15-740/18-740 Computer Architecture Lecture 29: Control Flow II

Prof. Onur Mutlu
Carnegie Mellon University
Fall 2011, 11/30/11

Announcements for This Week

- December 2:
 - Midterm II
 - Comprehensive
 - 2 letter-sized cheat sheets allowed
 - Sample exams and solutions are all posted
- December 1:
 - Homework 6 due

Project Schedule

- December 6 and 8
 - Milestone III meetings
 - Signup sheets will be posted online
- December 13
 - Project poster session
 - Location and time TBD
- December 18
 - Final project report due
 - Strive for a good conference paper (in terms of insight, explanations, writing, and formatting)
 - Address comments from the poster session and milestone III

Readings

Required:

- McFarling, "Combining Branch Predictors," DEC WRL TR, 1993.
- Carmean and Sprangle, "Increasing Processor Performance by Implementing Deeper Pipelines," ISCA 2002.

Recommended:

- Evers et al., "An Analysis of Correlation and Predictability: What Makes Two-Level Branch Predictors Work," ISCA 1998.
- Yeh and Patt, "Alternative Implementations of Two-Level Adaptive Branch Prediction," ISCA 1992.
- Jimenez and Lin, "Dynamic Branch Prediction with Perceptrons," HPCA 2001.
- Kim et al., "Diverge-Merge Processor (DMP): Dynamic Predicated Execution of Complex Control-Flow Graphs Based on Frequently Executed Paths," MICRO 2006.

Approaches to Conditional Branch Handling

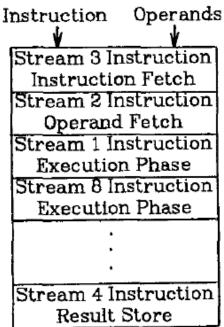
- Branch prediction
 - Static
 - Dynamic
- Eliminating branches
 - I. Predicated execution
 - Static
 - Dynamic
 - HW/SW Cooperative
 - II. Predicate combining (and condition registers)
- Multi-path execution
- Delayed branching (branch delay slot)
- Fine-grained multithreading

Fine-Grained Multithreading

- Idea: Hardware has multiple thread contexts. Each cycle, fetch engine fetches from a different thread.
 - By the time the fetched branch resolves, there is no need to fetch another instruction from the same thread

 Branch resolution latency overlapped with execution of other threads' instructions

- + No logic needed for branch prediction, (also for dependency checking)
- -- Single thread performance suffers
- -- Does not overlap latency if not enough threads to cover the whole pipeline
- -- Extra logic for keeping thread contexts

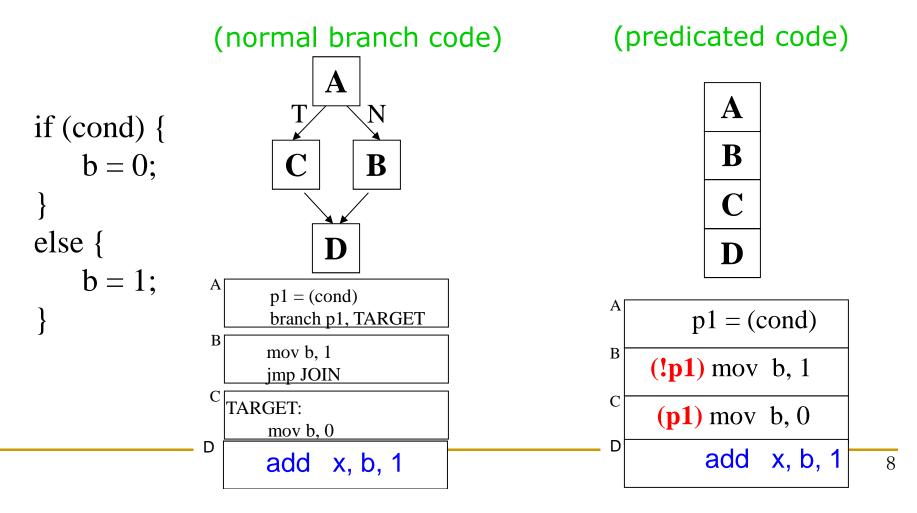


Approaches to Conditional Branch Handling

- Branch prediction
 - Static
 - Dynamic
- Eliminating branches
 - I. Predicated execution
 - Static
 - Dynamic
 - HW/SW Cooperative
 - II. Predicate combining (and condition registers)
- Multi-path execution
- Delayed branching (branch delay slot)
- Fine-grained multithreading

Predication (Predicated Execution)

- Idea: Compiler converts control dependency into a data dependency → branch is eliminated
 - Each instruction has a predicate bit set based on the predicate computation
 - Only instructions with TRUE predicates are committed (others turned into NOPs)

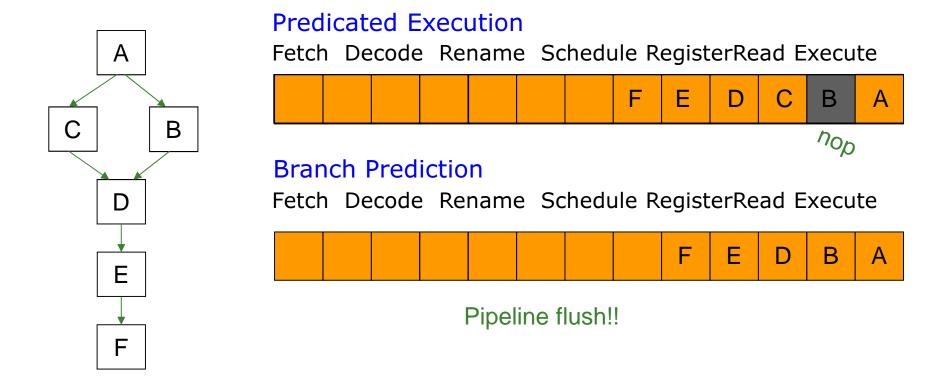


Conditional Move Operations

- Very limited form of predicated execution
- CMOV R1 ← R2
 - □ R1 = (ConditionCode == true) ? R2 : R1
 - Employed in most modern ISAs (x86, Alpha)

Predicated Execution (II)

 Predicated execution can be high performance and energyefficient



Predicated Execution (III)

Advantages:

- + Eliminates mispredictions for hard-to-predict branches
 - + No need for branch prediction for some branches
 - + Good if misprediction cost > useless work due to predication
- + Enables code optimizations hindered by the control dependency
 - + Can move instructions more freely within predicated code
 - + Vectorization with control flow
- + Reduces fetch breaks (straight-line code)

Disadvantages:

- -- Causes useless work for branches that are easy to predict
 - -- Reduces performance if misprediction cost < useless work
 - -- Adaptivity: Static predication is not adaptive to run-time branch behavior. Branch behavior changes based on input set, phase, control-flow path.
- -- Additional hardware and ISA support (complicates renaming and OOO)
- -- Cannot eliminate all hard to predict branches
 - -- Complex control flow graphs, function calls, and loop branches
- -- Additional data dependencies delay execution (problem esp. for easy branches)

Predicated Execution (III)

Advantages:

- + Eliminates mispredictions for hard-to-predict branches
 - + No need for branch prediction for some branches
 - + Good if misprediction cost > useless work due to predication
- + Enables code optimizations hindered by the control dependency
 - + Can move instructions more freely within predicated code
 - + Vectorization with control flow
- + Reduces fetch breaks (straight-line code)

Disadvantages:

- -- Causes useless work for branches that are easy to predict
 - -- Reduces performance if misprediction cost < useless work
 - -- Adaptivity: Static predication is not adaptive to run-time branch behavior. Branch behavior changes based on input set, phase, control-flow path.
- -- Additional hardware and ISA support (complicates renaming and OOO)
- -- Cannot eliminate all hard to predict branches
 - -- Complex control flow graphs, function calls, and loop branches
- -- Additional data dependencies delay execution (problem esp. for easy branches)

Approaches to Conditional Branch Handling

- Branch prediction
 - Static
 - Dynamic
- Eliminating branches
 - I. Predicated execution
 - Static
 - Dynamic
 - HW/SW Cooperative
 - II. Predicate combining (and condition registers)
- Multi-path execution
- Delayed branching (branch delay slot)
- Fine-grained multithreading

Multi-Path Execution

Idea: Execute both paths after a conditional branch

- For all branches: Riseman and Foster, "The inhibition of potential parallelism by conditional jumps," IEEE Transactions on Computers, 1972.
- For a hard-to-predict branch: Use dynamic confidence estimation

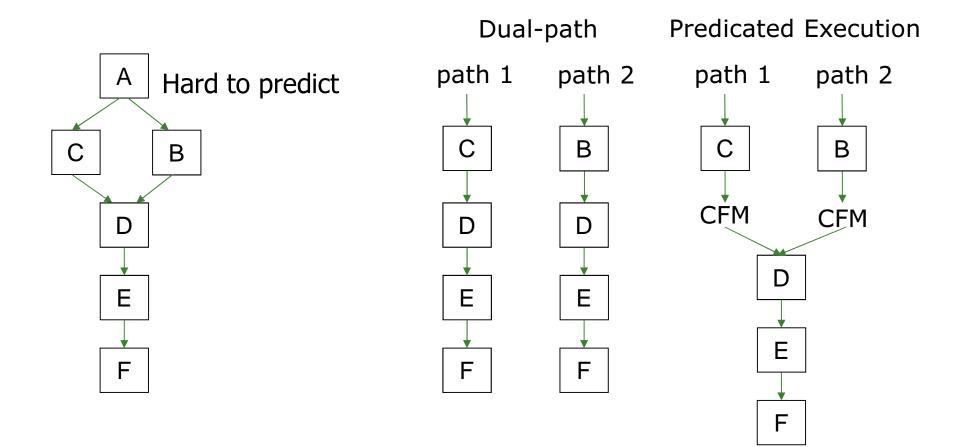
Advantages:

- + Improves performance if misprediction cost > useless work
- + No ISA change needed

Disadvantages:

- -- What happens when the machine encounters another hard-to-predict branch? Execute both paths again?
 - -- Paths followed quickly become exponential
- -- Each followed path requires its own register alias table, PC, GHR
- -- Wasted work (and reduced performance) if paths merge

Dual-Path Execution versus Predication

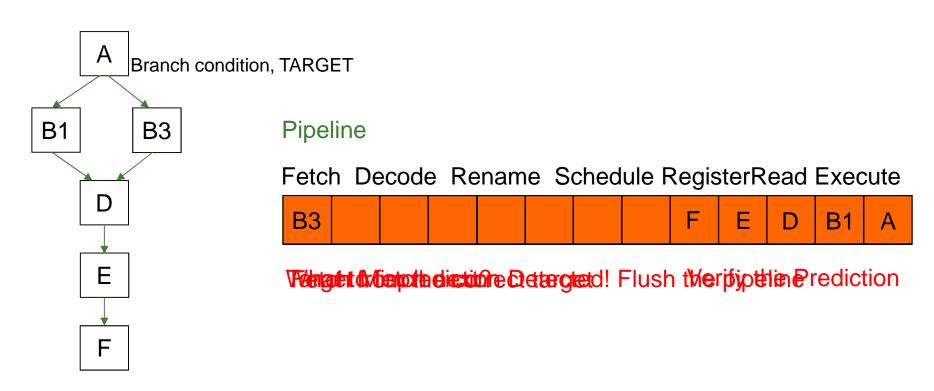


Approaches to Conditional Branch Handling

- Branch prediction
 - Static
 - Dynamic
- Eliminating branches
 - I. Predicated execution
 - Static
 - Dynamic
 - HW/SW Cooperative
 - II. Predicate combining (and condition registers)
- Multi-path execution
- Delayed branching (branch delay slot)
- Fine-grained multithreading

Branch Prediction

- Processors are pipelined to increase concurrency
- How do we keep the pipeline full in the presence of branches?
 - Guess the next instruction when a branch is fetched
 - Requires guessing the direction and target of a branch

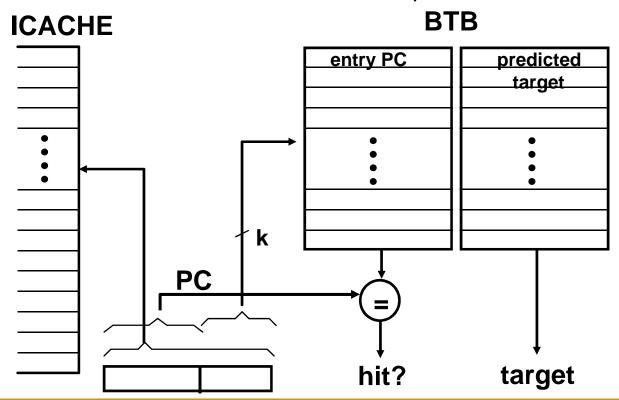


Branch Prediction

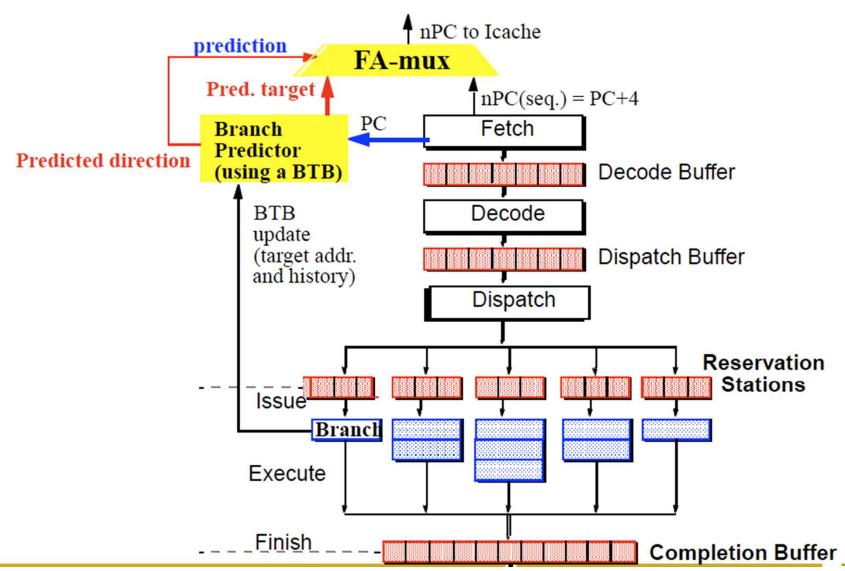
- Idea: Predict the next fetch address (to be used in the next cycle) when the branch is fetched
- Requires three things to be predicted:
 - Whether the fetched instruction is a branch
 - Conditional branch direction
 - Branch target address (if taken)
- Target addresses remain the same for conditional direct branches across dynamic instances
 - Idea: Cache the target address from previous instance
 - Called Branch Target Buffer (BTB) or Branch Target Address
 Cache

Branch Target Buffer

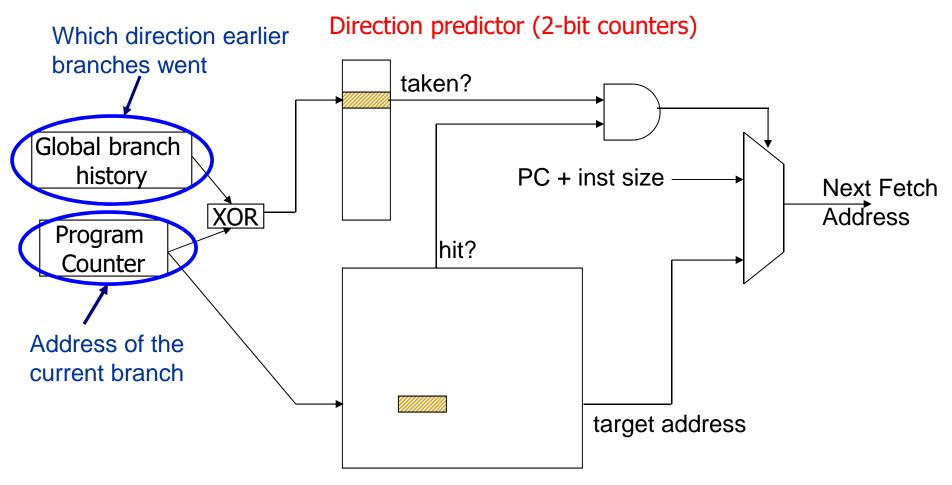
- Cache of branch target addresses accessed in parallel with the I-cache in the fetch stage
- Updated only by taken branches
- If BTB hit and the instruction is a predicted-taken branch
 - target from the BTB (assuming hit) is used as fetch address in the next cycle
- If BTB miss or the instruction is a predicted-not-taken branch
 - PC+N is used as the next fetch address in the next cycle



Branch Target Buffer in Fetch Stage



A Frontend with BTB and Direction Prediction



Cache of Target Addresses (BTB: Branch Target Buffer)

Direction Prediction

- 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)
 - Last time (single-bit)
 - Two-bit counter based
 - Two-level (global vs. local)
 - Hybrid

Static Branch Prediction (I)

- Always not-taken
 - Simple to implement: no need for BTB, no direction prediction
 - Low accuracy: ~40%
 - Compiler can layout code such that the likely path is the "not-taken" path: Good for wide fetch as well!
- Always taken
 - No direction prediction
 - □ Better accuracy: ~60%
 - Backward branches (i.e. loop branches) are usually taken
 - Backward branch: target address lower than branch PC
- Backward taken, forward not taken (BTFN)
 - □ Predict backward (loop) branches as taken, others not-taken

Static Branch Prediction (II)

Profile-based

- Idea: Compiler determines likely direction for each branch using profile run. Encodes that direction as a hint bit in the branch instruction format.
- + Per branch prediction (more accurate than schemes in previous slide)
- -- Requires hint bits in the branch instruction format
- Accuracy depends on the representativeness of profile input set

Static Branch Prediction (III)

Program-based

- Idea: Use heuristics based on program analysis to determine statically-predicted direction
- Opcode heuristic: Predict BLEZ as NT (negative integers used as error values in many programs)
- Loop heuristic: Predict a branch guarding a loop execution as taken (i.e., execute the loop)
- Pointer and FP comparisons: Predict not equal
- + Does not require profiling
- -- Heuristics might be not representative or good
- -- Requires ISA support
- Ball and Larus, "Branch prediction for free," PLDI 1993.
 - 20% misprediction rate

Dynamic Branch Prediction

 Idea: Predict branches based on dynamic information (collected at run-time)

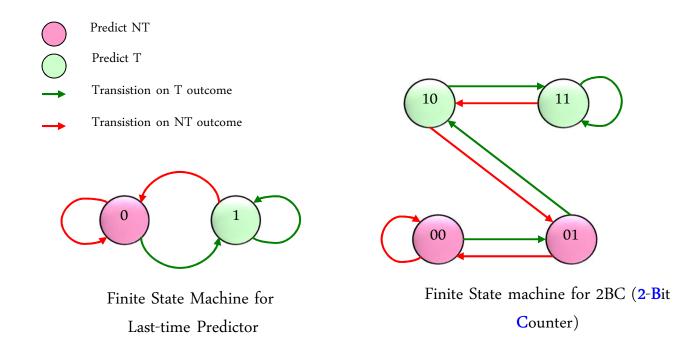
Advantages

- + No need for profiling: input set representativeness problem goes away
- + Prediction based on history of the execution of branches
 - + It can adapt to dynamic changes in branch behavior
- Disadvantages
 - -- More complex (requires additional hardware)

Last Time Predictor

- Last time predictor
 - Single bit per branch (stored in BTB)
- Always mispredicts the last iteration and the first iteration of a loop branch
 - \square Accuracy for a loop with N iterations = (N-2)/N
- + Loop branches for loops with large number of iterations

Two-Bit Counter Based Prediction



- Counter using saturating arithmetic
 - □ There is a symbol for maximum and minimum values

Two-Bit Counter Based Prediction

- Each branch associated with a two-bit counter
- One more bit provides hysteresis
- A strong prediction does not change with one single different outcome
- Accuracy for a loop with N iterations = (N-1)/N
 TNTNTNTNTNTNTNTNTN → 50% accuracy

(assuming init to weakly taken)

- + Better prediction accuracy
- -- More hardware cost (but counter can be part of a BTB entry)

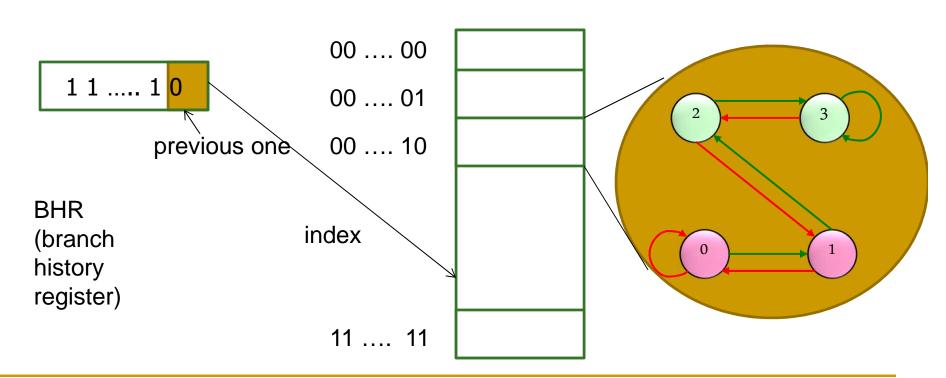
Can We Do Better?

If the loop test is done at the end of the body, the corresponding branch will execute the pattern $(1110)^n$, where 1 and 0 represent taken and not taken respectively, and n is the number of times the loop is executed. Clearly, if we knew the direction this branch had gone on the previous three executions, then we could always be able to predict the next branch direction.

 McFarling, "Combining Branch Predictors," DEC WRL TR 1993.

Two Level Branch Predictors

- First level: Branch history register (N bits)
 - The direction of last N branches
- Second level: Table of saturating counters for each history entry
 - The direction the branch took the last time the same history was seen?
 Pattern History Table (PHT)



Prediction and Update Functions

Prediction

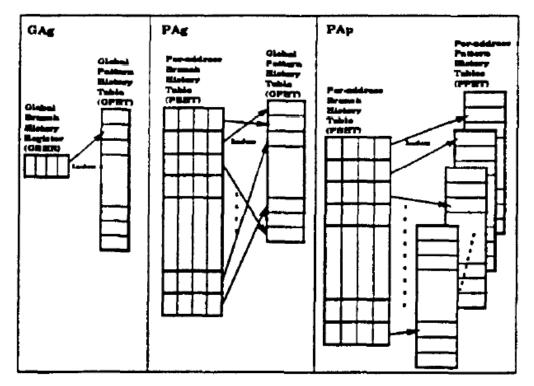
- Pattern History Table accessed at fetch time to generate a prediction
- Top bit of the 2-bit counter determines predicted direction

Update

- Pattern History Table accessed when the branch is retired to update the counters that generated the prediction
- If branch
 - actually taken: increment the counter
 - actually not-taken: decrement the counter

Two-Level Predictor Variations

- BHR can be global (G), per set of branches (S), or per branch (P)
- PHT counters can be adaptive (A) or static (S)
- PHT can be global (g), per set of branches (s), or per branch (p)



Yeh and Patt, "Two-Level Adaptive Training Branch Prediction," MICRO 1991.

Global Branch Correlation (I)

- GAg: Global branch predictor (commonly called)
- Exploits global correlation across branches
- Recently executed branch outcomes in the execution path is correlated with the outcome of the next branch

```
if (cond1)
...
if (cond1 AND cond2)
```

If first branch not taken, second also not taken

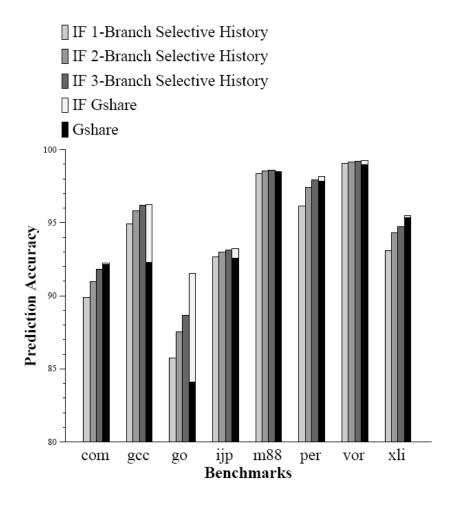
```
branch Y: if (cond1) a = 2;
...
branch X: if (a == 0)
```

If first branch taken, second definitely not taken

Global Branch Correlation (II)

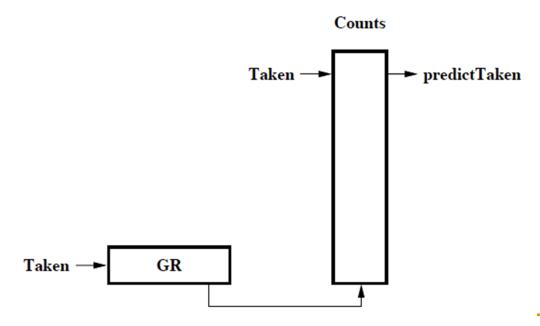
```
branch Y: if (cond1)
...
branch Z: if (cond2)
...
branch X: if (cond1 AND cond2)
```

- If Y and Z both taken, then X also taken
- If Y or Z not taken, then X also not taken
- Only 3 past branches' directions really matter (not necessarily the last 3 past branches)
- Evers et al., "An Analysis of Correlation and Predictability: What Makes Two-Level Branch Predictors Work," ISCA 1998.



Global Two-Level Prediction

- Idea: Have a single history register for all branches (called global history register)
- + Exploits correlation between different branches (as well as the instances of the same branch)
- -- Different branches interfere with each other in the history register → cannot separate the local history of each branch



How Does the Global Predictor Work?

After the initial startup time, the conditional branches have the following behavior, assuming GR is shifted to the left:

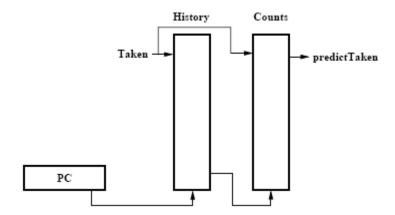
| test | value | GR | result |
|-------|-------|------|---------------|
| j<3 | j=1 | 1101 | taken |
| j<3 | j=2 | 1011 | taken |
| j<3 | j=3 | 0111 | not taken |
| i<100 | | 1110 | usually taken |

Pentium Pro Branch Predictor

- GAs
- 4-bit global history register
- Multiple pattern history tables (of 2 bit counters)
 - PHT determined by lower order bits of the branch address

Local Two-Level Prediction

- PAg, Pas, PAp
- Global history register produces interference
 - Different branches can go different ways for the same history
- Idea: Have a per-branch history register



- + No interference in the history register between branches
- -- Cannot exploit global branch correlation

Hybrid Branch Predictors

- Idea: Use more than one type of predictors (i.e., algorithms) and select the "best" prediction
 - E.g., hybrid of 2-bit counters and global predictor

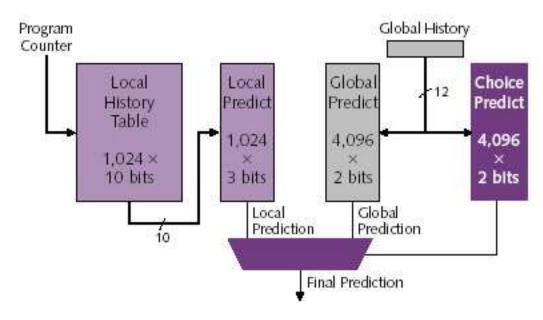
Advantages:

- + Better accuracy: different predictors are better for different branches
- + Reduced warmup time (faster-warmup predictor used until the slower-warmup predictor warms up)

Disadvantages:

- -- Need "meta-predictor" or "selector"
- -- Longer access latency
- McFarling, "Combining Branch Predictors," DEC WRL Tech Report, 1993.

Alpha 21264 Tournament Predictor



- Minimum branch penalty: 7 cycles
- Typical branch penalty: 11+ cycles
- 48K bits of target addresses stored in I-cache
- 32-entry return address stack
- Predictor tables are reset on a context switch

Effect on Prediction Accuracy

Bimodal: table of 2bc indexed by branch address

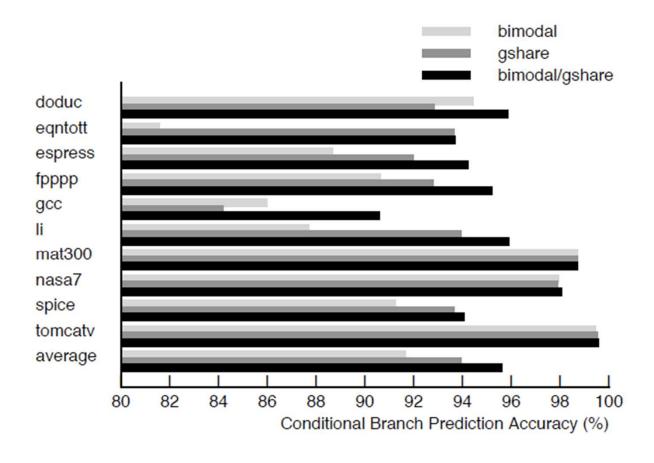
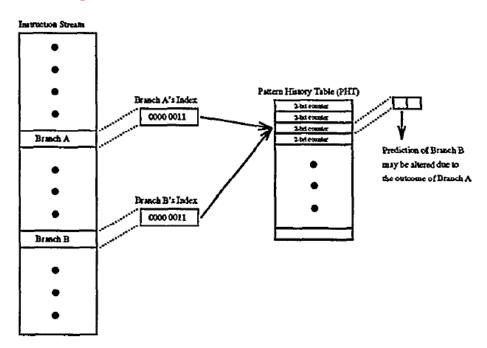


Figure 13: Combined Predictor Performance by Benchmark

The remaining slides are not covered in lecture. They are for your benefit.

Interference in the PHTs

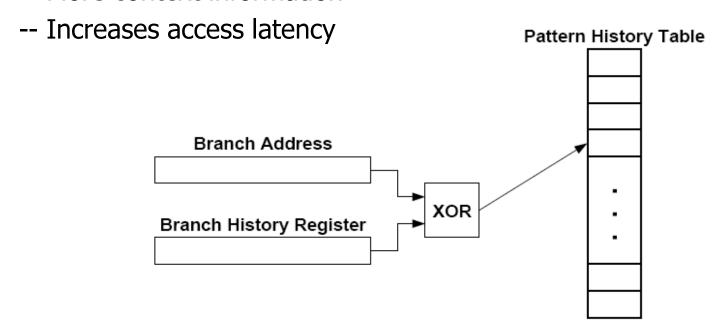
- Sharing the PHTs between histories/branches leads to interference
 - Different branches map to the same PHT entry and modify it
 - Can be positive, negative, or neutral



- Interference can be eliminated by dedicating a PHT per branch
 - -- Too much hardware cost
- How else can you eliminate interference?

Reducing Interference in PHTs (II)

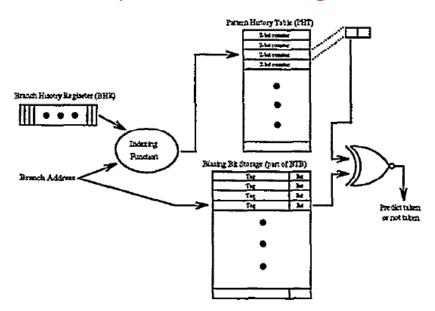
- Idea 1: Randomize the indexing function into the PHT such that probability of two branches mapping to the same entry reduces
 - Gshare predictor: GHR hashed with the Branch PC
 - + Better utilization of PHT
 - + More context information



McFarling, "Combining Branch Predictors," DEC WRL Tech Report, 1993.

Reducing Interference in PHTs (III)

- Idea 2: Agree prediction
 - Each branch has a "bias" bit associated with it in BTB
 - Ideally, most likely outcome for the branch
 - High bit of the PHT counter indicates whether or not the prediction agrees with the bias bit (not whether or not prediction is taken)
 - + Reduces negative interference (Why???)
 - -- Requires determining bias bits (compiler vs. hardware)



Sprangle et al., "The Agree Predictor: A Mechanism for Reducing Negative Branch History Interference," ISCA 1997.

Why Does Agree Prediction Make Sense?

- Assume two branches (b1, b2) have taken rates of 85% and 15%.
- Assume they conflict in the PHT
- Probability they have opposite outcomes
 - Baseline predictor:
 - P (b1 T, b2 NT) + P (b1 NT, b2 T) = (85%*85%) + (15%*15%) = 74.5%
 - Agree predictor:
 - Assume bias bits are set to T (b1) and NT (b2)
 - P (b1 agree, b2 disagree) + P (b1 disagree, b2 agree)
 - = (85%*15%) + (15%*85%) = 25.5%
- Agree prediction reduces the probability that two branches have opposite predictions in the PHT entry
 - Works because most branches are biased (not 50% taken)

Improved Branch Prediction Algorithms

Perceptron predictor

- Learns the correlations between branches in the global history register and the current branch using a perceptron
- Past branches that are highly correlated have larger weights and influence the prediction outcome more
- Jimenez and Lin, "Dynamic Branch Prediction with Perceptrons," HPCA 2001.

Enhanced hybrid predictors

- Multi-hybrid with different history lengths
- Seznec, "Analysis of the O-GEometric History Length Branch Predictor," ISCA 2005.

Pre-execution

- Similar to pre-execution based prefetching
- Chappell et al., "Difficult-Path Branch Prediction Using Subordinate Microthreads," ISCA 2002.

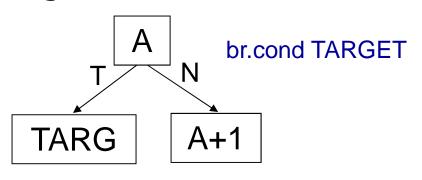
Call and Return Prediction

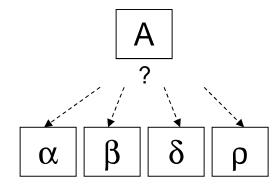
- Direct calls are easy to predict
 - Always taken, single target
 - Call marked in BTB, target predicted by BTB
- Returns are indirect branches
 - A function can be called from many points in code
 - A return instruction can have many target addresses
 - Next instruction after each call point for the same function
 - Observation: Usually a return matches a call
 - Idea: Use a stack to predict return addresses (Return Address Stack)
 - A fetched call: pushes the return (next instruction) address on the stack
 - A fetched return: pops the stack and uses the address as its predicted target
 - Accurate most of the time: 8-entry stack \rightarrow > 95% accuracy



Indirect Branch Prediction (I)

Register-indirect branches have multiple targets





R1 = MEM[R2] branch R1

Conditional (Direct) Branch

Indirect Jump

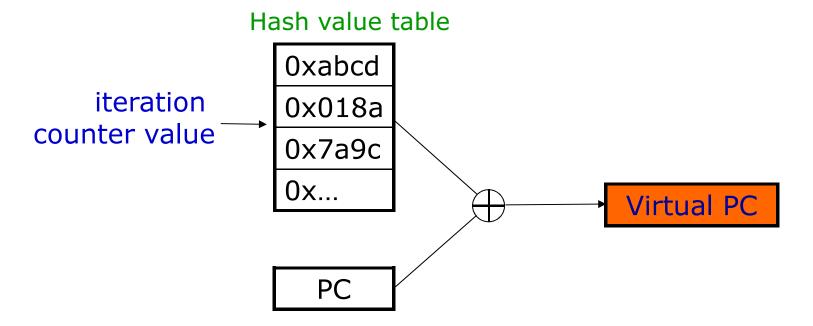
- Used to implement
 - Switch-case statements
 - Virtual function calls
 - Jump tables (of function pointers)
 - Interface calls

Indirect Branch Prediction (II)

- No direction prediction needed
- Idea 1: Predict the last resolved target as the next fetch address
 - + Simple: Use the BTB to store the target address
 - -- Inaccurate: 50% accuracy (empirical). Many indirect branches switch between different targets
- Idea 2: Use history based target prediction
 - E.g., Index the BTB with GHR XORed with Indirect Branch PC
 - Chang et al., "Target Prediction for Indirect Jumps," ISCA 1997.
 - + More accurate
 - -- An indirect branch maps to (too) many entries in BTB
 - -- Conflict misses with other branches (direct or indirect)
 - -- Inefficient use of space if branch has few target addresses

Indirect Branch Prediction (III)

- Idea 3: Treat an indirect branch as "multiple virtual conditional branches" in hardware
 - Only for prediction purposes
 - Predict each "virtual conditional branch" iteratively
 - □ Kim et al., "VPC prediction," ISCA 2007.



VPC Prediction (I)

Real Instruction

call R1

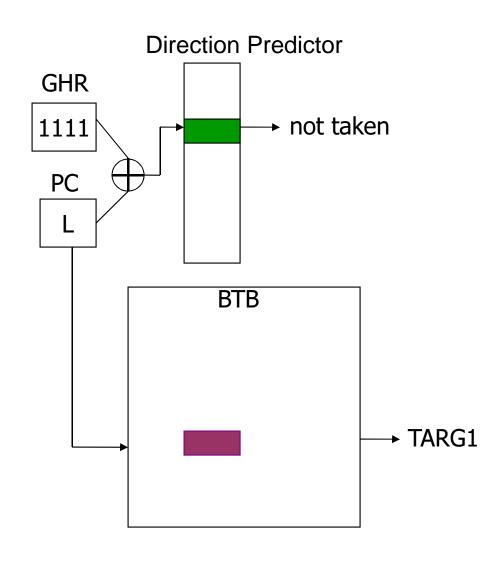
// PC: L

Virtual Instructions

cond. jump TARG1 // VPC: L cond. jump TARG2 // VPC: VL2 cond. jump TARG3 // VPC: VL3

cond. jump TARG4 // VPC: VL4

Next iteration



VPC Prediction (II)

Real Instruction

call R1

// PC: L

Virtual Instructions

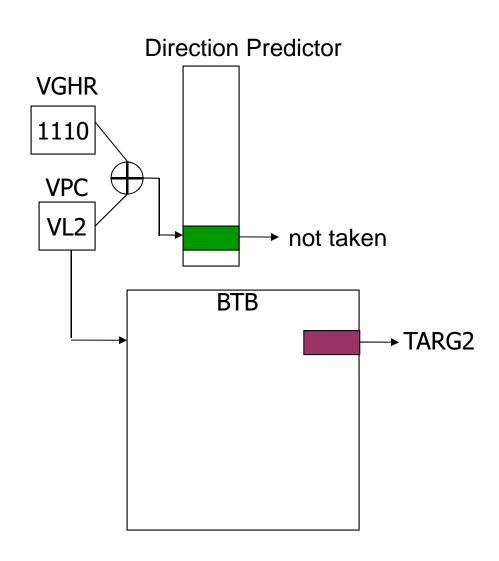
cond. jump TARG1 // VPC: L

cond. jump TARG2 // VPC: VL2

cond. jump TARG3 // VPC: VL3

cond. jump TARG4 // VPC: VL4

Next iteration



VPC Prediction (III)

Real Instruction

call R1

// PC: L

Virtual Instructions

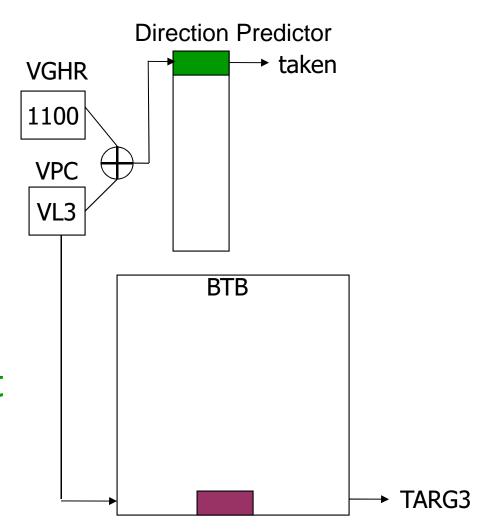
cond. jump TARG1 // VPC: L

cond. jump TARG2 // VPC: VL2

cond. jump TARG3 // VPC: VL3

cond. jump TARG4 // VPC: VL4

Predicted Target = TARG3



VPC Prediction (IV)

Advantages:

- + High prediction accuracy (>90%)
- + No separate indirect branch predictor
- + Resource efficient (reuses existing components)
- + Improvement in conditional branch prediction algorithms also improves indirect branch prediction
- + Number of locations in BTB consumed for a branch = number of target addresses seen

Disadvantages:

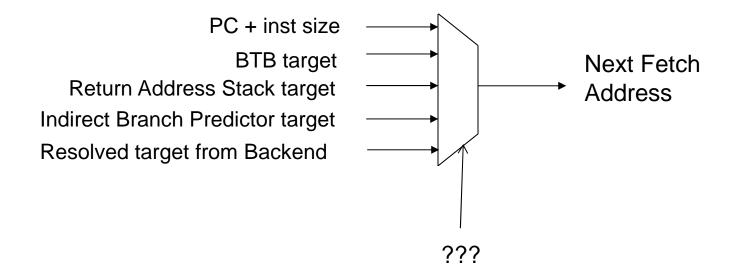
- -- Takes multiple cycles (sometimes) to predict the target address
- -- More interference in direction predictor and BTB

Issues in Branch Prediction (I)

- Need to identify a branch before it is fetched
- How do we do this?
 - □ BTB hit → indicates that the fetched instruction is a branch
 - BTB entry contains the "type" of the branch
- What if no BTB?
 - Bubble in the pipeline until target address is computed
 - E.g., IBM POWER4

Issues in Branch Prediction (II)

- Latency: Prediction is latency critical
 - Need to generate next fetch address for the next cycle
 - Bigger, more complex predictors are more accurate but slower



Issues in Branch Prediction (III)

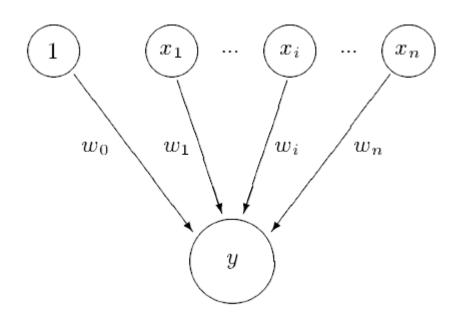
- State recovery upon misprediction
 - Misprediction detected when branch executes
 - Need to flush all instructions younger than the branch
 - Easy to invalidate instructions not yet renamed
 - Need to invalidate instructions in reservation stations and reorder buffer
 - Need to recover the Register Alias Table
 - Pentium 4: Retirement RAT copied to Frontend RAT
 - + Simple
 - -- Increases recovery latency (Branch has to be the oldest instruction in the machine!) ← Why is this not as bad???
 - Alpha 21264: Checkpoint RAT when branch is renamed, recover to checkpoint when misprediction detected
 - + Immediate recovery of RAT
 - -- More expensive (multiple RATs)

Open Research Issues in Branch Prediction

- Better algorithms
 - Machine learning techniques?
 - Needs to be low cost and *fast*
- Progressive evaluation of earlier prediction for a branch
 - As branch moves through the pipeline, more information becomes available → can we use this to override earlier prediction?
 - □ Falcon et al., "Prophet-critic hybrid branch prediction," ISCA 2004.

Perceptron Branch Predictor (I)

- Idea: Use a perceptron to learn the correlations between branch history register bits and branch outcome
- A perceptron learns a target Boolean function of N inputs



Each branch associated with a perceptron

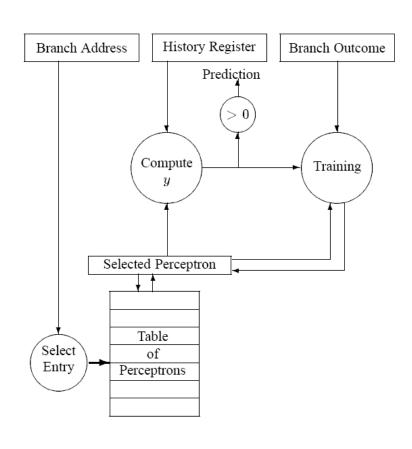
A perceptron contains a set of weights wi

- → Each weight corresponds to a bit in the GHR
- →How much the bit is correlated with the direction of the branch
- → Positive correlation: large + weight
- → Negative correlation: large weight

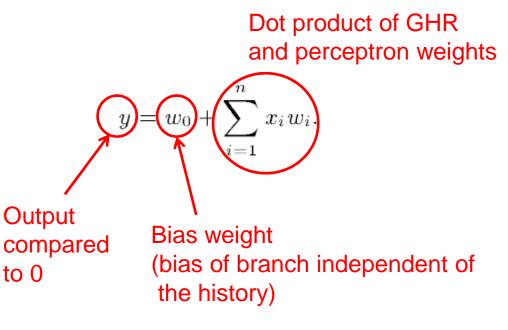
Prediction:

- → Express GHR bits as 1 (T) and -1 (NT)
- → Take dot product of GHR and weights
- \rightarrow If output > 0, predict taken
- Jimenez and Lin, "Dynamic Branch Prediction with Perceptrons," HPCA 2001.
- Rosenblatt, "Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms," 1962

Perceptron Branch Predictor (II)



Prediction function:



Training function:

if
$$\mathrm{sign}(y_{out}) \neq t \text{ or } |y_{out}| \leq \theta$$
 then for $i \coloneqq 0$ to n do $w_i \coloneqq w_i + tx_i$ end for end if

Perceptron Branch Predictor (III)

Advantages

+ More sophisticated learning mechanism → better accuracy

Disadvantages

- -- Hard to implement (adder tree to compute perceptron output)
- -- Can learn only linearly-separable functions
 e.g., cannot learn XOR type of correlation between 2 history
 bits and branch outcome