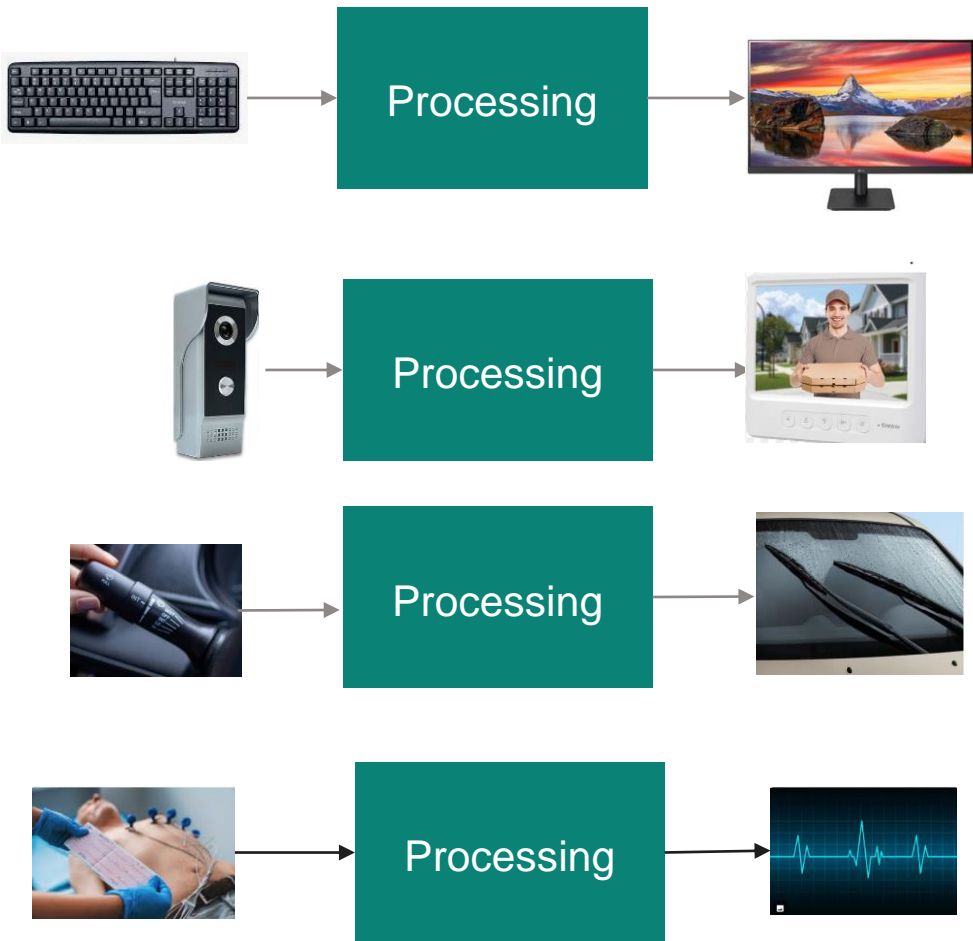# Computer architecture

Prakash Balasubramanian

06/April/2024
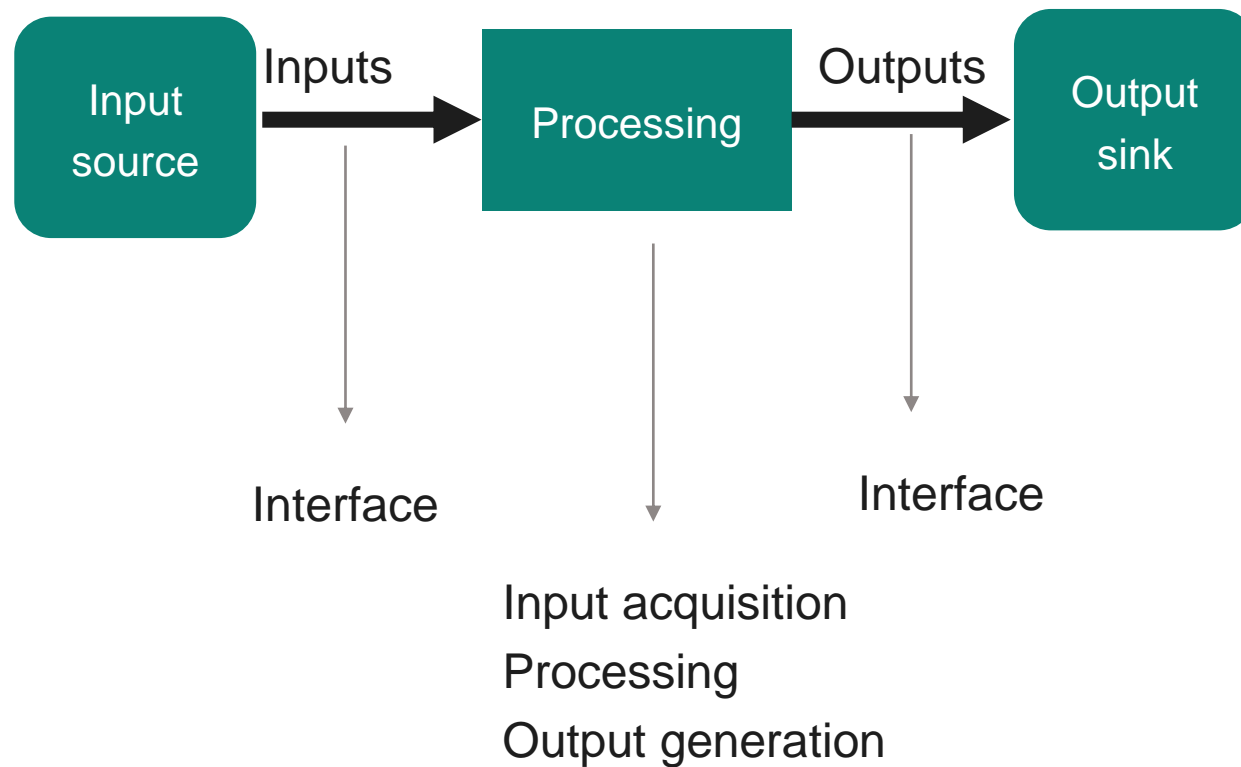
# Systems – 1/3
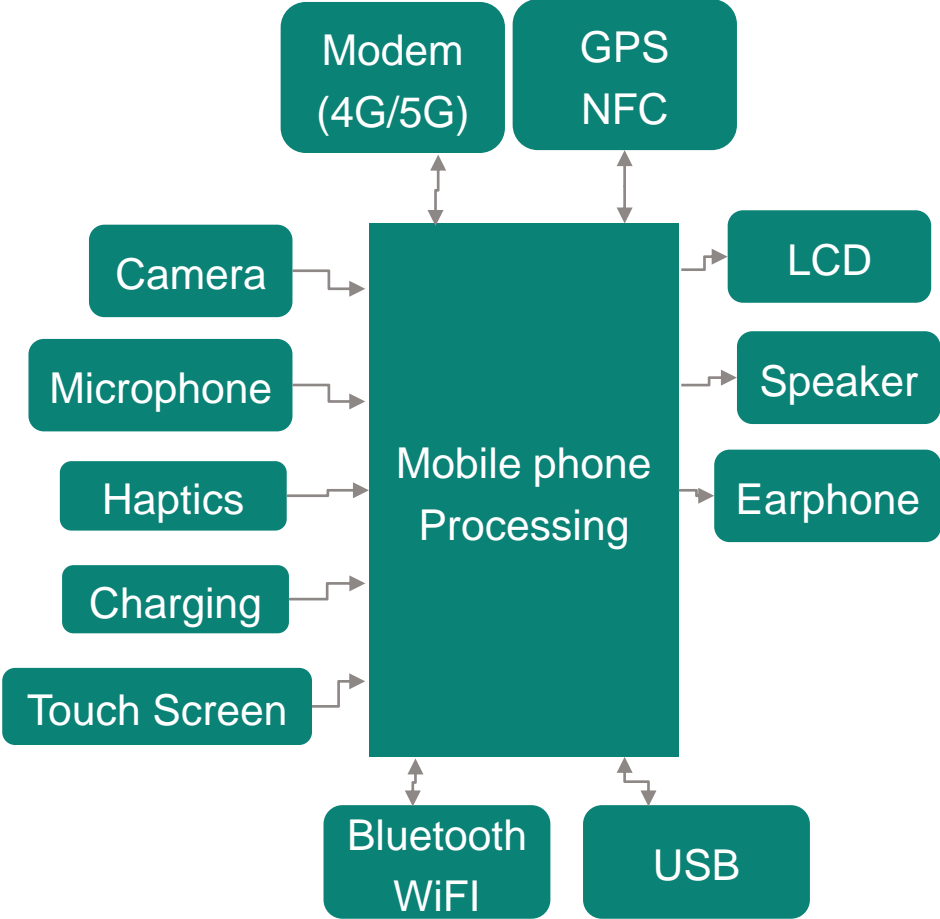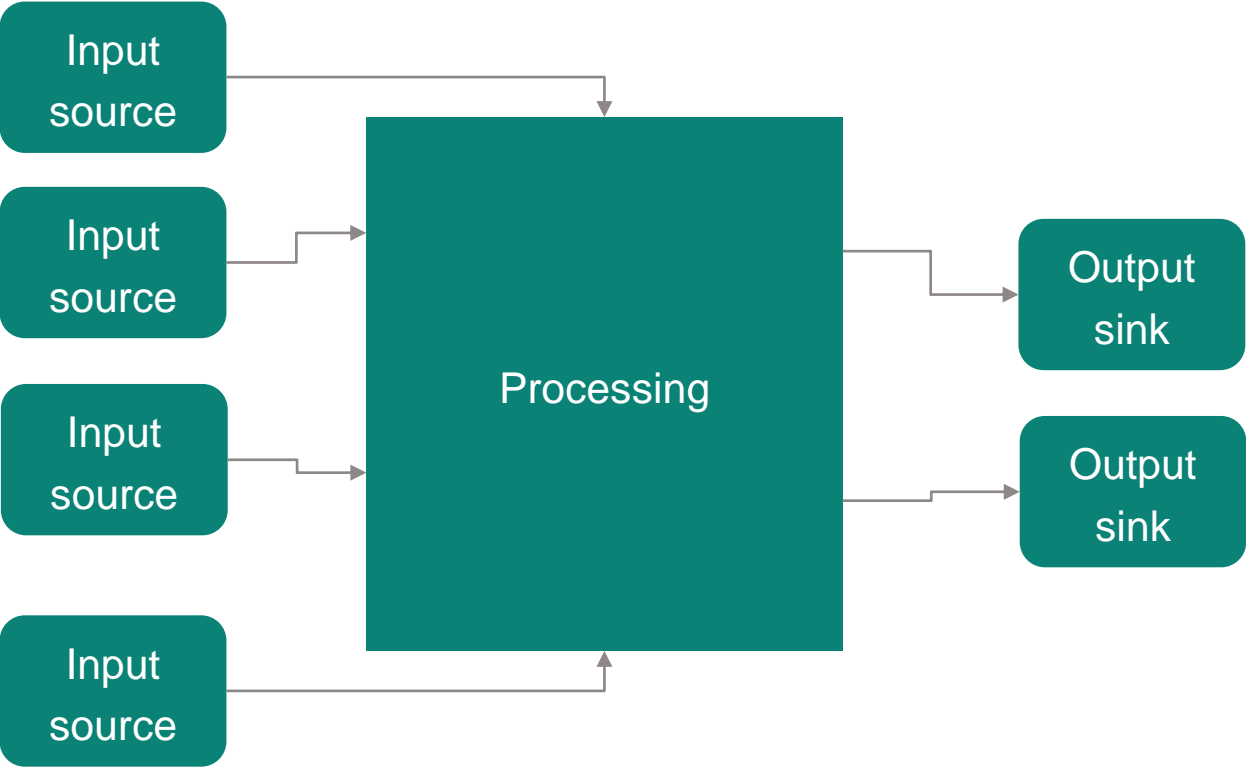
**Examples**
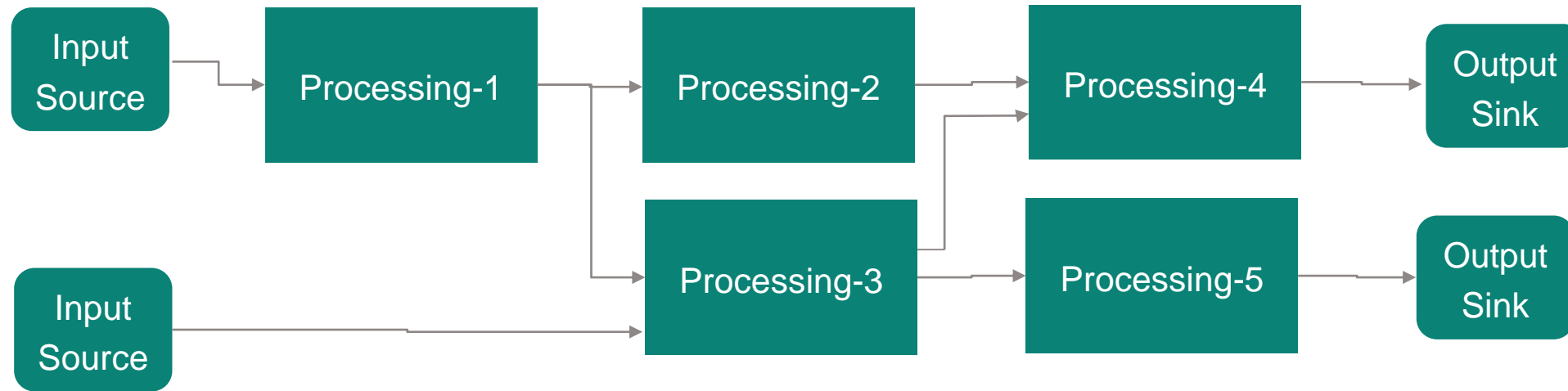
**Generalization**

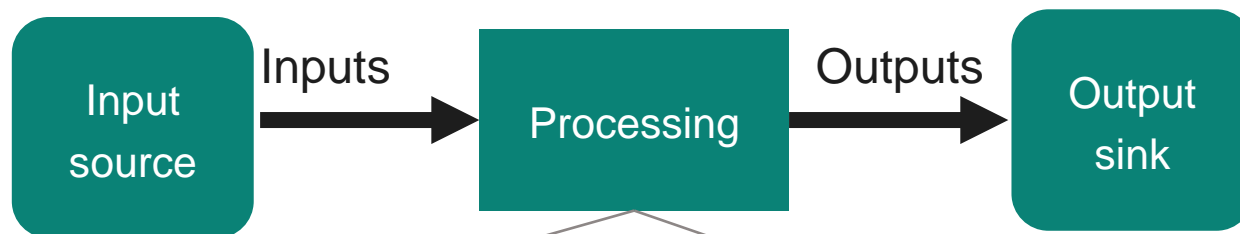## Multiple Input Multiple Output

# Systems – 3/3

## Hierarchy of systems, cascaded systems

# Real time system – 1/2



**General purpose processing**

- Working on an excel sheet on a Windows/Linux PC

- Watching a YouTube video

*Objectives of the system are met even if there are lapses and delays in processing.*
*These are Non-Real time systems*

**Real Time processing**

- Acquiring and processing an image frame before the next frame arrives
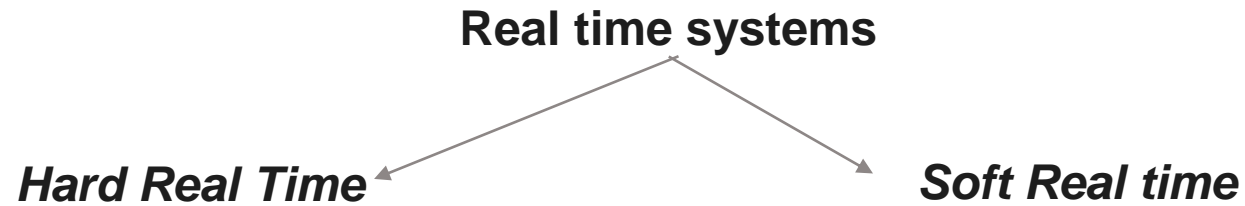- Reading out a sensor value when a new input signal becomes available

*Time is the essence*
*Processing must be performed at intended time*
*Processing must be completed within specified duration*
*MOST EMBEDDED SYSTEMS ARE REAL TIME SYSTEMS!*

# Real time system 2/2

**Real time systems**

*Hard Real Time*

*Soft Real time*

Deadlines cannot be missed.
Deadlines are Non-Negotiable!

An occasional miss in deadline is acceptable

- The pressing of the brake pedal must be recognized
- Once recognized, braking action must be applied within 20ms.
- If not, there is a good probability that an accident may happen and perhaps even result in fatality!

- A camera captures images at a rate of 30 frames per second
- Time gap between two frames is 1/30 = 33 ms
- Processing of a frame (e.g. Demosaicing) should therefore be done within 33 ms of its attival
- If every 472nd frame takes a little longer than 33ms for processing, there is barely any impact on end user
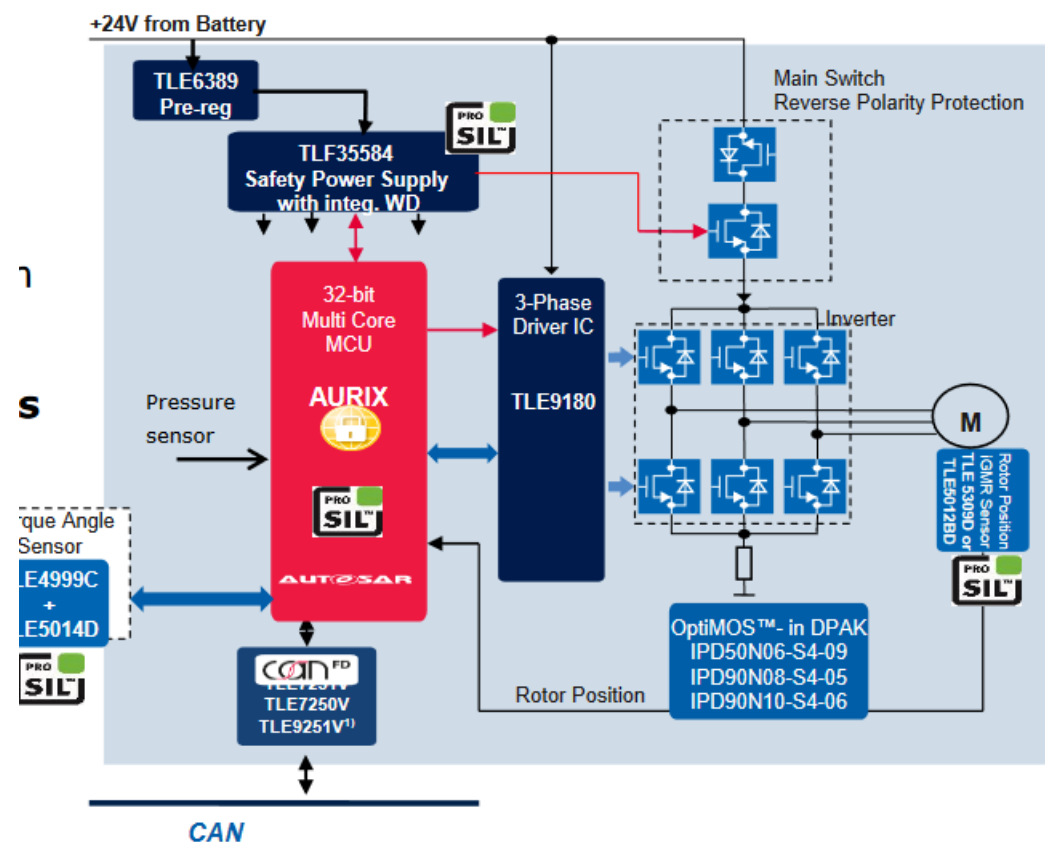
# General purpose vs Embedded systems
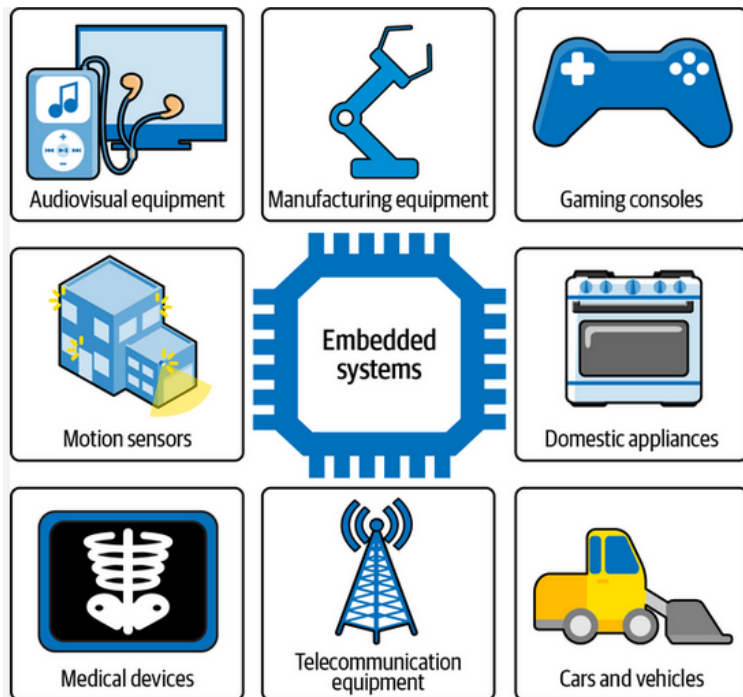
## General purpose applications



- Games, Video
- Browsing
- E-Mail
- Desktop publishing
- Star gazing app
- Stock market app

## Specific purpose system based on a Microprocessor/Microcontroller/SoC

# Pervasiveness of embedded systems



QUIZ: How would you classify this?

# Embedded systems – 1/3

**1970s**

```
Input-1          Input
source           peripheral          CPU
                 controller-1

Input-2          Input               Memories        Output         Output
source           peripheral                          peripheral     sink
                 controller-2                         controller-3
```

```
Keyboard         8255                6800            4511           7-Segment
                                                                    Display
Temp             741                 6116
Sensor
```

Reference Voltage

# Embedded systems – 2/3



**1980s - Today**

**Microcontroller**

- Input-1 source
- Input-2 source
- Input peripheral controller-1
- Input peripheral controller-2
- CPU
- Memories
- Output peripheral controller-3
- Output sink

*Early 2000s until today*

**Massive integration**

**Single Board Computers**

# Typical software stack



**Bare Metal Application**

Application

Middleware

Device Drivers

SoC/Microcontroller with CPU, Peripherals and Memory

Inputs

Outputs

**OS based Application**

Application

Middleware | OS

Device Drivers

SoC/Microcontroller with CPU, Peripherals and Memory

Inputs

Outputs

# Example of Peripherals and their device drivers

| OpAmp Driver | UART Driver | Timer Driver | GPIO Driver |
|:---:|:---:|:---:|:---:|

CPU

OpAmp    UART    Timer    GPIO

# A typical Microcontroller



CPU-1  CPU-2  Memory-1  Memory-2  Other Bus Masters

BUS or INTERCONNECT

Input Peripheral-1  Output Peripheral-2  Output Peripheral-3  Input-Output Peripheral-N

**Legend**

Bus Masters

Bus Slaves

# Von-Neumann CPU – 1/4

*All CPUs used during the last 70+ years are Von-Neumann architectures!*
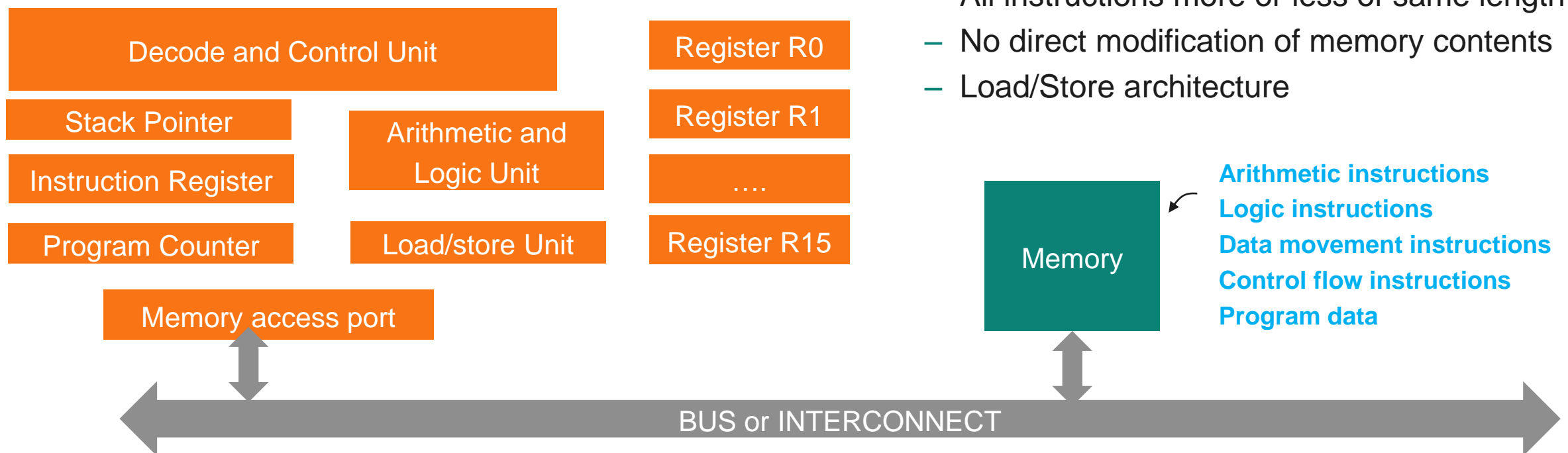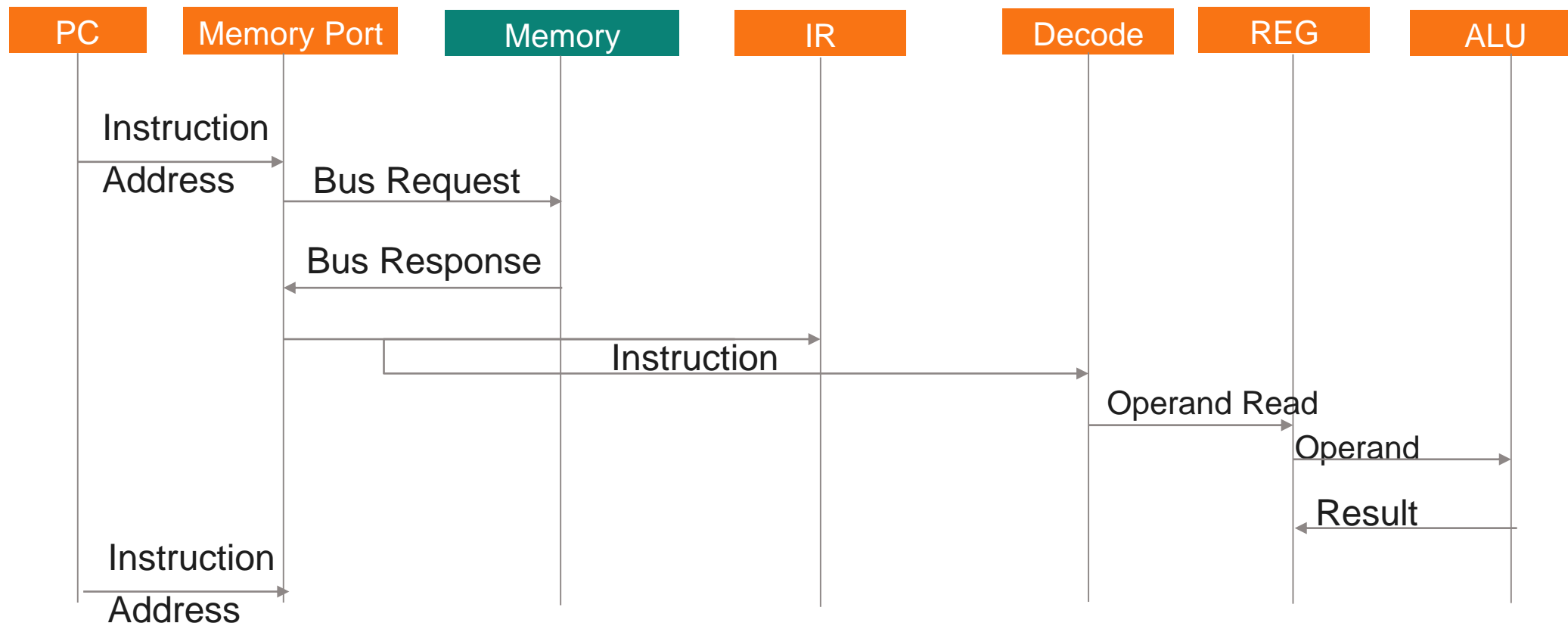
*All CPUs of today are improvisations of the basic RISC architecture*

– All instructions more or less of same length
– No direct modification of memory contents
– Load/Store architecture

| Decode and Control Unit | | Register R0 |
|---|---|---|
| Stack Pointer | Arithmetic and Logic Unit | Register R1 |
| Instruction Register | | .... |
| Program Counter | Load/store Unit | Register R15 |
| Memory access port | | |

**Memory**

**Arithmetic instructions**
**Logic instructions**
**Data movement instructions**
**Control flow instructions**
**Program data**

BUS or INTERCONNECT

# Von-Neumann CPU – 2/4

## Handling of ALU instructions
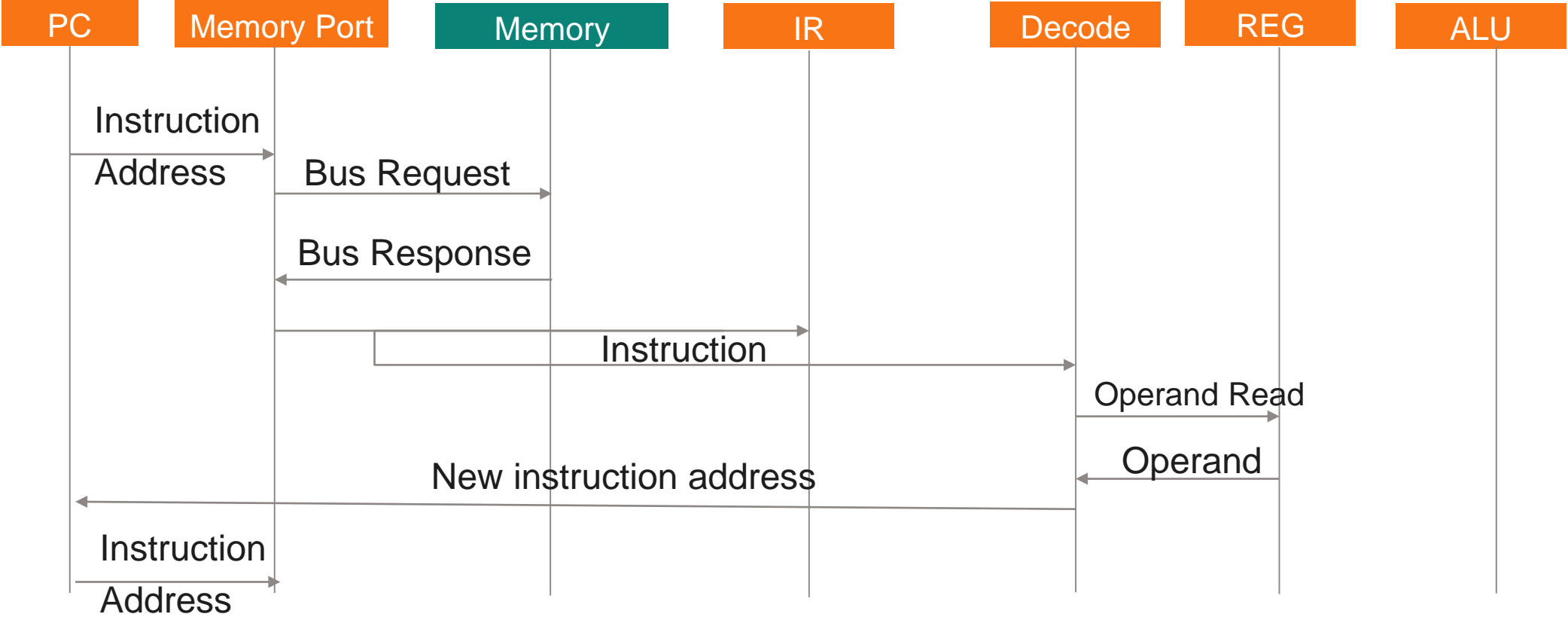
# Von-Neumann CPU – 3/4

**Handling of Load/Store instructions**

# Von-Neumann CPU – 4/4

**Handling of Control instructions**



| PC | Memory Port | Memory | IR | Decode | REG | ALU |
|----|-------------|--------|----|--------|-----|-----|

Instruction
Address

Bus Request

Bus Response

Instruction

Operand Read

Operand

New instruction address

Instruction
Address

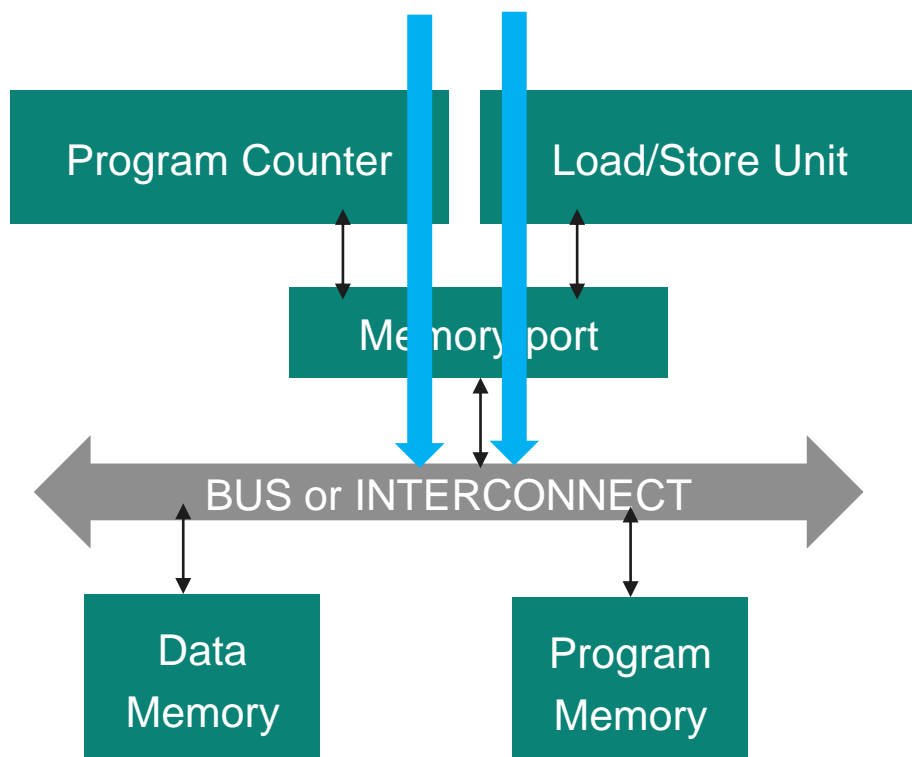Instruction Fetch          Instruction Decode          Instruction Execute

# Von-Neumann vs Harvard architecture – 1/2

Harvard architecture IS essentially an improvisation of Von-Neumann architecture.

Harvard architecture removes a Von-Neumann bottleneck. That's it!

| Program Counter | Load/Store Unit |
|---|---|

Memory port

BUS or INTERCONNECT

Data Memory

Program Memory

The same memory port is used for both fetching instructions during the Instruction Fetch Cycle and Data fetch/data store during the Execute cycle of a Load/Store instruction.
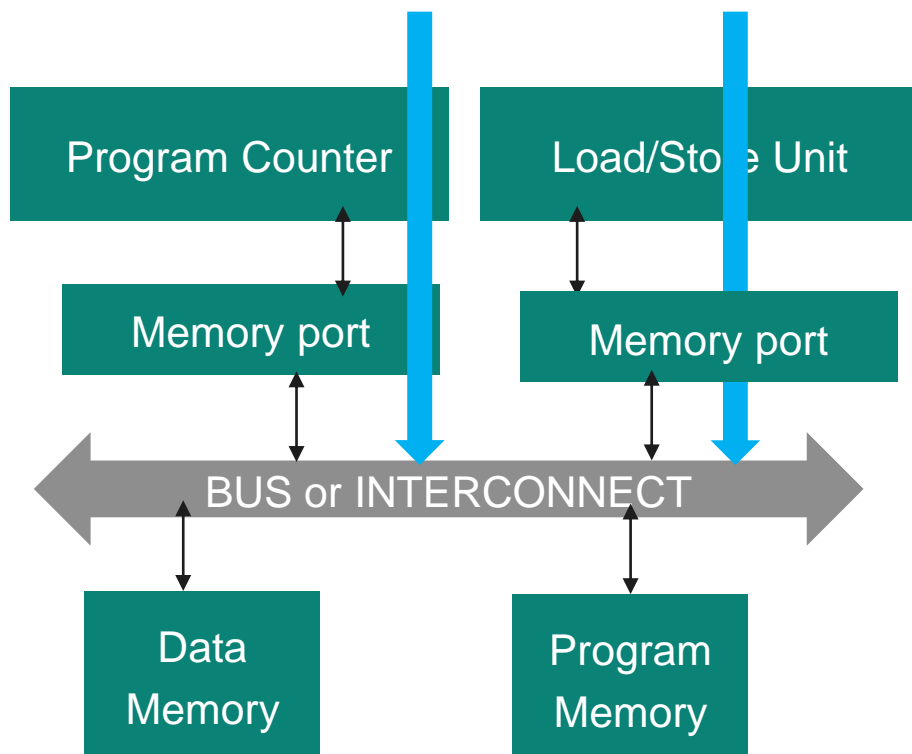
Only one can happen at a time and the other will have to wait!

This is the Von-Neumann bottleneck (among many others!)

# Von-Neumann vs Harvard architecture – 2/2

Harvard architecture IS essentially an improvisation of Von-Neumann architecture.

Harvard architecture removes a Von-Neumann bottleneck. That's it!



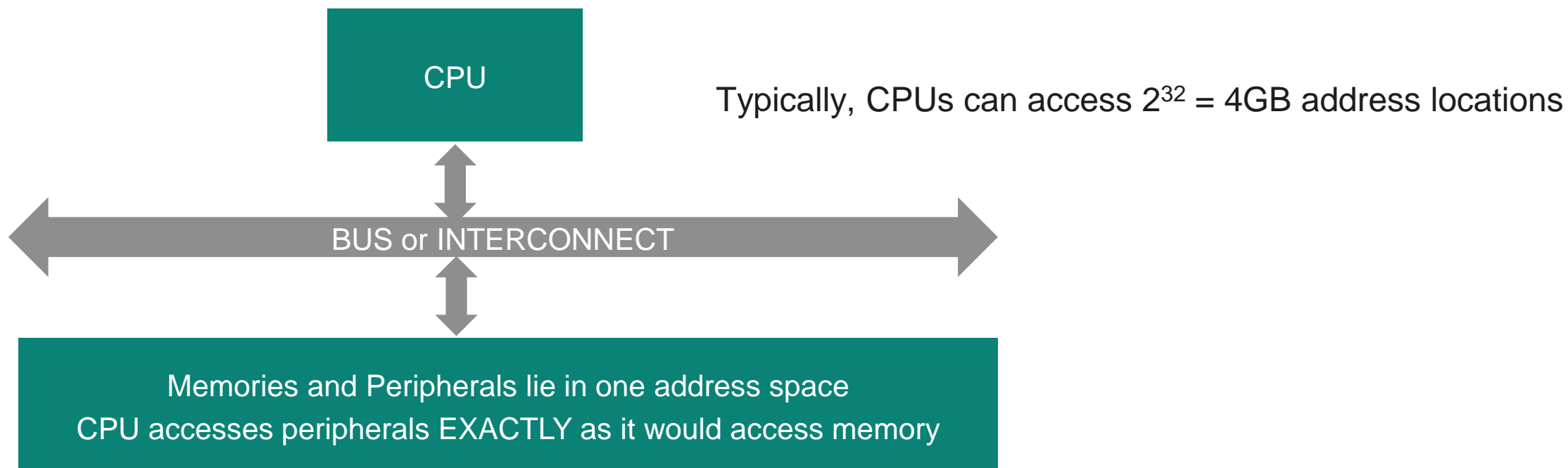By providing two memory ports, both instruction fetch and data load/store can happen concurrently!

With this, we have eliminated the bottleneck!

**Question:**

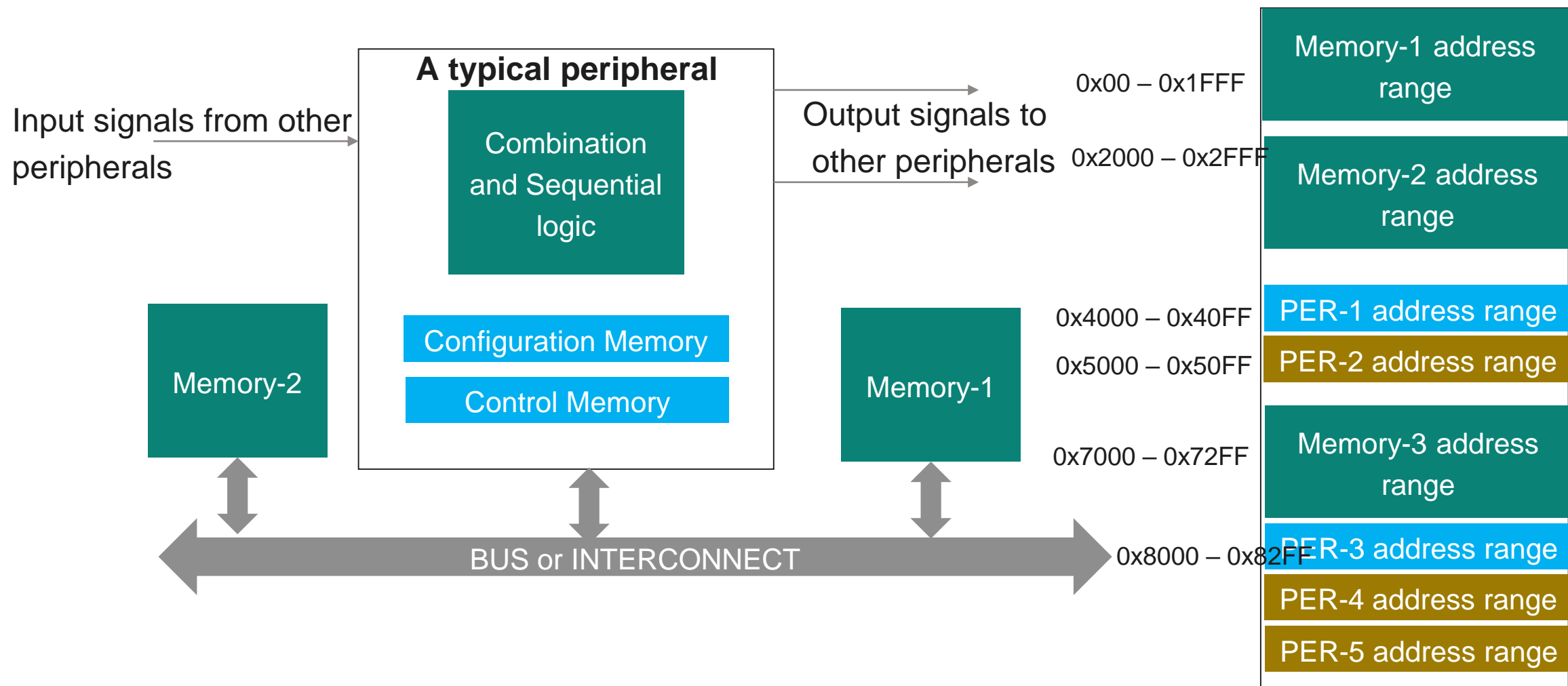Is the bottleneck truly eliminated? What is the problem with the picture?

# Address space of user program and peripherals – 1/2

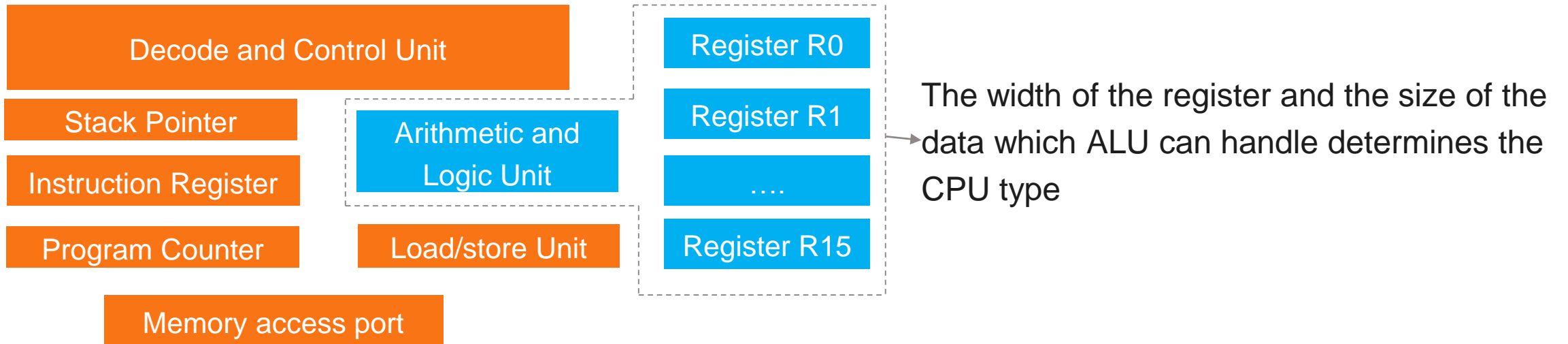REMEMBER THIS => **All modern CPUs follow Memory Mapped IO addressing scheme**

CPU

Typically, CPUs can access $2^{32}$ = 4GB address locations

BUS or INTERCONNECT

Memories and Peripherals lie in one address space
CPU accesses peripherals EXACTLY as it would access memory

# Address space of user program and peripherals – 2/2



*Unified address space of memory and peripherals*

**A typical peripheral**

Input signals from other peripherals

Combination and Sequential logic

Output signals to other peripherals

Configuration Memory

Control Memory

Memory-2

Memory-1

BUS or INTERCONNECT

0x00 – 0x1FFF

0x2000 – 0x2FFF

0x4000 – 0x40FF

0x5000 – 0x50FF

0x7000 – 0x72FF

0x8000 – 0x82FF

Memory-1 address range

Memory-2 address range

PER-1 address range

PER-2 address range

Memory-3 address range

PER-3 address range

PER-4 address range

PER-5 address range

# What type is my CPU?

| Decode and Control Unit |
| Stack Pointer |
| Instruction Register |
| Program Counter |

| Arithmetic and Logic Unit |
| Load/store Unit |

Memory access port

| Register R0 |
| Register R1 |
| …. |
| Register R15 |

The width of the register and the size of the data which ALU can handle determines the CPU type

If the register width is 32bit, then ALU is typically designed to operate on 32b operand at a time

If the register width is 16bit, then ALU is typically designed to operate on 16b operand at a time

# Memories – 1/2

### A toy application

```
const int LOOP_COUNT = 5;

int b[5];

int a[5] = { 20, -19, 773, 47, -146};

void main(void)

{

  int *ptr;

  ptr = b;


 while(port->portA==0)

 {

  for(int i=0;i<LOOP_COUNT;i++)

  {

    ptr[i] = a[i];

  }

 }

}
```
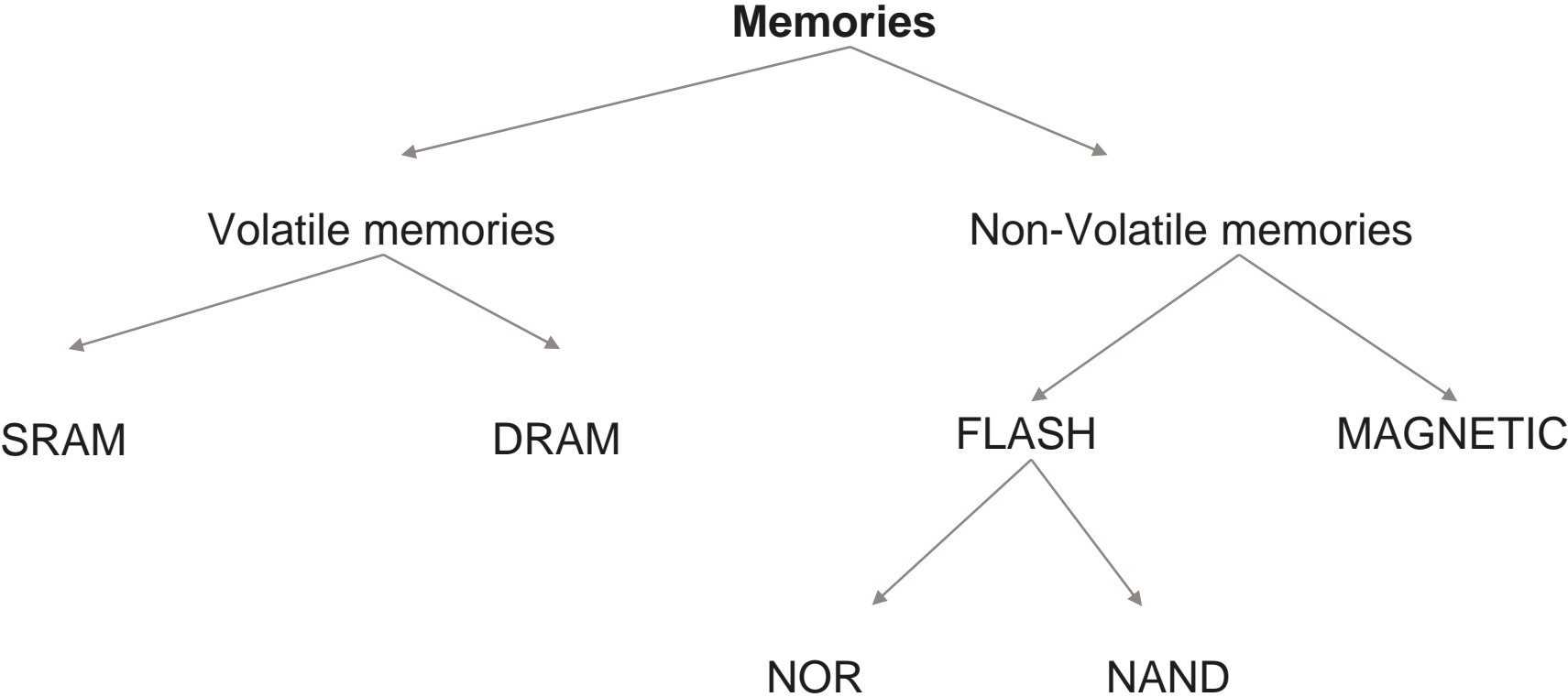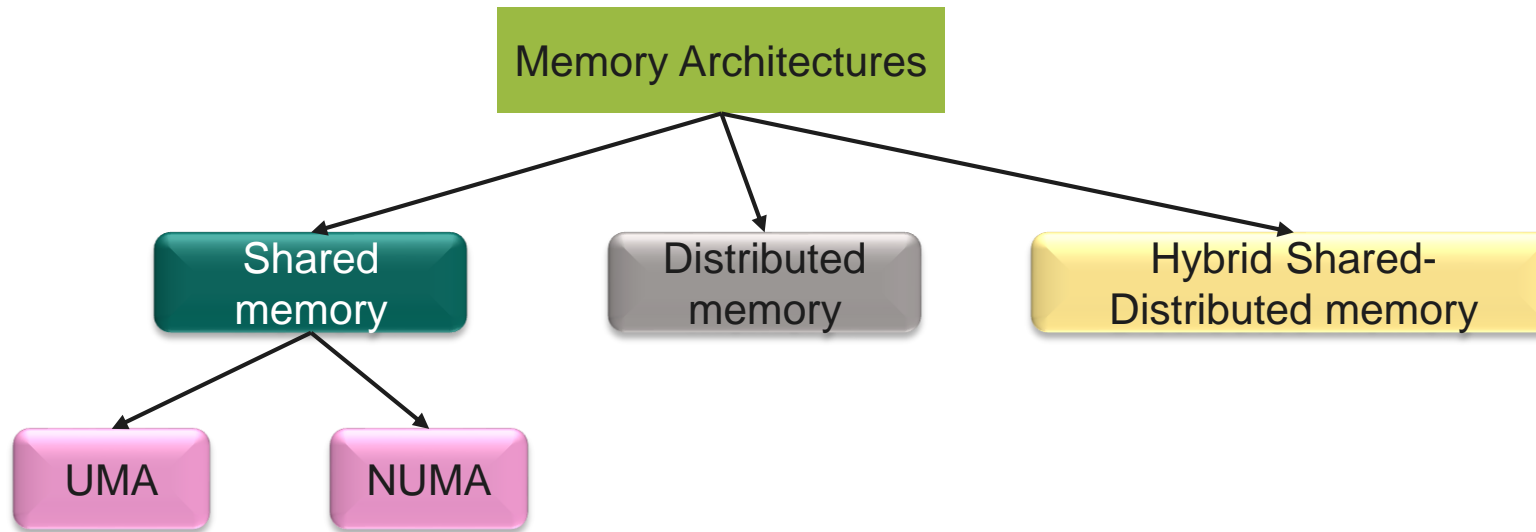
### Corresponding Program sections

– Constant data section / RODATA section

– BSS section

– DATA section

– TEXT section
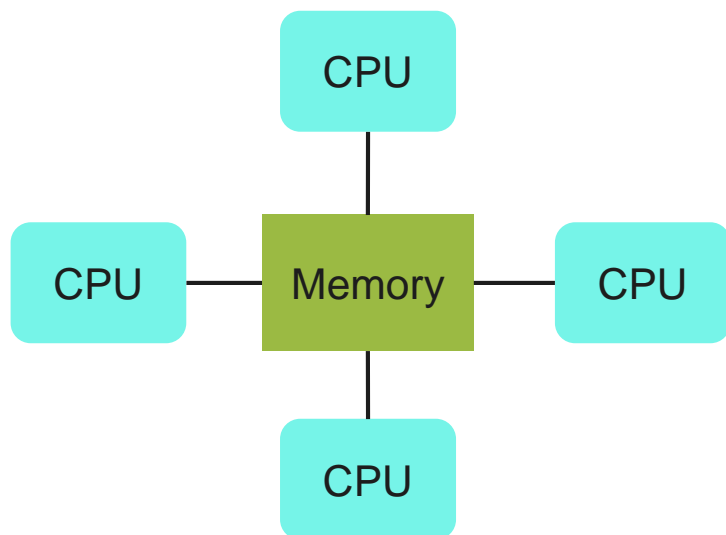
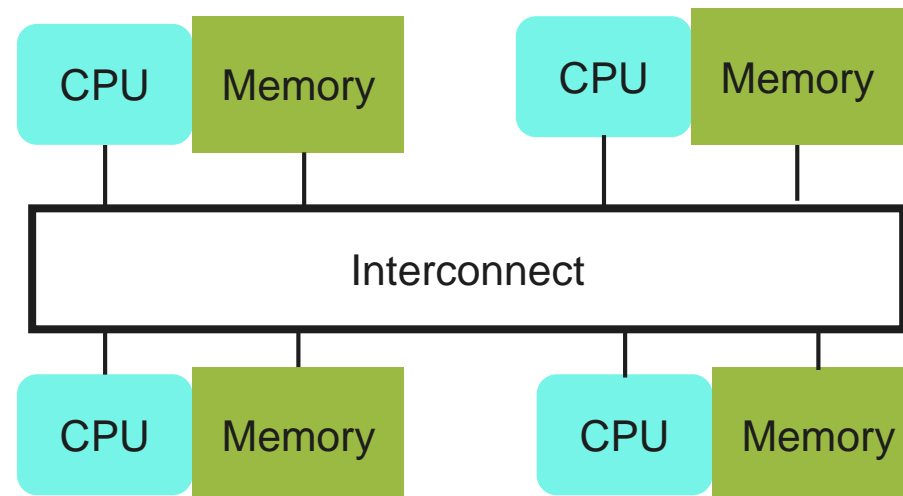### QUESTION: Where should these sections exist?

# Memory architectures

```
                      ┌─────────────────────────┐
                      │  Memory Architectures   │
                      └─────────────────────────┘
                        │           │         │
              ┌─────────┘           │         └──────────┐
              ▼                     ▼                     ▼
     ┌──────────────┐      ┌──────────────┐      ┌──────────────────┐
     │   Shared     │      │ Distributed  │      │  Hybrid Shared-  │
     │   memory     │      │   memory     │      │ Distributed memory│
     └──────────────┘      └──────────────┘      └──────────────────┘
         │      │
      ┌──┘      └──┐
      ▼           ▼
   ┌──────┐   ┌──────┐
   │ UMA  │   │ NUMA │
   └──────┘   └──────┘
```

# Shared Memory – Global address space

**Uniform Memory Access**



All CPUs have equal memory access latency

**Non-Uniform Memory Access**



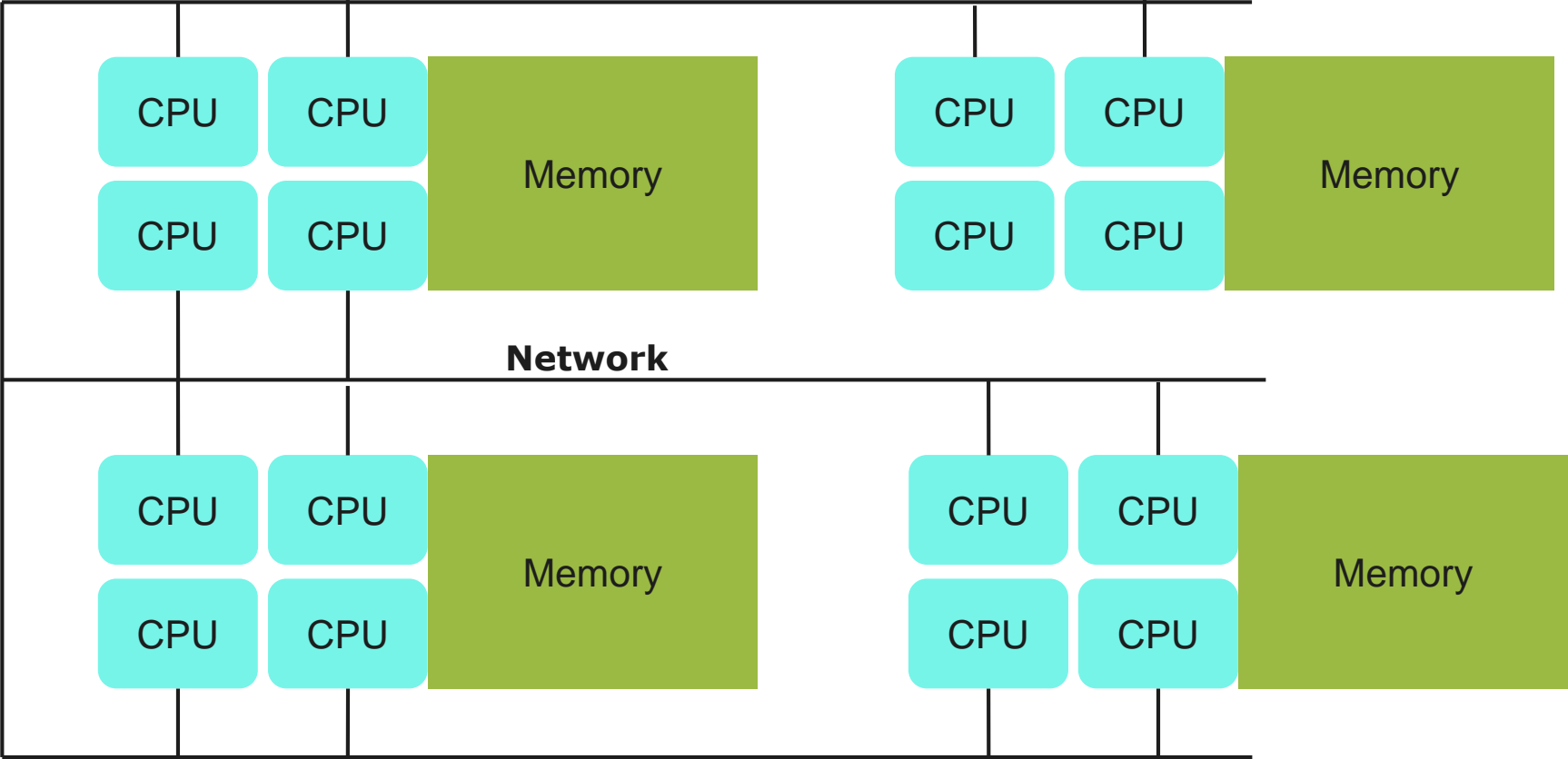Access latencies of local and remote memories are different
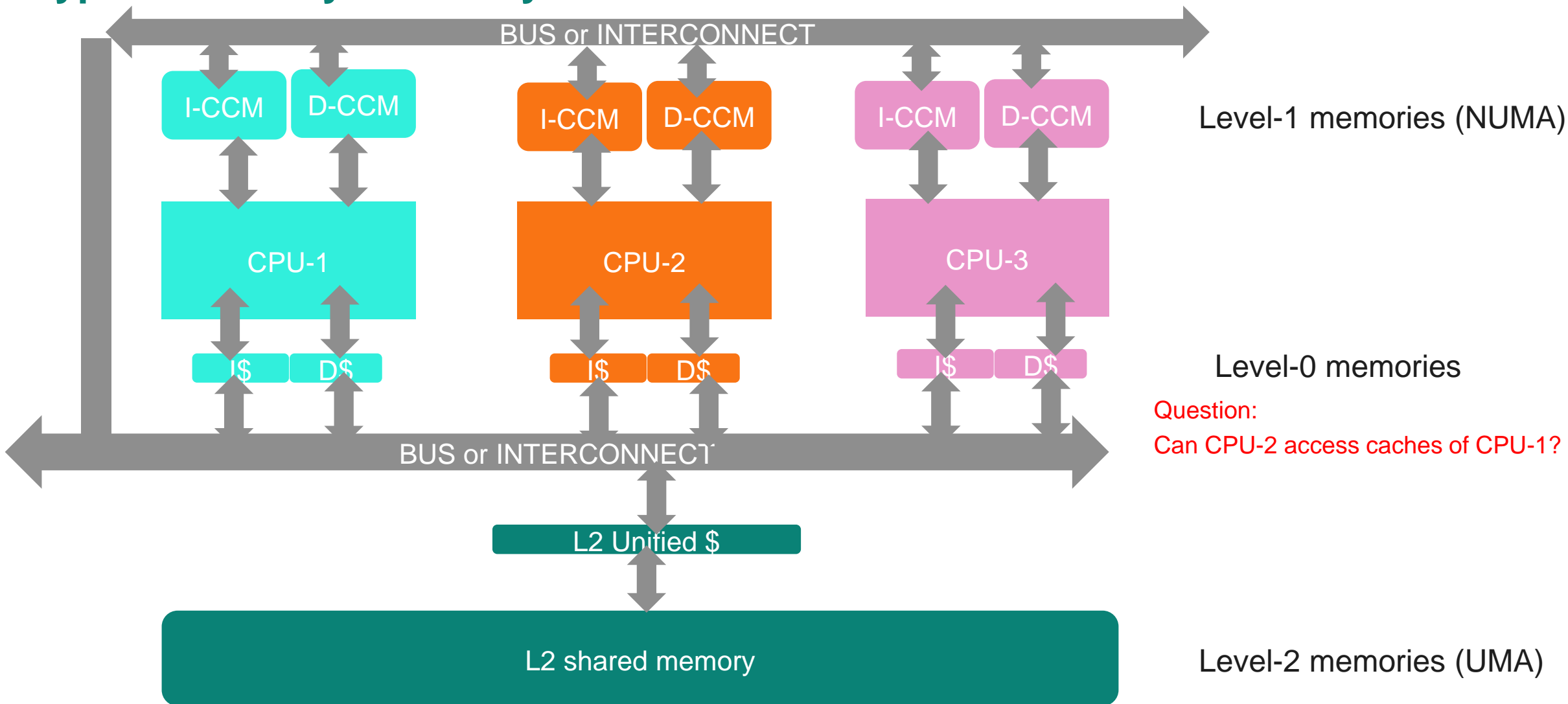
# Distributed Memory – Private address space



A CPU cannot access the memory space of another CPU
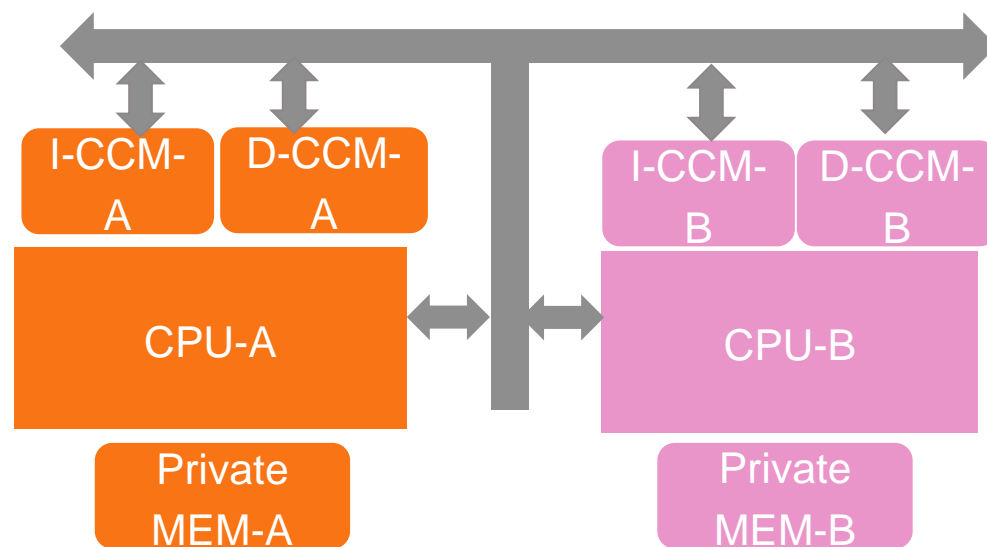
Typical HPC setup

# Hybrid memory model
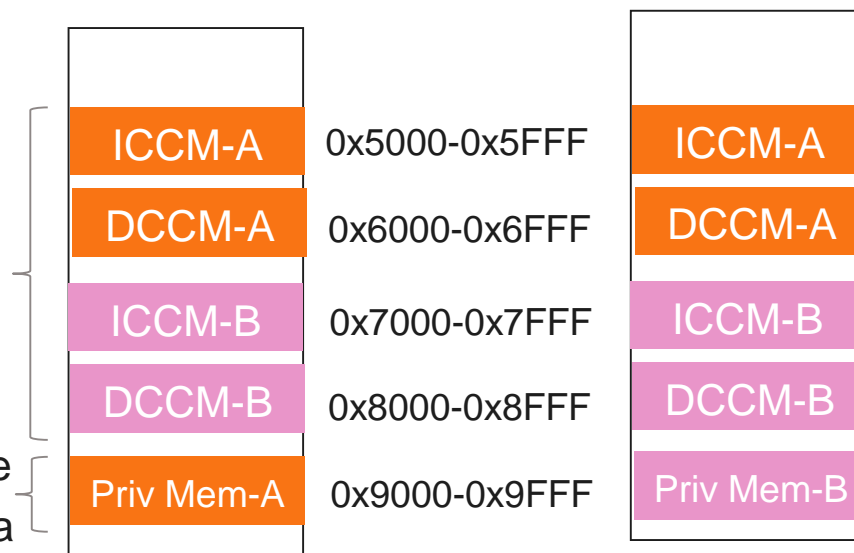
# Typical memory hierarchy on Microcontrollers



BUS or INTERCONNECT

I-CCM   D-CCM
I-CCM   D-CCM
I-CCM   D-CCM

Level-1 memories (NUMA)

CPU-1   CPU-2   CPU-3

I$   D$   I$   D$   I$   D$

Level-0 memories

Question:
Can CPU-2 access caches of CPU-1?

BUS or INTERCONNECT

L2 Unified $

L2 shared memory

Level-2 memories (UMA)

# Mapping of memories to CPU address space



These 4 memories are accessible to both CPUs.

CPU-A's private memory cannot be accessed by CPU-B and vice-versa

| | |
|---|---|
| ICCM-A | 0x5000-0x5FFF |
| DCCM-A | 0x6000-0x6FFF |
| ICCM-B | 0x7000-0x7FFF |
| DCCM-B | 0x8000-0x8FFF |
| Priv Mem-A | 0x9000-0x9FFF |

Both CPUs see ICCM-A at 0x5000, ICCM-B at 0x7000

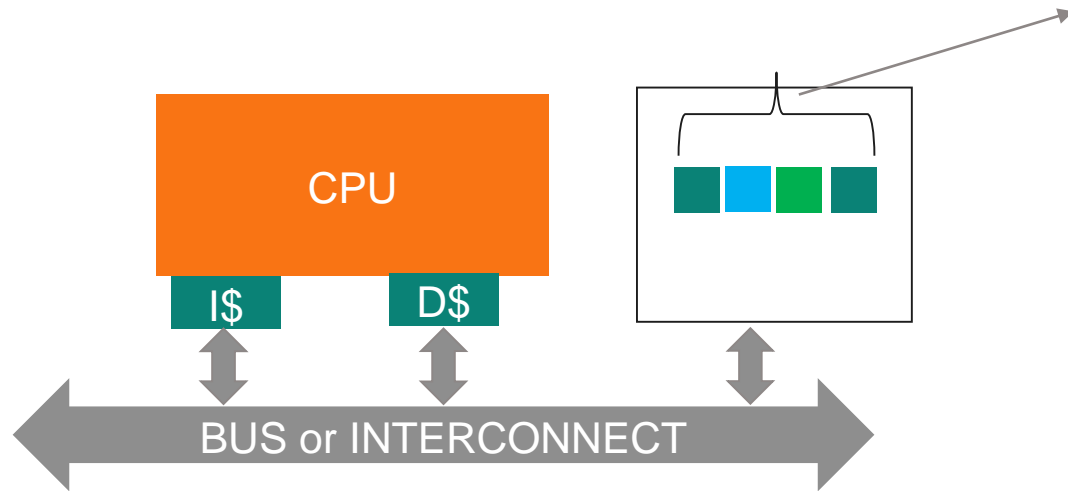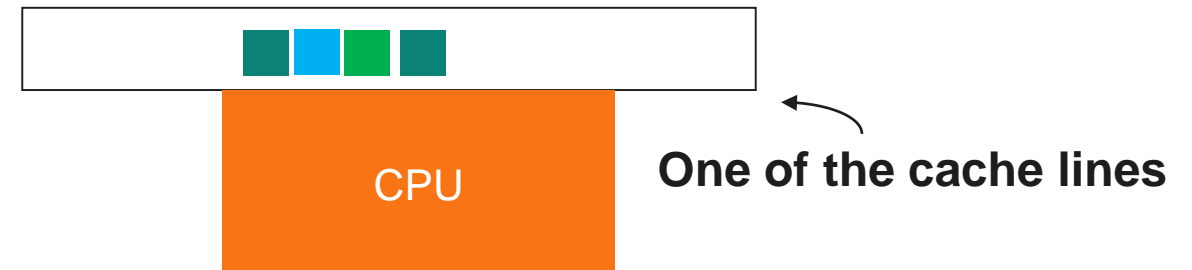When an address of 0x9000 is accessed, the two CPUs end up fetching stuff from their private memories only!

1. Memories are slow, CPUs are fast

2. The pipeline of a CPU will stall waiting for data to be written to/returned from memory

3. A lot of energy is wasted in moving data between CPU and memory

4. If your program has accessed a certain location, it likely may access it yet again soon

5. If your program has accessed a certain location, it likely may access surrounding locations too

**6. What can we do to reduce the pains of slow memory access and energy?**

**ANSWER: Cache memory**

# Cache memories – 2/2



1. When your program attempts to access any of these four words, all four words are brought into the cache

2. This is called CACHE LINE FILL

3. The requested word is now forwarded to CPU by cache

**One of the cache lines**

4. From now on, any access of the 4 address locations will result in fetches from the cache line
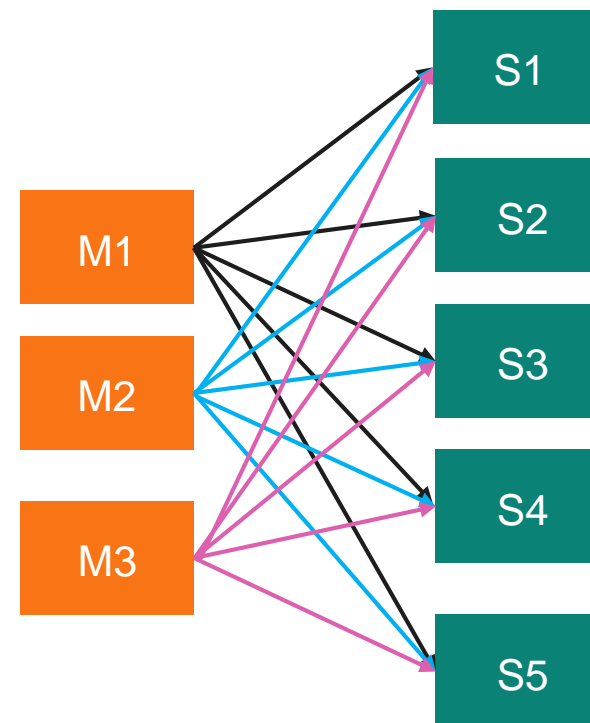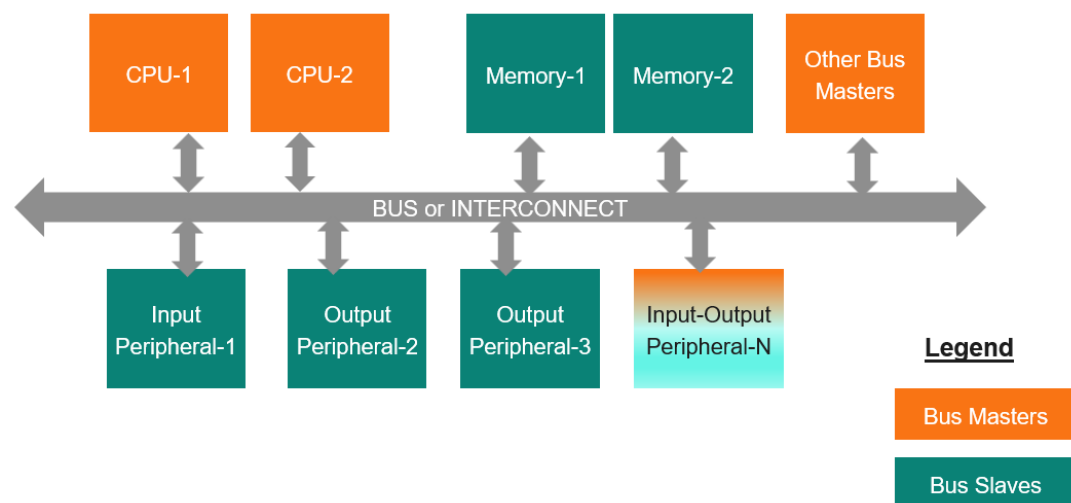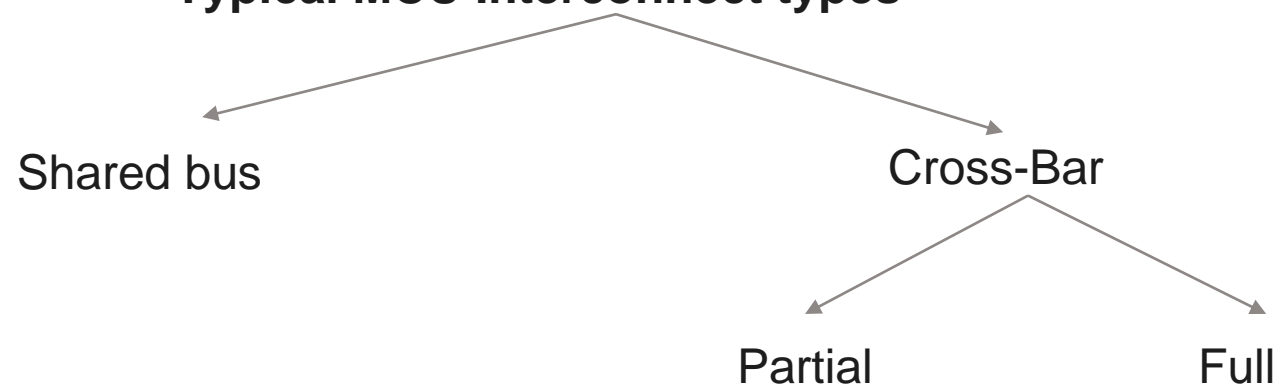
**Homework topics:**

Cache policies (Write Through, Write Back)

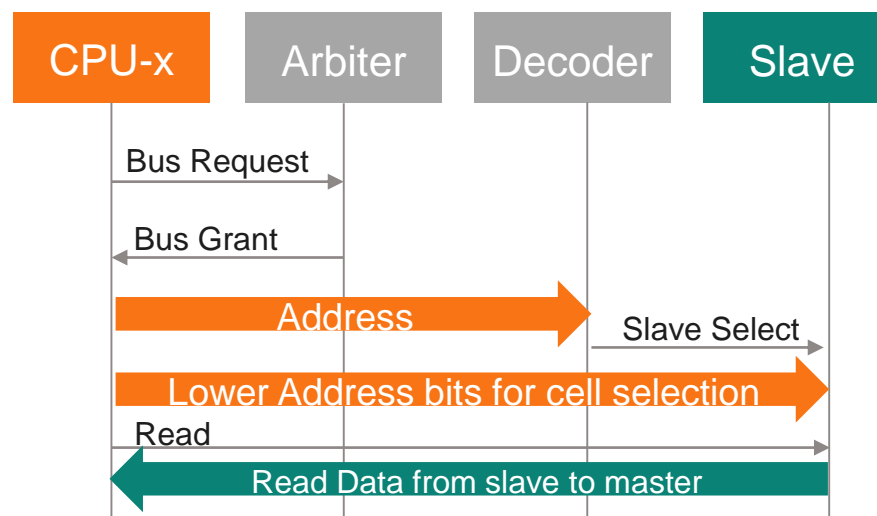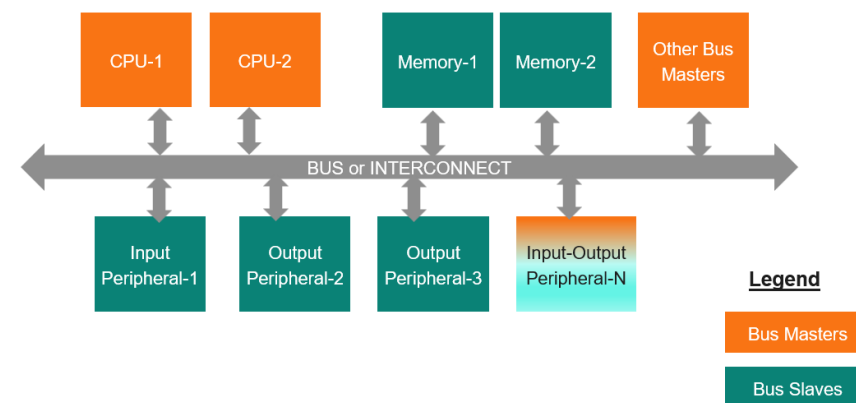Cache line mapping (associativity), Eviction and refill policies,

Cache coherence

**Typical MCU Interconnect types**
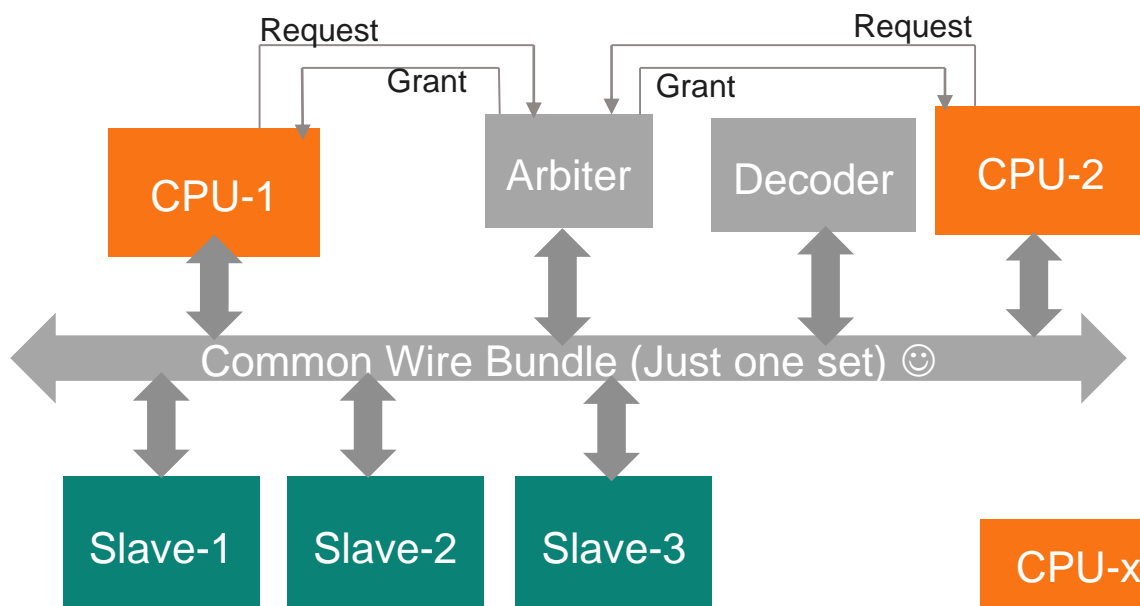
Shared bus

Cross-Bar

Partial

Full



CPU-1  CPU-2  Memory-1  Memory-2  Other Bus Masters

BUS or INTERCONNECT

Input Peripheral-1  Output Peripheral-2  Output Peripheral-3  Input-Output Peripheral-N

**Legend**

Bus Masters
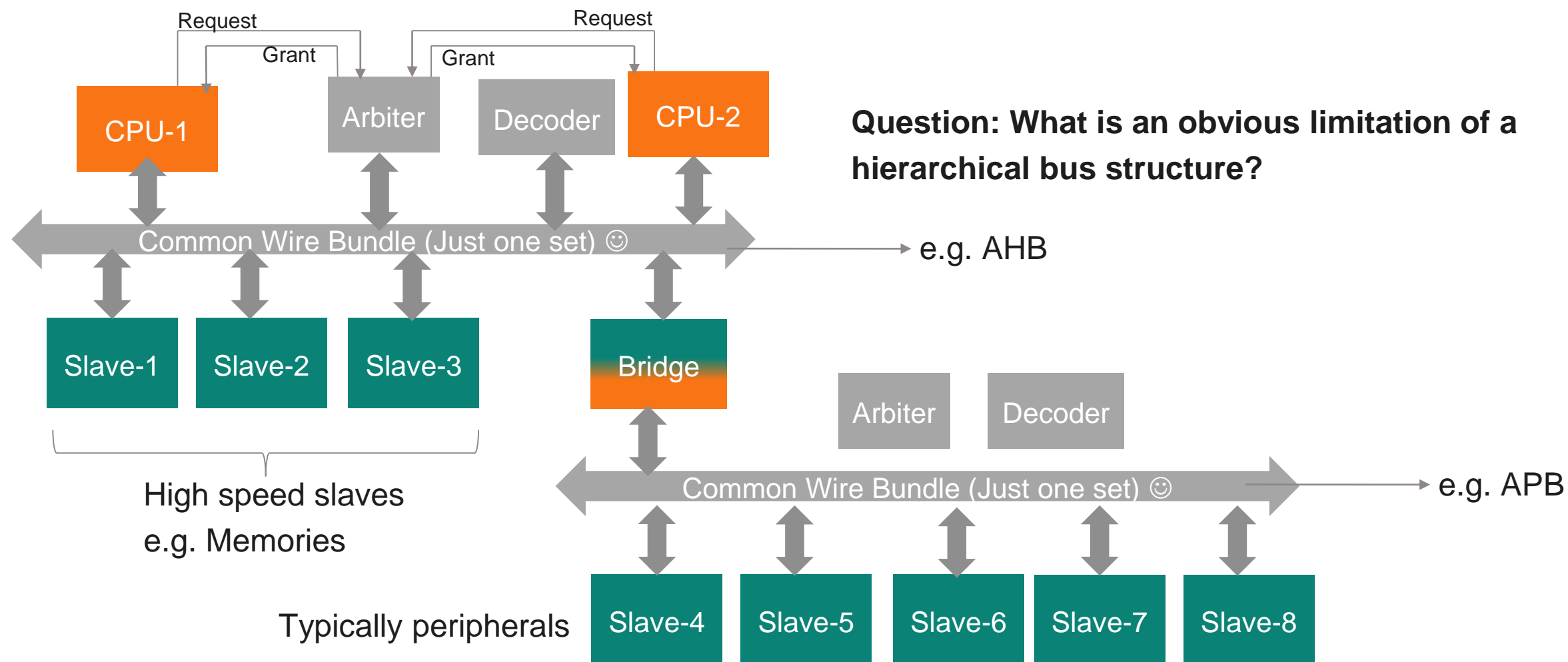
Bus Slaves

M1  M2  M3

S1  S2  S3  S4  S5

15 wire bundles for connecting 3 masters to 5 slaves!

That is already quite a lot!

Need a mechanism to reduce the bundle count.
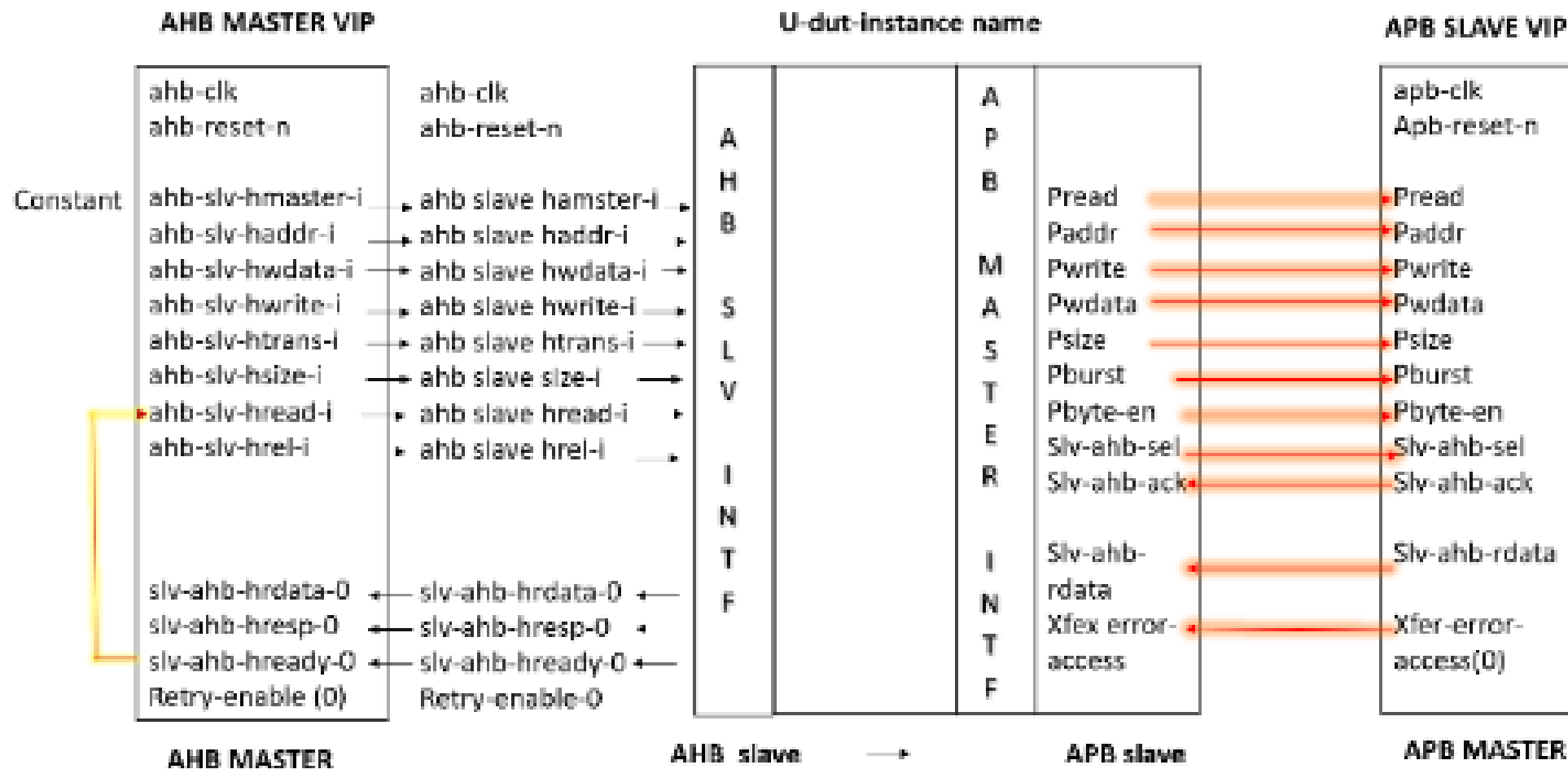
# Interconnects – 2/4

**On most MCUs, this is a typical multi-level bus structure!**

Request

Grant

Request

Grant

CPU-1

Arbiter

Decoder

CPU-2

**Question: What is an obvious limitation of a hierarchical bus structure?**

Common Wire Bundle (Just one set) ☺ → e.g. AHB

Slave-1 Slave-2 Slave-3

Bridge

Arbiter Decoder

High speed slaves
e.g. Memories

Common Wire Bundle (Just one set) ☺ → e.g. APB

Typically peripherals

Slave-4 Slave-5 Slave-6 Slave-7 Slave-8
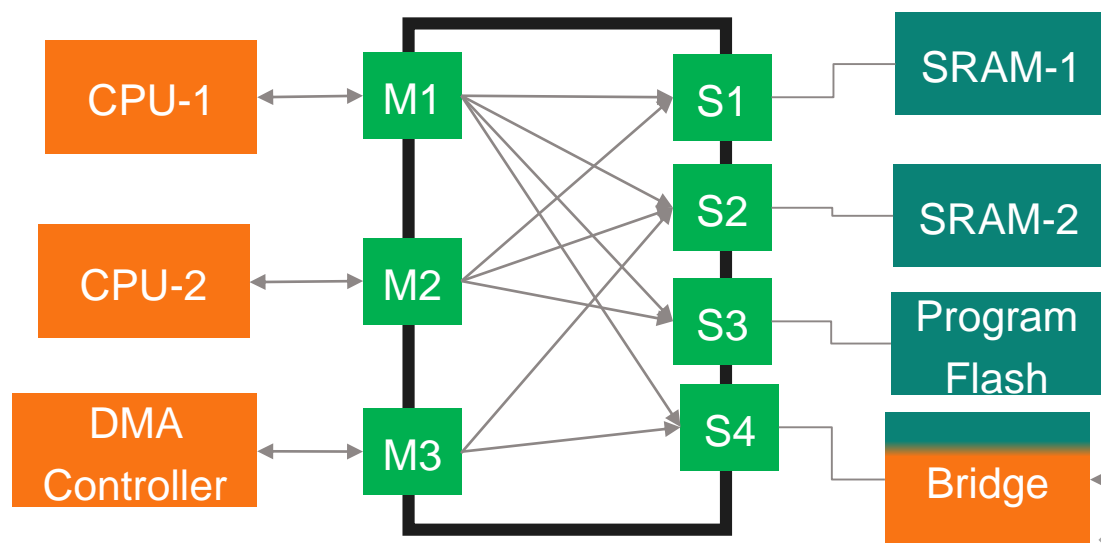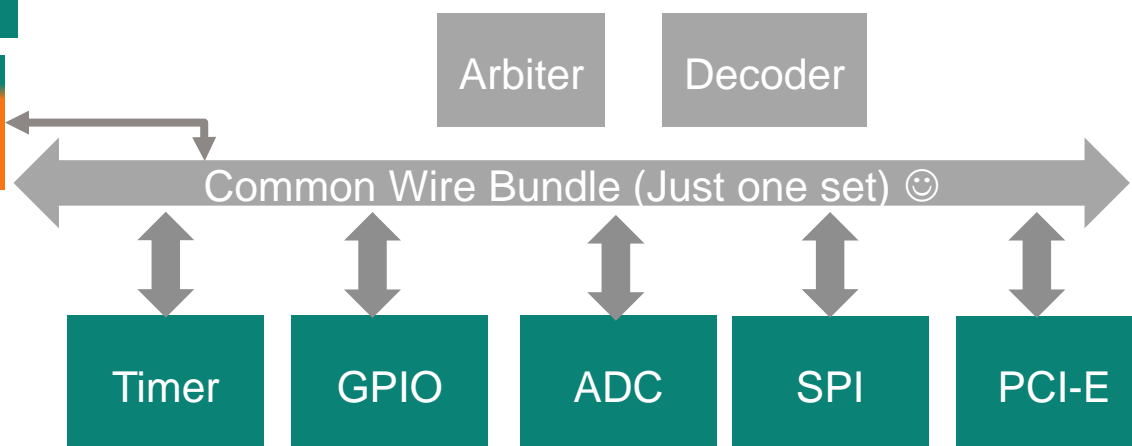
# AHB, APB, Bridge

What we want is multiple bus masters and slaves to be transacting concurrently!



While CPU-1 is writing to SRAM-1, CPU-2 may be reading from Flash.

At the same time, DMA controller may be moving data from PCI-E memory to SRAM-2!

**Question:** Are all masters ports connected to all slave ports?

# PSoC, Finally!