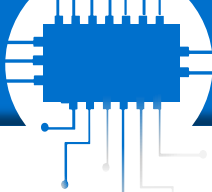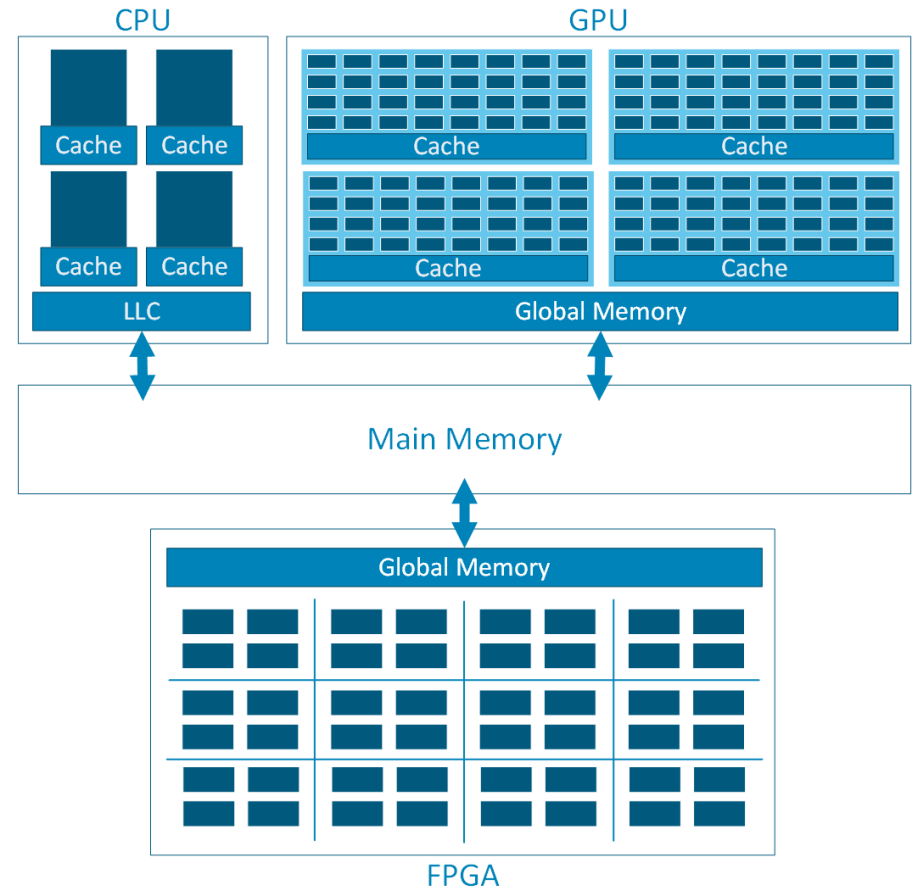# AI for EDGE

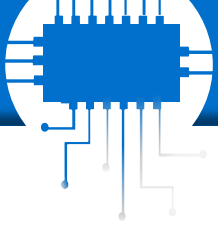# Compute architectures

o Heterogeneous compute environment

o Architecture comparison

- CPU
- GPU
- FPGA

o How to match the workload to compute device

o Alternative for edge computing
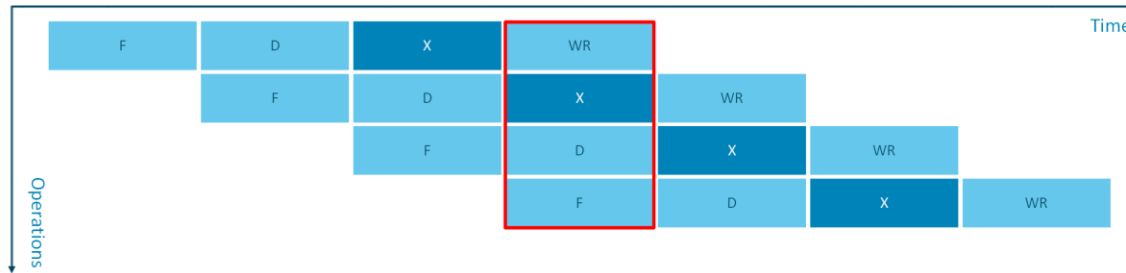
# Shared-memory computing system

- o One or more CPU (each core with own local cache)
- o Shared LLC (last level cache)
- o Set of accelerators:
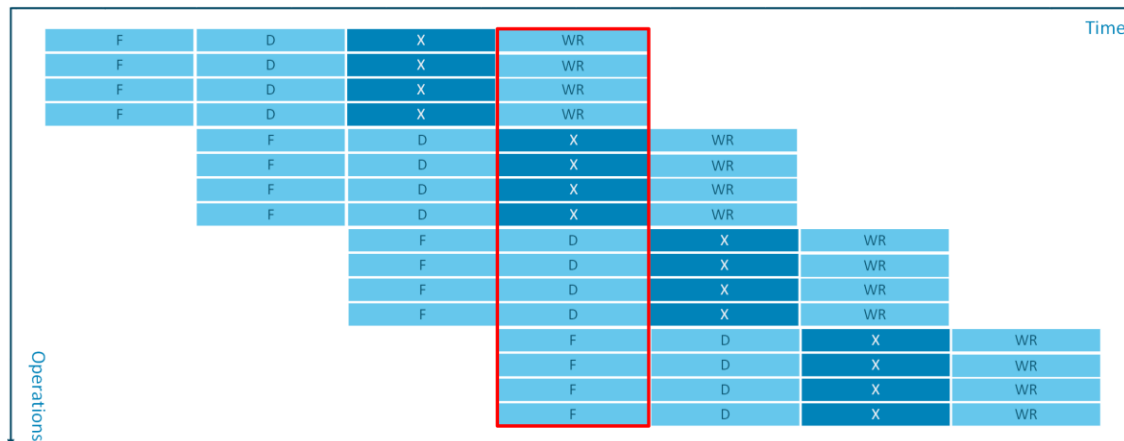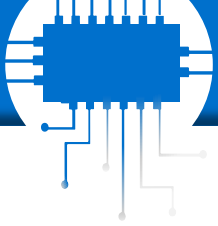  - • GPUs
  - • FPGAs
  - • other specialized hardware

## Scalar Pipelined Execution

Time →

Operations ↓

| F | D | X | WR | | | |
| | F | D | X | WR | | |
| | | F | D | X | WR | |
| | | | F | D | X | WR |

F – fetch
D – decode
X – execute
WR – write to mem/reg

Operation-Level Parallelism

## Superscalar Execution

Time →

Operations ↓

| F | D | X | WR |
| F | D | X | WR |
| F | D | X | WR |
| F | D | X | WR |
| | F | D | X | WR |
| | F | D | X | WR |
| | F | D | X | WR |
| | F | D | X | WR |
| | | F | D | X | WR |
| | | F | D | X | WR |
| | | F | D | X | WR |
| | | F | D | X | WR |
| | | | F | D | X | WR |
| | | | F | D | X | WR |
| | | | F | D | X | WR |
| | | | F | D | X | WR |

F – fetch
D – decode
X – execute
WR – write to mem/reg

Classroom materials
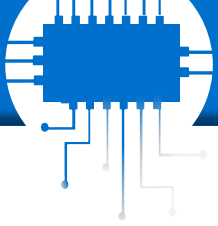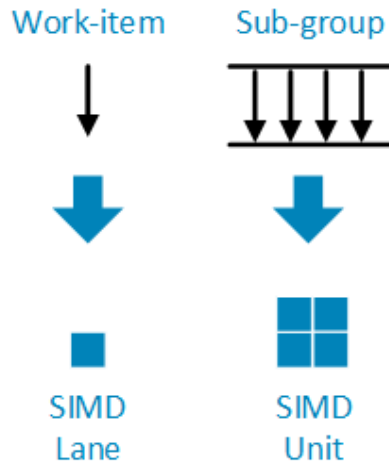
o Advantages:

- automatic parallelism on sequential code
- lower latency when compared to offload acceleration
- accurate branch prediction
- out-of-order superscalar execution
- better energy efficiency than GPU configuration

o Several types of parallelism to achieve performance:

- SIMD (single instruction, multiple data) data parallelism
- Thread-level parallelism (multiple threads, different logical cores)
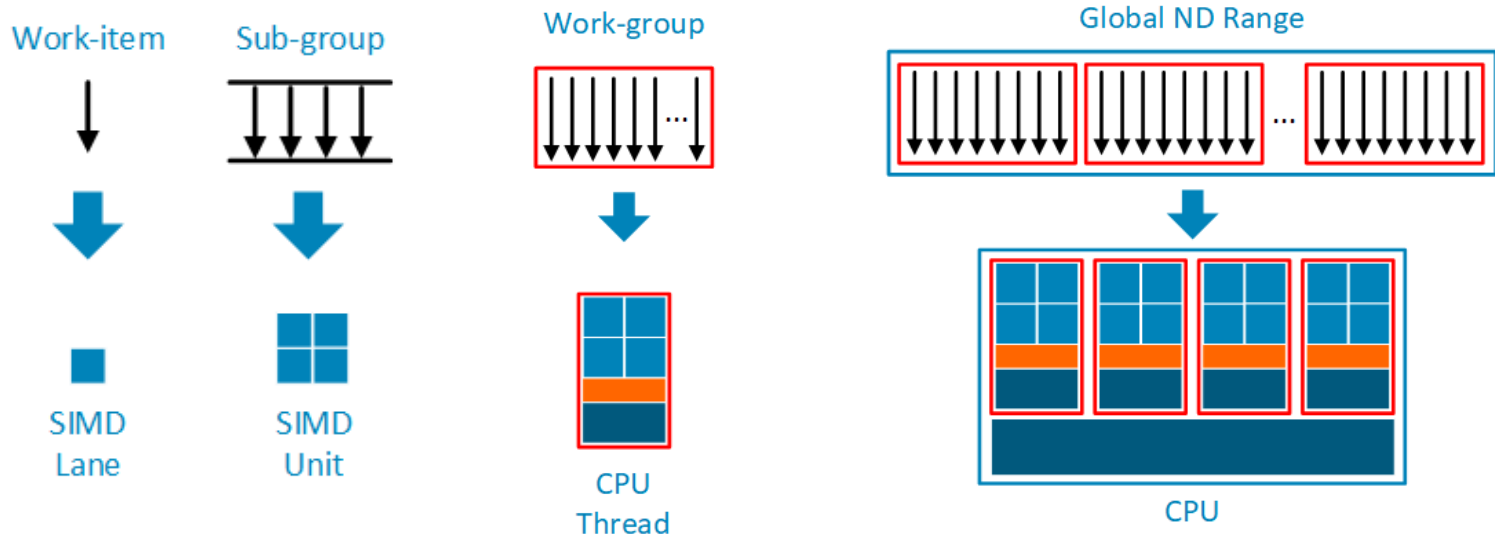- Instruction-level parallelism

o SIMD data parallelism

- each work-item can map to CPU SIMD lane
- vector data tapes are used to explicitly specify SIMD operations
- compiler can perform loop vectorization to generate SIMD code
  - one loop iteration maps to a CPU SIMD lane
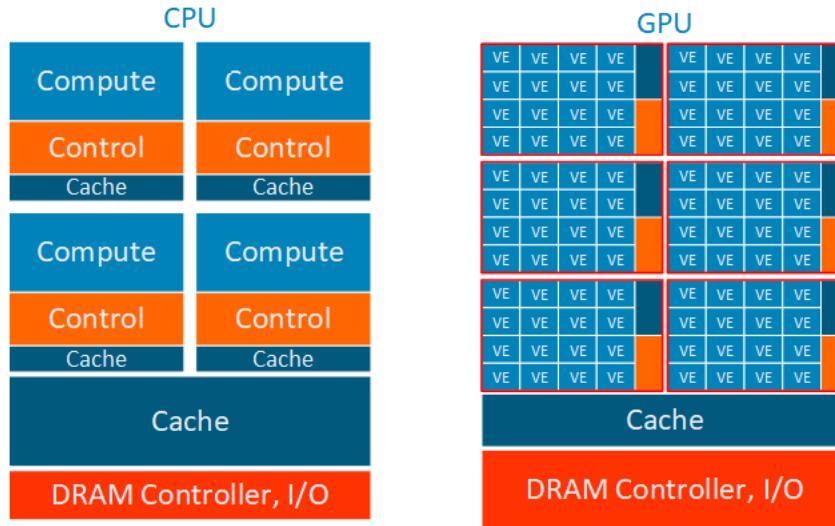  - multiple loop iterations execute together in SIMD fashion



Classroom materials

o Thread-level parallelism

- CPU core and hyper-thread parallelism (a machine with 2 cores and 4 hyper-threads can execute 4 work-groups in parallel)

- different work-groups can execute on different logical cores in parallel

o Massively-parallel, more specialized cores than CPUs

o Optimized for aggregate throughput across all cores

o Vector architecture (efficiently processes vector data)

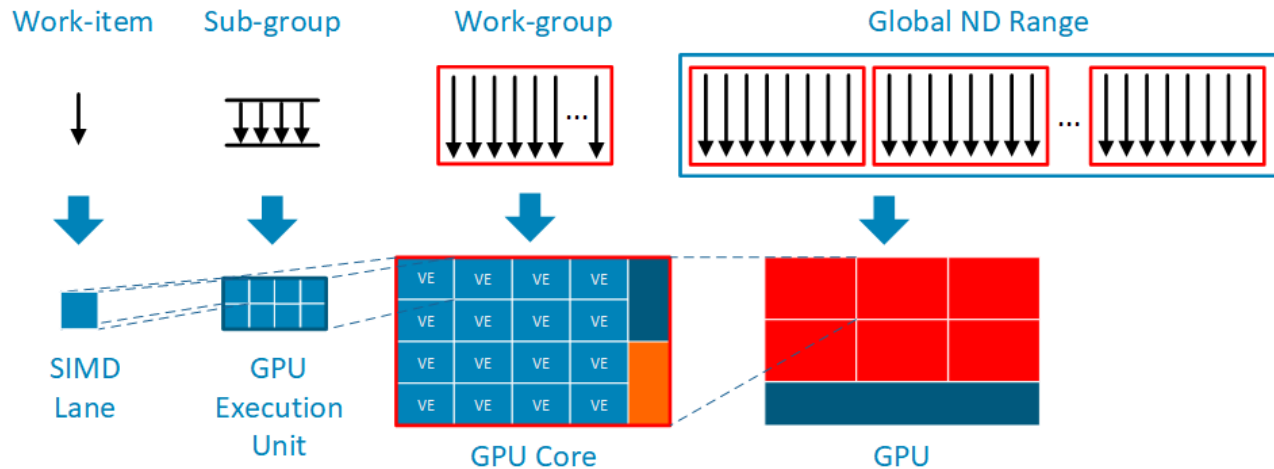o More silicon space to compute, less to cache and control



VE – vector engine (each VE can process multiple SIMD instruction streams)

SM – streaming multiprocessor (multiple VEs combine to form a compute unit with shared local memory and synchronization mechanisms)
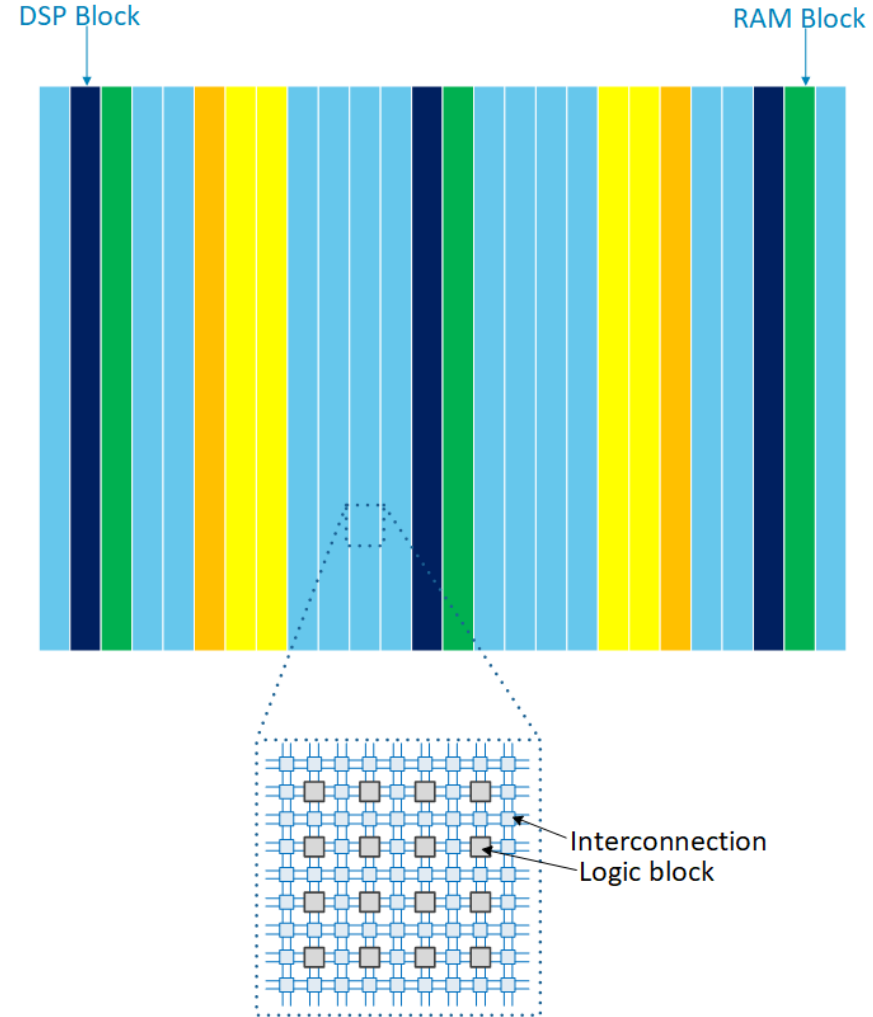
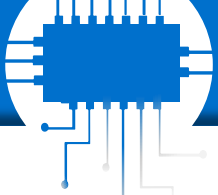GPUs rely on large data-parallel workloads to achieve performance

- every work-item is mapped to a SIMD lane
- sub-groups are mapped to the vector engine
- work-groups, which include work-items that can synchronize and share local data, are assigned for execution on streaming multiprocessors
- single-task kernels are rarely utilized (NDRange kernels are needed to fully populate deep execution pipeline)

- o Massive array of small processing units
- o Resources are connected by a mesh of programmable wires
- o Compute engines are defined by the user
- o Data flows through customized deep pipelines

DSP Block
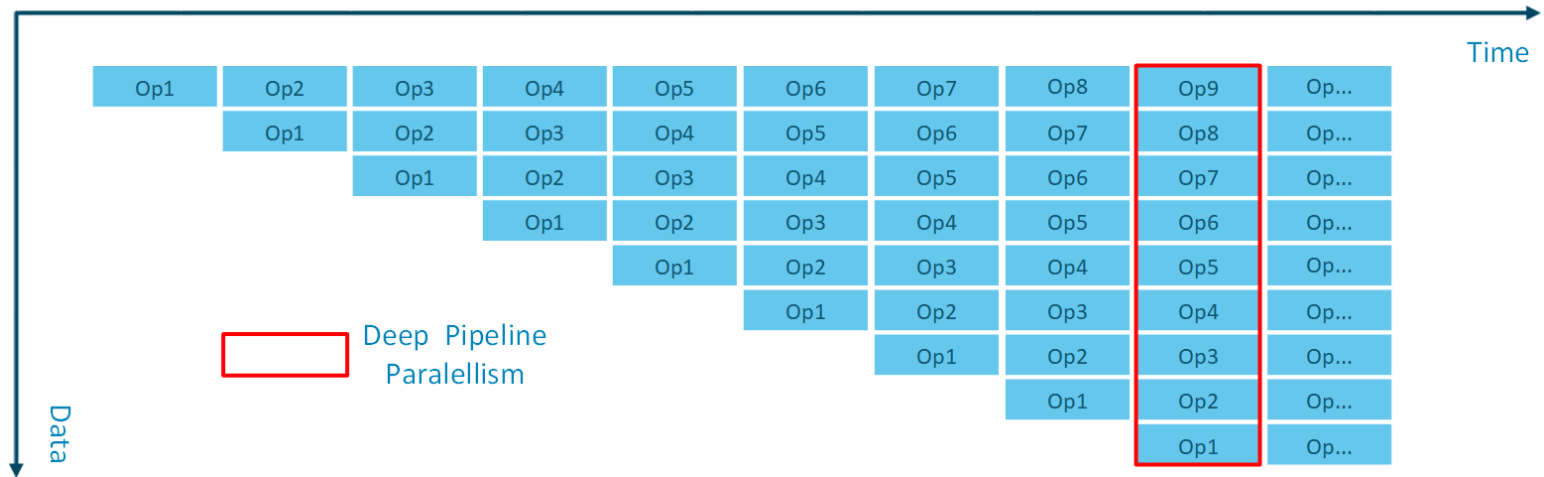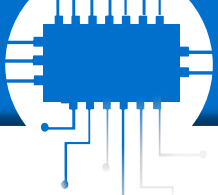
RAM Block

Interconnection
Logic block

Pipeline parallelism can be combined with other types of parallelism:

o task parallelism (multiple pipelines)

o superscalar execution (multiple independent instructions executing in parallel)

FPGA Operation Execution Parallelism
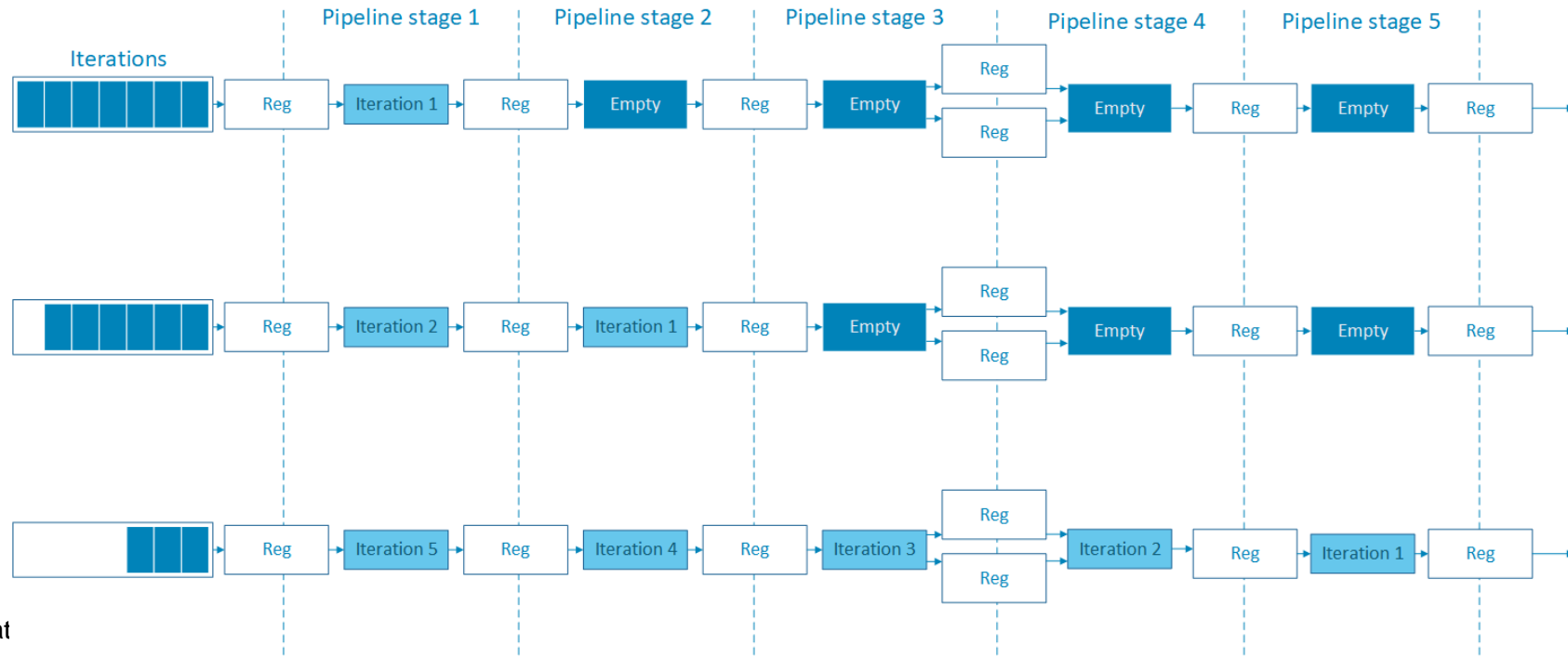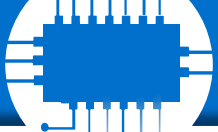
FPGA Advantages:

- Efficiency: no need for control units, instruction fetch and other execution overhead
- Flexibility: can be reconfigured to accommodate different functions and custom data type
- Custom Instructions: instructions not supported by CPUs can be easily implemented and executed on FPGAs
- Rich I/O: FPGA core can interact directly with various memory and custom interfaces

- The operations in the kernel are laid out spatially
- The key to performance is to keep the deep pipeline fully occupied
- With single-task kernels, the FPGA attempts to pipeline loop execution
- Every clock cycle, successive iterations of the loop enter the first stage of the pipeline



Classroom mat

- CPUs:
  - the most widely used generic processors
    (not as compute-dense as GPUs, and not as compute-efficient as FPGAs)
  - modern CPUs support SIMD instructions
  - the most flexible with the broadest library support
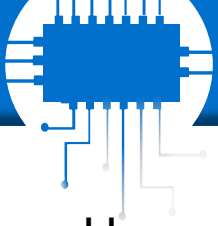- GPUs:
  - the most compute dense, massively data-parallel accelerators
  - employ a single instruction, multiple threads execution model (SIMT)
  - optimized for continuous reads and writes
- FPGAs:
  - compute engines defined by the user
  - reconfigurable, fine-grained datapath (1-bit resolution)
  - dependent kernels are deeply pipelined, independent are executed in parallel

- GPU: image processing, deep learning, data analysis
  - Little or no dependency across processed data
  - Simple control flow (minimal branching and loop divergence)
  - Matches data-types supported by the GPU

- FPGA: genomics sequencing, machine learning, image lossless compression
  - Suitable for algorithms that are easily expressed in serial code and may have dependencies across data elements

- CPU: task orchestration for heterogeneous computing system
  - Can still have superior performance in compute applications when vector, memory, and thread optimizations are applied
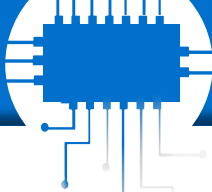
How many transistors does it cost to implement a single multiplication in a network?

1. As a dedicated digital circuit:
   - m.in. 38 transistors/bit

2. As a reconfigurable digital circuit:
   - 32k transistors/weight
     (in 40 LUTs)

o The average consumer CPU will draw between 65 to 85 watts of power, while the average GPU consumes anywhere between 200 to 500 watts.

o A typical microcontroller draws power in the order of milliwatts or microwatts, which is a thousand times less power consumption. This energy efficiency enables the TinyML devices to run on battery power while running ML applications on the edge.

o TinyML with its support for frameworks that include TensorFlow Lite, uTensor, and Arm's CMSIS-NN, brings together AI and small connected devices.
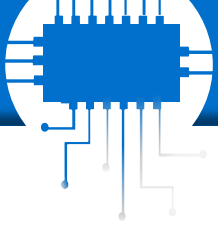
Classroom materials

**Definition:** *This is a concept of implementing <u>mainly deep neural networks</u> directly on embedded devices with <u>highly limited resources</u>.*
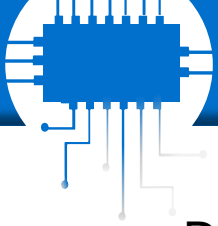
**This approach may include:**

- adapting (reducing) the network architecture to reduced hardware resources
- reduction of power consumption during processing (e.g. for use in battery powered applications)

TinyML can be understood as a network algorithm compression method.

Classroom materials

- Energy efficiency: Microcontrollers consume very little power, which delivers benefits in remote installations and mobile devices.

- Low latency: By processing data locally at the edge, data doesn't need to be transmitted to the cloud for inference. This greatly reduces device latency.

- Privacy: Data can be stored locally, not on cloud servers.

- Reduced bandwidth: With decreased dependency on the cloud for inference, bandwidth concerns are minimized.


- The future of TinyML using MCUs is promising for small edge devices and modest applications where an FPGA, GPU or CPU are not viable options.

Classroom materials

- Data Parallel C++ (programming of heterogeneous systems using SYCL)
- Intel oneAPI programming guide
- Intel DevCloud
- TinyML is bringing deep learning models to microcontrollers
  https://thenextweb.com/news/tinyml-deep-learning-microcontrollers-syndication
- ESP32-CAM: TinyML Image Classification
  https://mjrobot.org/2022/02/10/esp32-cam-tinyml-image-classification-fruits-vs-veggies/

Classroom materials