Dynamic Branch Prediction

Special Thanks to Dr. Muhamed Mudawar (KFUPM) and Dr. Onur Mutlu (ETHZ) for most of the Slides

STUDENTS-HUB.com

Recall: Control Dependence

- Question: What should the fetch PC be in the next cycle?
- Answer: The address of the next instruction
 - ✤ All instructions are control dependent on previous ones. Why?
- ✤ 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
- 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?

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)

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 → 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
 - ♦ 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

- Modern processors have deep pipelines
 - ♦ Branch penalty limits performance of deep pipelines
- ✤ Want to execute instructions beyond a branch even before that branch is resolved → use speculative execution
 - ♦ Branch prediction: dynamic vs. static
- What to predict?

IBM POWER4 Processor Pipeline



STUDENTS-HUB.com https://www.ixbt.com/cpu/ibm/power4/p4-pipe.png Uploaded By: Jibreel Bornat

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



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	2	3	4
Pred	Т	Т	Т	Т	Т
Act	Т	Т	Т	Т	NT
	\checkmark	\checkmark	\checkmark	\checkmark	®×
Pred	Т	Т	NT	Т	
Act	Т	NT	Т	NT	
		₿X	®× X	®×	
	Pred Act Pred Act	ame BHT initializ	ume BHT initialized as TaO1PredTActTV√PredTActTV×Second Colspan="2">NT√§	ume BHT initialized as Taken012PredTTActTTPredT√V√PredTActTNTNTActTNT\$\$\$	Jume BHT initialized as Taken0123PredTTTActTTTV√√√PredTT√ActTNTTActTNTT√<<√<<<<<

STUDENTS-HUB.com

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

						-	int a=0.
Assume BHT initialized as Weakly Taken (10)							
i		0	1	2	3	4	For(i=0; a<4; i++)
For 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	₿×	
lf Br	Pred	(10) T	(11) T	(10) T	(11) T		EISE
	Act	т	NT	т	NT		{}
		\checkmark	₽ ×	\checkmark	₽×		
							a++;
							}

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 → 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*!

BHT predicts this

Local vs. Global History Predictor Structure



Source: Combining Branch Predictors by Scott McFarling

STUDENTS-HUB.com

Branches Are Correlated!

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



Correlating Predictor: McFarling's 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 2m-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
 - Decremented if the branch was not taken. Finally, the branch outcome is shifted into the most significant bit of the global history register:

Correlating Predictor: McFarling's gshare predictor

- + Better utilization of PHT
- Increases access latency



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
 - ♦ 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

Example: Evaluating Branch Alternatives

Branch Scheme	Jump	Branch Not Taken	Branch Taken
Predict not taken	Penalty = 1 cycle	Penalty = 0 cycles	Penalty = 2 cycles
Delayed	Penalty = 0 cycles	Penalty = 0 cycles	Penalty = 1 cycle
BTB Prediction	Penalty = 1 cycle	Penalty = 2 cycles	Penalty = 2 cycles

- ✤ Assume: Jump = 2%, Branch-Not-Taken = 5%, Branch-Taken = 15%
- ✤ Assume a branch target buffer with hit rate = 90% for jump & branch
- Prediction accuracy for jump = 100%, for conditional branch = 95%
- ✤ What is the impact on the CPI? (Ideal CPI = 1 if no control hazards)

Evaluating Branch Alternatives- Solution

Branch Scheme	Jump = 2%	Branch NT = 5%	Branch Taken = 15%	СРІ
Predict not taken	0.02 × 1	0	$0.15 \times 2 = 0.30$	1+0.32
Delayed	0	0	0.15 × 1 = 0.15	1+0.15
BTB Prediction	0.02×0.1×1	0.05×0.9×0.05×2	0.15×(0.1+0.9×0.05)×2	1+0.05