

Sure! Let's break it all down in a simple way. First, I'll explain the **concept of vector addition**, then I'll walk you through the code line-by-line.

What is Vector Addition?

Imagine you have two lists (arrays) of numbers:

A = [1, 1, 1]

B = [2, 2, 2]

Vector addition means adding them **element by element**:

C = [1+2, 1+2, 1+2] = [3, 3, 3]

This is easy for a computer — and **very easy to parallelize**, because each addition is independent of the others. That's where **CUDA** (GPU programming) comes in: we can run many additions at the same time on a GPU!

What Is This CUDA Code Doing?

This code:

1. Creates two big arrays (**A** and **B**) of size 1,000,000.
 2. Adds them together using the **GPU**.
 3. Measures the time taken by the GPU.
 4. Verifies the result.
-

Code Breakdown (Step by Step)

1. Include Libraries

```
#include <iostream>    // For input/output
#include <cuda_runtime.h> // CUDA API
#include <ctime>        // For timing
```

2. Constants

```
constexpr int VECTOR_SIZE = 1'000'000;    // Size of arrays
constexpr int THREADS_PER_BLOCK = 256;    // GPU threads per block
```

⚙️ 3. GPU Kernel: The GPU Function

```
__global__ void vectorAdd(const int *A, const int *B, int *C, int n) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if (idx < n) {
        C[idx] = A[idx] + B[idx];
    }
}
```

This is the function that runs **on the GPU**. It:

- Calculates a **unique index** for each thread (`idx`)
- Adds `A[idx] + B[idx]` and stores it in `C[idx]`

Each thread runs one element of the addition.

🧠 4. Main Function: Allocate Memory

Host (CPU) Memory

```
cudaMallocHost(&h_A, bytes); // Pinned memory (faster transfer)
cudaMallocHost(&h_B, bytes);
cudaMallocHost(&h_C, bytes);
```

Device (GPU) Memory

```
cudaMalloc(&d_A, bytes);    // Allocate on GPU
cudaMalloc(&d_B, bytes);
cudaMalloc(&d_C, bytes);
```

📝 5. Initialize Data on CPU

```
for (int i = 0; i < VECTOR_SIZE; ++i) {
    h_A[i] = 1;
    h_B[i] = 2;
}
```

📦 6. Copy Data to GPU

```
cudaMemcpy(d_A, h_A, bytes, cudaMemcpyHostToDevice);
cudaMemcpy(d_B, h_B, bytes, cudaMemcpyHostToDevice);
```

7. Launch Kernel on GPU

```
int blocks = (VECTOR_SIZE + THREADS_PER_BLOCK - 1) / THREADS_PER_BLOCK;
```

This calculates how many **blocks** are needed.

Then we time and run it:

```
clock_t start = clock();  
vectorAdd<<<blocks, THREADS_PER_BLOCK>>>(d_A, d_B, d_C, VECTOR_SIZE);  
cudaDeviceSynchronize();  
clock_t end = clock();
```

8. Copy Result Back to CPU

```
cudaMemcpy(h_C, d_C, bytes, cudaMemcpyDeviceToHost);
```

9. Validate the Result

```
for (int i = 0; i < VECTOR_SIZE; ++i) {  
    if (h_C[i] != h_A[i] + h_B[i]) {  
        valid = false;  
        break;  
    }  
}
```

10. Show Execution Time

```
double ms = 1000.0 * (end - start) / CLOCKS_PER_SEC;  
std::cout << "CUDA execution time: " << ms << " ms\n";
```

11. Free Memory

```
cudaFree(...);    // Free GPU memory  
cudaFreeHost(...); // Free CPU memory
```

Output Example

You might see:

CUDA execution time: 0.55 ms

Result is correct!

Why Use GPU for Vector Addition?

While it seems simple, vector operations are **common in graphics, physics, AI, and simulations** — and GPUs can do **thousands of operations in parallel**, making them super fast.

Let me know if you want a CPU-only version of this code, or want to test the speed difference between CPU and GPU!