## COMP4211 05s1 Seminar 4:
## Branch Prediction

Slides due to
David A. Patterson, 2001

---

## Review Tomasulo

- Reservations stations: *implicit register renaming* to larger set of registers + buffering source operands
  - Prevents registers as bottleneck
  - Avoids WAR, WAW hazards of Scoreboard
  - Allows loop unrolling in HW
- Not limited to basic blocks (integer units gets ahead, beyond branches)
- Today, helps cache misses as well
  - Don't stall for L1 Data cache miss (insufficient ILP for L2 miss?)
- Lasting Contributions
  - Dynamic scheduling
  - Register renaming
  - Load/store disambiguation
- 360/91 descendants are Pentium III; PowerPC 604; MIPS R10000; HP-PA 8000; Alpha 21264

---

## Tomasulo Algorithm and Branch Prediction

- 360/91 predicted branches, but did not speculate: pipeline stopped until the branch was resolved
  - No speculation; only instructions that can complete
- Speculation with Reorder Buffer allows execution past branch, and then discard if branch fails
  - just need to hold instructions in buffer until branch can commit

---

## Case for Branch Prediction when Issue N instructions per clock cycle

1. Branches will arrive up to $n$ times faster in an $n$-issue processor
2. Amdahl's Law => relative impact of the control stalls will be larger with the lower potential CPI in an $n$-issue processor

## 7 Branch Prediction Schemes

1. 1-bit Branch-Prediction Buffer
2. 2-bit Branch-Prediction Buffer
3. Correlating Branch Prediction Buffer
4. Tournament Branch Predictor
5. Branch Target Buffer
6. Integrated Instruction Fetch Units
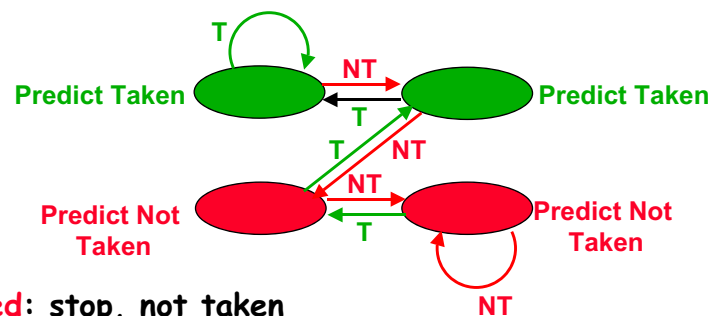7. Return Address Predictors

---

## Dynamic Branch Prediction

- Performance = $f$(accuracy, cost of misprediction)
- Branch History Table: Lower bits of PC address index table of 1-bit values
  - Says whether or not branch taken last time
  - No address check (saves HW, but may not be right branch)
- Problem: in a loop, 1-bit BHT will cause 2 mispredictions (avg is 9 iterations before exit):
  - End of loop case, when it exits instead of looping as before
  - First time through loop on *next* time through code, when it predicts *exit* instead of looping
  - Only 80% accuracy even if loop 90% of the time

---

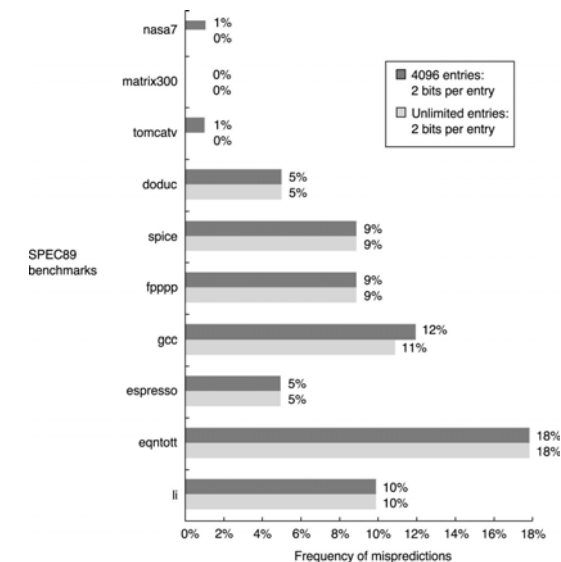## Dynamic Branch Prediction
### (Jim Smith, 1981)

- Solution: 2-bit scheme where change prediction only if get misprediction *twice:* (Figure 3.7, p. 198)



- **Red**: stop, not taken
- **Green**: go, taken
- Adds *hysteresis* to decision making process

---

## Prediction accuracy: 4K-entry 2-bit table vs infinite table size



| SPEC89 benchmarks | 4096 entries: 2 bits per entry | Unlimited entries: 2 bits per entry |
|---|---|---|
| nasa7 | 1% | 0% |
| matrix300 | 0% | 0% |
| tomcatv | 1% | 0% |
| doduc | 5% | 5% |
| spice | 9% | 9% |
| fpppp | 9% | 9% |
| gcc | 12% | 11% |
| espresso | 5% | 5% |
| eqntott | 18% | 18% |
| li | 10% | 10% |

Frequency of mispredictions

# Correlating Predictors

- 2-bit prediction uses a small amount of (hopefully) local information to predict behaviour

- Sometimes behaviour is correlated, and we can do better by keeping track of direction of related branches, for example consider the following code:

```
if (d==0)
    d = 1;
if (d==1) {
```

- If the first branch is not taken, neither is the second. Predictors that use the behaviour of other branches to make a prediction are called *correlating predictors* or *two-level predictors*
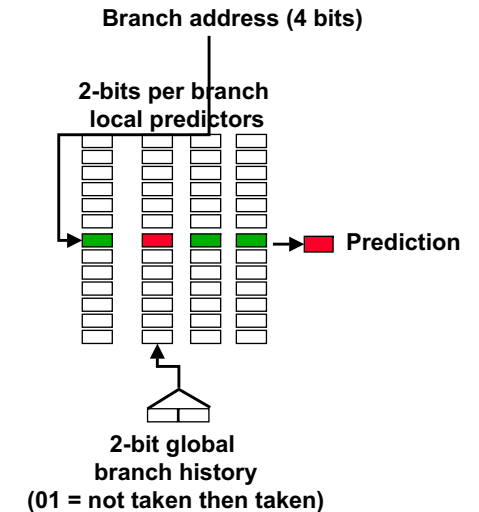
---

# Correlating Branches

Idea: taken/not taken of recently executed branches is related to behavior of next branch (as well as the history of that branch behavior)

- Then behavior of recent branches selects between, say, 4 predictions of next branch, updating just that prediction
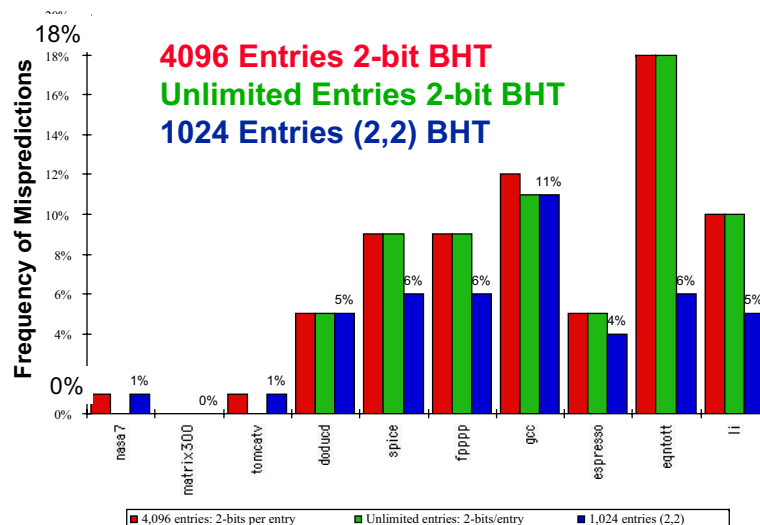
- (2,2) predictor: 2-bit global, 2-bit local



Branch address (4 bits)

2-bits per branch local predictors

Prediction

2-bit global branch history
(01 = not taken then taken)

---

# Accuracy of Different Schemes
(Figure 3.15, p. 206)

**4096 Entries 2-bit BHT**
**Unlimited Entries 2-bit BHT**
**1024 Entries (2,2) BHT**



Frequency of Mispredictions

- 4,096 entries: 2-bits per entry
- Unlimited entries: 2-bits/entry
- 1,024 entries (2,2)

---

# Re-evaluating Correlation

- Several of the SPEC benchmarks have less than a dozen branches responsible for 90% of taken branches:

| program | branch % | static | # = 90% |
|---------|----------|--------|---------|
| compress | 14% | 236 | 13 |
| eqntott | 25% | 494 | 5 |
| gcc | 15% | 9531 | 2020 |
| mpeg | 10% | 5598 | 532 |
| real gcc | 13% | 17361 | 3214 |

- Real programs + OS more like gcc

- Small benefits beyond benchmarks for correlation? problems with branch aliases?

# BHT Accuracy

- **Mispredict because either:**
  - Wrong guess for that branch
  - Got branch history of wrong branch when index the table
- **4096 entry table  programs vary from 1% misprediction (nasa7, tomcatv) to 18% (eqntott), with spice at 9% and gcc at 12%**
- **For SPEC92,
4096 about as good as infinite table**

# Tournament Predictors

- **Motivation for correlating branch predictors is 2-bit predictor failed on important branches; by adding global information, performance improved**
- **Tournament predictors: use 2 predictors, 1 based on global information and 1 based on local information, and combine with a selector**
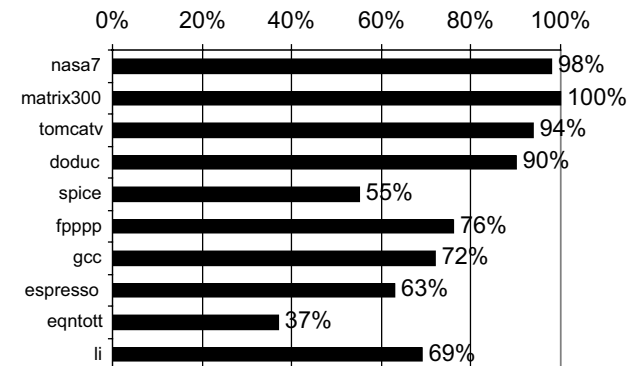- **Hopes to select right predictor for right branch**

# Tournament Predictor in Alpha 21264

- **4K 2-bit counters to choose from among a global predictor and a local predictor**
- **Global predictor** also has 4K entries and is indexed by the history of the last 12 branches; each entry in the global predictor is a standard 2-bit predictor
  - 12-bit pattern: ith bit 0 => ith prior branch not taken;
ith bit 1 => ith prior branch taken;
- **Local predictor** consists of a 2-level predictor:
  - **Top level** a local history table consisting of 1024 10-bit entries; each 10-bit entry corresponds to the most recent 10 branch outcomes for the entry. 10-bit history allows patterns 10 branches to be discovered and predicted.
  - **Next level** Selected entry from the local history table is used to index a table of 1K entries consisting a 3-bit saturating counters, which provide the local prediction
- **Total size: 4K*2 + 4K*2 + 1K*10 + 1K*3 = 29K bits!**

  **(~180,000 transistors)**

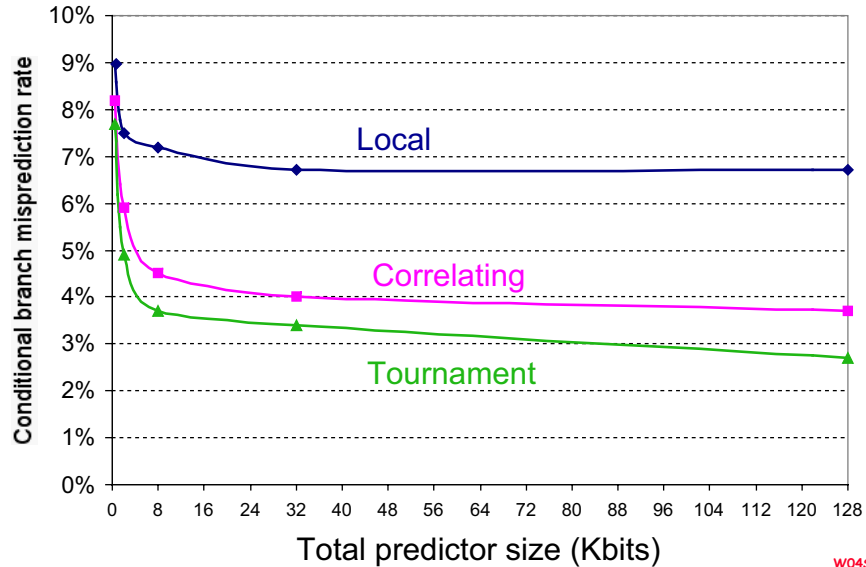# % of predictions from local predictor in Tournament Prediction Scheme

| | |
|---|---|
| nasa7 | 98% |
| matrix300 | 100% |
| tomcatv | 94% |
| doduc | 90% |
| spice | 55% |
| fpppp | 76% |
| gcc | 72% |
| espresso | 63% |
| eqntott | 37% |
| li | 69% |

## Accuracy v. Size (SPEC89)



Y-axis: Conditional branch misprediction rate (0% to 10%)

X-axis: Total predictor size (Kbits) — 0, 8, 16, 24, 32, 40, 48, 56, 64, 72, 80, 88, 96, 104, 112, 120, 128

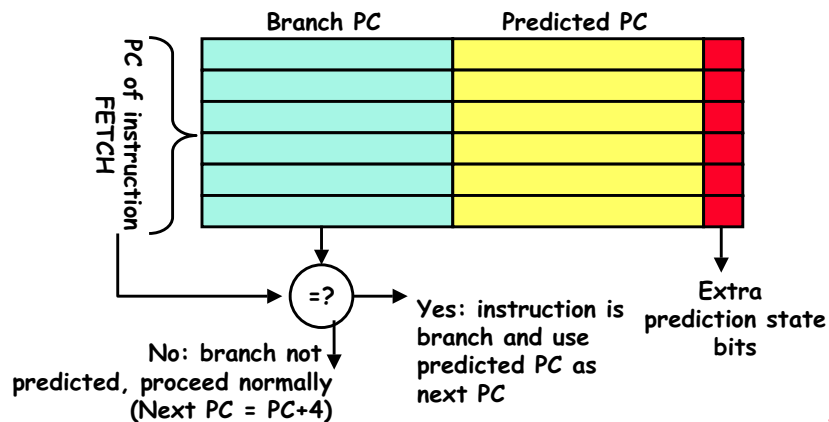Local
Correlating
Tournament

## Pitfall: Sometimes bigger and dumber is better

- **21264 uses tournament predictor (29 Kbits)**
- **Earlier 21164 uses a simple 2-bit predictor with 2K entries (or a total of 4 Kbits)**
- **SPEC95 benchmarks, 21264 outperforms**
  - **21264 avg. 11.5 mispredictions per 1000 instructions**
  - **21164 avg. 16.5 mispredictions per 1000 instructions**
- **Reversed for transaction processing (TP) !**
  - **21264 avg. 17 mispredictions per 1000 instructions**
  - **21164 avg. 15 mispredictions per 1000 instructions**
- **TP code much larger & 21164 hold 2X branch predictions based on local behavior (2K vs. 1K local predictor in the 21264)**
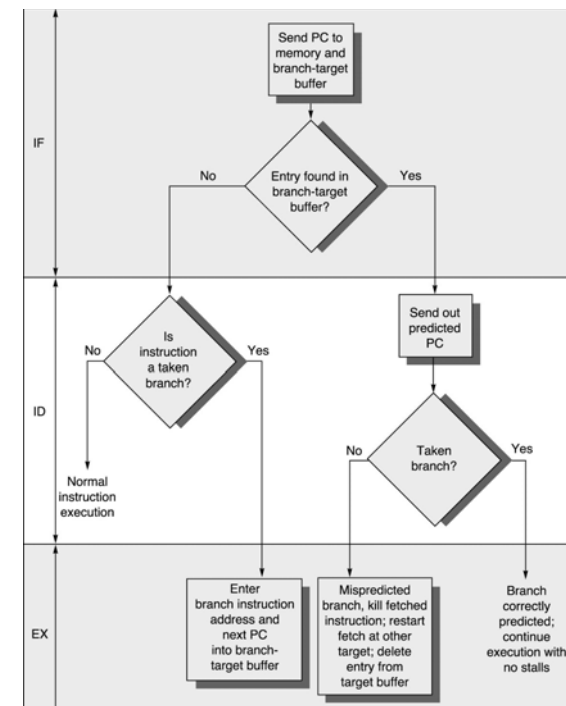
## Need Address at Same Time as Prediction

- **Branch Target Buffer (BTB): Address of branch index to get prediction AND branch address (if taken)**
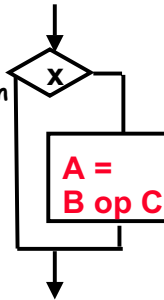  - **Note: must check for branch match now, since can't use wrong branch address (Figure 3.19, p. 210)**



Branch PC    Predicted PC

PC of instruction FETCH

=?

No: branch not predicted, proceed normally (Next PC = PC+4)

Yes: instruction is branch and use predicted PC as next PC

Extra prediction state bits

IF — Send PC to memory and branch-target buffer

Entry found in branch-target buffer? No / Yes

ID — Is instruction a taken branch? No / Yes

Send out predicted PC

Normal instruction execution

Taken branch? No / Yes

EX — Enter branch instruction address and next PC into branch-target buffer

Mispredicted branch, kill fetched instruction; restart fetch at other target; delete entry from target buffer

Branch correctly predicted; continue execution with no stalls

# Predicated Execution

- **Avoid branch prediction by turning branches into conditionally executed instructions:**

  **if (x) then A = B op C else NOP**
  - If false, then neither store result nor cause exception
  - Expanded ISA of Alpha, MIPS, PowerPC, SPARC have conditional move; PA-RISC can annul any following instr.
  - IA-64: 64 1-bit condition fields selected so conditional execution of any instruction
  - This transformation is called "if-conversion"
- **Drawbacks to conditional instructions**
  - Still takes a clock even if "annulled"
  - Stall if condition evaluated late
  - Complex conditions reduce effectiveness; condition becomes known late in pipeline

A =
B op C

# Special Case Return Addresses

- **Register Indirect branch hard to predict address**
- **SPEC89 85% such branches for procedure return**
- **Since stack discipline for procedures, save return address in small buffer that acts like a stack: 8 to 16 entries has small miss rate**

# Dynamic Branch Prediction Summary

- **Prediction becoming important part of scalar execution**
- **Branch History Table: 2 bits for loop accuracy**
- **Correlation: Recently executed branches correlated with next branch.**
  - Either different branches
  - Or different executions of same branches
- **Tournament Predictor: more resources to competitive solutions and pick between them**
- **Branch Target Buffer: include branch address & prediction**
- **Predicated Execution can reduce number of branches, number of mispredicted branches**
- **Return address stack for prediction of indirect jump**