

Christof Teuscher

ECE 410/510: Hardware for AI and ML

GPUs + CUDA

Portland State University
Department of Electrical and Computer Engineering (ECE)

www.teuscher-lab.com
teuscher@pdx.edu



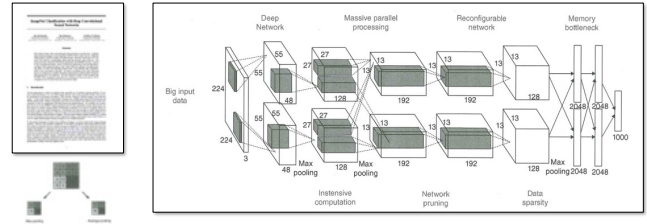
teuscher:Lab
teuscher-lab.com

Portland State
UNIVERSITY

Christof Teuscher teuscher@pdx.edu

www.teuscher-lab.com/teaching

AlexNet: A classical deep convolutional neural network



AlexNet paper: https://papers.nips.cc/paper_files/paper/2012/file/c39862d3b9d6b76c8436e924a68c45b-Paper.pdf

Liv & Law, Artificial Intelligence Hardware Design, 2021

teuscher:Lab
teuscher-lab.com

Portland State
UNIVERSITY

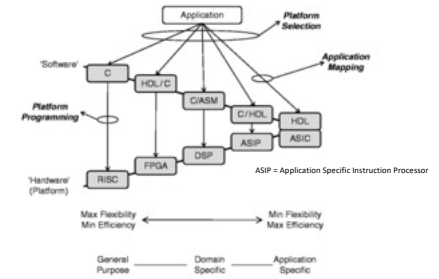
Flexibility vs performance

1. Why is running a neural network simulation on a general purpose processor inefficient?
2. How can we improve performance?
3. What is a better processor architecture?

teuscher:Lab
teuscher-lab.com

Portland State
UNIVERSITY

The HW/SW co-design space

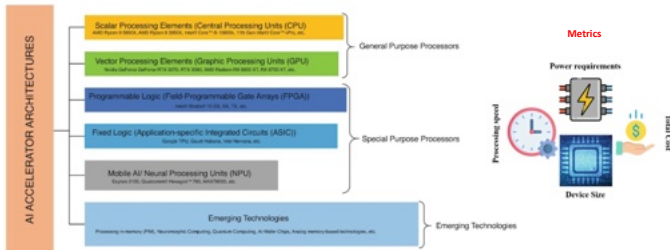


Schaumont, 2013

teuscher:Lab
teuscher-lab.com

Portland State
UNIVERSITY

Hardware for AI/ML

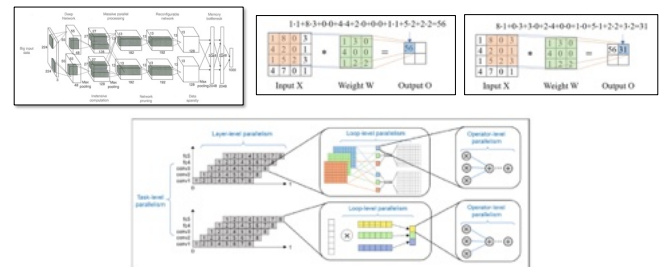


Mahra et al., 2023

teuscher:Lab
teuscher-lab.com

Portland State
UNIVERSITY

Exploiting parallelism

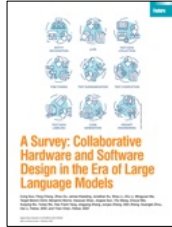


teuscher:Lab
teuscher-lab.com

Portland State
UNIVERSITY

Paper discussion

- What did you learn?



Paper discussion (1)

- The difference between these papers shows how the field of co-design has evolved **from general principles to specialized approaches** for complex AI systems over nearly three decades.
- LLMs face unique challenges compared to traditional deep learning systems, requiring extensive computational resources, high energy consumption, and complex software optimizations.
- LLMs differ from CNNs in their memory requirements and processing patterns, requiring unique co-design methodologies.

Paper discussion (2)

- Co-design is about a systems perspective: meeting system-level objectives.
- Co-design is the key to designing efficient systems.
- Designing HW is increasingly like designing SW.
- Flexibility vs performance trade-off.
 - E.g., reconfigurability increases usability, but does not increase performance.
- Co-design tools: modeling, validation, implementation.
- Identify critical segments of SW and run them on special HW.
- LLMs require unique optimization strategies.
 - GPT-3: 175 billion parameters, 350GB of GPU memory just for parameters
- GPUs are the cornerstone of modern deep learning infrastructure.

Paper discussion (3)

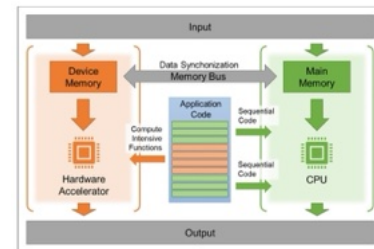
Table 1
Top 2 accelerators for efficiency

Name	Platform	Model	Energy efficiency (TOPS/W)	Quantization/Sparsity	Year
TruSOM [18]	ASIC (approx 28nm)	BERT	35.5 (INT8)	Sparsity	2022
SPRINT [19]	FPGA	GPT-2	-	Sparsity	2022
8-Former [16]	FPGA (approx 32 nm)	BERT	13.44 (INT8)	-	2022
SPRINT [19]	FPGA (approx 32 nm)	GPT-2, BERT	-	Quantization/Sparsity	2022
TruSOM [18]	FPGA (approx 32 nm)	GPT-2, BERT	18.84	Sparsity	2022
Moby [19]	ASIC (approx 32 nm)	BERT	688.8x RTX 2080Ti	Quantization/Sparsity	2022
LoCoNet [16]	ASIC (approx 32 nm)	BERT	8x GEMV (FP16)	Quantization/Sparsity	2022
STP [18]	ASIC (approx 32 nm)	BERT	18.1 (FP16)	Sparsity	2023
ANIMA [18]	FPGA (approx 42 nm)	BERT	-	Sparsity	2023
TP-RAFT [18]	ASIC (approx 42 nm)	BERT	-	Sparsity	2023
TIC-SAT [18]	ASIC (approx 42 nm)	BERT	0.48 (FP16)	Sparsity	2023
Transformer-SPR [18]	FPGA	BERT	-	Sparsity	2023
ANCT [18]	ASIC (approx 42 nm)	BERT	84.84x V100	Quantization/Sparsity	2023
TruSOM [18]	ASIC (approx 42 nm)	BERT	16.25x NVIDIA A100	Sparsity	2023
SPRINT [19]	ASIC (approx 42 nm)	BERT	4x GEMV (FP16)	Quantization/Sparsity	2023
C-Transformer [18]	ASIC (approx 42 nm)	GPT-2	33.4 (INT8)	-	2024
SPRINT [19]	FPGA (approx 42 nm)	LLAMA-GPT	6.67x A100 (FP16)	-	2024
ANIMA [18]	FPGA (approx 42 nm)	LLAMA-GPT	~100x	Sparsity	2024
ANIMA [18]	FPGA (approx 42 nm)	LLAMA-GPT	2.67x NVIDIA A100 (FP16)	-	2024
ANIMA [18]	FPGA (approx 42 nm)	LLAMA-GPT	2.67x NVIDIA A100 (FP16)	-	2024

Guo et al., 2025

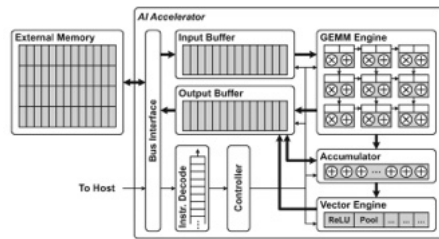
Hardware for AI/ML

General co-processing mechanism



Mahra, Ashutosh; Cha, Jaekwang; Park, Hyunbin; Kim, Shih; Artificial Intelligence and Hardware Accelerators, 2025

General AI accelerator architecture



GEMM - General Matrix Multiply

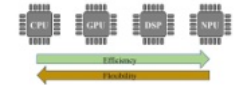
Tensor cores!

A vector is a 1st-order tensor
A matrix is a 2nd-order tensor.

Mahra, Ashutosh; Cha, Jaekwang; Park, Hyunbin; Kim, Shihy. Artificial Intelligence and Hardware Accelerators, 2025

Hardware for AI/ML

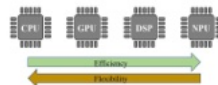
Processor	Power Consumption	Strengths	Limitations
CPU			
GPU			
FPGA			
ASIC			
Vision Processing Unit (VPU)			
Tensor Processing Unit (TPU)			



Mahra et al., 2023

Hardware for AI/ML

Processor	Power Consumption	Strengths	Limitations
CPU	High	<ul style="list-style-type: none"> Flexible General purpose processing Complex instructions and code System management 	<ul style="list-style-type: none"> Flexible memory access bottlenecks Low cores (2-16)
GPU	High	<ul style="list-style-type: none"> Parallel cores > 1000s of cores High performance AI processing 	<ul style="list-style-type: none"> Power consumption Large footprint
FPGA	Medium	<ul style="list-style-type: none"> Configurable logic gate Flexible In-field re-programmability 	<ul style="list-style-type: none"> Programming complexity
ASIC	Low	<ul style="list-style-type: none"> Custom logic designed with libraries Raster processing Small footprint 	<ul style="list-style-type: none"> Fixed function Expensive custom design
Vision Processing Unit (VPU)	Ultra low	<ul style="list-style-type: none"> Custom logic designed with libraries Raster processing Small footprint 	<ul style="list-style-type: none"> Limited dataset and batch size Limited network support
Tensor Processing Unit (TPU)	Low to medium	<ul style="list-style-type: none"> Specialized tool support Proprietary design Limited framework support 	



Mahra et al., 2023

Graphics Processing Units (GPUs)

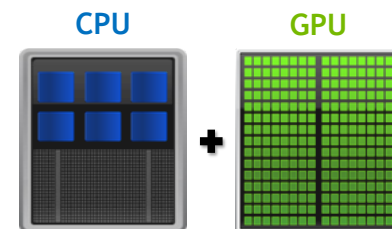
A bit of history

- Over several decades, GPUs evolved from a single core, fixed function hardware used for graphics solely, to a set of programmable parallel cores.
- The history of modern GPUs starts in 1995 with the introduction of the first 3D add-in cards.
- 1996: 3DFx's Voodoo graphics card took over about 85% of the market. Cards that could only render 2D became obsolete very fast.
- 1999: NVIDIA's GeForce 256, the "world's first GPU."



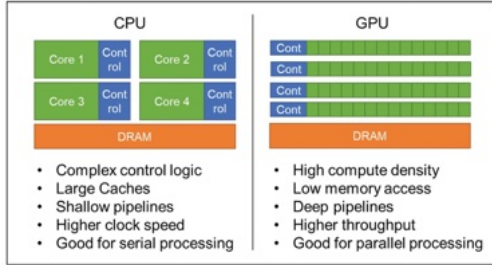
3dfx

Multi-core vs GPU



NVIDIA

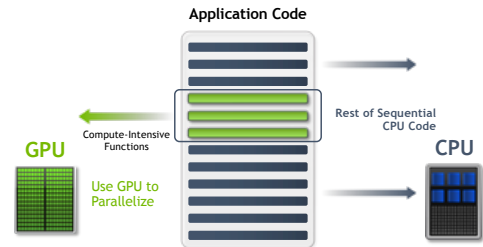
Multi-core vs GPU



SIMT (Single Instruction, Multiple Threads) is a parallel processing model used in GPUs where multiple threads execute the same instruction simultaneously, but on different data.

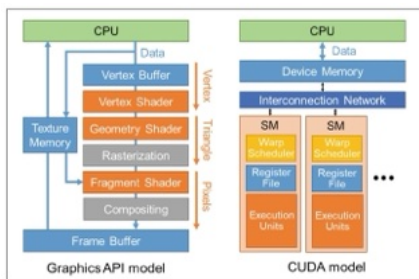
Mahra, Ashutosh; Cha, Jaekwang; Park, Hyunbin; Kim, Shihoo. Artificial Intelligence and Hardware Accelerators, 2025

Small Changes, Big Speed-up



NVIDIA

CUDA model



Mahra, Ashutosh; Cha, Jaekwang; Park, Hyunbin; Kim, Shihoo. Artificial Intelligence and Hardware Accelerators, 2025

NVIDIA Volta Microarchitecture (2013-2017)



GPU features	Nvidia Tesla P100	Nvidia Tesla V100	Nvidia A100
GPU codename	GP100	GV100	GA100
GPU architecture	Nvidia Pascal	Nvidia Volta	Nvidia Ampere
Compute capability	6.0	7.0	8.0
Threads / warp	32	32	32
Max warp / SM	64	64	64
Max threads / SM	2048	2048	2048
Max thread blocks / SM	20	32	32
Max 32-bit registers / SM	65536	65536	65536
Max registers / block	65536	65536	65536
Max registers / thread	255	255	255
Max thread block size	1024	1024	1024
FP32 cores / SM	64	64	64
Ratio of SM registers to FP32 cores	1024	1024	1024
Shared Memory Size / SM	64 KB	Configurable up to 96 KB	Configurable up to 164 KB

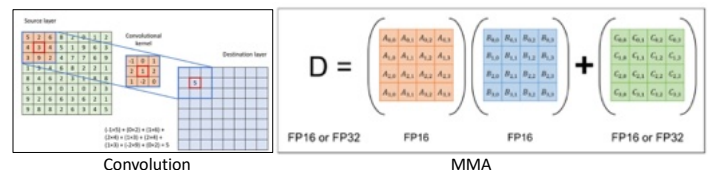
- First chip to feature tensor cores
- A GV100 SM is partitioned into four processing blocks

Warp

- A "warp" is a group of threads that execute instructions simultaneously on a GPU (Graphics Processing Unit).
- The warp scheduler is a component of the GPU architecture that manages the execution of these thread groups.
- A warp typically consists of 32 threads in NVIDIA GPUs (AMD calls similar groupings "wavefronts"). These threads execute the same instruction at the same time in a **SIMT** (Single Instruction, Multiple Thread) architecture.
- The warp scheduler is responsible for:
 - Deciding which warps are ready to execute
 - Allocating compute resources to warps
 - Managing warp context switching when threads need to wait for memory operations.
 - Hiding memory latency by switching between warps

Tensor cores

- Each tensor core executes MMA (matrix multiplication and accumulation), $D = A \times B + C$ on a 4×4 floating-point matrix.
- They are built to enhance deep learning by boosting matrix arithmetic.
- Low-precision operations often supported.



Mahra et al., 2023

Streaming vs tensor cores

SM Cores (Streaming Multiprocessors)

- These are the traditional GPU computing units
- Handle general-purpose parallel computing tasks
- Process standard graphics rendering workloads
- Execute regular FP32 (single precision) and INT32 operations
- Responsible for most traditional GPU workloads

Tensor Cores

- Specialized hardware accelerators designed specifically for matrix operations
- Optimized for deep learning and AI workloads
- Dramatically accelerate matrix multiply-accumulate operations (critical for neural networks)
- Can perform mixed-precision calculations (like FP16 and FP32)
- Provide order-of-magnitude speedups for deep learning training and inference
- Introduced with NVIDIA's Volta architecture and improved in subsequent generations

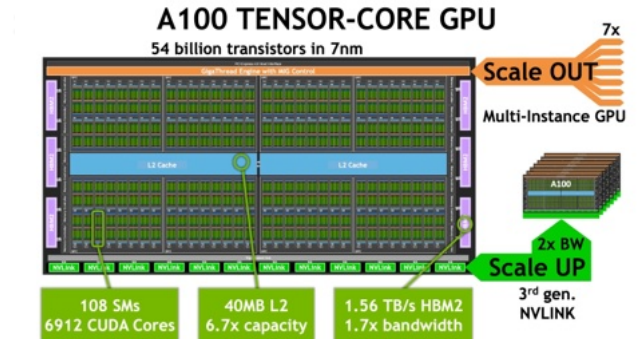
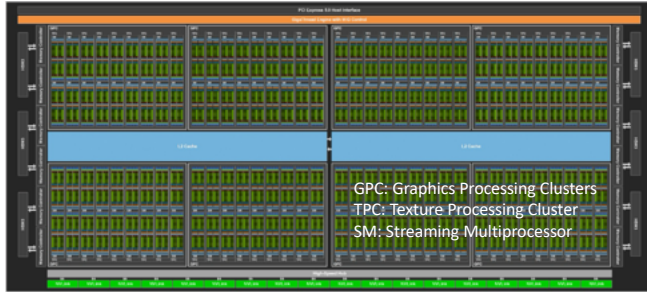
Consumer-grade GPUs

Graphics Card	GeForce RTX 3090	GeForce RTX 4090	GeForce RTX 5090
GPU Codename	GA102	AD102	GB502
GPU Architecture	NVIDIA Ampere	NVIDIA Ada Lovelace	NVIDIA Blackwell
GPUs	7	11	11
TPCs	41	84	95
SMs	82	128	170
CUDA Cores / SM	128	128	128
CUDA Cores / GPU	10496	16384	21760
Tensor Cores / SM	4 (3rd Gen)	4 (4th Gen)	4 (5th Gen)
Tensor Cores / GPU	328 (3rd Gen)	512 (4th Gen)	680 (5th Gen)
GPU Boost Clock (MHz)	1895	2520	2407
RT Cores	82 (3rd Gen)	128 (3rd Gen)	170 (4th Gen)
RT Tensors	496	192	312
Tensor Buffer Memory Size and Type	24 GB GDDR6X	24 GB GDDR6X	32 GB GDDR7
Memory Interface	384-bit	384-bit	512-bit
Memory Clock (GHz)	19.5 Gbps	21 Gbps	28 Gbps
Memory Bandwidth	936 GB/sec	1008 GB/sec	1760 GB/sec
ROPs	112	176	176

Pixel Fill-rate (Biganels/sec)	189.8	443.5	423.6
Texture Units	328	512	680
Text Fill-rate (Biganels/sec)	555.96	1290.2	1636.76
L1 Data Cache/Shared Memory	10496 KB	16384 KB	21760 KB
L2 Cache Size	6144 KB	73728 KB	98304 KB
TDP (Total Graphics Power)	350 W	450 W	575 W
Manufacturing Process	Samsung 8 nm 8N NVIDIA Custom Process	TSMC 4nm 4N NVIDIA Custom Process	TSMC 4nm 4N NVIDIA Custom Process
PCI Express Interface	Gen 4	Gen 4	Gen 5

GPC: Graphics Processing Clusters
TPC: Texture Processing Cluster
SM: Streaming Multiprocessor
RT: Ray Tracing
ROP: Render Output Unit

RTX 5090



V100

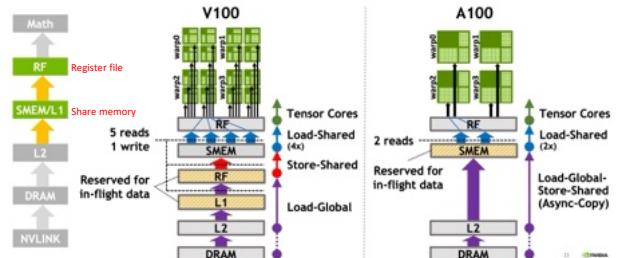


A100



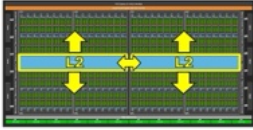
A100 data movement efficiency

3x SMEM/L1 bandwidth, 2x in-flight capacity



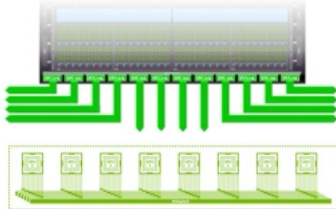
A100 data movement efficiency

Split L2 with hierarchical crossbar - 2.3x increase in bandwidth over V100, lower latency

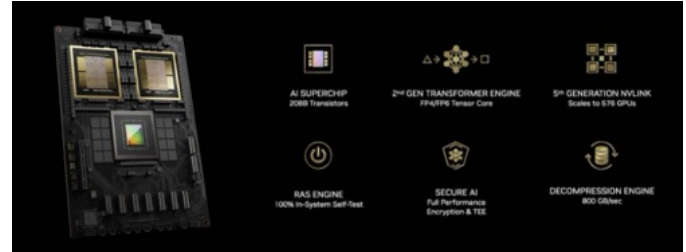


Third Generation NVLink

50 Gbit/sec per signal pair
12 links, 25 GB/s in/out, 600 GB/s total
2x vs. V100



GB100 \$40k, GB200: \$70k



RAS (Reliability, Availability, and Serviceability)

Datasetter

Die	GB100-1
Variant(s)	
Release date	Dec 2024
CUDA Cores	
Tensor Cores	
RT Cores	
Tensor Cores	
Streaming Multiprocessors	
Cache L1	
Cache L2	
Memory Interface	B102-G8
Die size	
Transistor count	104 Bn
Transistor density	503M
Package socket	
Products	B100 B200

Die	GB100-1	GB100-2	GB100-3	GB100-4	GB100-5
Variant(s)	GB100-100-A1 GB100-100-B1 GB100-100-C1	GB100-100-D1 GB100-100-E1 GB100-100-F1	GB100-100-G1 GB100-100-H1 GB100-100-I1	GB100-100-J1 GB100-100-K1 GB100-100-L1	GB100-100-M1 GB100-100-N1 GB100-100-O1
Release date	Jan 10, 2025	Jan 10, 2025	Jan 10, 2025	Jan 10, 2025	Jan 10, 2025
CUDA Cores	24,000	24,000	24,000	24,000	24,000
Tensor Cores	192	192	192	192	192
RT Cores	192	192	192	192	192
Tensor Cores	192	192	192	192	192
Cache L1	192	192	192	192	192
Cache L2	192	192	192	192	192
Memory Interface	B102-G8	B102-G8	B102-G8	B102-G8	B102-G8
Die size	104 Bn	104 Bn	104 Bn	104 Bn	104 Bn
Transistor count	104 Bn	104 Bn	104 Bn	104 Bn	104 Bn
Transistor density	503M	503M	503M	503M	503M
Products	B100 B200	B100 B200	B100 B200	B100 B200	B100 B200
Variant(s)	GB100-100-A1 GB100-100-B1 GB100-100-C1	GB100-100-D1 GB100-100-E1 GB100-100-F1	GB100-100-G1 GB100-100-H1 GB100-100-I1	GB100-100-J1 GB100-100-K1 GB100-100-L1	GB100-100-M1 GB100-100-N1 GB100-100-O1
Release date	Jan 10, 2025	Jan 10, 2025	Jan 10, 2025	Jan 10, 2025	Jan 10, 2025
CUDA Cores	24,000	24,000	24,000	24,000	24,000
Tensor Cores	192	192	192	192	192
RT Cores	192	192	192	192	192
Tensor Cores	192	192	192	192	192
Cache L1	192	192	192	192	192
Cache L2	192	192	192	192	192
Memory Interface	B102-G8	B102-G8	B102-G8	B102-G8	B102-G8
Die size	104 Bn	104 Bn	104 Bn	104 Bn	104 Bn
Transistor count	104 Bn	104 Bn	104 Bn	104 Bn	104 Bn
Transistor density	503M	503M	503M	503M	503M
Products	B100 B200	B100 B200	B100 B200	B100 B200	B100 B200
Variant(s)	GB100-100-A1 GB100-100-B1 GB100-100-C1	GB100-100-D1 GB100-100-E1 GB100-100-F1	GB100-100-G1 GB100-100-H1 GB100-100-I1	GB100-100-J1 GB100-100-K1 GB100-100-L1	GB100-100-M1 GB100-100-N1 GB100-100-O1
Release date	Jan 10, 2025	Jan 10, 2025	Jan 10, 2025	Jan 10, 2025	Jan 10, 2025
CUDA Cores	24,000	24,000	24,000	24,000	24,000
Tensor Cores	192	192	192	192	192
RT Cores	192	192	192	192	192
Tensor Cores	192	192	192	192	192
Cache L1	192	192	192	192	192
Cache L2	192	192	192	192	192
Memory Interface	B102-G8	B102-G8	B102-G8	B102-G8	B102-G8
Die size	104 Bn	104 Bn	104 Bn	104 Bn	104 Bn
Transistor count	104 Bn	104 Bn	104 Bn	104 Bn	104 Bn
Transistor density	503M	503M	503M	503M	503M
Products	B100 B200	B100 B200	B100 B200	B100 B200	B100 B200