# ADVANCED SYSTEM ARCHITECTURES

Khoa Khoa học và Kỹ thuật Máy tính
BM Kỹ thuật Máy tính

Trần Ngọc Thịnh
http://www.cse.hcmut.edu.vn/~tnthinh
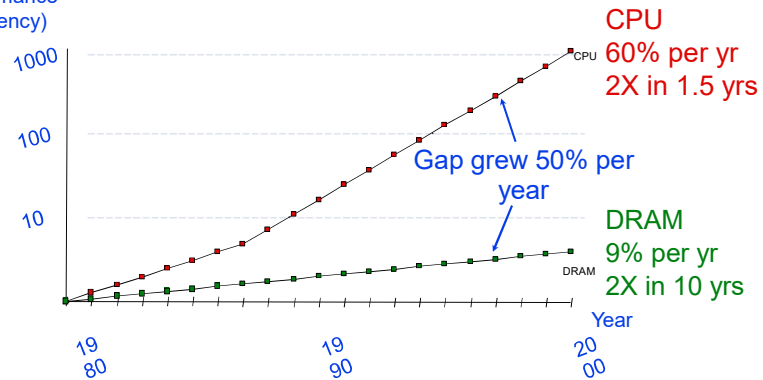
©2019, dce

---

# Memory Hierarchy Design

---

## Since 1980, CPU has outpaced DRAM ...

Four-issue 2GHz superscalar accessing 100ns DRAM could execute 800 instructions during time for one memory access!



Performance (1/latency)

CPU
CPU 60% per yr
2X in 1.5 yrs

Gap grew 50% per year

DRAM
DRAM 9% per yr
2X in 10 yrs

Year
1980 1990 2000

---

## Processor-DRAM Performance Gap Impact

- To illustrate the performance impact, assume a single-issue pipelined CPU with CPI = 1 using non-ideal memory.
- Ignoring other factors, the minimum cost of a full memory access in terms of number of wasted CPU cycles:

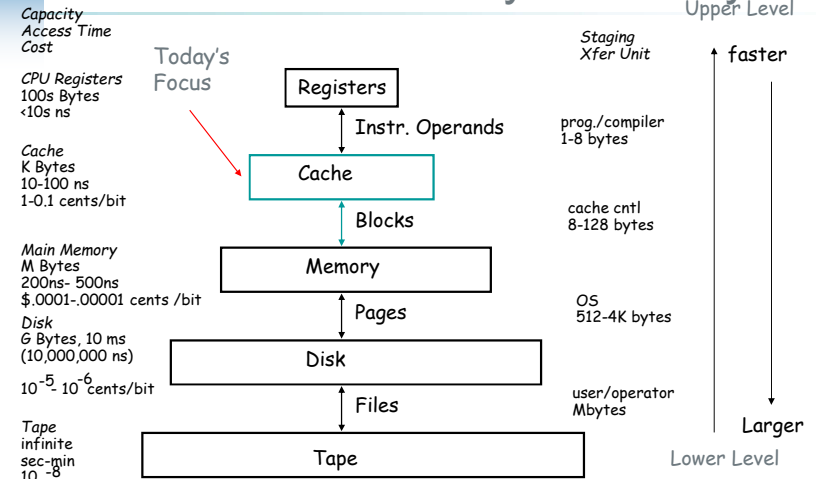| Year | CPU speed MHZ | CPU cycle ns | Memory Access ns | Minimum CPU memory stall cycles or instructions wasted | | |
|------|------|------|------|------|------|------|
| 1986: | 8 | 125 | 190 | 190/125 - 1 | = | 0.5 |
| 1989: | 33 | 30 | 165 | 165/30 -1 | = | 4.5 |
| 1992: | 60 | 16.6 | 120 | 120/16.6 -1 | = | 6.2 |
| 1996: | 200 | 5 | 110 | 110/5 -1 | = | 21 |
| 1998: | 300 | 3.33 | 100 | 100/3.33 -1 | = | 29 |
| 2000: | 1000 | 1 | 90 | 90/1 - 1 | = | 89 |
| 2002: | 2000 | .5 | 80 | 80/.5 - 1 | = | 159 |
| 2004: | 3000 | .333 | 60 | 60.333 - 1 | = | 179 |

## Addressing the Processor-Memory Performance GAP

- **Goal:** Illusion of large, fast, cheap memory. Let programs address a memory space that scales to the disk size, at a speed that is usually as fast as register access

- **Solution:** Put smaller, faster "cache" memories between CPU and DRAM. Create a "memory hierarchy".

## Levels of the Memory Hierarchy



Capacity
Access Time
Cost

Today's Focus

Staging Xfer Unit — faster — Upper Level

CPU Registers
100s Bytes
<10s ns

Registers

Instr. Operands — prog./compiler 1-8 bytes

Cache
K Bytes
10-100 ns
1-0.1 cents/bit

Cache

Blocks — cache cntl 8-128 bytes

Main Memory
M Bytes
200ns- 500ns
$.0001-.00001 cents /bit

Memory

Pages — OS 512-4K bytes

Disk
G Bytes, 10 ms
(10,000,000 ns)
$10^{-5}$- $10^{-6}$ cents/bit

Disk

Files — user/operator Mbytes

Tape
infinite
sec-min
$10^{-8}$

Tape

Larger — Lower Level

## Common Predictable Patterns

Two predictable properties of memory references:

- **Temporal Locality:** If a location is referenced, it is likely to be referenced again in the near future (e.g., loops, reuse).

- **Spatial Locality:** If a location is referenced it is likely that locations near it will be referenced in the near future (e.g., straightline code, array access).
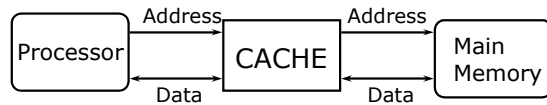
## Caches

Caches exploit both types of predictability:

- Exploit temporal locality by remembering the contents of recently accessed locations.

- Exploit spatial locality by fetching blocks of data around recently accessed locations.

## Simple view of cache

```
            Address         Address
Processor  ──────►  CACHE  ──────►   Main
           ◄──────         ◄──────   Memory
            Data            Data
```

- The processor accesses the cache first
- Cache hit: Just use the data
- Cache miss: replace a block in cache by a block from main memory, use the data
- The data transferred between cache and main memory is in blocks, and controlled by independent hardware
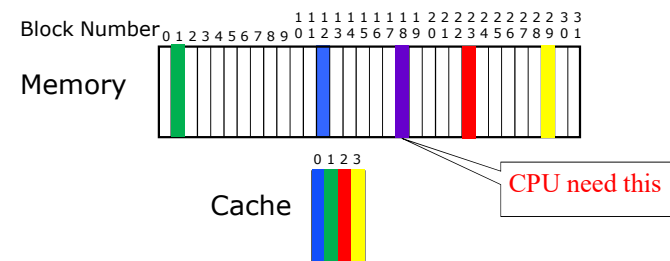
9

---

## Simple view of cache

- Hit rate: fraction of cache hit
- Miss rate: 1 – Hit rate
- Miss penalty: Time to replace a block + time to deliver the data to the processor

10

---

## Simple view of cache

- Example: For(i = 0; i < 10; i++) S = S + A[i];
- No cache: At least 12 accesses to main memory (10 A[i] and Read S, write S)
- With Cache: if A[i] and S is in a single block (ex 32-bytes), 1 access to load block to cache, and 1 access to write block to main memory
- Access to S: Temporal Locality
- Access to A[i]: Spatial Locality (A[i])

11

---

## Replacement

Block Number 0 1 2 3 4 5 6 7 8 9 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 3 3
                          0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1

Memory

0 1 2 3

Cache

CPU need this

- Cache cannot hold all blocks
- Replace a block by another that is currently needed by CPU
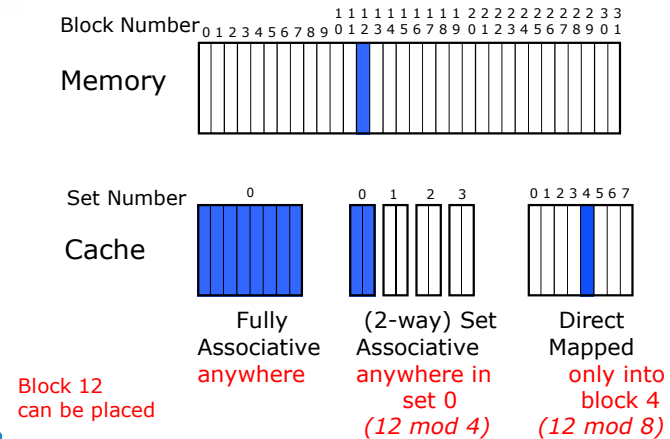
12

---

3

## Basic Cache Design & Operation Issues

- Q1: Where can a block be placed cache?
  *(Block placement strategy & Cache organization)*
  - Fully Associative, Set Associative, Direct Mapped.
- Q2: How is a block found if it is in cache?
  *(Block identification)*
  - Tag/Block.
- Q3: Which block should be replaced on a miss?
  *(Block replacement)*
  - Random, LRU, FIFO.
- Q4: What happens on a write?
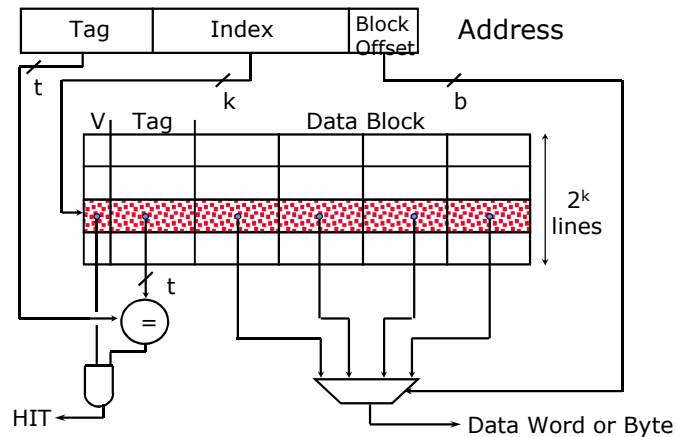  *(Cache write policy)*
  - Write through, write back.

13

---

## Q1: Where can a block be placed?



Block 12 can be placed

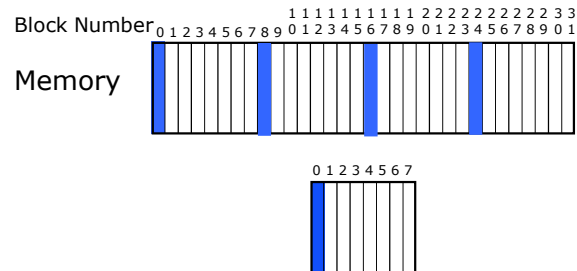| Fully Associative | (2-way) Set Associative | Direct Mapped |
|---|---|---|
| anywhere | anywhere in set 0 *(12 mod 4)* | only into block 4 *(12 mod 8)* |

14

---

## Direct-Mapped Cache



15

---

## Direct-mapped Cache

- Address: N bits ($2^N$ bytes)
- Cache has $2^k$ lines (blocks)
- Each line has $2^b$ bytes
- Block M is mapped to the line $M \% 2^k$
- Need t = N-k-bTag bits to identify mem. block
- Advantage: Simple
- Disadvantage: High miss rate
- What if CPU accesses block N0, N1 and N0 % $2^k$ = N1 % $2^k$ ?

16

---

4

## Direct-mapped Cache

Block Number
1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 3 3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1

Memory

0 1 2 3 4 5 6 7

- Access N0, N1 where N0 % $2^k$ = N1 % $2^k$
- Replace a block while there are many rooms available!

17

## 4KB Direct Mapped Cache  Example

**Index field**

Address (showing bit positions)
31 30 ... 13 12 11 ... 2 1 0

Tag field

Byte offset
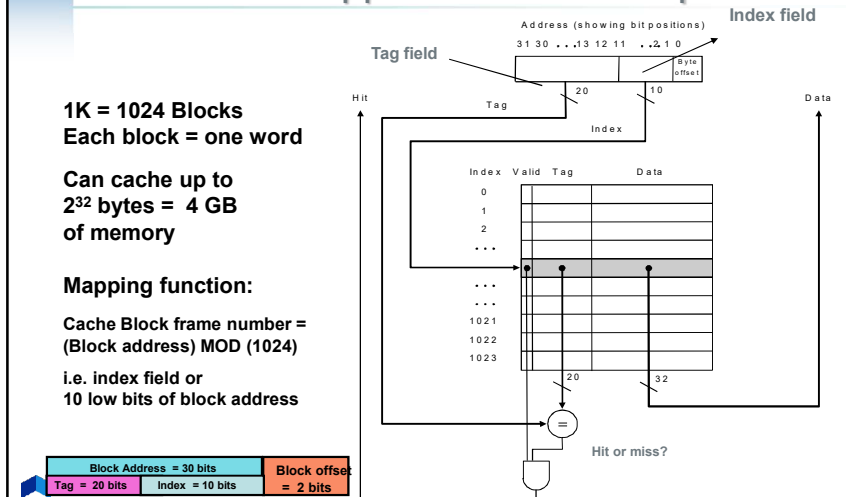
Hit

Tag

20

10

Data

Index

**1K = 1024 Blocks**
**Each block = one word**

**Can cache up to $2^{32}$ bytes =  4 GB of memory**

**Mapping function:**

Cache Block frame number =
(Block address) MOD (1024)

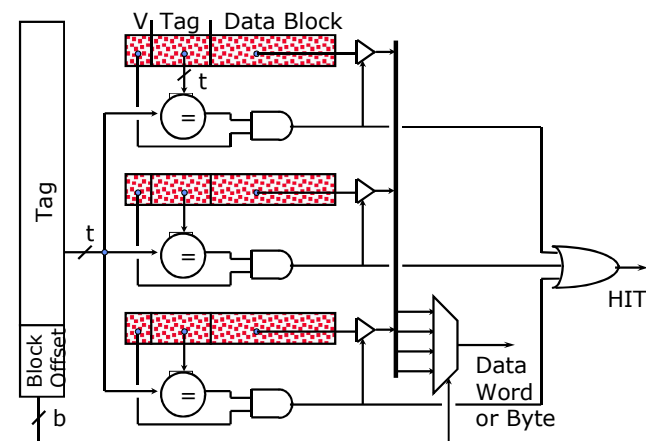i.e. index field or
10 low bits of block address

Index  Valid  Tag       Data
0
1
2
...
...
...
1021
1022
1023

20          32

=

Hit or miss?

| Block Address = 30 bits | | Block offset = 2 bits |
|---|---|---|
| Tag  = 20 bits | Index  = 10 bits | |

18

## 64KB Direct Mapped Cache Example

Tag field

Address (showing bit positions)
31 ... 16 15 ... 4 3 2 1 0

Index field

**4K= 4096 blocks**
**Each block =  four words =  16 bytes**

Hit

Tag

16

12

2 Byte offset

Index

Word select

Data

Block offset

**Can cache up to $2^{32}$ bytes = 4 GB of memory**

16 bits        128 bits
V  Tag         Data

4K entries

16    32    32    32    32

=

Hit or miss?

Mux

32

**Larger cache blocks take better advantage of spatial locality and thus may result in a lower miss rate**

| Block Address  =  28 bits | | Block offset = 4 bits |
|---|---|---|
| Tag = 16 bits | Index  = 12 bits | |

**Mapping Function:    Cache Block frame number  =  (Block address) MOD (4096)**
**i.e. index field or 12 low bit of block address**

19

## Fully Associative Cache

V  Tag   Data Block

t

Tag

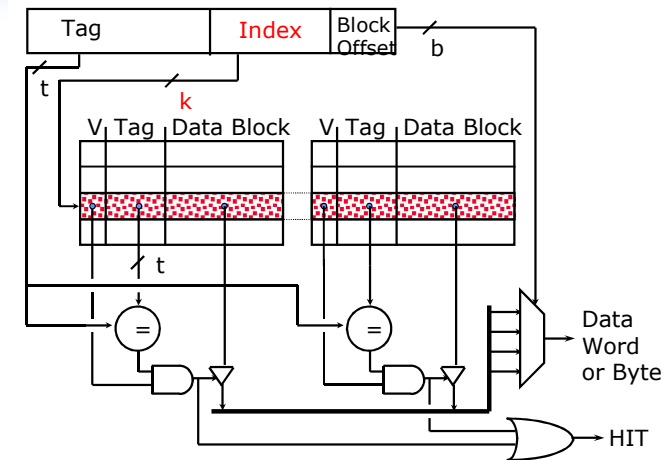t

=

=

=

Block Offset

b

HIT

Data Word or Byte

20

5

## Fully associative cache

- CAM: Content Addressable Memory
- Each block can be mapped to any lines in cache
- Tag bit: t = N-b. Compared to Tag of all lines
- Advantage: replacement occurs only when no rooms available
- Disadvantage: resource consumption, delay by comparing many elements
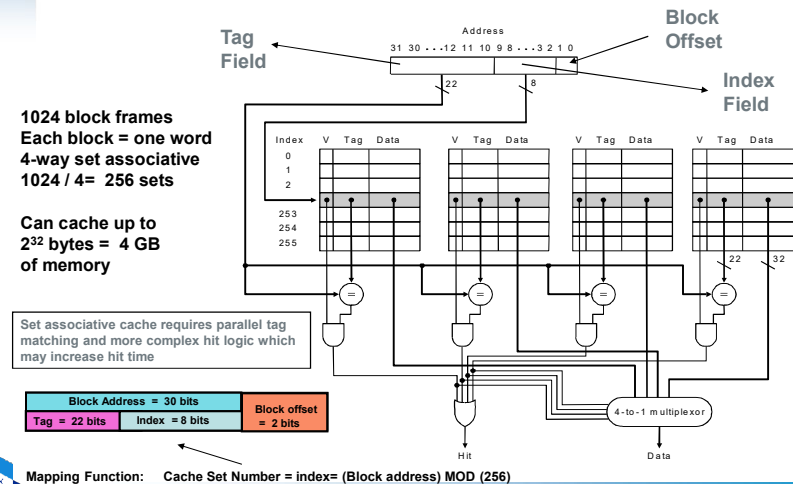
21

## Set-Associative Cache



22

## W-way Set-associative Cache

- Balancing: Direct mapped cache vs Fully associative cache
- Cache has $2^k$ sets
- Each set has *w* lines
- Block M is mapped to one of *w* lines in set M % $2^k$
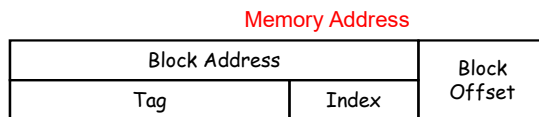- Tag bit: t = N-k-b
- Currently: widely used (Intel, AMD, …)

23

## 4K Four-Way Set Associative Cache: MIPS Implementation Example



1024 block frames
Each block = one word
4-way set associative
1024 / 4= 256 sets

Can cache up to
$2^{32}$ bytes = 4 GB
of memory

Set associative cache requires parallel tag matching and more complex hit logic which may increase hit time

| Block Address = 30 bits | | Block offset = 2 bits |
| Tag = 22 bits | Index = 8 bits | |

Mapping Function:    Cache Set Number = index= (Block address) MOD (256)

24

6

## Q2: How is a block found?

- Index selects which set to look in
- Compare Tag to find block
- Increasing associativity shrinks index, expands tag. Fully Associative caches have no index field.
- Direct-mapped: 1-way set associative?
- Fully associative: 1 set?

Memory Address

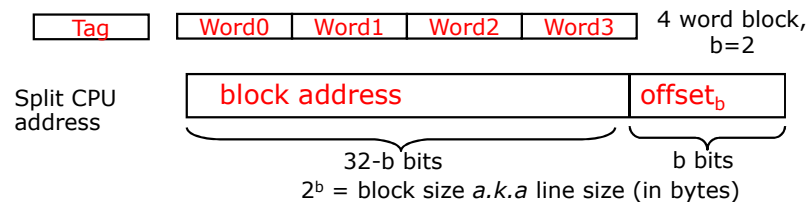| Block Address | | Block Offset |
|:---:|:---:|:---:|
| Tag | Index | |

---

## What causes a MISS?

- Three Major Categories of Cache Misses:
  - Compulsory Misses: first access to a block
  - Capacity Misses: cache cannot contain all blocks needed to execute the program
  - Conflict Misses: block replaced by another block and then later retrieved - (affects set assoc. or direct mapped caches)
    - Nightmare Scenario: ping pong effect!

---

## Block Size and Spatial Locality

Block is unit of transfer between the cache and memory

| Tag | | Word0 | Word1 | Word2 | Word3 | 4 word block, b=2 |

Split CPU address

| block address | offset$_b$ |

32-b bits        b bits

$2^b$ = block size *a.k.a* line size (in bytes)

Larger block size has distinct hardware advantages
- less tag overhead
- exploit fast burst transfers from DRAM
- exploit fast burst transfers over wide busses

*What are the disadvantages of increasing block size?*
   *Fewer blocks => more conflicts.  Can waste bandwidth.*

---

## Q3: Which block should be replaced on a miss?

- Easy for Direct Mapped
- Set Associative or Fully Associative:
  - Random
  - Least Recently Used (LRU)
    - LRU cache state must be updated on every access
    - true implementation only feasible for small sets (2-way, 4-way)
    - pseudo-LRU binary tree often used for 4-8 way
  - First In, First Out (FIFO) a.k.a. Round-Robin
    - used in highly associative caches
- Replacement policy has a second order effect since replacement only happens on misses

## Q4: What happens on a write?

- Cache hit:
    - *write through:* write both cache & memory
        - generally higher traffic but simplifies cache coherence
    - *write back:* write cache only
      (memory is written only when the entry is evicted)
        - a dirty bit per block can further reduce the traffic
- Cache miss:
    - *no write allocate:* only write to main memory
    - *write allocate (aka fetch on write):* fetch into cache

- Common combinations:
    - write through and no write allocate
    - write back with write allocate

## Reading assignment 1

- Cache performance
    - Replacement policy  (algorithms)
    - Optimization (Miss rate, penalty, …)
- Reference
    - Hennessy - Patterson - Computer Architecture. A Quantitative
    - *www2.lns.mit.edu/~avinatan/research/cache.pdf*
    - *More on internet*

## Reading Log1

- Listing some algorithms of replacement policy? Can it affect cache performance?

- What is cache performance?

- What criteria can effect cache performance?

- Explain optimization technique on each criterion

- Multi-level Caches