

Christof Teuscher

ECE 410/510: Hardware for AI and ML

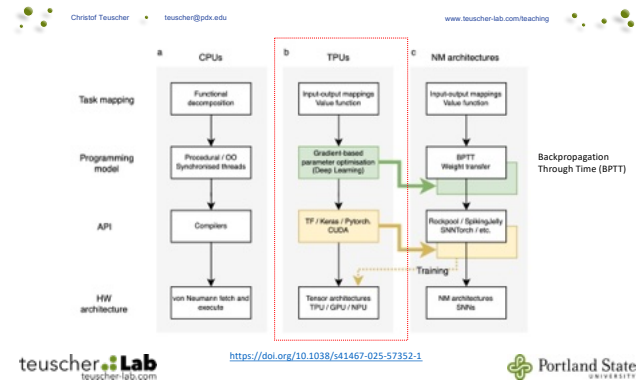
GPUs + CUDA + CNN/DNN

Portland State University
Department of Electrical and Computer Engineering (ECE)
www.teuscher-lab.com
teuscher@pdx.edu



teuscher Lab
teuscher-lab.com

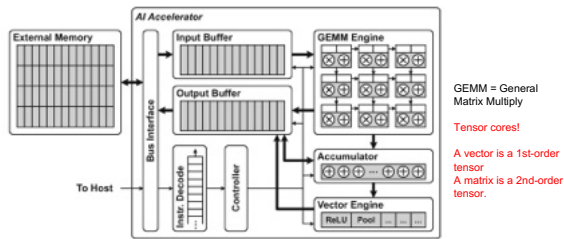
Portland State
UNIVERSITY



teuscher Lab
teuscher-lab.com

Portland State
UNIVERSITY

General AI accelerator architecture

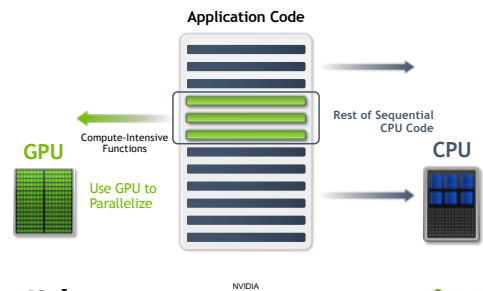


Mahra, Ashutosh; Cha, Jaekwang; Park, Hyunbin; Kim, Shih. Artificial Intelligence and Hardware Accelerators, 2025

teuscher Lab
teuscher-lab.com

Portland State
UNIVERSITY

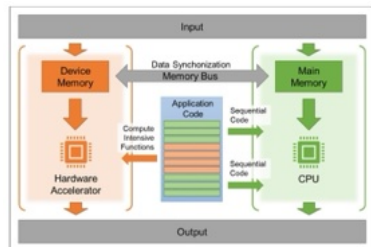
General co-processing mechanism



teuscher Lab
teuscher-lab.com

Portland State
UNIVERSITY

General co-processing mechanism

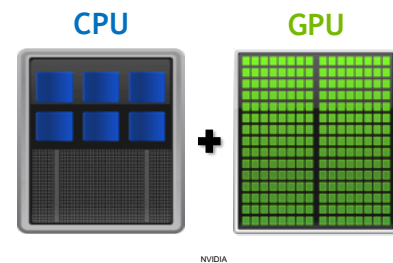


Mahra, Ashutosh; Cha, Jaekwang; Park, Hyunbin; Kim, Shih. Artificial Intelligence and Hardware Accelerators, 2025

teuscher Lab
teuscher-lab.com

Portland State
UNIVERSITY

Multi-core vs GPU

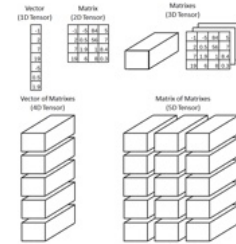


teuscher Lab
teuscher-lab.com

Portland State
UNIVERSITY

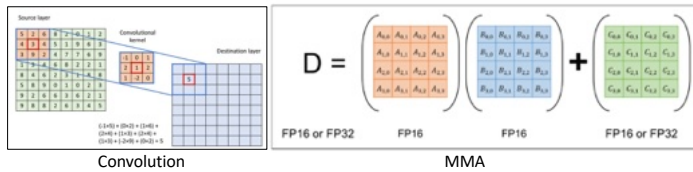
What is a tensor?

What is a tensor?



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.



Tensor cores (1)

Each Tensor Core operates on a 4×4 matrix and performs the following operation:

$$D = A \times B + C$$

where A, B, C, and D are 4×4 matrices (Figure 8). The matrix multiply inputs A and B are FP16 matrices, while the accumulation matrices C and D may be FP16 or FP32 matrices (see Figure 8).

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32 FP16 FP16 FP16 or FP32

Figure 8. Tensor Core 4×4 Matrix Multiply and Accumulate

Tensor cores (2)

Tensor Cores operate on FP16 input data with FP32 accumulation. The FP16 multiply results in a full precision product that is then accumulated using FP32 addition with the other intermediate products for a 4×4 matrix multiply (see Figure 9). In practice, Tensor Cores are used to perform much larger 2D or higher dimensional matrix operations, built up from these smaller elements.

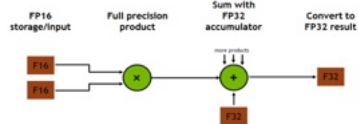


Figure 9. Mixed Precision Multiply and Accumulate in Tensor Core

Tensor cores (3)

Figure 10 shows the 4×4 matrix multiplication (using the two source 4×4 matrices outside the cube) requiring 64 operations (represented by the cube) to generate a 4×4 output matrix (shown below the cube). The Volta-based V100 accelerator with Tensor Cores can perform such calculations at 12x faster rate than Pascal-based Tesla P100.

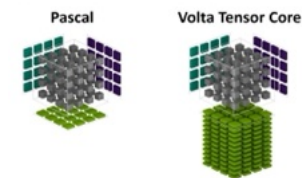


Figure 10. Pascal and Volta 4×4 Matrix Multiplication

Christof Teuscher

teuscher@psk.edu

www.teuscher-lab.com/teaching

Tensor cores (4)

New Turing Tensor Cores Provide Multi-Precision for AI Inference

<https://developer.nvidia.com/blog/nvidia-turing-architecture-in-depth>

teuscher:Lab

teuscher-lab.com

Portland State UNIVERSITY

Christof Teuscher

teuscher@psk.edu

www.teuscher-lab.com/teaching

CUDA vs tensor cores (1)

Metric	CUDA Cores	Tensor Cores
Primary Function	General-purpose Computing	Specialized for deep learning matrix computations
Arithmetic Operations	Executes parallel tasks efficiently	Optimized for matrix multiplication
Ideal Applications	Scientific research, machine learning, gaming, and more	Deep-learning models, neural network training
Raw Power	Generally, a higher number of cores per GPU	Fewer cores, but highly optimized for specific tasks
Performance Optimization	Efficient for complex algorithms in various fields	Accelerates matrix operations frequently used in AI
Use Case Suitability	A broad range of applications beyond deep learning	Specific tasks like training large-scale neural networks
Advantage in Deep Learning	Less efficient for matrix-heavy computations	Significant speedups and improved efficiency
Choice Based on Requirements	Suitable for diverse computing needs	Ideal for deep learning projects requiring extensive matrix computations
Future Prospects	Continuous improvements expected	Continuous improvements expected

<https://acecloud.ai/resources/blog/cuda-cores-vs-tensor-cores>

teuscher:Lab

teuscher-lab.com

Portland State UNIVERSITY

Christof Teuscher

teuscher@psk.edu

www.teuscher-lab.com/teaching

CUDA vs tensor cores (2)

CUDA Cores	Tensor Cores
CUDA cores deliver very high accuracy but at the cost of computing speed.	Tensor cores take the hit in accuracy but deliver accelerated computational speed.
CUDA cores are adequate for typical Machine Learning training. But enterprises predominantly prefer to use these for graphics-rich visual content and accelerated Ray-Tracing tasks.	Tensor cores are perfect for both low-end and high-end enterprise-grade AI development, data-rich ML training, and Deep Learning model development with multiple layers in neural networks.
CUDA cores are general-purpose cores that serve various tasks, such as graphics rendering, video editing, data analysis, machine learning, cryptographic operations, etc.	Tensor cores are specialized cores that perform high-end computations and allow mixed-precision training.
Enterprises that require small-scale/individual-level visual/video editing, game development tasks, data analytics operations, and indie research prefer CUDA-Cores GPUs.	Enterprises that demand robust, dedicated, and energy-efficient operation or research, such as managing data centers, massive AI training projects, High Performance Computing (HPC), etc., prefer Tensor core GPUs.

<https://acecloud.ai/resources/blog/cuda-cores-vs-tensor-cores>

teuscher:Lab

teuscher-lab.com

Portland State UNIVERSITY

Christof Teuscher

teuscher@psk.edu

www.teuscher-lab.com/teaching

CUDA vs tensor cores (2)

CUDA Cores	Tensor Cores
CUDA cores deliver very high accuracy but at the cost of computing speed.	Tensor cores take the hit in accuracy but deliver accelerated computational speed.
CUDA cores are adequate for typical Machine Learning training. But enterprises predominantly prefer to use these for graphics-rich visual content and accelerated Ray-Tracing tasks.	Tensor cores are perfect for both low-end and high-end enterprise-grade AI development, data-rich ML training, and Deep Learning model development with multiple layers in neural networks.
CUDA cores are general-purpose cores that serve various tasks, such as graphics rendering, video editing, data analysis, machine learning, cryptographic operations, etc.	Tensor cores are specialized cores that perform high-end computations and allow mixed-precision training.
Enterprises that require small-scale/individual-level visual/video editing, game development tasks, data analytics operations, and indie research prefer CUDA-Cores GPUs.	Enterprises that demand robust, dedicated, and energy-efficient operation or research, such as managing data centers, massive AI training projects, High Performance Computing (HPC), etc., prefer Tensor core GPUs.

<https://acecloud.ai/resources/blog/cuda-cores-vs-tensor-cores>

teuscher:Lab

teuscher-lab.com

Portland State UNIVERSITY

Christof Teuscher

teuscher@psk.edu

www.teuscher-lab.com/teaching

CUDA vs tensor cores (3)

CUDA Cores	Tensor Cores
CUDA Cores cannot reduce the use cycles needed for mixed multiplication and addition operations.	Tensor cores can reduce the number of cycles needed for mixed multiplication and addition operations.
CUDA Cores are slower via-via Tensor cores. Thus, Data scientists and ML Engineers do not prefer them for extensive ML training.	Tensor cores have significantly high computational power and are a perfect choice for Data Science, ML development, and algorithm training.
CUDA Cores perform a single operation in one clock cycle. 1x1 per GPU clock	Tensor cores can perform multiple operations in one clock cycle. (11 x 11) per GPU clock (11 x 11) per GPU clock
These are excellent for high-end graphics-based operations like Ray tracing, rendering physics engines, and other editing purposes that require accuracy.	These are excellent for handling Machine Learning operations that feed on massive data collection or extensive Artificial Neural Networks (ANN) layers.
Accelerated rendering of 3D graphics in Cloud-based games can leverage CUDA Cores.	Cloud services such as GPU-as-a-Service will prefer Tensor cores over CUDA cores.

<https://acecloud.ai/resources/blog/cuda-cores-vs-tensor-cores>

teuscher:Lab

teuscher-lab.com

Portland State UNIVERSITY

Christof Teuscher

teuscher@psk.edu

www.teuscher-lab.com/teaching

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	85
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 (2nd Gen)	128 (3rd Gen)	170 (4th Gen)
RT Trilangs	465	191	317.5
Frame Buffer Memory Size and Type	24 GB GDDR6X	24 GB GDDR6X	32 GB GDDR7
Memory Interface	384-bit	384-bit	512-bit
Memory Clock (MHz)	19.5 Gbps	21 Gbps	28 Gbps
Memory Bandwidth	936 GB/sec	1008 GB/sec	1792 GB/sec
ROPs	112	176	176

	Gen 4	Gen 4	Gen 5
Pixel Fill-rate (Bilions/sec)	189.8	443.5	423.6
Texture Units	328	512	680
Text Fill-rate (Bilions/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

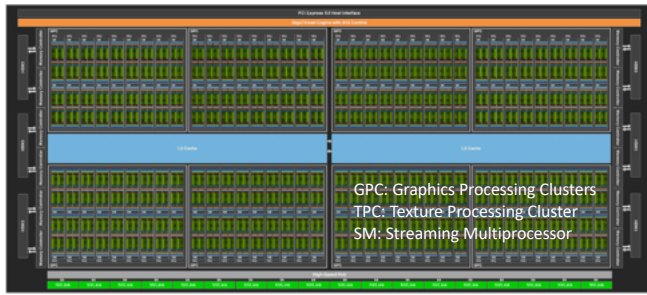
GPC: Graphics Processing Clusters
TPC: Texture Processing Cluster
SM: Streaming Multiprocessor
RT: Ray Tracing
ROP: Render Output Unit

teuscher:Lab

teuscher-lab.com

Portland State UNIVERSITY

RTX 5090



A100 TENSOR-CORE GPU

54 billion transistors in 7nm



V100

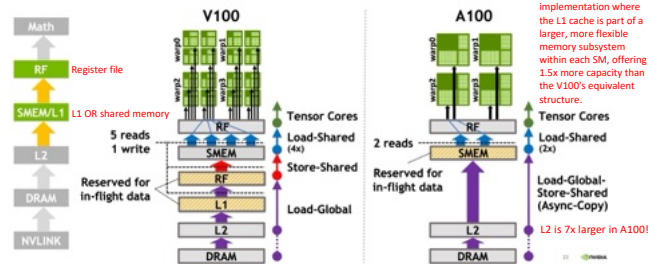


A100



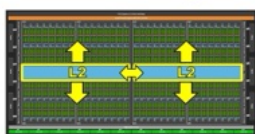
A100 data movement efficiency

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



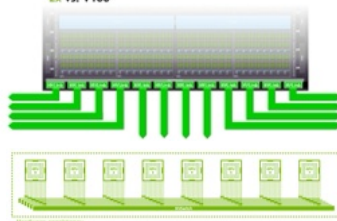
A100 data movement efficiency

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

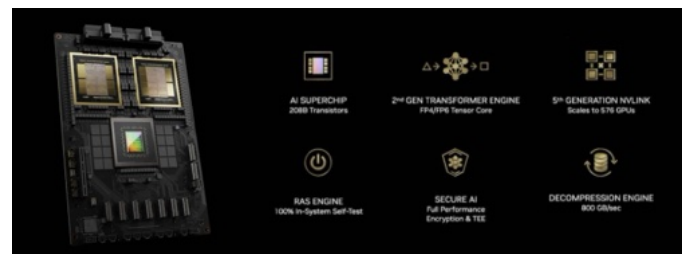


Third Generation NVLink

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)

Datacenter	
Die	GB102-2
Variant(s)	
Release date	Dec 2024
CUDA Cores	
Tensor Cores	
RT Cores	
Tensor Cores	
Streaming Multiprocessors	
Cache	
L1	
L2	
Memory Interface	B152-GB
Die size	
Transistor count	108 bn
Transistor density	
Package socket	SXM8
Products	B150 B150

Computer	
Die	GB102-2
Variant(s)	
Release date	Jan 10, 2025
CUDA Cores	96,768
Tensor Cores	192,512
RT Cores	192,512
Tensor Cores	192,512
Streaming Multiprocessors	192,512
Cache	
L1	
L2	
Memory Interface	B152-GB
Die size	
Transistor count	108 bn
Transistor density	
Package socket	SXM8
Products	B150 B150

CUDA

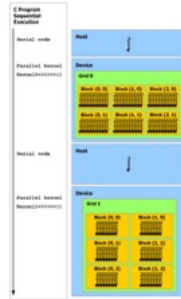
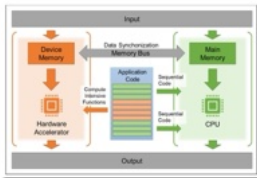
- **CUDA** (Compute Unified Device Architecture) is a parallel computing platform and programming model developed by NVIDIA.
- It enables developers to use NVIDIA graphics processing units (GPUs) for general-purpose computing, a practice known as **GPGPU** (General-Purpose computing on Graphics Processing Units).
- **CUDA programs to be written without knowing the exact hardware configuration in advance.**



- Good resource: <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

Programming and execution model

A GPU is built around an array of **Streaming Multiprocessors (SMs)**. A multithreaded program **partitioned into blocks of threads that execute independently from each other**, so that a GPU with more multiprocessors will automatically execute the program in less time than a GPU with fewer multiprocessors.



CUDA programming model assumes that the **CUDA threads execute on a physically separate device** that operates as a coprocessor to the host running the C++ program.

Serial code executes on the host while parallel code executes on the device.

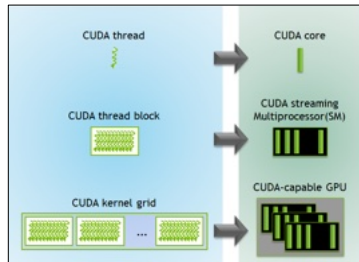
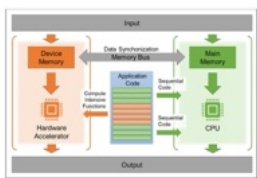
Some terminology

- A **warp** in CUDA is a fundamental execution unit consisting of 32 threads that execute in lockstep. All 32 threads in a warp execute the same instruction at the same time on different data (**SIMD**: Single Instruction, Data).
- A **CUDA kernel** is a function that runs on the GPU.
- A **block** in CUDA (also called a **thread block**) is a fundamental organizational unit of threads that execute a kernel function.
- The number of **threads per block** is limited (typically 1024 threads in modern GPUs), as they must reside on the same streaming multiprocessor (**SM**).
- **Thread blocks are composed of multiple warps**. For example, a block of 128 threads would consist of 4 warps.
- In CUDA programming, a **kernel grid** is the highest-level organization of threads when executing a CUDA kernel.

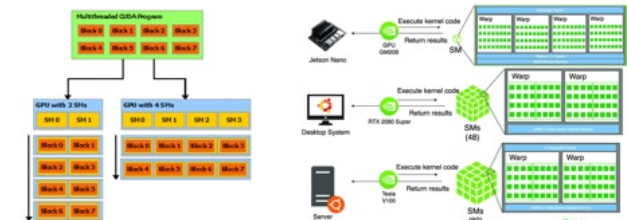


Programming and execution model

A GPU is built around an array of **Streaming Multiprocessors (SMs)**. A multithreaded program is **partitioned into blocks of threads that execute independently from each other**, so that a GPU with more multiprocessors will automatically execute the program in less time than a GPU with fewer multiprocessors.



Automatic scalability



A GPU is built around an array of **Streaming Multiprocessors (SMs)**. A multithreaded program is **partitioned into blocks of threads that execute independently from each other**, so that a GPU with more multiprocessors will automatically execute the program in less time than a GPU with fewer multiprocessors.

Kernels (4)

What happens if you ask for more blocks than are physically available on the GPU?

- The GPU doesn't reject your kernel launch. Instead, it schedules the execution of blocks over time. This is a key aspect of CUDA's execution model.
- The GPU hardware scheduler automatically queues the blocks and executes them as resources become available.
- Blocks waiting in the queue don't make progress until scheduled.
- The total execution time will increase as blocks must be processed sequentially.
- Memory for all blocks must still be allocated, potentially creating memory pressure.
- GPUs have a metric called "**occupancy**" which represents how many threads can be active simultaneously relative to the theoretical maximum. The scheduler tries to maximize occupancy based on resource availability.
- CUDA programs can be written without knowing the exact hardware configuration.



<https://developer.nvidia.com/blog/using-cuda-rtx-to-improve-cuda-programming-model/>

Another example

SAXPY stands for "Single-precision A Times Y"

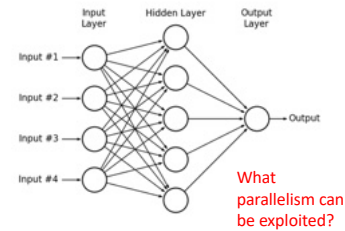
```
__global__ void saxpy(int N, float *x, float *y, float *z) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if (i < N) y[i] = x[i] + y[i];
}

int main() {
    int N = 1000000;
    float *x, *y, *z;
    x = (float *) malloc(N * sizeof(float));
    y = (float *) malloc(N * sizeof(float));
    z = (float *) malloc(N * sizeof(float));
    // Initialize x, y, and z
    for (int i = 0; i < N; i++) {
        x[i] = 1.0f;
        y[i] = 2.0f;
    }
    // Launch saxpy on 1024 threads
    cudaMemset(x, 0, N * sizeof(float));
    cudaMemset(y, 0, N * sizeof(float));
    cudaMemset(z, 0, N * sizeof(float));
    // Launch saxpy on 1024 threads
    cudaLaunchKernel(saxpy, 1024, N, x, y, z);
    // Wait for the kernel to finish
    cudaDeviceSynchronize();
    // Print the result
    for (int i = 0; i < N; i++) {
        printf("x[%d] = %f, y[%d] = %f, z[%d] = %f\n", i, x[i], y[i], z[i]);
    }
}
```

<https://developer.nvidia.com/blog/using-cuda-rtx-to-improve-cuda-programming-model/>

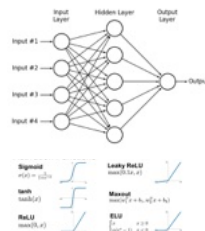
How to map a neural network on a GPU?

Multi-layer perceptron (MLP)



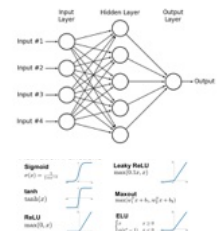
Multi-layer perceptron (MLP)

- Data parallelism:** You can process multiple training examples simultaneously. Since each example is processed independently through the forward pass, you can distribute batches across different computing units.
- Model parallelism:** The network itself can be partitioned, with different parts of the model assigned to different processors. This works particularly well for very large networks.
- Layer-wise parallelism:** Within each layer, the computations can be parallelized since neurons in the same layer operate independently.

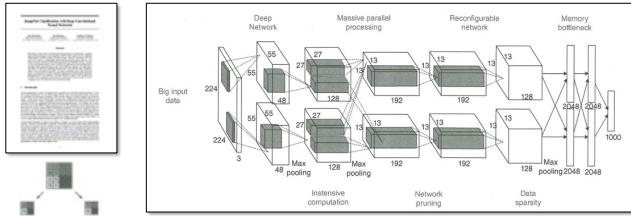


Multi-layer perceptron (MLP)

- Matrix multiplications:** The core operation in MLPs ($Wx + b$) is highly parallelizable. Each output neuron's computation is independent.
- Activation functions:** These are applied element-wise and can be computed in parallel across all neurons.
- Backpropagation:** During training, gradient computations can be parallelized both within and across layers.



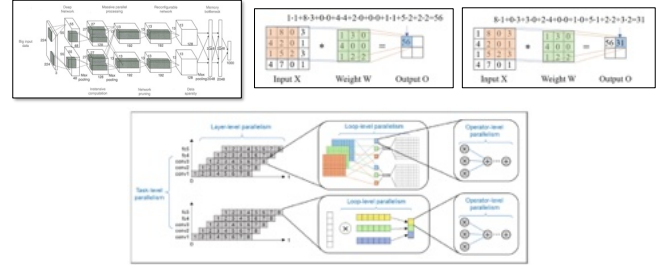
AlexNet: A classical deep convolutional neural network



AlexNet paper: https://papers.nips.cc/paper_files/paper/2012/file/c39862d3b9d6b76c8436e924a6845b-Paper.pdf

Liu & Law, Artificial Intelligence Hardware Design, 2021

Exploiting parallelism

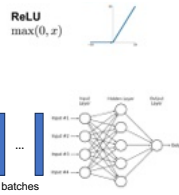


CUDA MLP (1)

```
#include <stdio.h>
#include <stdlib.h>
#include <cuda_runtime.h>
#include <time.h>

// Dimensions of our network
#define INPUT_SIZE 4
#define HIDDEN_SIZE 5
#define OUTPUT_SIZE 1
#define BATCH_SIZE 16 // Process 16 samples at once

// Activation function (ReLU)
__device__ float relu(float x) {
    return x > 0 ? x : 0;
}
```



CUDA C++ Language Extensions

7.1. Function Execution Space Specifiers

Function execution space specifiers denote whether a function executes on the host or on the device and whether it is callable from the host or from the device.

7.1.1. __global__

The `__global__` execution space specifier declares a function as being a kernel. Such a function is:

- Executed on the device,
- Callable from the host,
- Callable from the device for devices of compute capability 5.0 or higher (see [CUDA Dynamic Parallelism](#) for more detail).

Any call to a `__global__` function must have void return type, and cannot be a member of a class.

A call to a `__global__` function must specify its execution configuration as described in [Execution Configuration](#).

A call to a `__global__` function is asynchronous, meaning it returns before the device has completed its execution.

7.1.2. __device__

The `__device__` execution space specifier declares a function that is:

- Executed on the device,
- Callable from the device only.

The `__global__` and `__device__` execution space specifiers cannot be used together.

7.1.3. __host__

The `__host__` execution space specifier declares a function that is:

- Executed on the host,
- Callable from the host only.

<https://docs.nvidia.com/cuda/cuda-c-programming-guide/>

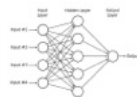
CUDA MLP (2)

```
// Forward pass kernel for first layer (input -> hidden)
__global__ void layer1_forward(float *input, float *weights, float *output, int batch_size) {
    int batch_idx = blockIdx.x * blockDim.x + threadIdx.x; // Global thread ID (x coordinate)
    int neuron_idx = blockIdx.y * blockDim.y + threadIdx.y; // Global y coordinate

    if (batch_idx < batch_size && neuron_idx < HIDDEN_SIZE) {
        float sum = bias[neuron_idx];
        for (int i = 0; i < INPUT_SIZE; i++) {
            sum += input[batch_idx * INPUT_SIZE + i] * weights[i * HIDDEN_SIZE + neuron_idx];
        }
        output[batch_idx * HIDDEN_SIZE + neuron_idx] = relu(sum);
    }
}

// Forward pass kernel for second layer (hidden -> output)
__global__ void layer2_forward(float *input, float *weights, float *output, int batch_size) {
    int batch_idx = blockIdx.x * blockDim.x + threadIdx.x;

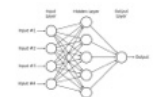
    if (batch_idx < batch_size) {
        float sum = bias[0];
        for (int i = 0; i < HIDDEN_SIZE; i++) {
            sum += input[batch_idx * HIDDEN_SIZE + i] * weights[i];
        }
        output[batch_idx] = sum; // No activation on output layer in this example
    }
}
```



CUDA MLP (3)

Why separate kernels for each layer?

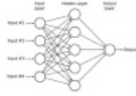
- Memory optimization:** Each layer has different memory access patterns and working set sizes. Separate kernels allow you to optimize memory usage specifically for each layer's requirements rather than trying to find a one-size-fits-all approach.
- Parallelism control:** Different layers may benefit from different thread configurations. For example, the first layer in our implementation handles 4 inputs to 5 neurons, while the output layer handles 5 inputs to 1 neuron - these naturally map to different parallelization strategies.
- Resource utilization:** GPU resources like shared memory, registers, and thread counts can be tailored specifically to each layer's computational needs, which improves overall efficiency.
- Kernel fusion limitations:** While fusing operations into a single kernel can reduce launch overhead, there are practical limits to how much computation can be packed into one kernel before it becomes inefficient. Separate kernels help manage this complexity.
- Synchronization points:** Having distinct kernels creates natural synchronization points between layers, ensuring all computations from one layer are complete before the next layer begins processing.
- Debugging and profiling:** With separate kernels, it's easier to identify performance bottlenecks or errors in specific parts of the network.
- Specialized optimizations:** Different mathematical operations may benefit from specialized implementations - convolutions, fully-connected layers, and activation functions all have different optimal implementations on GPUs.



CUDA MLP (4)

How much should be in a kernel?

- The optimal kernel size involves balancing several trade-offs:**
- Kernel launch overhead:** Each kernel launch incurs some overhead. Too many small kernels can result in performance degradation due to this overhead.
- Register pressure:** Larger kernels with more operations typically require more registers per thread, which can reduce occupancy (the number of threads that can run simultaneously on the GPU).
- Instruction cache:** GPUs have limited instruction cache. Very large kernels might not fit well in the cache, causing additional latency.
- Memory bandwidth utilization:** Kernels that load data, perform a small operation, and then store results can be memory-bound. Combining multiple operations in one kernel can improve arithmetic intensity and make better use of compute resources.
- Thread divergence:** If your kernel contains many conditional branches that cause different threads to follow different execution paths, performance will suffer.



As a general guideline:

- Good candidates for fusion:** Operations that work on the same data and have similar parallelization patterns (like a matrix multiplication followed by an element-wise activation function).
- Better kept separate:** Operations with fundamentally different access patterns or parallelization strategies.

CUDA MLP (5)

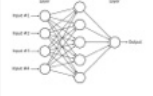
```
// Initialization weights with random values between -0.5 and 0.5
void initialize_weights(float* w_weights, int size) {
    for (int i = 0; i < size; i++) {
        w_weights[i] = ((float)rand() / RAND_MAX) - 0.5f;
    }
}

int main() {
    // Seed random number generator
    srand(time(0));

    // Host memory allocation
    float h_input[BATCH_SIZE * INPUT_SIZE];
    float h_output[BATCH_SIZE * HIDDEN_SIZE];
    float h_weights1[INPUT_SIZE * HIDDEN_SIZE];
    float h_weights2[HIDDEN_SIZE * OUTPUT_SIZE];
    float h_bias1[HIDDEN_SIZE];
    float h_bias2[OUTPUT_SIZE];

    // Initialize weights and biases
    initialize_weights(h_weights1, INPUT_SIZE * HIDDEN_SIZE);
    initialize_weights(h_weights2, HIDDEN_SIZE * OUTPUT_SIZE);
    initialize_weights(h_bias1, HIDDEN_SIZE);
    initialize_weights(h_bias2, OUTPUT_SIZE);

    // Initialize some example input data
    for (int i = 0; i < BATCH_SIZE * INPUT_SIZE; i++) {
        h_input[i] = ((float)rand() / RAND_MAX);
    }
}
```



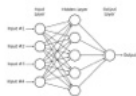
CUDA MLP (6)

```
// Host memory allocation
float h_input[BATCH_SIZE * INPUT_SIZE];
float h_output[BATCH_SIZE * HIDDEN_SIZE];
float h_weights1[INPUT_SIZE * HIDDEN_SIZE];
float h_weights2[HIDDEN_SIZE * OUTPUT_SIZE];
float h_bias1[HIDDEN_SIZE];
float h_bias2[OUTPUT_SIZE];

// Device memory allocation
float d_input, d_output;
float d_weights1, d_weights2, d_bias1, d_bias2;

cudaMalloc(&h_input, BATCH_SIZE * INPUT_SIZE * sizeof(float));
cudaMalloc(&h_output, BATCH_SIZE * HIDDEN_SIZE * sizeof(float));
cudaMalloc(&h_weights1, INPUT_SIZE * HIDDEN_SIZE * sizeof(float));
cudaMalloc(&h_weights2, HIDDEN_SIZE * OUTPUT_SIZE * sizeof(float));
cudaMalloc(&h_bias1, HIDDEN_SIZE * sizeof(float));
cudaMalloc(&h_bias2, OUTPUT_SIZE * sizeof(float));

// Copy data from host to device
cudaMemcpy(d_input, h_input, BATCH_SIZE * INPUT_SIZE * sizeof(float), cudaMemcpyHostToDevice);
cudaMemcpy(d_weights1, h_weights1, INPUT_SIZE * HIDDEN_SIZE * sizeof(float), cudaMemcpyHostToDevice);
cudaMemcpy(d_weights2, h_weights2, HIDDEN_SIZE * OUTPUT_SIZE * sizeof(float), cudaMemcpyHostToDevice);
cudaMemcpy(d_bias1, h_bias1, HIDDEN_SIZE * sizeof(float), cudaMemcpyHostToDevice);
cudaMemcpy(d_bias2, h_bias2, OUTPUT_SIZE * sizeof(float), cudaMemcpyHostToDevice);
```



CUDA MLP (7)

How will we use threads for this MLP?



A GPU is built around an array of Streaming Multiprocessors (SMs). A multithreaded program is partitioned into blocks of threads that execute independently from each other.

CUDA MLP (8)

```
// Set up execution configuration for first layer
dim3 threadsPerBlock(16, 16, 1); // 16 threads per block
dim3 blocksPerGrid(BATCH_SIZE * INPUT_SIZE / threadsPerBlock.x,
                  HIDDEN_SIZE / threadsPerBlock.y, 1);

// Set up execution configuration for second layer
dim3 threadsPerBlock2(16, 16, 1); // 16 threads per block
dim3 blocksPerGrid2(BATCH_SIZE * HIDDEN_SIZE / threadsPerBlock2.x,
                   OUTPUT_SIZE / threadsPerBlock2.y, 1);

// Launch kernels
layer1_forward=blocksPerGrid, threadsPerBlock==d_input, d_weights1, d_bias1, d_hidden, BATCH_SIZE);
layer2_forward=blocksPerGrid2, threadsPerBlock2==d_output, d_weights2, d_bias2, d_hidden, BATCH_SIZE);

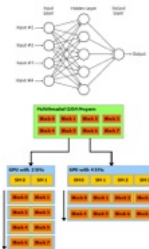
// Check for errors
cudaError_t err = cudaGetLastError();
if (err != cudaSuccess) {
    printf("Error: %s\n", cudaGetLastError().message);
    return -1;
}

// Copy results back to host
cudaMemcpy(h_output, d_output, BATCH_SIZE * OUTPUT_SIZE * sizeof(float), cudaMemcpyDeviceToHost);

// Print some outputs
printf("Sample outputs from MLP:\n");
for (int i = 0; i < BATCH_SIZE * OUTPUT_SIZE; i++) {
    printf("Sample %d: %f\n", i, h_output[i]);
}
```

Why?

kernel <= 0.75 * 1024
The kernel launches with a grid of M thread blocks.
Each thread block has 1 parallel threads.



A GPU is built around an array of Streaming Multiprocessors (SMs). A multithreaded program is partitioned into blocks of threads that execute independently from each other.

CUDA MLP (9)

```
// Free GPU memory
cudaFree(d_input);
cudaFree(d_hidden);
cudaFree(d_output);
cudaFree(d_weights1);
cudaFree(d_weights2);
cudaFree(d_bias1);
cudaFree(d_bias2);
```

