Computer Architecture

Lecture 08– Branch Prediction

Tian Xia

Institute of Artificial Intelligence and Robotics Xi'an Jiaotong University

http://gr.xjtu.edu.cn/web/pengjuren

Recap: Phases of Instruction Execution of OoO



Jump/Branch Instructions (in RISC-V)



Jump/Branch Instructions (in RISC-V)

| | PC-Relative Address | Absolute Address |
|---------------------------------|---|---|
| Unconditional Jump | lf Jump? <mark>YES</mark> | If Jump? YES |
| | Jump Where? Known Fast | Jump Where? Known Slow (depend on regs) |
| Conditional Jump (Branch) | If Jump? Conditional (depend on regs) | Jump Destination Prediction |
| | Jump Where? Known Fast | Branch Direction Prediction |

Incorrect Control-Flow Penalty



Importance of Branch Prediction



- Consider 4-way superscalar with 8 pipeline stages from fetch to dispatch, and 80-entry ROB, and 3 cycles from issue to branch resolution
- On a misprediction, could throw away 8*4+(80-1)=111 instructions (in worst case)
- Improving from 90% to 95% prediction accuracy, removes 50% of branch mispredictions
 - If 1/6 instructions are branches, then move from 60 instructions between mispredictions, to 120 instructions between mispredictions

Misprediction Penalty in SuperScalars



Misprediction Penalty in OoO Pipeline



- Steady instruction issue throughput (Insn-Per-Cycle, IPC) is broken
- Drain: linearly drop to 0
- **Refill**: wait for useful instructions reach issue port
- **Recover**: issue bandwidth grows in curve
- Wasted cycles = unused_slots / issue_throughput
 Typically 10—20 cycles in Intel Xeon CPU (4 issue bandwidth)

Average Run-Length between Branches

What is the average run length between braches ?

| | SPEC(int) | SPEC(fp) |
|----------|-----------|----------|
| Branches | 19% | 11% |
| Loads | 24% | 26% |
| Stores | 10% | 7% |
| Others | 47% | 56% |

Roughly 1 jump for every 5-10 instructions

- Essential in modern processors to mitigate branch delay latency
 - Two types of Prediction:
 - **Predict Branch Direction** (Taken or not ?)
 - Predict Branch/Jump Address (Target address ? PC related and Absolute)
- Typical Ideal Prediction: >95%

Reducing Control-Flow Penalty

- Software solutions
 - Eliminate branches loop unrolling
 - Increases the run length
 - Reduce resolution time instruction scheduling
 - Compute the branch condition as early as possible (of limited value because branches often in critical path through code)
- Hardware solutions
 - Find something else to do (delay slots)
 - Replaces pipeline bubbles with useful work (requires software cooperation) – quickly see diminishing returns
 - Speculate, i.e., branch prediction
 - Speculative execution of instructions beyond the branch
 - Many advances in accuracy, widely used

Branch Prediction

Motivation:

Branch penalties limit performance of deeply pipelined processors Modern branch predictors have high accuracy (>95%) and can reduce branch penalties significantly

Required hardware support:

Prediction structures:

- <u>Branch history tables (BHT), branch target buffers(BTB),</u> etc.
 Misprediction recovery mechanisms:
 - Keep result computation separate from commit
 - Kill instructions following branch in pipeline (a.k.a. flush)
 - Restore state to that following branch

Static Branch Prediction

Overall probability a branch is taken is ~60-70% but:



ISA can attach preferred direction semantics to branches, e.g., Motorola MC88110

- bne0 (preferred taken) beq0 (not taken)
- typically reported as ~80% accurate

Speculating Both Directions?

- An alternative to branch prediction is to execute both directions of a branch speculatively (Predicated Execution)
 - Resource requirement is proportional to the number of concurrent speculative executions
 - Only half the resources engage in useful work when both directions of a branch are executed speculatively



With accurate branch prediction, it is more cost effective to dedicate all resources to the predicted direction!

Dynamic Branch Prediction Learning based on past behavior

- Branch behavior is monitored during program execution
 - History data can influence prediction of future execution of the branch instruction
- Temporal correlation
 - The way a branch resolves may be a good predictor of the way it will resolve at the next execution

- Spatial correlation
 - Several branches may resolve in a highly correlated manner (a preferred path of execution)

Prediction based on Instruction



- Fast decode instructions belonging to one fetch group, only detect whether it is branch or not ?
- Some design doing pre-decode when moving Instructions from L2\$ to L11\$ (with Branch Flag)

Prediction based on PC



- The Physical Address of an instruction of a program is changeable (Depends on OS Page Scheduler)
- The <u>Virtual Address of an instruction of a program is fixed after</u> <u>compilation</u>
- ASID is used together with PC for prediction (Otherwise, Branch predictor need to be clear after context switch)

One-Bit Branch History Predictor

- For each branch, remember last way branch went
- Has problem with loop-closing backward branches, as two mispredictions occur on every loop execution
 - first iteration predicts loop backwards branch not-taken (loop was exited last time)
 - 2. last iteration predicts loop backwards branch taken (loop continued last time)



1-bit Saturating Counter

1-Bit Saturating Counter



What happens on loop branches ? At best, two mispredictions for every use of loop

2-Bit Branch Prediction

- Assume 2 BP bits per instruction instead of one
- Change the prediction after *two consecutive mistakes*!



BP state:(predict take/¬take) x (last prediction right/wrong)

Pattern History Table (PHT) PHT 2-bit saturating Fetch PC 0,0 counter I-Cache Hash Update Instruction **FSM** offset Opcode Update Logic **Mispredict?** Predict Taken/¬Taken? **Target PC** Branch?

- 4K-entry PHT (*Why the low bits of PC?*), 2 bits/entry (~80-90% correct predictions)
- PHT is updated at the **commit stage**

Prediction Based on Local Branch History



- Local branch history: latest branch results of one branch instruction
- In many codes, loops may have fixed iteration period
- Use local history buffer to capture loop patterns, need to be larger than iteration period

Prediction Based on Local Branch History



- For each branch instruction, record its latest K branch results
- For total 2^K possible history record, put one 2-bit saturating counter for each case (into a PHT)
- Prediction reach **100%** if K > Loop Period
- **Too expensive** to hold all branch instructions!

Two-Level Branch Predictor



•

Two-Level Branch Predictor (Hash)



- Avoid conflict of **different branches** with the **same history**
- Less hardware resource requirement

Exploiting Spatial Correlation *Yeh and Patt, 1992*

- If **first** condition false, **second** condition also false
- Prediction must be a function of own branch as well as recent outcomes of other branches.

Global branch history: the direction of the last N branches executed by the processor

Two-Level Branch Predictor



Two-Level Branch Predictor (Hash) (Global history based prediction)



- All branch instructions share one PHT
- Can blend branch PC with GHR bits to avoid conflicts

Limitations of BHT/PHTs

- If predict branch **taken**, need to know **WHERE TO JUMP**!
- BHT/PHT only predicts branch direction. Therefore, cannot redirect fetch stream until after branch target is determined.



UltraSPARC-III fetch pipeline

Branch Target Buffer (BTB)



- Keep both the branch PC and **target PC** in the BTB
- Only taken branches and jumps held in BTB (to save space)
- If PC doesn't exist in BTB: PC+4 is fetched (continue as non-taken)
- Next PC determined *before* the branch instruction is fetched
 Useful for both absolute branch and PC-relative branch

Combining BTB and BHT

- BTB entries are considerably more expensive than BHT, but can redirect fetches at earlier stage in pipeline and can accelerate indirect branches (JR)
- BHT can hold many more entries and is more accurate



BTB/BHT only updated after branch resolves in E stage

Absolute Address Jump/Branch with uses of Jump Register (JALR)

- Switch statements (jump to address of matching case)
 BTB works well if same case used repeatedly
- Dynamic function call (jump to run-time function address)

BTB works well if same function usually called, (e.g., in C++ programming, when objects have same type in virtual function call)

 Subroutine returns (jump to return address) BTB works well if usually return to the same place
 Often one function called from many distinct call sites! How well does BTB work for each of these cases?

Absolute Address Jump/Branch (Function calls and returns)



- CALL for some shared subroutine
- The Target address of call is fixed
- Can be solved using BTB

- Return from callee to **different callers**
- The target(return) address is unfixed
- Can't be solved using BTB

Subroutine Return Address Stack (RAS)

Small structure to accelerate JR for subroutine returns, typically much more accurate than BTBs.



Absolute Branch Prediction with RAS



Decode Call and Return Instructions

- Subroutine Call and Return instructions are usually pseudo instructions
- Automatically translated into fixed-format real instructions by compiler

RISC-V Pseudo Call and Return

| Pseudo Code | Instructions | Description |
|-------------|---|--------------------------------|
| CALL | <pre>auipc x1,offset[31:12] jair x0,x1,offset[11:0]</pre> | Call far-away subroutine |
| RET | jalr x0, x1, 0 | Return from subroutine call |

 Decoder recognizes the fixed encoding to find subroutine Call and Return branch instructions.

Return Address Stack in Pipeline

- How to use return address stack (RAS) in deep fetch pipeline?
- Only know if subroutine call/return at decode stage



Return Address Stack prediction checked

Return Address Stack in Pipeline

- Can remember whether PC is call/return in extended
 BTB structure
- Instead of target-PC, just store push/pop bit

PC Generation/Mu RAS Ρ Instruction Fetch Stage 1 Push/Pop before instructions decoded! F Instruction Fetch Stage 2 Branch Address Calc/Begin Decode B **Complete Decode** Steer Instructions to Functional units R **Register File Read** Integer Execute E

Return Address Stack prediction checked

Branch Prediction for CALL/Return



If RAS is **full**, drop the new one? Or the oldest?

- Drop the oldest in a First-in-First-out fashion (FIFO)
- Use it in a Last-in-First-out fashion (Stack)

Absolute Address Jump/Branch

(Switch case, function pointer array, ...)

Switch Case

Function Pointer Array



Absolute Branch Prediction with Target Cache



- Use local branch history to help predict target address
- Use hash to combine branch PC and branch history
- Like BHT, but stores target address instead of direction₄₀

Temporal and Spatial Correlation

If *first condition* false, *second condition* also false.

How about the *first condition* is true, does this help make the prediction of the *second 'if"* statement ?

No!

Prediction must be a function of **own branch** as well as recent outcomes of **other branches**.

History register records the direction of the last N branches executed by the processor

Local branch history

Global branch history

Tournament Branch Predictor (Alpha 21264)

- Choice predictor learns whether best to use local or global branch history in predicting next branch (Different schemes work better for different branches)
- Global history is speculatively updated but restored on mispredict
- Claim 90-100% success on range of applications





2-Bit Choice Prediction



Tournament Predictors – Combine approaches (local and global)

In-Order vs. Out-of-Order Branch Prediction



- Speculative fetch but not speculative execution - branch resolves before later instructions complete
- Completed values held in bypass network until commit

- Speculative execution, with branches resolved after later instructions complete
- Completed values held in rename registers in ROB or unified physical register file until commit

Both styles of machine can use **same branch predictors** in front-end fetch pipeline, and both can execute multiple instructions per cycle

Common to have **10-30 pipeline stages** in either style of design

InO vs. OoO Mispredict Recovery

- In-order execution?
 - Design so no instruction issued after branch can write-back before branch resolves
 - Kill all instructions in pipeline behind mispredicted branch
- Out-of-order execution?
 - Multiple instructions following branch in program order can complete before branch resolves
 - A simple solution would be to handle like **precise exception**
 - Kill following instructions in pipeline and ROB
 - **Restore rename register** status to the branch time point
 - How about **multiple branches** are encountered?

Branch Misprediction in OoO Pipeline



- Can have multiple unresolved branches in ROB
- Can resolve branches out-of-order by killing all the instructions in ROB that follow a mispredicted branch
- Use mask bits to tag instructions that are dependent on different speculative branches
- Mask bits cleared as branch resolves, and reused for next branch

Rename Table Recovery

- Have to quickly recover rename table on branch mispredictions
- MIPS R10K only has four snapshots for each of four outstanding speculative branches
- Alpha 21264 has 80 snapshots, one per ROB instruction



Load-Store Queue Design

- After control hazards, data hazards through memory are probably next most important bottleneck to superscalar performance
- Modern superscalars use very sophisticated load-store reordering techniques to reduce effective memory latency by allowing loads and stores to be speculatively issued
 - Speculative load \rightarrow ROB and physical register (rename)
 - Speculative store \rightarrow Store Buffer

Recap: I2OI (OoO superscalar)



Speculative Store Buffer



L1 Data Cache

- Just like register updates, stores should not modify the memory until after the instruction is committed. A speculative store buffer is a structure introduced to hold speculative store data.
- During decode, store buffer slot allocated in program order
- Stores split into "store address" and "store data" micro-operations
- "Store address" execution writes tag
- "Store data" execution writes data
- Store commits when oldest instruction and both address and data available:
 - clear speculative bit and eventually move data to cache
 - On store abort:, clear valid bit

Load bypass from speculative store buffer



- If data in both store buffer and cache, which should we use?
 Speculative store buffer
- If same address in store buffer twice, which should we use?
 Youngest store older than load

Summary: a Full Branch Prediction method



Next Lecture : Very Long Instruction Words (Data Level Parallel)