

Reduction of Control Hazards (Branch) Stalls with Dynamic Branch Prediction

➔ So far we have dealt with control hazards in instruction pipelines by:

- 1 – Assuming that the branch is not taken (i.e stall when branch is taken).
- 2 – Reducing the branch penalty by resolving the branch early in the pipeline
 - Branch penalty if branch is taken = stage resolved - 1 In IF ? $1-1 = 0?$
- 3 – Branch delay slot and canceling branch delay slot. (**ISA support needed**)
- 4 – Compiler-based static branch prediction encoded in branch instructions ➔
 - Prediction is based on program profile or branch direction
 - **ISA support needed.**

ISA
Support
Needed

How to further reduce the impact of branches on pipeline processor performance ?

Dynamic Branch Prediction:

Why? Better branch prediction accuracy than static prediction and thus fewer branch stalls

Why?

- Hardware-based schemes that utilize run-time behavior of branches to make dynamic predictions: + *No ISA support needed* ←

By CPU

How?

- Information about outcomes of previous occurrences of branches are used to dynamically predict the outcome of the current branch.

Branch Target Buffer (BTB): (Goal: zero stall taken branches) ←

- A hardware mechanism that aims at reducing the stall cycles resulting from ➔ **correctly predicted taken branches to zero cycles (No stalls).**

Static Conditional Branch Prediction

- ➔ Branch prediction schemes can be classified into static (at compilation time) and dynamic (at runtime) schemes.
- ➔ Static methods are carried out by the compiler. They are static because the prediction is already known before the program is executed.
- ➔ Static Branch prediction is encoded in branch instructions using one prediction (or branch direction hint) bit = **0 = Not Taken**, = **1 = Taken**
- ➔ **Must be supported by ISA**, Ex: HP PA-RISC, PowerPC, UltraSPARC

X = Static
Prediction bit
X=0 Not Taken
X=1 Taken

Branch Encoding

X

How?

Two basic methods to statically predict branches at compile time:

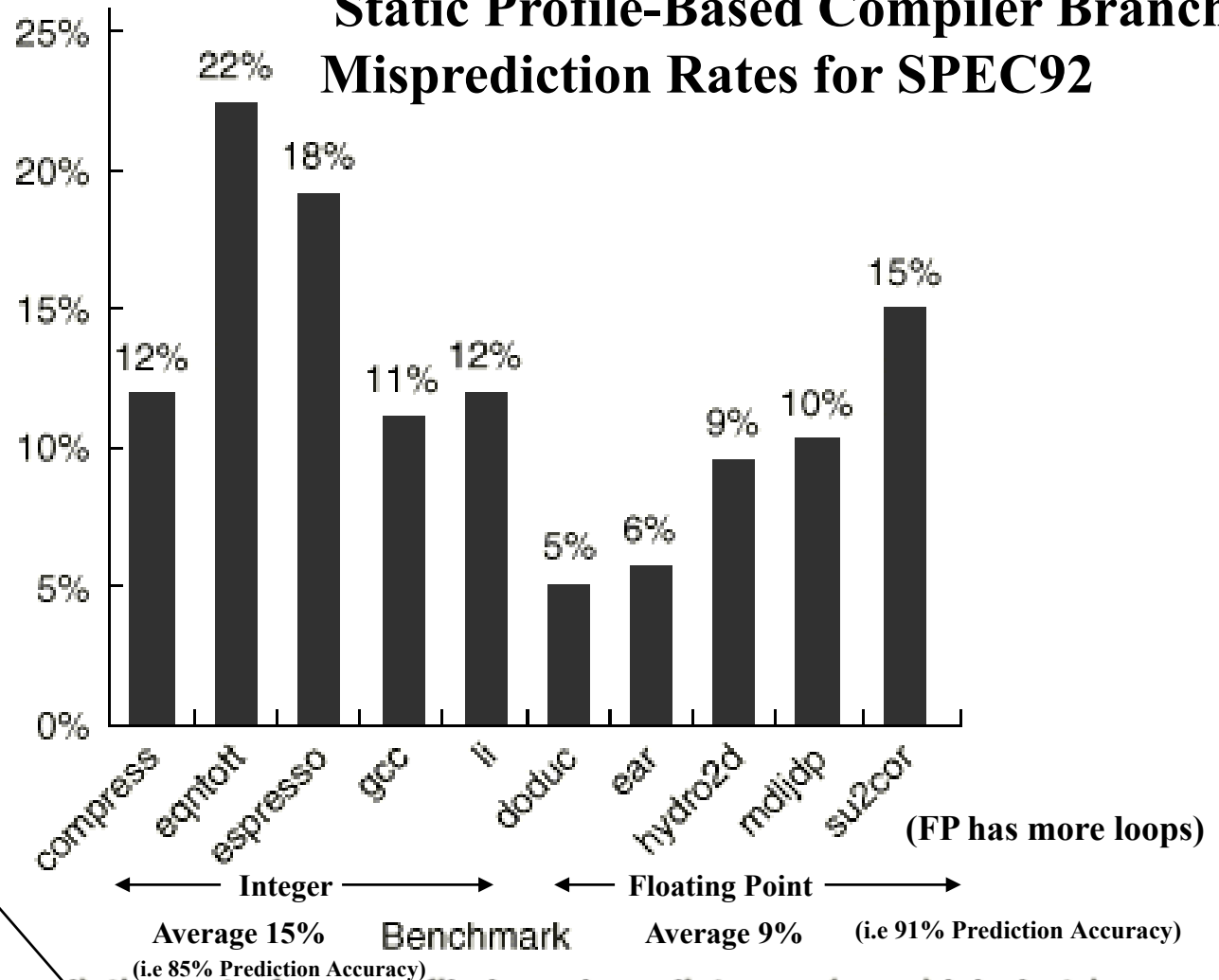
- 1 – Use the direction of a branch to base the prediction on. Predict **backward** branches (branches which decrease the PC) to be **taken** (e.g. branches used in loops) and **forward** branches (branches which increase the PC) **not to be taken**.
- 2 – Profiling can also be used to predict the outcome of a branch.
 - ➔ A number runs of the program are used to collect program behavior information (i.e. if a given branch is likely to be taken or not)
 - This information is included in the opcode of the branch (one bit branch direction hint) as the static prediction.

Static prediction was previously discussed in lecture 2
4th edition: in Chapter 2.3, 3rd Edition: In Chapter 4.2

CMPE550 - Shaaban

Static Profile-Based Compiler Branch Misprediction Rates for SPEC92

Misprediction rate



More Loops
In FP Code

Thus lower
misprediction rate
In FP code

Misprediction rate for a profile-based predictor varies widely but is generally better for the FP programs, which have an average misprediction rate of 9% with a standard deviation of 4%, than for the integer programs, which have an average misprediction rate of 15% with a standard deviation of 5%.

(repeated here from lecture2)

CMPE550 - Shaaban

Dynamic Conditional Branch Prediction

No ISA Support Needed

➔ • Dynamic branch prediction schemes are different from static mechanisms because they utilize hardware-based mechanisms that use the run-time behavior of branches to make more accurate predictions than possible using static prediction.

Why?

• Usually information about outcomes of previous occurrences of branches (branching history) is used to dynamically predict the outcome of the current branch. The two main types of dynamic branch prediction are:

➔ 1 – One-level or Bimodal: Usually implemented as a Pattern History Table (PHT), a table of usually two-bit saturating counters (predictors) which is indexed by a portion of the branch address (low bits of address). (First proposed mid 1980s)

• Also called *non-correlating* dynamic branch predictors. ←

➔ 2 – Two-Level Adaptive Branch Prediction. (First proposed early 1990s).

• Also called *correlating* dynamic branch predictors. ←

+ • To potentially reduce the stall cycles resulting from correctly predicted taken branches to zero cycles, a Branch Target Buffer (BTB) that includes the addresses of conditional branches that were taken along with their targets is added to the instruction fetch (IF) stage.

BTB discussed next →

BTB

4th Edition: Dynamic Prediction in Chapter 2.3, BTB in Chapter 2.9
(3rd Dynamic Prediction in Chapter 3.4, BTB in Chapter 3.5)

CMPE550 - Shaaban

Branch Target Buffer (BTB)

Why?

- Effective branch prediction requires the target of the branch at an early pipeline stage. (resolve the branch early in the pipeline)
 - One can use additional adders to calculate the target, as soon as the branch instruction is decoded. This would mean that one has to wait until the ID stage before the target of the branch can be fetched, taken branches would be fetched with a one-cycle penalty (this was done in the enhanced MIPS pipeline Fig A.24).

In IF ? $1-1 = 0?$

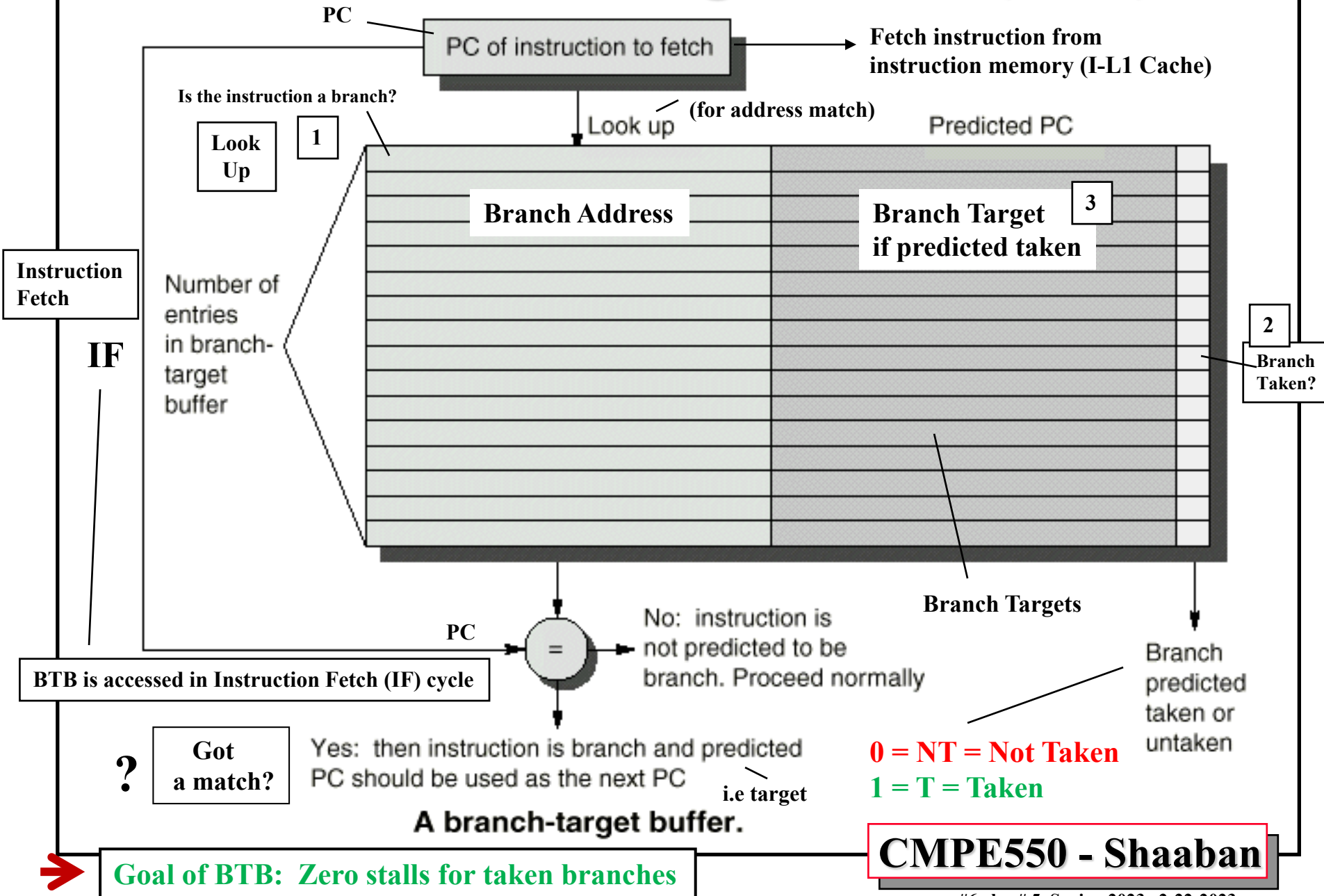
BTB Goal

How?

- To avoid this problem and to potentially achieve zero stall cycles for taken branches, one can use a Branch Target Buffer (BTB).
- A typical BTB is an associative memory where the addresses of taken branch instructions are stored together with their target addresses.
- The BTB is accessed in Instruction Fetch (IF) cycle and provides answers to the following three questions while the current instruction is being fetched:
 - 1 – Is the instruction a branch?
 - 2 – If yes, is the branch predicted taken?
 - 3 – If yes, what is the branch target?
- Instructions fetched from the target are stored in the BTB in case the branch is predicted-taken and thus found in BTB.
- After the branch has been resolved the BTB is updated. If a branch is encountered for the first time a new entry is created once it is resolved as taken.

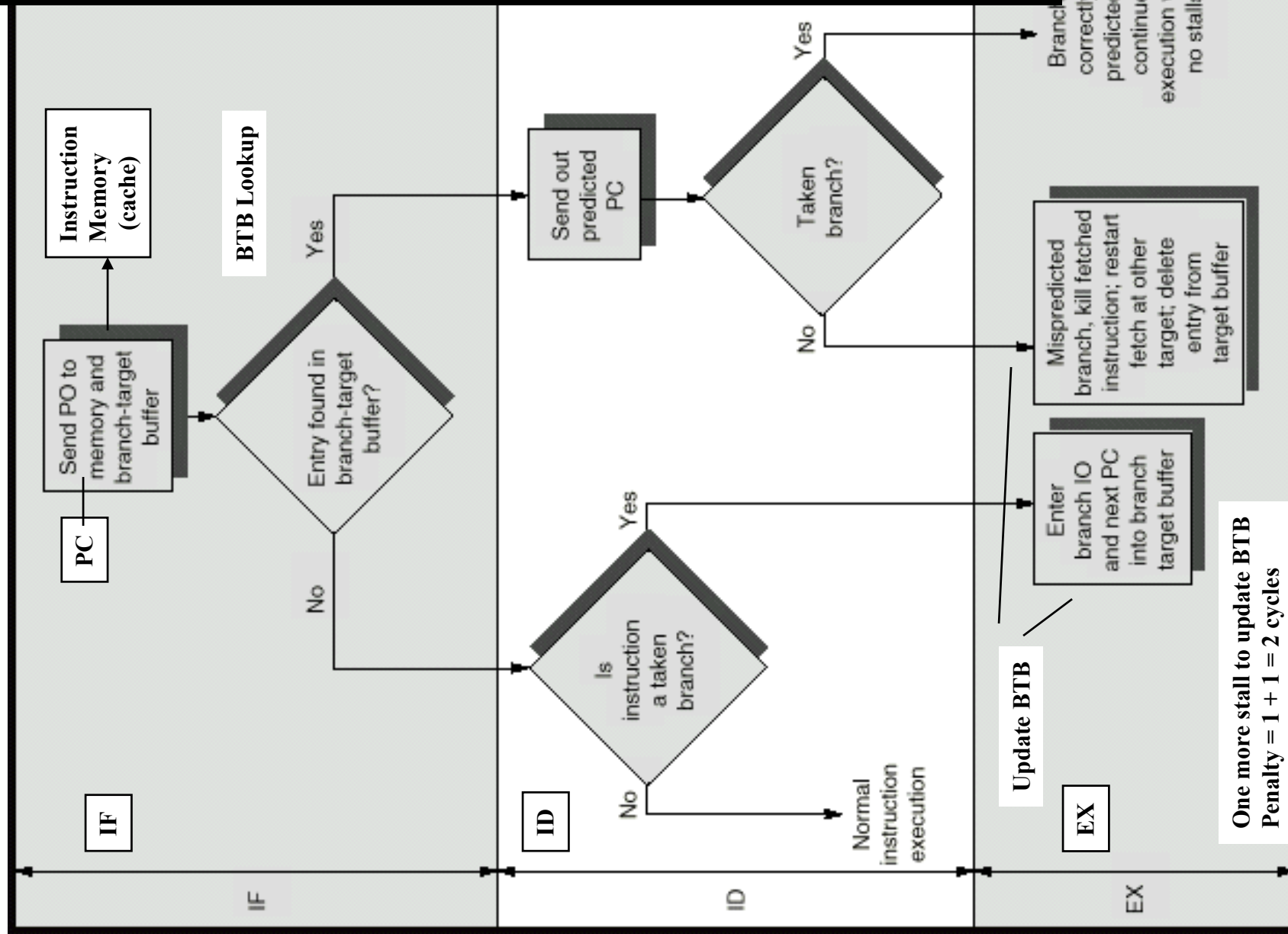
Goal of BTB: Zero stalls for taken branches

Basic Branch Target Buffer (BTB)



Here, branches are assumed to be resolved in ID

BTB Operation



Branch Penalty Cycles Using A Branch-Target Buffer (BTB)

➔ Base Pipeline Taken Branch Penalty = 1 cycle (*i.e. branches resolved in ID*)

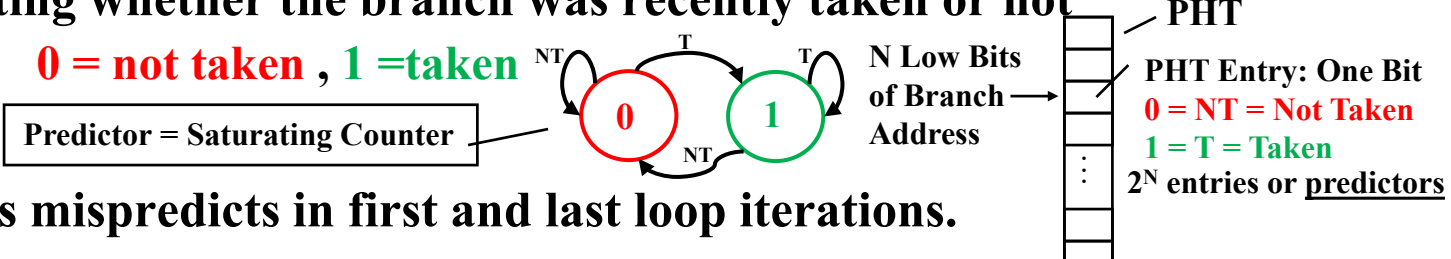
i.e In BTB?		Or Not A Branch		
No	Not Taken	Not Taken	0	
Instruction in buffer	Prediction	Actual branch	Penalty cycles	
Yes	Taken	Taken	0	
Yes	Taken	Not taken	2 = 1 + 1	
No		Taken	2 = 1 + 1	

BTB Goal: Taken Branches with zero stalls

➔ Assuming one more stall cycle to update BTB
Penalty = 1 + 1 = 2 cycles

Penalties for all possible combinations of whether the branch is in the buffer and what it actually does, assuming we store only taken branches in the buffer.

Basic Dynamic Branch Prediction

- ➔ • Simplest method: (One-Level or Non-Correlating) ←
- A branch prediction buffer or Pattern History Table (PHT) indexed by N low address bits of the branch instruction.
 - Each buffer location (or PHT entry or predictor) contains one bit indicating whether the branch was recently taken or not
 - e.g. **0 = not taken**, **1 = taken**
 - Always mispredicts in first and last loop iterations.
- 

- ➔ • To improve prediction accuracy, two-bit prediction is used:
- ➔ – A prediction must miss twice before it is changed. ➔ (Smith Algorithm, 1985)

Why 2-bit Prediction?

- Thus, a branch involved in a loop will be mispredicted only once when encountered the next time as opposed to twice when one bit is used.

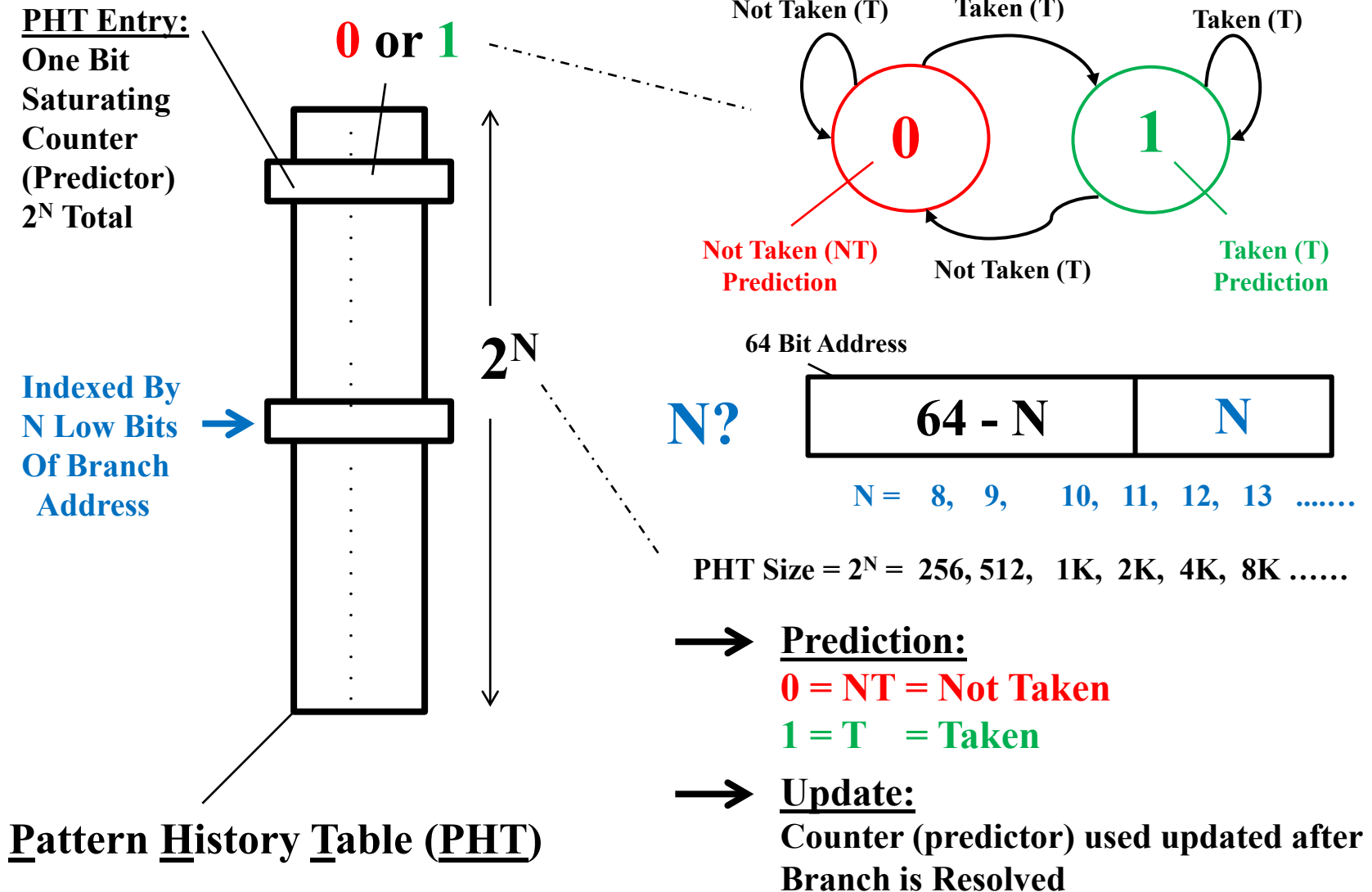
- Two-bit prediction is a specific case of n-bit saturating counter incremented when the branch is taken and decremented when the branch is not taken.

The counter (predictor) used is updated after the branch is resolved

- Two-bit saturating counters (predictors) are usually always used based on observations that the performance of two-bit PHT prediction is comparable to that of n-bit predictors.

Smith Algorithm

One Level (Bimodal) Dynamic Branch Prediction



One-Level Bimodal Branch Predictors

Pattern History Table (PHT)

Most common one-level implementation

Sometimes referred to as
Decode History Table (DHT)
or
Branch History Table (BHT)

2-bit saturating counters (predictors)

Indexed by

N Low Bits of Branch Address

Table (PHT) has 2^N entries
(also called predictors).

2-bit saturating counters

Example:

For $N=12$

Table has $2^N = 2^{12}$ entries
= 4096 = 4k entries

Number of bits needed = $2 \times 4k = 8k$ bits

High bit determines
branch prediction

0 = NT = Not Taken

1 = T = Taken

Prediction Bits

0	0	Not Taken (NT)
0	1	
1	0	Taken (T)
1	1	

When to
update

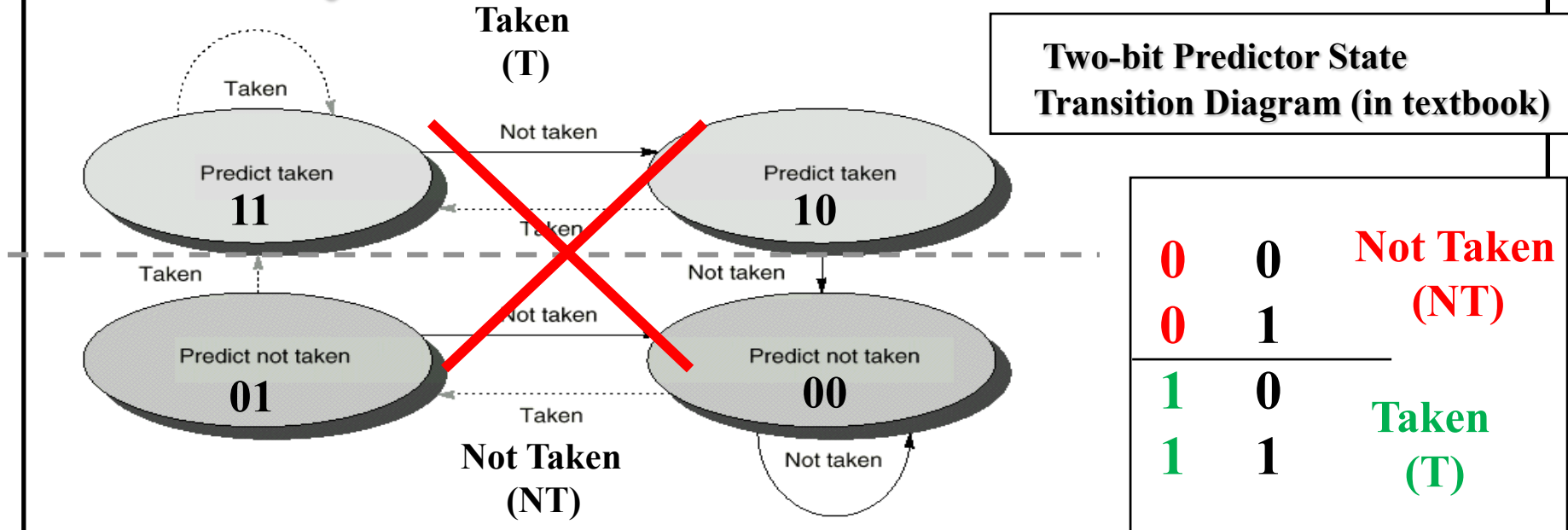
Update counter after branch is resolved:
- Increment counter used if branch is taken
- Decrement counter used if branch is not
taken

What if different branches map to the same predictor (counter)?

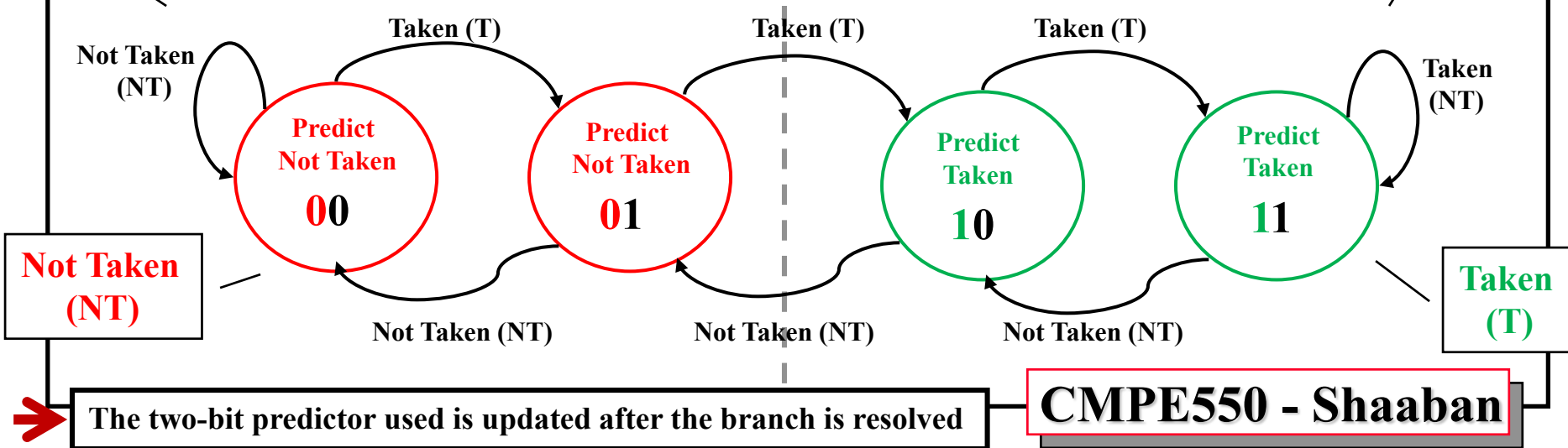
This is called branch address aliasing and leads to interference with current branch prediction by other branches and may lower branch prediction accuracy for programs with aliasing.

CMPE550 - Shaaban

Basic Dynamic Two-Bit Branch Prediction:



Or Two-bit saturating counter predictor state transition diagram (Smith Algorithm):



→ FP has more branches involved in loops than integer code

FP

$N=12$
 $2^N = 4096$

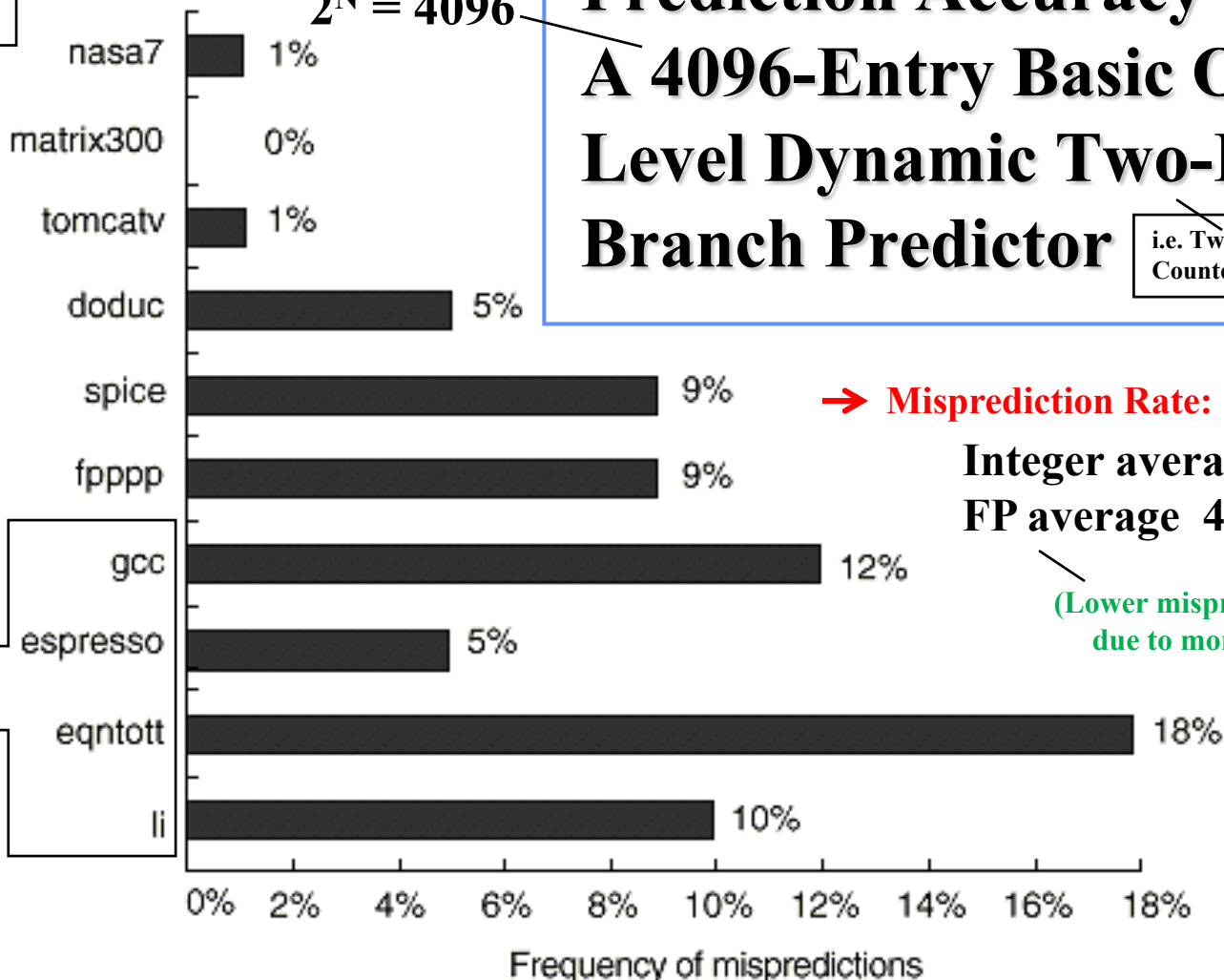
Prediction Accuracy of A 4096-Entry Basic One-Level Dynamic Two-Bit Branch Predictor

i.e. Two-bit Saturating Counters (Smith Algorithm)

SPEC89 benchmarks

Integer

→ Has more branches involved in IF-Then-Else constructs than FP



→ Misprediction Rate:

Integer average 11%

FP average 4%

(Lower misprediction rate due to more loops)

→ Prediction accuracy of a 4096-entry two-bit prediction buffer for the SPEC89 benchmarks.

CMPE550 - Shaaban

From The Analysis of Static Branch Prediction :

MIPS Performance Using Canceling Delay Branches

Benchmark	% conditional branches	% conditional branches with empty slots	% conditional branches that are cancelling	% cancelling branches that are cancelled	% branches with cancelled delay slots	Total % branches with empty or cancelled delay slot
compress	14%	18%	31%	43%	13%	31%
eqntott	24%	24%	50%	24%	12%	36%
espresso	15%	29%	19%	21%	4%	33%
gcc	15%	16%	33%	34%	11%	27%
li	15%	20%	55%	48%	26%	46%
Integer average	17%	21%	38%	34%	13%	35%
doduc	8%	33%	12%	62%	8%	41%
ear	10%	37%	36%	14%	5%	42%
hydro2d	12%	0%	69%	24%	16%	17%
mdljdp2	9%	0%	86%	10%	8%	8%
su2cor	3%	7%	17%	57%	10%	17%
FP average	8%	16%	44%	34%	9%	25%
Overall average	12%	18%	41%	34%	11%	30%

Delayed and cancelling delay branches for MIPS allow branch hazards to be hidden 70% of the time on average for these 10 SPEC benchmarks.

➔ 70% Static Branch Prediction Accuracy

(repeated here from lecture2)



+ ISA Support Needed

CMPE550 - Shaaban

Prediction Accuracy of Basic One-Level Two-Bit Branch Predictors:

SPEC89 benchmarks

$N=12$ $2^N = 4096 = 4K$ $N = \text{All branch address bits}$

■ 4096 entries: 2 bits per entry
□ Unlimited entries: 2 bits per entry

FP

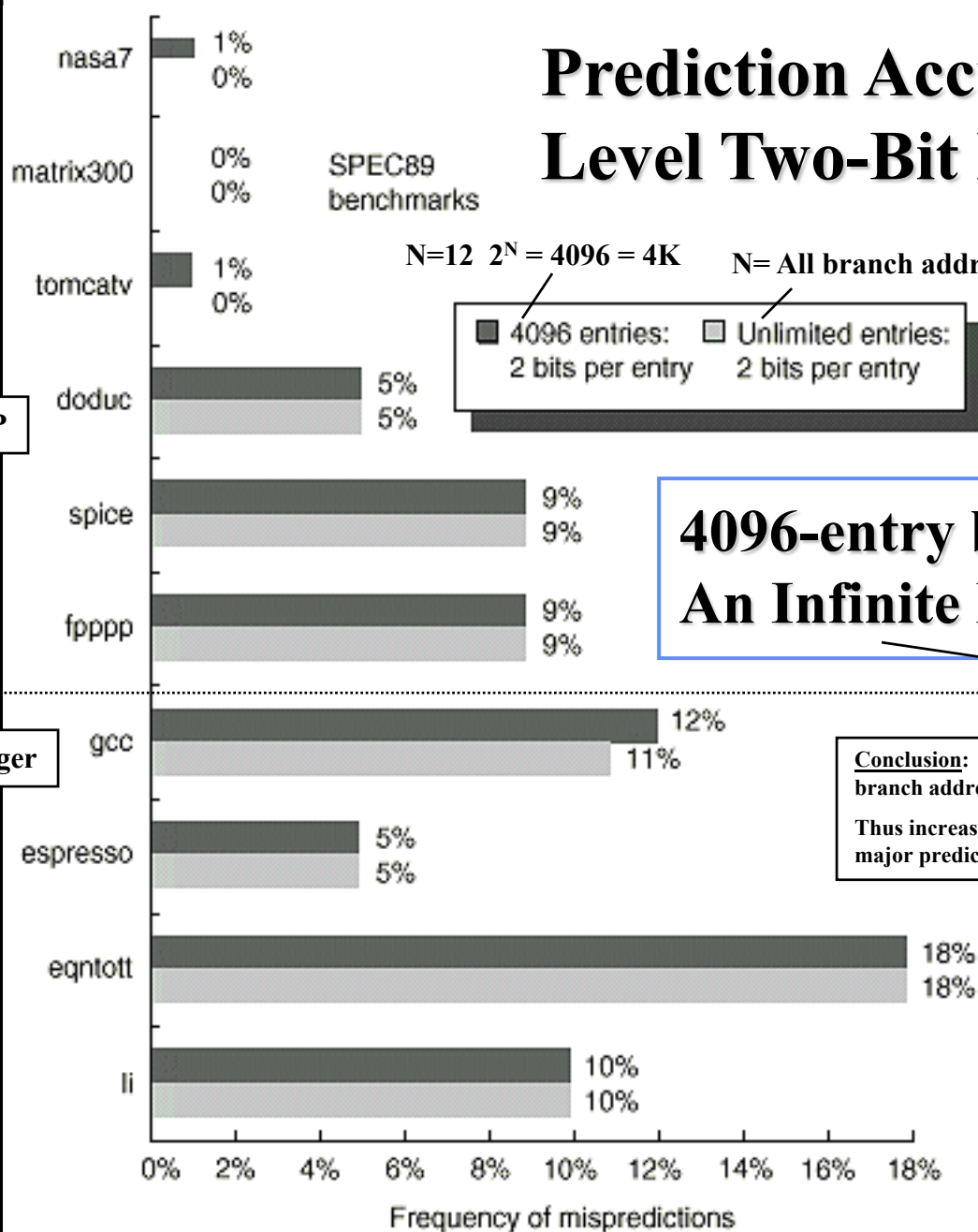
**4096-entry buffer (PHT) Vs.
An Infinite Buffer Under SPEC89**

i.e. All branch address bit are used

Conclusion: SPEC89 programs do not have many branches that suffer from branch address aliasing (interference) when using a 4096-entry PHT.

Thus increasing PHT size (which usually lowers aliasing) did not result in major prediction accuracy improvement.

Integer



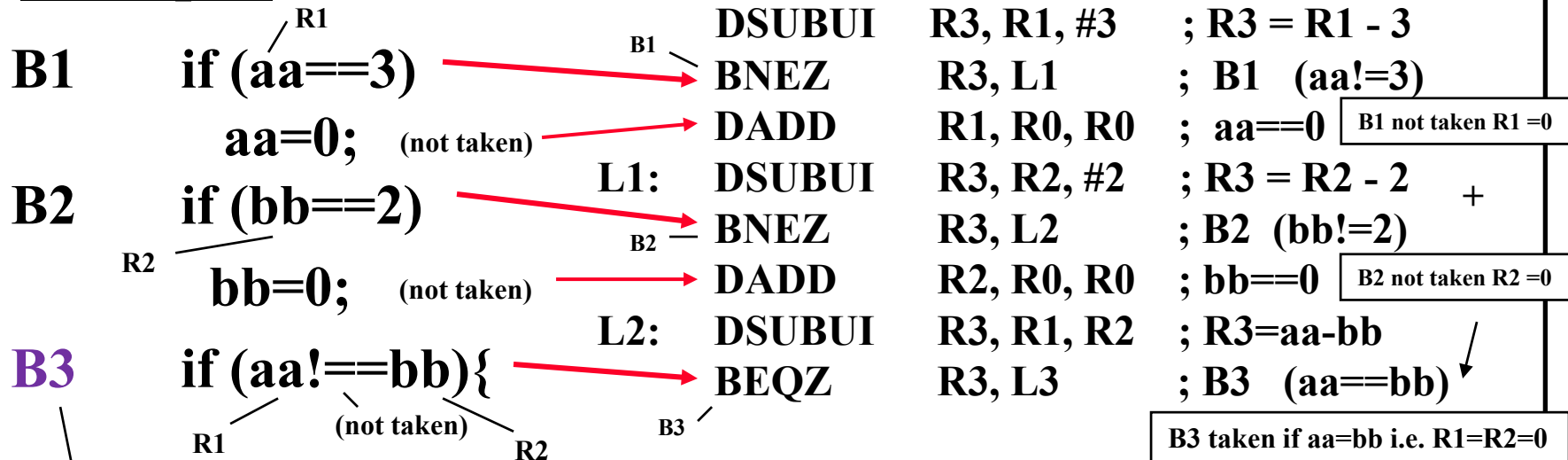
Correlating Branches

➔ Recent branches are possibly correlated: The behavior of recently executed branches affects outcome of current branch.

➔ Occur in branches used to implement if-then-else constructs Which are more common in integer than floating point code

US B3 in the example below

Example:



Branch B3 is correlated with branches B1, B2. If B1, B2 are both not taken, then B3 will be taken. Using only the behavior of one branch cannot detect this behavior.

aa=3 bb=2

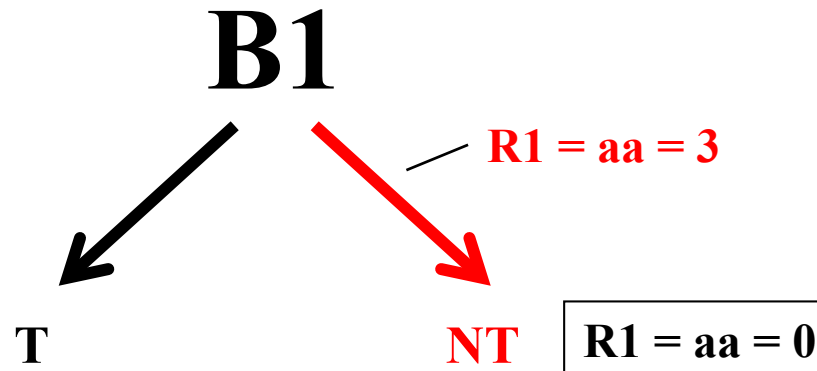
B3 in this case

➔ Both B1 and B2 Not Taken → B3 Taken

CMPE550 - Shaaban

0 = NT = Not Taken

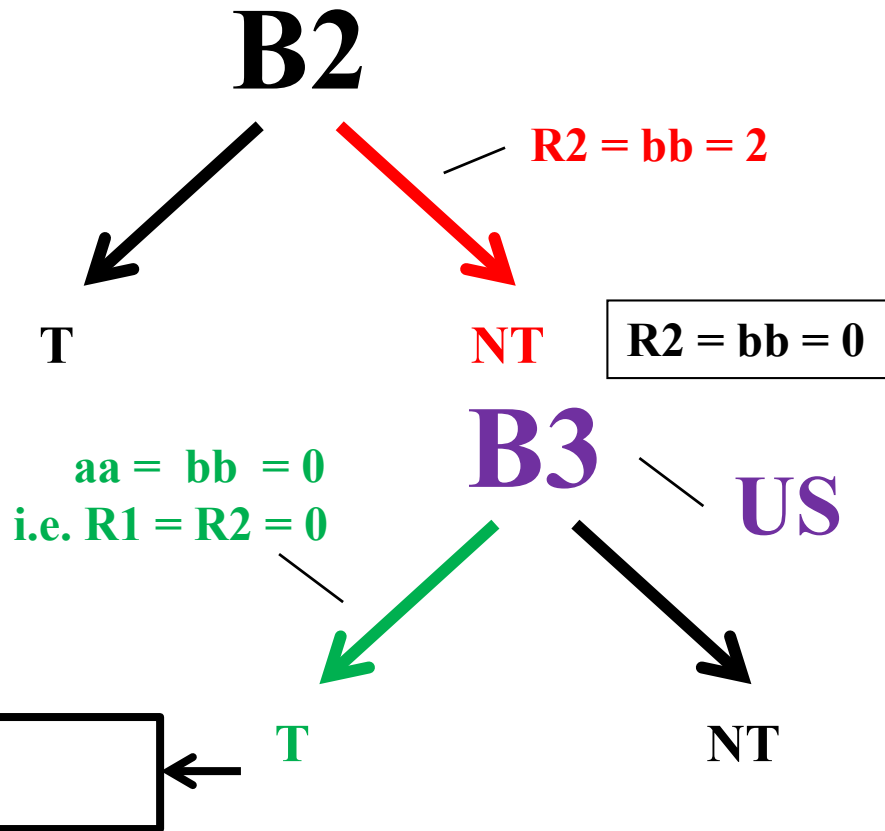
1 = T = Taken



➔ The prediction in one level dynamic branch prediction only depends on the behavior of the current branch being predicted.

➔ How to take into account that recently executed branches (B1, B2 here) may affect outcome of current branch, **B3** here (correlated branches)?

➔ How to record the outcome recently executed branches? **?**



Both B1 and B2 Not Taken → B3 Taken

CMPE550 - Shaaban

0 = NT = Not Taken

1 = T = Taken

B1

$R1 = aa = 3$

T

0

NT

$R1 = aa = 0$

B2

$R2 = bb = 2$

T

0

NT

$R2 = bb = 0$

B3

US

1

T

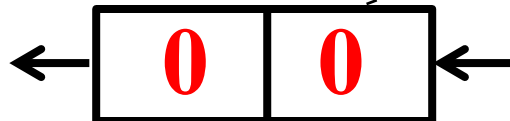
NT

$aa = bb = 0$
i.e. $R1 = R2 = 0$

How to record the outcome recently executed branches? ?

Branch History Register (BHR):

An m bit shift register that records the outcome of m recently executed branches



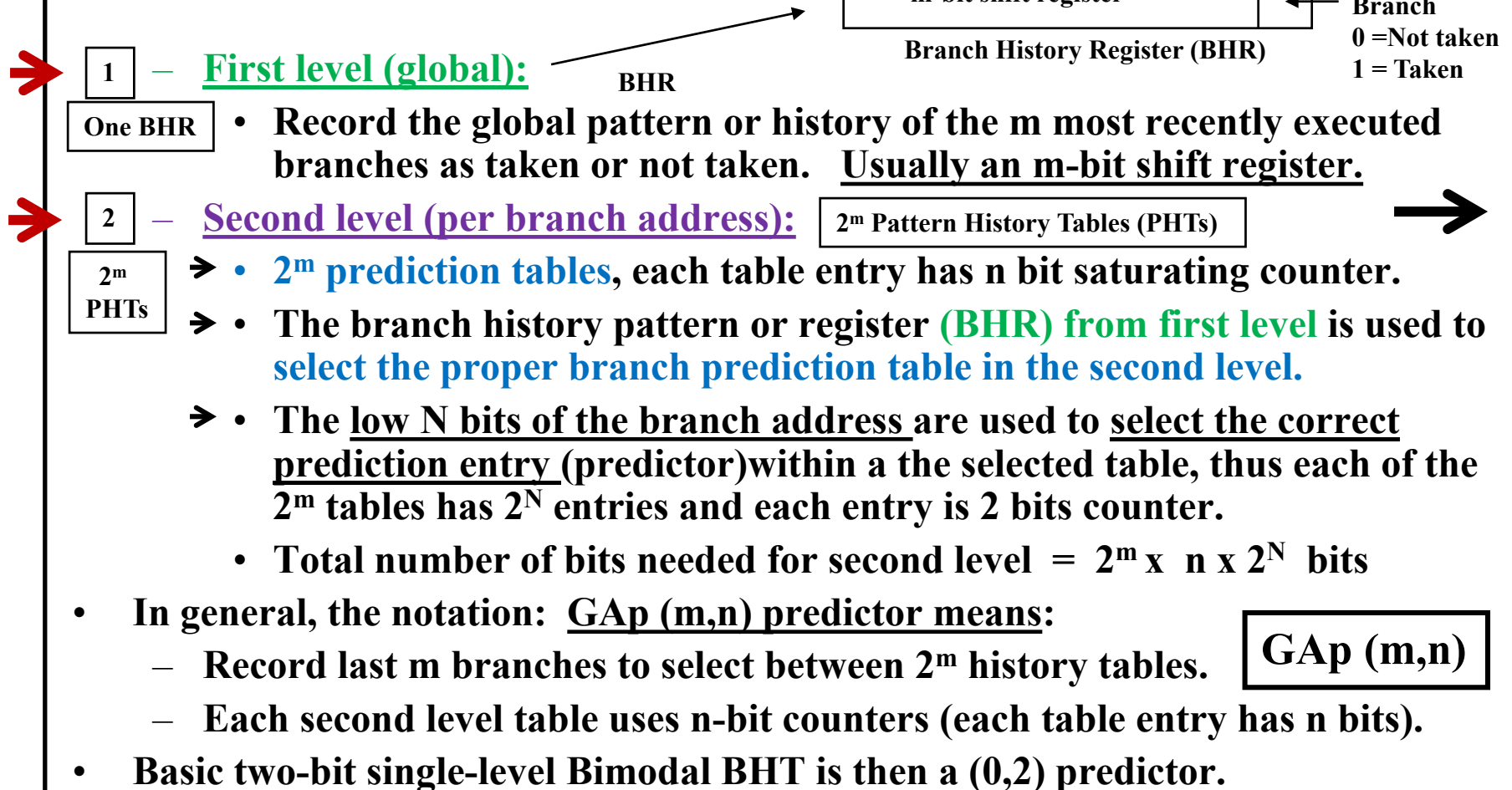
BHR of size $m = 2$ bits

Both B1 and B2 Not Taken → B3 Taken

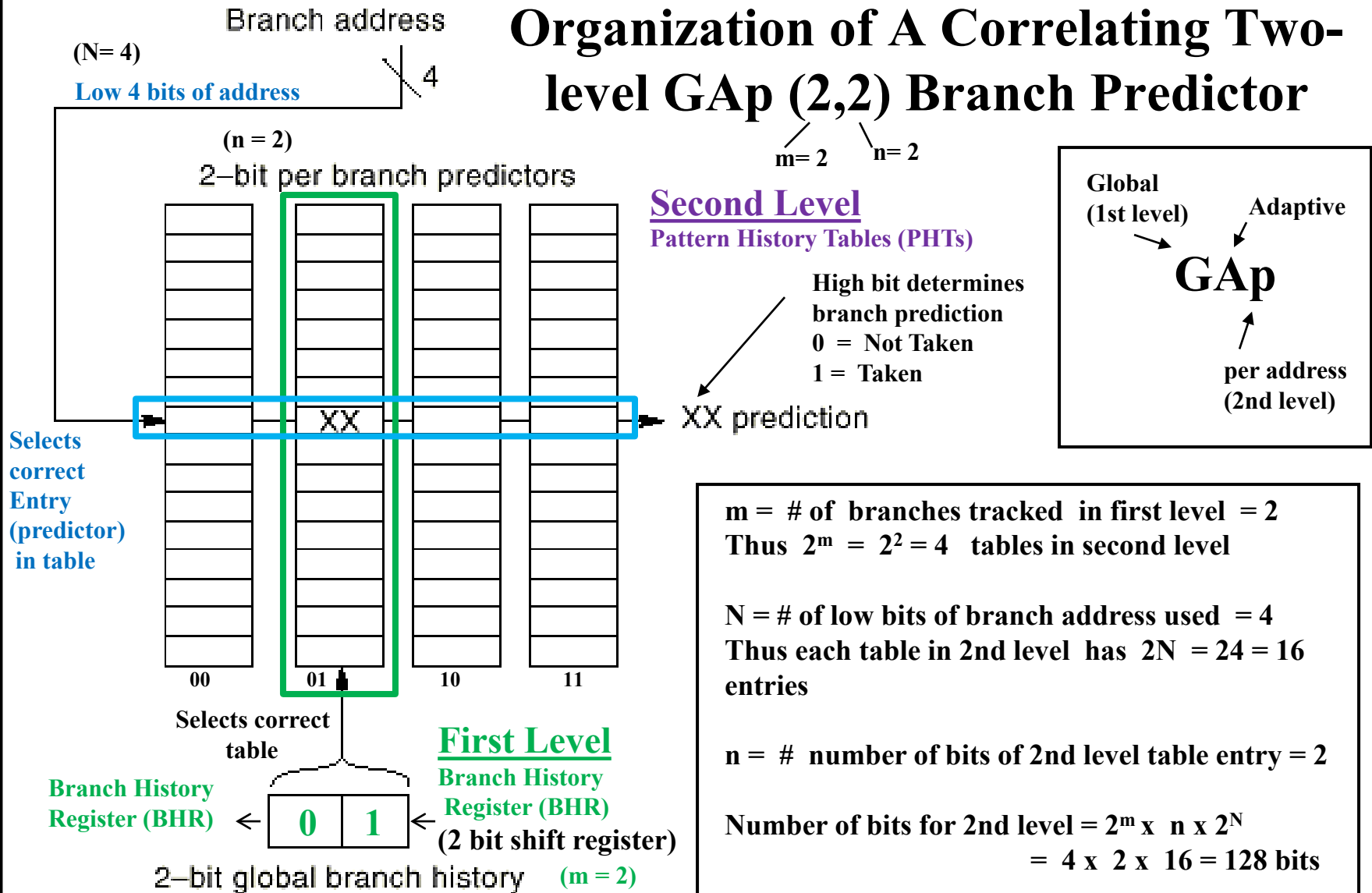
CMPE550 - Shaaban

Correlating Two-Level Dynamic GAP Branch Predictors

- 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.
- Uses two levels of branch history:



Organization of A Correlating Two-level Gap (2,2) Branch Predictor



A (2,2) branch-prediction buffer uses a two-bit global history to choose from among four predictors for each branch address.

Dynamic Branch Prediction: Example

b1 Not Taken ($d = R1 = 0$) $\Rightarrow d = R1 = 1$

if ($d==0$)

Not Taken — $d=1$;

if ($d==1$)

$d = R1 = 1 \Rightarrow$ b2 Not Taken

b1

BNEZ R1, L1 ; branch b1 ($d!=0$)

DADDIU R1, R0, #1 ; $d==0$, so $d=1$

L1:

DADDIU R3, R1, # -1

b2

BNEZ R3, L2 ; branch b2 ($d!=1$)

...

L2:

b1 Taken

b2 Taken

Possible execution sequences for a code fragment.

Initial value of d	$d==0?$	b1		Value of d before b2	$d==1?$	b2
0	Yes	Not taken	\longrightarrow	1	Yes	Not taken
1	No	Taken	\longrightarrow	1	Yes	Not taken
2	No	Taken	\longrightarrow	2	No	Taken

One Level

Behavior of a one-bit predictor initialized to not taken.

d=?	b1 prediction	b1 action	New b1 prediction	b2 prediction	b2 action	New b2 prediction
2	NT	T	T	NT	T	T
0	T	NT	NT	T	NT	NT
2	NT	T	T	NT	T	T
0	T	NT	NT	T	NT	NT

One Level with one-bit table entries (predictors) : NT = 0 = Not Taken
T = 1 = Taken

CMPE550 - Shaaban

Dynamic Branch Prediction: Example (continued)

b1 Not Taken ($d = R1 = 0$) $\Rightarrow d = R1 = 1$

if ($d==0$)

Not Taken — $d=1$;

if ($d==1$)

$d = R1 = 1 \Rightarrow$ b2 Not Taken

b1

BNEZ R1, L1 ; branch b1 ($d!=0$)

DADDIU R1, R0, #1 ; $d==0$, so $d=1$

L1:

DADDIU R3, R1, # -1

b2

BNEZ R3, L2 ; branch b2 ($d!=1$)

...

L2:

b1 Taken

b2 Taken

Combinations and meaning of the taken/not taken prediction bits.

Initial value of d	$d==0?$	b1	Value of d before b2	$d==1?$	b2
0	Yes	Not taken	1	Yes	Not taken
1	No	Taken	1	Yes	Not taken
2	No	Taken	2	No	Taken

The action of the one-bit predictor with one bit of correlation, initialized to not taken/not taken.

d=?	b1 prediction	b1 action	New b1 prediction	b2 prediction	b2 action	New b2 prediction
2	NT/NT	T	T/NT	NT/NT	T	NT/T
0	T/NT	NT	T/NT	NT/T	NT	NT/T
2	T/NT	T	T/NT	NT/T	T	NT/T
0	T/NT	NT	T/NT	NT/T	NT	NT/T

Two level GAp(1,1)

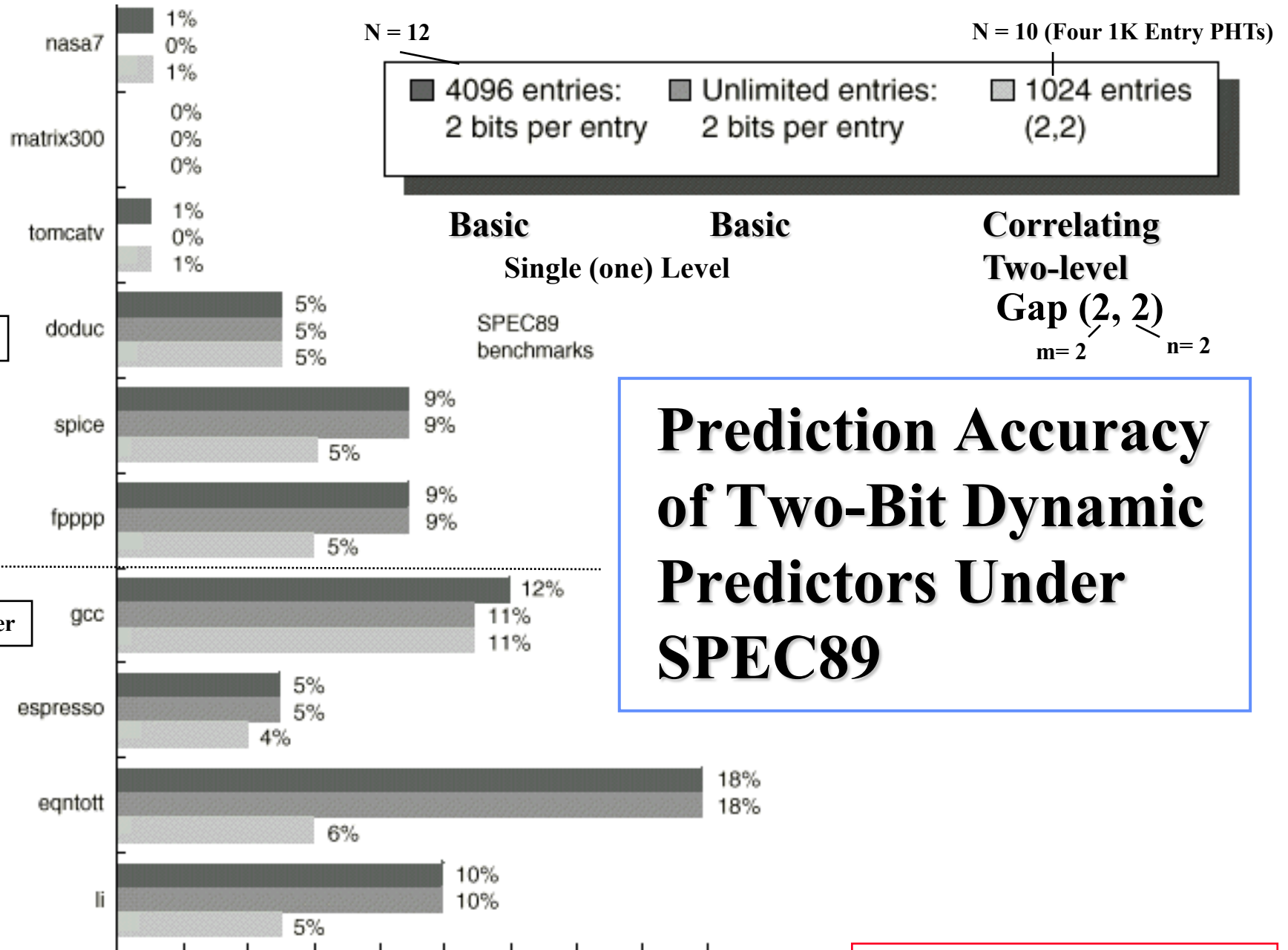
m=1

n=1

CMPE550 - Shaaban

FP

Integer



McFarling's gshare Predictor

gshare = global history with index sharing

- McFarling noted (1993) that using global history information might be less efficient than simply using the address of the branch instruction, especially for small predictors.
- ➔ • He suggests using both global history (BHR) and branch address by **hashing them together**. He proposed using the XOR of global branch history register (BHR) and branch address since he expects that this value has more information than either one of its components. The result is that this mechanism outperforms Gap scheme by a small margin. Bitwise XOR ➔
- This mechanism uses less hardware than Gap, since both branch history (first level) and pattern history (second level) are kept globally.
- The hardware cost for k history bits is $k + 2 \times 2^k$ bits, neglecting costs for logic.

➔ gshare is one of the most widely implemented two level dynamic branch prediction schemes

gshare Predictor

Branch and pattern history are kept globally. History and branch address are XORed and the result is used to index the pattern history table.

Here:
 $m = N = k$ bits

(BHR)

Branch History
Shift Register

Branch Address

$m = N = k$ bits

$N = k$ bits

k bits

First Level:

$m = k$ bits

k bits

XOR (bitwise XOR)

2-bit saturating counters (predictors)

Index the second level

Second Level:

(One PHT)

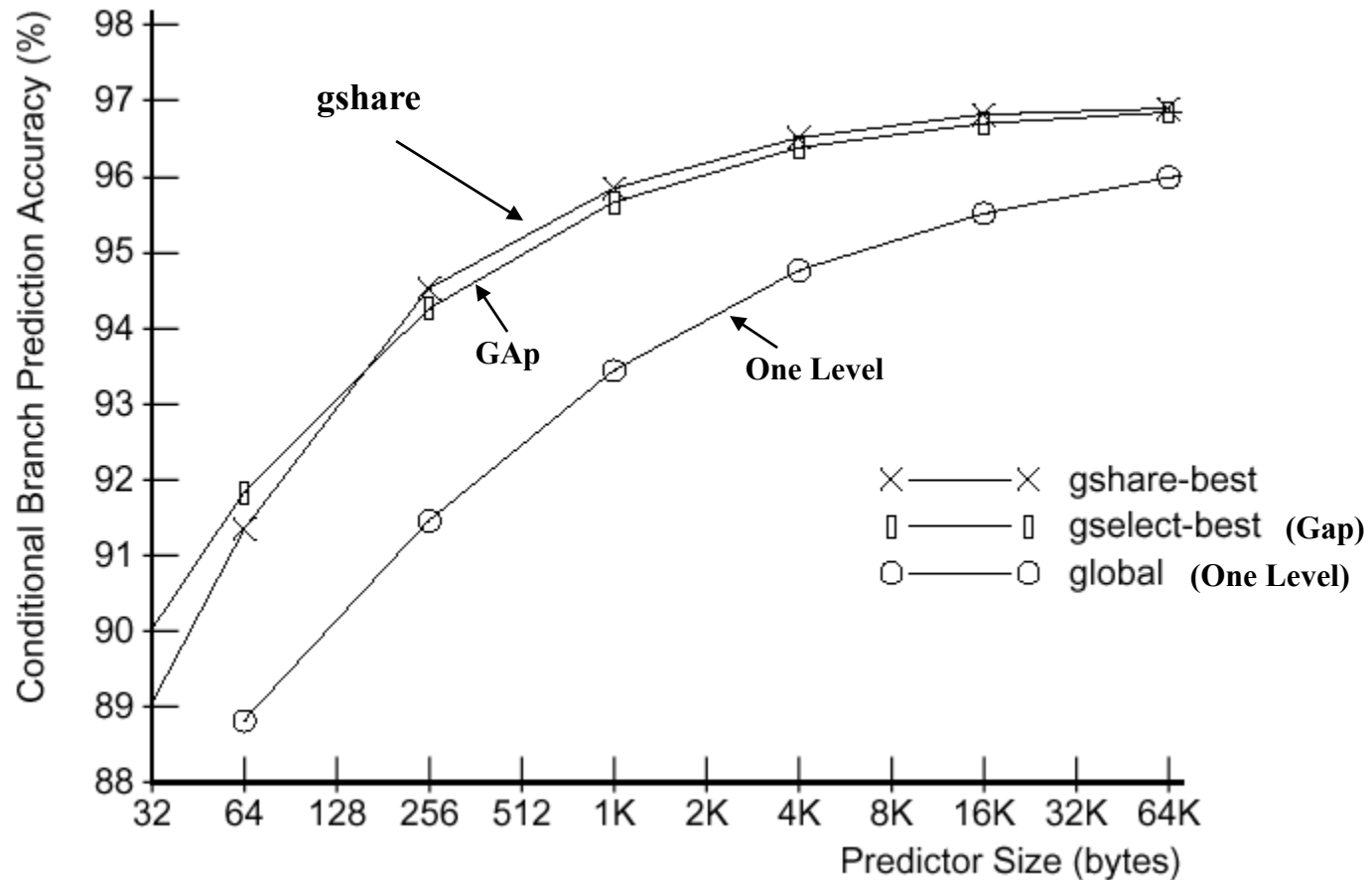
One Pattern History Table (PHT) with 2^k entries (predictors)

Prediction

gshare = global history with index sharing

CMPE550 - Shaaban

gshare Performance



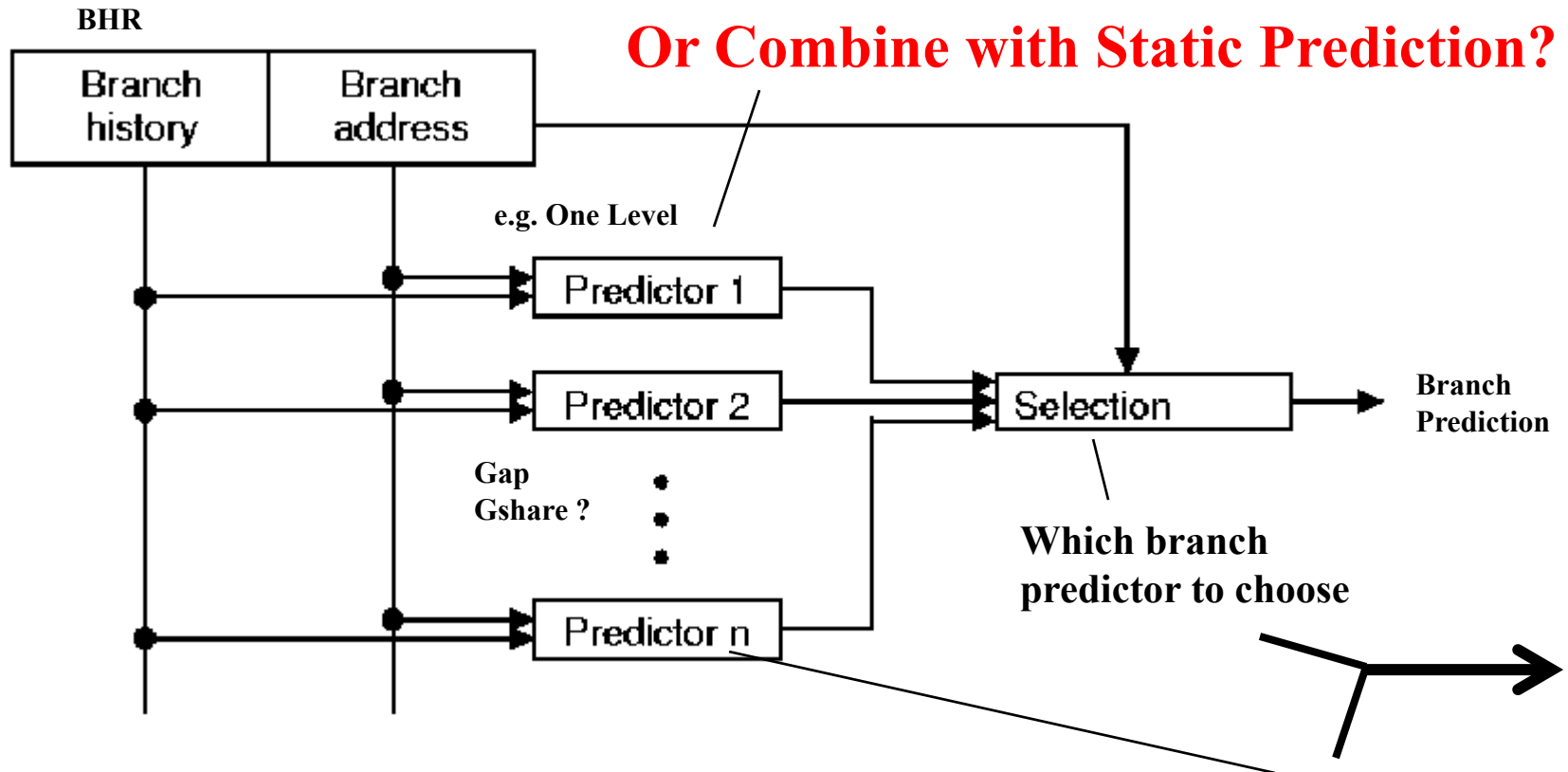
Hybrid Dynamic Branch Predictors

(Also known as tournament or combined predictors)

- ➔ • Hybrid predictors are simply combinations of two (most common) or more branch prediction mechanisms.
- This approach takes into account that different mechanisms may perform best for different branch scenarios.
- ➔ • McFarling presented (1993) a number of different combinations of two branch prediction mechanisms.
- ➔ • He proposed to use an additional 2-bit counter selector array which serves to select the appropriate predictor for each branch.
- One predictor is chosen for the higher two counts, the second one for the lower two counts. The selector array counter used is updated as follows:
 1. If the first predictor is wrong and the second one is right the selector counter used counter is decremented,
 2. If the first one is right and the second one is wrong, the selector counter used is incremented.
 3. No changes are carried out to selector counter used if both predictors are correct or wrong.

Predictor
Selector
Array
Counter
Update

A Generic Hybrid Predictor



➔ Usually only two predictors are used (i.e. $n=2$)
e.g. As in Alpha, IBM POWER 4 - 9 ...

McFarling's Hybrid Predictor Structure

The hybrid predictor contains an additional counter array (selector array) with 2-bit up/down saturating counters. Which serves to select the best predictor to use. Each counter in the selector array keeps track of which predictor is more accurate for the branches that share that counter. Specifically, using the notation P1c and P2c to denote whether predictors P1 and P2 are correct respectively, the selector counter is incremented or decremented by P1c-P2c as shown.

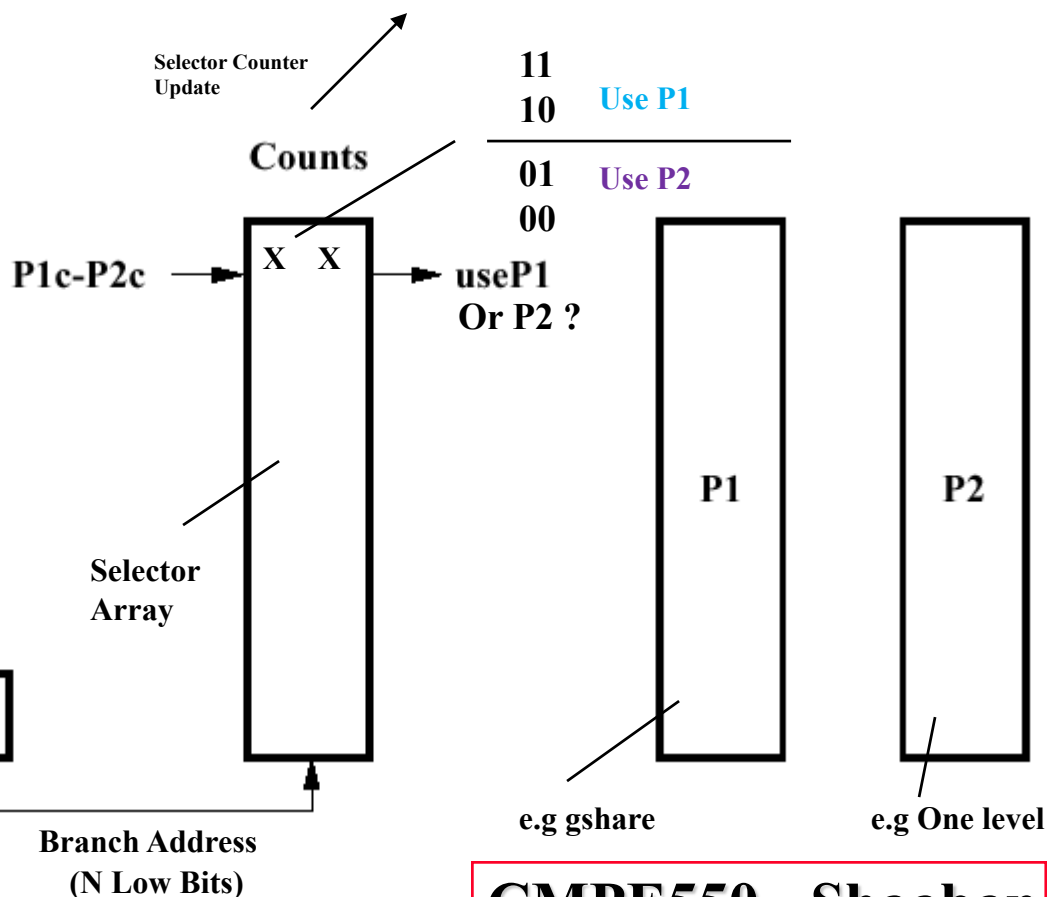
Both wrong

P2 correct

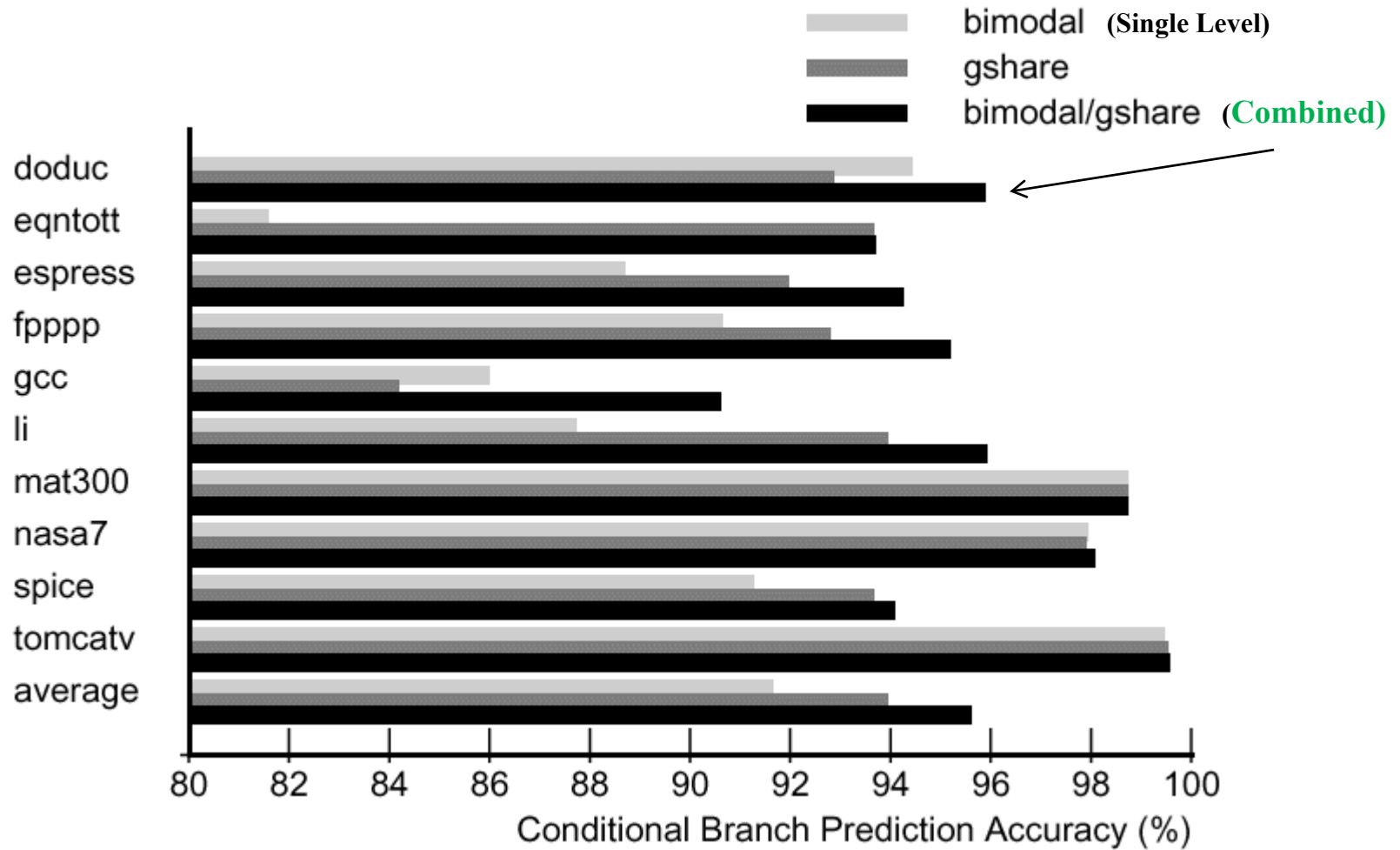
P1 correct

Both correct

P1c	P2c	P1c-P2c	
0	0	0	(no change)
0	1	-1	(decrement counter)
1	0	1	(increment counter)
1	1	0	(no change)



MCFarling's Hybrid Predictor Performance by Benchmark



Processor Branch Prediction Examples

<u>Processor</u>	<u>Released</u>	<u>Accuracy</u>	<u>Prediction Mechanism</u>
Cyrix 6x86	early '96	ca. 85%	PHT associated with BTB
Cyrix 6x86MX	May '97	ca. 90%	PHT associated with BTB
AMD K5	mid '94	80%	PHT associated with I-cache
AMD K6	early '97	95%	2-level adaptive associated with BTIC and ALU
Intel Pentium	late '93	78%	PHT associated with BTB
Intel P6	mid '96	90%	2 level adaptive with BTB
S+D PowerPC750	mid '97	90%	PHT associated with BTIC
MC68060	mid '94	90%	PHT associated with BTIC
DEC Alpha	early '97	95%	Hybrid 2-level adaptive associated with I-cache
S+D HP PA8000	early '96	80%	PHT associated with BTB
S+D SUN UltraSparc	mid '95	88%int 94%FP	PHT associated with I-cache

S+D : Uses both static (ISA supported) and dynamic branch prediction

PHT = One Level

CMPE550 - Shaaban