# CUDA with numba python

1. What is CUDA?
2. What is Numba?
3. How do CUDA and Numba work together?
4. Basic CUDA concepts for parallel programming
5. Simple coding examples with Numba + CUDA
6. When to use it

## 1. What is CUDA?

CUDA (Compute Unified Device Architecture) is a parallel computing platform and application programming interface (API) model created by NVIDIA . It allows developers to use NVIDIA GPUs for general purpose processing (GPGPU).

With CUDA, you can offload compute-intensive tasks from the CPU to the GPU, leading to significant speedups — especially for data-parallel operations.

## 2. What is Numba?

Numba is a just-in-time (JIT) compiler in Python that translates a subset of Python and NumPy code into fast machine code using the LLVM compiler infrastructure.

It supports:

- CPU-based JIT compilation

- GPU acceleration via CUDA

The key feature here is that you can write Python functions and compile them to run on the GPU using CUDA , with minimal changes and no need to write C/C++ or CUDA C code directly.

## 3. How do CUDA and Numba Work Together?

Numba provides a high-level interface to program NVIDIA GPUs using CUDA, directly from Python.

You can write a function in Python, decorate it with special Numba decorators like `@cuda.jit`, and Numba will compile and execute it on your GPU.

This makes it very convenient for Python users who want to harness GPU power without diving deep into CUDA C.

## 4. Basic CUDA Concepts

Here are some fundamental CUDA concepts:

| CONCEPT | DESCRIPTION |
|---|---|
| Thread | The smallest unit of execution. Each thread runs the same kernel function. |
| Block | A group of threads. Threads in a block can communicate/synchronize. |
| Grid | A collection of blocks. Blocks are independent and cannot synchronize. |
| Kernel | A function that runs in parallel on the GPU. You launch kernels from the host (CPU). |

Example analogy:

- Think of a thread as a worker
- A block is a team of workers
- A grid is a collection of teams

## 5. Simple Coding Examples with Numba + CUDA

Let's start with a basic example: adding two arrays element-wise using CUDA via Numba.

### ☑ Prerequisites:

Make sure you have:

```
pip install numba numpy
```

Also ensure you have an NVIDIA GPU and the appropriate drivers installed.

```python
import numpy as np
from numba import cuda

# Define the CUDA kernel
@cuda.jit
def add_kernel(a, b, c):
    i = cuda.grid(1)  # Get index in 1D grid
    if i < a.shape[0]:  # Prevent out-of-bounds access
        c[i] = a[i] + b[i]

# Host code
def main():
    n = 100000
    a = np.arange(n, dtype=np.float32)
    b = np.arange(n, dtype=np.float32)
    c = np.zeros(n, dtype=np.float32)

    # Allocate device memory and copy data
    d_a = cuda.to_device(a)
    d_b = cuda.to_device(b)
    d_c = cuda.device_array_like(c)

    # Set up grid/block sizes
    threads_per_block = 256
    blocks_per_grid = (n + threads_per_block - 1) // threads_per_block

    # Launch the kernel
    add_kernel[blocks_per_grid, threads_per_block](d_a, d_b, d_c)

    # Copy result back to host
    c = d_c.copy_to_host()

    print("Result (first 5):", c[:5])

if __name__ == '__main__':
```

Let's compare the performance of a simple array addition using:

1. A CPU-based Numba-optimized version
2. A GPU-based Numba CUDA version

You've already provided the CUDA GPU-based code (array addition with `@cuda.jit`). Now, we'll implement an equivalent CPU version using Numba's `@jit` decorator and measure both versions.

## 🧪 Goal

Compare the execution time of:

- CPU-based array addition (`@jit`)
- GPU-based array addition (`@cuda.jit`)

```python
1    from numba import jit, cuda
2    import numpy as np
3    from time import perf_counter
4
5    # CUDA kernel
6    @cuda.jit
7    def add_kernel_gpu(a, b, c):
8        i = cuda.grid(1)
9        if i < a.size:
10           c[i] = a[i] + b[i]
11
12   # CPU kernel with JIT optimization
13   @jit(nopython=True)
14   def add_kernel_cpu(a, b, c):
15       for i in range(a.size):
16           c[i] = a[i] + b[i]
17       return c
18
19   def main():
20       n = 10_000_000
21       print(f"Array size: {n}")
22
23       # Host arrays
24       a = np.arange(n).astype(np.float32)
25       b = 2 * np.arange(n).astype(np.float32)
26       c_cpu = np.zeros_like(a)
27
28       # CPU timing
29       start = perf_counter()
30       add_kernel_cpu(a, b, c_cpu)
31       end = perf_counter()
32       print(f"CPU Time: {end - start:.4f} seconds")
33
34       # GPU setup
35       d_a = cuda.to_device(a)
36       d_b = cuda.to_device(b)
37       d_c = cuda.device_array_like(a)
38
39
```

```python
40        # Kernel-only timing
41        start = perf_counter()
42        threads_per_block = 256
43        blocks_per_grid = (n + threads_per_block - 1) // threads_per_block
44        add_kernel_gpu[blocks_per_grid, threads_per_block](d_a, d_b, d_c)
45        cuda.synchronize()  # Wait for GPU to finish
46        end = perf_counter()
47        print(f"GPU (kernel only): {end - start:.4f} seconds")
48
49        # Optional: Copy result back
50        # c_gpu = d_c.copy_to_host()
51
52    if __name__ == '__main__':
53        main()
```

## Sample Output

```
Array size: 10000000
CPU Time: 0.0957 seconds
GPU (kernel only): 0.0155 seconds
```

## 🧠 Key Takeaways

| FEATURE | CPU ( `@JIT` ) | GPU ( `@CUDA.JIT` ) |
|---|---|---|
| **Execution Unit** | Single-core / Multi-core CPU | Massively parallel GPU |
| **Memory** | RAM | Device memory (VRAM) |
| **Overhead** | Low (no transfer cost) | Includes device memory transfers |
| **Performance** | Good for small data | Excels at large-scale parallelism |

💡 For small arrays (<10,000 elements), CPU may outperform GPU due to memory copy overhead.

## 🚀 Try It With Larger Data

Try increasing `n` to `100_000_000` and observe how the GPU scales better than the CPU.

# 🚀 Matrix Multiplication with CUDA (via Numba)

```python
import numpy as np
from numba import cuda, float32
from time import perf_counter

# Define block size for CUDA
TPB = 16  # Threads per block (Typically 16x16 for matrices)

# CPU version for comparison
def matmul_cpu(A, B, C):
    for row in range(C.shape[0]):
        for col in range(C.shape[1]):
            tmp = 0.
            for k in range(A.shape[1]):
                tmp += A[row, k] * B[k, col]
            C[row, col] = tmp
    return C

# CUDA kernel using Numba
@cuda.jit
def matmul_gpu(A, B, C):
    row = cuda.threadIdx.y + cuda.blockIdx.y * cuda.blockDim.y
    col = cuda.threadIdx.x + cuda.blockIdx.x * cuda.blockDim.x

    if row < C.shape[0] and col < C.shape[1]:
        tmp = 0.
        for k in range(A.shape[1]):
            tmp += A[row, k] * B[k, col]
        C[row, col] = tmp

# Helper function to run GPU version
def run_matmul_gpu(A, B):
    rows, cols = A.shape[0], B.shape[1]
    C_gpu = np.zeros((rows, cols), dtype=np.float32)
    threads_per_block = (TPB, TPB)
    blocks_per_grid_x = (cols + TPB - 1) // TPB
    blocks_per_grid_y = (rows + TPB - 1) // TPB
    blocks_per_grid = (blocks_per_grid_x, blocks_per_grid_y)

    d_A = cuda.to_device(A)
```

```python
40        d_B = cuda.to_device(B)
41        d_C = cuda.to_device(C_gpu)
42
43        matmul_gpu[blocks_per_grid, threads_per_block](d_A, d_B, d_C)
44
45        C_result = d_C.copy_to_host()
46        return C_result
47
48  # Main test function
49  def main():
50        # You can change this size to see different performance impacts
51        N = 1024
52        A = np.random.rand(N, N).astype(np.float32)
53        B = np.random.rand(N, N).astype(np.float32)
54        C_cpu = np.zeros((N, N), dtype=np.float32)
55
56        print(f"Matrix size: {N}x{N}")
57
58        # CPU Timing
59        start = perf_counter()
60        matmul_cpu(A, B, C_cpu)
61        end = perf_counter()
62        print(f"CPU Time: {end - start:.4f} seconds")
63
64        # GPU Timing
65        start = perf_counter()
66        C_gpu = run_matmul_gpu(A, B)
67        end = perf_counter()
68        print(f"GPU Time: {end - start:.4f} seconds")
69
70        # Check result accuracy
71        assert np.allclose(C_cpu, C_gpu, atol=1e-4), "CPU and GPU results
72  do not match!"
73
74  if __name__ == '__main__':
75        main()
```

## 📊 Sample Output

Depending on your hardware, you may get something like:

```
Matrix size: 1024x1024
CPU Time: 15.2345 seconds
GPU Time: 0.2134 seconds
```

This shows a significant speedup using the GPU!