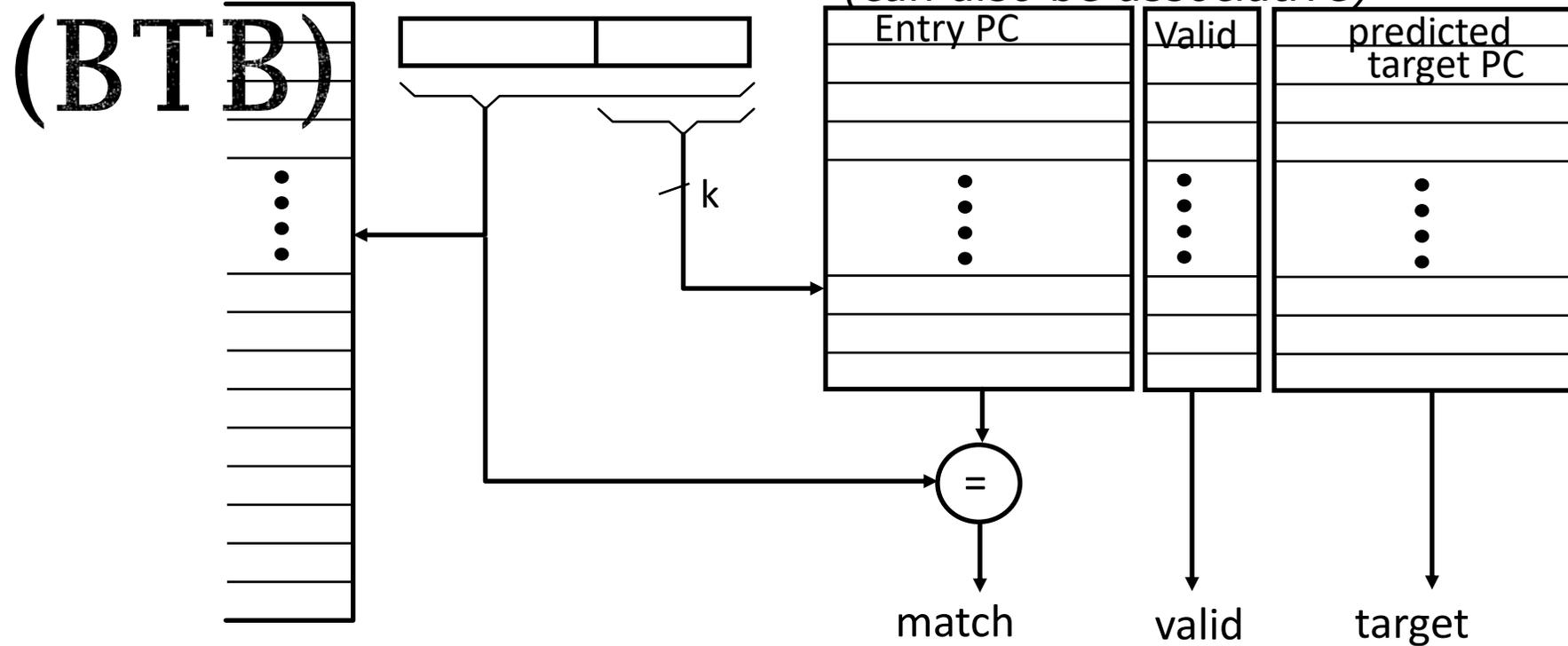
BTB contains useful information for branch and jump instructions only

❓ Do not update it for other instructions

- For all other instructions the next PC is PC+4 !
- *How to achieve this effect without decoding the instruction?*

BRANCH TARGET BUFFER (BTB)

2^k-entry direct-mapped BTB
(can also be associative)



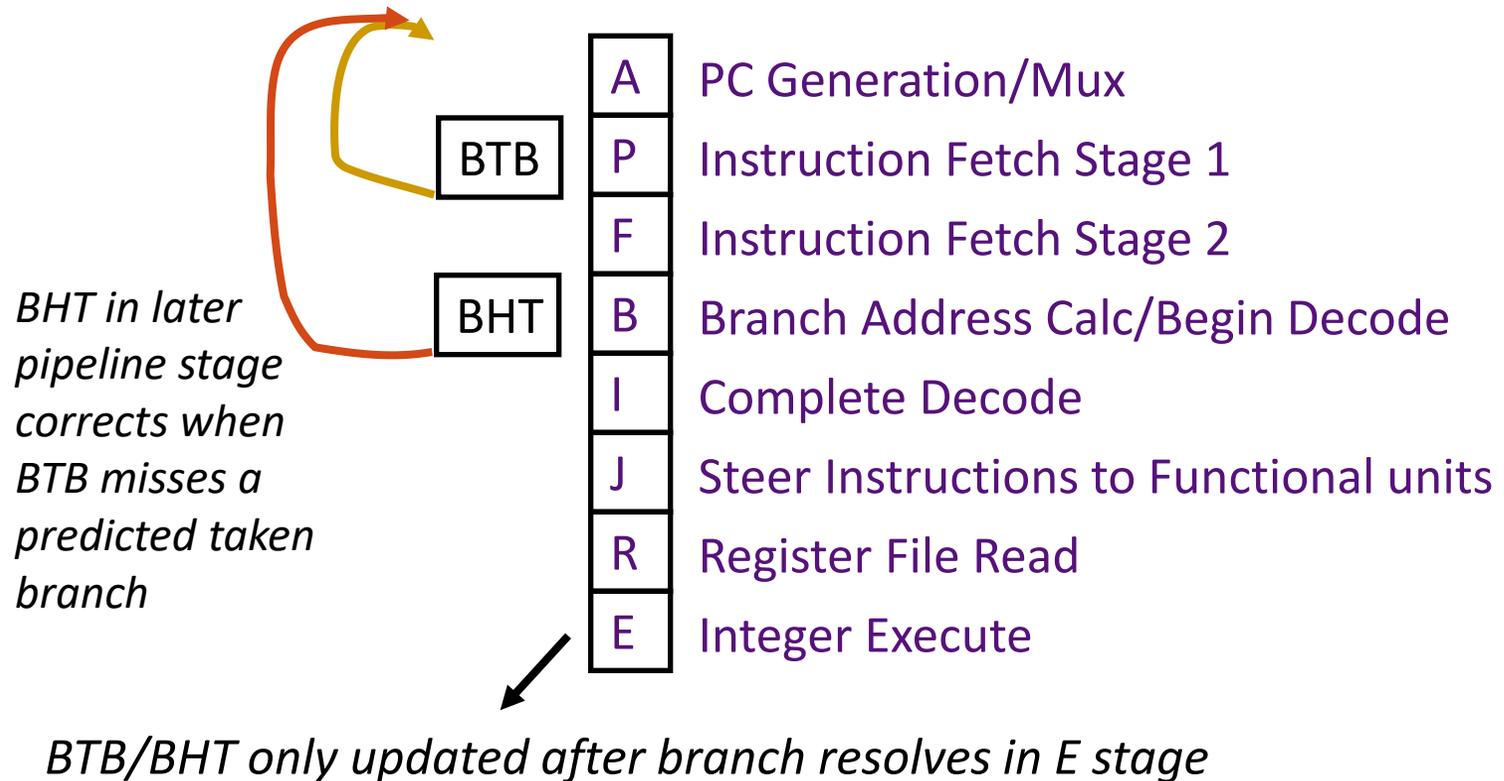
- Keep both the branch PC and target PC in the BTB
- PC+4 is fetched if match fails
- Only *taken* branches and jumps held in BTB
- Next PC determined *before* branch fetched and decoded

COMBINING BTB AND

BHT

BTB entries are considerably more expensive than BHT, but can redirect fetches at earlier stage in pipeline and can accelerate indirect branches (JR)

- BHT can hold many more entries and is more accurate



USES OF JUMP REGISTER (JR)

BTB works well if same case used repeatedly

- Switch statements (jump to address of matching case)

BTB works well if same function usually called, (e.g., in

- Dynamic function call (jump to run-time function address)
C++ programming, when objects have same type in virtual function call)

BTB works well if usually return to the same place

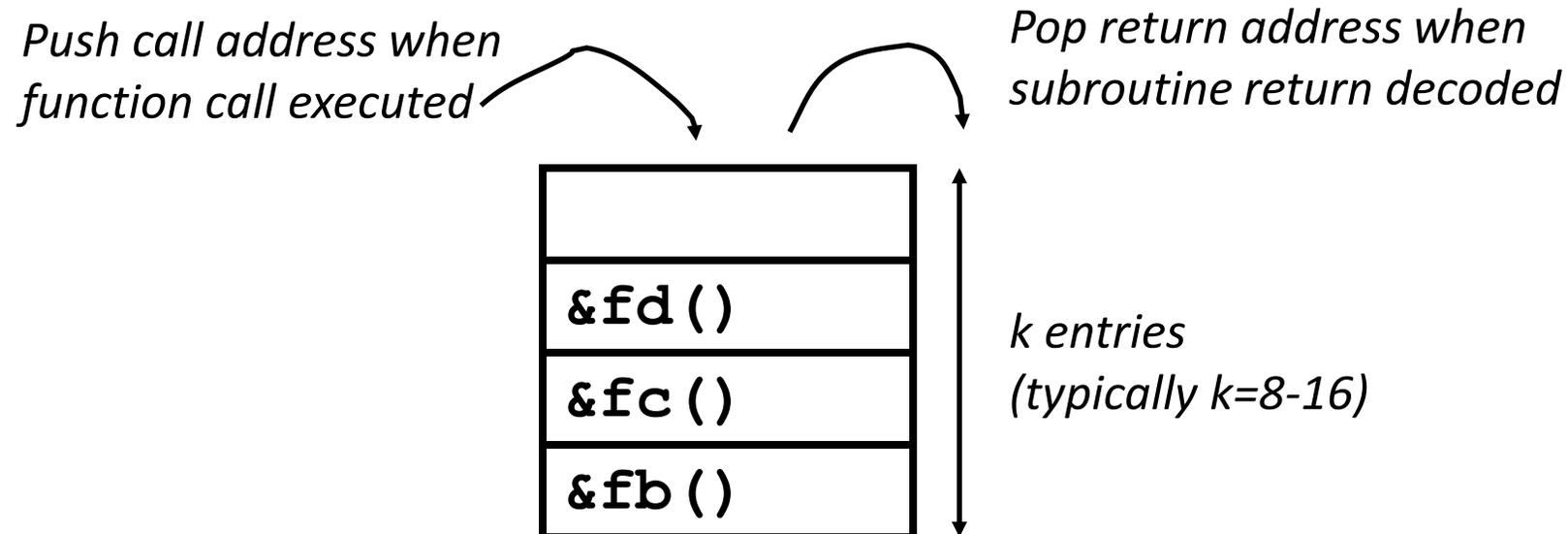
- Subroutine returns (jump to return address)
⇒ Often (one function called from) many distinct call sites!

How well does BTB work for each of these cases?

SUBROUTINE RETURN STACK

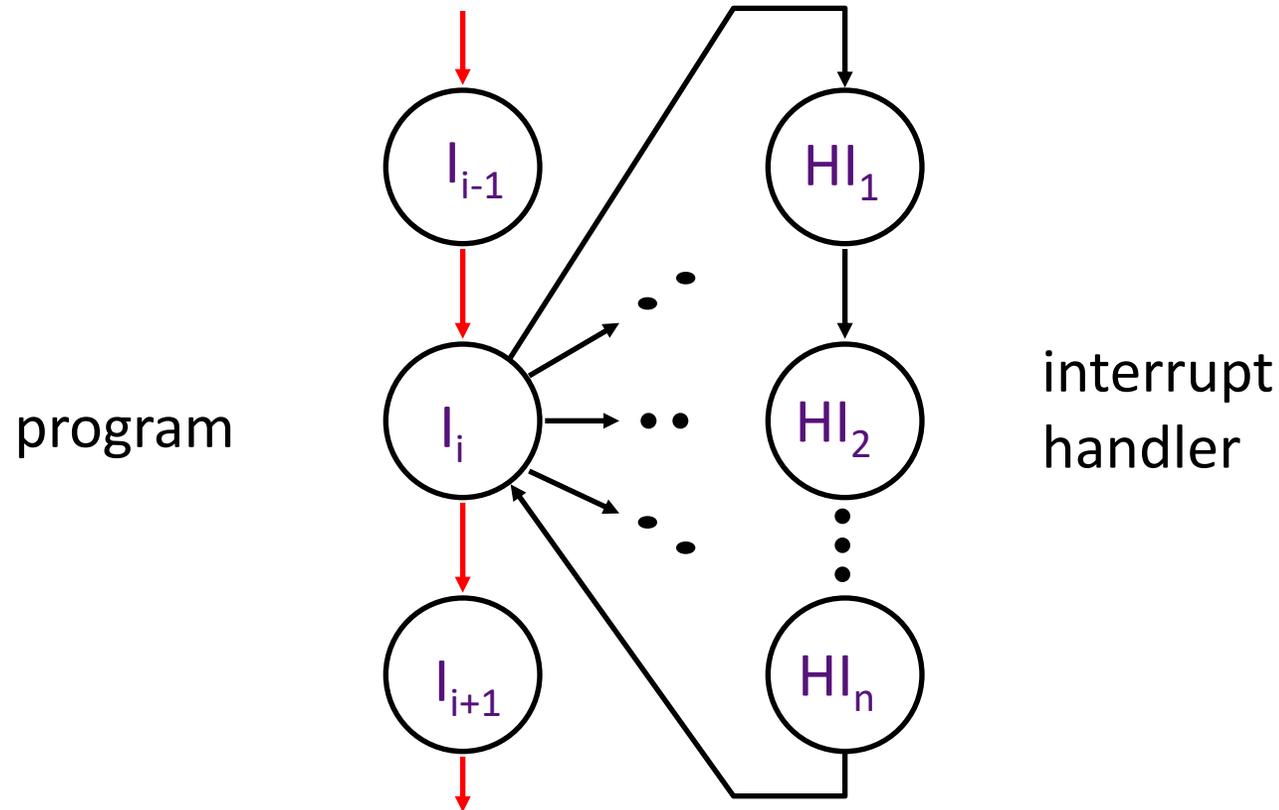
Small structure to accelerate JR for subroutine returns,
typically much more accurate than BTBs.

```
fa () { fb (); }  
fb () { fc (); }  
fc () { fd (); }
```



INTERRUPTS:

ALTERING THE NORMAL FLOW OF CONTROL



An *external or internal event* that needs to be processed by another (system) program. The event is usually unexpected or rare from program's point of view.

CAUSES OF INTERRUPTS

Interrupt: an *event* that requests the attention of the processor

- Asynchronous: an *external event*
 - input/output device service-request
 - timer expiration
 - power disruptions, hardware failure
- Synchronous: an *internal event (a.k.a. traps or exceptions)*
 - undefined opcode, privileged instruction
 - arithmetic overflow, FPU exception
 - misaligned memory access
 - *virtual memory exceptions*: page faults, TLB misses, protection violations
 - system calls, e.g., jumps into kernel

HISTORY OF EXCEPTION HANDLING

▪ First system with exceptions was Univac-I, 1951

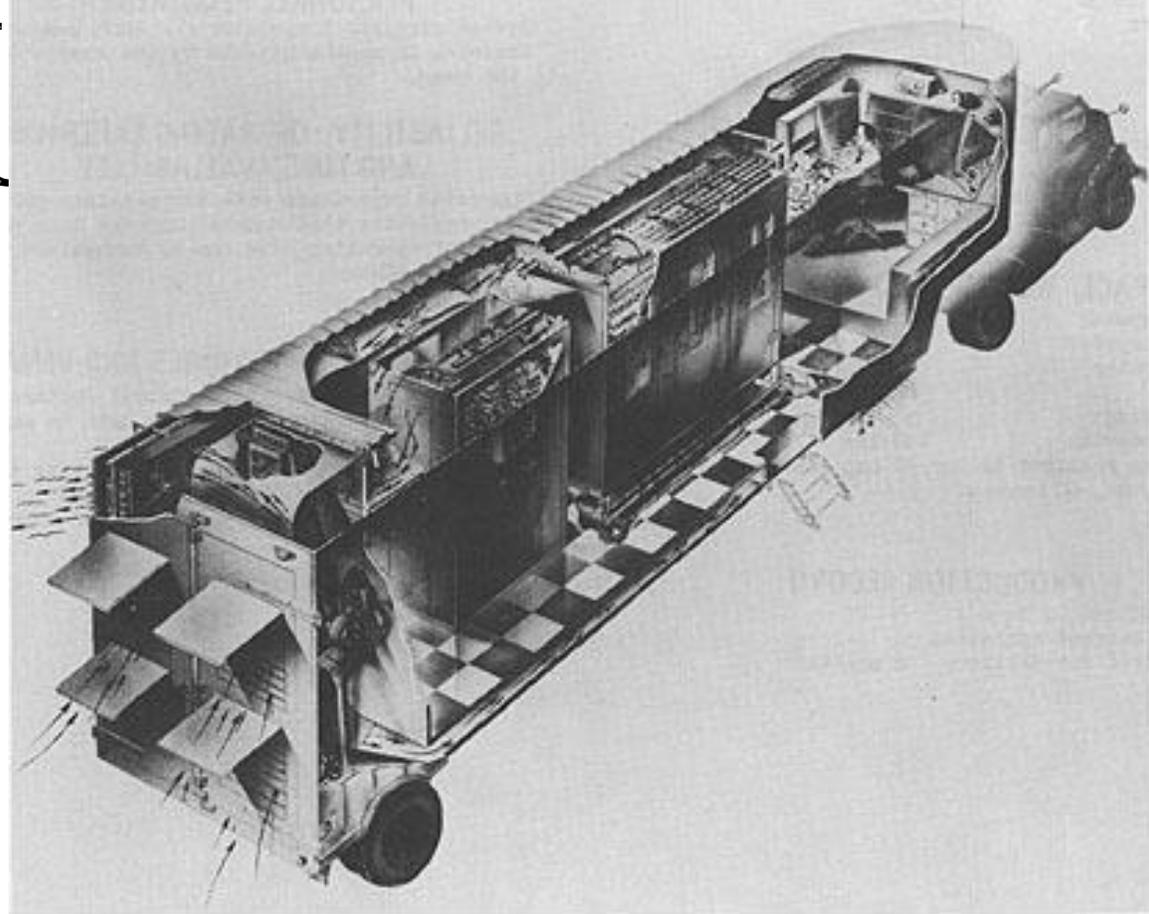
- Arithmetic overflow would either
 - 1. trigger the execution a two-instruction fix-up routine at address 0, or
 - 2. at the programmer's option, cause the computer to stop
- Later Univac 1103, 1955, modified to add external interrupts
 - Used to gather real-time wind tunnel data

▪ First system with I/O interrupts was DYSEAC, 1954

- Had two program counters, and I/O signal caused switch between two PCs
- Also, first system with DMA (direct memory access by I/O device)

[Courtesy Mark Smotherman]

DYSEAC FIRST MOBILE COMPU



- Carried in two tractor trailers, 12 tons + 8 tons
- Built for US Army Signal Corps

[Courtesy Mark Smotherman]

ASYNCHRONOUS

INTERRUPTS:

- An I/O device requests attention by asserting one of the prioritized interrupt request lines

INVOKING THE INTERRUPT HANDLER

- When the processor decides to process the interrupt
 - It stops the current program at instruction I_j , completing all the instructions up to I_{j-1} (*precise interrupt*)
 - It saves the PC of instruction I_j in a special register (EPC)
 - It disables interrupts and transfers control to a designated interrupt handler running in the kernel mode

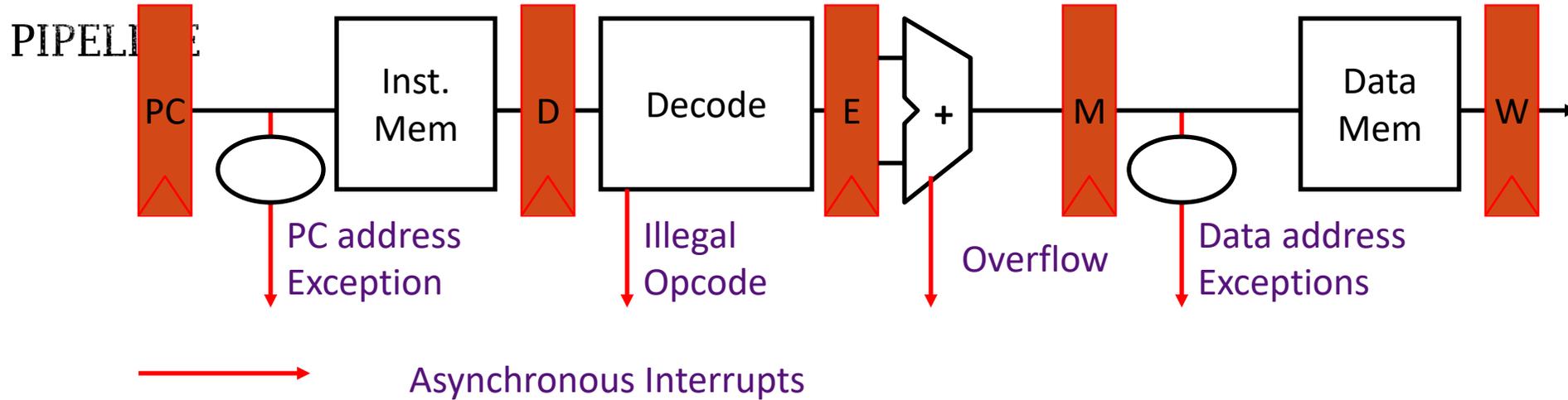
INTERRUPT HANDLER

- Saves EPC before enabling interrupts to allow nested interrupts ⇒
 - need an instruction to move EPC into GPRs
 - need a way to mask further interrupts at least until EPC can be saved
- Needs to read a *status register* that indicates the cause of the interrupt
- Uses a special indirect jump instruction RFE (*return-from-exception*) which
 - enables interrupts
 - restores the processor to the user mode
 - restores hardware status and control state

SYNCHRONOUS INTERRUPTS

- A synchronous interrupt (exception) is caused by a *particular instruction*
- In general, the instruction cannot be completed and needs to be *restarted* after the exception has been handled
 - requires undoing the effect of one or more partially executed instructions
- In the case of a system call trap, the instruction is considered to have been completed
 - a special jump instruction involving a change to privileged kernel mode

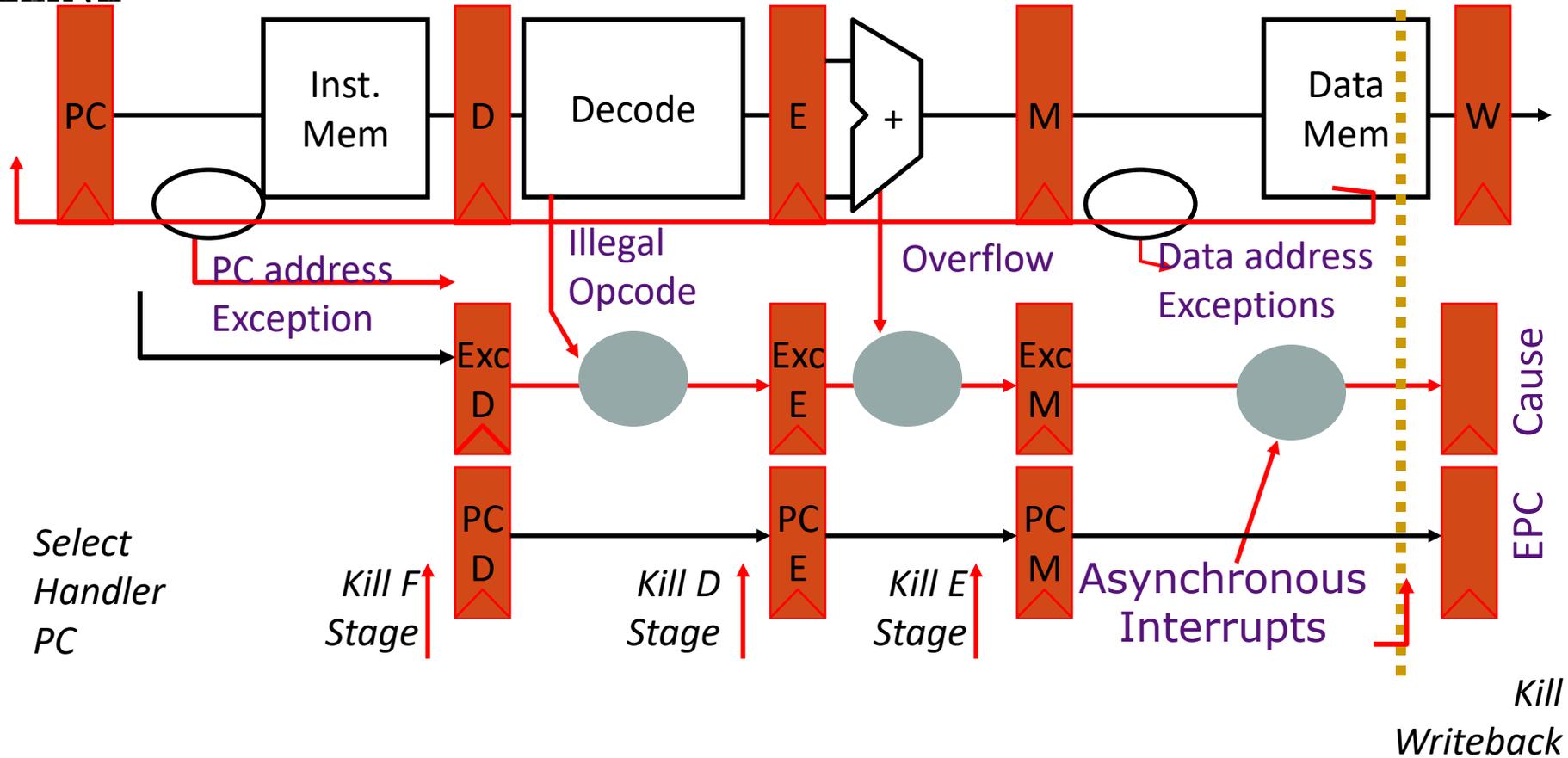
EXCEPTION HANDLING 5-STAGE



- How to handle multiple simultaneous exceptions in different pipeline stages?
- How and where to handle external asynchronous interrupts?

EXCEPTION HANDLING

PIPELINE



EXCEPTION HANDLING

5-STAGE PIPELINE

- Hold exception flags in pipeline until commit point (M stage)
- Exceptions in earlier pipe stages override later exceptions *for a given instruction*
- Inject external interrupts at commit point (override others)
- If exception at commit: update Cause and EPC registers, kill all stages, inject handler PC into fetch stage

SPECULATING ON

EXCEPTIONS

- Prediction mechanism
- Exceptions are rare, so simply predicting no exceptions is very accurate!

- Check prediction mechanism

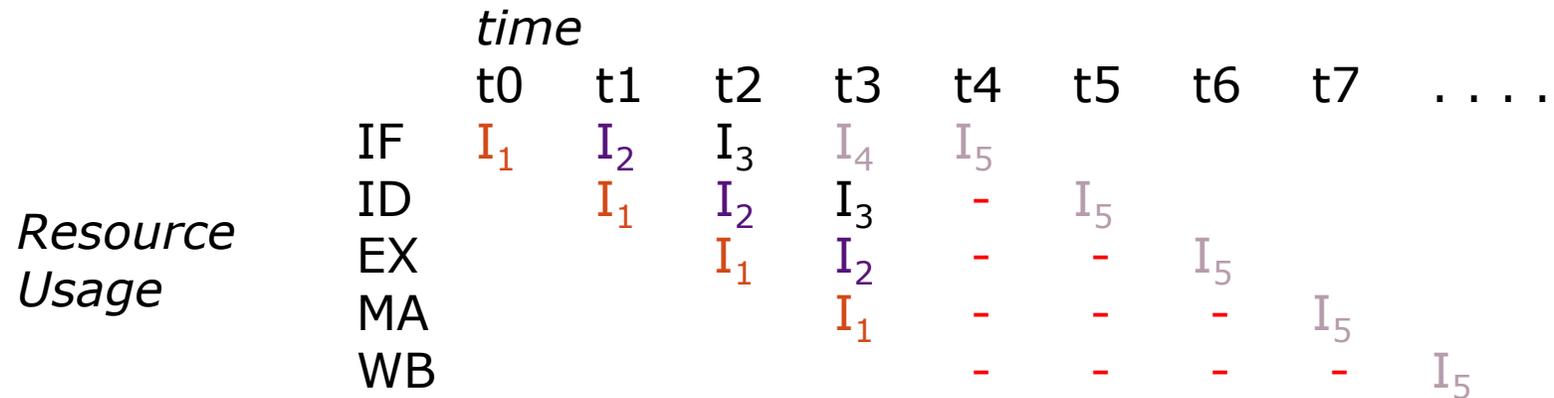
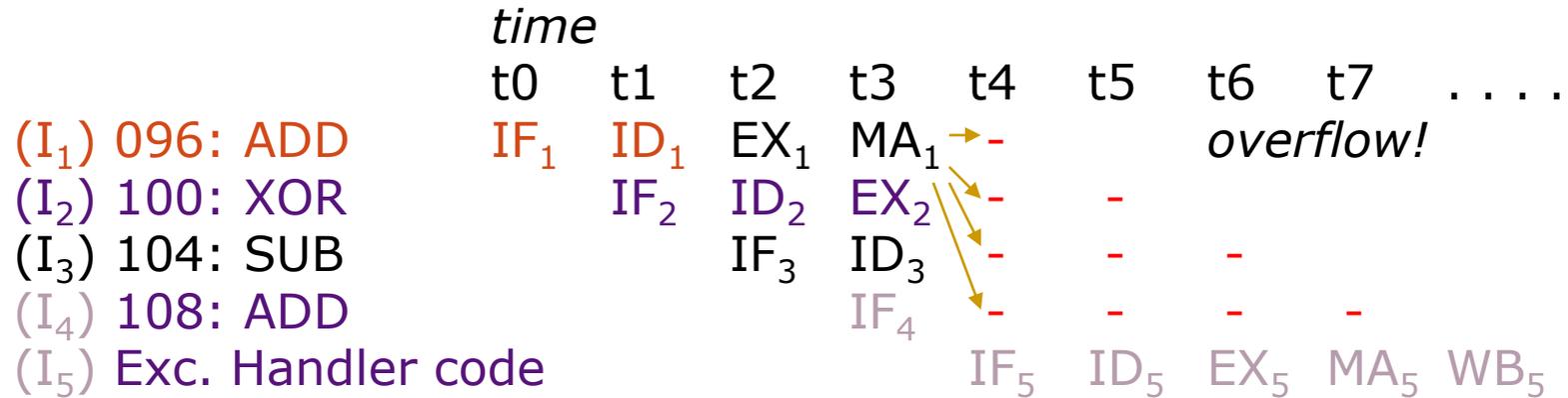
- Exceptions detected at end of instruction execution pipeline, special hardware for various exception types

- Recovery mechanism

- Only write architectural state at commit point, so can throw away partially executed instructions after exception
- Launch exception handler after flushing pipeline

- Bypassing allows use of uncommitted instruction results by following instructions

EXCEPTION PIPELINE DIAGRAM



ACKNOWLEDGEMENTS

- These slides contain material developed and copyright by:
 - Arvind (MIT)
 - Krste Asanovic (MIT/UCB)
 - Joel Emer (Intel/MIT)
 - James Hoe (CMU)
 - John Kubiatowicz (UCB)
 - David Patterson (UCB)
- MIT material derived from course 6.823
- UCB material derived from course CS252