





# EURO<sup>4SEE</sup>

Optimizing Deep Learning Systems for Hardware Assoc. Prof. Erdem AKAGÜNDÜZ, METU







• PII.a : Existing Solutions: CPU, GPU, TPU, FPGA, ASIC basics

PII.b : Memory hierarchy, bandwidth bottlenecks, movement

costs

• PII.c : Precision







- Now that we understand the fundamental computation types and bottlenecks from Part I, let's look at
  - o the major hardware solutions
  - o and see how they map to those categories.







- Now that we understand the fundamental computation types and bottlenecks from Part I, let's look at
  - o the major hardware solutions
  - o and see how they map to those categories.
- Hardware types are designed to excel at specific clusters of computations.
- Choosing the right hardware depends on the dominant operations in workloads
- We will review CPUs, GPUs, TPUs, FPGAs, and ASICs in this context.

### **CPUs (Central Processing Units)**

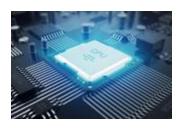






- CPU is a general-purpose processor designed for flexibility, capable of handling diverse tasks and complex control flow, but limited in parallel throughput.
  - O <u>Strengths</u>: Flexible, good at logic + branching + control flow, memory access, OS-level tasks.
  - <u>Limitations</u>: Relatively few cores → less efficient for massive matrix multiplications.
  - Best at: Control-heavy tasks, preprocessing, small inference workloads.

#### **MultiCore CPU Systems**







- Modern CPUs include multiple cores, allowing parallel execution of tasks.
  - Strengths: Excellent for logic, branching, control flow, and sequential operations.
  - <u>Limitations</u>: Parallelism is limited compared to GPUs/TPUs; not as efficient for large-scale matrix multiplications.
  - O Best at: very complex control flow, multiple interface systems

### **GPUs (Graphics Processing Units)**







- GPU is a massively parallel processor optimized for matrix and tensor operations, offering very high throughput but less efficiency for branching-heavy tasks.
  - $\circ$  Strengths: Thousands of cores  $\rightarrow$  excellent at parallel arithmetic (matrix/vector operations).
  - <u>Limitations</u>: Higher latency in single-threaded tasks, less efficient for branching.
  - Best at: Training deep networks, batch inference, dense linear algebra.

#### **GPUs (Graphics Processing Units)**







- Is GPU only an NVIDIA thing?
  - o No, AMD, Intel, and others also produce GPUs.
  - GPUs are a class of hardware (parallel processors optimized for graphics/matrix math), not a brand.
- Why Deep Learning Platforms need only NVIDIA GPU then?
  - o CUDA Ecosystem. NVIDIA developed CUDA, a mature programming framework tightly integrated with deep learning libraries.
  - Popular DL frameworks (TensorFlow, PyTorch, JAX, etc.) are built and optimized first for CUDA.
  - NVIDIA has maintained a strong lead in ML hardware adoption, making it the de facto standard.

## **TPUs (Tensor Processing Units)**







- TPU is a domain-specific processor tailored for deep learning workloads, highly efficient for matrix multiplications and convolutions but less versatile outside ML.
  - Strengths: Specialized for tensor/matrix multiplications common in deep learning.
  - O <u>Limitations</u>: Less flexible than GPUs, tailored to DL workloads only.
  - O Best at: Large-scale training and inference in tensor-heavy models.

#### **FPGA**







- FPGA is a reconfigurable processor that allows custom hardware pipelines, offering low latency and energy efficiency at the cost of programming complexity.
  - Strengths: Reconfigurable hardware, can be optimized for specific arithmetic or logic patterns, low power requirements
  - o <u>Limitations</u>: Programming complexity, slower development cycle.
  - Best at: Custom accelerations (low-latency inference, energy-efficient specialized operations).





- ASIC is a fixed-function processor designed for maximum efficiency in a specific workload, achieving the best performance-per-watt but with no flexibility once fabricated.
  - Strengths: Fully customized for one set of operations, highest efficiency and energy savings.
  - <u>Limitations</u>: No flexibility "locked" to one task.
  - Best at: Dedicated inference accelerators, production environments with fixed workloads.

## Comparison



HW	Key Strengths	Limitations	Typical Computation Type	Precision Trends
CPU	General-purpose, flexible, strong control	Low parallelism, lower	Scalar/vector ops, branching-	FP32 / FP64
	flow, large ecosystem	FLOPs/Watt	heavy tasks	
GPU	High parallelism, excellent for matrix/tensor	High power use, less	Matrix multiplications, tensor	FP32, FP16, BF16,
	ops, optimized libraries (CUDA, ROCm)	efficient for branching	ops, data-parallel workloads	INT8
TPU	Optimized for deep learning (matmul +	Limited to ML workloads,	Tensor ops (matmul, conv),	FP32, BF16, INT8
	conv), high throughput for inference &	less general	low precision inference	
	training			
FPGA	Customizable, low latency, energy efficient	Harder to program, long	Pipelined ops, fixed-function	Any (user-defined),
		development cycle	acceleration	often INT8 or lower
ASIC	Maximum efficiency for a fixed workload,	Inflexible, huge upfront	Highly repetitive, domain-	Any (user-defined),
	very low energy cost	design cost	specific ops	often INT8 or lower





Case Study	Best Hardware	Why
Edge Devices		Low energy, compact form factor, efficient for small/medium matmul, inference focus
Cloud / Data Centers		High throughput, flexible for multiple tenants, strong software ecosystem
Supercomputers / HPC		HPC tasks often require control flow (CPU) + massive tensor ops (GPU), also support for FP64 precision in scientific computing

#### **Next: Part II.b**



• PII.a : Existing Solution: CPU, GPU, TPU, FPGA, ASIC basics

• PII.b : Memory hierarchy, bandwidth bottlenecks, movement

costs

PII.c : Precision (FP32 →INT8) and energy efficiency implications



## Thanks!





This project has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement No 101191697. The JU receives support from the Digital Europe Programme and Germany, Türkiye, Republic of North Macedonia, Montenegro, Serbia, Bosnia and Herzegovina.