

ECE 4750 Computer Architecture, Fall 2021

T13 Advanced Processors: Branch Prediction

School of Electrical and Computer Engineering
Cornell University

revision: 2021-10-26-18-28

1	Branch Prediction Overview	3
2	Software-Based Branch Prediction	4
2.1.	Static Software Hints	5
2.2.	Branch Delay Slots	6
2.3.	Predication	7
3	Hardware-Based Branch Prediction	8
3.1.	Fixed Branch Predictor	8
3.2.	Branch History Table (BHT) Predictor	10
3.3.	Two-Level Predictor For Temporal Correlation	14
3.4.	Two-Level Predictor For Spatial Correlation	15
3.5.	Generalized Two-Level Predictors	16
3.6.	Tournament Predictors	17
3.7.	Branch Target Buffers (BTBs) Predictor	18

Copyright © 2018 Christopher Batten, Christina Delimitrou. All rights reserved. This hand-out was originally prepared by Prof. Christopher Batten at Cornell University for ECE 4750 / CS 4420. It has since been updated by Prof. Christina Delimitrou. Download and use of this handout is permitted for individual educational non-commercial purposes only.

1. Branch Prediction Overview

Assume incorrect branch prediction in dual-issue I2OL processor.

bne													
opA													
opB													
opC													
opD													
opE													
opF													
opG													
opTARG													

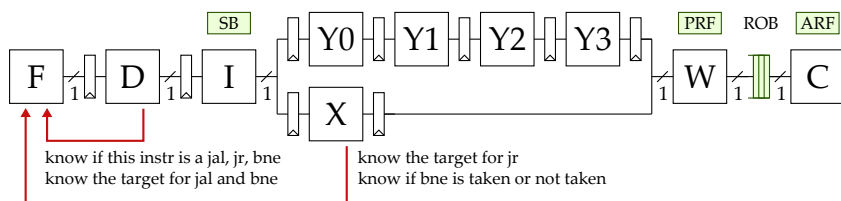
Assume correct branch prediction in dual-issue I2OL processor.

bne													
opA													
opTARG													
opX													
opY													
opZ													

Three critical pieces of information we need to predict control flow:

- (1) Is this instruction a control flow instruction?
- (2) What is the target of this control flow instruction?
- (3) Do we redirect control flow to the target or next instr?

When do we know these critical pieces of information?



	jal	jr	bne
(1) Is this instruction a control flow instruction?	D	D	D
(2) What is the target of this control flow instruction?	D	X	D
(3) Do we redirect ctrl flow to the target or next instr?	D	D	X

What do we need to predict in F stage vs. D stage?

	jal	jr	bne
F stage	predict 1,2,3	predict 1,2,3	predict 1,2,3
D stage	no prediction	predict 2	predict 3

2. Software-Based Branch Prediction

- Static software hints
- Branch delay slots
- Predication

2.1. Static Software Hints

Software provides hints about whether a control flow instruction is likely to be taken or not taken. These hints are part of the instruction and thus are available earlier in the pipeline (e.g., in the D stage).

bne.t														
opA														
opTARG														
bne.nt														
opY														
opZ														

What if the hint is wrong?

bne.t														
opA														
opTARG														
bne.nt														
opA														
opB														

2.2. Branch Delay Slots

Without branch delay slots must squash fall through instructions if branch is taken.

bne														
opA														
opB														
targ														

With branch delay slots compiler can put useful work in the slots. Instructions in the delay slots are always executed regardless of branch condition.

bne														
opA														
opB														
targ														

2.3. Predication

Not really “prediction”. Idea is to turn control flow into dataflow completely eliminating the control hazard.

Conditional move instructions conditionally move a source register to a destination register.

`movn rd, rs1, rs2` if ($R[rs2] \neq 0$) $R[rd] \leftarrow R[rs1]$

`movz rd, rs1, rs2` if ($R[rs2] == 0$) $R[rd] \leftarrow R[rs1]$

Pseudocode	w/o Predication	w/ Predication
if (a < b)	<code>slt x1, x2, x3</code>	<code>slt x1, x2, x3</code>
x = a	<code>beq x1, x0, L1</code>	<code>movz x4, x2, x1</code>
else	<code>addi x4, x2, x0</code>	<code>movn x4, x3, x1</code>
x = b	<code>jal x0, L2</code>	
	L1:	
	<code>addi x4, x3, x0</code>	
	L2:	

Full predication enables almost all instructions to be executed under a predicate. If predicate is false, instruction should turn into a NOP.

Pseudocode	w/ Predication
if (a < b)	<code>slt.p p1, x2, x3</code>
opA	<code>(p1) opA</code>
opB	<code>(p1) opB</code>
else	<code>(!p1) opC</code>
opC	<code>(!p1) opD</code>
opD	

- What if both sides of branch have many instructions?
- What if one side of branch has many more than the other side?

3. Hardware-Based Branch Prediction

- Fixed branch predictor
- Branch history table (BHT) predictor
- Two-level predictor for temporal correlation
- Two-level predictor for temporal correlation
- Generalized two-level predictors
- Tournament predictor
- Branch target buffer (BTB) predictor

3.1. Fixed Branch Predictor

- Always predict not taken
 - What we have been assuming so far
 - Simple to implement and can perform prediction in F
 - Poor accuracy, especially on very important backwards branch in loops
- Always predict taken
 - Difficult to implement: we don't know if this is a branch until D
 - Difficult to implement: we don't know target until at least D
 - Could predict not taken in F, and then adjust in D
 - Poor accuracy, especially on if/then/else
- Predict taken for backward branches and predict not taken for forward branches
 - Difficult to implement: we don't know if this is a branch until D
 - Difficult to implement: we don't know target until at least D
 - Could predict not taken in F, and then adjust in D
 - Better accuracy

```

loop:          <----- .
    lw  x1, 0(x2)          | backward
    lw  x3, 0(x4)          | branches
    slt  x5, x1, x3        | taken on avg
    beq  x5, x0, L1        --. forward | 90%
    addi x6, x1, x0        | branches |
    jal  x0, L2            | taken on avg |
L1:          <-' 50%      |
    addi x6, x3, x0        |
L2:          |
    sw   x6, 0(x7)         |
    addi x2, x2, 4          |
    addi x4, x4, 4          |
    addi x7, x7, 4          |
    addi x8, x8, -1         |
    bne  x8, x0, loop -----'

```

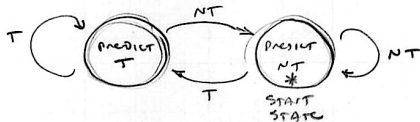
- For now let's focus on conditional branches as opposed to unconditional jumps
- Let's assume we always predict not-taken in the F stage
- In the D stage, we know if the instruction is a branch and we know the target of the branch
- So key goal is to predict whether or not we need to redirect the control flow, i.e., to predict the branch outcome in the D stage instead of waiting until the X stage
- By doing this prediction in the D stage we can reduce the branch misprediction penalty by several cycles although it is still not zero if we predict the branch is taken

3.2. Branch History Table (BHT) Predictor

How can we do better? Exploit structure in the program, namely **temporal correlation**: the outcomes of specific static branch in the past may be a good indicator of the outcomes of future dynamic instances of the same static branch.

One-Bit Saturating Counter

Remember the previous outcome of a specific static branch and predict the outcome will be the same for the next dynamic instance of the same branch.



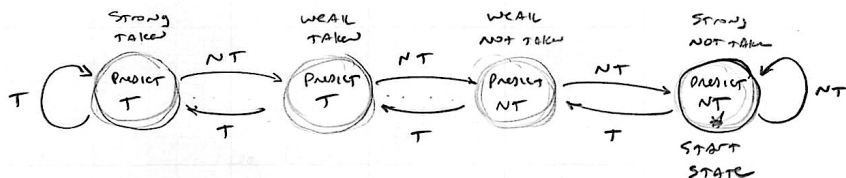
Consider how this saturating counter would be have for a backwards branch in a loop with four iterations. Assume the entire loop is executed several times.

Iteration	Prediction	Actual	Mispredict?
1			
2			
3			
4			
1			
2			
3			
4			

Exploiting temporal correlation works well, but a one-bit saturating counter will always mispredicts the backwards branch in a loop twice. Loops are *very* common!

Two-Bit Saturating Counter

Remember the last *two* outcomes of a specific static branch. Require two consecutive “counter examples” before changing the prediction.

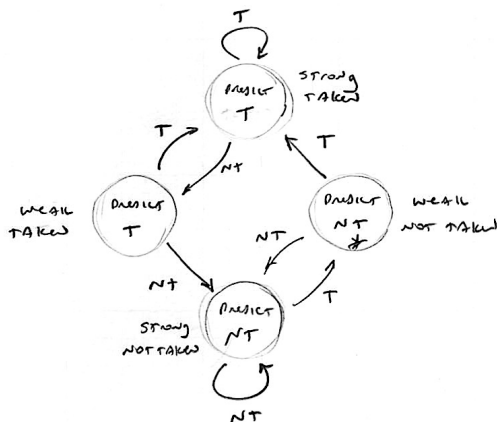


Consider how this saturating counter would behave for a backwards branch in a loop with four iterations. Assume the entire loop is executed several times.

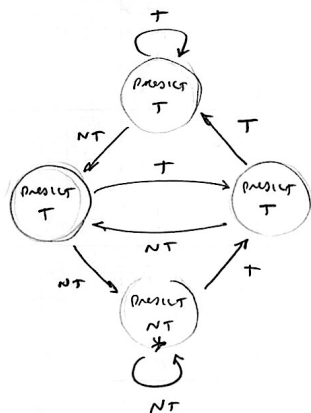
Iteration	Prediction	Actual	Mispredict?	ST	WT	WNT	SNT
1				○	○	○	○
2				○	○	○	○
3				○	○	○	○
4				○	○	○	○
1				○	○	○	○
2				○	○	○	○
3				○	○	○	○
4				○	○	○	○

What if start state is strongly taken?

Other Two-Bit FSM Branch Predictors



JUMP DIRECTLY TO
STRONG FROM WEAK

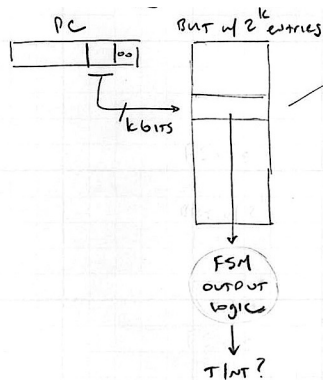


BIASED TOWARDS
PREDICTING BRANCHES
AS TAKEN

See Fig 5.8 in SHEN + LIPASTI for other alternatives

Branch History Table

- So far we have focused on a simple FSM that exploits temporal correlation to make a prediction for a *specific static branch*
- To make predictions for many different static branches, we need to keep track of a *dedicated* FSM per static branch
- A **branch history table** (BHT) is a table where each entry is the state of the FSM for a different static branch.



- Two PC's can "alias" to the same entry in BHT
- Aliasing is similar to a cache conflict
- We could store the PC as a tag along with the FSM state to make sure we don't mix up the FSM state across two static branches
- Storing the PC is too expensive though, so we can just let branches alias and this just reduces the branch prediction accuracy
- Can reduce aliasing with larger BHT

BHT with 4k entries and 2bits/entry = 80–90% accuracy

How do we continue to improve prediction accuracy? Exploit even more complicated temporal correlation.

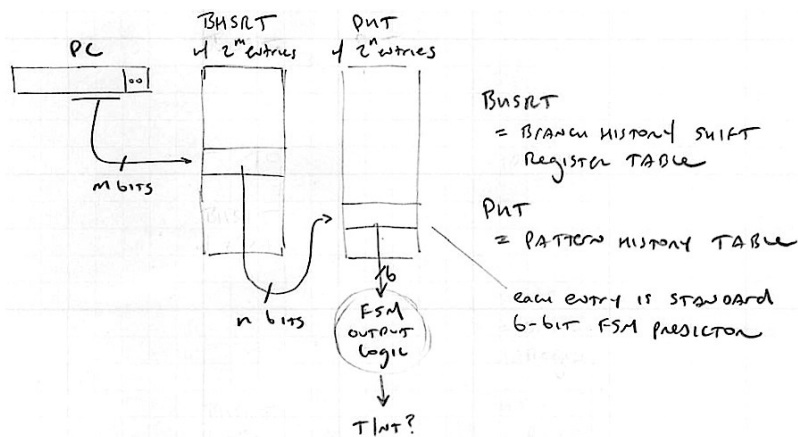
MORE COMPLICATED TEMPORAL CORRELATION

OFTEN A BRANCH EXHIBITS MORE COMPLICATED PATTERN THAN JUST "ALWAYS TAKEN" OR "ALWAYS NOT TAKEN". COULD DEVELOP A MORE COMPLICATED FSM, BUT THEN PATTERNS VARY PER BRANCH. WE WANT PER BRANCH CUSTOMIZED FSMs.

```
void convolve( int B[], int A[], int size ) {
    for (int i = 2; i < size - 2; i++)
        for (int j = 0; j < 5; j++)
            B[i - (2-j)] = A[i] * COEFF[j]
}
```

CAN WE PREDICT THAT EVERY FIFTH DYNAMIC INSTANCE OF THE BACKWARDS LOOP BRANCH WILL BE NOT TAKEN?

3.3. Two-Level Predictor For Temporal Correlation



When a branch is taken or not taken we shift in either a one (taken) or a zero (not taken) into the least significant bit of the corresponding BHSR.

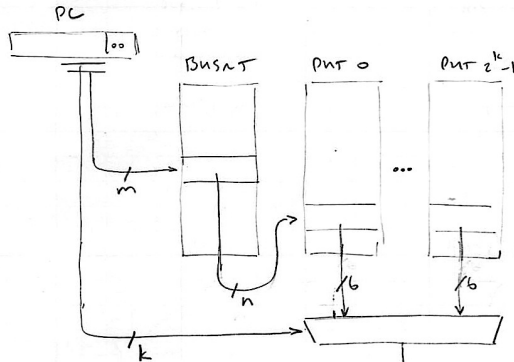
Index	Value
...	
0111	ST
1000	WT
1001	WT
1010	WT
1011	ST
1100	WT
1101	ST
1110	ST
1111	SNT

- BHSR captures temporal pattern for that branch
- We use the BHSR to index into the PHT. A BHT has an entry per branch, but a PHT has an entry per branch pattern.
- The PHT says for a given pattern over the past n executions of a branch, should I take or not take the next execution of this branch?
- Once the two-level predictor is warmed up for previous nested loop example, the state of the PHT would be what is shown on the left
- Need at least four bits of “history” to learn this pattern and perfectly predict this branch

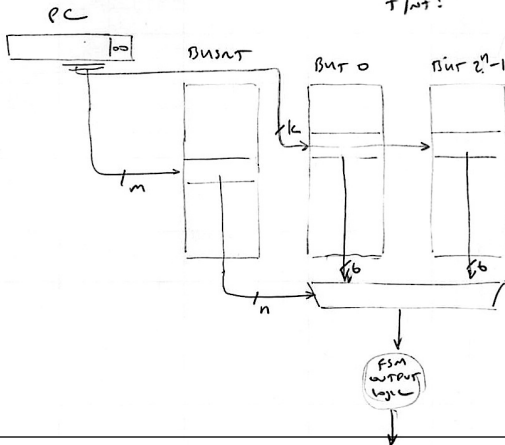
3. Hardware-Based Branch Prediction 3.3. Two-Level Predictor For Temporal Correlation

PROBLEM: MULTIPLE BRANCHES WITH SAME HISTORY MIGHT NEED DIFFERENT PREDICTIONS. IN OTHER WORDS, A LIASING IN THE PUT CAN REDUCE ACCURACY

SOLUTION: ADD MULTIPLE PUTS, USE BITS FROM PC TO CHOOSE WHICH PUT TO USE



ISOMORPHIC - TWO DIFFERENT WAYS OF DRAWING THE SAME TWO-LEVEL STRUCTURE.



15/INT?

HW BPNED: EXPLOITING SPATIAL CORRELATION

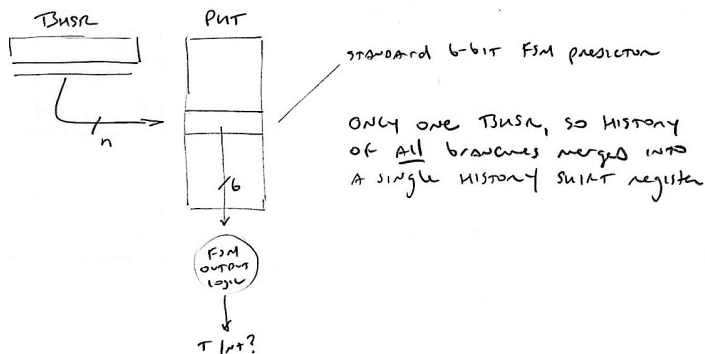
The way one branch is resolved may be a good indication of the way a later (different) branch will resolve

```

if (x < 7)          slt r2, r1, 7
    y++             beq r2, L1          branch 0
if (x < 5)          addi r3, r3, 1
    z++             L1:
                    slt r2, r1, 5
                    beq r2, L2          branch 1
                    addi r4, r4, 1
                    L2:
  
```

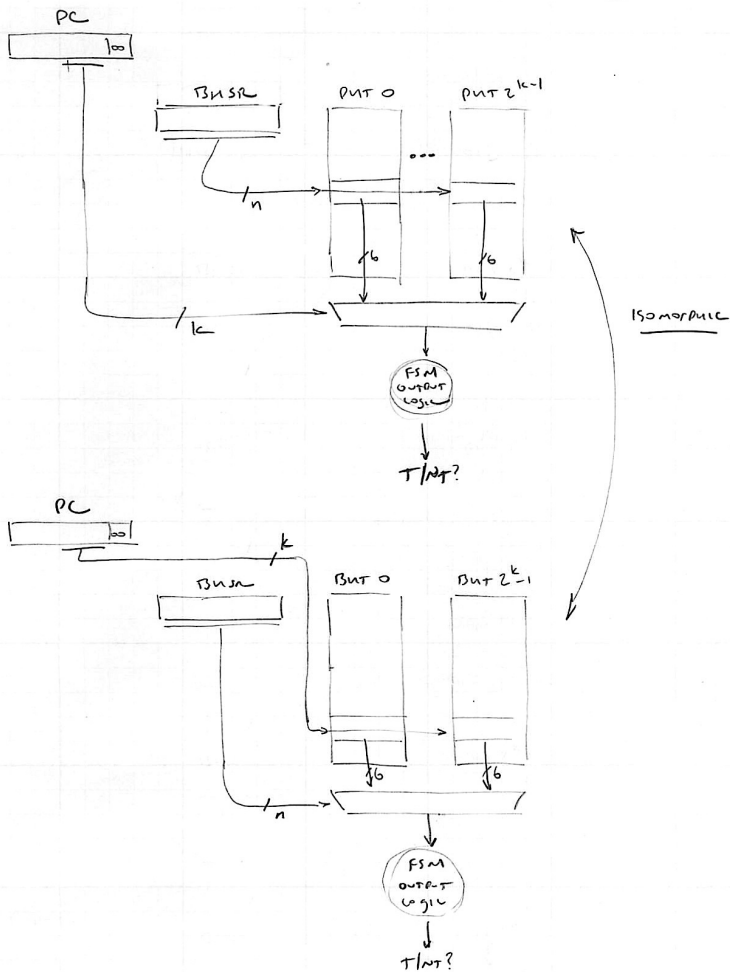
if branch 0 is taken (ie $x \geq 7$)
then branch 1 is always taken (ie x must be ≥ 5)

so whether branch 0 is taken or not taken can be used to predict if we should take branch 1



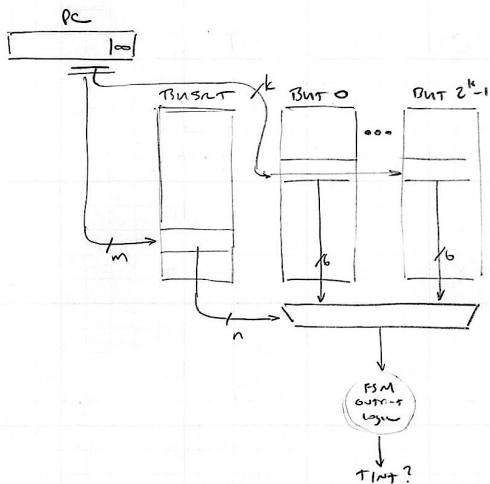
For above example, BPNR will capture history - so we will know when first branch is taken, and that value (r) of the BPNR will point to an entry in the PNR that predicts taken.

As before, multiple PNTs CAN HELP AVOID ALIASING W PNT



GENERALIZED TWO-LEVEL BUTS

COMBINED APPROACH TO EXPLOIT BOTH COMPLEX TEMPORAL CORRELATION AND SPATIAL CORRELATION



Difference from discussed on complex temporal correlation is that we purposely choose a smaller m to cause aliasing in the BUTs. Since this aliasing allows us to capture spatial correlation.

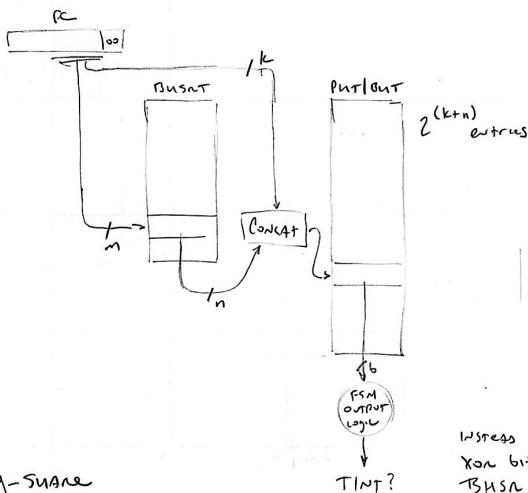
(choose higher order m bits)

		put for each PC			
		one put $k = 0$	$0 < k < 30$	$k = 30$	
ONE BULK	$m = 0$	GA _g	GA _s	GA _p	
	$0 < m < 30$	PA _g	PA _s	PA _p	
ONE BULK for each PC		$m = 30$	SA _g	SA _s	SA _p

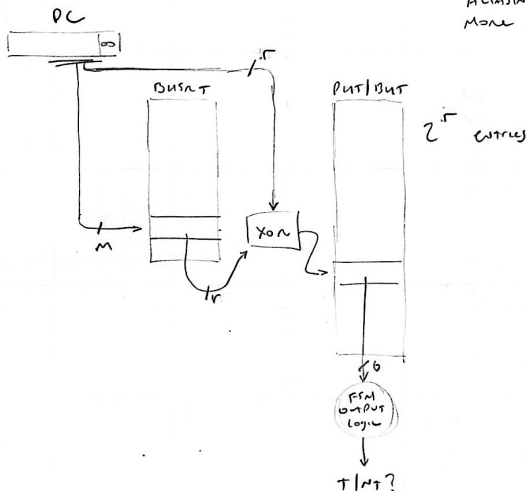
97% ACCURACY

g-select

ISOMORPHIC TO PREVIOUS FIGURE



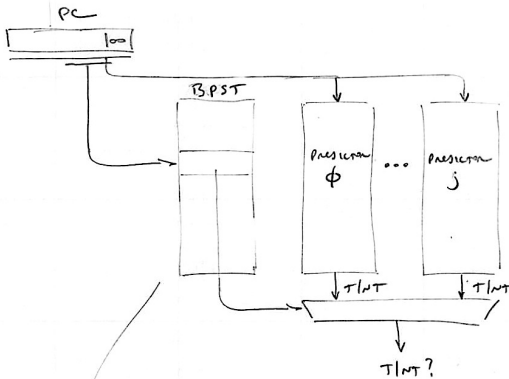
Wasteful of concatenating
XOR bits from PC with
BHSN, helps AVOID
ALIASING in the PUT/BUT
MORE EFFECTIVELY

g-square

TOURNAMENT PREDICTIONS

Different Predictors Are better At Predicting Different Types of Branches

- one-level 2-bit saturating counter
- two-level gsnare
- loops
- irregular code



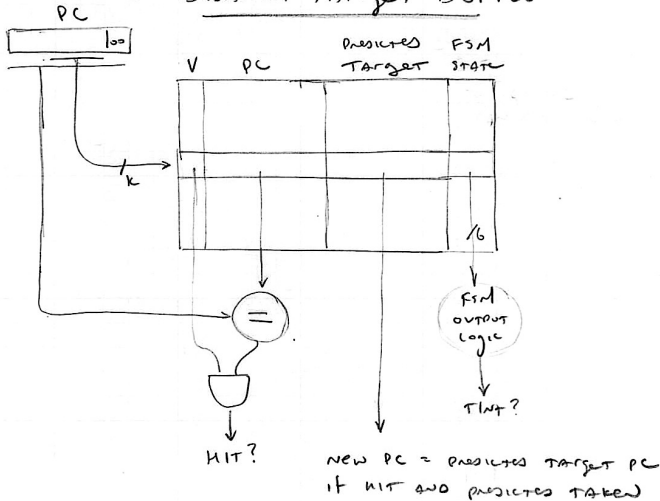
Branch Predictor Selection Table

- predicts which branch predictor we should use

HW BTB: PREDICTING TARGET ADDRESS

EVEN WITH BEST POSSIBLE PREDICTION OF BRANCH OUTCOME
STILL NEED TO WAIT FOR TARGET ADDRESS TO BE DETERMINED.

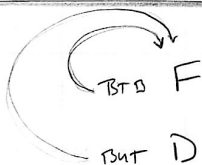
BRANCH TARGET BUFFER



CAN PUT BTB IN FETCH STAGE

- PREDICTING IF PC POINTS TO A BRANCH
- PREDICTING TARGET OF BRANCH
- PREDICTING IF BRANCH IS TAKEN.

SOMETIMES $b=0$, if hit then ASSUME PREDICT TAKEN

COMBINE BTB AND BUT

BUT D UPDATE BTB for J/branches

BTB is much more expensive than BUT, BUT BTB earlier in pipeline + can accelerate JR

I

X UPDATE BUT for branches
UPDATE BTB for JR

COMBINE BTB w/ few entries with BUT w/ many entries

M

W

RETURN ADDRESS STACK PREDICTION

BTB only works for JR function call returns it
Always call function from same place (NOT realistic)

STACK PREDICTION

- PUSH Target Address on stack for JAL/JALR
- POP off Target Address for JR to predict target

Move stack prediction into fetch AND predict which PC's are JR.

USE TOURNAMENT PREDICTION TO CHOOSE BETWEEN BTB AND STACK PREDICTION.