# Multicore Systems
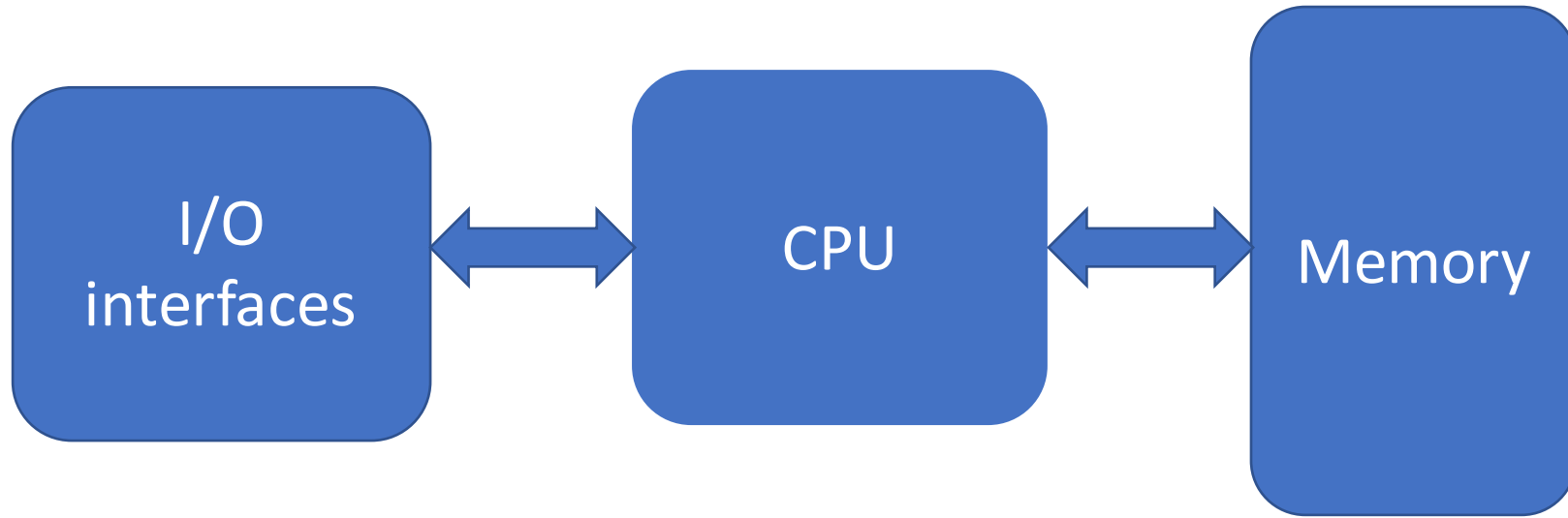
Sujoy Sinha Roy

School of Computer Science

University of Birmingham
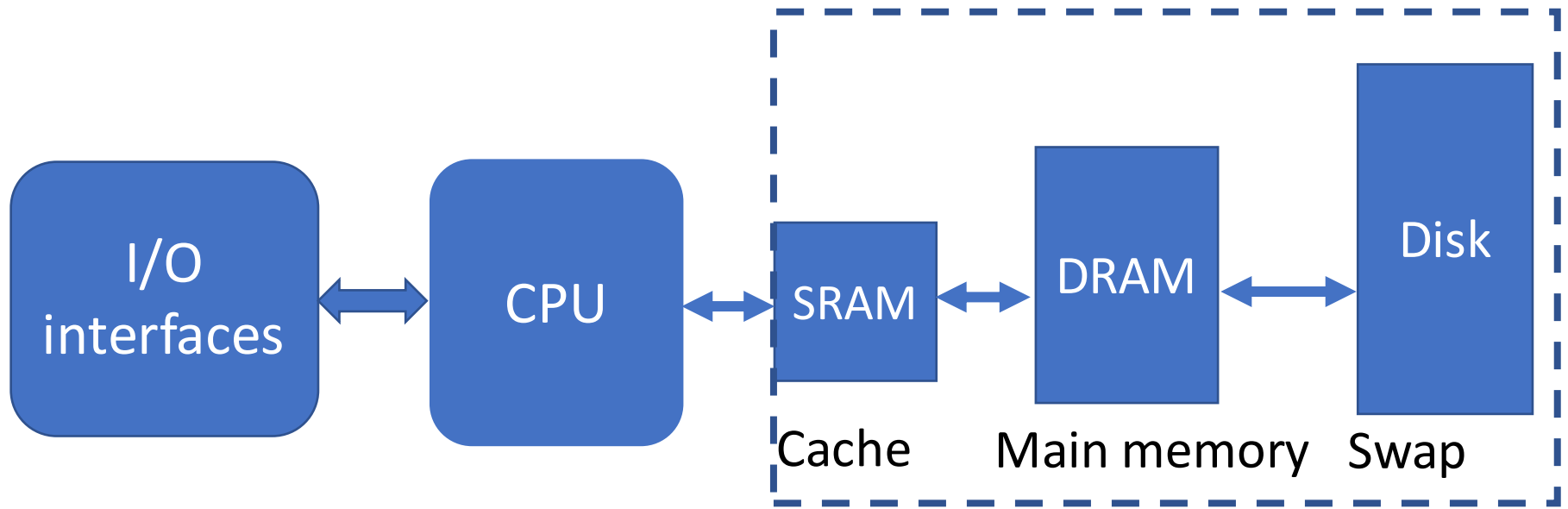
# Recap: Simple von Neumann Computer

```
┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│     I/O      │ <──> │     CPU      │ <──> │    Memory    │
│  interfaces  │      │              │      │              │
└──────────────┘      └──────────────┘      └──────────────┘
```

We started C programming assuming a simple von Neumann computation model.

→ We executed program on a single CPU (or processor) step-by-step.

→ We stored data and program instructions in the Memory.

# Recap: Memory Hierarchy



We studied memory hierarchy from a programmer's perspective
→ Learnt how to write efficient C code to exploit memory hierarchy
→ Programs are still executed on a single CPU (or processor)
   step-by-step but more efficiently

| Produced | From May 7, 1997 to December 26, 2003[1] |
|---|---|
| Common manufacturer(s) | Intel |
| Max. CPU clock rate | 233 MHz to 450 MHz |
| FSB speeds | 66 MHz to 100 MHz |
| Min. feature size | 0.35 μm to 0.18 μm |
| Instruction set | IA-32, MMX |

| | |
|---|---|
| **Produced** | From February 26, 1999 to May 18, 2007[1] |
| **Common manufacturer(s)** | Intel |
| **Max. CPU clock rate** | 450 MHz to 1.4 GHz |
| **FSB speeds** | 100 MHz to 133 MHz |
| **Min. feature size** | 0.25 μm to 0.13 μm |
| **Instruction set** | IA-32, MMX, SSE |

Clock frequency almost doubled w.r.t Pentium II

| Produced | From November 20, 2000 to August 8, 2008 |
| --- | --- |
| Max. CPU clock rate | 1.3 GHz to 3.8 GHz |
| FSB speeds | 400 MT/s to 1066 MT/s |
| Instruction set | x86 (i386), x86-64 (only some chips), MMX, SSE, SSE2, SSE3 |

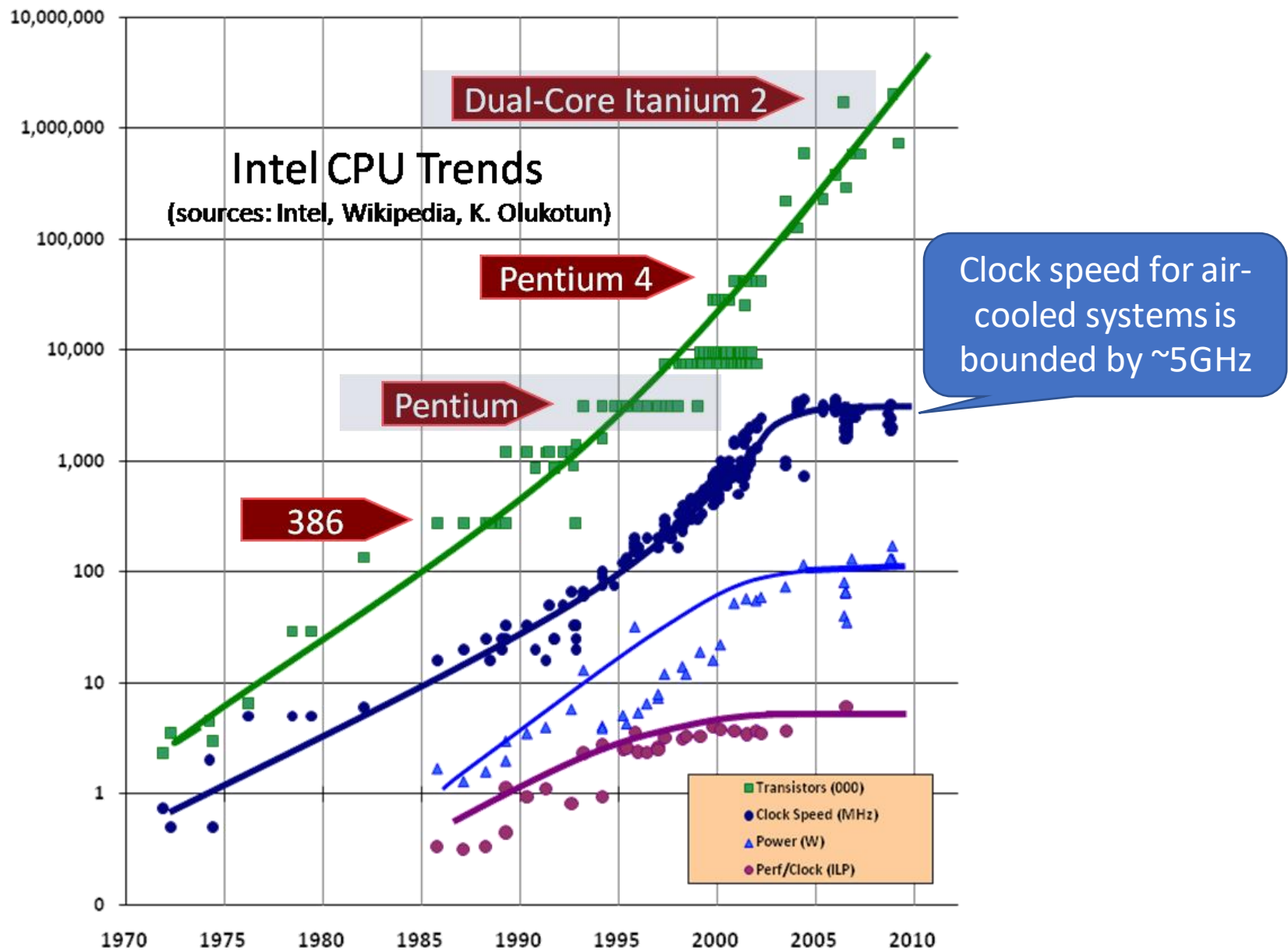Clock frequency almost doubled w.r.t Pentium III
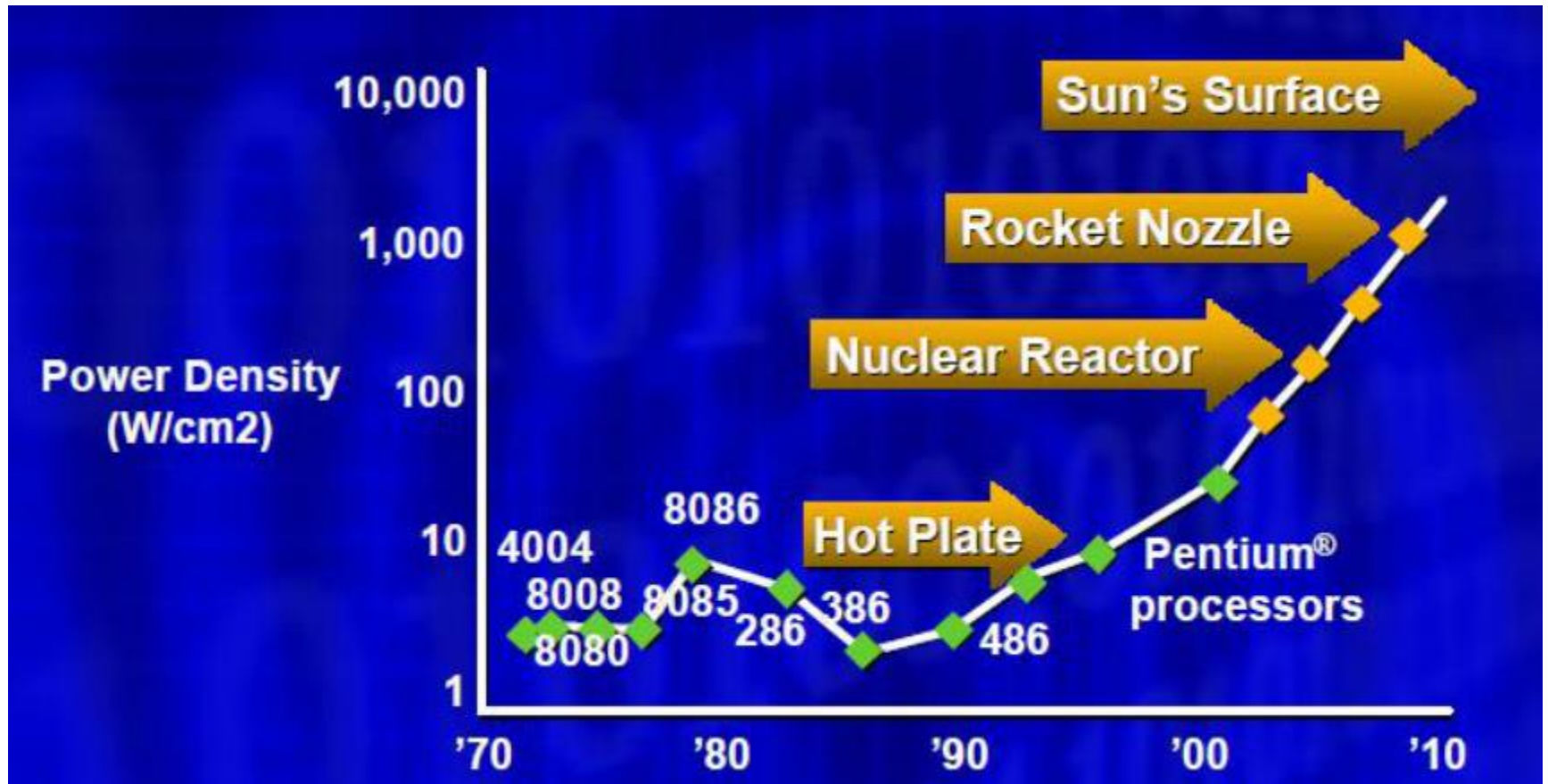
**What is the clock frequency in 2020?**

Intel® Core™ i7-10510U Processor (8M Cache, Up to 4.80 GHz)

- 8 MB Intel® Smart Cache Cache
- 4 Cores
- 8 Threads
- 4.90 GHz Max Turbo Frequency
- U - Ultra-low power
- 10th Generation

Little increase in clock frequency after 10 years since Pentium 4

Source: Herb Sutter "The Free Lunch Is Over: A Fundamental Turn Toward Concurrency in Software" available at http://www.gotw.ca/publications/concurrency-ddj.htm

Source: Patrick Gelsinger, Intel Developer's Forum, Intel Corporation, 2004.

# Applications are demanding more resources!

Alternative solution:

- Put many processing cores on the microprocessor chip.

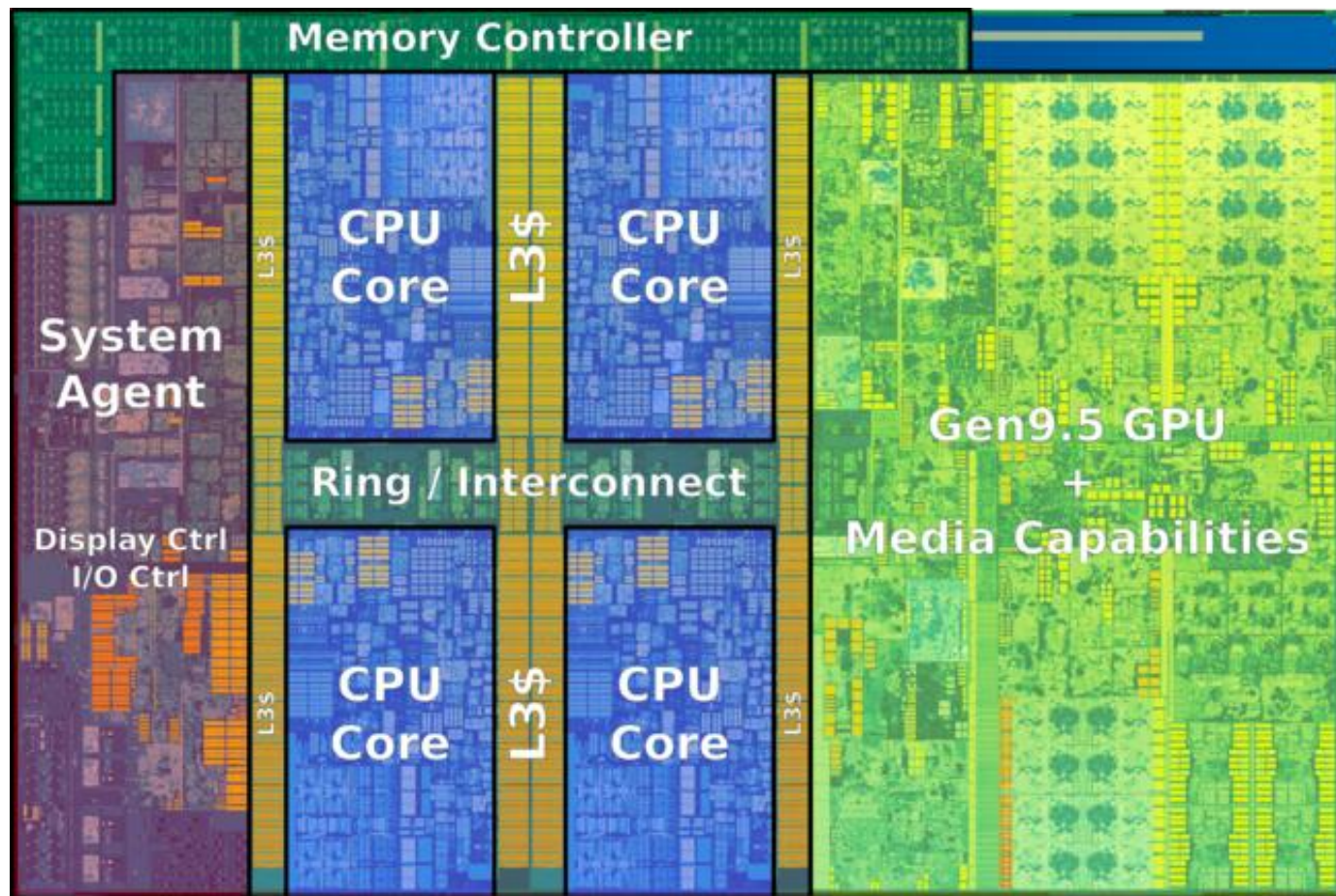- The number of cores doubles with each generation.

Photo of the quad-core GT2 Kaby Lake processor chip. Used in mainstream desktop computers.

It has 4 CPU cores to compute in parallel.

## Samsung Galaxy S10+

Exynos 9 Series (9820)

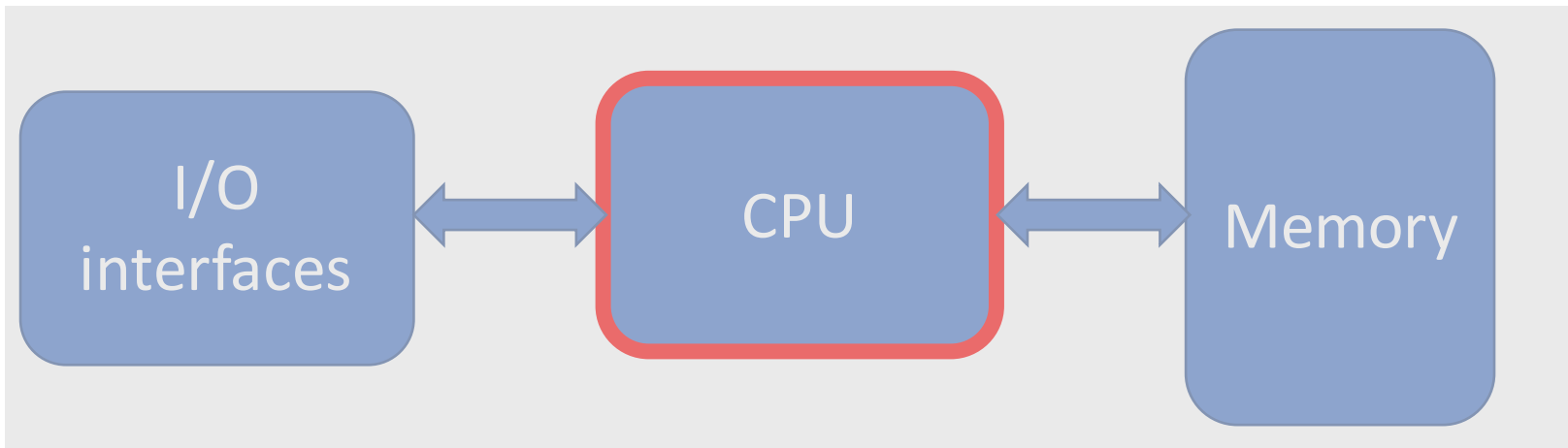6.4" Quad HD+ Curved Dynamic AMOLED (3040x1440)

Android 9 (Pie)
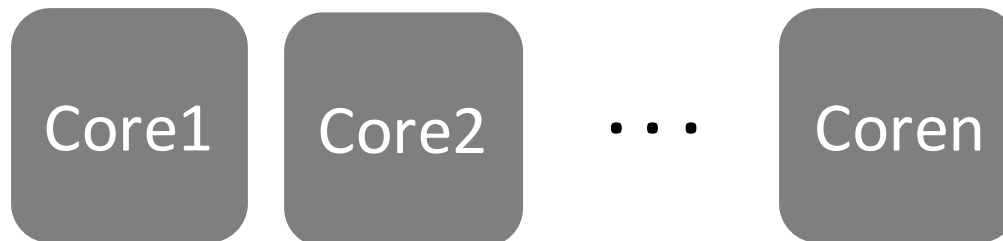
8/12GB RAM, 128/256GB/1TB Storage, microSD(up to 512GB)

Rear: 12MP AF(Wide)+12MP AF(Tele)+16MP FF(Ultra Wide),
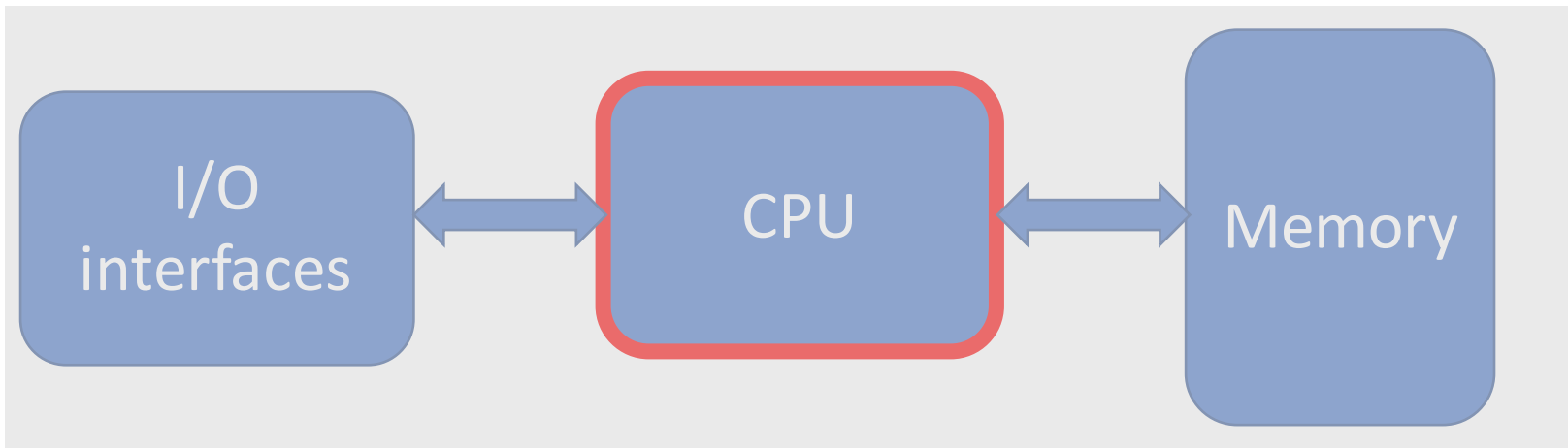Front: 10MP AF+8MP FF(RGB Depth)

# Superior Performance for
# Seamless Multi-tasking

Featuring a 4[th] generation custom CPU, the Exynos 9820's innovative tri-cluster architecture delivers premium processing power. The CPU consists of two custom cores for ultimate processing power, two Cortex-A75 cores for optimal performance, and four Cortex-A55 cores for greater efficiency, resulting in superior performance that lasts. Tri-cluster with intelligent task scheduler boosts multi-core performance by 15 percent when compared to the Exynos 9810, while the 4[th] generation custom CPU with enhanced memory access capability and cutting-edge architecture design improves single core performance by up to 20 percent or boosts power efficiency by up to 40 percent.[1]
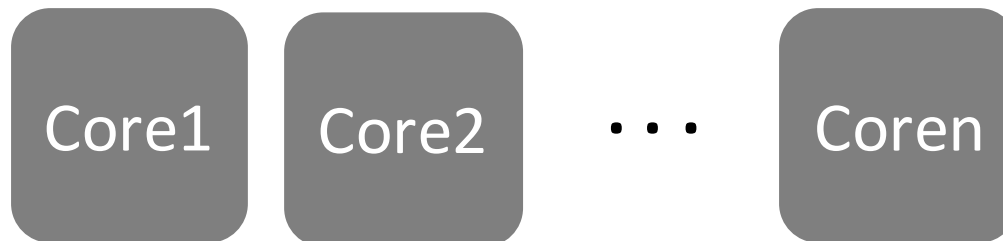
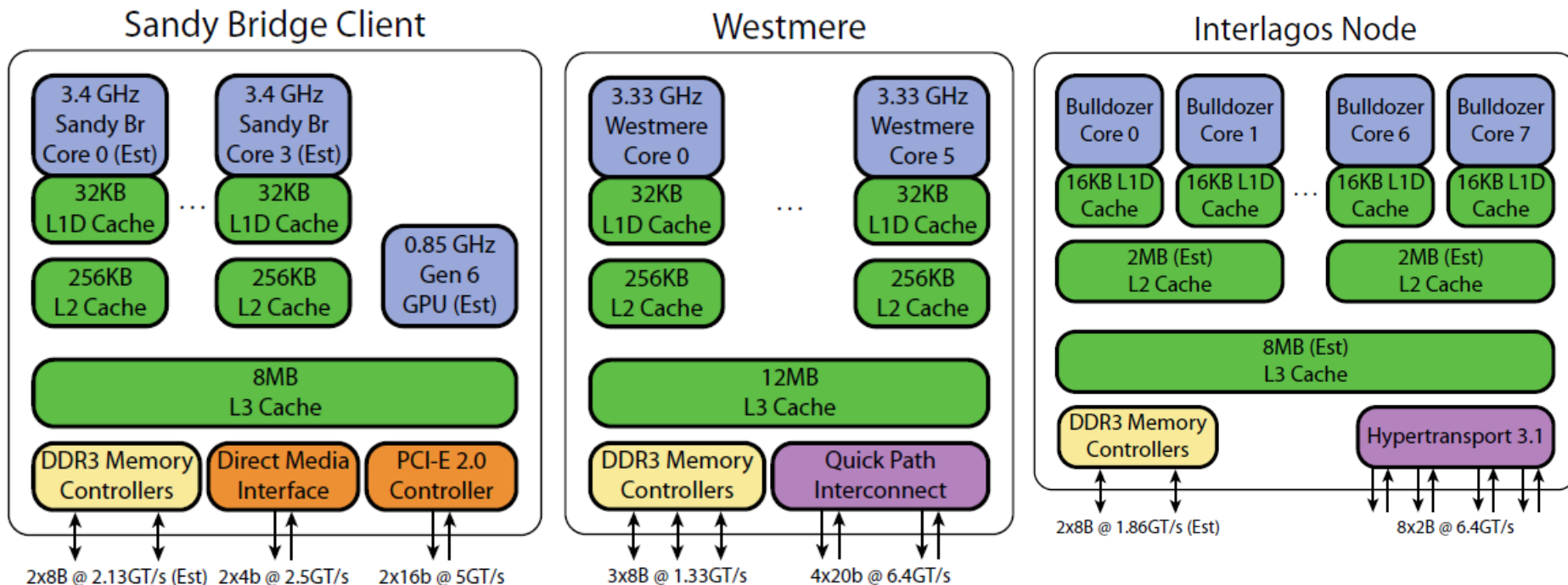https://www.samsung.com/semiconductor/minisite/exynos/products/mobileprocessor/exynos-9-series-9820/

I/O interfaces ⟷ CPU ⟷ Memory

From Single Core
To
Multi Core

Core1　Core2　...　Coren

# Examples of Intel Multicore Processors

## Sandy Bridge Client

| 3.4 GHz Sandy Br Core 0 (Est) | 3.4 GHz Sandy Br Core 3 (Est) | |
| 32KB L1D Cache | ... 32KB L1D Cache | 0.85 GHz Gen 6 GPU (Est) |
| 256KB L2 Cache | 256KB L2 Cache | |

8MB L3 Cache

| DDR3 Memory Controllers | Direct Media Interface | PCI-E 2.0 Controller |

2x8B @ 2.13GT/s (Est)   2x4b @ 2.5GT/s   2x16b @ 5GT/s

## Westmere

| 3.33 GHz Westmere Core 0 | 3.33 GHz Westmere Core 5 |
| 32KB L1D Cache | ... 32KB L1D Cache |
| 256KB L2 Cache | 256KB L2 Cache |

12MB L3 Cache

| DDR3 Memory Controllers | Quick Path Interconnect |

3x8B @ 1.33GT/s   4x20b @ 6.4GT/s

## Interlagos Node

| Bulldozer Core 0 | Bulldozer Core 1 | Bulldozer Core 6 | Bulldozer Core 7 |
| 16KB L1D Cache | 16KB L1D Cache | ... 16KB L1D Cache | 16KB L1D Cache |
| 2MB (Est) L2 Cache | | 2MB (Est) L2 Cache | |

8MB (Est) L3 Cache

| DDR3 Memory Controllers | Hypertransport 3.1 |

2x8B @ 1.86GT/s (Est)   8x2B @ 6.4GT/s
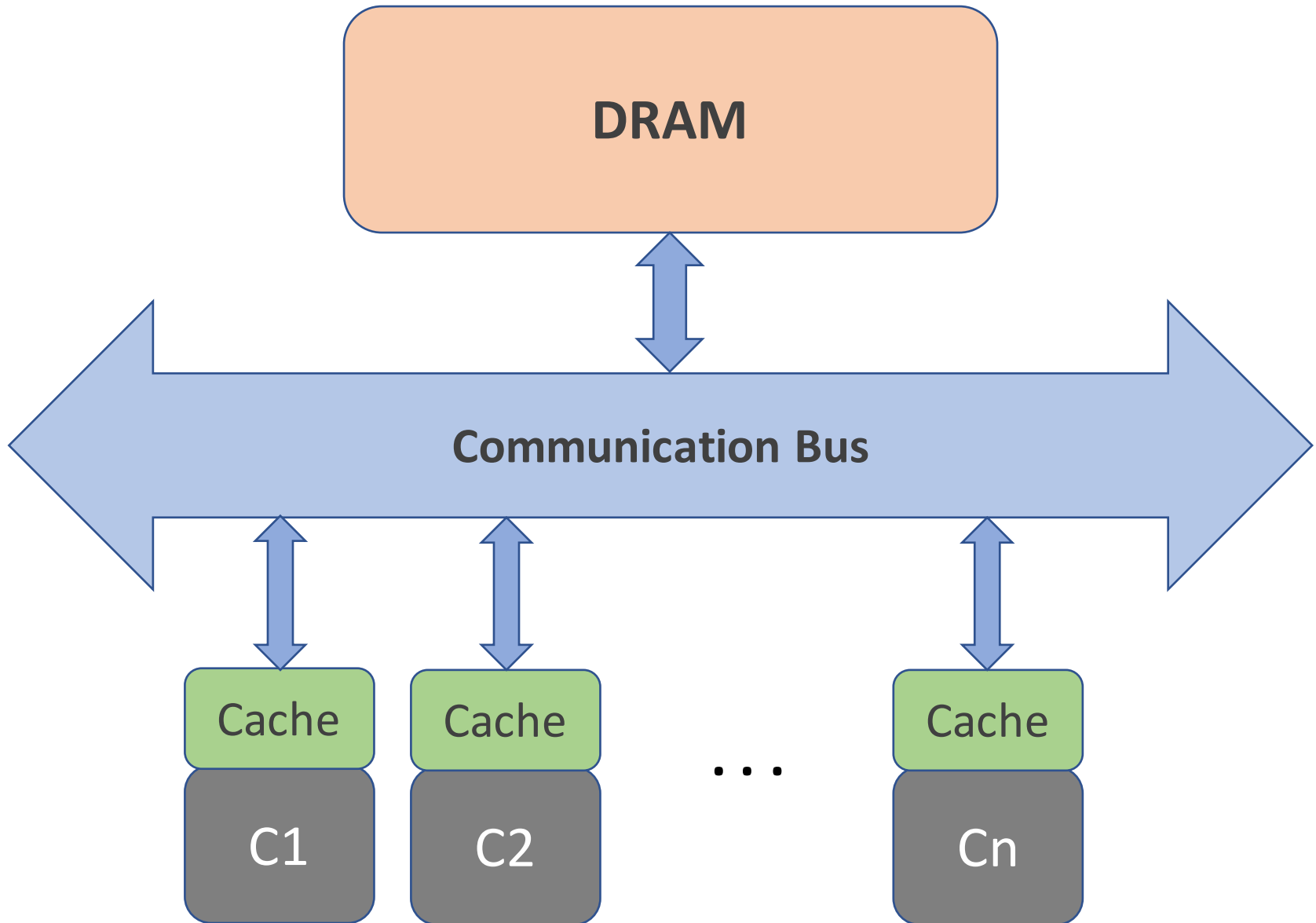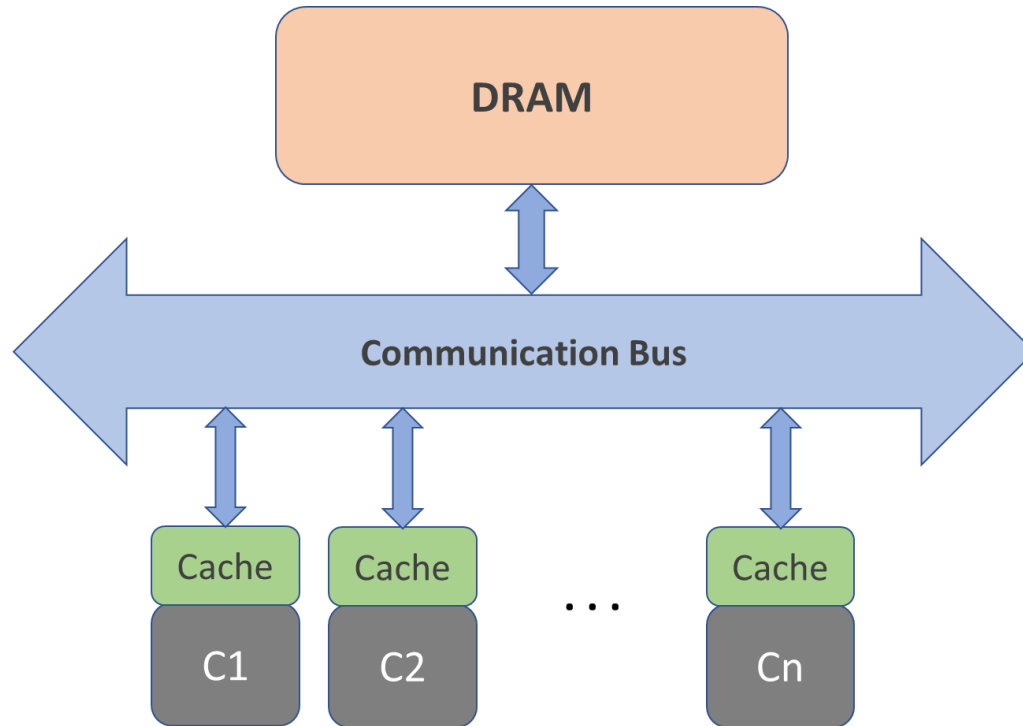
- Each Core has its own L1 Cache
- In some processors, each core has its own L2 Cache
- All cores share one L3 cache
- All cores share one DDR3 DRAM (main memory)

Image source https://www.realworldtech.com/sandy-bridge/2/
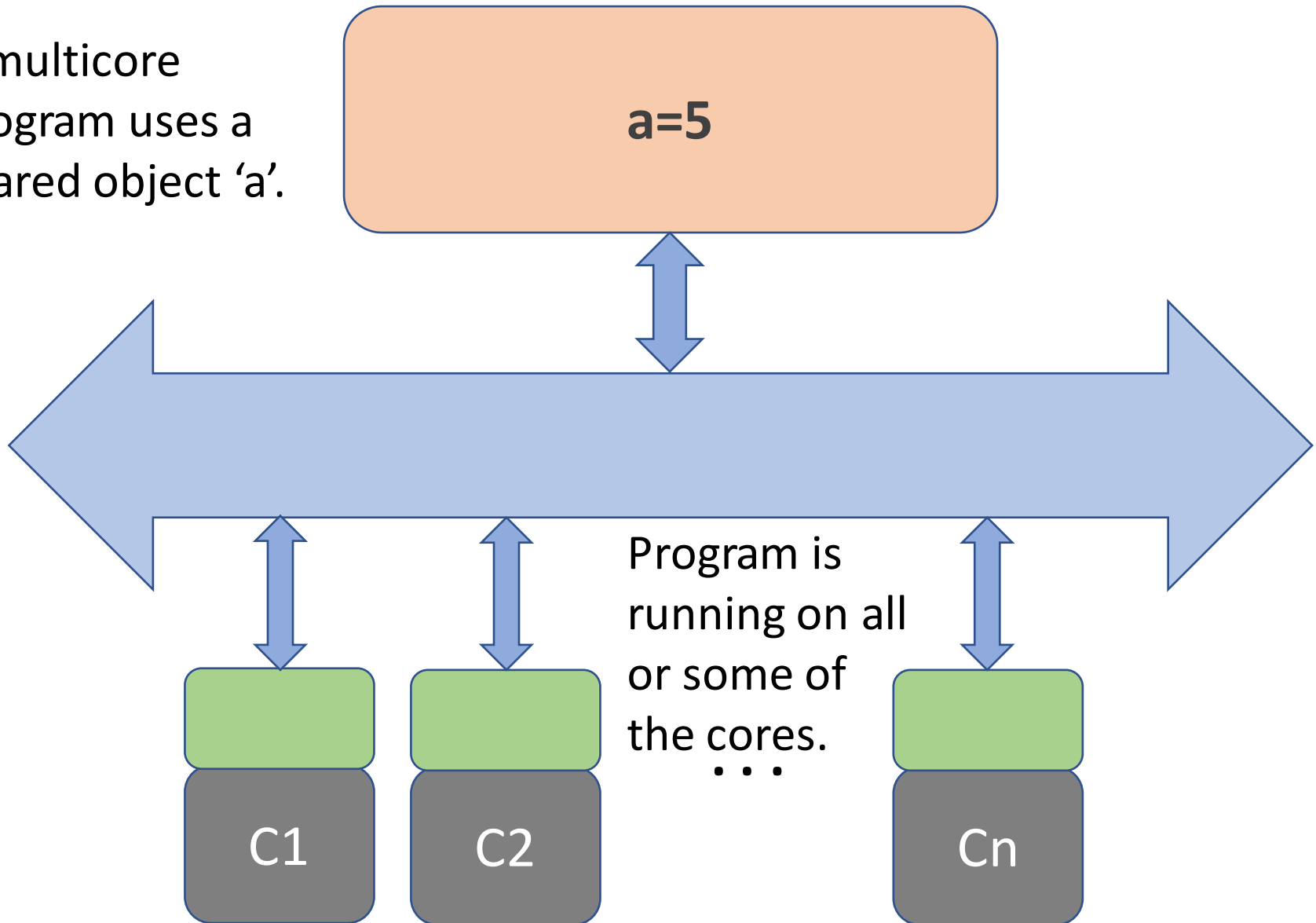
16

# Simplified Multicore Architecture
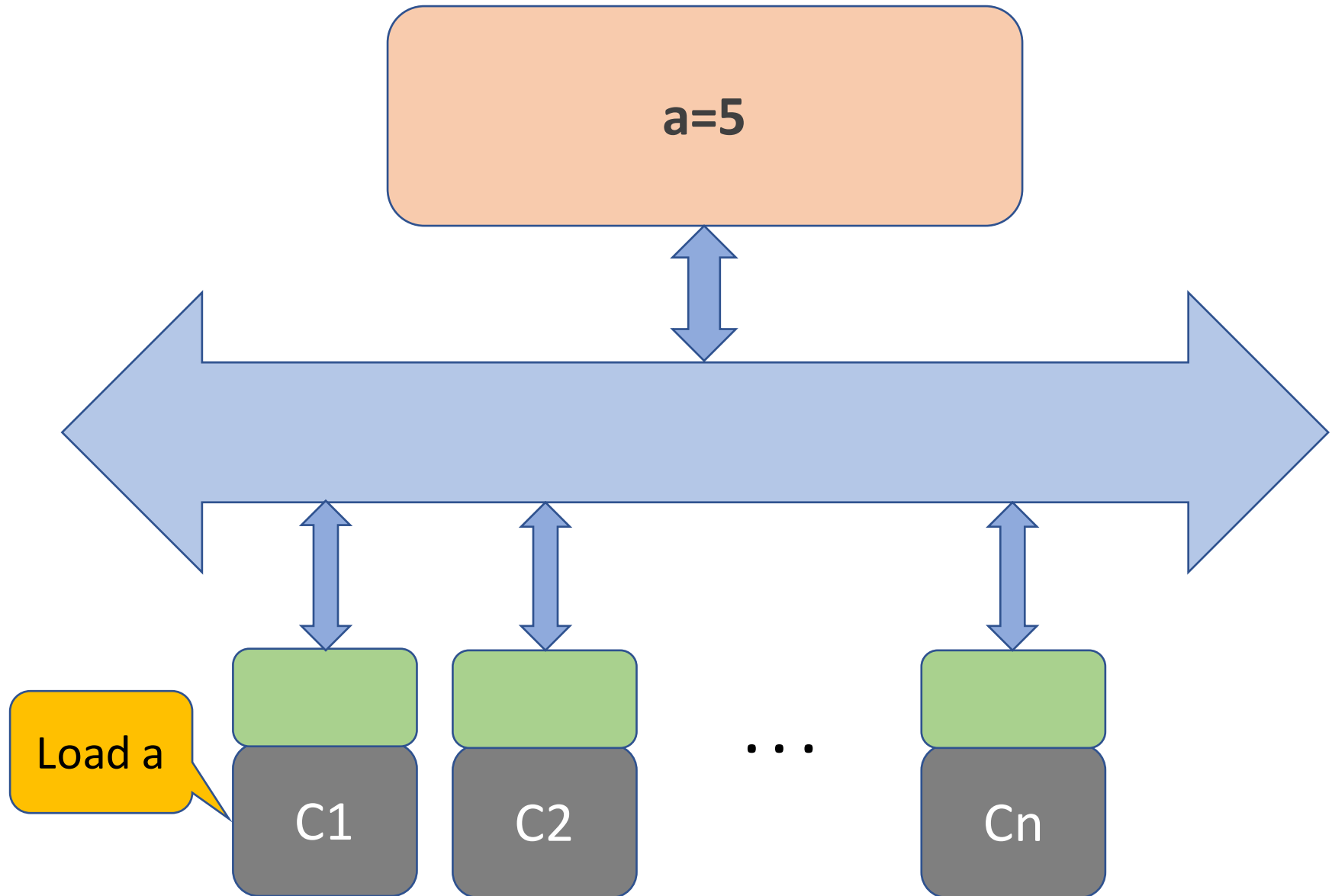
# Coherence of data



- 'Coherence' is the quality of being logical and consistent. During program execution, data must remain coherent.

- On a multicore system, a data variable may reside in multiple caches and might get updated 'locally' by its local core
  → this causes coherence problem
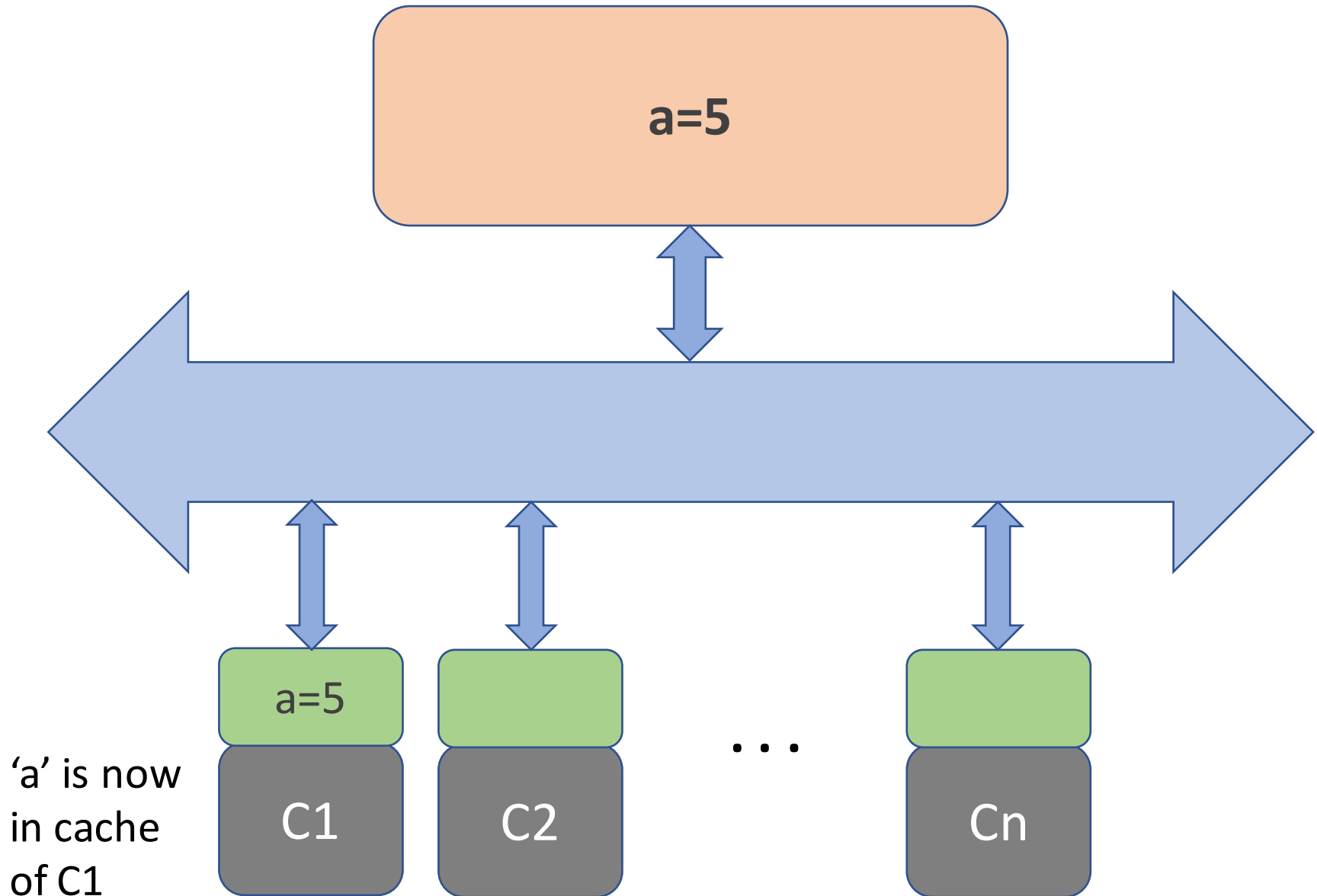
# Example: Cache coherence problem
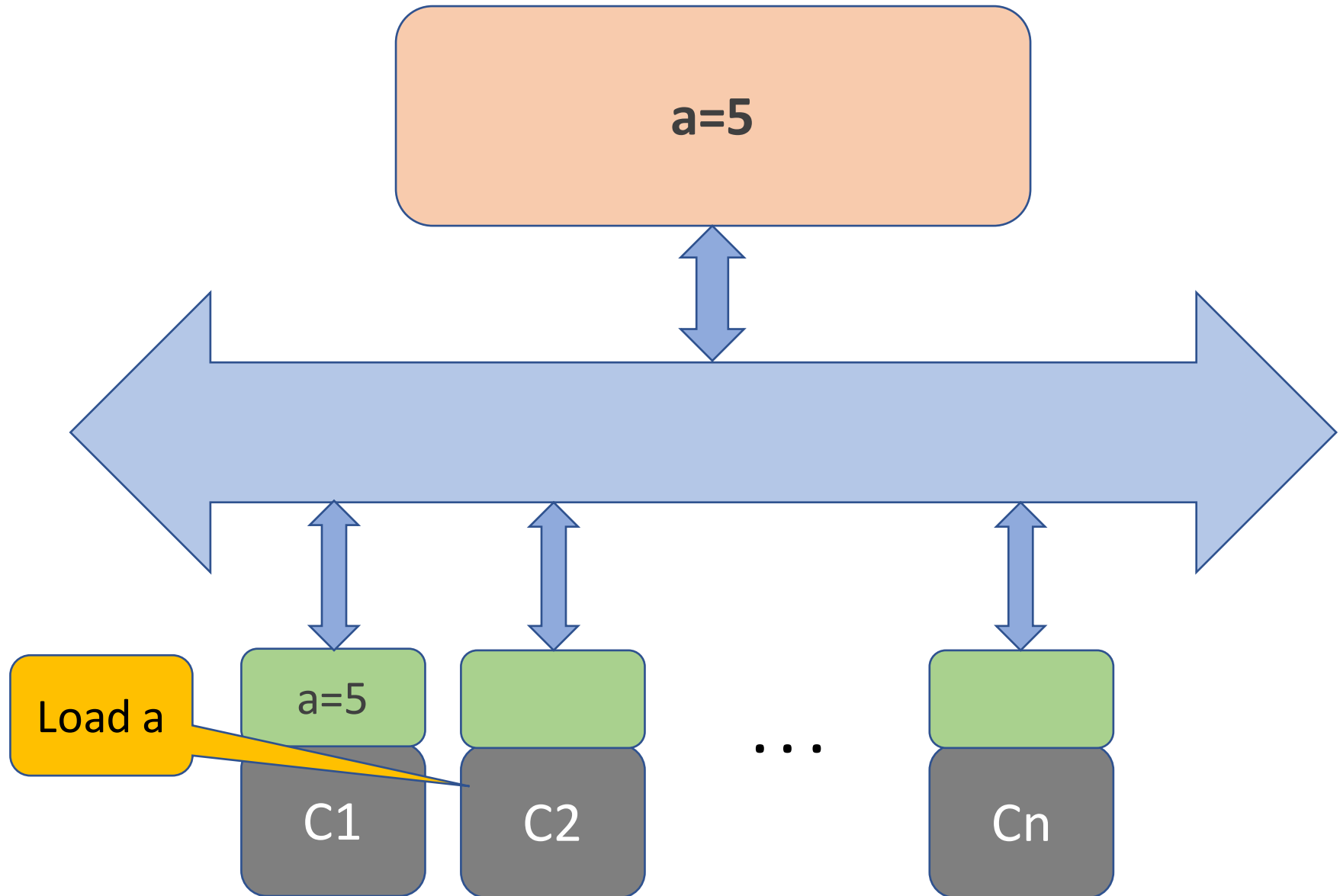
A multicore program uses a shared object 'a'.

a=5

Program is running on all or some of the cores.

C1    C2    • • •    Cn

# Example: Cache coherence problem

# Example: Cache coherence problem

a=5

a=5

'a' is now in cache of C1

C1          C2          . . .          Cn

# Example: Cache coherence problem

# Example: Cache coherence problem

a=5

Multiple copies of 'a' are now in caches of C1 and C2.

a=5

a=5

C1

C2

. . .

Cn

# Example: Cache coherence problem

a=5

a=5

a=5

Load a

C1

C2

. . .

Cn

# Example: Cache coherence problem

a=5

a=5    a=5    . . .    a=5

C1     C2              Cn

Multiple copies of 'a' are now in caches of C1, C2, Cn

# Example: Cache coherence problem

# Example: Cache coherence problem



a=5

a=5    a=5    . . .    **a=10**

C1     C2              **Cn**
                       **a = 10**

Only local copy of 'a' has the latest value. All other copies are outdated.

# Example: Cache coherence problem



a=5

Load a

a=5      a=5      ...      a=10

C1       C2                Cn
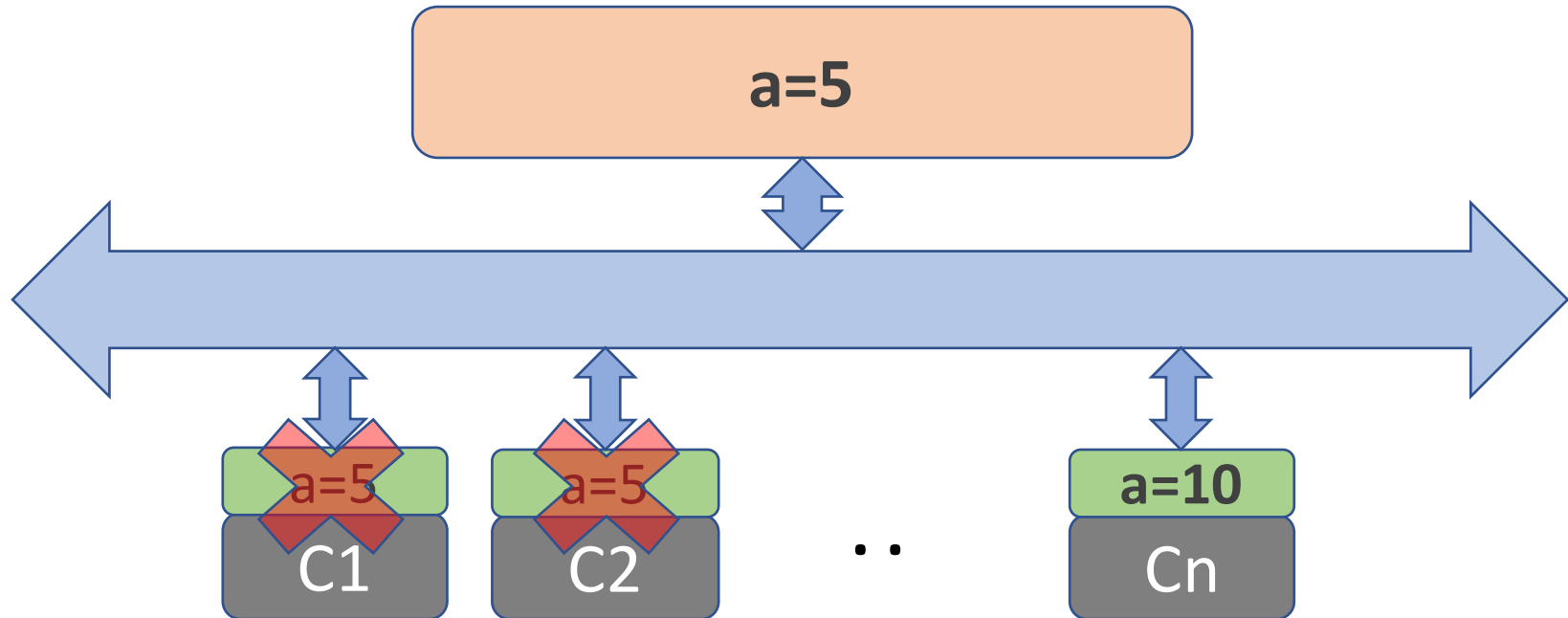                           a = 10

C1 gets an outdated 'a' from its Cache.

# Coherence mechanisms

To achieve coherence:
a **write to a memory location** must cause all other copies of this location to be **removed from the caches** they are in.

a=5

a=5    C1

a=5    C2

. .

a=10    Cn

Do: As soon as 'a' gets updated by a core, invalidate or delete all other copies of it.

# Coherence protocols

- Coherence protocols apply cache coherence in multicore systems.

- Goal is that two cores **must never see different values** for the same shared data.

How this is achieved?
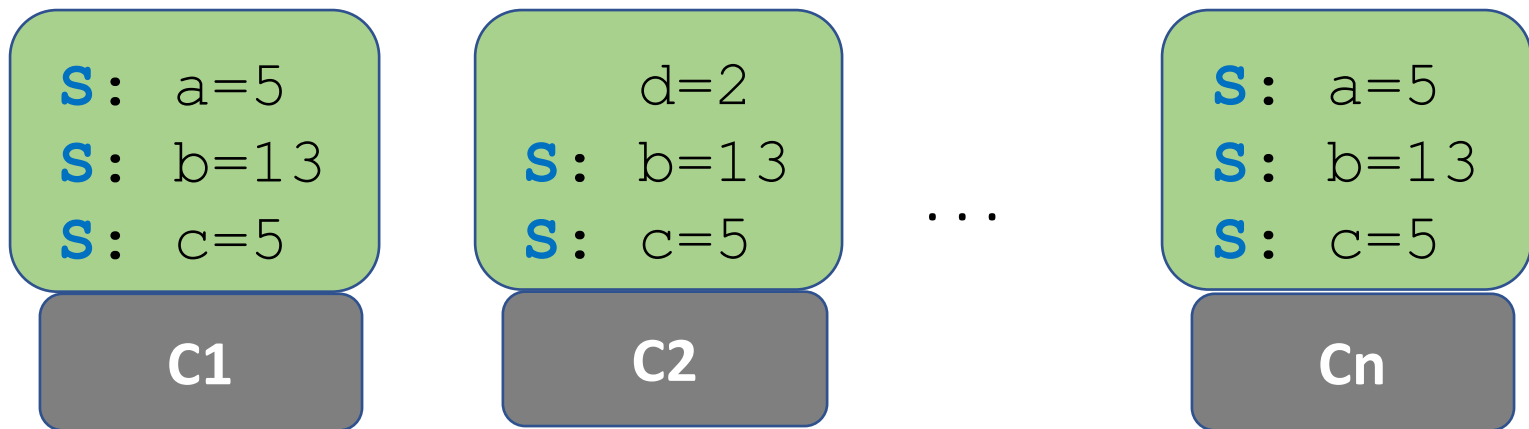
Additional information is stored and passed

- Is this data stored in multiple caches?

- Has this cache line been modified?

# MSI protocols

MSI protocol is a simple cache coherence protocol.
In this protocol, each cache line is labeled with a state:

- **M**: cache block has been modified.

- **S**: other caches may be sharing this block.

- **I**: cache block is invalid

## Initial state of caches before computation

```
S:  a=5          d=2          S:  a=5
S:  b=13     S:  b=13         S:  b=13
S:  c=5      S:  c=5    ...   S:  c=5
```

**C1**          **C2**                **Cn**

# MSI protocols

MSI protocol is a simple cache coherence protocol.
In this protocol, each cache line is labeled with a state:

- **M**: cache block has been modified.

- **S**: other caches may be sharing this block.

- **I**: cache block is invalid

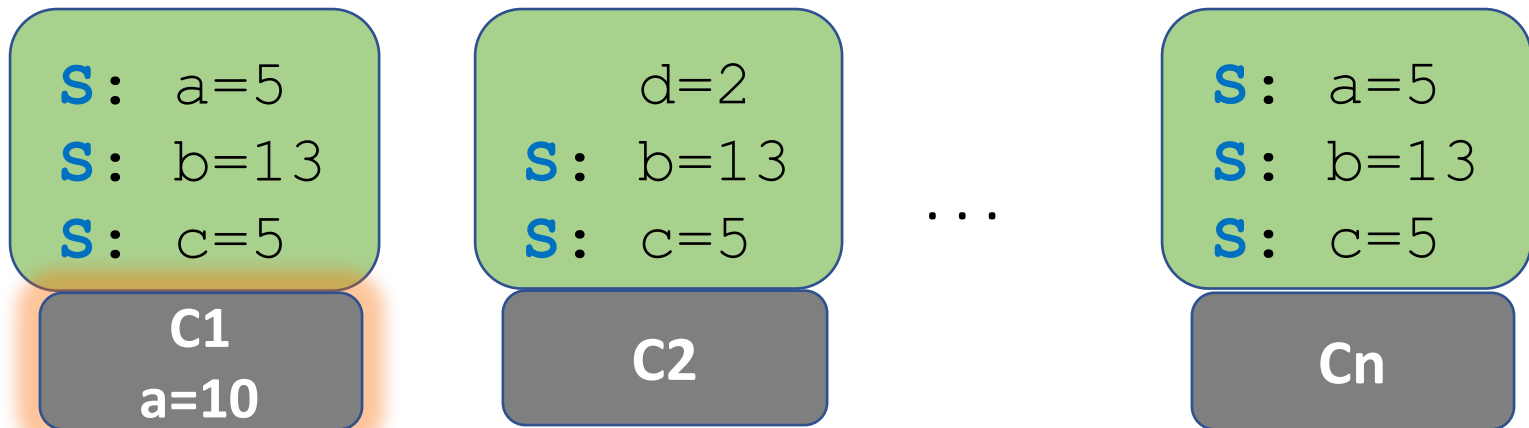**C1 computes new value a=10 in register**

# MSI protocols

MSI protocol is a simple cache coherence protocol.
In this protocol, each cache line is labeled with a state:

- **M**: cache block has been modified.

- **S**: other caches may be sharing this block.

- **I**: cache block is invalid

**Hardware invalidates all other copies of a**

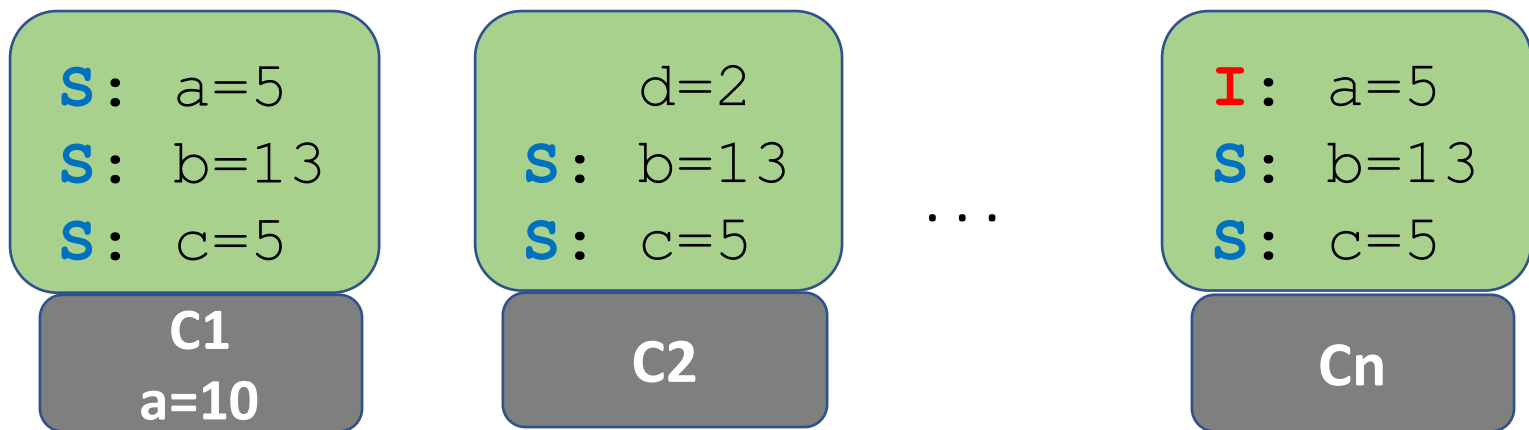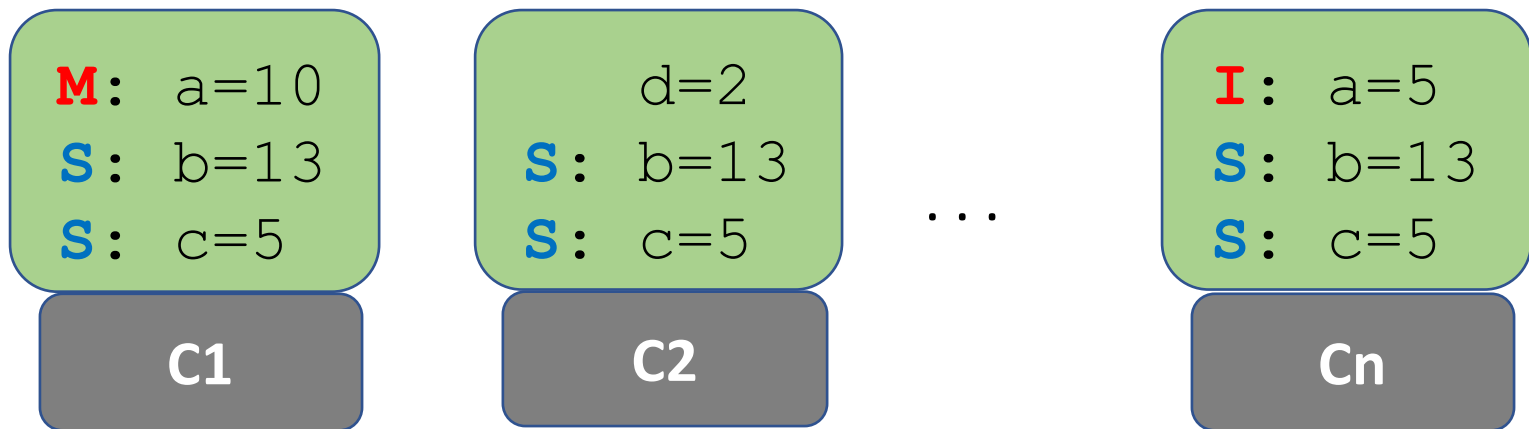| S: a=5 | d=2 | | I: a=5 |
|---|---|---|---|
| S: b=13 | S: b=13 | ... | S: b=13 |
| S: c=5 | S: c=5 | | S: c=5 |
| **C1** | **C2** | | **Cn** |
| **a=10** | | | |

# MSI protocols

MSI protocol is a simple cache coherence protocol.
In this protocol, each cache line is labeled with a state:

- **M**: cache block has been modified.

- **S**: other caches may be sharing this block.

- **I**: cache block is invalid

**Hardware updates a in cache of C1 and marks it with 'M'**

| M: a=10 |
| S: b=13 |
| S: c=5 |
| **C1** |

| d=2 |
| S: b=13 |
| S: c=5 |
| **C2** |

. . .

| I: a=5 |
| S: b=13 |
| S: c=5 |
| **Cn** |

# MSI protocols

MSI protocol is a simple cache coherence protocol.
In this protocol, each cache line is labeled with a state:

- `M`: cache block has been modified.

- `S`: other caches may be sharing this block.

- `I`: cache block is invalid

**Hardware also updates 'a' in the shared memory**

```
M:  a=10          d=2          I:  a=5
S:  b=13      S:  b=13         S:  b=13
S:  c=5       S:  c=5    ...   S:  c=5
```
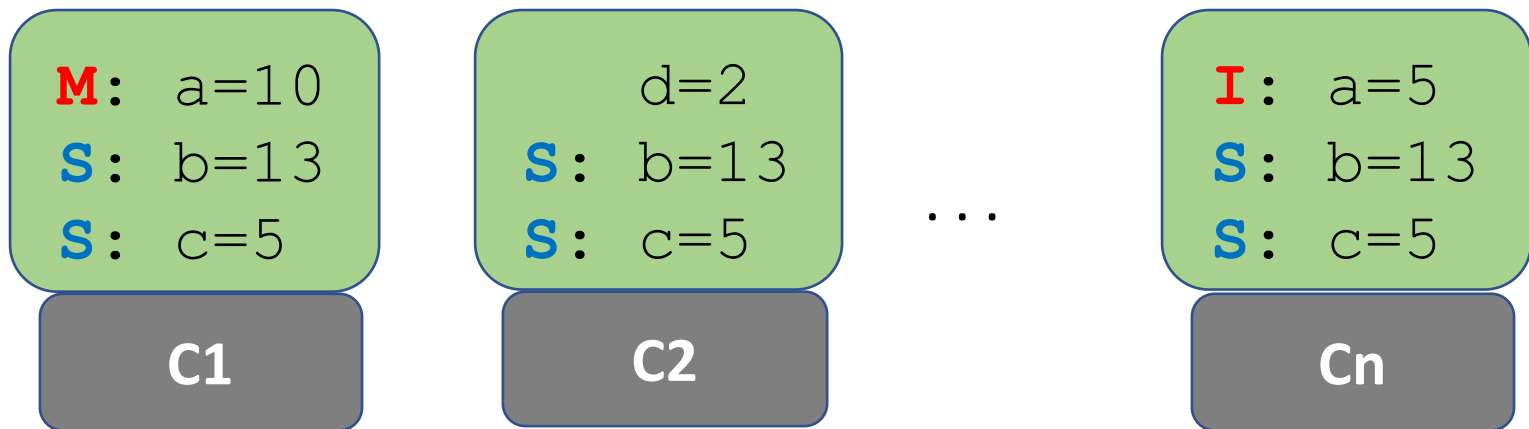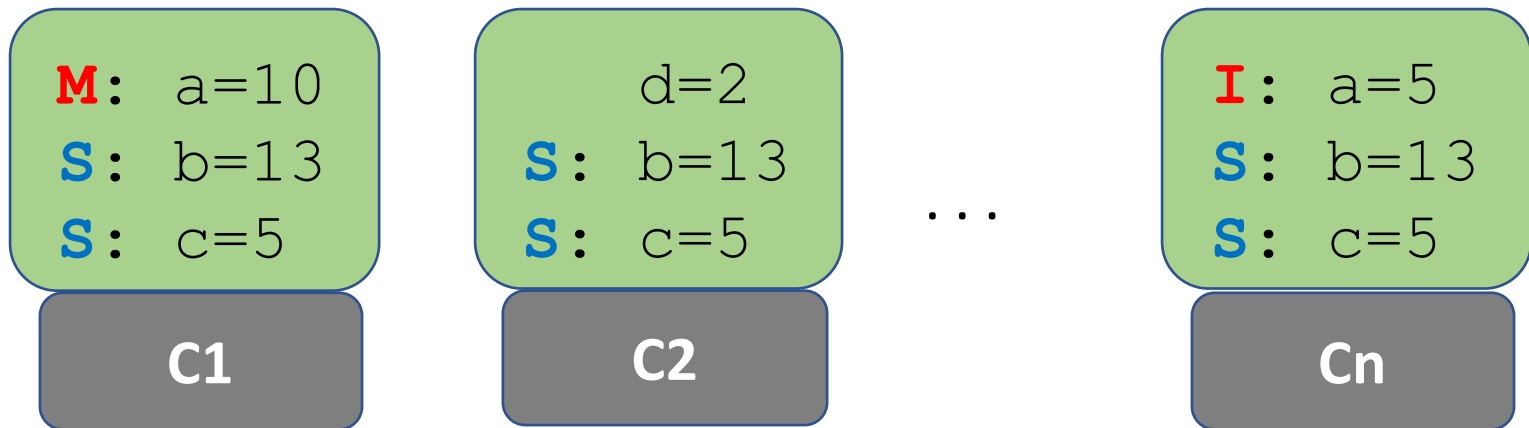
| C1 | C2 | Cn |

# MSI protocols

MSI protocol is a simple cache coherence protocol.
In this protocol, each cache line is labeled with a state:

- **M**: cache block has been modified.

- **S**: other caches may be sharing this block.

- **I**: cache block is invalid

**If Cn wants to load 'a' then 'a' will be brought from the shared memory to the local cache of Cn.**

```
M:  a=10          d=2          I:  a=5
S:  b=13      S:  b=13         S:  b=13
S:  c=5       S:  c=5   ...    S:  c=5
```

**C1**  **C2**  **Cn**

# Programming multicore platforms

- Option 1: Program directly targeting processor cores
  - ➢ Programmer takes care of synchronization
  - ➢ Painful and error-prone


- Option 2: Use a *concurrency platform*
  - ➢ It abstracts processor cores, handles synchronization and communication protocols, and performs load balancing.

  - ➢ Hence offers much easier multicore programming environment

  - ➢ Examples:
    - **Pthreads** and WinAPI threads
    - OpenMP