



TRƯỜNG ĐẠI HỌC
BÁCH KHOA HÀ NỘI
HANOI UNIVERSITY
OF SCIENCE AND TECHNOLOGY

KIẾN TRÚC MÁY TÍNH

Computer Architecture

Course ID: IT3283

Nguyễn Kim Khánh

Nội dung học phần

Chương 1. Giới thiệu chung

Chương 2. Hệ thống máy tính

Chương 3. Số học và logic máy tính

Chương 4. Kiến trúc tập lệnh

Chương 5. Bộ xử lý

Chương 6. Bộ nhớ máy tính

Chương 7. Hệ thống vào-ra

Chương 8. Các kiến trúc song song

Chương 6

BỘ NHỚ MÁY TÍNH

Nội dung của chương 6

6.1. Tổng quan hệ thống nhớ

6.2. Bộ nhớ chính

6.3. Bộ nhớ đệm (cache)

6.4. Bộ nhớ ngoài

6.1. Tổng quan hệ thống nhớ

1. Các đặc trưng của bộ nhớ

■ Vị trí

- Bên trong CPU:
 - tập thanh ghi
- Bộ nhớ trong:
 - bộ nhớ chính
 - bộ nhớ đệm (cache)
- Bộ nhớ ngoài:
 - các thiết bị lưu trữ

■ Dung lượng

- Độ dài từ nhớ (tính bằng bit)
- Số lượng từ nhớ

Các đặc trưng của bộ nhớ (tiếp)

- Đơn vị truyền
 - Từ nhớ
 - Khối nhớ
- Phương pháp truy nhập
 - Truy nhập tuần tự (băng từ)
 - Truy nhập trực tiếp (các loại đĩa)
 - Truy nhập ngẫu nhiên (bộ nhớ bán dẫn)
 - Truy nhập liên kết (cache)

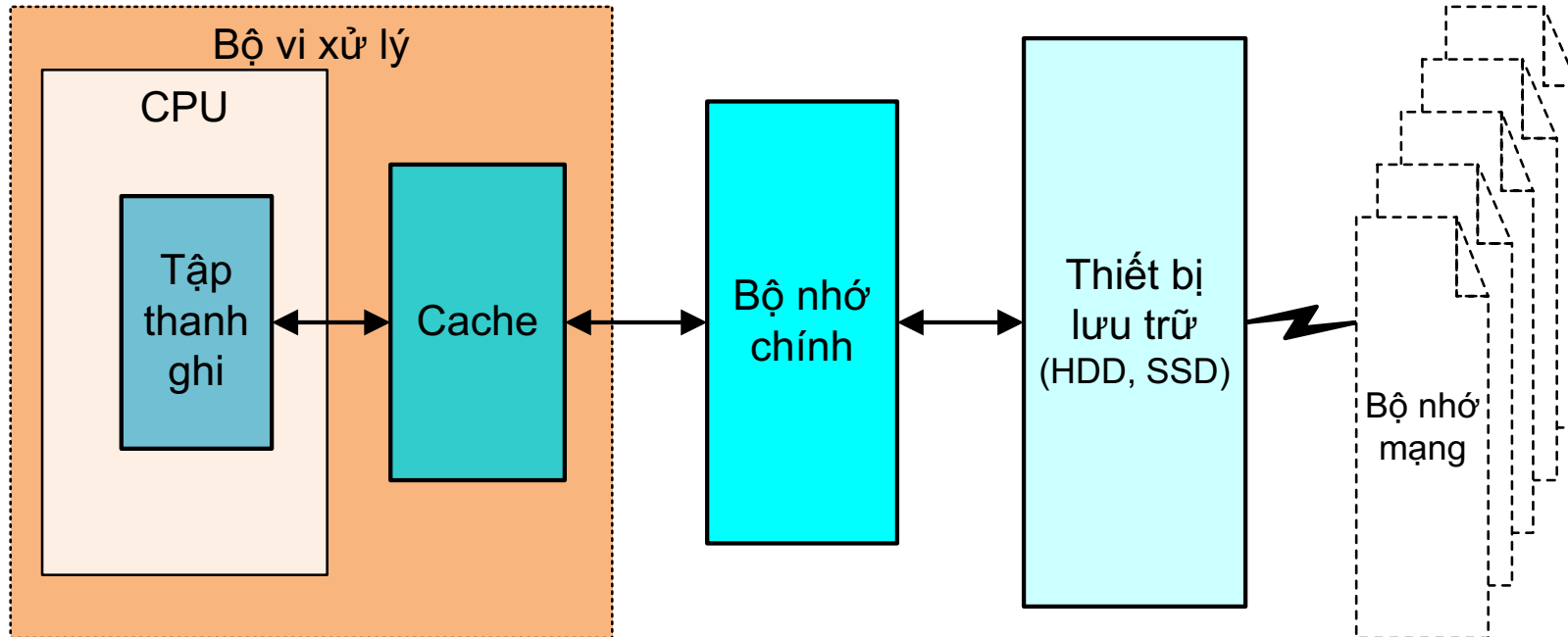
Các đặc trưng của bộ nhớ (tiếp)

- Hiệu năng (performance)
 - Thời gian truy nhập
 - Chu kỳ nhớ
 - Tốc độ truyền
- Kiểu vật lý
 - Bộ nhớ bán dẫn
 - Bộ nhớ từ
 - Bộ nhớ quang

Các đặc trưng của bộ nhớ (tiếp)

- Các đặc tính vật lý
 - Khả biến / Không khả biến (volatile / nonvolatile)
 - Xoá được / không xoá được
- Tổ chức

2. Phân cấp bộ nhớ



Từ trái sang phải:

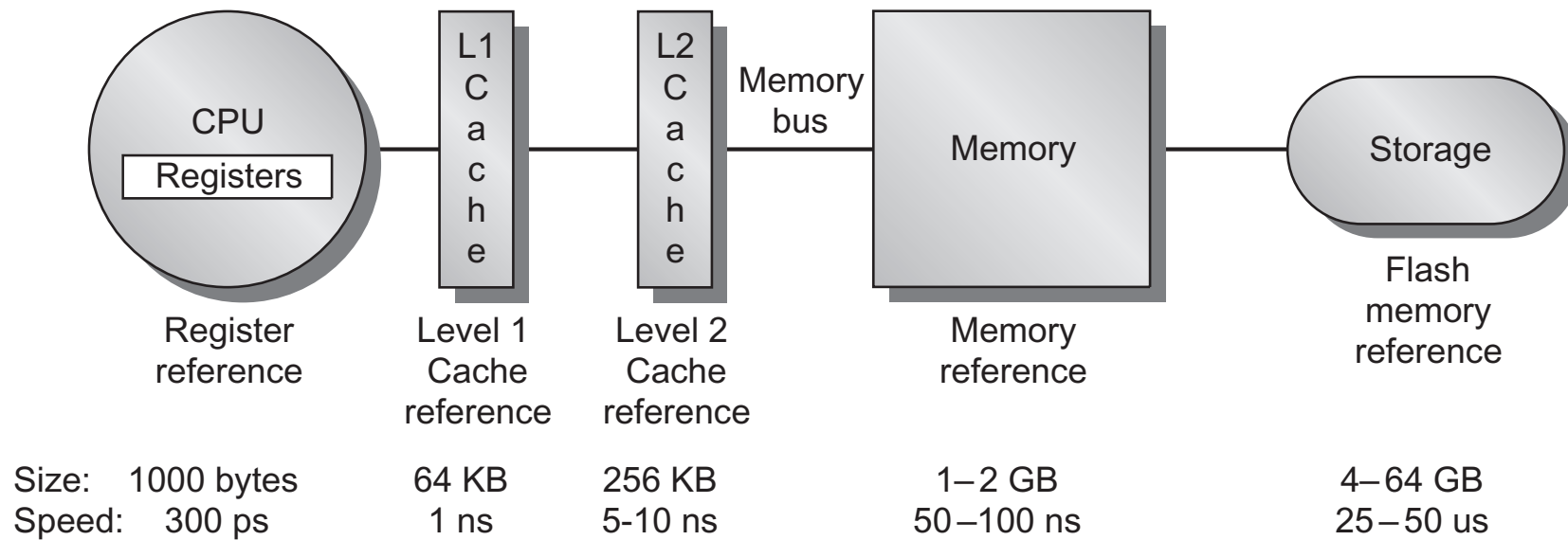
- dung lượng tăng dần
- tốc độ giảm dần
- giá thành cùng dung lượng giảm dần

Công nghệ bộ nhớ

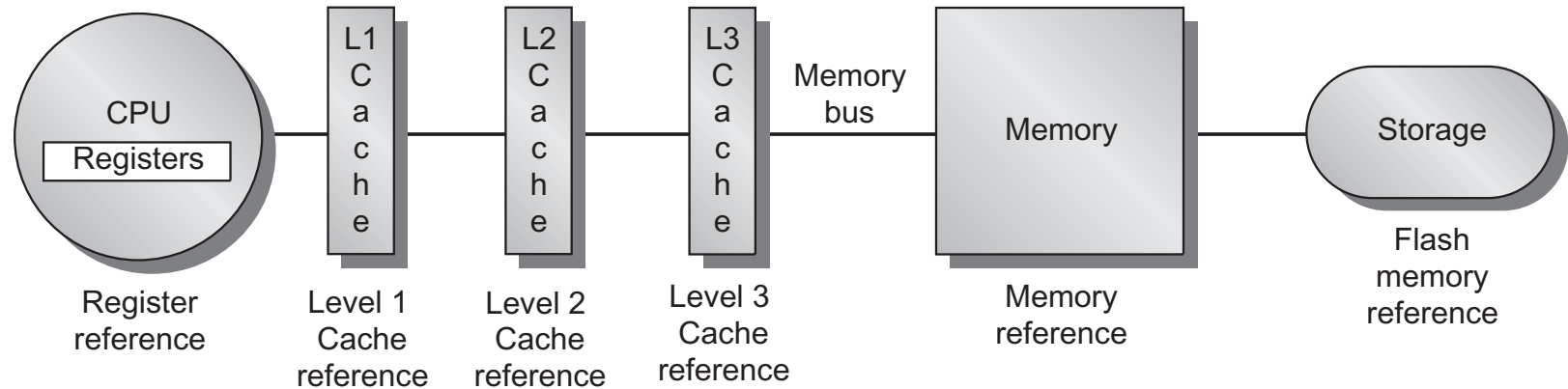
Công nghệ bộ nhớ	Thời gian truy nhập	Giá thành/GiB (2012)
SRAM	0,5 – 2,5 ns	\$500 – \$1000
DRAM	50 – 70 ns	\$10 – \$20
Flash memory	5000 – 50 000 ns	\$0,75 – \$1
HDD	5 – 20 ms	\$0,05 – \$0,1

- Bộ nhớ lý tưởng
 - Thời gian truy nhập như SRAM
 - Dung lượng và giá thành như ổ đĩa cứng

Phân cấp bộ nhớ cho thiết bị di động

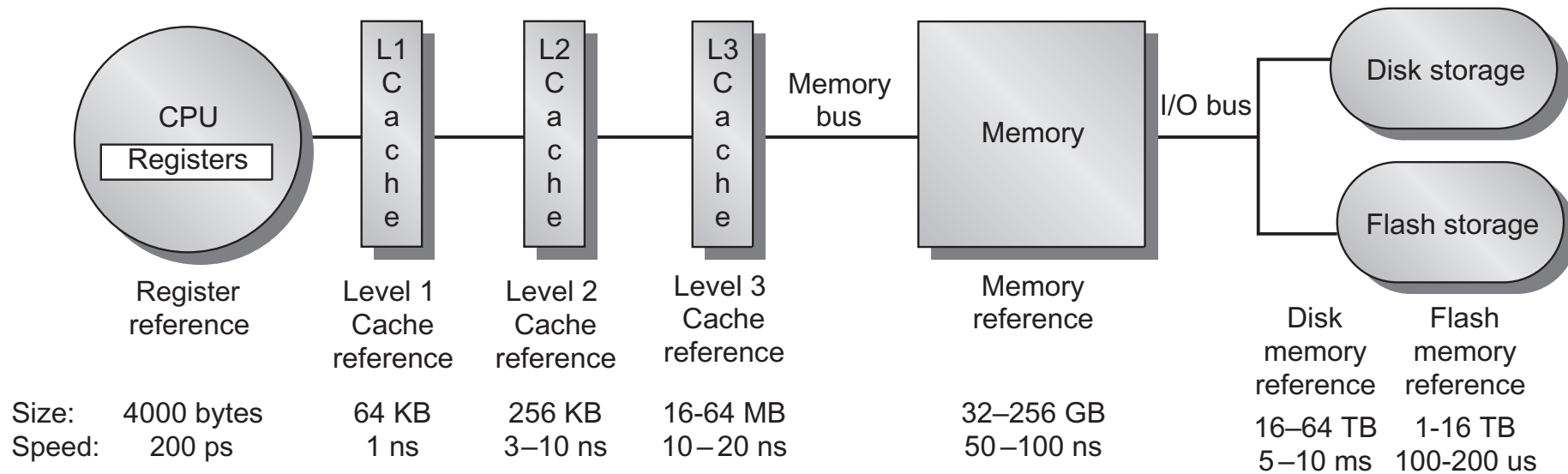


Phân cấp bộ nhớ cho máy tính PC



Laptop	Size: 1000 bytes Speed: 300 ps	64 KB 1 ns	256 KB 3–10 ns	4-8 MB 10–20 ns	4–16 GB 50–100 ns	256 GB-1 TB 50-100 uS
Desktop	Size: 2000 bytes Speed: 300 ps	64 KB 1 ns	256 KB 3–10 ns	8-32 MB 10–20 ns	8–64 GB 50–100 ns	256 GB-2 TB 50-100 uS

Phân cấp bộ nhớ cho máy chủ



Nguyên lý cục bộ hoá tham chiếu bộ nhớ

- Trong một khoảng thời gian đủ nhỏ CPU thường chỉ tham chiếu các thông tin trong một khối nhớ cục bộ
- Ví dụ:
 - Cấu trúc chương trình tuần tự
 - Vòng lặp có thân nhỏ
 - Cấu trúc dữ liệu mảng

6.2. Bộ nhớ chính

1. Bộ nhớ bán dẫn

Kiểu bộ nhớ	Tiêu chuẩn	Khả năng xóa	Cơ chế ghi	Tính khả biến
Read Only Memory (ROM)	Bộ nhớ chỉ đọc	Không xóa được	Mặt nạ	Không khả biến
Programmable ROM (PROM)			Bảng điện	
Erasable PROM (EPROM)	bằng tia cực tím, cả chip			
Electrically Erasable PROM (EEPROM)	bằng điện, mức từng byte			
Flash memory	bằng điện, từng khối			
Random Access Memory (RAM)	Bộ nhớ đọc-ghi	bằng điện, mức từng byte	Bảng điện	Khả biến

ROM (Read Only Memory)

- Bộ nhớ không khả biến
- Lưu trữ các thông tin sau:
 - Thư viện các chương trình con
 - Các chương trình điều khiển hệ thống (BIOS)
 - Các bảng chức năng
 - Vi chương trình

Các kiểu ROM

- ROM mặt nạ:
 - thông tin được ghi khi sản xuất
- PROM (Programmable ROM)
 - Cần thiết bị chuyên dụng để ghi
 - Chỉ ghi được một lần
- EPROM (Erasable PROM)
 - Cần thiết bị chuyên dụng để ghi
 - Xóa được bằng tia tử ngoại
 - Ghi lại được nhiều lần
- EEPROM (Electrically Erasable PROM)
 - Có thể ghi theo từng byte
 - Xóa bằng điện

Bộ nhớ Flash

- Ghi theo khối
- Xóa bằng điện
- Dung lượng lớn

RAM (Random Access Memory)

- Bộ nhớ đọc-ghi (Read/Write Memory)
- Khả biến
- Lưu trữ thông tin tạm thời
- Có hai loại: SRAM và DRAM
(Static and Dynamic)

SRAM (Static) – RAM tĩnh

- Các bit được lưu trữ bằng các Flip-Flop
→ thông tin ổn định
- Cấu trúc phức tạp
- Dung lượng chip nhỏ
- Tốc độ nhanh
- Đắt tiền
- Dùng làm bộ nhớ cache

DRAM (Dynamic) – RAM động

- Các bit được lưu trữ trên tụ điện
→ cần phải có mạch làm tươi
- Cấu trúc đơn giản
- Dung lượng lớn
- Tốc độ chậm hơn
- Rẻ tiền hơn
- Dùng làm bộ nhớ chính

Một số DRAM tiên tiến thông dụng

- Cải tiến để tăng tốc độ
- Synchronous DRAM (SDRAM): làm việc được đồng bộ bởi xung clock
- DDR-SDRAM (Double Data Rate SDRAM)
- DDR3, DDR4

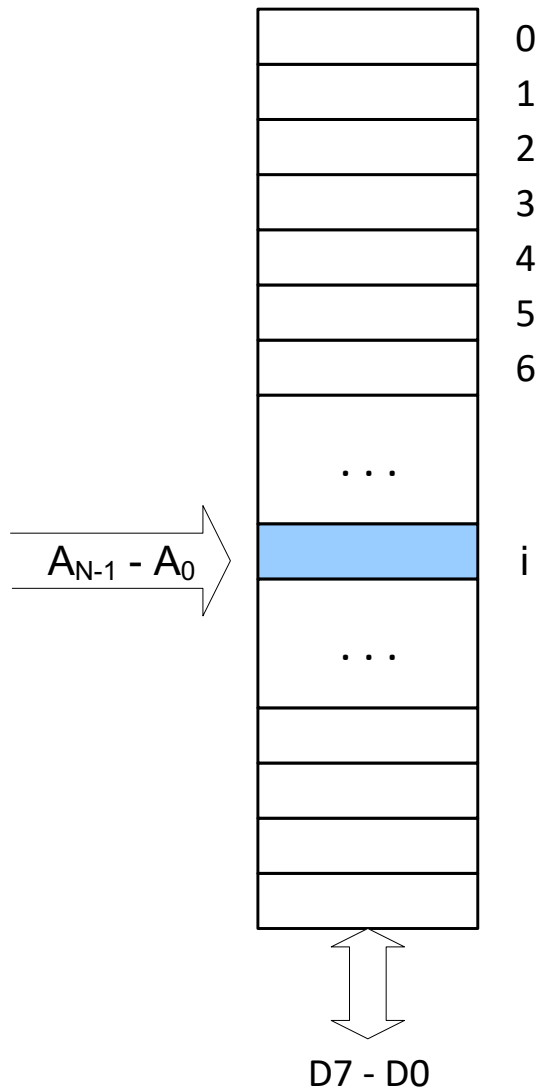
2. Các đặc trưng cơ bản của bộ nhớ chính

- Chứa các chương trình đang thực hiện và các dữ liệu đang được sử dụng
- Tồn tại trên mọi hệ thống máy tính
- Bao gồm các ngăn nhớ được đánh địa chỉ trực tiếp bởi CPU
- Dung lượng của bộ nhớ chính nhỏ hơn không gian địa chỉ bộ nhớ mà CPU quản lý.
- Việc quản lý logic bộ nhớ chính tùy thuộc vào hệ điều hành

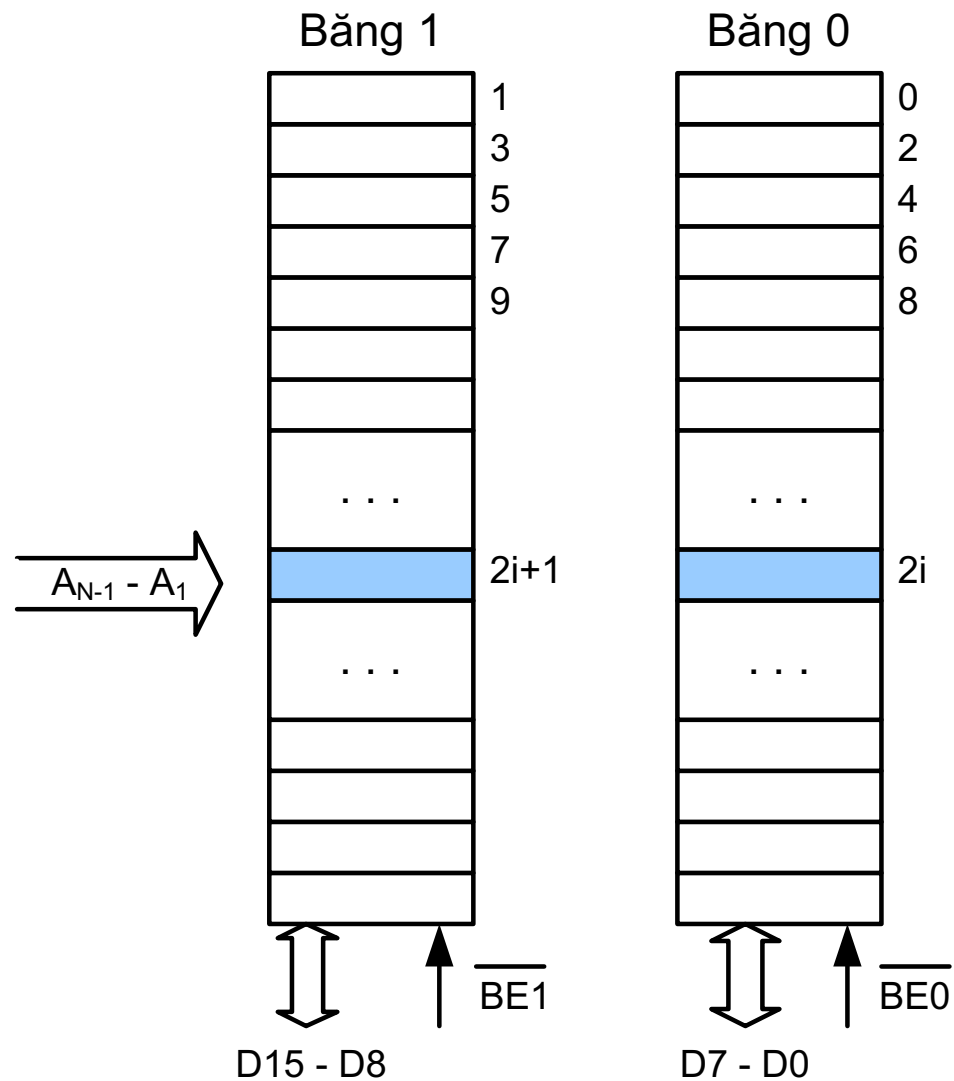
Tổ chức bộ nhớ đan xen (interleaved memory)

- Độ rộng của bus dữ liệu để trao đổi với bộ nhớ: m
= 8, 16, 32, 64, 128 ... bit
- Các ngăn nhớ được tổ chức theo byte
→ tổ chức bộ nhớ vật lý khác nhau

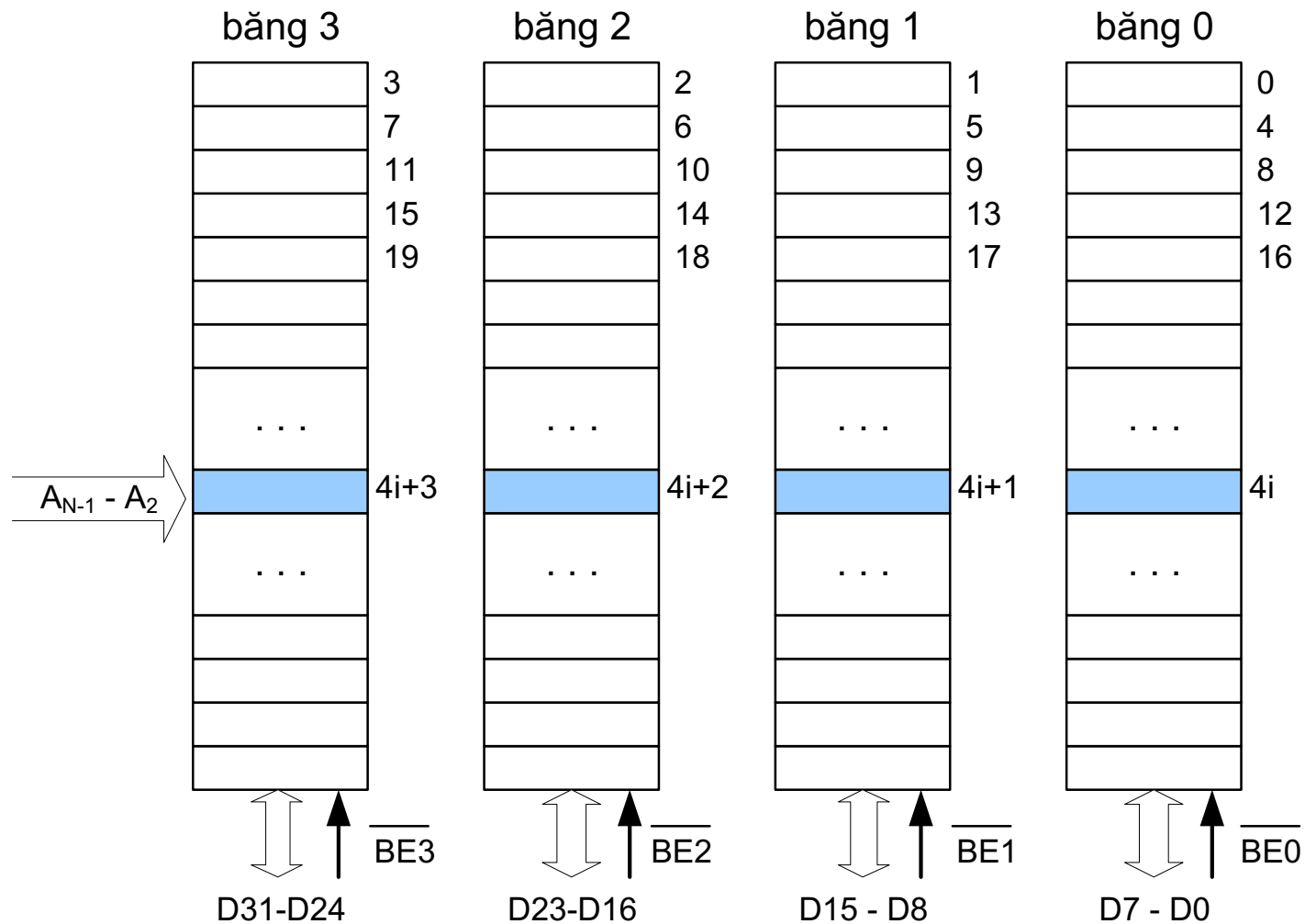
$m=8\text{bit} \rightarrow$ một bảng nhớ tuyến tính



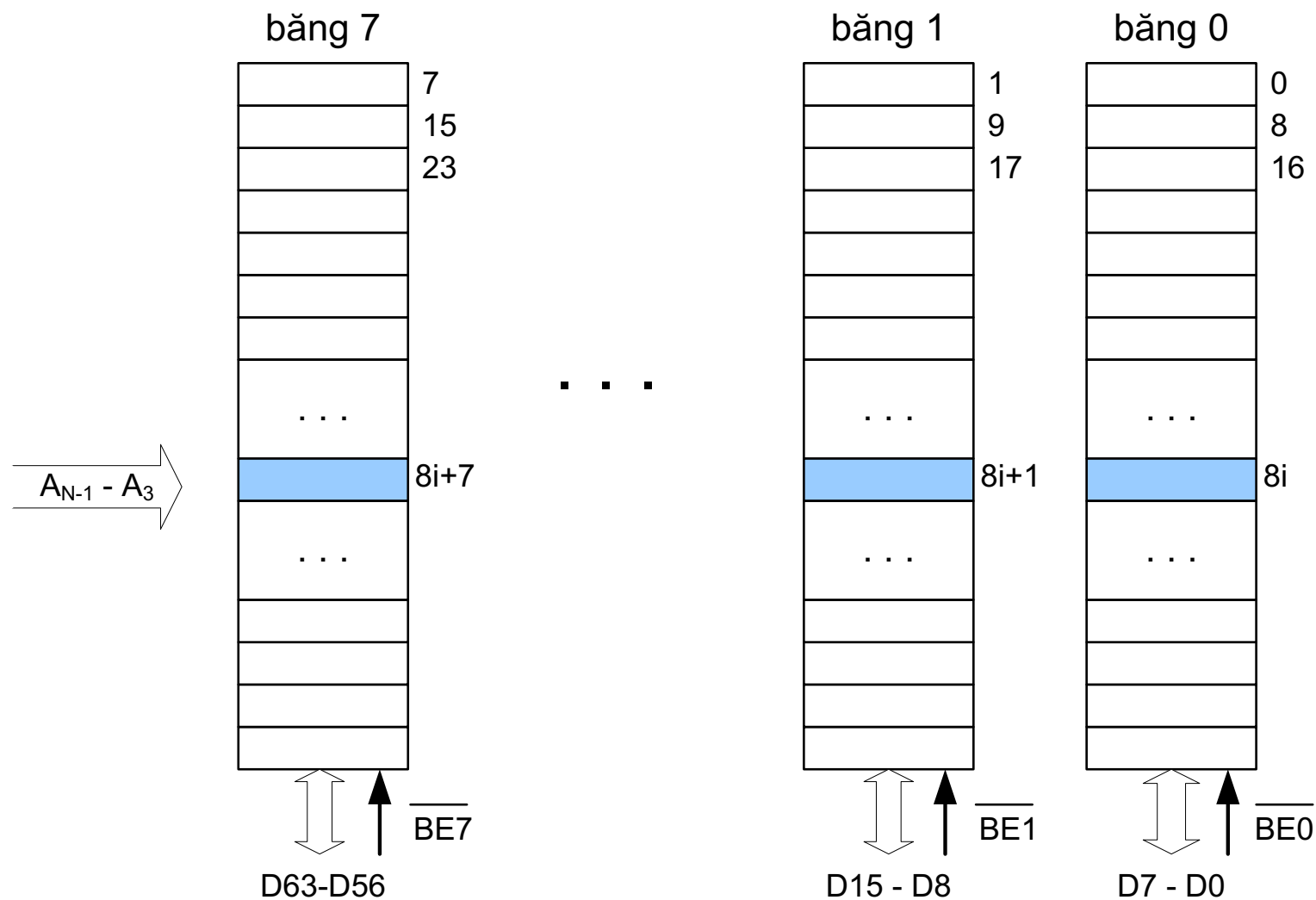
$m = 16\text{bit} \rightarrow$ hai bảng nhớ đơn xen



$m = 32\text{bit} \rightarrow$ bốn bảng nhớ đơn xen



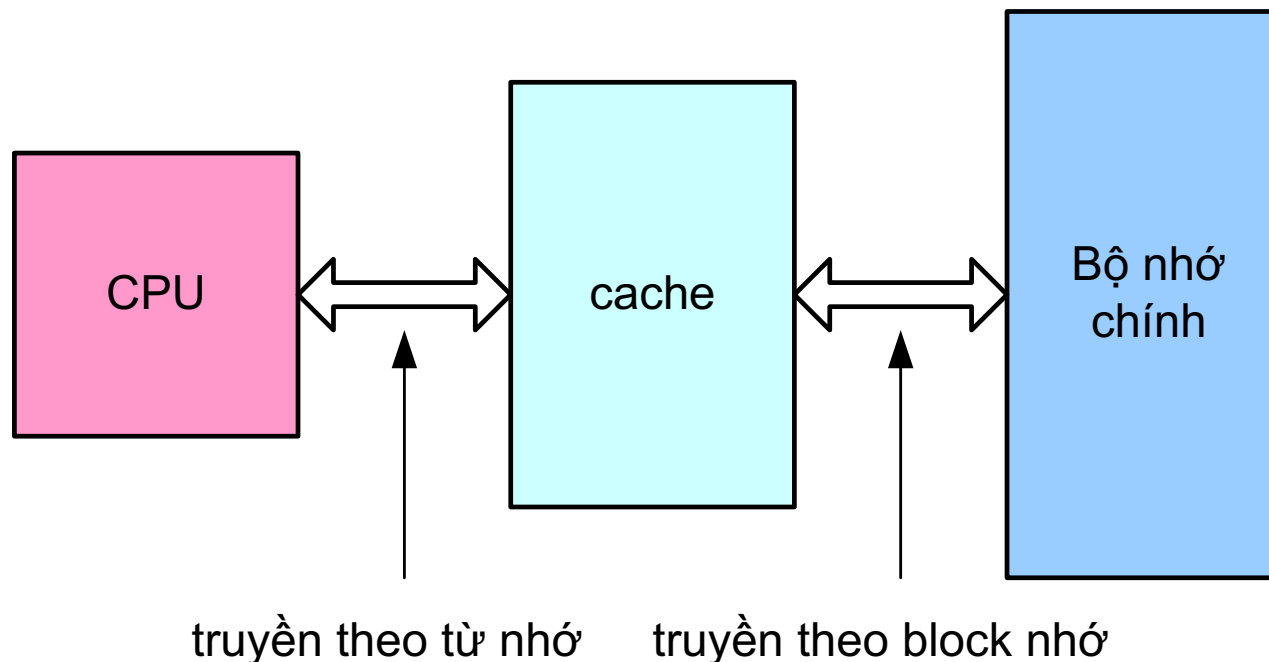
$m = 64\text{bit} \rightarrow$ tám bảng nhớ đơn xen



6.3. Bộ nhớ đệm (cache memory)

1. Nguyên tắc chung của cache

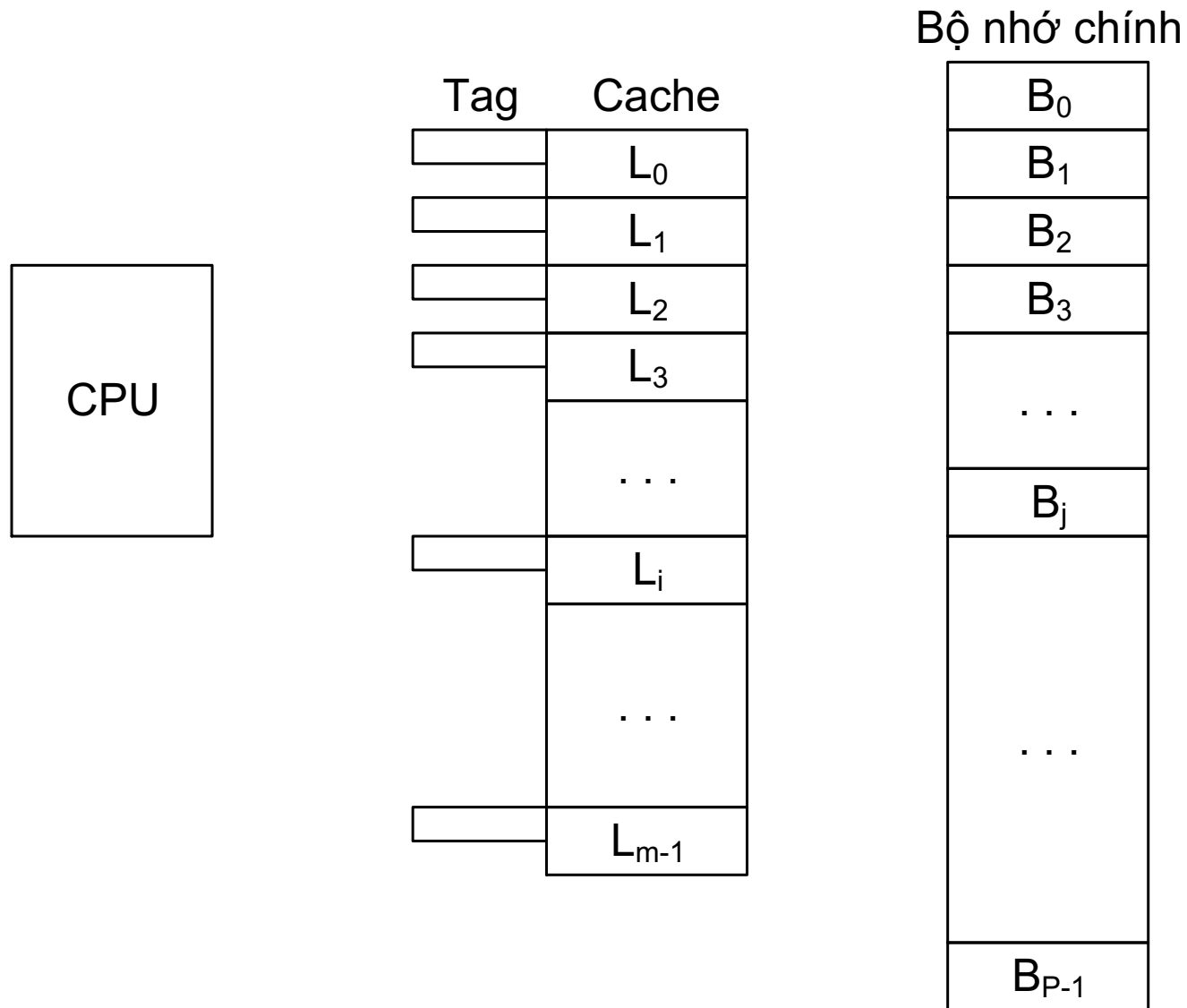
- Cache có tốc độ nhanh hơn bộ nhớ chính
- Cache được đặt giữa CPU và bộ nhớ chính nhằm tăng tốc độ CPU truy cập bộ nhớ
- Cache có thể được đặt trên chip CPU



Ví dụ về thao tác của cache

- CPU yêu cầu nội dung của ngăn nhớ
- CPU kiểm tra trên cache với dữ liệu này
- Nếu có, CPU nhận dữ liệu từ cache (nhanh)
- Nếu không có, đọc Block nhớ chứa dữ liệu từ bộ nhớ chính vào cache
- Tiếp đó chuyển dữ liệu từ cache vào CPU

Cấu trúc chung của cache / bộ nhớ chính



Cấu trúc chung của cache / bộ nhớ chính (tiếp)

- Bộ nhớ chính có 2^N byte nhớ
- Bộ nhớ chính và cache được chia thành các khối có kích thước bằng nhau
 - Bộ nhớ chính: $B_0, B_1, B_2, \dots, B_{p-1}$ (p Blocks)
 - Bộ nhớ cache: $L_0, L_1, L_2, \dots, L_{m-1}$ (m Lines)
 - Kích thước của Block (Line) = 8,16,32,64,128 byte
- Mỗi Line trong cache có một thẻ nhớ (Tag) được gắn vào

Cấu trúc chung của cache / bộ nhớ chính (tiếp)

- Một số Block của bộ nhớ chính được nạp vào các Line của cache
- Nội dung Tag (thẻ nhớ) cho biết Block nào của bộ nhớ chính hiện đang được chứa ở Line đó
- Nội dung Tag được cập nhật mỗi khi Block từ bộ nhớ chính nạp vào Line đó
- Khi CPU truy nhập (đọc/ghi) một từ nhớ, có hai khả năng xảy ra:
 - Từ nhớ đó có trong cache (cache hit)
 - Từ nhớ đó không có trong cache (cache miss).

2. Các phương pháp ánh xạ

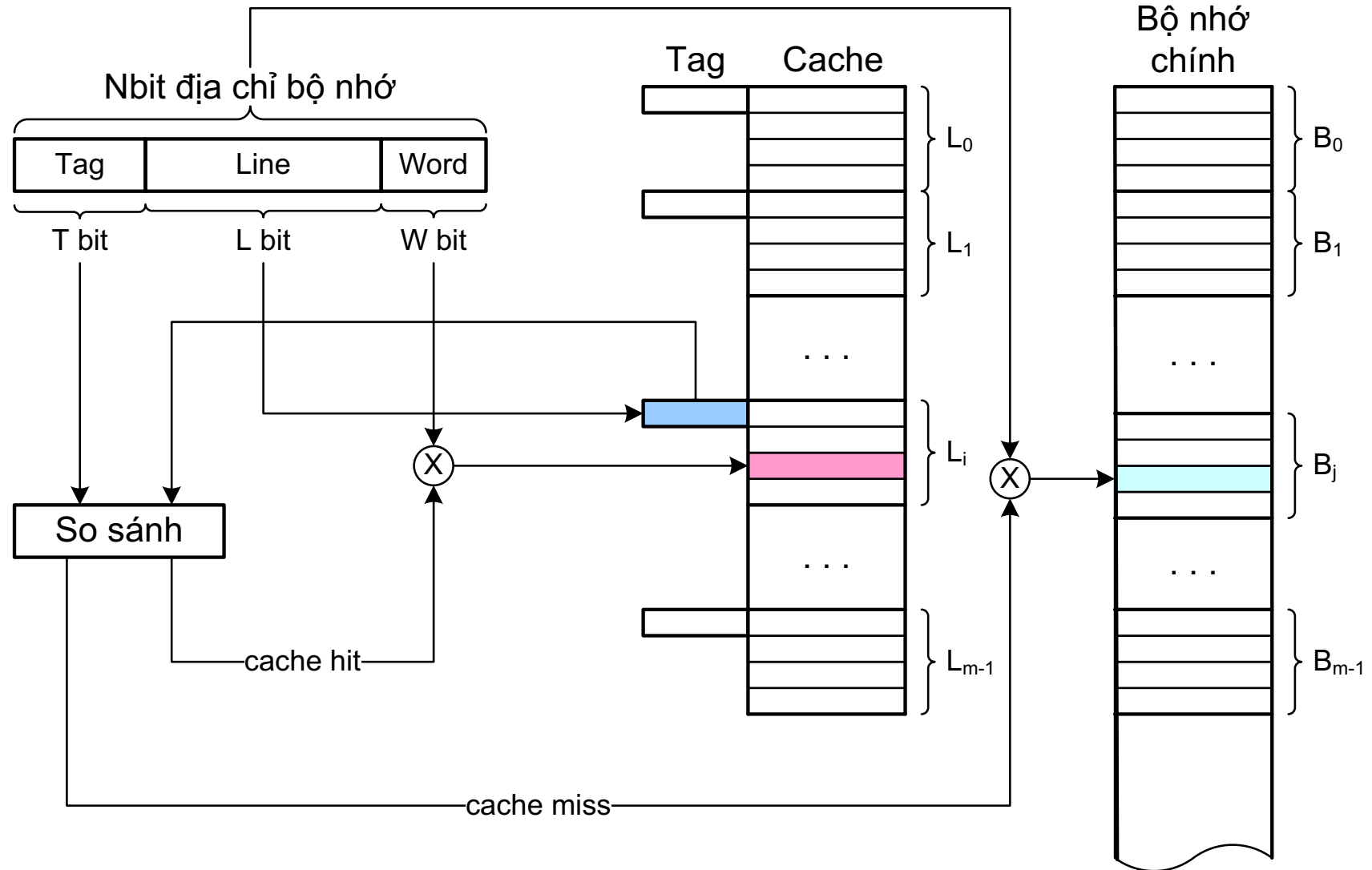
(Chính là các phương pháp tổ chức bộ nhớ cache)

- Ánh xạ trực tiếp
(Direct mapping)
- Ánh xạ liên kết toàn phần
(Fully associative mapping)
- Ánh xạ liên kết tập hợp
(Set associative mapping)

Ánh xạ trực tiếp

- Mỗi Block của bộ nhớ chính chỉ có thể được nạp vào một Line của cache:
 - $B_0 \rightarrow L_0$
 - $B_1 \rightarrow L_1$
 -
 - $B_{m-1} \rightarrow L_{m-1}$
 - $B_m \rightarrow L_0$
 - $B_{m+1} \rightarrow L_1$
 -
- Tổng quát
 - B_j chỉ có thể nạp vào $L_{j \bmod m}$
 - m là số Line của cache.

Ánh xạ trực tiếp (tiếp)



Ánh xạ trực tiếp (tiếp)

- Địa chỉ N bit của bộ nhớ chính chia thành ba trường:
 - Trường **Word** gồm W bit xác định một từ nhớ trong Block hay Line:
$$2^W = \text{kích thước của Block hay Line}$$
 - Trường **Line** gồm L bit xác định một trong số các Line trong cache:
$$2^L = \text{số Line trong cache} = m$$
 - Trường **Tag** gồm T bit:
$$T = N - (W+L)$$

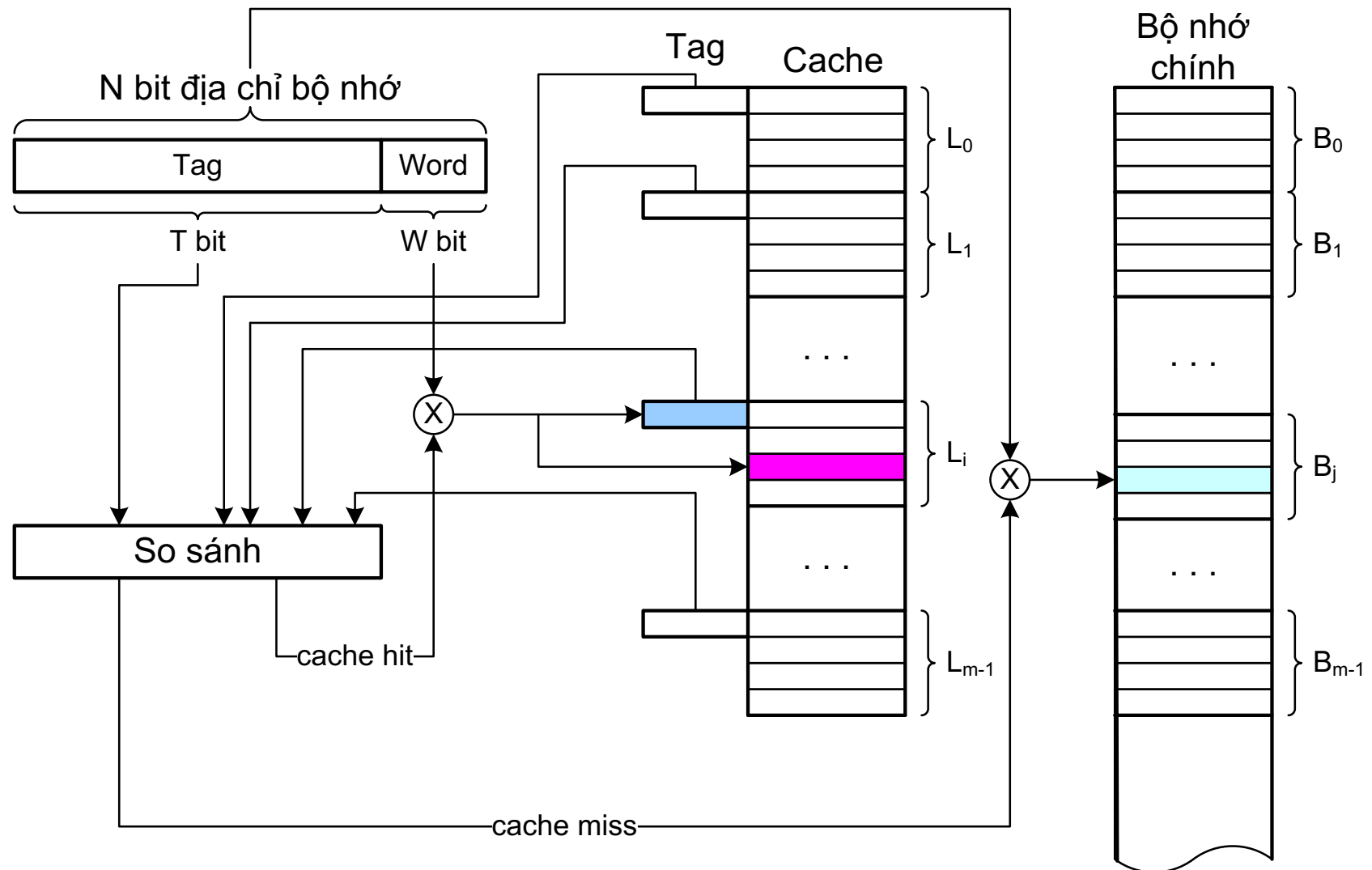
Ánh xạ trực tiếp (tiếp)

- Mỗi thẻ nhớ (Tag) của một Line chứa được T bit
- Khi Block từ bộ nhớ chính được nạp vào Line của cache thì Tag ở đó được cập nhật giá trị là T bit địa chỉ bên trái của Block đó
- Khi CPU muốn truy nhập một từ nhớ thì nó phát ra một địa chỉ N bit cụ thể
 - Nhờ vào giá trị L bit của trường Line sẽ tìm ra Line tương ứng
 - Đọc nội dung Tag ở Line đó (T bit), rồi so sánh với T bit bên trái của địa chỉ vừa phát ra
 - Giống nhau: cache hit
 - Khác nhau: cache miss
- Ưu điểm: Bộ so sánh đơn giản
- Nhược điểm: Xác suất cache hit thấp

Ánh xạ liên kết toàn phần

- Mỗi Block có thể nạp vào bất kỳ Line nào của cache
- Địa chỉ của bộ nhớ chính chia thành hai trường:
 - Trường Word
 - Trường Tag dùng để xác định Block của bộ nhớ chính
- Tag xác định Block đang nằm ở Line đó

Ánh xạ liên kết toàn phần (tiếp)



Ánh xạ liên kết toàn phần (tiếp)

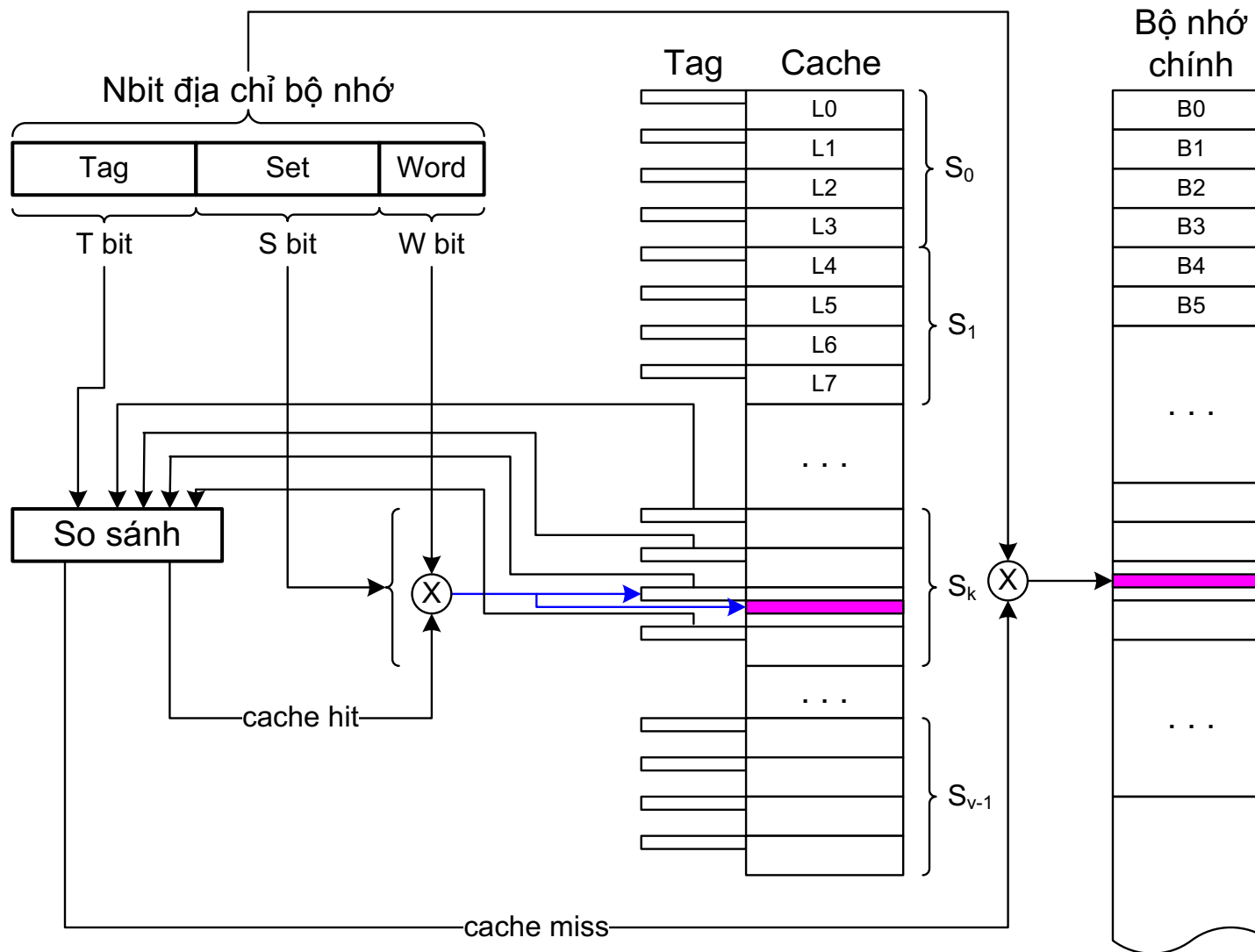
- Mỗi thẻ nhớ (Tag) của một Line chứa được T bit
- Khi Block từ bộ nhớ chính được nạp vào Line của cache thì Tag ở đó được cập nhật giá trị là T bit địa chỉ bên trái của Block đó
- Khi CPU muốn truy nhập một từ nhớ thì nó phát ra một địa chỉ N bit cụ thể
 - So sánh T bit bên trái của địa chỉ vừa phát ra với lần lượt nội dung của các Tag trong cache
 - Nếu gặp giá trị bằng nhau: cache hit xảy ra ở Line đó
 - Nếu không có giá trị nào bằng: cache miss
- Ưu điểm: Xác suất cache hit cao
- Nhược điểm:
 - So sánh đồng thời với tất cả các Tag → mất nhiều thời gian
 - Bộ so sánh phức tạp



Ánh xạ liên kết tập hợp

- Dung hòa cho hai phương pháp trên
- Cache được chia thành các Tập (Set)
- Mỗi một Set chứa một số Line
- Ví dụ:
 - 4 Line/Set \rightarrow 4-way associative mapping
- Ánh xạ theo nguyên tắc sau:
 - $B_0 \rightarrow S_0$
 - $B_1 \rightarrow S_1$
 - $B_2 \rightarrow S_2$
 -

Ánh xạ liên kết tập hợp (tiếp)



Ánh xạ liên kết tập hợp (tiếp)

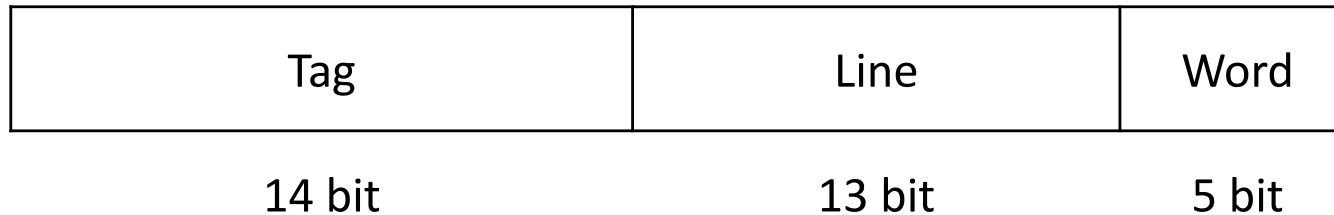
- Kích thước Block = 2^W Word
- Trường Set có S bit dùng để xác định một trong số các Set trong cache. $2^S =$ Số Set trong cache
- Trường Tag có T bit: $T = N - (W+S)$
- Khi CPU muốn truy nhập một từ nhớ thì nó phát ra một địa chỉ N bit cụ thể
 - Nhờ vào giá trị S bit của trường Set sẽ tìm ra Set tương ứng
 - So sánh T bit bên trái của địa chỉ vừa phát ra với lần lượt nội dung của các Tag trong Set đó
 - Nếu gặp giá trị bằng nhau: cache hit xảy ra ở Line tương ứng
 - Nếu không có giá trị nào bằng: cache miss
- Tổng quát cho cả hai phương pháp trên
- Thông dụng với: 2,4,8,16Lines/Set

Ví dụ về ánh xạ địa chỉ

- Giả sử máy tính đánh địa chỉ cho từng byte
- Không gian địa chỉ bộ nhớ chính = 4GiB
- Dung lượng bộ nhớ cache là 256KiB
- Kích thước Line (Block) = 32byte.
- Xác định số bit của các trường địa chỉ cho ba trường hợp tổ chức:
 - Ánh xạ trực tiếp
 - Ánh xạ liên kết toàn phần
 - Ánh xạ liên kết tập hợp 4 đường

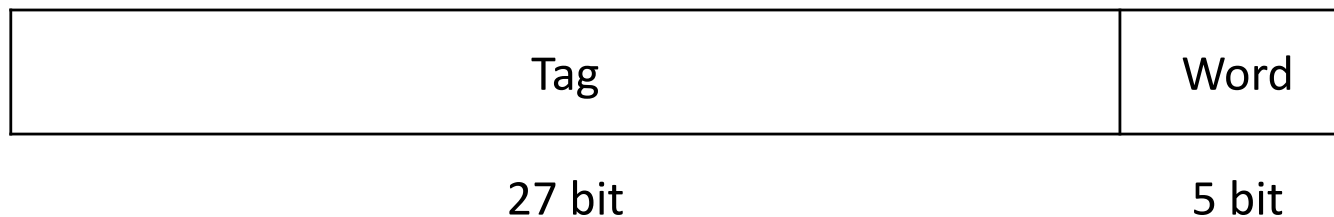
Với ánh xạ trực tiếp

- Bộ nhớ chính = 4GiB = 2^{32} byte \rightarrow Số bit địa chỉ của bộ nhớ chính là: **$N = 32$ bit**
- Cache = 256 KiB = 2^{18} byte
- Kích thước Line = 32 byte = 2^5 byte \rightarrow số bit địa chỉ của trường Word là: **$W = 5$ bit**
- Số Line trong cache = $2^{18} / 2^5 = 2^{13}$ Line \rightarrow số bit địa chỉ trường Line là: **$L = 13$ bit**
- Số bit địa chỉ của trường Tag là:
 $T = 32 - (13 + 5) = 14$ bit



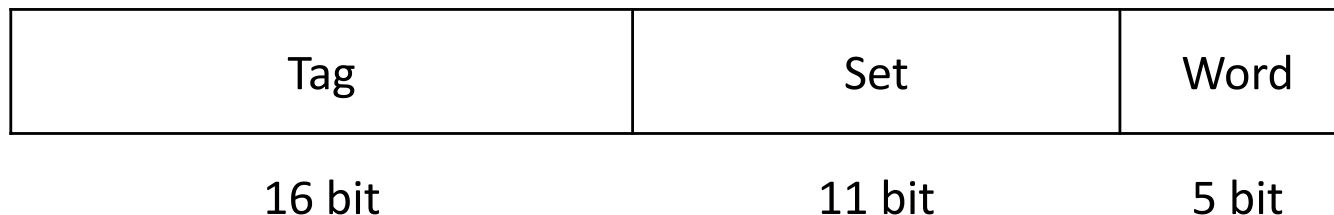
Với ánh xạ liên kết toàn phần

- Bộ nhớ chính = 4GiB = 2^{32} byte \rightarrow số bit địa chỉ của bộ nhớ chính là: $N = 32$ bit
- Kích thước Line = 32 byte = 2^5 byte \rightarrow số bit địa chỉ của trường Word là: $W = 5$ bit
- Số bit địa chỉ của trường Tag là:
 $T = 32 - 5 = 27$ bit



Với ánh xạ liên kết tập hợp 4 đường

- Bộ nhớ chính = 4GiB = 2^{32} byte \rightarrow số bit địa chỉ của bộ nhớ chính là: $N = 32$ bit
- Kích thước Line = 32 byte = 2^5 byte \rightarrow số bit địa chỉ của trường Word là: $W = 5$ bit
- Số Line trong cache = $2^{18} / 2^5 = 2^{13}$ Line
- Một Set có 4 Line = 2^2 Line
 \rightarrow số Set trong cache = $2^{13} / 2^2 = 2^{11}$ Set
 \rightarrow số bit địa chỉ của trường Set là: $S = 11$ bit
- Số bit địa chỉ của trường Tag là:
 $T = 32 - (11 + 5) = 16$ bit



3. Thay thế block trong cache

Với ánh xạ trực tiếp:

- Không phải lựa chọn
- Mỗi Block chỉ ánh xạ vào một Line xác định
- Thay thế Block ở Line đó

Thay thế block trong cache (tiếp)

Với ánh xạ liên kết: cần có thuật giải thay thế:

- **Random**: Thay thế ngẫu nhiên
- **FIFO** (First In First Out): Thay thế Block nào nằm lâu nhất ở trong Set đó
- **LFU** (Least Frequently Used): Thay thế Block nào trong Set có số lần truy nhập ít nhất trong cùng một khoảng thời gian
- **LRU** (Least Recently Used): Thay thế Block ở trong Set tương ứng có thời gian lâu nhất không được tham chiếu tới
- **Tối ưu nhất: LRU**

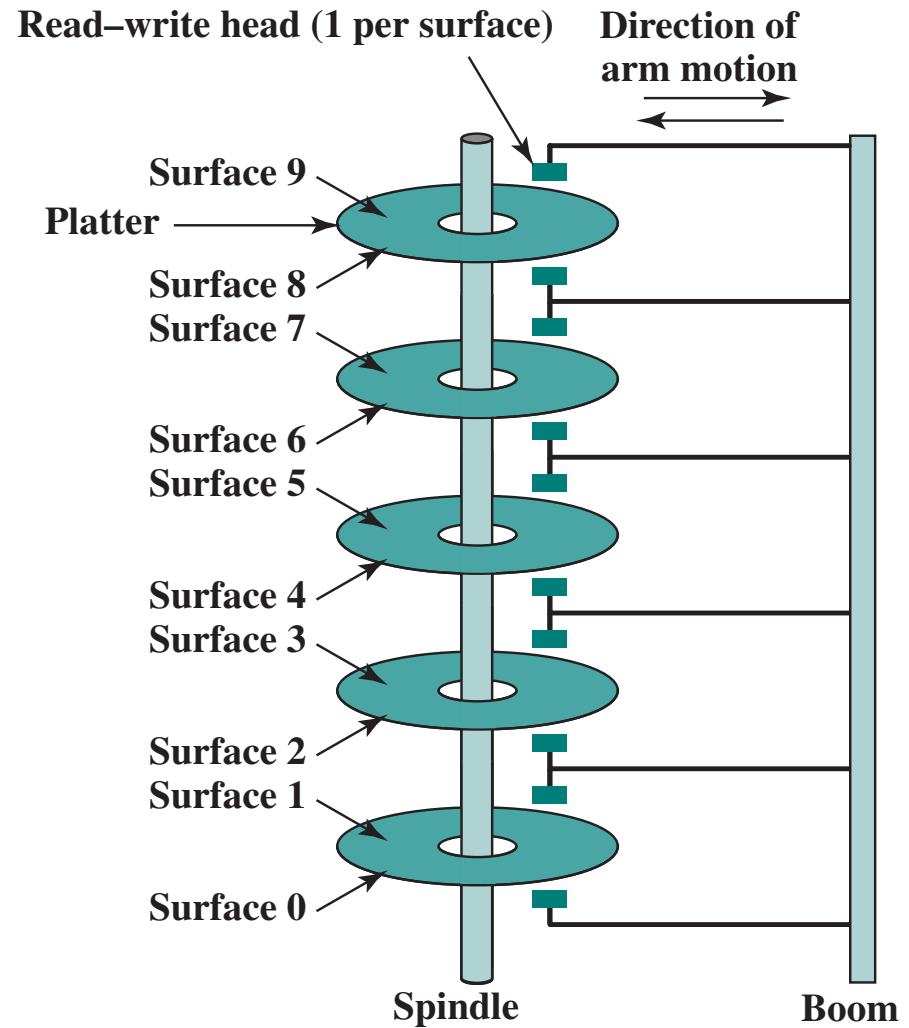
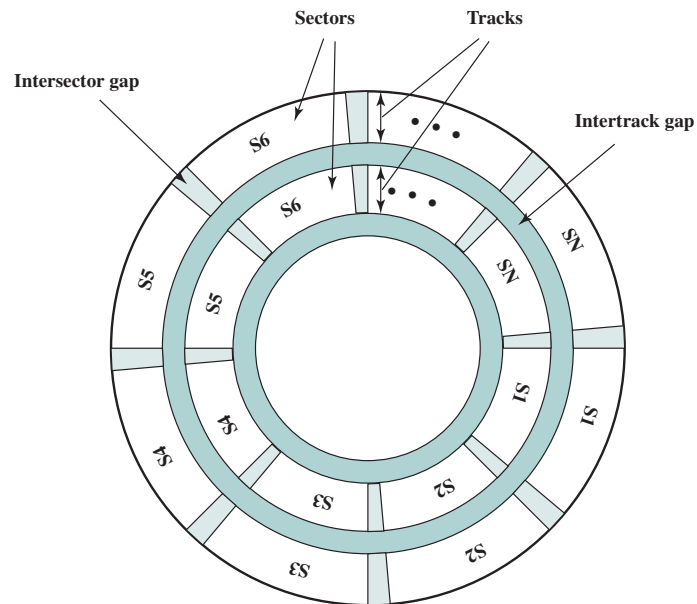
4. Phương pháp ghi dữ liệu khi cache hit

- Ghi xuyên qua (Write-through):
 - ghi cả cache và cả bộ nhớ chính
 - tốc độ chậm
- Ghi trả sau (Write-back):
 - chỉ ghi ra cache
 - tốc độ nhanh
 - khi Block trong cache bị thay thế cần phải ghi trả cả Block về bộ nhớ chính

6.4. Bộ nhớ ngoài

- Tồn tại dưới dạng các thiết bị lưu trữ
- Các kiểu bộ nhớ ngoài
 - Băng từ: ít sử dụng
 - Đĩa từ: Ổ đĩa cứng HDD (Hard Disk Drive)
 - Đĩa quang: CD, DVD
 - Bộ nhớ Flash:
 - Ổ nhớ thể rắn SSD (Solid State Drive)
 - USB flash
 - Thẻ nhớ

Ổ đĩa cứng (HDD – Hard Disk Drive)



- Dung lượng lớn
- Tốc độ đọc/ghi chậm
- Tốn năng lượng
- Dễ bị lỗi cơ học
- Rẻ tiền

Ổ SSD (Solid State Drive)

- Bộ nhớ bán dẫn flash
- Không khả biến
- Tốc độ nhanh
- Tiêu thụ năng lượng ít
- Gồm nhiều chip nhớ flash và cho
- song song
- Ít bị lỗi
- Đắt tiền



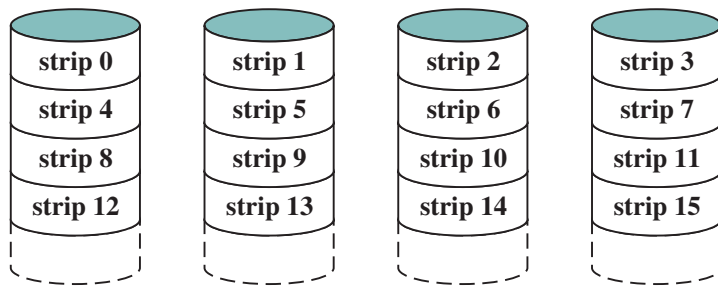
Đĩa quang

- CD (Compact Disc)
 - Dung lượng thông dụng 650MB
- DVD
 - Digital Video Disc hoặc Digital Versatile Disk
 - Ghi một hoặc hai mặt
 - Một hoặc hai lớp trên một mặt
 - Thông dụng: 4,7GB/lớp

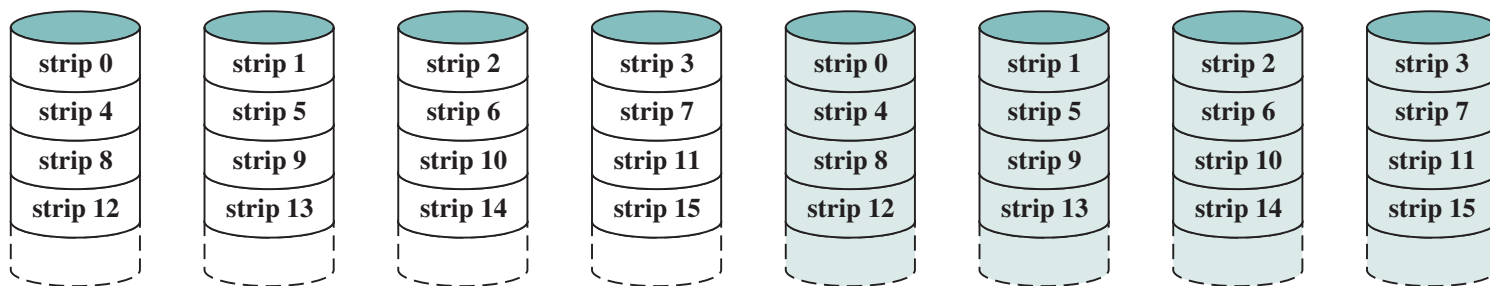
Hệ thống lưu trữ dung lượng lớn: RAID

- Redundant Array of Inexpensive Disks
- (Redundant Array of Independent Disks)
- Tập các ổ đĩa cứng vật lý được OS coi như một ổ logic duy nhất → dung lượng lớn
- Dữ liệu được lưu trữ phân tán trên các ổ đĩa vật lý → truy cập song song (nhanh)
- Lưu trữ thêm thông tin dư thừa, cho phép khôi phục lại thông tin trong trường hợp đĩa bị hỏng → an toàn thông tin
- 7 loại phổ biến (RAID 0 – 6)

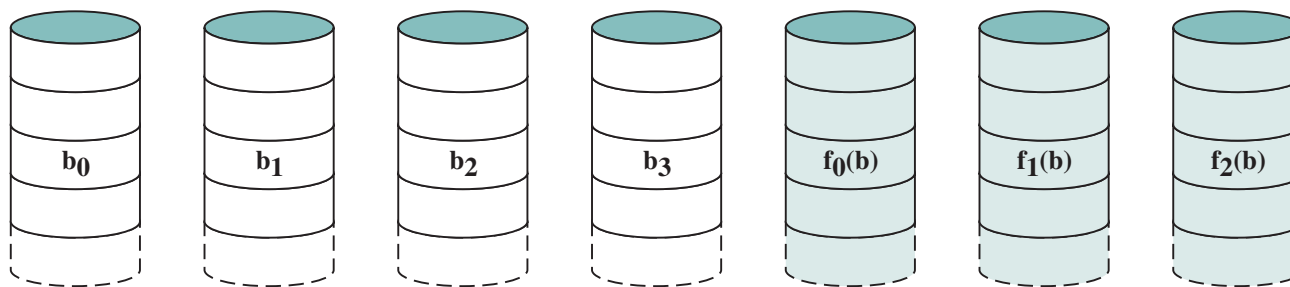
RAID 0, 1, 2



(a) RAID 0 (Nonredundant)

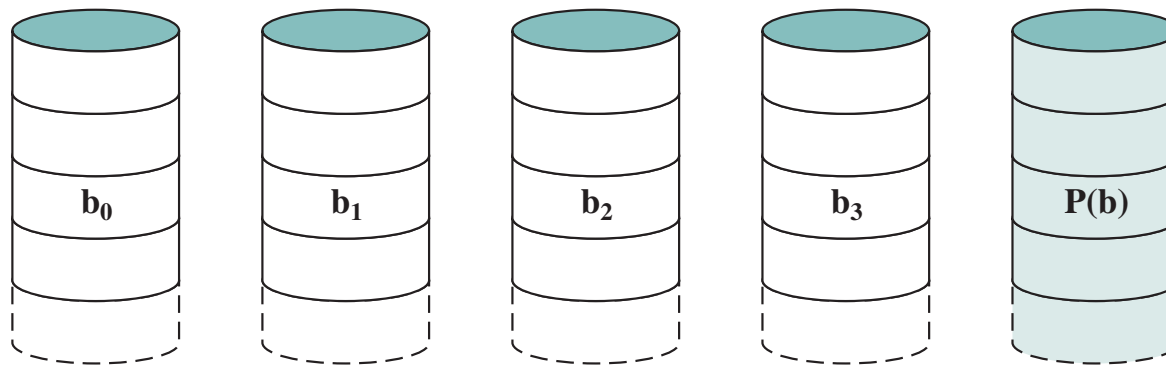


(b) RAID 1 (Mirrored)

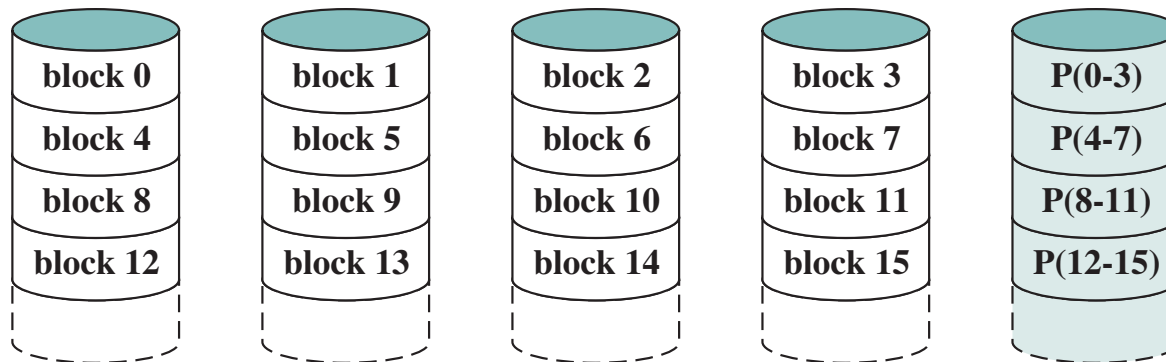


(c) RAID 2 (Redundancy through Hamming code)

RAID 3 & 4

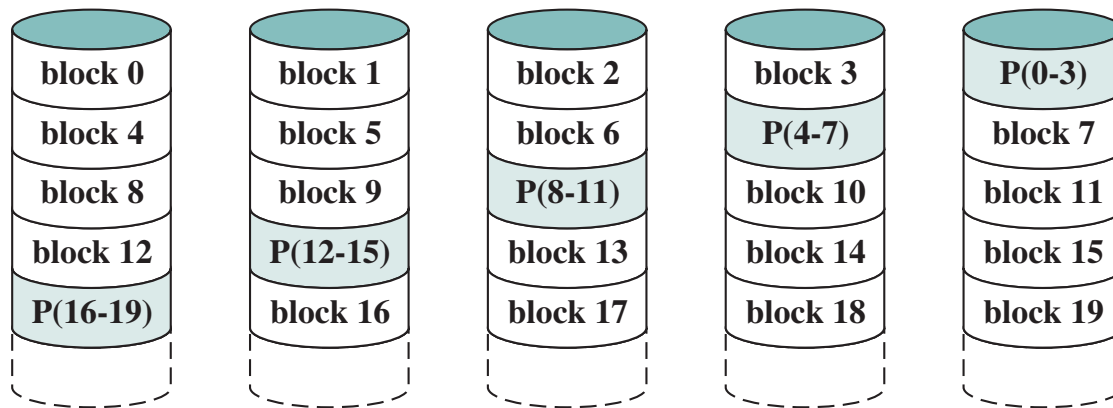


(d) RAID 3 (Bit-interleaved parity)

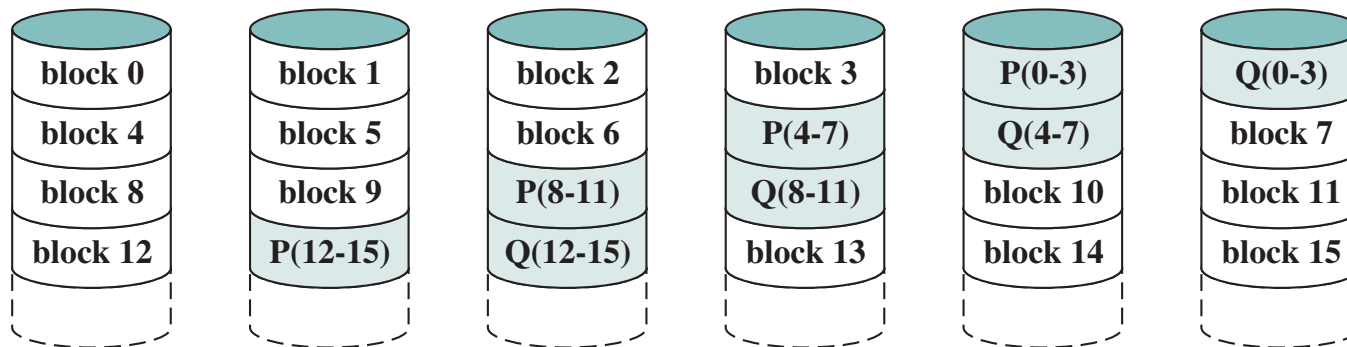


(e) RAID 4 (Block-level parity)

RAID 5 & 6

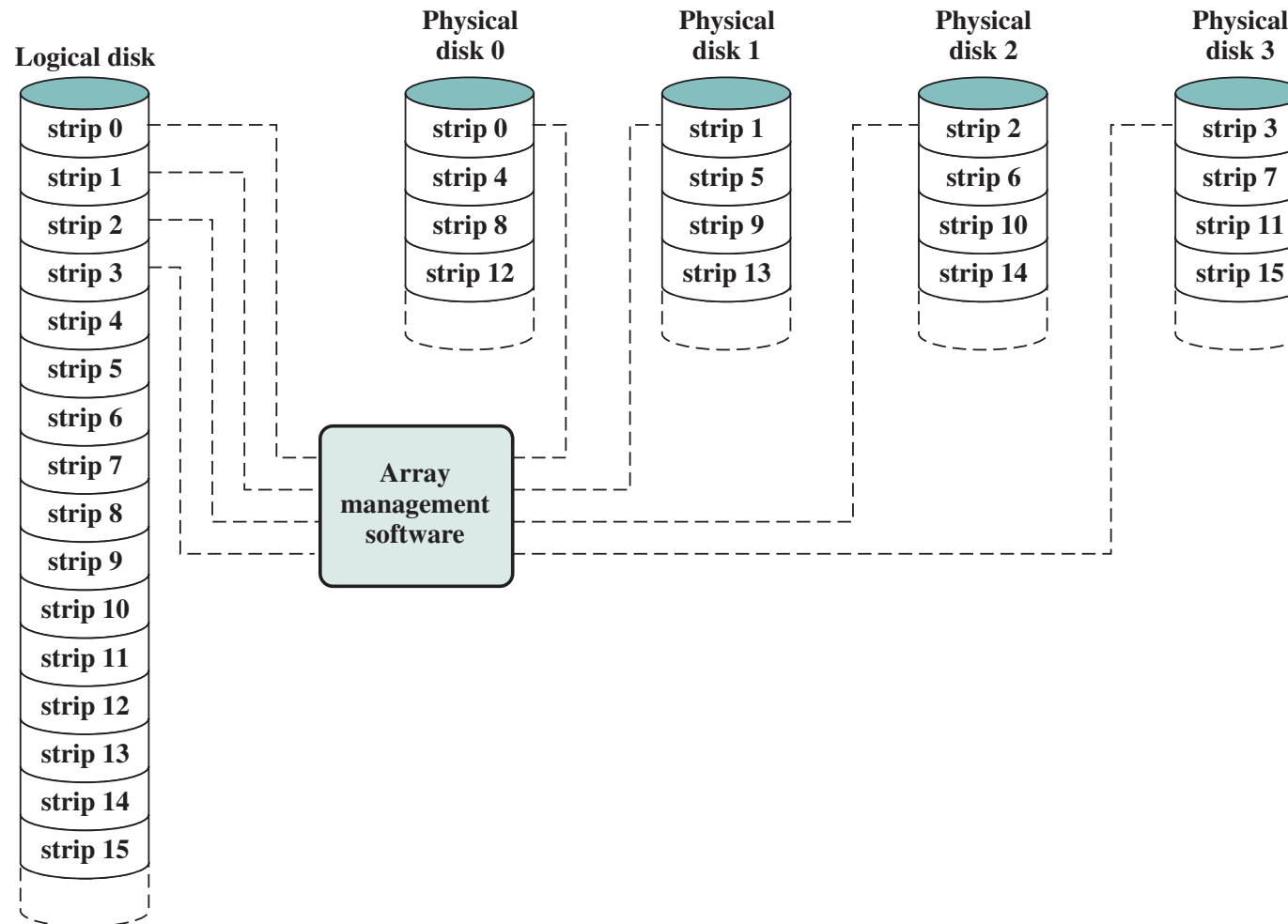


(f) RAID 5 (Block-level distributed parity)



(g) RAID 6 (Dual redundancy)

Ảnh xạ dữ liệu của RAID 0



Hết chương 6