

Compute Unified Device Architecture (CUDA)

What is CUDA used for?

↳ Graphics and rate tracing

↳ Deep learning

↳ and more

* The point of writing GPU Kernel is to run something faster

Chapter 1 Deep Learning Ecosystem

* Research: PyTorch, TensorFlow, MLX (apple devices)

* Production: Training and Inference

↳ vLLM (inference only)

↳ Tensor RT

↳ Triton (OpenAI)

* Low Level: CUDA, ROCm

* Inference For Edge Computing & Embedded Systems:

↳ Edge Computing: low-latency and highly efficient local computing in the context of real-world distributed systems

↳ PyTorch Mobile, TensorFlow Lite

* Cloud Computing:

↳ AWS, Google Cloud, Azure, VASTAI

* Compilers:

↳ XLA, LLVM, NVCC

Quick C++ Review

Basic Structure

Using namespace std; // standard library

int main()
{

// code here

data-type example; // variable declaration

cin >> example; // input

cout << example << endl; // output (print)

if (condition) // if and nested if

if (condition 2){

}

3

switch (expression){

case value:

// code here if matches value

break;

case ...

default:

break;

3

for (initialization; test; update){

3

dataType arrge-name [size]; // multidimensional datatype array-name [size]... [size]

vector <data-type> vector-name // resize automatically, #include <vector>

push-back() insert elem @ end, pop-back() removes elem from end, clear() remove all elem, empty() checks if empty,

at(i) access elem at index i, front() access first elem, back() last elem, erase() removes specific elem at position

x/

return 0;

3

Quick C++ Review Continued . . .

References: points to an object of a given class, letting you access the value of an object

```
int var = 12;
```

```
int & ref = var;
```

ref is now a reference to var

Pointers: stores memory address of another var

```
int i = 3;
```

```
int *ptr = &i; // pointer to i (stores address of i)
```

Functions:

```
return-type func-name (param); // function declaration
```

```
return-type func-name (param) {
```

// code here for func definition)

}

String Functions:

- * substring: string substr (size_t first, char copy, size_t length);

- * append() adds to end of string

- * compare()

- * empty() checks if empty

OOP

- * Class & objects:

↓
blueprint
for objects
behavior

↓
instance or
var of class

- * Encapsulation: wrap up data & methods together within a single entity (classes)

- * Abstraction: shows only necessary details & hides internal

- * Inheritance: deriving properties of Parent class to Child class

- * Polymorphism: different functionalities to functions with same name (overload, override)

- * File handling: open(), getline(), << to write file after opening

GPUs Background Info

* CPU: Central Processing Unit → host executes functions

- ↳ general purpose
- ↳ high clock speed
- ↳ few cores
- ↳ high cache
- ↳ low latency - (Complete it fast) (Can't do as much comparatively operations per second)
- ↳ low throughput

- minimize time of one
- metric: latency in seconds

* Graphics Processing Unit → device executes kernels

- ↳ specialized
- ↳ low clock speed
- ↳ many cores
- ↳ low caches
- ↳ high latency
- ↳ high throughput

↳ maximize throughput

↳ metric: throughput in tasks per second

• Tensor Processing Unit (TPU)

↳ Specialized GPU for deep learning algorithms (matrix multiplication)

• FPGA: Field Programmable Gate Array

↳ specialized hardware that can be reconfigured to perform specific tasks

↳ ^{very} low latency

↳ very high throughput

↳ very high power consumption

↳ very high cost

Typical CUDA Program

1. CPU allocates CPU memory
2. CPU copies data to GPU
3. CPU launches Kernel on GPU (processing is done here)
4. CPU copies results from GPU back to CPU to do something useful with it

* GPU Kernel: → smallest unit of execution → mini program that runs independently but shares resources

↳ a function that runs on the GPU

↳ runs many threads in parallel (1000s at once) (each thread executes same func. but on different data)

↳ - global — void add(int *a, int *b, int *c) ↳ name and parameters of function
int index = blockIdx.x * blockDim.x + threadIdx.x; ↳ computes global thread ID, so each thread knows which element it should work on
keyword tells CUDA this func. runs on GPU and can be called from CPU

$c[\text{index}] = a[\text{index}] + b[\text{index}];$ ↳ num of threads per block

↳ Vector addition in parallel
of 1000s of GPU threads

3

↳ Block → a group of threads that can cooperate using shared memory

↳ Grid → full set of all blocks launched for one kernel call

↳ Launch: call kernel from CPU side

ex: `int N=1000;
add<<(N+255)/256, 256>>(d_a, d_b, d_c);`

launches: $(N+255)/256$ blocks

with 256 threads