

CUDA Programming



Outline

- What is CUDA?
- When to use CUDA?
- Principles of CUDA
- How CUDA works?



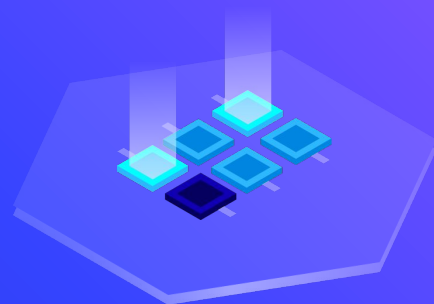
The background features a dense bundle of fiber optic cables that glow with a vibrant blue light. The cables are arranged in a fan-like pattern, creating a sense of depth and movement. Overlaid on this is a complex digital network of white lines and nodes, resembling a data flow or circuit board. Several small, semi-transparent boxes containing binary code (001, 011, 010) are scattered throughout the network, adding a technical and futuristic feel to the overall composition.

What is Cuda ?

What is Cuda



- 📌 CUDA = Cuda Architecture + Cuda Programming Model
- 📌 Developed by NVIDIA
- 📌 Parallel computing on graphics processing units (GPUs)





Programming language
based on C++



CUDA



Massively parallel
hardware designed

Software development kit

“ CUDA language is based on C++ with a few additional keywords and concepts, which makes it fairly easy for non-GPU programmers to pick up

```
void c_hello(){  
    printf("Hello World!\n");  
}  
  
int main() {  
    c_hello();  
    return 0;  
}
```

```
__global__ void cuda_hello(){  
    printf("Hello World from GPU!\n");  
}  
  
int main() {  
    cuda_hello<<<1,1>>>();  
    return 0;  
}
```

History of CUDA

Compute Unified Device Architecture

1989

Pixel Machine



2007

CUDA

Compute Unified Device
Architecture
By NVIDIA



GPU was invented

purely fixed-function devices
GeForce 256

1999



OpenCL

CUDA competitor
Apple and the Khronos Group
Not limited to Intel/AMD
CPUs with Nvidia GPUs

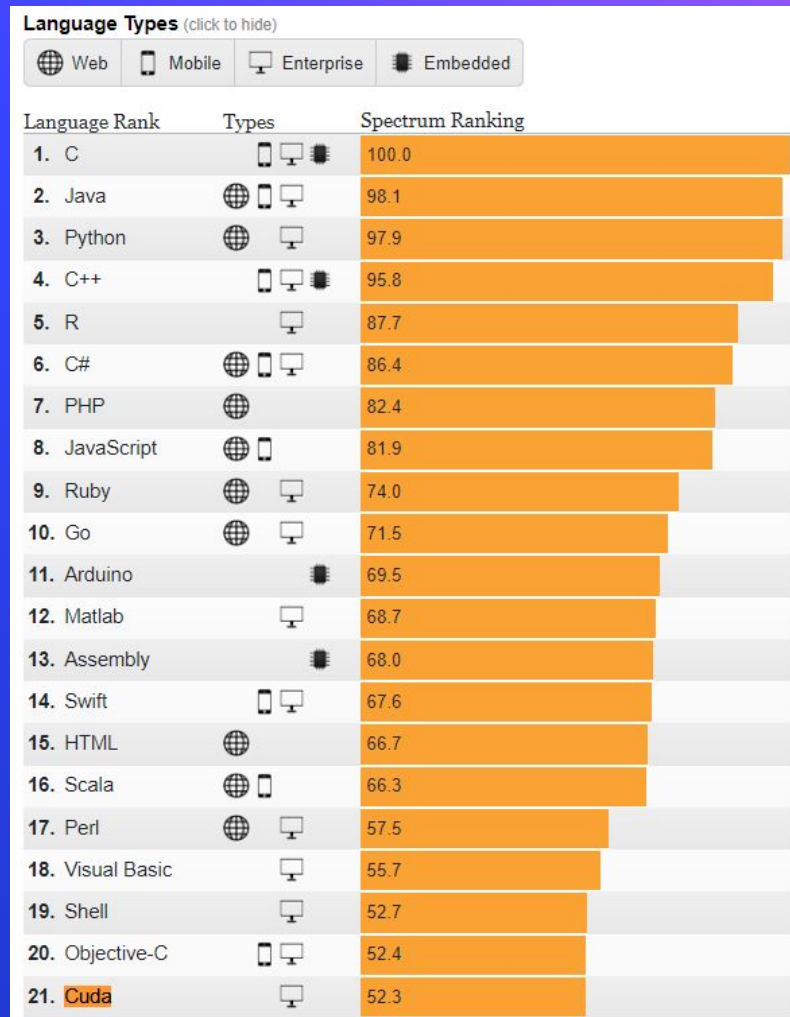
2009



Popularity

y

The Top Programming Languages 2016





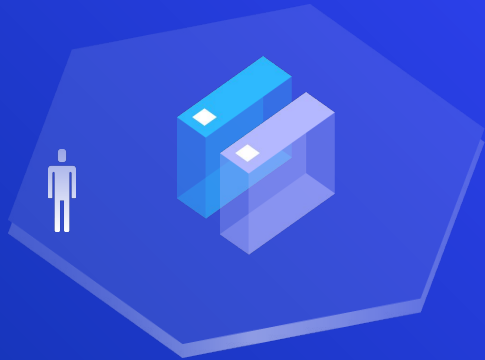
When to use Cuda

?

When to use Cuda ?

⬡ Lots of data

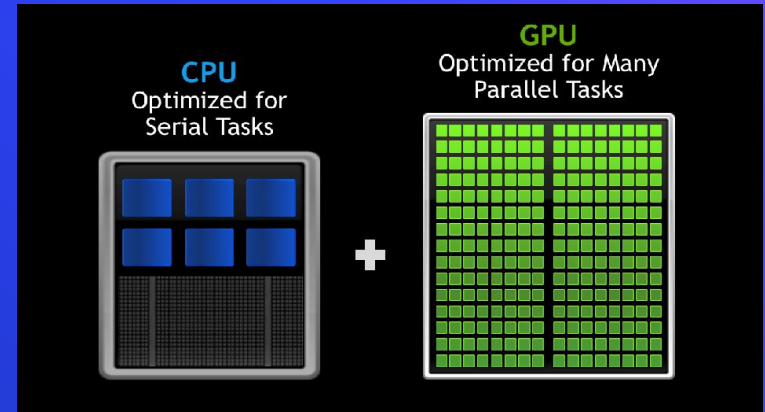
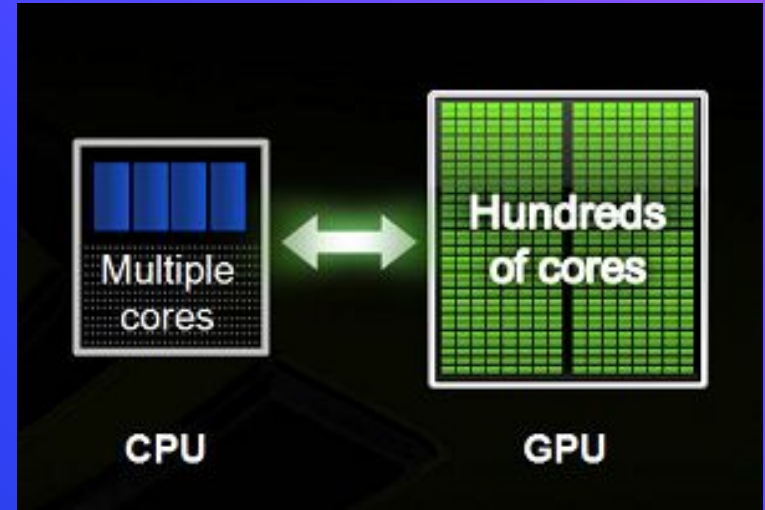
⬡ Lots of computations



Benefits of Cuda

?

- Harnesses the power of the GPU by using parallel processing;
- Running thousands of simultaneous reads instead of single, dual, or quad reads on the CPU.



Benefits of Cuda

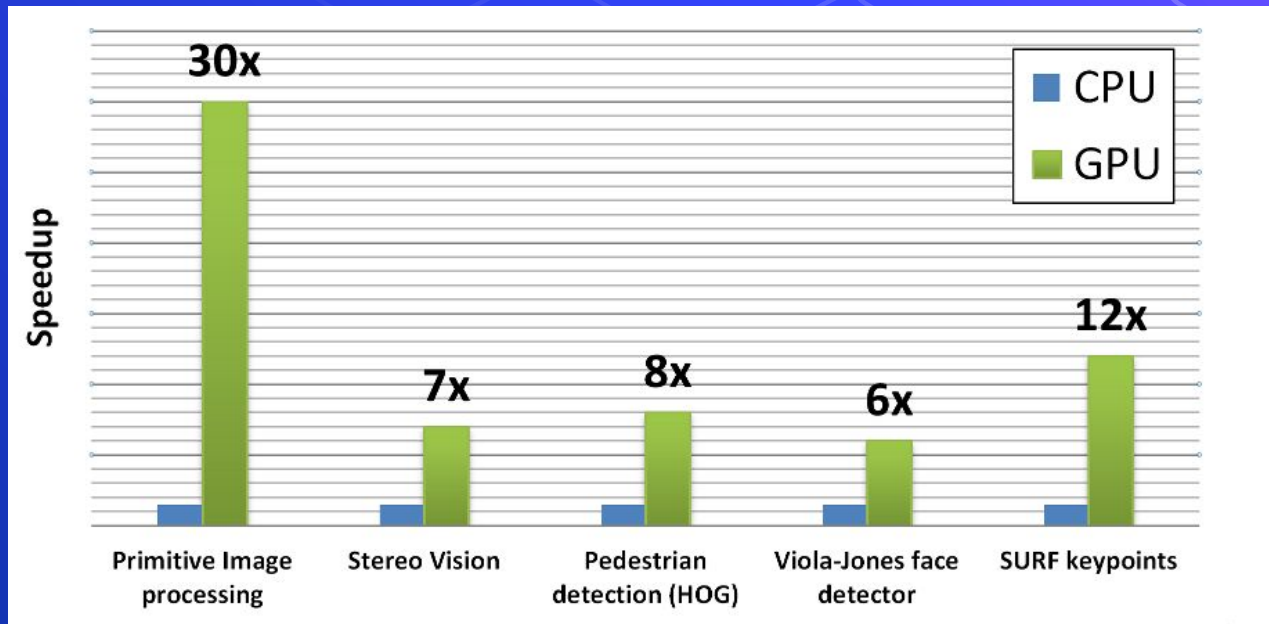
?

- ⬡ C/C++ is widely used, easy to learn how to program for CUDA.
- ⬡ Most of graphics cards of NVIDIA support CUDA.



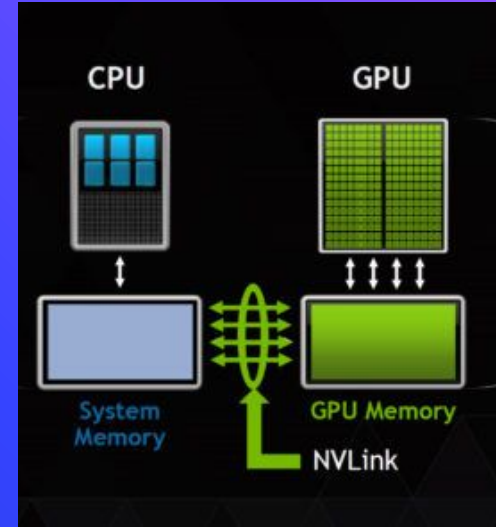
5x to 200x

Huge increase in processing power over conventional CPU processing. Early reports suggest speed increases of 5x to 200x over CPU processing speed.



Limitations of CUDA

Over traditional general purpose computation on GPUs



⬡ Latency between the CPU and GPU

Copying between host and device memory may incur a performance hit due to system bus bandwidth

⬡ A single process must run spread across multiple memory spaces

Threads should be running in groups of at least 32 for best performance

⬡ CUDA-enabled GPUs are only available from NVIDIA

Unlike OpenCL

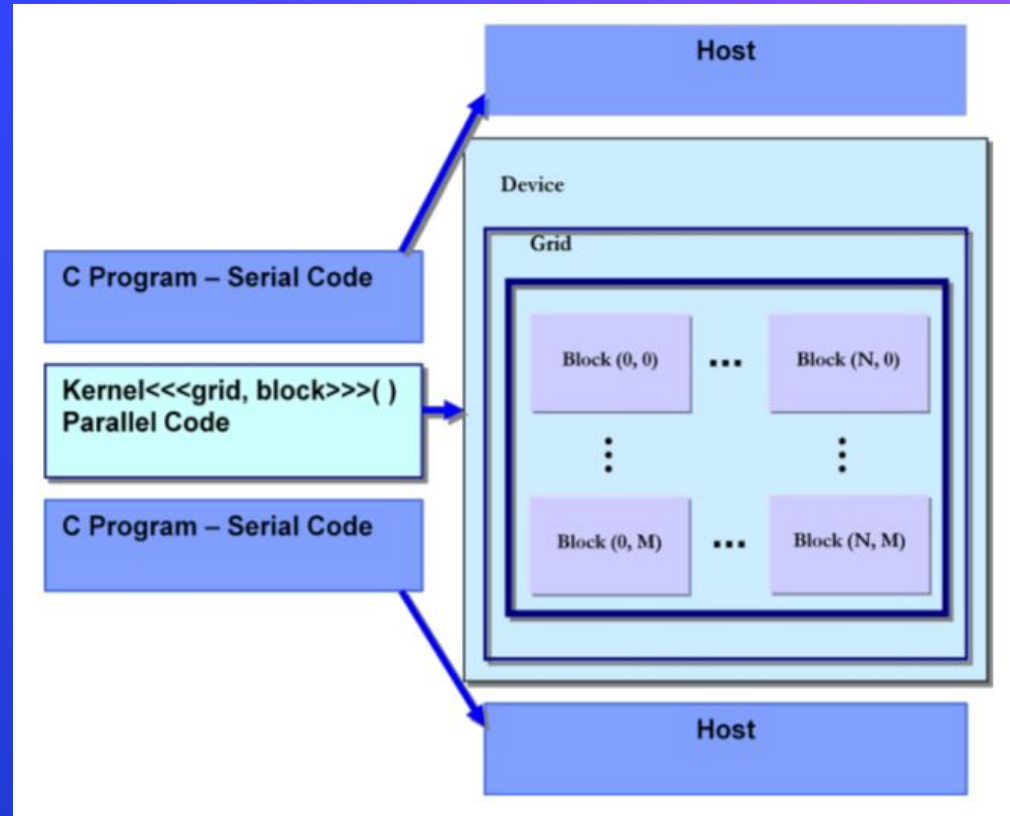


Principles of CUDA

Parallel structure

Paradigm

- ⬡ Parallel programming
- ⬡ Combination of serial and parallel executions



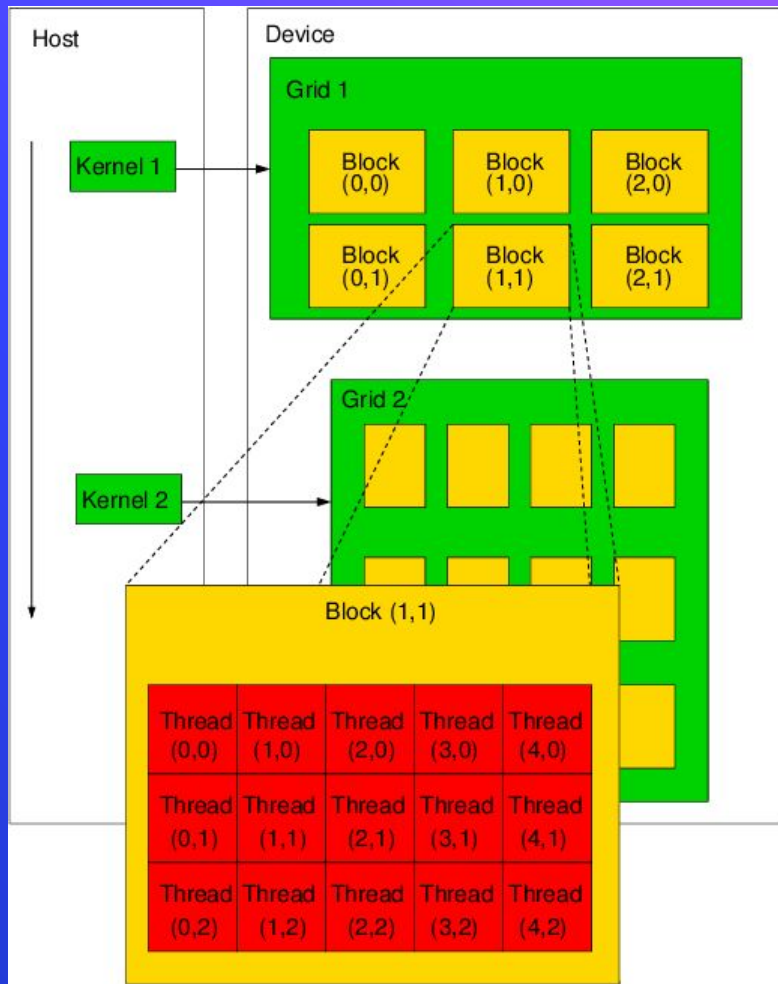
Programming Model

- Device = GPU
- Host = CPU
- Kernel = Functions run on Device



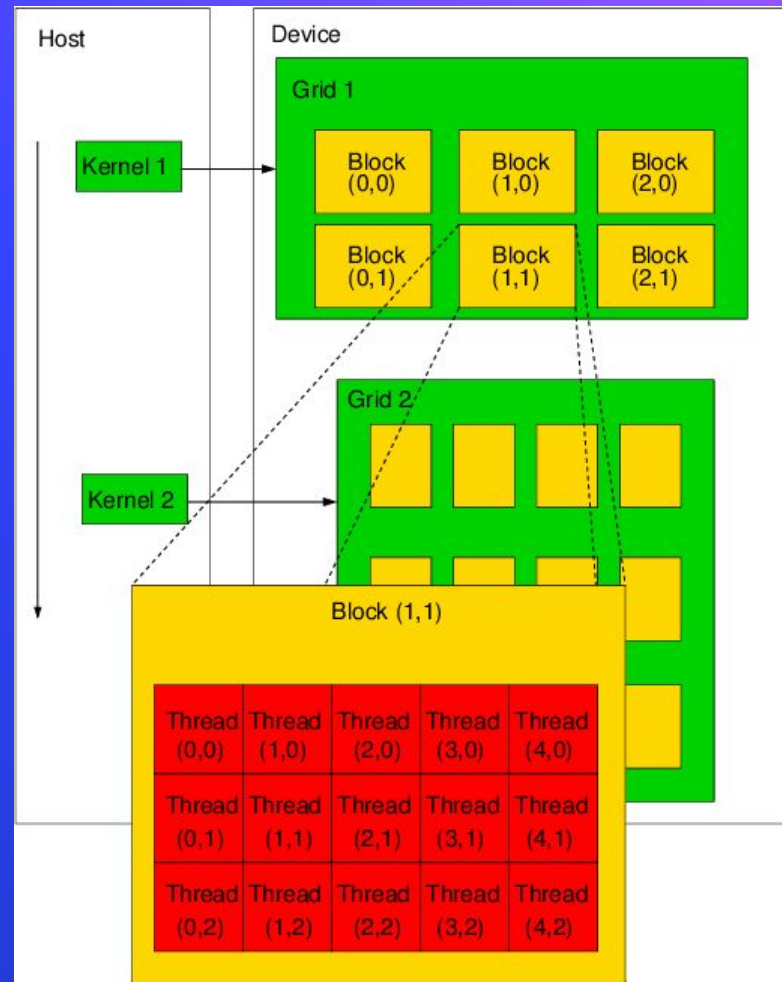
More about CUDA kernel

specific functions in CUDA, A kernel can be a function or a full program invoked by the CPU. It is executed N number of times in parallel on GPU by using N number of threads



Programming Model

- 512, 1024 or 2048 threads in one block
- All blocks define a grid
- All block execute same program (kernel)
- Blocks are independent
- Only one kernel at a time



Programming Model

Technical specifications	Compute capability (version)				
	1.0	1.1	1.2	1.3	2.x
Maximum x- or y- dimensions of a grid of thread blocks	65535				
Maximum number of threads per block	512			1024	

In our GPU using for experiment

Max of block = 65535

Max of threads each block = 1024 (Compute Capability: 3.5)

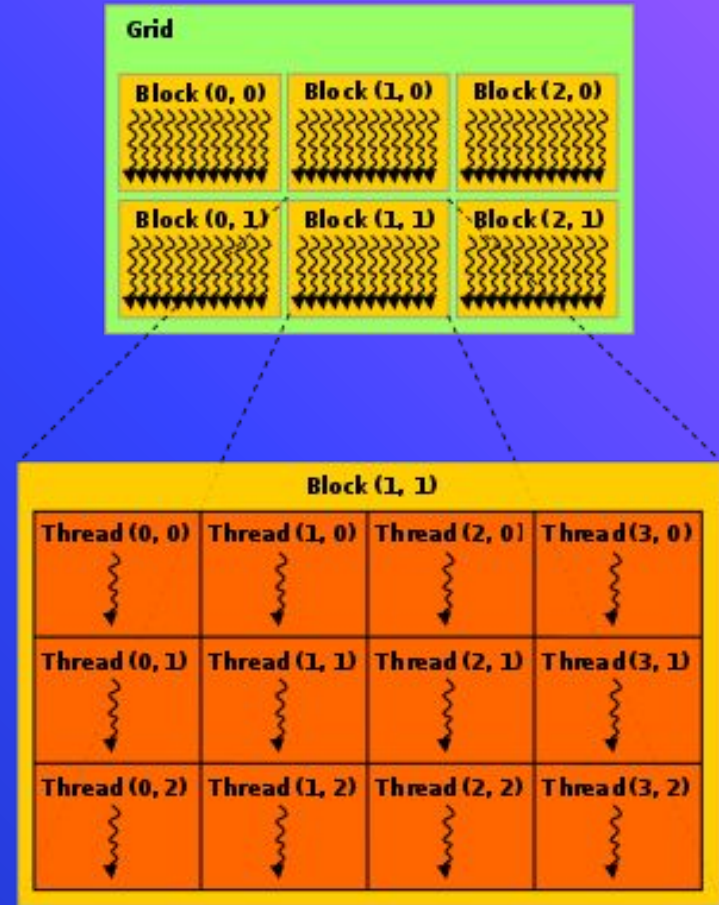
CUDA Thread

❖ Thread Cooperation

- Share memory > powerful feature of CUDA
- All threads run same code

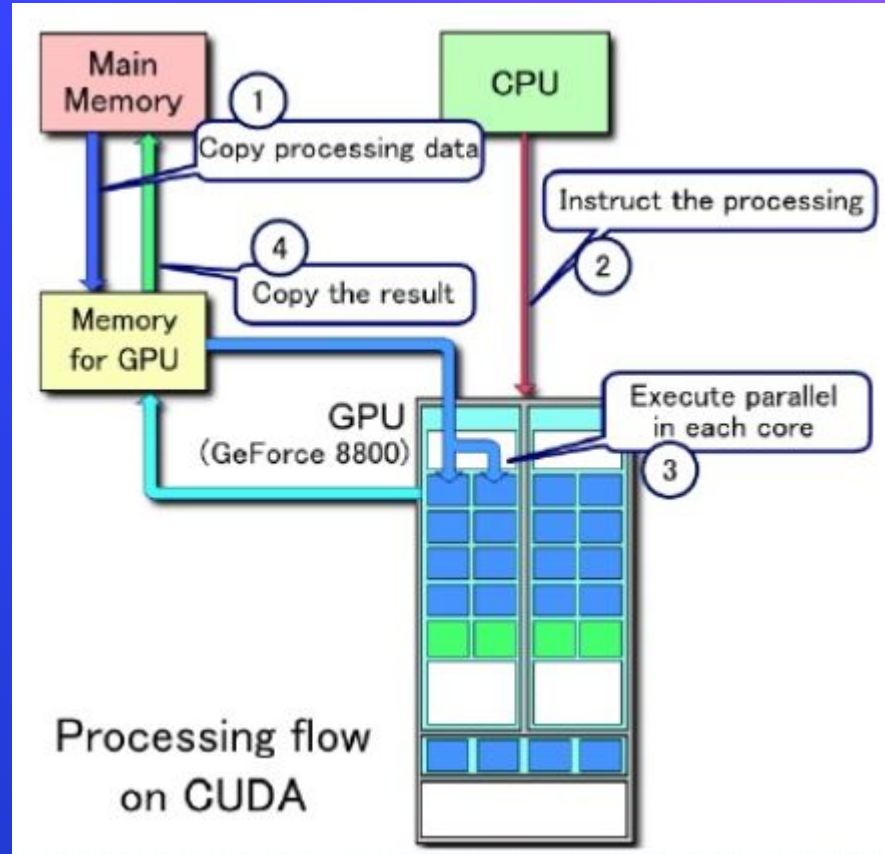
❖ CUDA threads vs CPU threads

- CUDA thread is lightweight
- CUDA use 1000s threads, CPU single, dual, or quad



Processing Flow

1. Copy data from Main mem to GPU mem
2. CPU instructs the process to GPU
3. GPU execute parallel in each core
4. Copy result GPU mem to Main mem



Processing Flow

- Copy data from Main memory to GPU mem
- GPU execute parallel
- Copy result GPU mem to Main mem

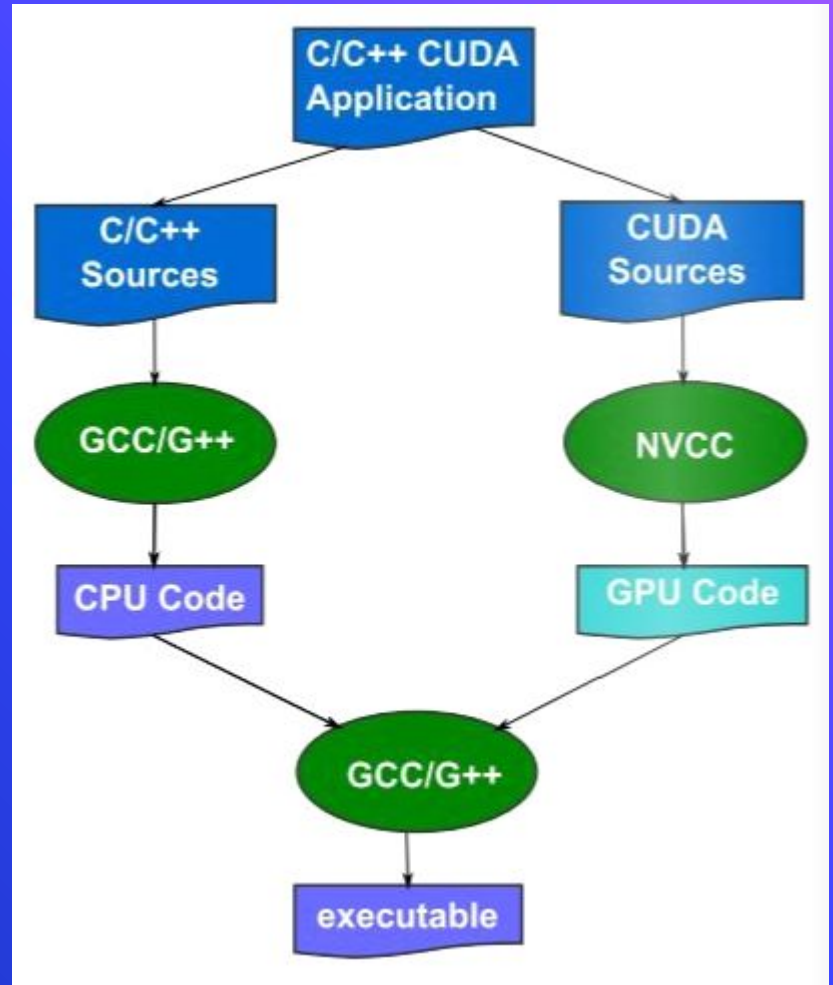
```
cudaMemcpy( d_a, a, SIZE*sizeof(int), cudaMemcpyHostToDevice );
cudaMemcpy( d_b, b, SIZE*sizeof(int), cudaMemcpyHostToDevice );
cudaMemcpy( d_c, c, SIZE*sizeof(int), cudaMemcpyHostToDevice );

VectorAdd<<< 1, SIZE >>>(d_a, d_b, d_c, SIZE);

cudaMemcpy( c, d_c, SIZE*sizeof(int), cudaMemcpyDeviceToHost );
```

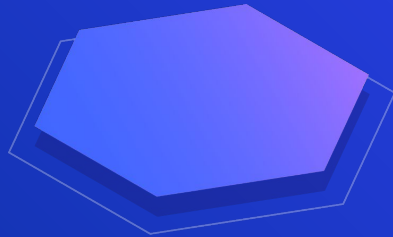
CUDA compiler

- ⬡ GCC/G++ Compiler for Host code
- ⬡ NVCC Compiler for device code
- ⬡ GCC/G++ for linking



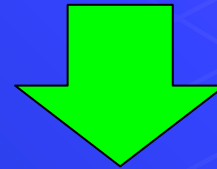
Syntax

Identifier



RULES FOR NAMING IDENTIFIER

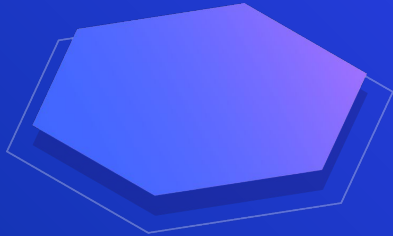
- Case sensitive
- Cannot start with a digit
- Underscore can be used as first character
- Other special characters are not allowed
- Cannot use keywords as identifier.



Followsym: **[a-zA-Z0-9_ \$]**

Identifier: **[a-zA-Z]{followsym}* | {_followsym}+**

Specification



Function Qualifiers

Function Call

CUDA Built-in device variable

Variable Qualifiers

Function Qualifiers

Example

Function Qualifiers

- __global__ - Invoked from within host (CPU) code, cannot be called from device (GPU) code - must return void
- __device__ - called from other GPU functions, cannot be called from host (CPU) code
- __host__ - can only be executed by CPU, called from host
- __host__ and __device__ qualifiers can be combined
- Example use overriding operators
- Example will generate both CPU and GPU code

```
#include <stdio.h>

__global__
void add(int *a, int *b, int *c, int n) {

    int i = threadIdx.x;

    __shared__ bool sharedArr[100];

    for (int i = 0; i < n; i++)
        c[i] = a[i] + b[i];
}
```

“ CUDA language is based on C++ with a few additional keywords and concepts, which makes it fairly easy for non-GPU programmers to pick up

```
void c_hello(){  
    printf("Hello World!\n");  
}  
  
int main() {  
    c_hello();  
    return 0;  
}
```

```
__global__ void cuda_hello(){  
    printf("Hello World from GPU!\n");  
}  
  
int main() {  
    cuda_hello<<<1,1>>>();  
    return 0;  
}
```

Function Qualifiers

- `__global__` : { invoked from within CPU code
can not be called from GPU code
must return void
- `__device__` : { called from other GPU functions
can not be called from CPU code
- `__host__` : { can only be executed by CPU
Called from host

Function Qualifiers

- `__global__` : Invoked from within host (CPU) code,
cannot be called from device (GPU) code
- must return void
- `__device__` : called from other GPU functions,
cannot be called from host (CPU) code
- `__host__` : can only be executed by CPU, called from
host
- `__host__` and `__device__` qualifiers can be combined
- Example use controlling memory
- Compiler will generate both CPU and GPU code

Function Call

Example

Function Qualifiers

- `__global__`: Invoked from within host (CPU) code.
- cannot be called from device (GPU) code
- must return void
- `__device__`: called from other GPU functions,
- cannot be called from host (CPU) code
- `__host__`: can only be executed by CPU, called from host
- `__host__` and `__device__` qualifiers can be combined
- Example use: controlling memory
- Example will generate both CPU and GPU code

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = i;  
    b[i] = i;  
}
```

```
add<<<1,1>>>(a,b,c,SIZE);
```

```
cudaDeviceSynchronize();
```

```
cudaFree(a);
```

```
cudaFree(b);
```

```
cudaFree(c);
```

Function Call

kernel_function<<<M,N>>>(list of arguments);

M is number of block

N is number of threads each block (size of block)

*M*N is number of times that function is executed.*

CUDA Built-in device variable

Example

```
__global__  
void add(int *a, int *b, int *c, int n) {  
    int i = threadIdx.x;  
    __shared__ bool sharedArr[100];  
    for (int i = 0; i < n; i++)  
        c[i] = a[i] + b[i];  
}  
  
int main() {  
    int *a, *b, *c;
```

Function Qualifiers

- __global__ - Invoked from within host (CPU) code.
- __device__ - Cannot be called from device (GPU) code.
- __host__ - Cannot be called from host (CPU) code.
- __host__ __device__ - Can only be executed by CPU, called from both.
- __host__ __device__ __force__ - Qualifiers can be combined.
- Example will generate both CPU and GPU code.

CUDA Built-in device variable

- ★ `dim3 blockDim` : Dimensions of the grid in blocks
- ★ `dim3 blockIdx` : Dimensions of the block in threads
- ★ `dim3 threadIdx` : Block index within the grid
- ★ `dim3 threadIdx` : Thread index within the block



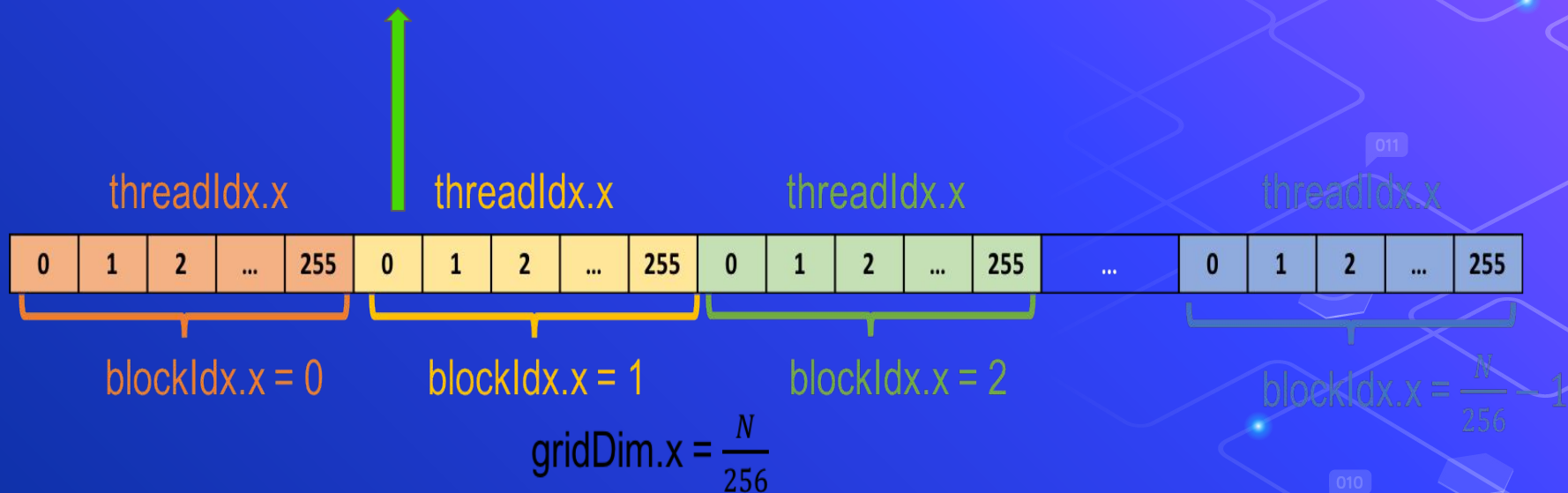
The thread index i is calculated by the following formula :

$$i = blockIdx.x * blockDim.x + threadIdx.x$$

CUDA Built-in Device Variables

- All `__global__` and `__device__` functions have access to these automatically defined variables
- `dim3 blockDim`;
 - Dimensions of the grid in blocks (at most 20)
- `dim3 blockIdx`;
 - Dimensions of the block in threads
- `dim3 blockIdx`;
 - Block index within the grid
- `dim3 threadIdx`;
 - Thread index within the block

$$\text{index} = \text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x} = 1 * 256 + 0 = 256$$



Variable Qualifiers

Example

Function Qualifiers

- __global__ - Invoked from within host (CPU) code.
- __device__ - Cannot be called from device (GPU) code.
- __host__ - Cannot be called from host (CPU) code.
- __host__ __device__ - Can only be executed by CPU, called from both.
- __restrict__ - Example use: controlling pointers.
- __constant__ - Example: will generate both CPU and GPU code.

```
__global__
void add(int *a, int *b, int *c, int n) {

    int i = threadIdx.x;

    __shared__ bool sharedArr[100];

    for (int i = 0; i < n; i++)
        c[i] = a[i] + b[i];
}

int main() {
    int *a, *b, *c;
```

Variable Qualifiers

- `__shared__` : accessible by all threads in the same block.
- `__device__` : accessible by all threads.

Function Qualifiers

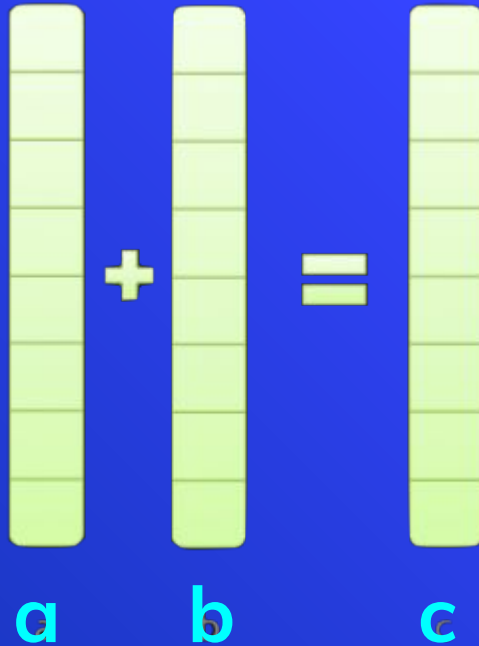
- `__global__` : Invoked from within host (CPU) code, cannot be called from device (GPU) code, must return void
- `__device__` : called from other GPU functions, cannot be called from host (CPU) code
- `__host__` : can only be executed by CPU, called from host
- `__host__` and `__device__` qualifiers can be combined
- Example use: controlling memory
- Compiler will generate both CPU and GPU code

A bundle of glowing blue fiber optic cables is the central focus, with light rays emanating from them. The background is dark, and there are faint, glowing blue circles and lines scattered around. On the right side, there is a white line-art network diagram with several nodes and connections. Some nodes are labeled with binary code: '001' at the top right, '011' in the middle right, and '010' at the bottom right. The overall aesthetic is futuristic and technological.

How to use?

Example:

Adding two vector **a** and **b** resulted in vector **c** with a, b and c have the same size.



No Parallel

```
int main() {
    int *a, *b, *c;
    int SIZE = 1000000;

    // Allocation memory in devices for a,b and c
    cudaMallocManaged(&a, SIZE*sizeof(int));
    cudaMallocManaged(&b, SIZE*sizeof(int));
    cudaMallocManaged(&c, SIZE*sizeof(int));

    // Simple initialization
    for (int i = 0; i < SIZE; i++) {
        a[i] = i;
        b[i] = i;
    }

    // Launch kernel on device
    add<<<1,1>>>(a,b,c,SIZE);

    // Free memory
    cudaFree(a);
    cudaFree(b);
    cudaFree(c);

    return 0;
}
```

```
#include <stdio.h>

__global__
void add(int *a, int *b, int *c, int n) {

    // Iteration to do addition
    for (int i = 0; i < n; i++)
        c[i] = a[i] + b[i];
}
```

- Using only one block with one thread.
- Using a loop to do addition on each index.

No Parallel

```
==12305== NVPROF is profiling process 12305, command: ./add_normal
==12305== Profiling application: ./add_normal
==12305== Profiling result:
```

Type	Time(%)	Time	Calls	Avg	Min
GPU activities:	100.00%	266.64ms	1	266.64ms	266.64ms
API calls:	50.29%	274.20ms	3	91.401ms	831.23us
	49.00%	267.15ms	1	267.15ms	267.15ms
	0.42%	2.3071ms	1	2.3071ms	2.3071ms
	0.19%	1.0519ms	3	350.62us	321.01us
	0.04%	236.00us	96	2.4580us	191ns
	0.04%	213.08us	1	213.08us	213.08us
	0.00%	22.577us	1	22.577us	22.577us
	0.00%	3.6050us	1	3.6050us	3.6050us
	0.00%	1.8500us	3	616ns	235ns
	0.00%	1.2910us	2	645ns	262ns
	0.00%	339ns	1	339ns	339ns

With Parallel

```
int main() {
    int *a, *b, *c;
    int SIZE = 1000000;
    int THREAD_SIZE = 256;

    // Allocation memory in devices for a,b and c
    cudaMallocManaged(&a, SIZE*sizeof(int));
    cudaMallocManaged(&b, SIZE*sizeof(int));
    cudaMallocManaged(&c, SIZE*sizeof(int));

    // Simple initialization
    for (int i = 0; i < SIZE; i++) {
        a[i] = i;
        b[i] = i;
    }

    // Launch kernel function with SIZE/THREAD_SIZE + 1 block
    // Each block has THREAD_SIZE threads
    add<<<SIZE/THREAD_SIZE + 1, THREAD_SIZE>>>(a,b,c,SIZE);

    // Free memory
    cudaFree(a);
    cudaFree(b);
    cudaFree(c);

    return 0;
}
```

```
#include <stdio.h>

__global__
void add(int *a, int *b, int *c, int n) {
    // Get thread index
    int index = threadIdx.x + blockIdx.x * blockDim.x;
    if (index < n)
        c[index] = a[index] + b[index];
}
```

- Using only THREAD_SIZE block, each block has $\text{SIZE/THREAD_SIZE} + 1$ threads.
- Each thread we do addition on one index.

With Parallel

```
==11831== NVPROF is profiling process 11831, command: ./add_parallel
```

```
==11831== Profiling application: ./add_parallel
```

```
==11831== Profiling result:
```

Type	Time(%)	Time	Calls	Avg	Min
GPU activities:	100.00%	57.537us	1	57.537us	57.537us
API calls:	98.42%	262.15ms	3	87.383ms	631.43us
	0.95%	2.5211ms	1	2.5211ms	2.5211ms
	0.39%	1.0363ms	3	345.43us	300.18us
	0.09%	249.13us	96	2.5950us	192ns
	0.08%	204.57us	1	204.57us	204.57us
	0.05%	145.72us	1	145.72us	145.72us
	0.01%	31.471us	1	31.471us	31.471us
	0.00%	4.1350us	1	4.1350us	4.1350us
	0.00%	2.2310us	3	743ns	227ns
	0.00%	1.4950us	2	747ns	402ns
	0.00%	324ns	1	324ns	324ns

No Parallel

```
==12305== NVPROF is profiling process 12305, command: ./add_normal
==12305== Profiling application: ./add_normal
==12305== Profiling result:
```

Type	Time(%)	Time	Calls	Avg	Min
GPU activities:	100.00%	266.64ms	1	266.64ms	266.64ms
API calls:	50.29%	274.20ms	3	91.401ms	831.23us
	49.00%	267.15ms	1	267.15ms	267.15ms
	0.42%	2.3071ms	1	2.3071ms	2.3071ms
	0.19%	1.0519ms	3	350.62us	321.01us
	0.04%	236.00us	96	2.4580us	191ns
	0.04%	213.08us	1	213.08us	213.08us
	0.00%	22.577us	1	22.577us	22.577us
	0.00%	3.6050us	1	3.6050us	3.6050us
	0.00%	1.8500us	3	616ns	235ns
	0.00%	1.2910us	2	645ns	262ns
	0.00%	339ns	1	339ns	339ns

With Parallel

```
==11831== NVPROF is profiling process 11831, command: ./add_parallel
==11831== Profiling application: ./add_parallel
==11831== Profiling result:
```

Type	Time(%)	Time	Calls	Avg	Min
GPU activities:	100.00%	57.537us	1	57.537us	57.537us
API calls:	98.42%	262.15ms	3	87.383ms	631.43us
	0.95%	2.5211ms	1	2.5211ms	2.5211ms
	0.39%	1.0363ms	3	345.43us	300.18us
	0.09%	249.13us	96	2.5950us	192ns
	0.08%	204.57us	1	204.57us	204.57us
	0.05%	145.72us	1	145.72us	145.72us
	0.01%	31.471us	1	31.471us	31.471us
	0.00%	4.1350us	1	4.1350us	4.1350us
	0.00%	2.2310us	3	743ns	227ns
	0.00%	1.4950us	2	747ns	402ns
	0.00%	324ns	1	324ns	324ns

No Parallel

```
#include <stdio.h>

__global__
void add(int *a, int *b, int *c, int n) {

    // Iteration to addition
    for (int i = 0; i < n; i++)
        c[i] = a[i] + b[i];
}

int main() {
    int *a, *b, *c;
    int SIZE = 1000000;

    // Allocation memory in devices for a,b and c
    cudaMallocManaged(&a, SIZE*sizeof(int));
    cudaMallocManaged(&b, SIZE*sizeof(int));
    cudaMallocManaged(&c, SIZE*sizeof(int));

    // Simple initialization
    for (int i = 0; i < SIZE; i++) {
        a[i] = i;
        b[i] = i;
    }

    // Launch kernel function with one block with one thread
    add<<<1,1>>>(a,b,c,SIZE);

    // Free memory
    cudaFree(a);
    cudaFree(b);
    cudaFree(c);

    return 0;
}
```

With Parallel

```
#include <stdio.h>

__global__
void add(int *a, int *b, int *c, int n) {
    // Get thread index
    int index = threadIdx.x + blockIdx.x * blockDim.x;
    if (index < n)
        c[index] = a[index] + b[index];
}

int main() {
    int *a, *b, *c;
    int SIZE = 1000000;
    int THREAD_SIZE = 256;

    // Allocation memory in devices for a,b and c
    cudaMallocManaged(&a, SIZE*sizeof(int));
    cudaMallocManaged(&b, SIZE*sizeof(int));
    cudaMallocManaged(&c, SIZE*sizeof(int));

    // Simple initialization
    for (int i = 0; i < SIZE; i++) {
        a[i] = i;
        b[i] = i;
    }

    // Launch kernel function with SIZE/THREAD_SIZE + 1 blocks
    // Each block has THREAD_SIZE threads
    add<<<SIZE/THREAD_SIZE + 1,THREAD_SIZE>>>(a,b,c,SIZE);

    // Free memory
    cudaFree(a);
    cudaFree(b);
    cudaFree(c);

    return 0;
}
```

001

Experiment: runtime: 126x

	Plus	Quadratic
GPU one thread	0.266s	
GPU Parallel	0.000057s	
CPU	0.0072s	

Thanks

!

Nguyễn Minh Trí
Đỗ Trí Nhân



References

Ref

[An Introduction to CUDA Programming](#)

<https://docs.nvidia.com/cuda/parallel-thread-execution/index.html>

<https://stackoverflow.com/questions/5211746/what-is-cuda-like-what-is-it-for-what-are-the-benefits-and-how-to-start>

https://www.tutorialspoint.com/cuda/cuda_key_concepts.htm

https://www.tutorialspoint.com/cuda/cuda_keywords_and_thread_organization.htm

<https://developer.nvidia.com/cuda-faq>

<https://www.infoworld.com/article/3299703/what-is-cuda-parallel-programming-for-gpus.html>

<https://www.sciencedirect.com/science/article/pii/S0167819119301759>

<https://spectrum.ieee.org/static/interactive-the-top-programming-languages-2016#index/2016/1/1/1/1/1/50/1/50/1/50/1/30/1/30/1/30/1/20/1/20/1/5/1/5/1/20/1/100/>

<https://medium.com/@inaccel/cpu-gpu-fpga-or-tpu-which-one-to-choose-for-my-machine-learning-training-948902f058e0>

<http://www.diva-portal.org/smash/get/diva2:447977/FULLTEXT01.pdf>

<http://proceeding.vap.ac.vn/index.php/proceedingvap/article/view/000214/203>

<https://www.slideserve.com/ryu/cuda-programming>

<https://www.myzhar.com/blog/tutorials/tutorial-nvidia-gpu-cuda-compute-capability/>

https://books.google.com.vn/books?id=CYxjDwAAQBAJ&pg=PA7&lpg=PA7&dq=pytorch+is+wrapper+of+cuda&source=bl&ots=43Smx_h6jw&sig=ACfU3U20IxlU8l4kSDSTS5RGQEFk6QMnuQ&hl=en&sa=X&ved=2ahUKewiz_tLSqa7pAhXOzIsBHU55DCMQ6AEwD3oECBIQAQ#v=onepage&q=pytorch%20is%20wrapper%20of%20cuda&f=false

<http://www.c4learn.com/cplusplus/cpp-variable-naming/>

<https://www.infoworld.com/article/3299703/what-is-cuda-parallel-programming-for-gpus.html>

<https://tatourian.blog/2013/09/03/nvidia-gpu-architecture-cuda-programming-environment/>

Main

<https://www.yumpu.com/en/document/read/50976433/cuda-parallel-programming-tutorial>

<http://www.diva-portal.org/smash/get/diva2:447977/FULLTEXT01.pdf>

<https://www.slideserve.com/ryu/cuda-programming>

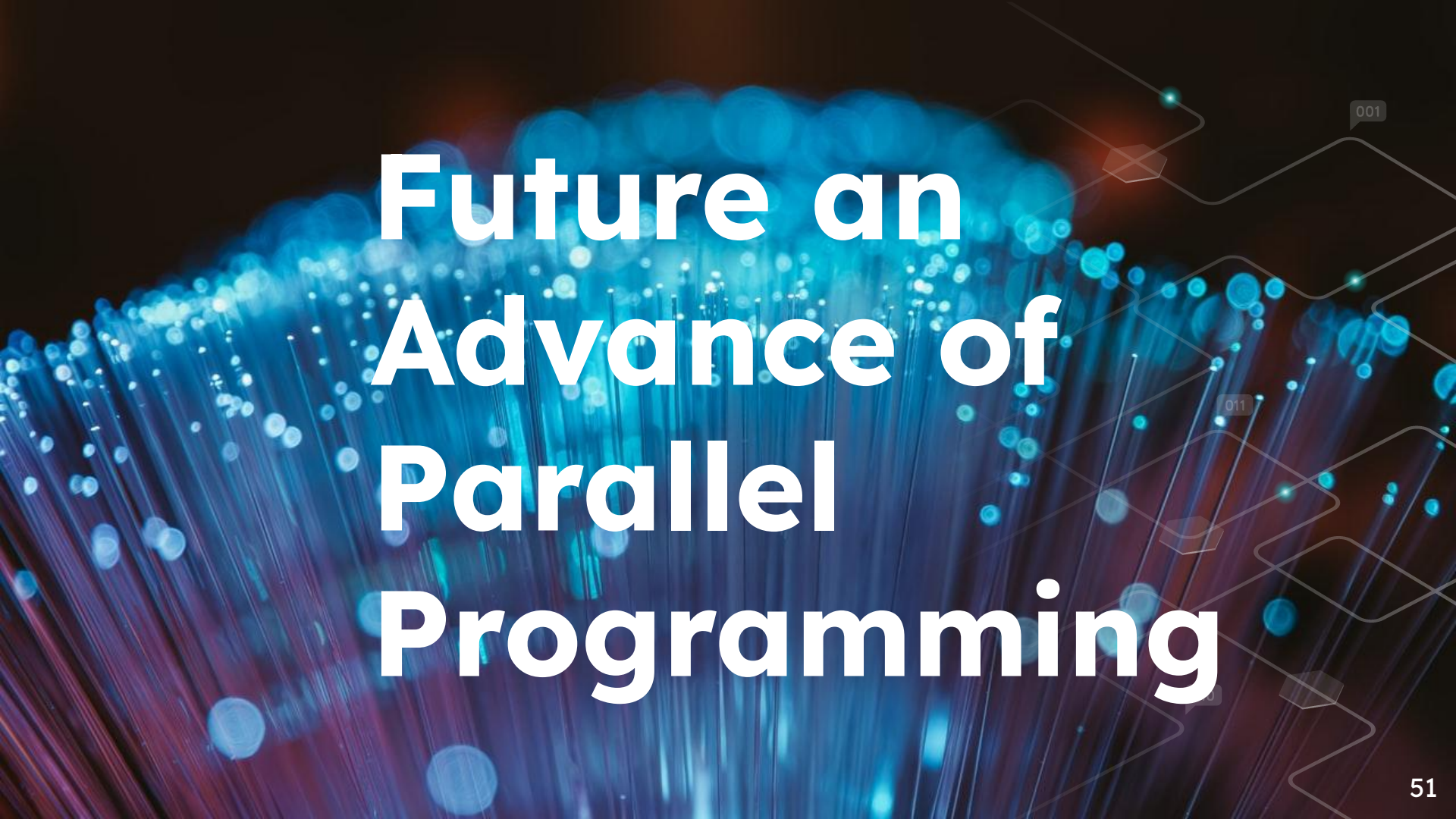
<https://www.nvidia.com/docs/IO/116711/sc11-cuda-c-basics.pdf>



Questions?

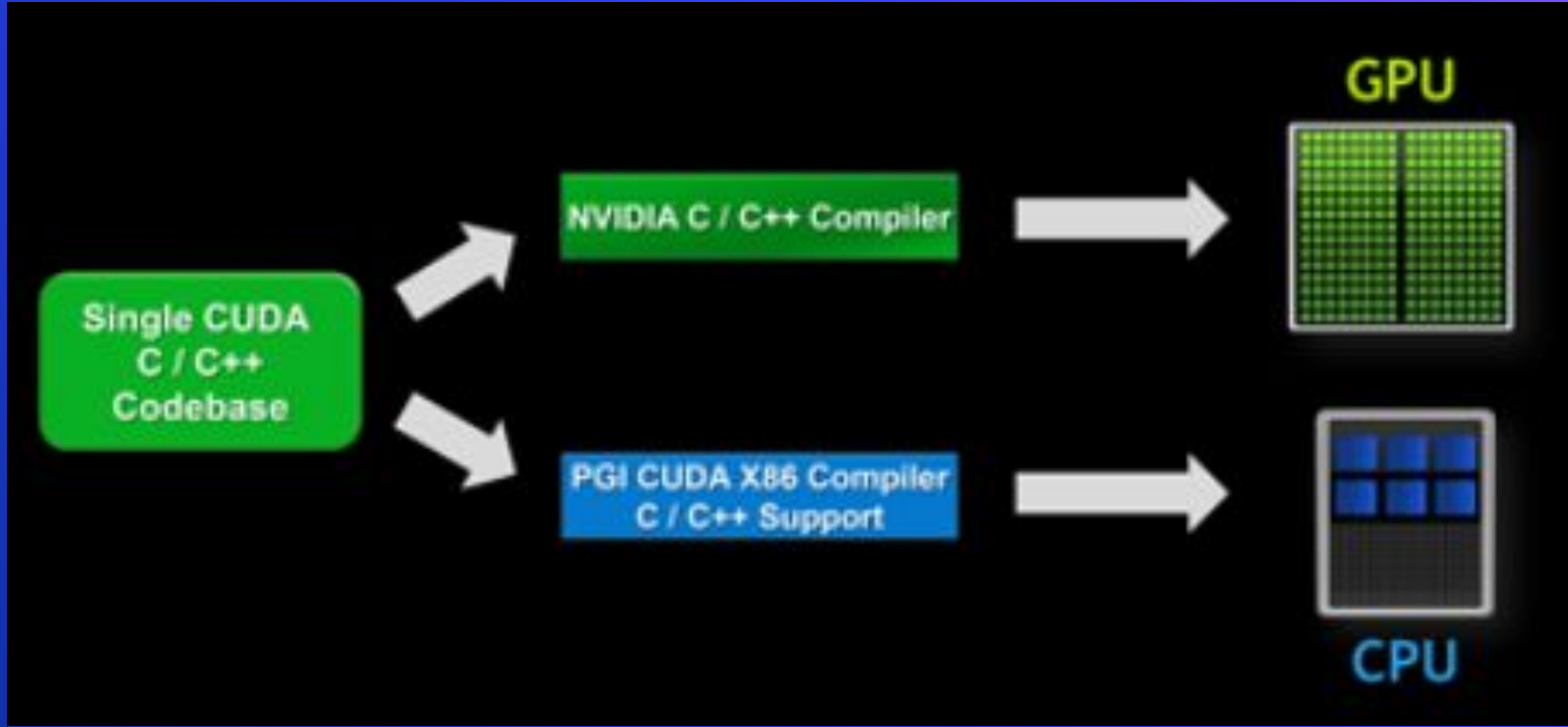
The image shows a laptop screen with a blue background. On the screen, the text "Solving 1000 quadratic equations" is written in white. The background of the screen features a faint, light blue network diagram with lines and nodes. There are also small, semi-transparent speech bubble icons containing binary code (001, 011, 010) scattered across the screen. The laptop itself is dark blue and is shown from a slightly elevated perspective.

Solving 1000 quadratic equations

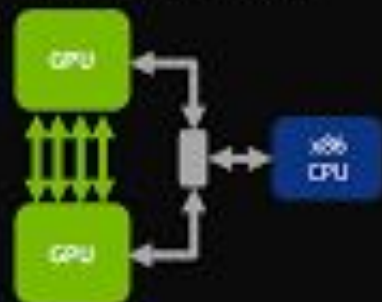


Future an Advance of Parallel Programming

CUDA compiler - Future (option CPU or GPU)



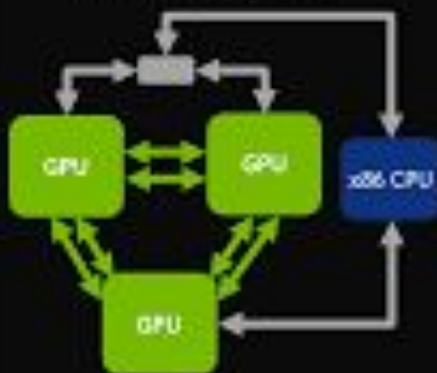
2 GPUs per Node

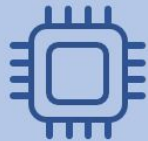


4 GPUs per Node



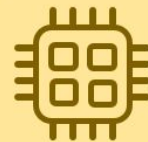
3 GPUs per Node





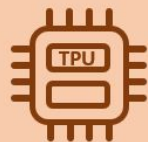
CPU

- Small models
- Small datasets
- Useful for design space exploration



GPU

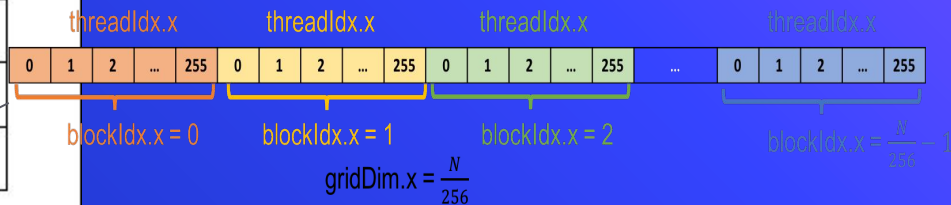
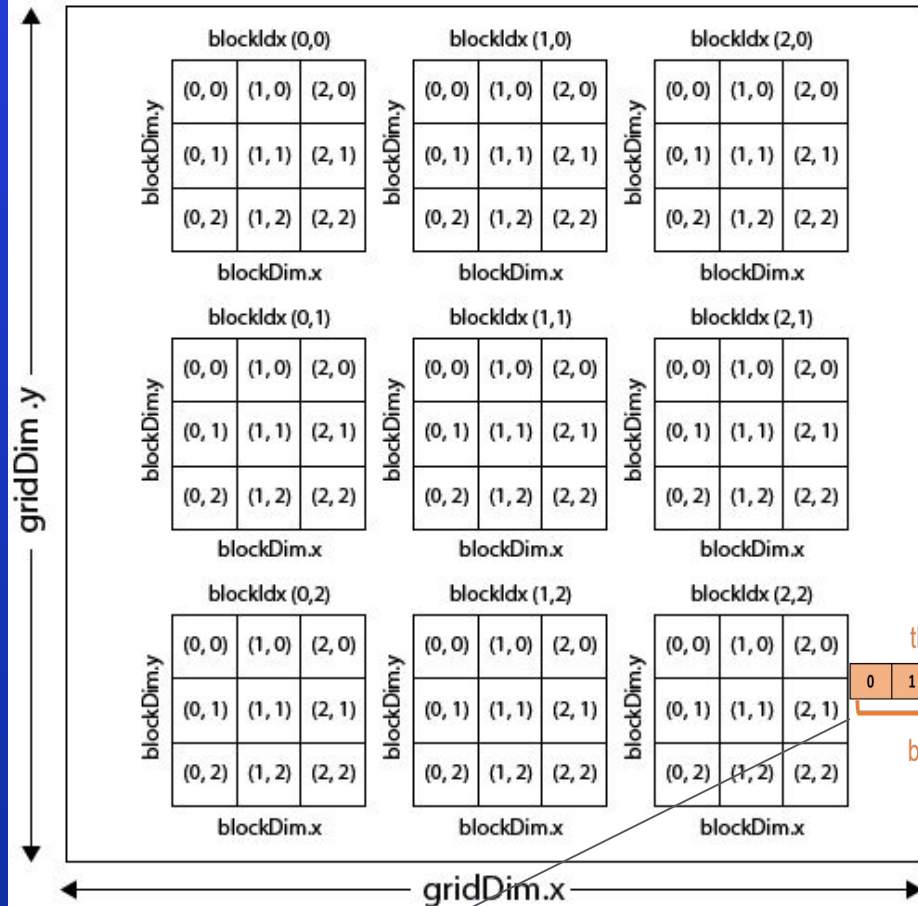
- Medium-to-large models, datasets
- Image, video processing
- Application on CUDA or OpenCL



TPU

- Matrix computations
- Dense vector processing
- No custom TensorFlow operations

CUDA Grid





NVIDIA CUDA



Open CL

Sample Code

Example: Increment Array Elements

CPU program

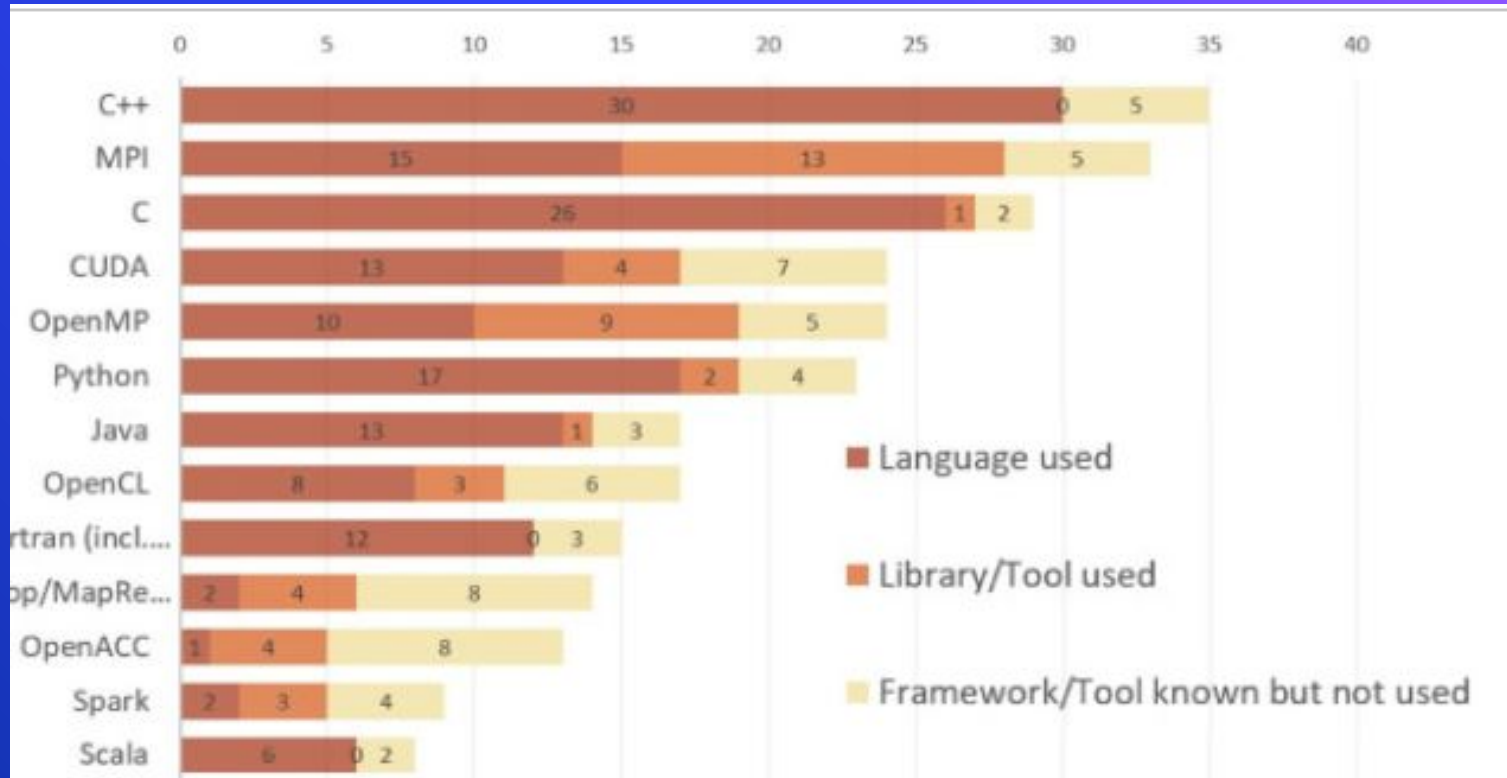
```
void increment_cpu(float *a, float b, int N)
{
    for (int idx = 0; idx < N; idx++)
        a[idx] = a[idx] + b;
}
```

```
void main()
{
    .....
    increment_cpu(a, b, N);
}
```

CUDA program

```
__global__ void increment_gpu(float *a, float b, int N)
{
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if (idx < N)
        a[idx] = a[idx] + b;
}
```

```
void main()
{
    ...
    dim3 dimBlock (blocksize);
    dim3 dimGrid( ceil( N / (float)blocksize ) );
    increment_gpu<<<dimGrid, dimBlock>>>(a, b, N);
}
```



<https://www.sciencedirect.com/science/article/pii/S0167819119301759>

How GPU Acceleration Works

