# Chapter 5

## Large and Fast: Exploiting Memory Hierarchy

# The Memory Hierarchy

## The BIG Picture

- Common principles apply at all levels of the memory hierarchy
  - Based on notions of caching
- At each level in the hierarchy
  - Block placement
  - Finding a block
  - Replacement on a miss
  - Write policy

# Block Placement

- Determined by associativity
  - Direct mapped (1-way associative)
    - One choice for placement
  - n-way set associative
    - n choices within a set
  - Fully associative
    - Any location
- Higher associativity reduces miss rate
  - Increases complexity, cost, and access time

# Finding a Block

| Associativity | Location method | Tag comparisons |
|---|---|---|
| Direct mapped | Index | 1 |
| n-way set associative | Set index, then search entries within the set | n |
| Fully associative | Search all entries | #entries |
| | Full lookup table | 0 |

- **Virtual memory**
  - ◆ Full table lookup makes full associativity feasible
  - ◆ Benefit in reduced miss rate
- **Cache and TLB**
  - ◆ Set-associative, some cache uses direct map

# Replacement

- Choice of entry to replace on a miss
  - Least recently used (LRU)
    - Complex and costly hardware for high associativity
  - Random
    - Close to LRU, easier to implement

- Virtual memory
  - LRU approximation with hardware support

- Cache
  - Both LRU and random is ok

# Write Policy

- Write-through
  - Update both upper and lower levels
  - Simplifies replacement, but may require write buffer
- Write-back
  - Update upper level only
  - Update lower level when block is replaced
  - Need to keep more state
- Virtual memory
  - Only write-back is feasible, given disk write latency

# 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 | | |
| Increase associativity | | |
| Increase block 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 | | |
| Increase block 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 | | |

# Cache Design Trade-offs

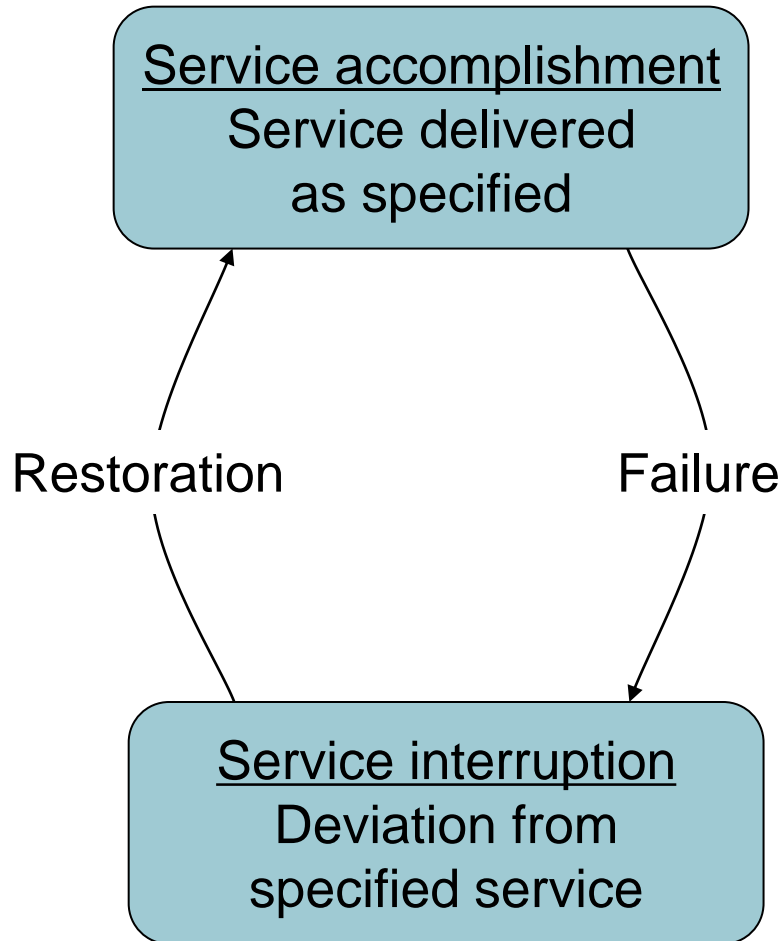| Design change | Effect on miss rate | Negative performance effect |
|---|---|---|
| Increase cache size | Decrease capacity misses | May increase access time |
| Increase associativity | Decrease conflict misses | May increase access time |
| Increase block size | Decrease compulsory misses | Increases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

# Multilevel On-Chip Caches

| Characteristic | ARM Cortex-A8 | Intel Nehalem |
|---|---|---|
| L1 cache organization | Split instruction and data caches | Split instruction and data caches |
| L1 cache size | 32 KiB each for instructions/data | 32 KiB each for instructions/data per core |
| L1 cache associativity | 4-way (I), 4-way (D) set associative | 4-way (I), 8-way (D) set associative |
| L1 replacement | Random | Approximated LRU |
| L1 block size | 64 bytes | 64 bytes |
| L1 write policy | Write-back, Write-allocate(?) | Write-back, No-write-allocate |
| L1 hit time (load-use) | 1 clock cycle | 4 clock cycles, pipelined |
| L2 cache organization | Unified (instruction and data) | Unified (instruction and data) per core |
| L2 cache size | 128 KiB to 1 MiB | 256 KiB (0.25 MiB) |
| L2 cache associativity | 8-way set associative | 8-way set associative |
| L2 replacement | Random(?) | Approximated LRU |
| L2 block size | 64 bytes | 64 bytes |
| L2 write policy | Write-back, Write-allocate (?) | Write-back, Write-allocate |
| L2 hit time | 11 clock cycles | 10 clock cycles |
| L3 cache organization | – | Unified (instruction and data) |
| L3 cache size | – | 8 MiB, shared |
| L3 cache associativity | – | 16-way set associative |
| L3 replacement | – | Approximated LRU |
| L3 block size | – | 64 bytes |
| L3 write policy | – | Write-back, Write-allocate |
| L3 hit time | – | 35 clock cycles |

# 2-Level TLB Organization

| Characteristic | ARM Cortex-A8 | Intel Core i7 |
|---|---|---|
| Virtual address | 32 bits | 48 bits |
| Physical address | 32 bits | 44 bits |
| Page size | Variable: 4, 16, 64 KiB, 1, 16 MiB | Variable: 4 KiB, 2/4 MiB |
| TLB organization | 1 TLB for instructions and 1 TLB for data<br><br>Both TLBs are fully associative, with 32 entries, round robin replacement<br><br>TLB misses handled in hardware | 1 TLB for instructions and 1 TLB for data per core<br><br>Both L1 TLBs are four-way set associative, LRU replacement<br><br>L1 I-TLB has 128 entries for small pages, 7 per thread for large pages<br><br>L1 D-TLB has 64 entries for small pages, 32 for large pages<br><br>The L2 TLB is four-way set associative, LRU replacement<br><br>The L2 TLB has 512 entries<br><br>TLB misses handled in hardware |

# Dependability

```
┌──────────────────────────────┐
│   Service accomplishment      │
│      Service delivered        │
│        as specified           │
└──────────────────────────────┘
        ↑              ↓
   Restoration       Failure
        ↑              ↓
┌──────────────────────────────┐
│    Service interruption       │
│      Deviation from           │
│     specified service         │
└──────────────────────────────┘
```
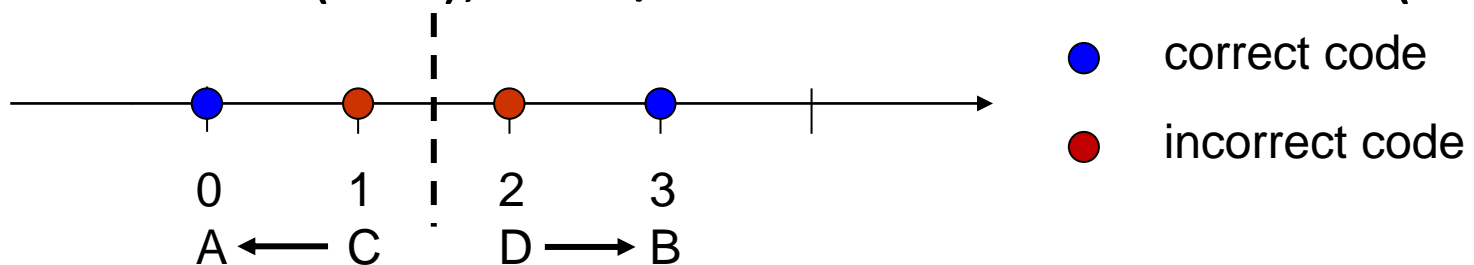
- Fault: failure of a component
  - May or may not lead to system failure

# Dependability Measures

- Reliability: mean time to failure (MTTF)

- Service interruption: mean time to repair (MTTR)

- Mean time between failures

  - MTBF = MTTF + MTTR

- Availability = MTTF / (MTTF + MTTR)

- Improving Availability

  - Increase MTTF: fault avoidance, fault tolerance, fault forecasting

  - Reduce MTTR: fault detection, fault diagnosis and fault repair

# The Hamming SEC Code

- Hamming distance
  - Number of bits that are different between two bit patterns
  - E.g. use 111 to represent 1, use 000 to represent 0, hamming distance (d) is 3, d=3.
- Minimum distance = 2 provides single bit error detection
  - E.g. odd-parity code: 10 → 101, 11 → 110, d = 2
- Minimum distance = 3 provides single error correction(SEC), 2 bit/ double error detection (DED)

```
         0        1     2        3
         A ← C         D → B
```

● correct code

● incorrect code

# Encoding SEC

- To calculate Hamming code:
  - ◆ Number bits from 1 on the left
  - ◆ All bit positions that are a power of 2 are parity bits (bit 1 2 4 8 are parity bits)
  - ◆ Each parity bit checks certain data bits:

| Bit position | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |

| Encoded date bits | | p1 | p2 | d1 | p4 | d2 | d3 | d4 | p8 | d5 | d6 | d7 | d8 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Parity bit coverate | p1 | X | | X | | X | | X | | X | | X | |
| | p2 | | X | X | | | X | X | | | X | X | |
| | p4 | | | | X | X | X | X | | | | | X |
| | p8 | | | | | | | | X | X | X | X | X |

# Decoding SEC

- Value of parity bits indicates which bits are in error

    - Use numbering from encoding procedure

    - E.g.

        - Parity bits = 0000 indicates no error

        - Parity bits = 0101 indicates bit 10 was flipped

| Bit position | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | | |
| Encoded date bits | | p1 | p2 | d1 | p4 | d2 | d3 | d4 | p8 | d5 | d6 | d7 | d8 | | |
| Parity bit coverate | p1 | X | | X | | X | | X | | X | | X | | √ | 0 |
| | p2 | | X | X | | | X | X | | | X | X | | X | 1 |
| | p4 | | | | X | X | X | X | | | | | X | √ | 0 |
| | p8 | | | | | | | | X | X | X | X | X | X | 1 |

# SEC/DED Code

- Add an additional parity bit for the whole word ($p_n$)

- Make Hamming distance = 4

- Decoding:

  - Let H = SEC parity bits

    - H even, $p_n$ even, no error

    - H odd, $p_n$ odd, correctable single bit error

    - H even, $p_n$ odd, error in $p_n$ bit

    - H odd, $p_n$ even, double error occurred

- Note: ECC DRAM uses SEC/DED with 8 bits protecting each 64 bits

# Summary

- Cache Performance

  - Mainly depends on miss rate and miss penalty

- To improve cache performance:

  - Fully associative cache

  - Set-associative cache

  - Replacement policy

  - Multilevel cache

- Dependability

  - MTTF, MTTR, reliability, availability

  - Hamming code: SEC/DED code

# Virtual Machines

- Host computer emulates guest operating system and machine resources

  - Improved isolation of multiple guests

  - Avoids security and reliability problems

  - Aids sharing of resources

- Virtualization has some performance impact

  - Feasible with modern high-performance computers

- Examples

  - IBM VM/370 (1970s technology!)

  - VMWare

  - Microsoft Virtual PC

# Virtual Machine Monitor

- Maps virtual resources to physical resources

  - Memory, I/O devices, CPUs

- Guest code runs on native machine in user mode

  - Traps to VMM on privileged instructions and access to protected resources

- Guest OS may be different from host OS

- VMM handles real I/O devices

  - Emulates generic virtual I/O devices for guest

# Example: Timer Virtualization

- In native machine, on timer interrupt

  - OS suspends current process, handles interrupt, selects and resumes next process

- With Virtual Machine Monitor

  - VMM suspends current VM, handles interrupt, selects and resumes next VM

- If a VM requires timer interrupts

  - VMM emulates a virtual timer

  - Emulates interrupt for VM when physical timer interrupt occurs

# Instruction Set Support

- User and System modes

- Privileged instructions only available in system mode

  - Trap to system if executed in user mode

- All physical resources only accessible using privileged instructions

  - Including page tables, interrupt controls, I/O registers

- Renaissance of virtualization support

  - Current ISAs (e.g., x86) adapting

# Concluding Remarks

- Fast memories are small, large memories are slow
  - We really want fast, large memories ☹
  - Caching gives this illusion ☺
- Principle of locality
  - Programs use a small part of their memory space frequently
- Memory hierarchy
  - L1 cache ↔ L2 cache ↔ … ↔ DRAM memory ↔ disk
- Virtual Memory and TLB