Lecture 21: BPred, OOO, Memory Hierarchy





	SIGN1: Look	at Branch PC
Bimodal Pre	dictor	Content of the second s
Code (Sign R=BEA Fr PC= 325 Front e.d i Capture the situation with a 14b number	the (situation) of end 14 bits 3ranch PC 14 bit 2,527 Number a 14 b signature (0 - 16,383) [1	Table descreations Backed 3 - 20 1 - 20 1 - 20 1 - 20 1 - 20 3 - 20 1 - 20 3 - 20

Local, Global, Tournament 146 hot 3 146 sign 146 PC **Bimodal Predictor** OPTZ's combo of PC+ History Concat, A XOR 7146 PC =T14b sign Create better 14bftish 14 bits Signame Table of **Branch PC 16K** entries OFTO: Pr of 2-bit OPT 1: Remember what the last saturating counters 14 branches did NT/O 11011011011011 Br History Rie 10101101101104

2-Bit Prediction

- For each branch, maintain a 2-bit saturating counter: if the branch is taken: counter = min(3,counter+1) if the branch is not taken: counter = max(0,counter-1) ... sound familiar?
- If (counter >= 2), predict taken, else predict not taken
- The counter attempts to capture the common case for each branch

Indexing functions Multiple branch predictors History, trade-offs

An Out-of-Order Processor Implementation

Issue Queue (IQ)

An Out-of-Order Processor Implementation

Issue Queue (IQ)

Example Code

Completion times	with in-order	with ooo
ADD R1, R2, R3	5	5
ADD R4, R1, R2	6	6
LW R5, 8(R4)	7	7
ADD R7, R6, R5	9	9
ADD R8, R7, R5	10	10
LW R9, 16(R4)	11	7
ADD R10, R6, R9	13	9
ADD R11, R10, R9	14	10

Cache Hierarchies

- Data and instructions are stored on DRAM chips DRAM is a technology that has high bit density, but relatively poor latency – an access to data in memory can take as many as 300 cycles today!
- Hence, some data is stored on the processor in a structure called the cache – caches employ SRAM technology, which is faster, but has lower bit density
- Internet browsers also cache web pages same concept

Memory Hierarchy

• As you go further, capacity and latency increase

Locality

- Why do caches work?
 - Temporal locality: if you used some data recently, you will likely use it again
 - Spatial locality: if you used some data recently, you will likely access its neighbors
- No hierarchy: average access time for data = 300 cycles
- 32KB 1-cycle L1 cache that has a hit rate of 95%: average access time = 0.95 x 1 + 0.05 x (301) = 16 cycles