

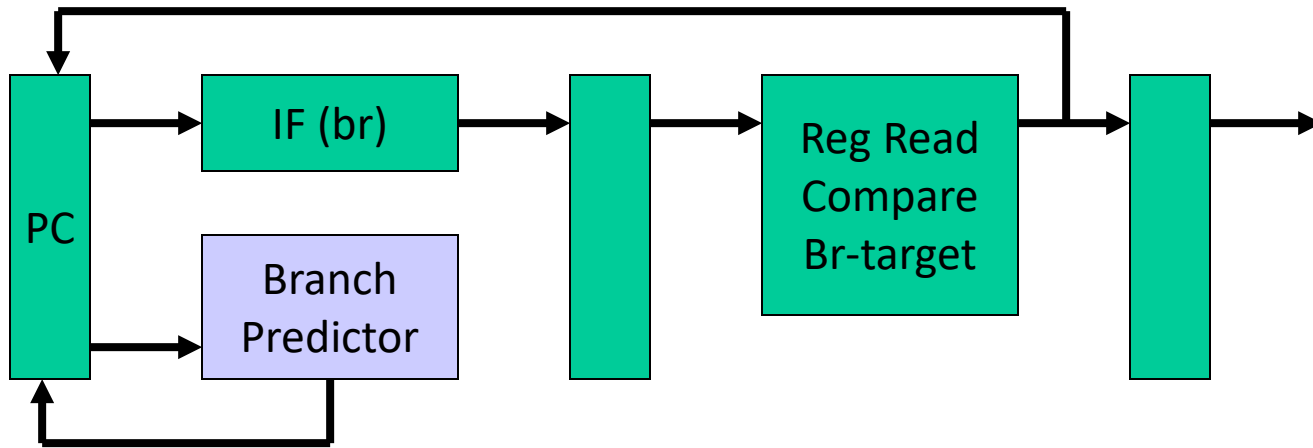
# Lecture 21: BPred, OOO, Memory Hierarchy

---

- Today's topics:
  - Branch Predictors
  - Out-of-order execution
  - Cache intro

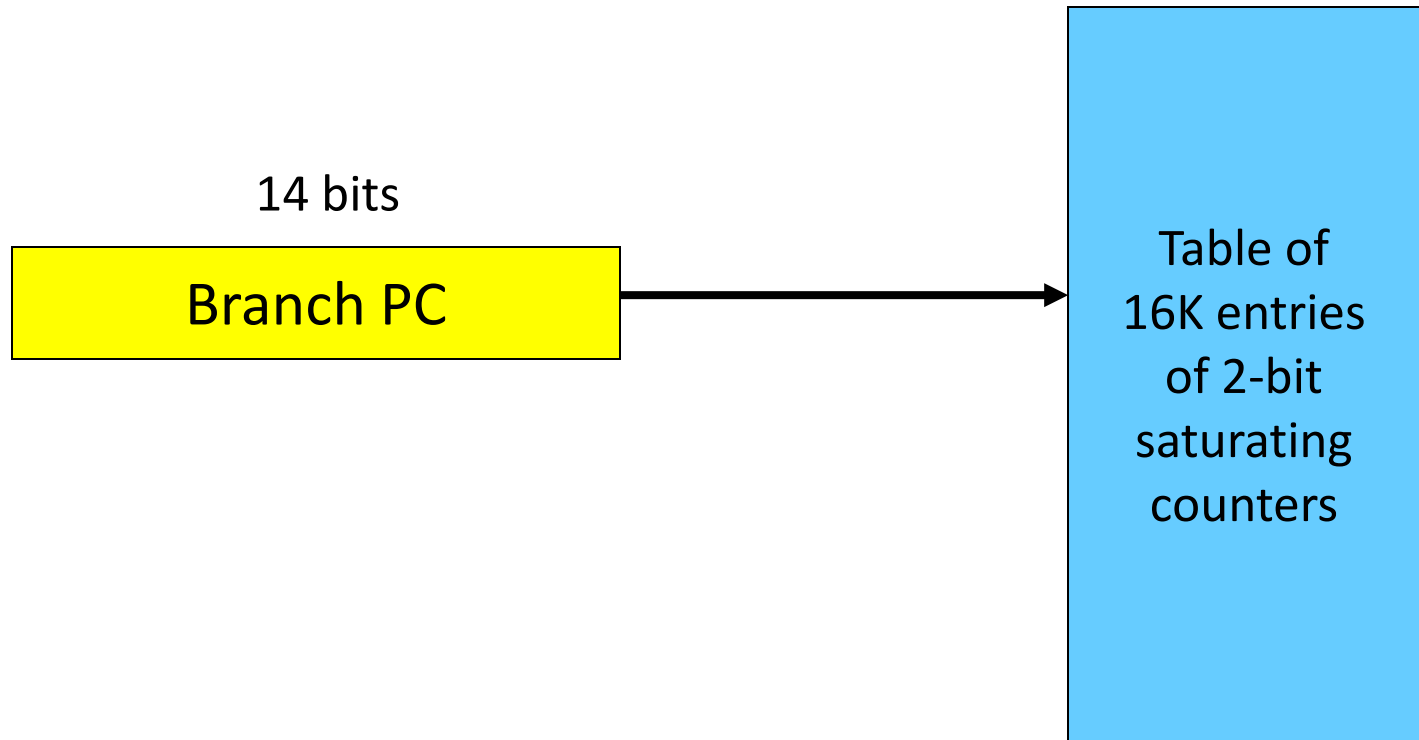
# Pipeline with Branch Predictor

---



# Bimodal Predictor

---



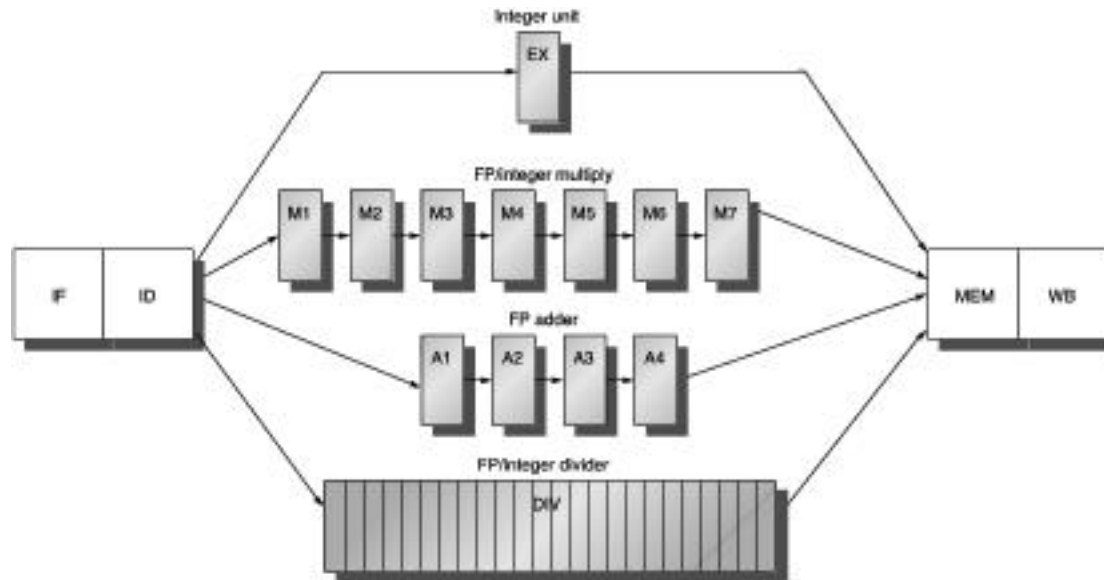
## 2-Bit Prediction

---

- For each branch, maintain a 2-bit saturating counter:  
if the branch is taken:  $\text{counter} = \min(3, \text{counter} + 1)$   
if the branch is not taken:  $\text{counter} = \max(0, \text{counter} - 1)$   
... sound familiar?
- If ( $\text{counter} \geq 2$ ), 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

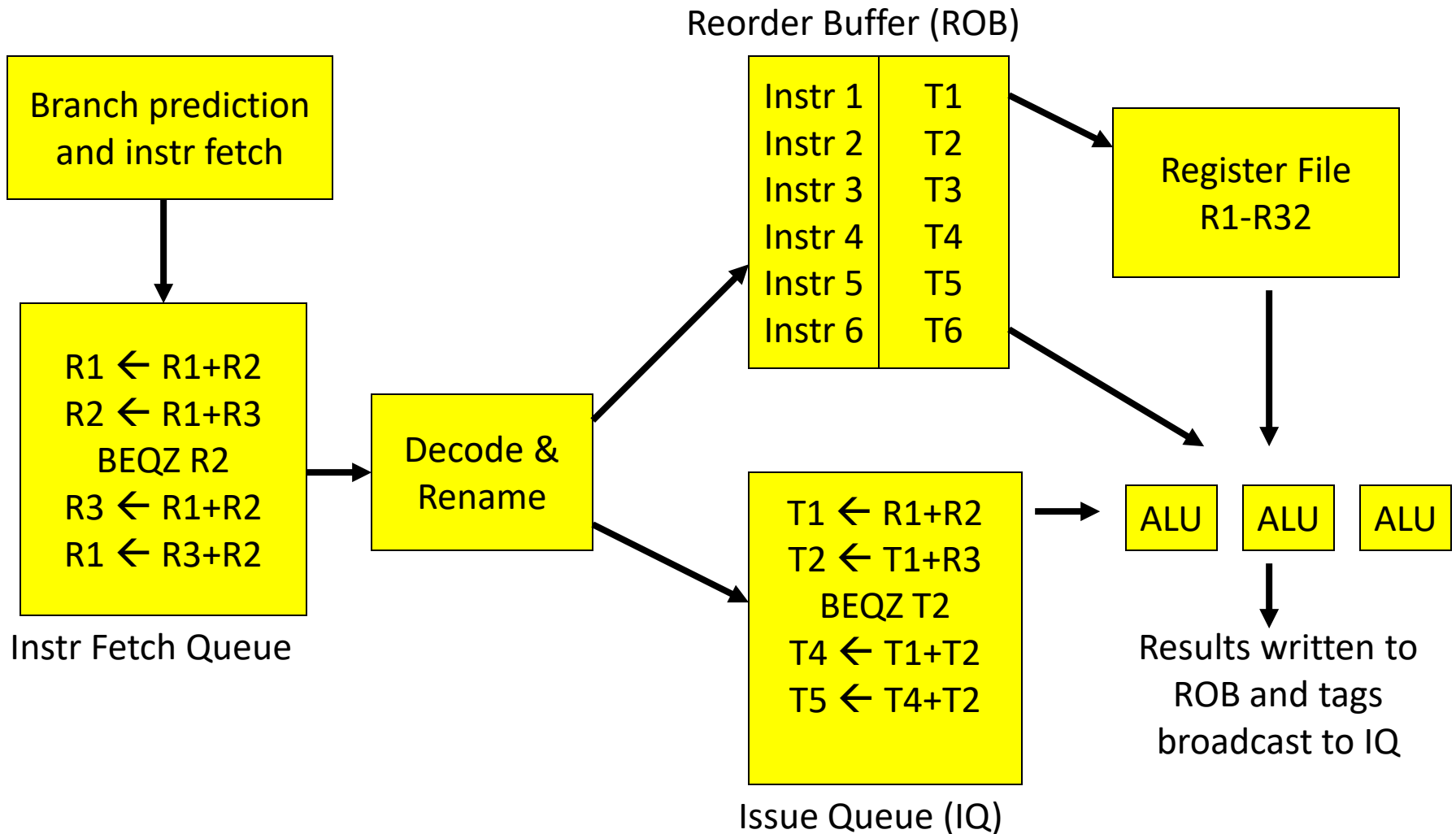
# Multicycle Instructions



© 2008 Elsevier Science (USA). All rights reserved.

- Multiple parallel pipelines – each pipeline can have a different number of stages
- Instructions can now complete out of order – must make sure that writes to a register happen in the correct order

# An Out-of-Order Processor Implementation



# 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

---

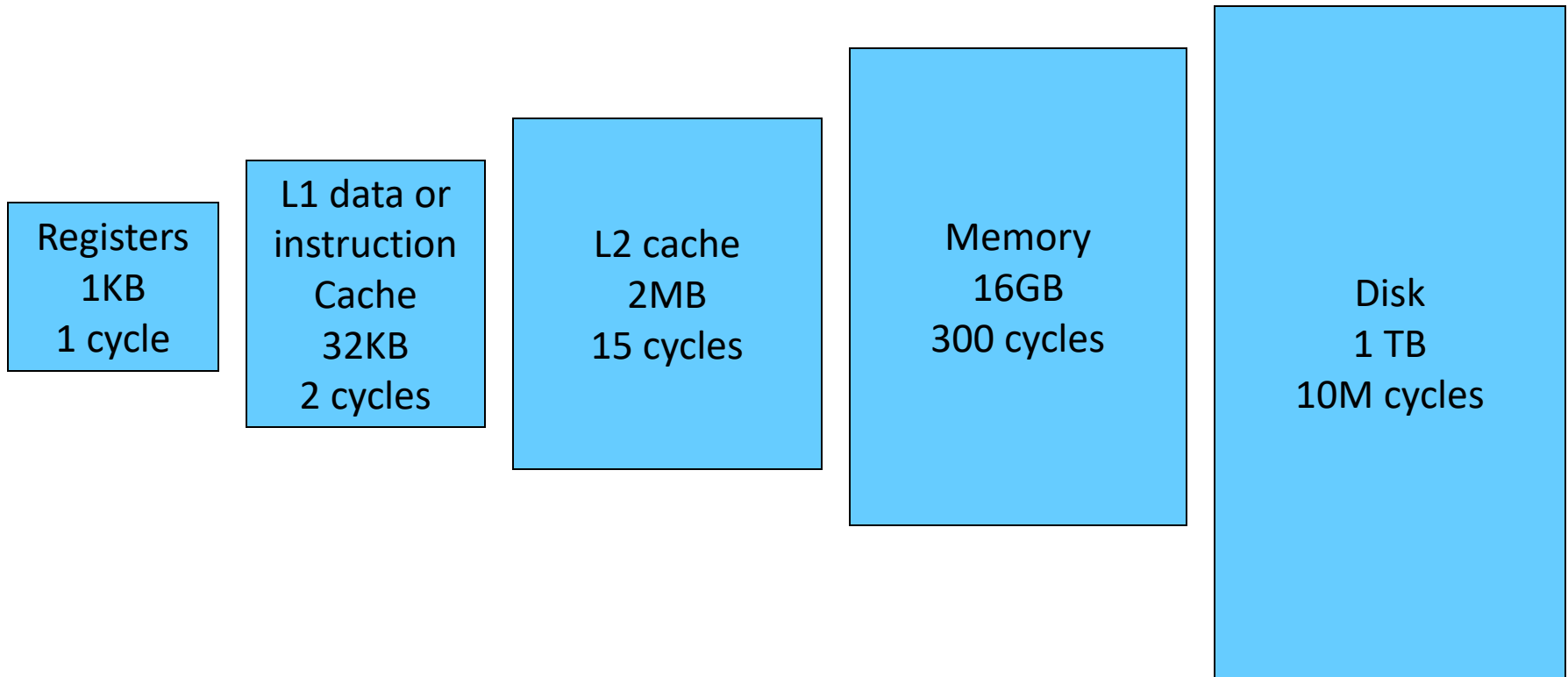
- Data and instructions are stored on DRAM chips – DRAM is a technology that has high bit density, but relatively poor latency – an access to data in memory can take as many as 300 cycles today!
- Hence, some data is stored on the processor in a structure called the cache – caches employ SRAM technology, 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 =  $0.95 \times 1 + 0.05 \times (301)$   
= 16 cycles