# Computer Architecture

Spring 2020

**Hamed Farbeh**

**farbeh@aut.ac.ir**

Department of Computer Engineering

Amirkabir University of Technology

# Copyright Notice

Lectures adopted from

- Computer Organization and Design: The Hardware/Software Interface, 5$^{th}$ edition, David A. Patterson, John L. Hennessy, MK pub., 2014
  - Chapter 5: Large and Fast: Exploiting Memory Hierarchy

# Multilevel Caches

- Primary cache attached to CPU
  - Small, but fast
- Level-2 cache services misses from primary cache
  - Larger, slower, but still faster than main memory
- Main memory services L2 cache misses
- Some high-end systems include L3 cache

# Multilevel Cache Example

- Given
  - CPU base CPI = 1, clock rate = 4GHz
  - Miss rate/instruction = 2%
  - Main memory access time = 100ns
- With just primary cache
  - Miss penalty = 100ns/0.25ns = 400 cycles
  - Effective CPI = 1 + 0.02 × 400 = 9

# Example (cont.)

- Now add L2 cache
  - Access time = 5ns
  - Global miss rate to main memory = 0.5%
- Primary miss with L2 hit
  - Penalty = 5ns/0.25ns = 20 cycles
- Primary miss with L2 miss
  - Extra penalty = 400 cycles
- CPI = 1 + 0.02 × 20 + 0.005 × 400 = 3.4
- Performance ratio = 9/3.4 = 2.6

# Multilevel Cache Considerations

- Primary cache
  - Focus on minimal hit time
- L2 cache
  - Focus on low miss rate to avoid main memory access
  - Hit time has less overall impact
- For multicores
  - Dedicated separate L1 Cache
  - Dedicated or shared unified L2 cache
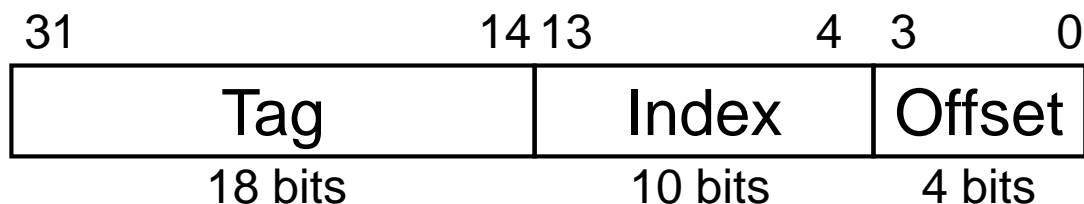  - Shared unified L3 cache

# Sources of Misses

- Compulsory misses (aka cold start misses)
  - First access to a block
- Capacity misses
  - Due to finite cache size
  - A replaced block is later accessed again
- Conflict misses (aka collision misses)
  - In a non-fully associative cache
  - Due to competition for entries in a set
  - Would not occur in a fully associative cache of the same total size

# Cache Design Trade-offs

| Design change | Effect on miss rate | Negative performance effect |
|---|---|---|
| Increase cache size | Decrease capacity misses | May increase access time |
| Increase associativity | Decrease conflict misses | May increase access time |
| Increase block size | Decrease compulsory misses | Increases miss penalty. For very large block size, may increase miss rate due to pollution. |

# Cache Control

- Example cache characteristics
  - Direct-mapped, write-back, write allocate
  - Block size: 4 words (16 bytes)
  - Cache size: 16 KB (1024 blocks)
  - 32-bit byte addresses
  - Valid bit and dirty bit per block
  - Blocking cache (vs. non-blocking)
    - CPU waits until access is complete

| 31 | 14 13 | 4 3 | 0 |
|---|---|---|---|
| Tag | Index | Offset | |
| 18 bits | 10 bits | 4 bits | |

# Interface Signals



```
          Read/Write                    Read/Write
          Valid                         Valid
                      32                            32
          Address    /                  Address    /
  CPU                        Cache                        Memory
                      32                            128
          Write Data /                  Write Data /
                      32                            128
          Read Data  /                  Read Data  /
          Ready                         Ready
```
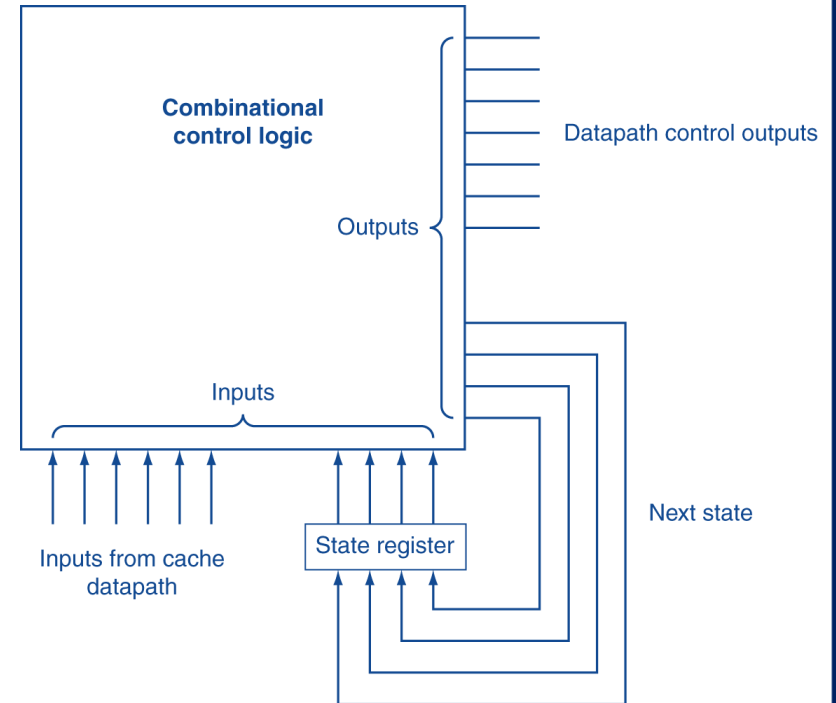
Multiple cycles per access

# Finite State Machines

- Use an FSM to sequence control steps

- Set of states, transition on each clock edge

  - State values are binary encoded

  - Current state stored in a register

  - Next state
    $= f_n$ (current state, current inputs)

- Control output signals
  $= f_o$ (current state)



Combinational control logic

Datapath control outputs

Outputs

Inputs

Inputs from cache datapath

State register

Next state

# Cache Controller FSM



Idle

Compare Tag
If Valid && Hit ,
Set Valid, SetTag,
if Write Set Dirty

Cache Hit
Mark Cache Ready

Valid CPU request

Could partition into separate states to reduce clock cycle time

Cache Miss and Old Block is Clean

Cache Miss and Old Block is Dirty

Memory Ready

Allocate
Read new block from Memory

Memory not Ready

Memory Ready

Write-Back
Write Old Block to Memory

Memory not Ready