

Assignment 16

CUDA programming

1. What is CUDA ?

A. CUDA = Compute Unified Device Architecture

CUDA is a parallel computing platform and programming model developed by NVIDIA for general computing on graphical processing units (GPUs).

It allows to access the raw computing power of CUDA GPUs to process data faster than with traditional CPUs. CUDA can achieve higher parallelism and efficiency than general-purpose CPU code using parallel processes. It enables parallel processing by breaking down a task into thousands of smaller "threads" executed independently.

2. What is the prerequisite for learning CUDA?

- Understanding parallel computing concepts such as **threads**, **blocks**, **grids**, and **synchronization** is crucial for CUDA programming.
- Many applications of CUDA involve numerical computations, such as **matrix operations** and **linear algebra**.
- CUDA enable or **capable GPU** devices, also have to installed **CUDA toolkit** and **nvidia developer driver**.

3. Which are the languages that support CUDA?

- #### A.
- C++, C, C#, Fortran, Python, Java

4. What do you mean by a CUDA ready architecture?

- #### A.
- CUDA Ready Architecture refers to a hardware architecture designed by NVIDIA to support CUDA and their parallel computing platform and programming model. In short the GPU architecture is designed and optimized to support CUDA

5. How CUDA works?

- #### A.
- When a processor is given a task, it will pass the instructions for that task to the GPU. The GPU will then do its work, following the instructions from the CPU. GPUs run one

kernel (a group of tasks) at a time. Each kernel consists of blocks, which are independent groups of ALUs. Each block contains threads, which are levels of computation. The threads in each block typically work together to calculate a value. Once the job is completed, the results from the GPU are given back to the CPU.

6. What are the benefits and limitations of CUDA programming?

A. Benefits :

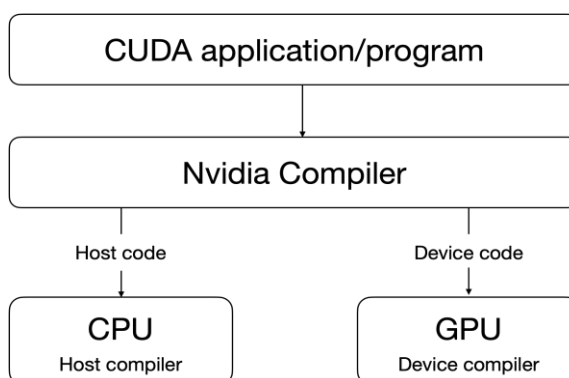
- Massive parallelism
- Performance boost
- Integrated memory and shared memory

Limitations :

- Works only on nvidia GPUs
- Highly Parallelizable Problems Only(Not all problems can be effectively parallelized for GPUs)
- Limited Portability (Porting to other GPU architectures requires significant effort.)

7. Understand and explain the CUDA program structure with an example.

- A. CUDA programming includes code for both the GPU and CPU. By default, a typical C programme is a CUDA programme that simply contains the host code. The CPU is referred to as the host, while the GPU is referred to as the device. While the host code can be compiled using a regular C compiler such as GCC, the device code requires a specific compiler to comprehend the API functions that are used. For Nvidia GPUs, the compiler is known as the NVCC (Nvidia C Compiler).



The GPU runs the device code, whereas the CPU runs the host code. The NVCC runs a CUDA programme and isolates the host code from the device code. To do this, specific CUDA keywords are searched for. The code that is intended to run on the GPU (device code) is identified by special CUDA keywords known as 'Kernels' that

label data-parallel functions. The NVCC further compiles the device code, which is then run on the GPU.

8. Explain CUDA thread organization.

- A.** Organizing CUDA threads is similar to forming a team to collaborate on a large project. Here's a simplified breakdown:

Grids: Think of a grid as the overall project. It is made up of several smaller sections known as blocks.

Blocks: Each block represents a smaller team inside the project. They collaborate on a certain aspect of the project and can interact with one another effortlessly.

Threads: Each block has individual workers known as threads. They function similarly to team members who conduct the real work. Each thread completes a little portion of the job allocated to the block.

Divide the work into grids, which are then divided into blocks, with each block containing many threads. This arrangement helps to organize and perform duties efficiently on the GPU.

9. Install and try CUDA sample program and explain the same. (installation steps).

- A.** CUDA installation :

Pre-requisite:

- anaconda
- Python
- Tensorflow
- Jupyter
- CUDA enabled graphic card

To install CUDA run this following command in conda power shell or in its environment.

```
conda install -c conda-forge cudatoolkit=11.2 cudnn=8.1.0 -y
```

```

Windows PowerShell
PS C:\Users\Meetradsinh> conda install -c conda-forge cudatoolkit=11.2 cudnn=8.1.0 -y
Retrieving notices: ...working... done
Collecting package metadata (current_repodata.json): done
Solving environment: unsuccessful initial attempt using frozen solve. Retrying with flexible solve.
Collecting package metadata (repodata.json): done
Solving environment: done

==> WARNING: A newer version of conda exists. <==
current version: 23.7.4
latest version: 24.3.0

Please update conda by running

    $ conda update -n base -c defaults conda

Or to minimize the number of packages updated during conda update use

    conda install conda=24.3.0

## Package Plan ##

environment location: C:\Users\Meetradsinh\anaconda3

added / updated specs:
- cudatoolkit=11.2
- cudnn=8.1.0

The following packages will be downloaded:

package | build | size | channel
-----|-----|-----|-----
archspec-0.2.3 | pyhd8ed1ab_0 | 48 KB | conda-forge
conda-23.9.0 | py311h1ea47a8_2 | 1.2 MB | conda-forge
cudatoolkit-11.2.2 | h933977f_10 | 879.9 MB | conda-forge
cudnn-8.1.0.77 | h3e0f4f4_0 | 610.8 MB | conda-forge
python_abi-3.11 | 2_cp311 | 5 KB | conda-forge
truststore-0.8.0 | pyhd8ed1ab_0 | 20 KB | conda-forge
-----|-----|-----|-----
Total: | | 1.46 GB |

```

```

Jupyter Assignment_16 Last Checkpoint: 10 minutes ago
File Edit View Run Kernel Settings Help Trusted
+ + + + + Code
Open in... Python 3 (ipykernel)

[1]: import tensorflow as tf
[2]: gpus = tf.config.list_physical_devices("GPU")
[6]: if gpus:
      for gpu in gpus:
          print("Found a GPU with the name:", gpu)
      else:
          print("Failed to detect a GPU.")

Found a GPU with the name: PhysicalDevice(name='/physical_device:GPU:0', device_type='GPU')

[7]: !nvidia-smi

Wed Apr 3 15:10:23 2024
+-----+
| NVIDIA-SMI 551.86 | Driver Version: 551.86 | CUDA Version: 12.4 |
+-----+
| GPU Name | TCC/WDDM | Bus-Id | Disp.A | Volatile Uncorr. ECC |
| Fan Temp Perf | Pwr:Usage/Cap | Memory-Usage | GPU-Util | Compute M. |
|=====+-----+
| 0 NVIDIA GeForce RTX 4060 ... | WDDM | 00000000:01:00:0 Off | | N/A |
| N/A 47C P8 | 4W / 35W | 0MiB / 8188MiB | 0% | Default |
|=====+-----+
| Processes: |
| GPU GI CI PID Type Process name | GPU Memory |
| ID ID | | | | | Usage |
|=====+-----+
| No running processes found |
+-----+

[ ]:

```

As shown in above image if CUDA was successfully installed then by running above code it will show available gpu and also by running nvidia-smi it will show all detail of gpu and also show CUDA version.

Code example :

```

import numpy as np
from numba import cuda

# Define the CUDA kernel function

@cuda.jit
# @cuda.jit(gridsize=(arr.size // threadsperblock + 1, 1), blocksize=(threadsperblock,))
def add_to_array(arr):
    """
    This kernel adds 1.0 to each element of the input array.
    Each thread processes one element of the array.
    """

    # Get the unique thread index within a block
    idx = cuda.threadIdx.x

    # Check if the thread index is within the array bounds
    if idx < arr.size:
        arr[idx] += 1.0

# Create a sample NumPy array on the CPU (host)
arr = np.array([1.0, 2.0, 3.0])

# Allocate memory for the array on the GPU (device)
d_arr = cuda.to_device(arr)

# Configure the number of threads per block (adjust as needed)
threadsperblock = 512

gridsize = (arr.size // threadsperblock + 1, 1)
add_to_array(gridsize, threadsperblock)(d_arr)

# Copy results back from GPU to CPU
host_arr = d_arr.copy_to_host() # Use copy_to_host instead of to_host
arr = host_arr # Assign the copied data to the original array
arr

```

Output :

```
array([3., 4., 5.])
```