#### 332 Advanced Computer Architecture Chapter 4

#### Part 2: Branch Target Prediction

October 2022 Paul H J Kelly

These lecture notes are partly based on the course text, Hennessy and Patterson's Computer Architecture, a quantitative approach (4-6<sup>th</sup> eds), and on the lecture slides of David Patterson's Berkeley course (CS252)

Course materials online on https://scientia.doc.ic.ac.uk/2223/modules/60001/materials and https://www.doc.ic.ac.uk/~phjk/AdvancedCompArchitecture/aca20/

#### **Branch Prediction - context**

• If we have a branch predictor....

- We want to fetch the correct (predicted) next instruction without any stalls
- We need the prediction before the preceding instruction has been decoded
- We need to predict conditional branches
  - Direction prediction
- And indirect branches
  - Target prediction

### **Branch Target Buffer**

- Need address at same time as prediction
- Especially for indirect branches and virtual method calls
- Note that we must check for branch match, since can't use wrong branch address



# Branch target prediction: BTBs

- re: "In order to predict a branch, we need to know that current instruction is branch instruction"
- This doesn't have to be true!



In parallel with every ifetch

Check whether the BTB predicts that the instruction we are fetching *will* be a taken branch

# Branch target prediction: BTBs

- re: "In order to predict a branch, we need to know that current instruction is branch instruction"
  - In parallel with every ifetch Control Hazard on Branches Check whether the BTB predicts that the instruction we are fetching will be a taken branch 10: beg r1,r3,36 When a **taken** branch is 14: and r2,r3,r5 committed, we update the BTB with the branch's 18: or r6,r1,r7 target address (and with the tag of the address of 22: add r8,r1,r9 the branch instruction). 36: xor r10.r1.r11 If we're not smart we risk a three-cycle stall
- This doesn't have to be true!

# Branch Target Buffer (BTB)

- Cache of branch target addresses accessed in parallel with the I-cache in the fetch stage
- Updated only by taken branches (the direction-predictor determines whether BTB is used)
- If BTB hit and the instruction is a predicted-taken branch
  - target from the BTB 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



**BTB** is

*indexed* with low-order PC address bits,

*tagged* with high-order bits

(Note: we could use an n-way set-associative design here)

*Credit: Onur Mutlu* 

#### Target prediction: recall the 5-stage MIPS pipeline



Where does branch prediction happen?



#### **Combining BTB with direction Prediction**



Credit: Onur Mutlu

#### Combining fast simple predictor with slower bigger predictor



(What if branch is predicted-taken but BTB miss?)

Where does branch prediction happen?



## Updating the branch prediction





- A function might be called from different places
- In each case it must return to the right place
- Address of next instruction must be saved and restored

Hennessy and Patterson 6<sup>th</sup> ed p232-233



- jsr must save return address somewhere
- On x86 jsr pushes return address onto stack
- ret jumps to the address on the top of the stack
- On MIPS, "jal F" (jump-and-link) jumps to F, and stashes the current PC in a special register \$ra.
- Function returns with an indirect jump "jr \$ra"
- If the function body has other calls, compiler must push \$ra to the stack



- Return addresses form a stack (even if they are stored in registers)
- They should be easy to predict!
- We need to add *another* branch target predictor
- That maintains a hardware stack of return addresses
- Presumably a small stack



#### **Return Address Predictor - mispredictions**





Fixed small number of entries

- What happens if the call stack is deeper than the RAP's stack?
  - On return, the RAP's stack will be empty!
- Why might the prediction from the RAP be wrong?
  - Maybe the return address was overwritten
  - Maybe the stack pointer was changed
  - Maybe because we switched to another thread



- Q: when should the RAS be updated?
- The BTB is updated when a branch is *committed*
- But if we wait for commit to update the RAS, we might not have a prediction for the return from H
- Or: if we mispredict that the conditional "IF(C)" is true

   We might have the wrong RAS prediction for the return from G

### Branch prediction and multi-issue

- In a processor that fetches, issues and dispatches multiple instructions per cycle.....
- What if we encounter two (or more) branches in one issue "packet"?

### Branch prediction and multi-issue

- In a processor that fetches, issues and dispatches multiple instructions per cycle.....
- What if we encounter two (or more) branches in one issue "packet"?
- But all the BTB needs is to predict the next instruction to fetch – it doesn't matter which branch is responsible
- Commonly, a bigger slower branch predictor may later *re-steer* the processor if it has a better prediction that should over-ride the BTB

### **Dynamic Branch Prediction Summary**

- Prediction seems essential (?)
- Two questions: branch *takenness*, branch *target*

#### Takenness:

- Branch History Table: 2 bits for loop accuracy
  - Saturating counter (bimodal) scheme handles highly-biased branches well
  - Some applications have highly dynamic branches
- Correlation: Recently executed branches correlated with next branch.
  - Either different branches
  - Or different executions of same branches
- Tournament Predictor: try two or more competitive solutions and pick between them
- Predicated Execution can reduce number of branches, number of mispredicted branches

#### Target:

- Branch Target Buffer: include branch address & prediction
- BTB update
- Return address stack for prediction of indirect jump

#### Beyond:

- Prediction mechanisms have many applications beyond branch prediction:
  - Way prediction, prefetching, store-to-load forwarding, value prediction, etc
    - George Z. Chrysos and Joel S. Emer. 1998. Memory dependence prediction using store sets. ISCA98
  - Predictors can increase performance, but make it *harder* to optimize programs



# Branch prediction resources

- Design tradeoffs for the Alpha EV8 Conditional Branch Predictor (André Seznec, Stephen Felix, Venkata Krishnan, Yiannakis Sazeides)
  - SMT: 4 threads, wide-issue superscalar processor, 8-way issue, 512 registers (cancelled June 2001 when Alpha dropped)
  - Paper: <u>http://citeseer.ist.psu.edu/seznec02design.html</u>
  - Talk: <u>http://ce.et.tudelft.nl/cecoll/slides/PresDelft0803.ppt</u>
- Branch prediction in the Pentium family (Agner Fog)
  - Reverse engineering Pentium branch predictors using direct access to BTB
  - <u>http://www.x86.org/articles/branch/branchprediction.htm</u>
- Championship Branch Prediction Competition (CBP), organised by the Journal of Instruction-level Parallelism

<u>http://www.jilp.org/cbp/</u>

 The CBP-1 winning entry: TAgged GEometric history length predictor (TAGE): for each branch, maintain a predictor for what history length (from a geometric progression) works best.

<u>http://www.irisa.fr/caps/people/seznec/JILP-COTTAGE.pdf</u>

#### Example: Branch prediction in Intel Atom, Silvermont and Knights Landing

- two-level adaptive predictor with a global history table,
- Branch history register has 12 bits
- The pattern history table on the Atom has 4096 entries and is shared between threads
- The branch target buffer has 128 entries, organized as 4 ways by 32 sets
  - (size on Silvermont unknown, but probably bigger, and not shared between threads)
- Unconditional jumps make no entry in the global history table, but always-taken and nevertaken branches do
- Silvermont has branch prediction both at the fetch stage and at the later decode stage in the pipeline, where the latter can correct errors in the former
- No special predictor for loops (as there is for some other Intel CPUs)
  - Loops are predicted in the same way as other branches
- Penalty for mispredicting a branch is 11-13 clock cycles.
- It often occurs that a branch has a correct entry in the pattern history table, but no entry in the branch target buffer, which is much smaller:
  - If a branch is correctly predicted as taken, but no target can be predicted because of a missing BTB entry, then the penalty will be approximately 7 clock cycles.
- Pattern prediction evident for indirect branches on Knights Landing but not on Silvermont.
  - Indirect branches are predicted to go to the same target as last time on Silvermont
- Return stack buffer with 8 entries on the Atom and 16 entries on Silvermont and Knights Landing

"The microarchitecture of Intel, AMD and VIA CPUs An optimization guide for assembly programmers and compiler makers" http://www.agner.org/optimize/microarchitecture.pdf

# Piazza question: better predictions for indirect branches

- As you say, a BTB should give you a prediction for an indirect branch.
- However it might not be a very good one the killer app is polymorphic calls in object-oriented languages (virtual calls where the target object has a different type on different invocations).
- For that we need to add global history to the branch target prediction. We did not cover this in the lectures.
- This paper evaluates three alternative schemes:
- Dharmawan, Tubagus & Jeyachandra, E & Rahmadhani, Andri. (2016). Techniques to Improve Indirect Branch Prediction. 10.13140/RG.2.2.24350.02884.
- The state of the art is perhaps represented by this article in the same ISCA2020 "Industry" track:
- <u>The IBM z15 High Frequency Mainframe Branch Predictor (computer.org)</u> (section VI], pg 35-6). Basically they use the branch history to index a special BTB (actually they expand the branch history concept to include a couple of bits of the PC address of each taken branch in the history).