# Arhitecturi Paralele
## Memoria Cache

Dr. Ing. Cristian Chilipirea – cristian.chilipirea@gmail.ro

# Ierarhia de memorii

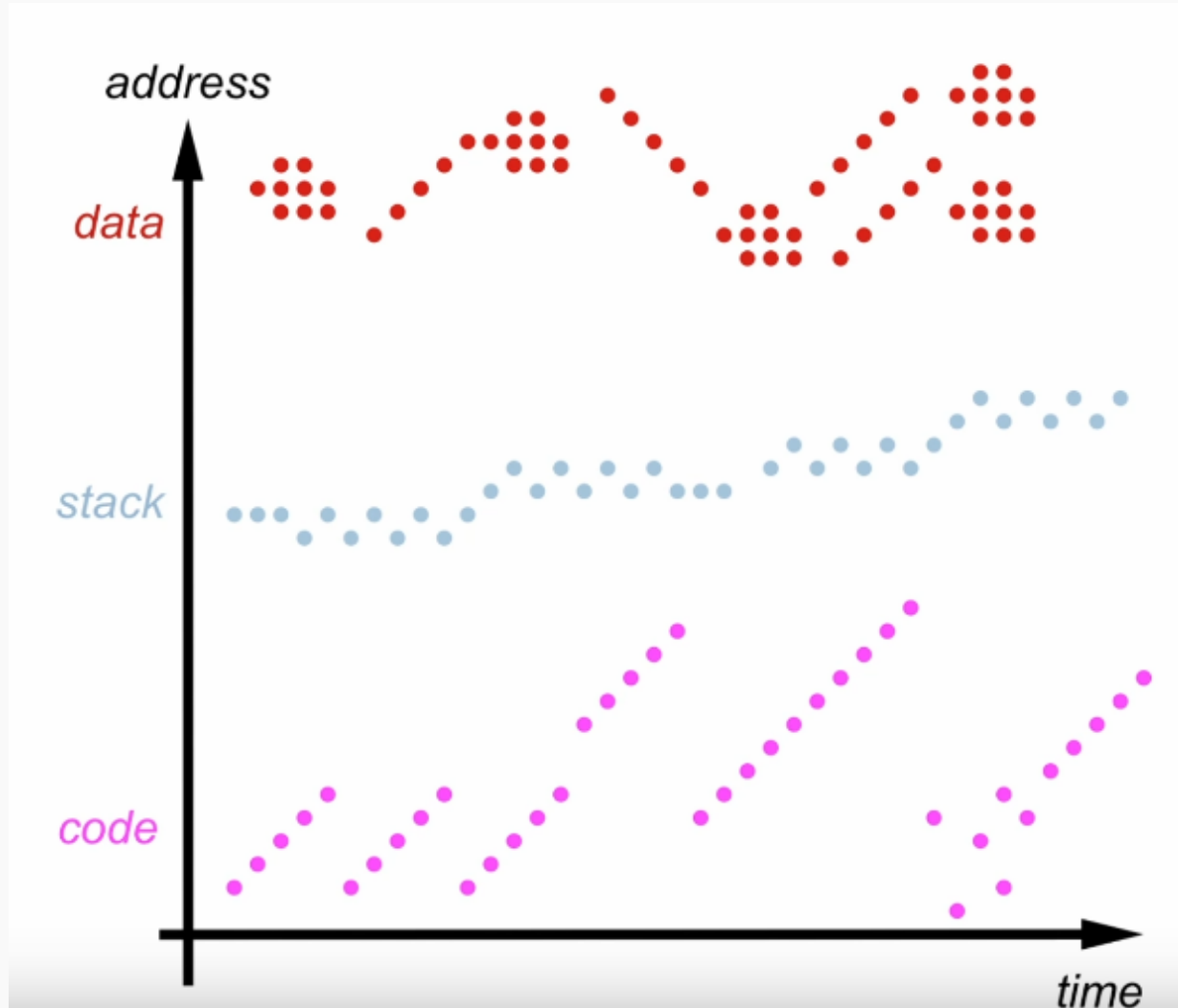| | Capacitate | Latență | Cost/GB | Controlat de |
|---|---|---|---|---|
| Regiștri | x 1 B | 1ns < | - | Compilator |
| Cache (SRAM) | x 10MB | 1-10ns | $5000 | Hardware |
| RAM (DRAM) | x 100GB | 70-100ns | $50 | Hardware/Kernel |
| SSD (Flash) | x 1TB | 7-150µs | $1 | Kernel |
| HDD (Magnetic) | x 19TB | 1-10ms | $0.1 | Kernel |

# Principiul "locality"

- Dacă avem un acces asupra unei date de la locația X, un acces la locația X+ΔX la un moment de timp t+Δt devine mai probabil cu cât ΔX și Δt se apropie de 0.

- **Localitate spațială**: Două zone de memorie apropiate vor fi accesate la un interval de timp apropiat.

- **Localitate temporală**: Dacă o zonă de memorie e accesată, sunt șanse mari să fie accesată din nou.
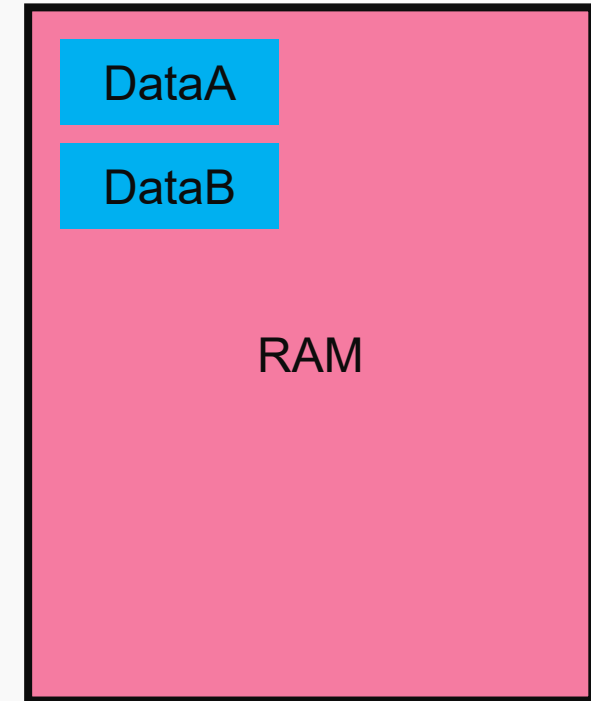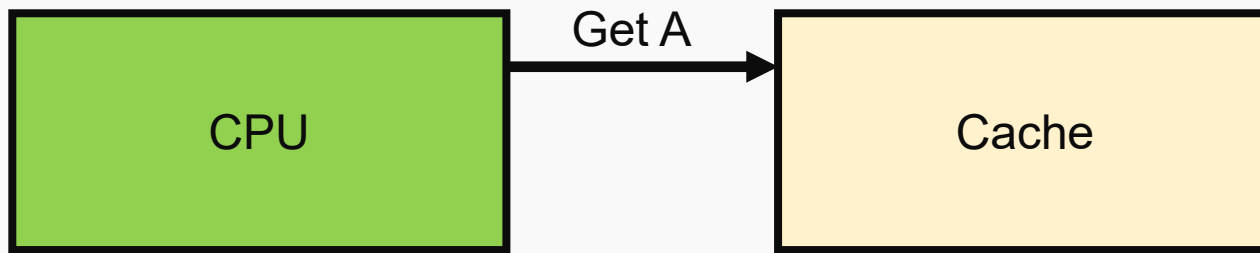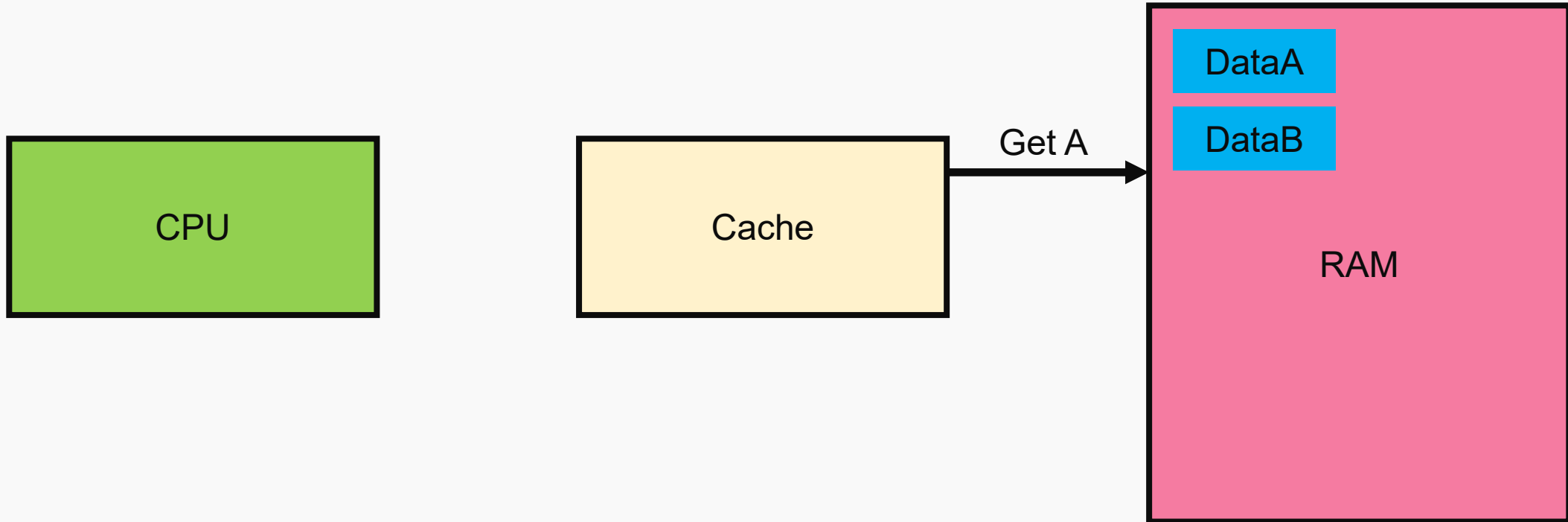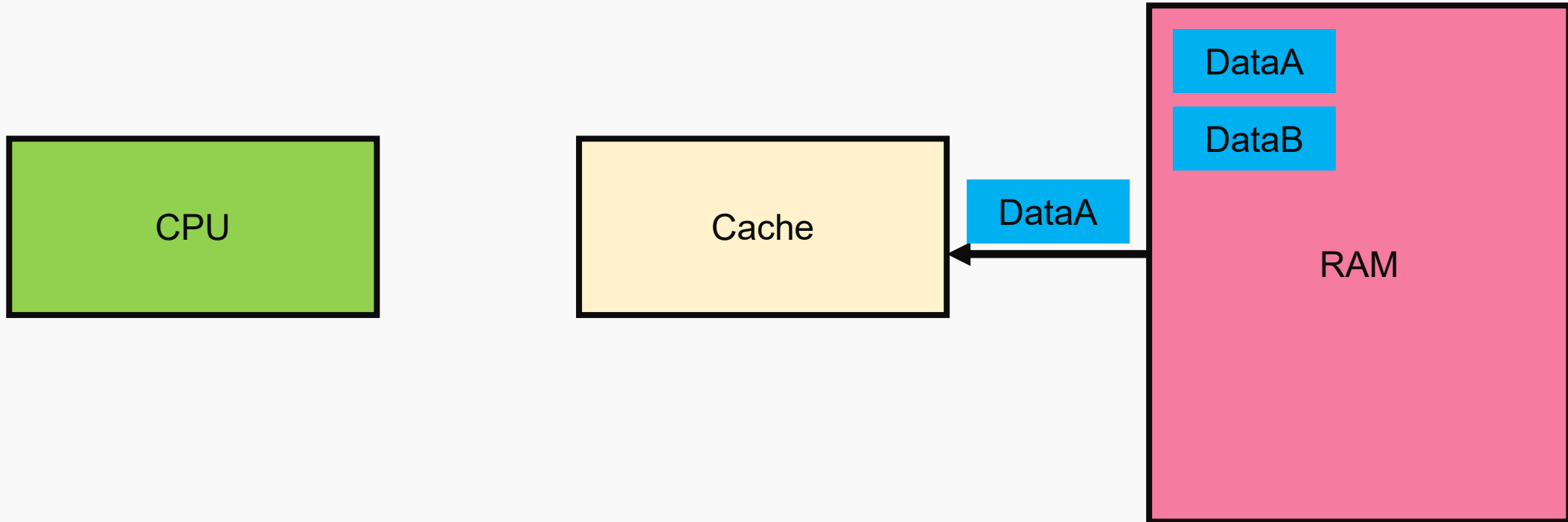
# Pattern acces la memorie

# Utilizare

CPU — Get A → Cache

RAM
- DataA
- DataB

# Utilizare

# Utilizare

# Utilizare
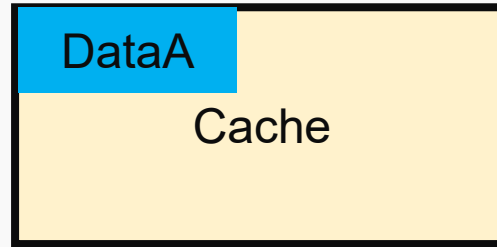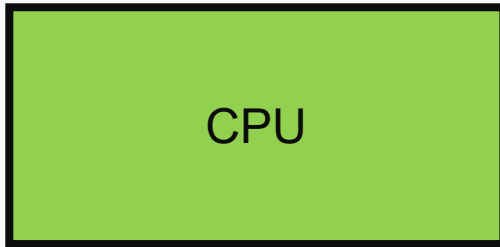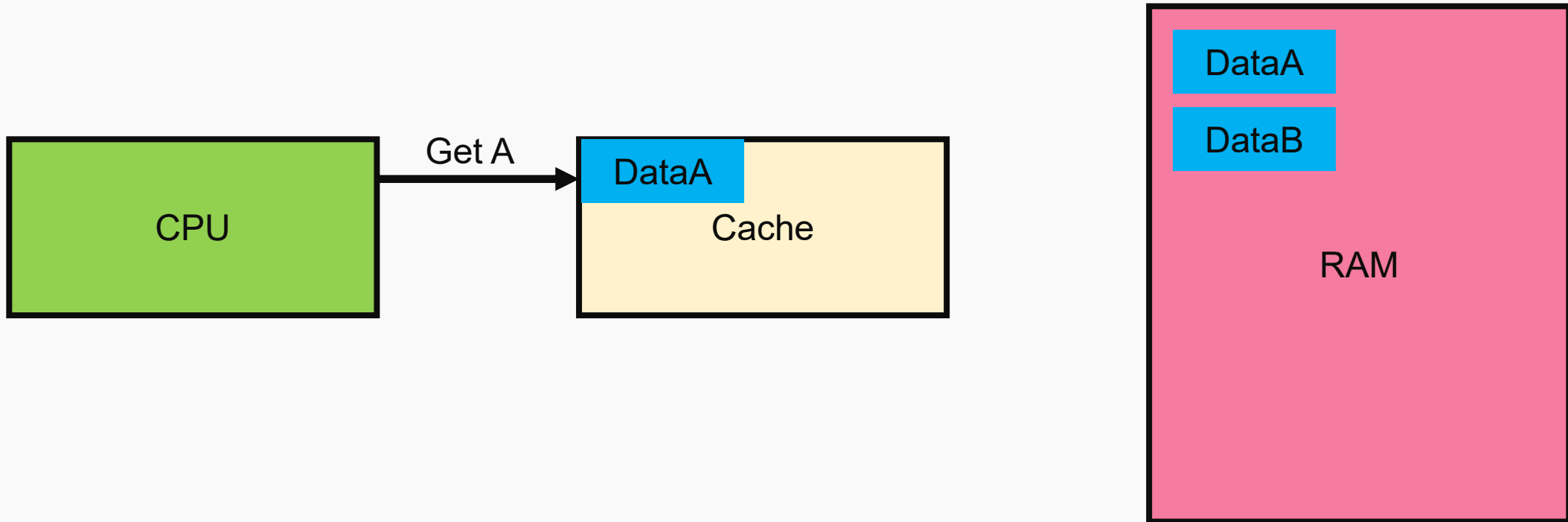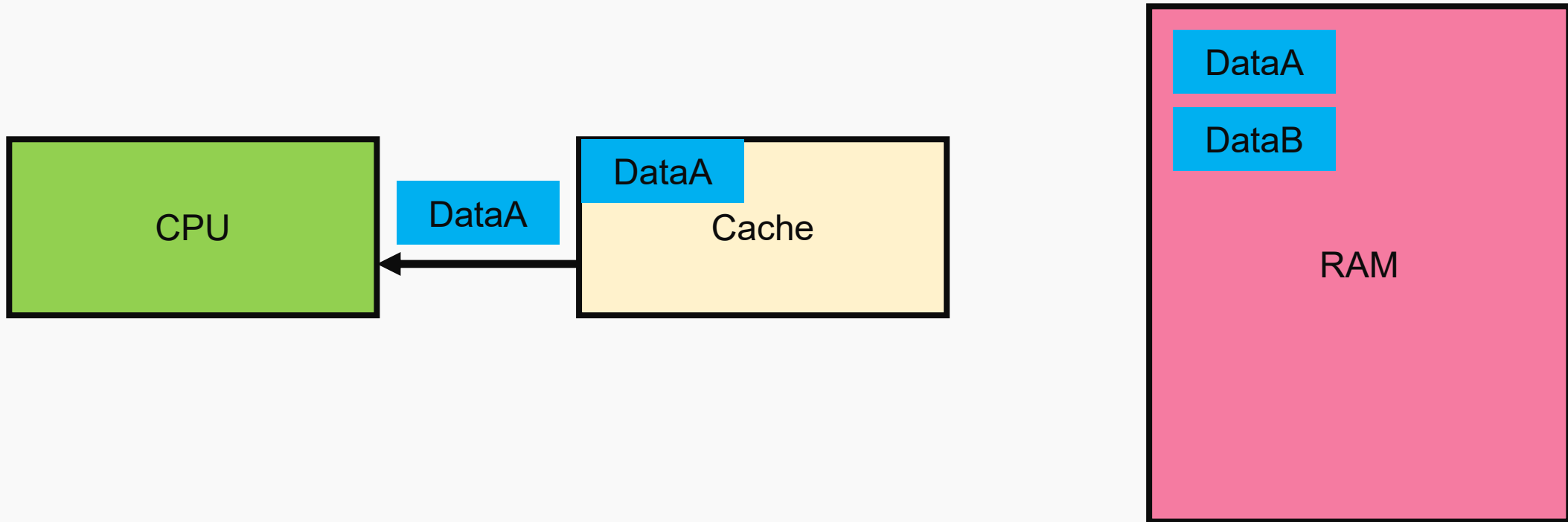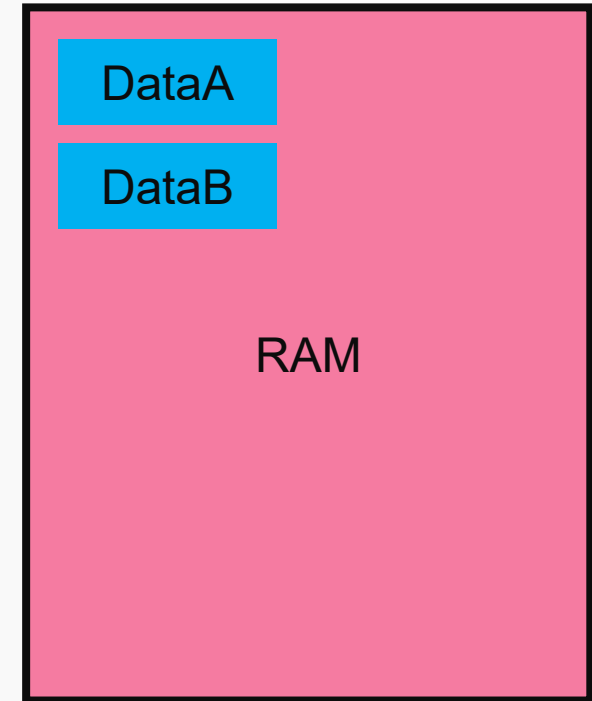


CPU

DataA
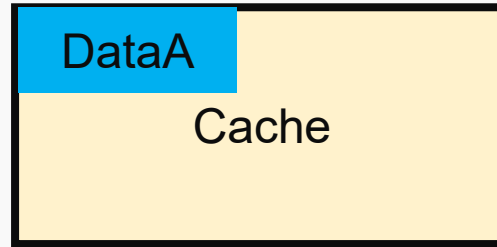Cache

DataA
DataB
RAM

# Utilizare

# Utilizare

# Utilizare – Cache Hit

# Utilizare – Cache Hit

# Utilizare

CPU

DataA
Cache

DataA
DataB
RAM

# Utilizare – Cache **Miss**

# Utilizare – Cache **Miss**

# Utilizare – Cache **Miss**



CPU

DataA
Cache

DataB

DataA
DataB

RAM
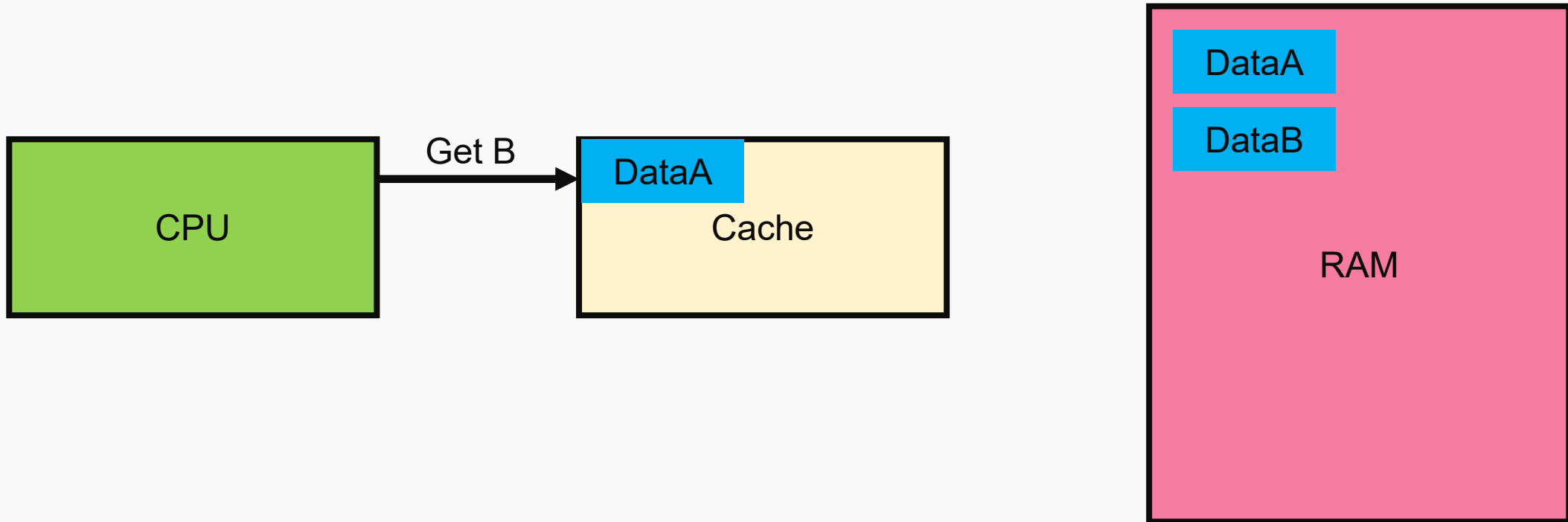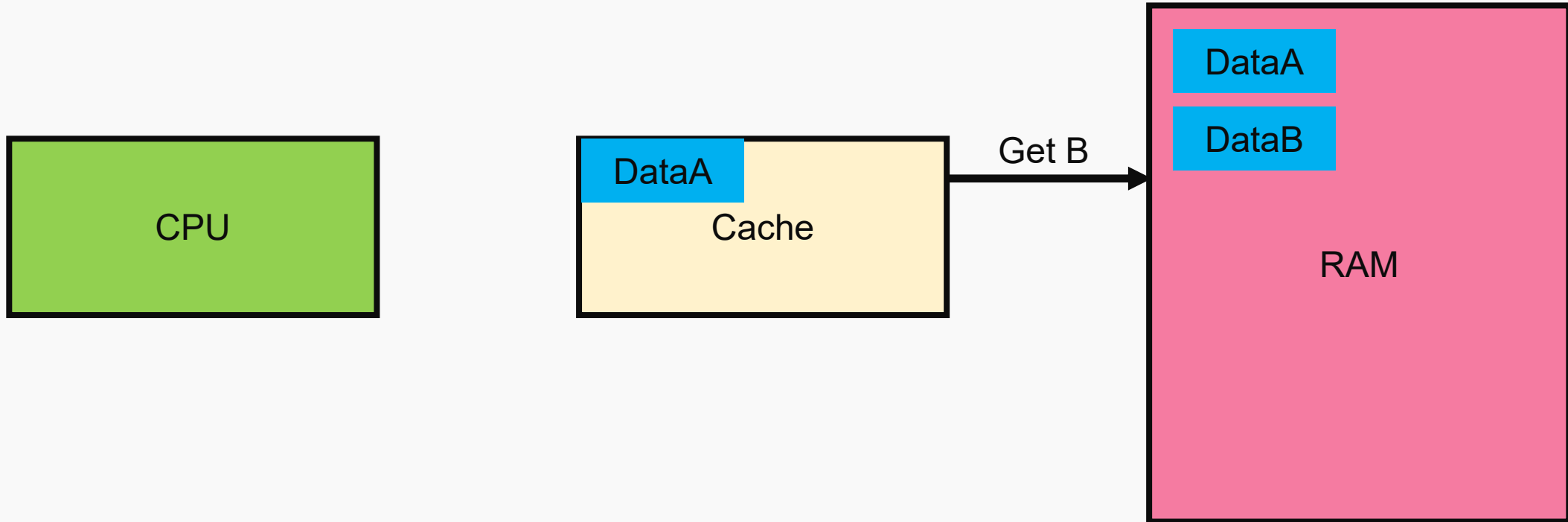
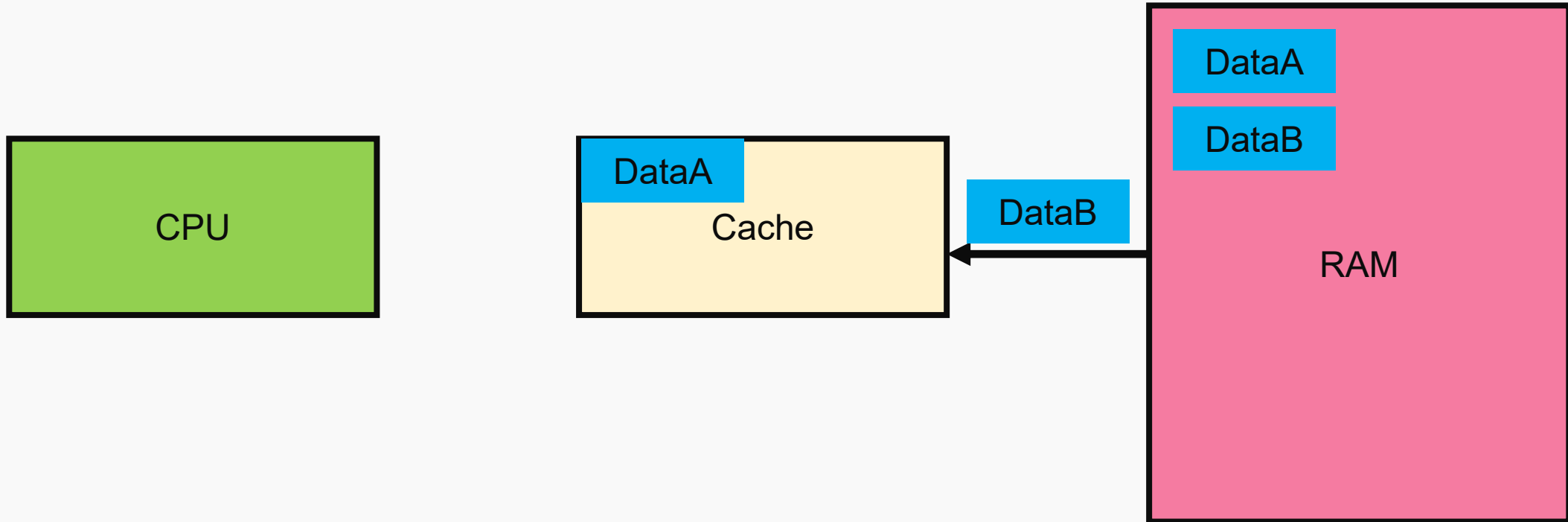# Utilizare – Cache **Miss**

# Utilizare – Cache **Miss**

# Hit/Miss Ratio

- Hit Ratio: $HR = \dfrac{hits}{hits+misses} = 1 - MR$

- Miss Ratio: $MR = \dfrac{misses}{hits+misses} = 1 - HR$

- Average Memory Access Time:
$$AMAT = HitTime + MissRatio * MissPenalty$$

# Direct-Mapped Cache

Valid bit      Tag (40 bits)      Data (64 byte - lines)

Tag bits (40)      Index bits (18)      Offset bits (6)

# Problemă: Cache Thrashing



Se vor scrie în aceeași zonă din Cache

# Matrici - Utilizare Cache Bună

# Fully-Associative Cache

Valid bit    Tag (58 bits)    Data (64 byte - lines)

Tag bits (58)

Offset bits (6)

# N (20) Way-Associative Cache



Tag (45 bits)  Data (64 byte - lines)  Tag (45 bits) Data (64 byte - lines)  Tag (45 bits)  Data (64 byte - lines)

Tag bits (45)

Index bits (13)

Offset bits (6)

# Associativity Implies Choices

**Direct-mapped**

address

- Compare addr with only one tag

- Location A can be stored in exactly one cache line

**N-way set-associative**

address

N

- Compare addr with N tags simultaneously

- Location A can be stored in exactly one set, but in any of the N cache lines belonging to that set

**Fully associative**

address

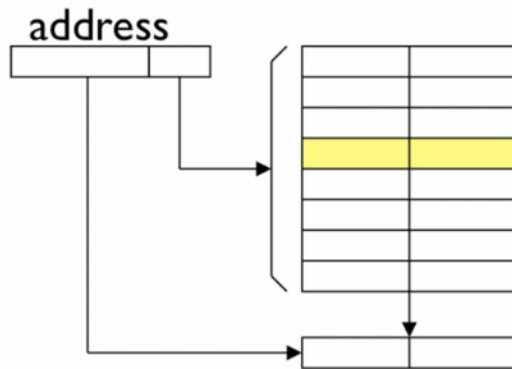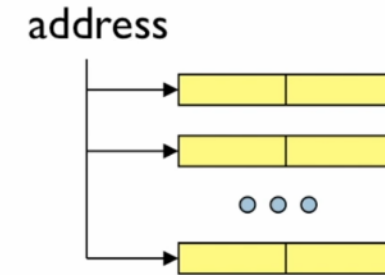- Compare addr with each tag simultaneously

- Location A can be stored in any cache line

# Cache replacement policies

- Random
- FIFO/LIFO
- Least Recently Used/Most recently used
- Least Frequently Used <<<<< LRU approximations
- Bélády (clairvoyant)

# Write Policy

Write-through: CPU writes are cached, but also written to main memory immediately (stalling the CPU until write is completed). Memory always holds current contents
- Simple, slow, wastes bandwidth

Write-behind: CPU writes are cached; writes to main memory may be buffered. CPU keeps executing while writes are completed in the background
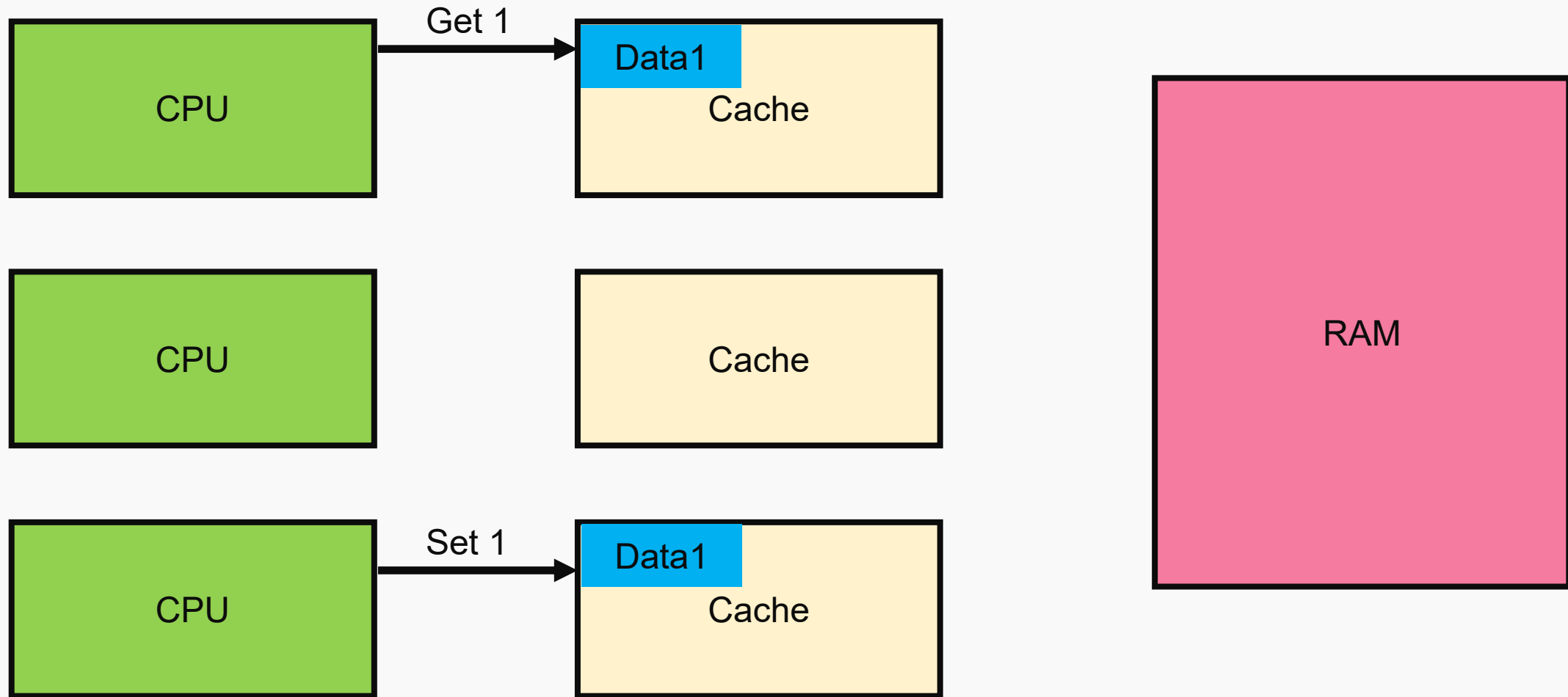- Faster, still uses lots of bandwidth

Write-back: CPU writes are cached, but not written to main memory until we replace the block. Memory contents can be "stale"
- Fastest, low bandwidth, more complex
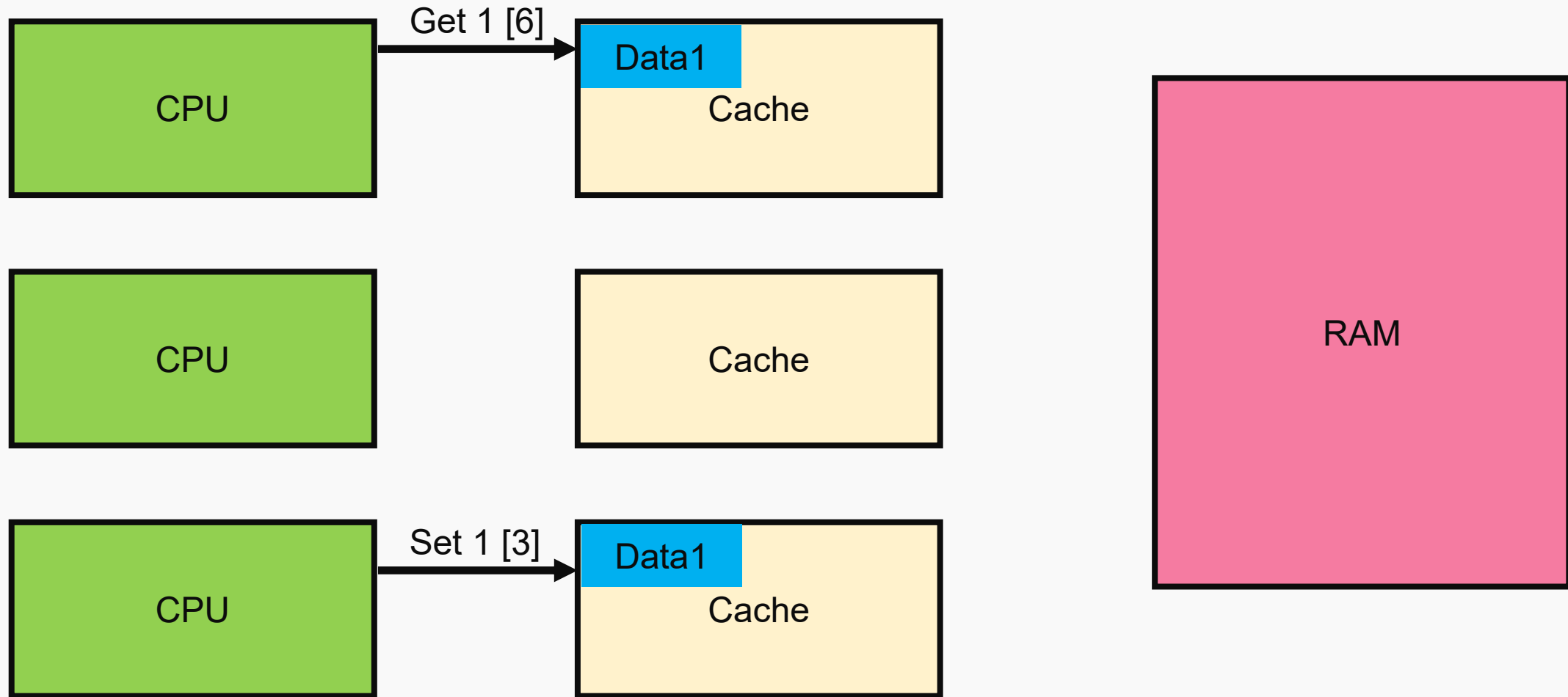- Commonly implemented in current systems

# False sharing

# False sharing – gets worse – cache lines

# False sharing