

Advanced Branch Prediction

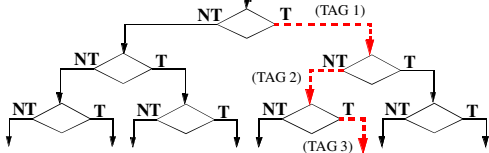
Prof. Mikko H. Lipasti
University of Wisconsin-Madison

Lecture notes based on notes by John P. Shen
Updated by Mikko Lipasti

Advanced Branch Prediction

- Control Flow Speculation
 - Branch Speculation
 - Mis-speculation Recovery
- Branch Direction Prediction
 - Static Prediction
 - Dynamic Prediction
 - Hybrid Prediction
- Branch Target Prediction
- High-bandwidth Fetch
- High-Frequency Fetch

Branch Speculation

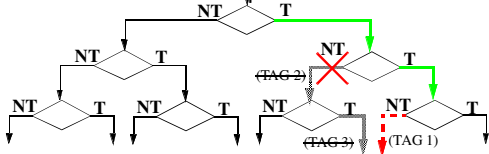


- **Leading Speculation**
 - Typically done during the Fetch stage
 - Based on potential branch instruction(s) in the current fetch group
- **Trailing Confirmation**
 - Typically done during the Branch Execute stage
 - Based on the next Branch instruction to finish execution

Branch Speculation

- **Leading Speculation**
 1. Tag speculative instructions
 2. Advance branch and following instructions
 3. Buffer addresses of speculated branch instructions
- **Trailing Confirmation**
 1. When branch resolves, remove/deallocate speculation tag
 2. Permit completion of branch and following instructions

Branch Speculation



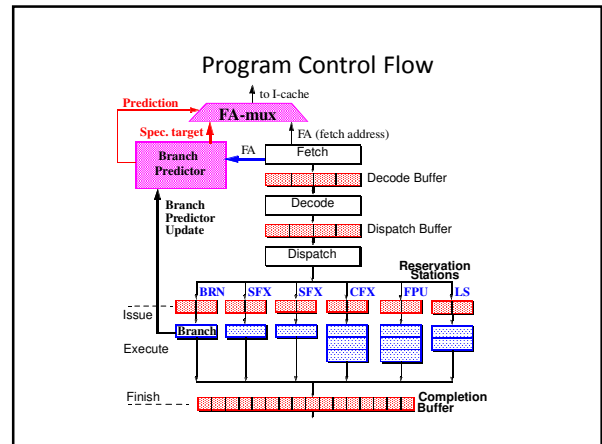
- **Start new correct path**
 - Must remember the alternate (non-predicted) path
- **Eliminate incorrect path**
 - Must ensure that the mis-speculated instructions produce no side effects

Mis-speculation Recovery

- **Start new correct path**
 1. Update PC with computed branch target (if predicted NT)
 2. Update PC with sequential instruction address (if predicted T)
 3. Can begin speculation again at next branch
- **Eliminate incorrect path**
 1. Use tag(s) to deallocate ROB entries occupied by speculative instructions
 2. Invalidate all instructions in the decode and dispatch buffers, as well as those in reservation stations

Tracking Instructions

- Assign branch tags
 - Allocated in circular order
 - Instruction carries this tag throughout processor
- Track instruction groups
 - Instructions managed in groups, max. one branch per group
 - ROB structured as groups
 - Leads to some inefficiency
 - Simpler tracking of speculative instructions



Static Branch Prediction

- Single-direction
 - Always not-taken: Intel i486
- Backwards Taken/Forward Not Taken
 - Loop-closing branches
 - Used as backup in Pentium Pro, II, III, 4
- Heuristic-based:


```
void * p = malloc (numBytes);
if (p == NULL)
    errorHandlerFunction( );
```

Static Branch Prediction

Heuristic Name	Description
Loop Branch	If the branch target is back to the head of a loop, predict taken.
Pointer	If a branch compares a pointer with NULL, or if two pointers are compared, predict in the direction that corresponds to the pointer being not NULL, or the two pointers not being equal.
Opcode	If a branch is testing that an integer is less than zero, less than or equal to zero, or equal to a constant, predict in the direction that corresponds to the test evaluating to false.
Guard	If the operand of the branch instruction is a register that gets used before being redefined in the successor block, predict that the branch goes to the successor block.
Loop Exit	If a branch occurs inside a loop, and neither of the targets is the loop head, then predict that the branch does not go to the successor that is the loop exit.
Loop Header	Predict that the successor block of a branch that is a loop header or a loop pre-header is taken.
Call	If a successor block contains a subroutine call, predict that the branch goes to that successor block.
Store	If a successor block contains a store instruction, predict that the branch does not go to that successor block.
Return	If a successor block contains a return from subroutine instruction, predict that the branch does not go to that successor block.

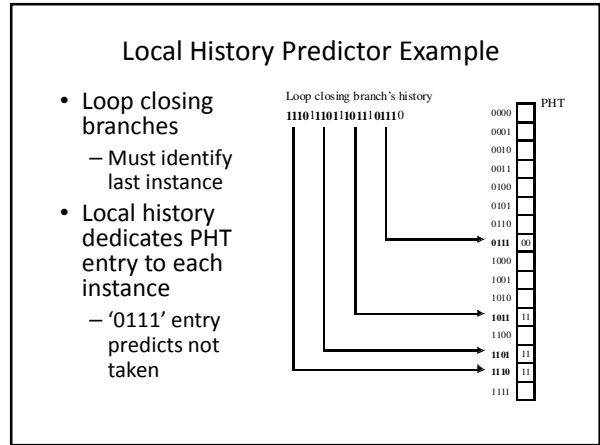
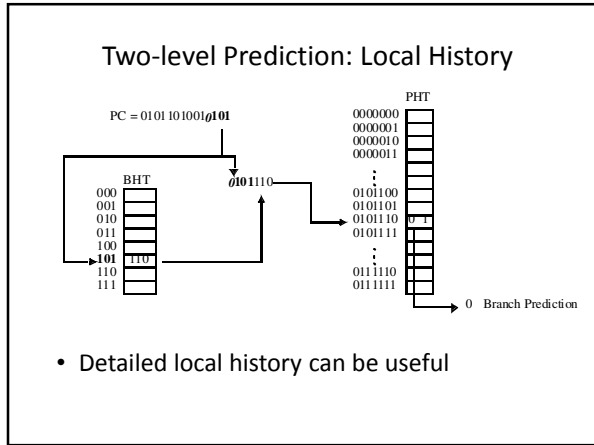
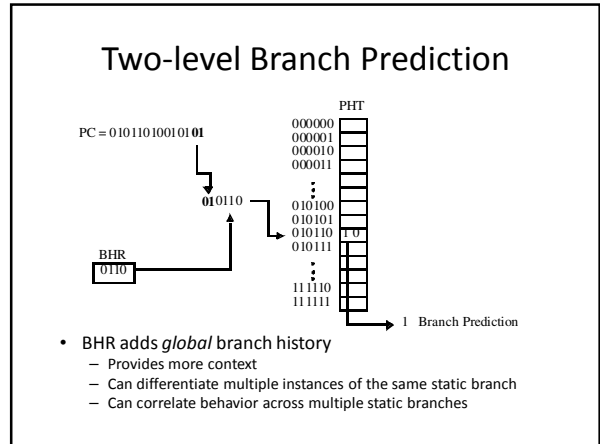
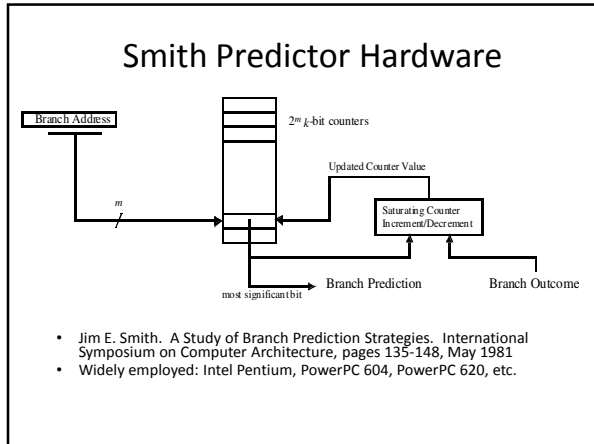
- Heuristic-based: Ball/Larus
 - Thomas Ball and James R. Larus. Branch Prediction for Free. ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, pages 300-313, May 1993.

Static Branch Prediction

- Profile-based
 1. Instrument program binary
 2. Run with representative (?) input set
 3. Recompile program
 - a. Annotate branches with hint bits, or
 - b. Restructure code to match predict not-taken
- Best performance: 75-80% accuracy

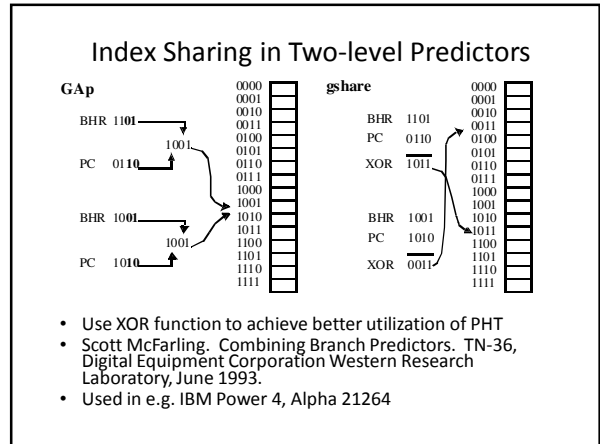
Dynamic Branch Prediction

- Main advantages:
 - Learn branch behavior autonomously
 - No compiler analysis, heuristics, or profiling
 - Adapt to changing branch behavior
 - Program phase changes branch behavior
- First proposed in 1980
 - US Patent #4,370,711, Branch predictor using random access memory, James. E. Smith
- Continually refined since then



Two-level Taxonomy

- Based on indices for branch history and pattern history
 - BHR: {G,P,S}: {Global, Per-address, Set}
 - PHT: {g,p,s}: {Global, Per-address, Set}
 - 9 combinations: GAg, GAp, GAs, PAg, PAp, PAs, SAg, SAp and SAs
- Tse-Yu Yeh and Yale N. Patt. Two-Level Adaptive Branch Prediction. International Symposium on Microarchitecture, pages 51-61, November 1991.



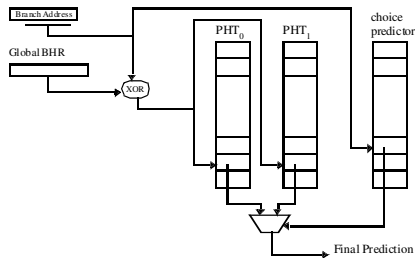
Sources of Mispredictions

- Lack of history (training time)
- Randomized behavior
 - Usually due to randomized input data
 - Surprisingly few branches depend on input data values
- BHR capacity
 - Correlate to branch that already shifted out
 - E.g. loop count > BHR width
- PHT capacity
 - Aliasing/interference
 - Positive
 - Negative

Reducing Interference

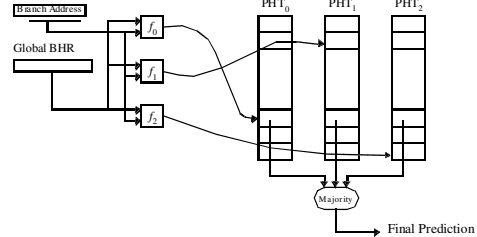
- Compulsory aliasing (*cold miss*)
 - Not important (less than 1%)
 - Only remedy is to set appropriate initial value
- Capacity aliasing (*capacity miss*)
 - Increase PHT size
- Conflict aliasing (*conflict miss*)
 - Change indexing scheme or partition PHT in a clever fashion

Bi-Mode Predictor



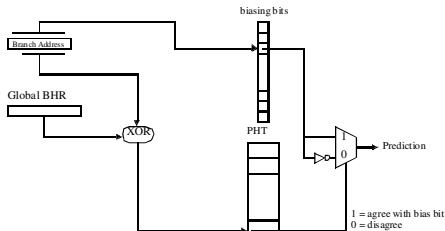
- PHT partitioned into T/NT halves
 - Selector chooses source
- Reduces negative interference, since most entries in PHT₀ tend towards NT, and most entries in PHT₁ tend towards T

gskewed Predictor



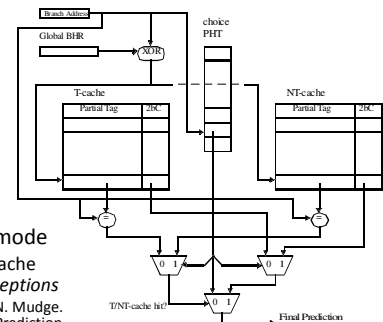
- Multiple PHT banks indexed by different hash functions
 - Conflicting branch pair unlikely to conflict in more than one PHT
- Majority vote determines prediction
- Used in Alpha EV8 (ultimately cancelled)
- P. Michaud, A. Seznec, and R. Uhlig. Trading Conflict and Capacity Aliasing in Conditional Branch Predictors. ISCA-24, June 1997

Agree Predictor



- Same principle as bi-mode
- PHT records whether branch bias matches outcome
 - Exploits 70-80% static predictability
- Used in in HP PA-8700
- E. Sprangle, R. S. Chappell, M. Alsop, and Y. N. Patt. The Agree Predictor: A Mechanism for Reducing Negative Branch History Interference. ISCA-24, June 1997.

YAGS Predictor



- Based on bi-mode
 - T/NT PHTs cache only the *exceptions*
- A. N. Eden and T. N. Mudge. The YAGS Branch Prediction Scheme. MICRO, Dec 1998.

Branch Filtering

- Highly-biased branches
 - e.g. '11111' history
 - Eliminated from PHT
- P-Y Chang, M. Evers, and Y. Patt. Improving Branch Prediction Accuracy by Reducing Pattern History Table Interference. PACT, October 1996.

Alloyed-History Predictors

- Local history vs. global history
- Kevin Skadron, Margaret Martonosi, and Douglas W. Clark. Alloyed Global and Local Branch History: A Robust Solution to Wrong-History Mispredictions. International Journal of Parallel Programming, 31(2), April 2003.

Path History

- Sometimes T/NT history is not enough
- Path history (PC values) can help

Path-Based Branch Predictor

- Ravi Nair. Dynamic Path-Based Branch Correlation. International Symposium on Microarchitecture, pages 15-23, December 1995.

Dynamic History Length

- Branch history length:
 - Some prefer short history (less training time)
 - Some require longer history (complex behavior)
- Vary history length
 - Choose through profile/compile-time hints
 - Or learn dynamically
- References
 - Maria-Dana Tarlescu, Kevin B. Theobald, and Guang R. Gao. Elastic History Buffer: A Low-Cost Method to Improve Branch Prediction Accuracy. ICCD, October 1996.
 - Toni Juan, Sanji Sanjeevan, and Juan J. Navarro. Dynamic History-Length Fitting: A Third Level of Adaptivity for Branch Prediction. ISCA, June 1998.
 - Jared Stark, Marius Evers, and Yale N. Patt. Variable Path Branch Prediction. ACM SIGPLAN Notices, 33(11):170-179, 1998

Loop Count Predictors

BHR entry:

- To predict last iteration's NT branch:
 - Must have length(BHR) > loop count
 - Not feasible for large loop counts
- Instead, BHR has mode bit
 - Once history == '111...11' or '000...00' switch to count mode
 - Now nth entry in PHT trains to NT and predicts nth iteration as last one
 - Now length(BHR) > log₂(loop count) is sufficient
- Used in Intel Pentium M/Core Duo/ Core 2 Duo

Understanding Advanced Predictors

- Four types of history
 - Local (bimodal) history (Smith predictor)
 - Table of counters summarizes local history
 - Simple, but only effective for biased branches
 - Local outcome history
 - Shift register of individual branch outcomes
 - Separate counter for each outcome history
 - Global outcome history
 - Shift register of recent branch outcomes
 - Separate counter for each outcome history
 - Path history
 - Shift register of recent (partial) block addresses
 - Can differentiate similar global outcome histories
- Can combine or “alloy” histories in many ways

Understanding Advanced Predictors

- History length
 - Short history—lower training cost
 - Long history—captures macro-level behavior
 - Variable history length predictors
- Really long history (long loops)
 - Loop count predictors
 - Fourier transform into frequency domain
- Limited capacity & interference
 - Constructive vs. destructive
 - Bi-mode, gskewed, agree, YAGS
 - Read sec. 9.3.2 carefully

Perceptron Branch Prediction

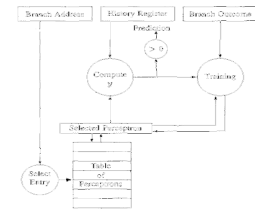
[Jimenez, Lin HPCA 2001]

$$y = w_0 + \sum_{i=1}^n x_i w_i$$

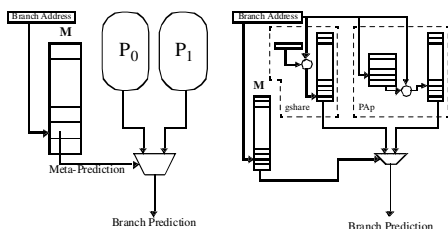
- Perceptron
 - Basis in AI concept [1962]
 - Computes boolean result based on multiple weighted inputs
- Adapted for branch prediction
 - x_i from branch history (1 T, -1 NT)
 - w_i incremented whenever branch outcome matches x_i
 - Finds correlation between current branch and any subset of prior branches

Perceptrons - Implementation

- Complex dot product must be computed for every prediction
 - Too slow
- Arithmetic tricks, pipelining:
 - Daniel A. Jimenez and Calvin Lin. Neural methods for dynamic branch prediction. ACM Transactions on Computer Systems, 20(4):369–397, November 2002.
- Key insight:
 - Not all branches in history are important
 - Perceptron weights learn this

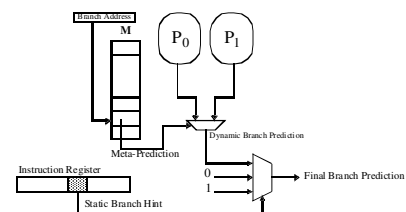


Combining or Hybrid Predictors

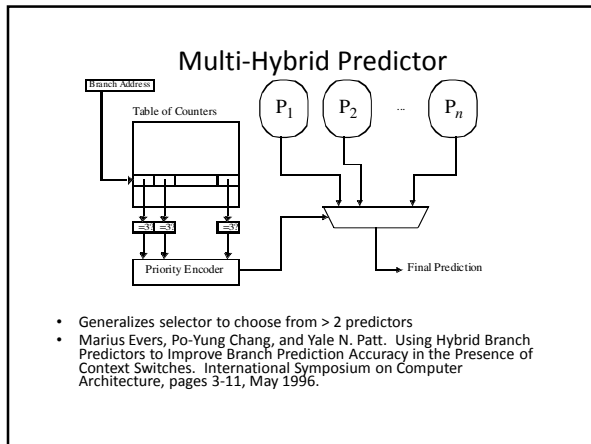


- Select “best” history
- Reduce interference w/partial updates
- Scott McFarling. Combining Branch Predictors. TN-36, Digital Equipment Corporation Western Research Laboratory, June 1993.

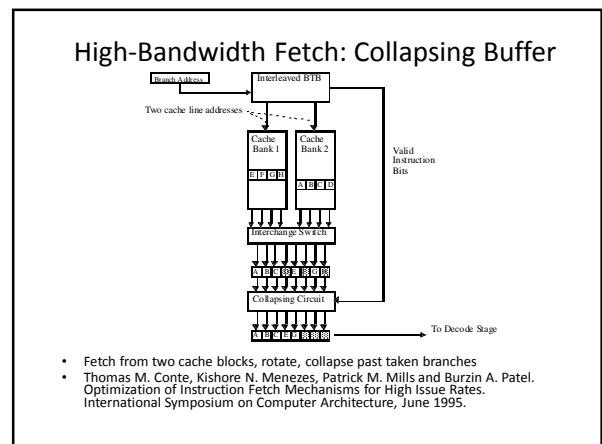
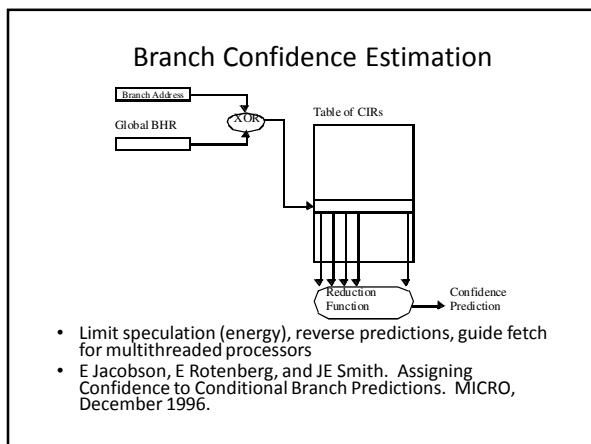
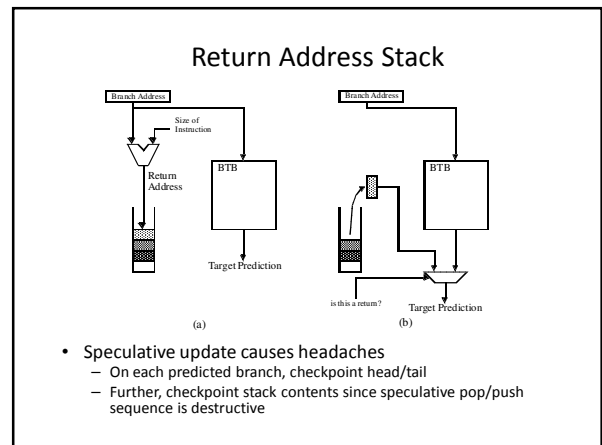
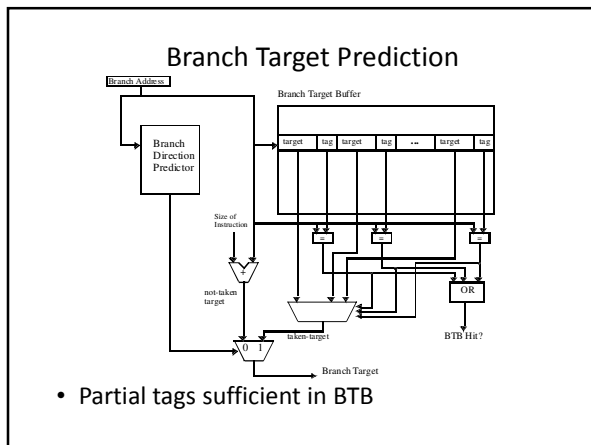
Branch Classification



- Static (profile-based) branch hints select which prediction to use
 - Static T/Static NT/Dynamic
- P-Y Chang, E Hao, TY Yeh, and Y Patt. Branch Classification: a New Mechanism for Improving Branch Predictor Performance. MICRO, Nov. 1994.
- D Grunwald, D Lindsay, and B Zorn. Static Methods in Hybrid Branch Prediction. PACT, October 1998



- ### Multiple History Lengths
- Championship Branch Prediction (CBP)
 - 2 contests, standardized methods and traces
 - Insight from perceptron BP:
 - Some branches need short history
 - Others need very long history
 - Geometric history length (O-GEHL)
 - Geometric series of history lengths
 - Tagged Geometric History Length (TAGE)
 - Choose longest matching history
- © Shen, Lipasti 38



High-Bandwidth Fetch: Trace Cache

- Fold out taken branches by *tracing* instructions as they commit into a *fill buffer*
- Eric Rotenberg, S. Bennett, and James E. Smith. Trace Cache: A Low Latency Approach to High Bandwidth Instruction Fetching. MICRO, December 1996.

Intel Pentium 4 Trace Cache

- No first-level instruction cache: trace cache only
- Trace cache BTB identifies next trace
- Miss leads to fetch from level two cache
- Trace cache instructions are decoded (uops)

High Frequency: Next-line Prediction

- Embed next fetch address in instruction cache
– Enables high-frequency back-to-back fetch
- Brad Calder and Dirk Grunwald. Next Cache Line and Set Prediction. International Symposium on Computer Architecture, pages 287-296, June 1995.

High Frequency: Overriding Predictors

If slow predictor agrees with fast predictor, do nothing
If predictors do not match, flush A, B, and C, and restart fetch at new predicted target

- Simple, fast predictor turns around every cycle
- Smarter, slower predictor can override
- Widely used: PowerPC 604, 620, Alpha 21264

Advanced Branch Prediction Summary

- Control Flow Speculation
 - Branch Speculation
 - Mis-speculation Recovery
- Branch Direction Prediction
 - Static Prediction
 - Dynamic Prediction
 - Hybrid Prediction
- Branch Target Prediction
- High-bandwidth Fetch
- High-Frequency Fetch