# **Dynamic Branch Prediction**

STUDENTS-HUB.com

# Control Dependences and Branch Prediction

- Modern processors have deep pipelines
  - ♦ Pipeline depth for i7 is 14 stages
  - Because of the combination of deep pipelining and multiple issues per clock, the i7 has many instructions in-flight at once (up to 256, and typically at least 30).
  - ♦ Branch penalty limits performance of deep pipelines.
- ♦ Want to execute instructions beyond a branch even before that branch is resolved → use speculative execution
- What to predict?
  - ♦ Answer: The address of the next instruction
  - $\diamond$  If the fetched instruction is a non-control-flow instruction:
    - Next Fetch PC is the address of the next-sequential instruction
    - Easy to determine if we know the size of the fetched instruction
  - $\diamond$  If the instruction that is fetched is a control-flow instruction:
    - How do we determine the next Fetch PC?
    - In fact, how do we even know whether or not the fetched instruction is a control-flow instruction?

STUDENTS-HUB.com

# Branch Types

Туре	Direction at fetch time	Number of possible next fetch addresses?	When is next fetch address resolved?
Conditional	Unknown	2	Execution (register dependent)
Unconditional	Always taken	1	Decode (PC + offset)
Call	Always taken	1	Decode (PC + offset)
Return	Always taken	Many	Execution (register dependent)
Indirect	Always taken	Many	Execution (register dependent)

Different branch types can be handled differently

STUDENTS-HUB.com

## How to Handle Control Dependences

- Critical to keep the pipeline full with correct sequence of dynamic instructions.
- Potential solutions if the instruction is a control-flow instruction:
  - ♦ Stall the pipeline until we know the next fetch address
  - $\diamond$  Always Guess Next PC = PC + 4
  - ♦ Reducing the Branch Misprediction Penalty
  - ♦ Employ delayed branching (branch delay slot)
  - ♦ Eliminate control-flow instructions (predicated execution)
  - ♦ Guess the next fetch address (branch prediction)
  - ♦ Do something else (fine-grained multithreading)
  - Fetch from both possible paths (if you know the addresses of both possible paths) (multipath execution)

### Stall Fetch Until Next PC is Known: Good Idea?



STUDENTS-HUB.com

### Always Guess NextPC = PC + 4

- Always predict the next sequential instruction is the next instruction to be executed
- This is a form of next fetch address prediction (and branch prediction)

### How can you make this more effective?

- Idea: Maximize the chances that the next sequential instruction is the next instruction to be executed
  - Software: Lay out the control flow graph such that the "likely next instruction" is on the not-taken path of a branch
    - Profile guided code positioning
  - ♦ Hardware: ??? (how can you do this in hardware...)
    - Cache traces of executed instructions  $\rightarrow$  Trace cache

STUDENTS-HUB.com

## Reducing the Branch Misprediction Penalty

### Resolve branch condition and target address early



STUDENTS-HUB.com

# Delayed Branch

- Define branch to take place after the next instruction
- MIPS defines one delay slot
  - ♦ Reduces branch penalty
- Compiler fills the branch delay slot
  - ♦ By selecting an independent instruction

from before the branch

 $\diamond$  Must be okay to execute instruction in the

delay slot whether branch is taken or not

- ✤ If no instruction is found
- Compiler fills delay slot with a NO-OP STUDENTS-HUB.com



### The Branch Problem

- Control flow instructions (branches) are frequent
  - $\diamond$  15-25% of all instructions
- Problem: Next fetch address after a control-flow instruction is not determined after N cycles in a pipelined processor
  - ♦ N cycles: (minimum) branch resolution latency
- If we are fetching W instructions per cycle (i.e., if the pipeline is W wide)
  - ♦ A branch misprediction leads to N x W wasted instruction slots

# Branch Prediction (A Bit More Enhanced)

- Idea: Predict the next fetch address (to be used in the next cycle) → no wasted cycle(s) on correct prediction
- Requires three things to be predicted at fetch stage:
  - ♦ Whether the fetched instruction is a branch
  - ♦ Direction (1-bit)
    - Single direction for unconditional jumps and calls/returns
    - Binary for conditional branches
  - ♦ Target (32-bit or 64-bit addresses)
    - Some are easy
      - One address: uni-directional jumps
      - Two: addresses: fall through (not taken) vs. taken
    - Many: function pointer or indirect jump (e.g. jr r31)
- Observation: Target address remains the same for a conditional direct branch across dynamic instances
  - ♦ Idea: Store the target address from previous instance and access it with the PC

STUDEN Called Branch Target Buffer (BTB) or Branch Target Address Cache By: Jibreel Bornat

### **Branch Prediction Techniques**

### Compile time (static)

- ♦ Always not taken
- ♦ Always taken
- ♦ BTFN (Backward taken, forward not taken)
- ♦ Profile based (likely direction)
- ♦ Program analysis based (likely direction)

### Run time (Dynamic Branch Prediction)

- ♦ Last time prediction (single-bit)
- ♦ Two-bit counter based prediction
- ♦ Two-level prediction (global vs. local)
- ♦ Hybrid

Common disadvantage of compile time methods?

Cannot adapt to dynamic changes in branch behavior

## **Dynamic Branch Prediction**

- Idea: Predict branches based on dynamic information (collected at run-time)
- Advantages
  - + Prediction based on history of the execution of branches. It can adapt to dynamic changes in branch behavior
  - + No need for static profiling: input set representativeness problem goes away

### Disadvantages

-- More complex (requires additional hardware)

### 1-bit Dynamic Branch Prediction Scheme

- Prediction is just a hint that is assumed to be correct
- If incorrect then fetched instructions are killed
- 1-bit prediction scheme is simplest to implement
  - $\diamond$  1 bit per branch instruction
  - ♦ Record last outcome of a branch instruction (Taken/Not taken)
  - ♦ Use last outcome to predict future behavior of a branch



STUDENTS-HUB.com

### 1-Bit Predictor: Shortcoming

- Inner loop branch 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



## 1-Bit Predictor: Example

Assume initially branch is Not Taken						
i		0	3	4		
For Br	Pred	NT	Т	NT	Т	NT
	Act	Т	Т	Т	Т	NT
		®×X		₽×X	$\checkmark$	$\checkmark$
lf Br	Pred	Т	Т	Т	Т	
	Act	Т	NT	Т	NT	
		$\checkmark$	®× ×	$\checkmark$	×	

## 1-Bit Predictor: Example

Ass	ume init					
		0	3	4		
For Br	Pred	Т	Т	NT	Т	NT
	Act	Т	Т	Т	Т	NT
		$\checkmark$	$\checkmark$	×	$\checkmark$	$\checkmark$
lf Br	Pred	Т	Т	Т	Т	
	Act	Т	NT	Т	NT	
		$\checkmark$	®× ×	$\checkmark$	×	

## Dynamic Branch Prediction - One-bit Branch History Table (BHT)

- Predict branch based on past history of branch
- One-bit Branch History Table (BHT)



## Dynamic Branch Prediction - One-bit Branch History Table (BHT)

One-bit Branch History Table (BHT)



STUDENTS-HUB.com

Uploaded By: Jibreel Bornat

## 1-bit Branch History Table (BHT) Example

Assume BHT initialized as Taken						
		0	1	3	4	
For Br	Pred	Т	Т	Т	Т	Т
	Act	Т	Т	Т	Т	NT
		$\checkmark$				₿×
lf Br	Pred	Т	Т	NT	Т	
	Act	Т	NT	Т	NT	
		$\checkmark$	®. ×	®× ×	®× ×	

### 2-bit Prediction Scheme

- 1-bit prediction scheme has a performance shortcoming
- 2-bit prediction scheme works better and is often used
  - ♦ 4 states: strong and weak predict taken / predict not taken
- Implemented as a saturating counter
  - ♦ Counter is incremented to max=3 when branch outcome is taken
  - ♦ Counter is decremented to min=0 when branch is not taken



#### STUDENTS-HUB.com

Uploaded By: Jibreel Bornat

### 2-bit Prediction Scheme with BHT

2-bit scheme change prediction only if we get two mispredictions



## 2-bit Branch History Table (BHT) Example

							7	
	Assume BHT initialized as Weakly Taken (10)					<pre>For(i=0; a&lt;4; i++)</pre>		
i			0	1	2	3	4	{
F	or Br	Pred	(10) T	(11) T	(11) T	(11) T	(11) T	
		Act	Т	Т	Т	Т	(10) T	if(a%2==0)
								{}
			$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	₿×	Flsp
lf	Br	Pred	(10) T	(11) T	(10) T	(11) T		{}
		Act	Т	NT	т	NT		
			$\checkmark$	₿×	$\checkmark$	₽×		a++;
								}

STUDENTS-HUB.com

### 2-bit Prediction Scheme with BHT Accuracy

### Mispredict because either:

- $\diamond$  Wrong guess for that branch
- $\diamond$  Got branch history of wrong branch when index the table

### ✤ 4096 entry table:



## For More Advanced Branch Prediction ...

- Hypothesis: recent branches are correlated; that is, behavior of recently executed branches affects prediction of current branch
- Two possibilities: current branch depends on
  - Local behavior: Last *m* outcomes of the same branch (*local* branch predictor), e.g., a loop of 3 iterations is executed repetitively
    - $\rightarrow$  a history record of the loop branch of the last 6 iterations should be able to predict the direction of that branch correctly
  - ♦ Global behavior: Last *m* most recently executed branches
     → because branches are often *correlated*!

## Local Correlation: Yeh-Patt predictor

- It is possible to do quite well considering only information about the current branch (local information).
- We do this by considering, What happened the last time a branch had the same history that the current branch does now?
  - $\diamond$  This diagram shows how it operates.



#### STUDENTS-HUB.com

# The Yeh-Patt predictor Example

 Let us work out an example, assuming that two branches have this history:

### A: T N T N T N T N B: T T T T T T T T



STUDENTS-HUB.com

PT entries 01 and 10 are "trained" for Alandades dtrained ifor Bel Bornat

## **Global Correlation**

Branch direction of multiple branches

♦ Not independent but correlated to the path taken

Example: path 1-1 of b3 can be known beforehand

```
if (aa==2) // b1
    aa = 0;
if (bb==2) // b2
    bb = 0;
if (aa!=bb) {// b3
    .....
}
```

#### How to capture global behavior?

Idea: record *m* most recently executed branches as taken or not taken, and use that pattern to select the proper *n*-bit branch history table



# Global Correlation: Gshare Predictor

- The Gshare architecture uses an m-bit global history register to keep track of the direction of the last m executed branches.
- To simplify the implementation, this global history register is xored with the (PC>>2) to create an index into a 2<sup>m</sup>-entry pattern history table of n-bit counters.
- The result of this index is the prediction for the current branch.
- The predictor then compares this prediction with the real branch direction to determine if the branch was correctly predicted or not, and updates the prediction statistics.
- The predictor then updates the n-bit counter used to perform the prediction.
- The counter is
  - $\diamond$  Incremented if the branch was taken, and
  - $\diamond$  Decremented if the branch was not taken.
- Finally, the branch outcome is shifted into the most significant bit of the global history register: STUDENTS-HUB.com
  Uploaded By: Jibreel Bornat

### **Gshare Predictor Architecture**



STUDENTS-HUB.com

### Two-Level Global Adaptive Branch Predictors (GAp)

Improve branch prediction by looking not only at the history of the branch in question but also at that of other branches using two levels of branch history.



 Record the global pattern or history of the m most recently executed branches as taken or not taken. Usually an m-bit shift register.

#### ♦ Second level (per branch address):

**Pattern History Tables (PHTs)** 

- 2<sup>m</sup> prediction tables, each table entry has n bit saturating counter.
- The branch history pattern from first level is used to select the proper branch prediction table in the second level.
- The low N bits of the branch address are used to select the correct prediction entry (predictor)within a the selected table, thus each of the 2<sup>m</sup> tables has 2<sup>N</sup> entries and each entry is 2 bits counter.
- Total number of bits needed for second level =  $2^m x n x 2^N$  bits
- ✤ In general, the notation: <u>GAp (m,n) predictor means</u>:
  - ♦ Record last m branches to select between 2<sup>m</sup> history tables.

♦ Each second level table uses n-bit counters (each table entry has n bits).

Basic two-bit single-level Bimodal BHT is then a (0,2) predictor. STUDENTS-HUB.com

### Two-Level Dynamic GAp Architecture Example



### An Example: GAp(1,1)

; b2

bnez R3, L2

. . .

. . .

Initial value of d	d==0?	b1	Value of d before b2	d==1?	b2
0	Yes	NT	1	Yes	NT
1	No	т	1	Yes	NT
2	No	Т	2	No	Т

=> if b1 is NT, then b2 is NT

## Behavior of one-bit Standard Predictor initialized to not taken; d alternates between 0 and 2.

d=?	b1 prediction	b1 action	New b1 prediction	b2 prediction	b2 action	new b2 prediction
2	N	Т	р т	NT	Т	Т
0	Т	NT	D NT	Т	NT	NT
2	NT	Т	Т	NT	Т	Т
0	Т	NT	NT	Т	NT	NT

=> All branches are mispredicted

STUDENTS-HUB.com

L1:

12.

### An Example: GAp(1,1)

### Introduce one bit of correlation

Each branch has two separate prediction bits: one prediction assuming the last branch executed was not taken, and another prediction assuming it was taken

Prediction bits	Prediction if last branch NT	Prediction if last branch T
NT/NT	NT	NT
NT/T	NT	Т
T/NT	т	NT
T/T	Т	Т

# Behavior of one-bit predictor with one bit of correlation initialized to NT/NT; Assume <u>last</u> branch NT

d=?	b1	b1	New b1	b2 prediction	b2	new b2		?	N
	prediction	action	prediction		action	prediction		b1	Т
2	NT/NT	т	T/NT	NT/NT	т	NT/T	lİ	b2	Т
0	T/NT	NT	T/NT	NT/T	NT	NT/T		b1	N
	ТЛІТ	<b>–</b>	ТЛІТ		<b>–</b>			b2	N
2		I	I/NI	N1/1	I	N1/1		b1	Т
0	T/NT	NT	T/NT	NT/T	NT	NT/T	ŀ	 h0	
							ļ	02	

#### => Only misprediction is on the first iteration

STUDENTS-HUB.com

Uploaded By: Jibreel Bornat

b1

b2

NT

NT

## An Example: GAp(2,2)



All entries in the matrix are initialized to 01. In the tables next slide, matrix entries are italicized if they were updated due to the execution of the previous branch. The entry with the yellow background is the entry used for the prediction.

STUDENTS-HUB.com

# An Example: GAp(2,2)



### Accuracy of Different Schemes



STUDENTS-HUB.com

## Hybrid Branch Predictors

- Idea: Use more than one type of predictor (i.e., multiple algorithms) and select the "best" prediction
  - $\diamond$  E.g., hybrid of 2-bit counters and global predictor
- ✤ Advantages:
  - + Better accuracy: different predictors are better for different branches
  - + Reduced warmup time (faster-warmup predictor used until the slowerwarmup predictor warms up)
- Disadvantages:
  - -- Need "meta-predictor" or "selector" to decide which predictor to use
  - -- Longer access latency

### Hybrid Branch Predictor - Basic Schema

- Some branches correlated to global history, some correlated to local history
  - ♦ Use more than one type of predictors and select "best"



STUDENTS-HUB.com

## Example: Tournament Predictor



STUDENTS-HUB.com

# Comparing Accuracy of Different Predictors

The misprediction rate for three different predictors on SPEC89 versus the size of the predictor in kilobits.



## Performance of Core i7 920 and core i7 6700

The misprediction rate for the integer SPECCPU2006 benchmarks on the Intel Core i7 920 and 6700.



# Zero-Delayed Branching

- How to achieve zero delay for a jump or a taken branch?
  - $\diamond\,$  Jump or branch target address is computed in the ID stage
  - ♦ Next instruction has already been fetched in the IF stage

### **Solution**

- Introduce a Branch Target Buffer (BTB) in the IF stage
  - ♦ Store the target address of recent branch and jump instructions
- Use the lower bits of the PC to index the BTB
  - ♦ Each BTB entry stores Branch/Jump address & Target Address
  - ♦ Check the PC to see if the instruction being fetched is a branch
  - ♦ Update the PC using the target address stored in the BTB

STUDENTS-HUB.com

# Branch Target Buffer (IF Stage)

- The branch target buffer is implemented as a small cache
  - $\diamond$  Stores the target address of recent branches and jumps
- We must also have prediction bits
  - $\diamond$  To predict whether branches are taken or not taken
  - $\diamond~$  The prediction bits are determined by the hardware at runtime



### Branch Target Buffer - cont'd

- Each Branch Target Buffer (BTB) entry stores:
  - ♦ Address of a recent jump or branch instruction
  - ♦ Target address of jump or branch
  - ♦ Prediction bits for a conditional branch (Taken or Not Taken)
  - To predict jump/branch target address and branch outcome before instruction is decoded and branch outcome is computed
- Use the lower bits of the PC to index the BTB
  - ♦ Check if the PC matches an entry in the BTB (jump or branch)
  - If there is a match and the branch is predicted to be Taken then Update the PC using the target address stored in the BTB
- The BTB entries are updated by the hardware at runtime STUDENTS-HUB.com
  Uploaded By: Jibreel Bornat

# Dynamic Branch Prediction with BTB



STUDENTS-HUB.com

### Gshare Predictor with BTB



Cache of Target Addresses (BTB: Branch Target Buffer)

### **Tournament Predictor with BTB**



# Branch Target Cache

- Similar to BTB, but we also want to get the *target instruction*!
  - Prediction returns not just the target address, but also the instruction stored there
  - Allows zero-cycle *unconditional* branches (branch-folding)
    - Send target-instruction to ID rather than branch
    - Branch is not even sent into pipe



© 2003 Elsevier Science (USA). All rights reserved.

### How about Subroutine Returns?

- Different call sites make return address hard to predict
  - ♦ printf() may be called by many callers
  - ♦ Target of "return" instruction in printf() is a moving target
- But return address is actually easy to predict
  - It is the address after the last call instruction that have not returned from yet
  - ♦ Can use a Return Address Stack (RAS)
- \* RAS:
  - $\diamond$  Call will push return address on the stack
  - ♦ Return uses the prediction of top-of-stack

### Return Address Stack



- May not know if it is a return instruction prior to decoding
  - Rely on BTB for speculation
  - Fix once recognize Return

STUDENTS-HUB.com