

RISC Design

Memory System

Virendra Singh
Professor

Computer Architecture and Dependable Systems Lab

Department of Electrical Engineering

Indian Institute of Technology Bombay

<http://www.ee.iitb.ac.in/~viren/>

E-mail: viren@ee.iitb.ac.in

EE-739: Processor Design



Lecture 28 (01 Apr 2021)

CADSL

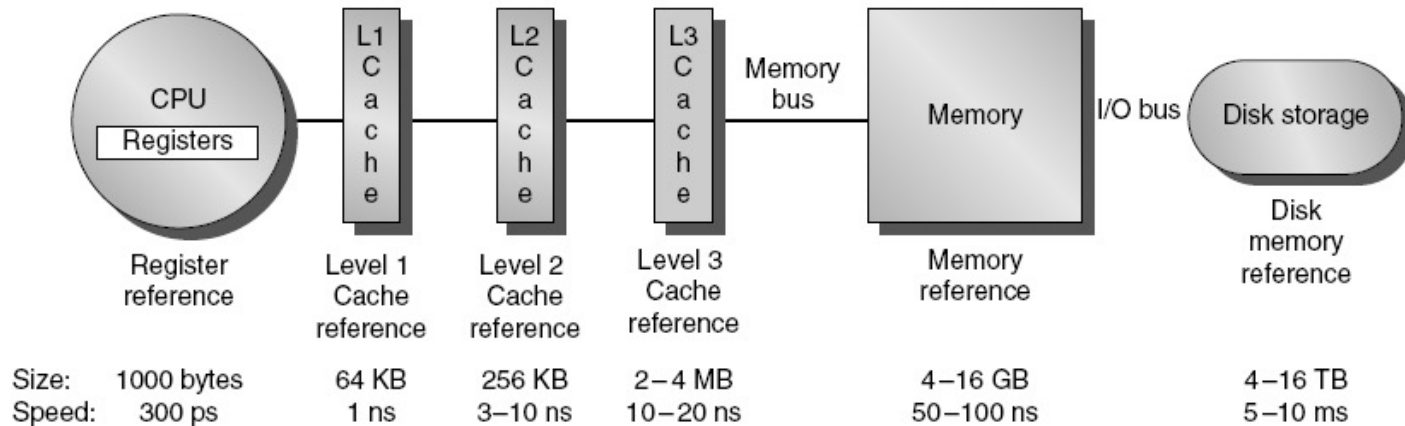
Memory Optimization

$$\frac{\text{Misses}}{\text{Instruction}} = \frac{\text{Miss rate} \times \text{Memory accesses}}{\text{Instruction count}} = \text{Miss rate} \times \frac{\text{Memory accesses}}{\text{Instruction}}$$

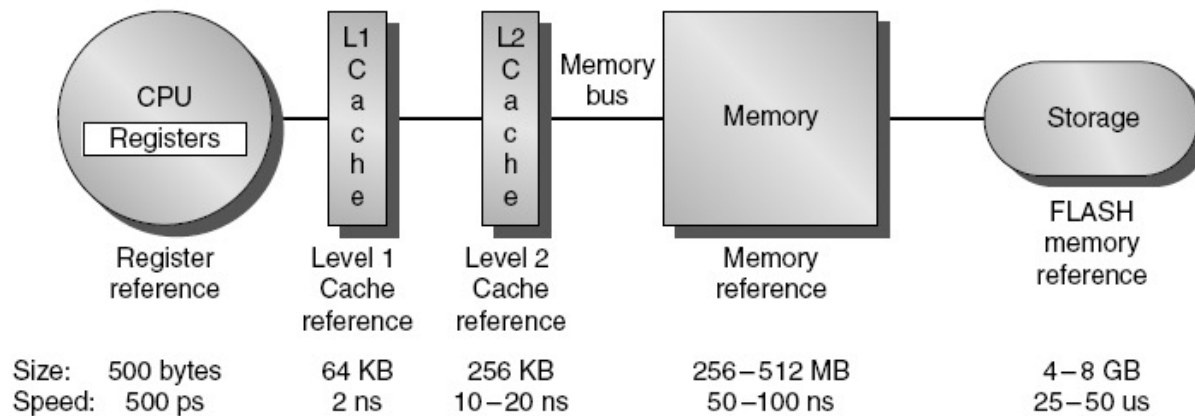
$$\text{Average memory access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$



Memory Hierarchy



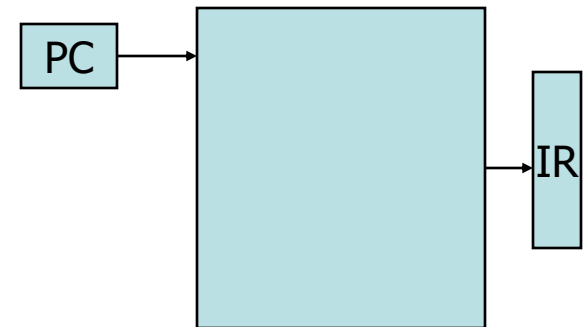
(a) Memory hierarchy for server



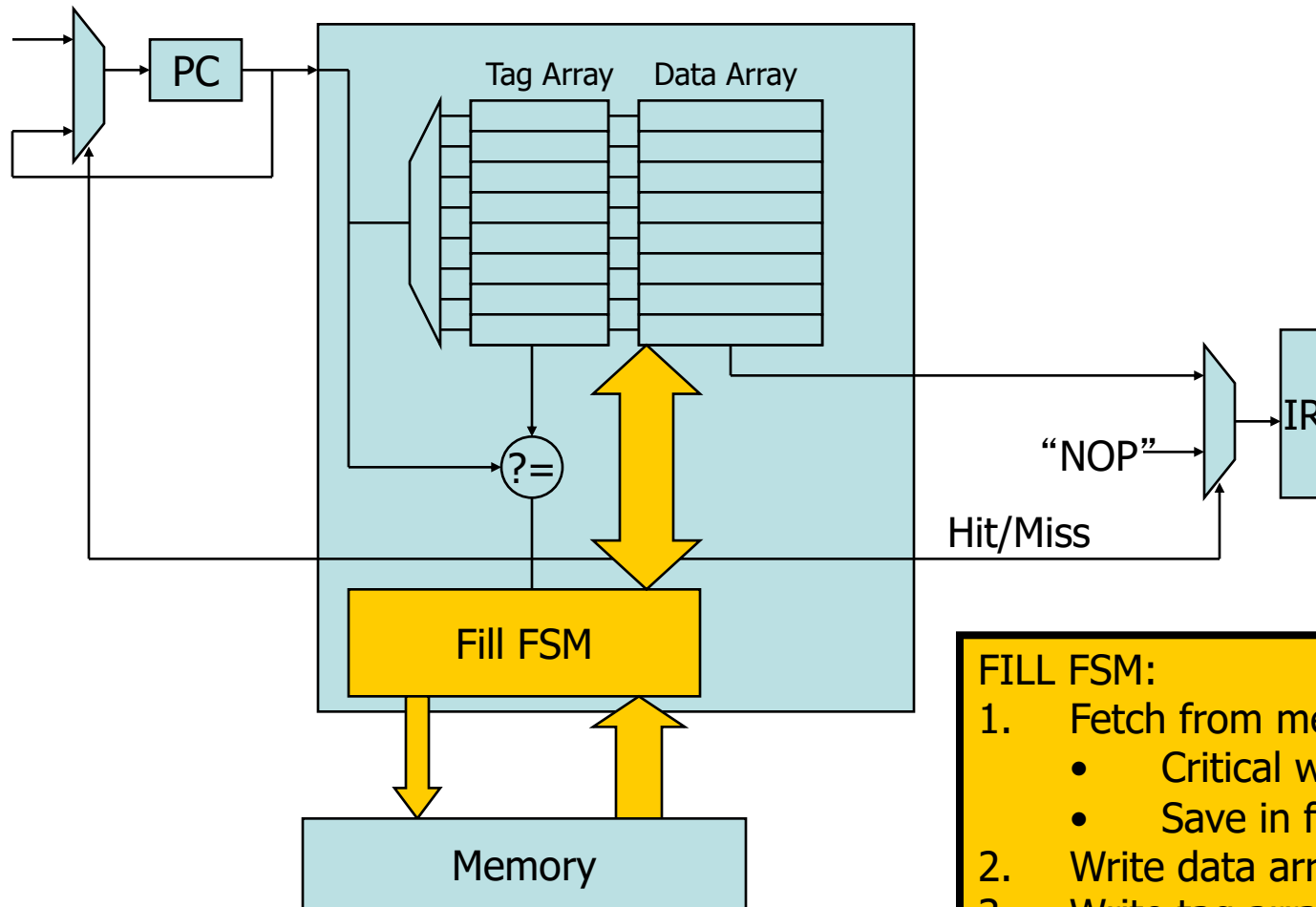
(b) Memory hierarchy for a personal mobile device

Caches and Pipelining

- Instruction cache
 - No writes, so simpler
- Interface to pipeline:
 - Fetch address (from PC)
 - Supply instruction (to IR)
- What happens on a miss?
 - Stall pipeline; inject nop
 - Initiate cache fill from memory
 - Supply requested instruction, end stall condition



I-Caches and Pipelining



FILL FSM:

1. Fetch from memory
 - Critical word first
 - Save in fill buffer
2. Write data array
3. Write tag array
4. Miss condition ends

D-Caches and Pipelining

- Pipelining loads from cache
 - Hit/Miss signal from cache
 - Stalls pipeline or inject NOPs?
 - Hard to do in current real designs, since wires are too slow for global stall signals
 - Instead, treat more like branch misprediction
 - Cancel/flush pipeline
 - Restart when cache fill logic is done



D-Caches and Pipelining

- Stores more difficult

- MEM stage:

- Perform tag check
 - Only enable write on a hit
 - On a miss, must not write (data corruption)

- Problem:

- Must do tag check and data array access sequentially
 - This will hurt cycle time or force extra pipeline stage
 - Extra pipeline stage delays loads as well: IPC hit!



Performance: Miss

- Miss rate
 - Fraction of cache access that result in a miss
- Causes of misses
 - Compulsory
 - First reference to a block
 - Capacity
 - Blocks discarded and later retrieved
 - Conflict
 - Program makes **repeated references** to multiple addresses from different blocks that map to the same location in the cache



Memory Optimization

$$\frac{\text{Misses}}{\text{Instruction}} = \frac{\text{Miss rate} \times \text{Memory accesses}}{\text{Instruction count}} = \text{Miss rate} \times \frac{\text{Memory accesses}}{\text{Instruction}}$$

$$\text{Average memory access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$

- Reducing miss rate
 - Larger block size, larger cache size, higher associativity
- Reducing miss penalty
 - Multi-level caches
- Reducing time to hit in the cache
 - Avoid address translation when indexing caches



Memory Hierarchy Basics

- **Five** basic cache optimizations:
 - Larger block size
 - Reduces compulsory misses
 - Increases capacity and conflict misses, increases miss penalty
 - Larger total cache capacity to reduce miss rate
 - Increases hit time, increases power consumption
 - Higher associativity
 - Reduces conflict misses
 - Increases hit time, increases power consumption



Memory Hierarchy Basics

- Five basic cache optimizations:
 - Higher number of cache levels
 - Reduces overall memory access time
 - Avoiding address translation in cache indexing
 - Reduces hit time



Summary

- Memory technology
- Memory hierarchy
 - Temporal and spatial locality
- Caches
 - Placement
 - Identification
 - Replacement
 - Write Policy
- Pipeline integration of caches



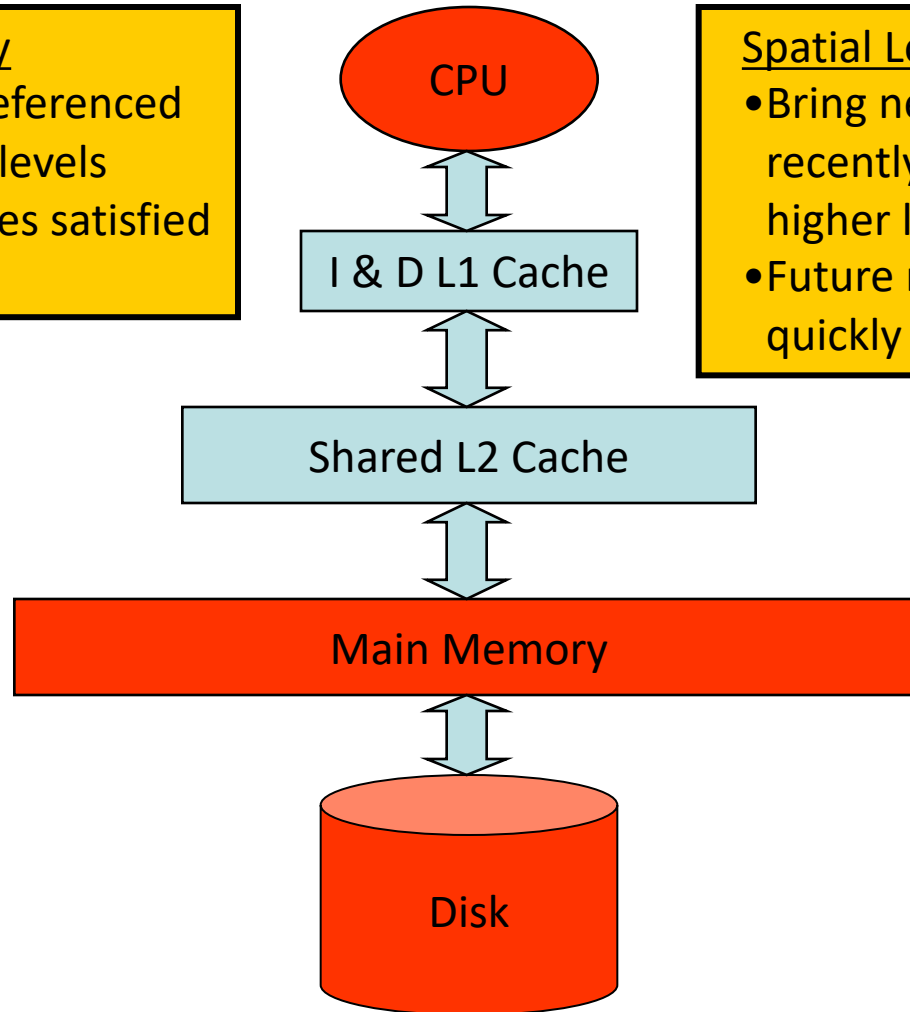
Memory Hierarchy

Temporal Locality

- Keep recently referenced items at higher levels
- Future references satisfied quickly

Spatial Locality

- Bring neighbors of recently referenced to higher levels
- Future references satisfied quickly



Thank You

