Dynamic Branch Prediction

Hung-Wei Tseng



Recap: Pipelining



Recap: Pipelining

add x1, x2, x3 ld x4, 0(x5) sub x6, x7, x8 sub x9, x10, x11 sd x1, 0(x12) xor x13, x14, x15 and x16, x17, x18 add x19, x20, x21 sub x22, x23, x24 ld x25, 4(x26) sd x27, 0(x28)

					_		
F	ID	EX	MEM	WB			
	IF	ID	EX	MEM	WB		
		IF	ID	EX	MEM	WB	
			IF	ID	EX	MEM	WB
				IF	ID	EX	MEM
					IF	ID	EX
						IF	ID
							IF
				Afte we inst	er thi are c ructi	s poi omp on ea	nt, letinț ach c



Recap: Three pipeline hazards

- Structural hazards resource conflicts cannot support simultaneous execution of instructions in the pipeline
- Control hazards the PC can be changed by an instruction in the pipeline
- Data hazards an instruction depending on a the result that's not yet generated or propagated when the instruction needs that



Recap: Tips of drawing a pipeline diagram

- Each instruction has to go through all 5 pipeline stages: IF, ID, EXE, MEM, WB in order — only valid if it's single-issue, RISC-V 5-stage pipeline
- An instruction can enter the next pipeline stage in the next cycle if
 - No other instruction is occupying the next stage
 - This instruction has completed its own work in the current stage
 - The next stage has all its inputs ready and it can retrieve those inputs
- Fetch a new instruction only if
 - We know the next PC to fetch
 - We can predict the next PC
 - Flush an instruction if the branch resolution says it's mis-predicted.

Recap: Solving Structural Hazards

- Stall can address the issue but slow
- Improve the pipeline unit design to allow parallel execution



Recap: The impact of control hazards

 Assuming that we have an application with 20% of branch instructions and the instruction stream incurs no data hazards. When there is a branch, we disable the instruction fetch and insert no-ops until we can determine the PC. What's the average CPI if we execute this program on the 5-stage RISC-V pipeline?

A. 1 B. 1.2 C. 1.4 D. 1.6 E. 1.8

add x1,	x2, x3	IF	ID	EX	MEM	WB						
ld x4,	0(x5)		IF	ID	EX	MEM	WB					
bne x0	x7, L			IF	ID	EX	MEM	WB				
add x0,	x0, x0				IF	ID	EX	MEM	WB			
add x0,	x0, x0					IF	ID	EX	MEM	WB		
sub x9	x10,x11						IF	ID	EX	MEM	WB	
cd v1	0(y12)							IF	ID	EX	MEM	W

Outline

- 2-bit predictor
- 2-level global predictor
- Hybrid predictors
- Perceptrons
- Branch and coding

Why can't we proceed without stalls/no-ops?

- How many of the following statements are true regarding why we have to stall for each branch in the current pipeline processor
 - The target address when branch is taken is not available for instruction fetch stage of the next cycleYou need a cheatsheet for that — branch target buffer
 - ② The target address when branch is not-taken is not available for instruction fetch
 - stage of the next cycle. You need to predict that history/states The branch outcome cannot be decided until the comparison result of ALU is not out
 - The next instruction needs the branch instruction to write back its result (4)
 - A. 0
 - B. 1
 - - D. 3

E. 4

A basic dynamic branch predictor





2-bit/Bimodal local predictor

- Local predictor every branch instruction has its own state
- 2-bit each state is described using 2 bits
- Change the state based on actual outcome
- If we guess right no penalty
- If we guess wrong flush (clear pipeline registers) for mis-predicted instructions that are currently in IF and ID stages and reset the PC **(**)



	branch PC	target PC	Stat
	0x400048	0x400032	10
Predict Taken	0x400080	0x400068	11
	0x401080	0x401100	00
	0x4000F8	0x400100	01





2-bit local predictor



predict state actual 10 Т Т 11 Т Т 11 Т Т 11 Т Т Т NT 11

90% accuracy! $CPI_{average} = 1 + 20\% \times 10\% \times 2 = 1.04$

https://www.pollev.com/hungweitseng close in 1:30

2-bit local predictor

• What's the overall branch prediction (include both branches) accuracy for this nested for loop?

```
i = 0;
do {
    if( i % 2 != 0) // Branch X, taken if i % 2 == 0
       a[i] *= 2;
    a[i] += i;
} while ( ++i < 100)// Branch Y</pre>
```

(assume all states started with 00)

A. ~25% B. ~33% C. ~50% D. ~67% E. ~75%

2-bit local

А



D

С

2-bit local predictor

• What's the overall branch prediction (include both branches) accuracy for this nested for loop?

```
i = 0;
                                                            b
do {
    if( i % 2 != 0) // Br Can We if i % P a[i] *= 2;
    a[i] += i;
} while ( ++i < 100)// Branch Vetter job?</pre>
                                                          3
(assume all states started with 00)
                                                          3
  A. ~25%
                                                          4
  B. ~33%
                                                          4
  C. ~50%
                                                          5
  D. ~67%
                                                          5
                                 For branch Y, almost 100%,
                                                          6
                                 For branch X, only 50%
    ~75%
```

6

ranch?	state	prediction	actual
Х	00	NT	Т
Y	00	NT	Т
Х	01	NT	NT
Y	01	NT	Т
Х	00	NT	Т
Y	10	Т	Т
Х	01	NT	NT
Y	11	Т	Т
Х	00	NT	Т
Y	11	Т	Т
Х	01	NT	NT
Y	11	Т	Т
Х	00	NT	Т
Y	11	Т	Т

Two-level global predictor

Marius Evers, Sanjay J. Patel, Robert S. Chappell, and Yale N. Patt. 1998. An analysis of correlation and predictability: what makes two-level branch predictors work. In Proceedings of the 25th annual international symposium on Computer architecture (ISCA '98).

2-bit local predictor

• What's the overall branch prediction (include both branches) accuracy for this nested for loop?

```
i = 0;
                                                    i b
do {
   if( i % 2 != 0) // Branch X, taken if i % 2 == 0
      a[i] *= 2;
                         This pattern
   a[i] += i;
} while ( ++i < 100)// Branch</pre>
(assume all states state peats all the tin
  A. ~25%
                                                    4
  B. ~33%
                                                    4
  C. ~50%
                                                    5
  D. ~67%
                                                    5
                             For branch Y, almost 100%,
                                                    6
                             For branch X, only 50%
  E. ~75%
                                                    6
```

ranch?	state	prediction	actual
Х	00	NT	Т
Y	00	NT	Т
Х	01	NT	NT
Υ	01	NT	Т
Х	00	NT	Т
ne	10	Т	Т
X	01	NT	NT
Y	11	Т	Т
Х	00	NT	Т
Y	11	Т	Т
Х	01	NT	NT
Y	11	Т	Т
Х	00	NT	Т
Y	11	Т	Т







Performance of GH predictor



Near perfect after this

i	branch?	GHR	state	prediction	actual
0	Х	000	00	NT	Т
1	Y	001	00	NT	Т
1	Х	011	00	NT	NT
2	Y	110	00	NT	т
2	Х	101	00	NT	Т
3	Y	011	00	NT	Т
3	Х	111	00	NT	NT
4	Y	110	01	NT	Т
4	Х	101	01	NT	Т
5	Y	011	01	NT	Т
5	Х	111	00	NT	NT
6	Y	110	10	Т	Т
6	Х	101	10	Т	Т
7	Y	011	10	Т	Т
7	Х	111	00	NT	NT
8	Y	110	11	Т	Т
8	Х	101	11	Т	Т
9	Y	011	11	Т	Т
9	Х	111	00	NT	NT
10	Y	110	11	Т	Т
10	Х	101	11	Т	Т
11	Y	011	11	Т	Т

Better predictor?

 Consider two predictors — (L) 2-bit local predictor with unlimited BTB entries and (G) 4-bit global history with 2-bit predictors. How many of the following code snippet would allow (G) to outperform (L)?



https://www.pollev.com/hungweitseng close in 1:30

i = 0;do { if(rand() %2 == 0) a[i] *= 2; a[i] += i; } while (++i < 100)</pre>

Global v.s. local

А

С

D

Ε

Better predictor?

 Consider two predictors — (L) 2-bit local predictor with unlimited BTB entries and (G) 4-bit global history with 2-bit predictors. How many of the following code snippet would allow (G) to outperform (L)? about the same about the same \Box_{i} i do { i = 0;do { = 0; do { if(i % 10 != 0) do { a[i] *= 2; a[i] += i; sum += A[i*2+j]; a[i] += i; } while (++i < 100); while (++i < 100); while(++j < 2);</pre> while (++i < 100); A. 0 Β. C. 2 D. 3 50% if you use L E. 4

i**L_could be better** do { if(rand() %2 == 0) a[i/] *= 2;

This branch mess up the global history!

while (++i < 100)

Hybrid predictors

gshare predictor







gshare predictor

 Allowing the predictor to identify both branch address but also use global history for more accurate prediction



Local History Predictor branch PC local history

x400048	1000
x400080	0110
x401080	1010
x4000F8	0110

Predict Taken

Tournament Predictor

- The state predicts "which predictor is better"
 - Local history
 - Global history
- The predicted predictor makes the prediction



TAGE

André Seznec. The L-TAGE branch predictor. Journal of Instruction Level Parallelism (http:// wwwjilp.org/vol9), May 2007.

Better predictor?

 Consider two predictors — (L) 2-bit local predictor with unlimited BTB entries and (G) 4-bit global history with 2-bit predictors. How many of the following code snippet would allow (G) to outperform (L)? about the same about the same G: i do { = 0; do { 1 = 0; do { if(i % 10 != 0) do { a[i] *= 2; a[i] += i; sum += A[i*2+j]; } while (++i < 100);</pre> a[i] += i; while (++i < 100): while(++j < 2); while (++i < 100); A. 0 different branch needs different length of history Β. C. 2 global predictor can work if the history is long enough! D. 3 E. 4

i**L_could be better** do { if(rand() %2 == 0) a[i] *= 2; a[i] += i; while (++i < 100)



Perceptron

Jiménez, Daniel, and Calvin Lin. "Dynamic branch prediction with perceptrons." Proceedings HPCA Seventh International Symposium on High-Performance Computer Architecture. IEEE, 2001.

The following slides are excerpted from https://www.jilp.org/cbp/Daniel-slides.PDF by Daniel Jiménez 35

Branch Prediction is Essentially an ML Problem

- The machine learns to predict conditional branches
- Artificial neural networks
 - Simple model of neural networks in brain cells
 - Learn to recognize and classify patterns

Mapping Branch Prediction to NN

- The inputs to the perceptron are branch outcome histories
 - Just like in 2-level adaptive branch prediction
 - Can be global or local (per-branch) or both (alloyed)
 - Conceptually, branch outcomes are represented as
 - +1, for taken
 - -1, for not taken
- The output of the perceptron is
 - Non-negative, if the branch is predicted taken
 - Negative, if the branch is predicted not taken
- Ideally, each static branch is allocated its own perceptron



Mapping Branch Prediction to NN (cont.)

- Inputs (x's) are from branch history and are -1 or +1
- n + 1 small integer weights (w's) learned by on-line training
- Output (y) is dot product of x's and w's; predict taken if y
- Training finds correlations between history and outcome





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

Training Algorithm $x_{1..n}$ is the *n*-bit history register, x_0 is 1. $w_{0..n}$ is the weights vector. t is the Boolean branch outcome. θ is the training threshold. if $|y| \leq \theta$ or $((y \geq 0) \neq t)$ then for each $0 \le i \le n$ in parallel if $t = x_i$ then $w_i := w_i + 1$ else $w_i := w_i - 1$ end if end for end if



Advanced Dynamic Predictors

- Which of the following predictor works the best when the processor has very limited hardware budget (e.g., 1K) for a branch predictor
 - A. gshare
 - B. bi-mode (2-bit local)
 - C. Tournament predictor
 - D. Perceptrons
 - E. TAGE

https://www.pollev.com/hungweitseng close in 1:30

Advanced Branch Predictors

А

В

Such the presentation to see live-content. For screen share software, share the entire screen. Get help at policy.com/spc

С

D

Ε

Advanced Dynamic Predictors

- Which of the following predictor works the best when the processor has very limited hardware budget (e.g., 1K) for a branch predictor
 - A. gshare
 - B. bi-mode (2-bit local)
 - C. Tournament predictor
 - D. Perceptrons
 - E. TAGE



How good is prediction using perceptrons?







How good is prediction using perceptrons?





Figure 4: Misprediction Rates at a 4K budget. The perceptron predictor has a lower misprediction rate than gshare for all benchmarks except for 186. crafty and 197. parser.

Figure 5: Misprediction Rates at a 16K budget. Gshare outperforms the perceptron predictor only on 186.crafty. The hybrid predictor is consistently better than the PHT schemes.



History/training for perceptrons



Hardware budget		ngth	
in kilobytes	gshare	bi-mode	perceptron
1	6	7	12
2	8	9	22
4	8	11	28
8	11	13	34
16	14	14	36
32	15	15	59
64	15	16	59
128	16	17	62
256	17	17	62
512	18	19	62

Table 1: Best History Lengths. This table shows the best amount of global history to keep for each of the branch prediction schemes.

Branch predictors in processors

- The Intel Pentium MMX, Pentium II, and Pentium III have local branch predictors with a local 4-bit history and a local pattern history table with 16 entries for each conditional jump.
- Global branch prediction is used in Intel Pentium M, Core, Core 2, and Silvermont-based Atom processors.
- Tournament predictor is used in DEC Alpha, AMD Athlon processors
- The AMD Ryzen multi-core processor's Infinity Fabric and the Samsung Exynos processor include a perceptron based neural branch predictor.



Branch and programming

Demo revisited

• Why the sorting the array speed up the code despite the increased instruction count?

```
if(option)
    std::sort(data, data + arraySize);
for (unsigned i = 0; i < 100000; ++i) {</pre>
    int threshold = std::rand();
    for (unsigned i = 0; i < arraySize; ++i) {</pre>
        if (data[i] >= threshold)
             sum ++;
    }
}
```

https://www.pollev.com/hungweitseng close in 1:30 Demo revisited

- Why the performance is better when option is not "0"
 - ① The amount of dynamic instructions needs to execute is a lot smaller
 - ② The amount of branch instructions to execute is smaller
 - ③ The amount of branch mis-predictions is smaller
 - ④ The amount of data accesses is smaller

```
A. 0 if(option)
    std::sort(data, data + arraySize);
B. 1
C. 2 for (unsigned i = 0; i < 100000; ++i) {
    int threshold = std::rand();
D. 3    for (unsigned i = 0; i < arraySize; ++i) {
        if (data[i] >= threshold)
        Sum ++;
        }
```

ot "O" ecute is a lot smaller smaller

Why sort

А

D

С

Demo revisited

- Why the performance is better when option is not "0"
 - ① The amount of dynamic instructions needs to execute is a lot smaller
 - ② The amount of branch instructions to execute is smaller
 - The amount of branch mis-predictions is smaller
 - ④ The amount of data accesses is smaller
 - A. O if(option)

B.

std::sort(data, data + arraySize);

ot "O" ecute is a lot smaller smaller

	Without sorting	With sorting
The ediction curacy of before reshold	50%	100%
The ediction curacy of X after reshold	50%	100%

Demo: Popcount

- The population count (or popcount) of a specific value is the number of set bits (i.e., bits in 1s) in that value.
- Applications
 - Parity bits in error correction/detection code
 - Cryptography
 - Sparse matrix
 - Molecular Fingerprinting
 - Implementation of some succinct data structures like bit vectors and wavelet trees.

Demo: pop count

- Given a 64-bit integer number, find the number of 1s in its binary representation.
- Example 1: Input: 9487 Output: 7 Explanation: 9487's binary representation is Ob10010100001111

```
int main(int argc, char *argv[]) {
     uint64_t key = 0xdeadbeef;
     int count = 1000000000;
     uint64_t sum = 0;
     for (int i=0; i < count; i++)</pre>
     printf("Result: %lu\n", sum);
     return sum;
```

}

- sum += popcount(RandLFSR(key));

https://www.pollev.com/hungweitseng close in 1:30

Four implementations

• Which of the following implementations will perform the best on modern pipeline processors? inline int popcount(uint64_t x) {





inline int popcount(uint64_t x) { int c = 0;int table $[16] = \{0, 1, 1, 2, 1,$ 2, 2, 3, 1, 2, 2, 3, 2, 3, 3, 4; for (uint64_t i = 0; i < 16; i++) c += table[(x & 0xF)];x = x >> 4;

return c;

Four implementations

• Which of the following implementations will perform the best on modern pipeline processors? inline int popcount(uint64_t x) {

```
int c = 0;
   inline int popcount(uint64_t x){
                                                   while(x)
                                                               {
     int c=0;
                                                     c += x \& 1;
     while(x) {
                                                     x = x >> 1;
           c += x \& 1;
                                                     c += x & 1;
           x = x >> 1;
                                                     x = x >> 1;
        }
                                            \mathbf{m}
                                                     c += x & 1;
       return c;
                                                     x = x >> 1;
   }
                                                     c += x \& 1;
                                                     x = x >> 1;
                                                                      in
   inline int popcount(uint64_t x) {
         int c = 0;
                                                   return c;
         int table [16] = \{0, 1, 1, 2, 1, 
                                               }
   2, 2, 3, 1, 2, 2, 3, 2, 3, 3, 4;
         while(x)
\mathbf{O}
                                                                   c += table[(x \& 0xF)];
             x = x >> 4;
         }
         return c;
   }
                                                                        }
                                                      63
```



fe int popcount(uint64_t x) { int c = 0;int table $[16] = \{0, 1, 1, 2, 1,$ 2, 2, 3, 1, 2, 2, 3, 2, 3, 3, 4; for (uint64_t i = 0; i < 16; i++) c += table[(x & 0xF)];x = x >> 4;

return c;

Announcements

- Assignment #3 due next Monday
- Reading Quiz due this Wednesday
- Project is released
 - Please check website to the link of GitHub repo
 - You may discuss, but each needs an individual/distinguishable version of code
 - You need to write a brief report
 - Grading rubrics
 - $\cdot 20\%$ report
 - 20% if you code can compile and run
 - 60% performance based. The sample prefetcher is the baseline. We calculate your score at this part using min(Speedup-1, 1). If you can speedup by 2, you score full credits in this part
 - Due 11/29 no extension

Computer Science & Engineering





