

Assignment 2

- Title: Vector and Matrix Operations using CUDA.

- Problem Statement:

Design a parallel algo to

1. Add two large vectors
2. Multiply vector & matrix
3. Multiply 2 $N \times N$ arrays $\approx N^2$ processors

- Objectives:

- To learn abt CUDA architecture
- To learn CUDA programming concepts

- Outcomes:

We will be able to learn about CUDA architecture & programming

- Requirements:

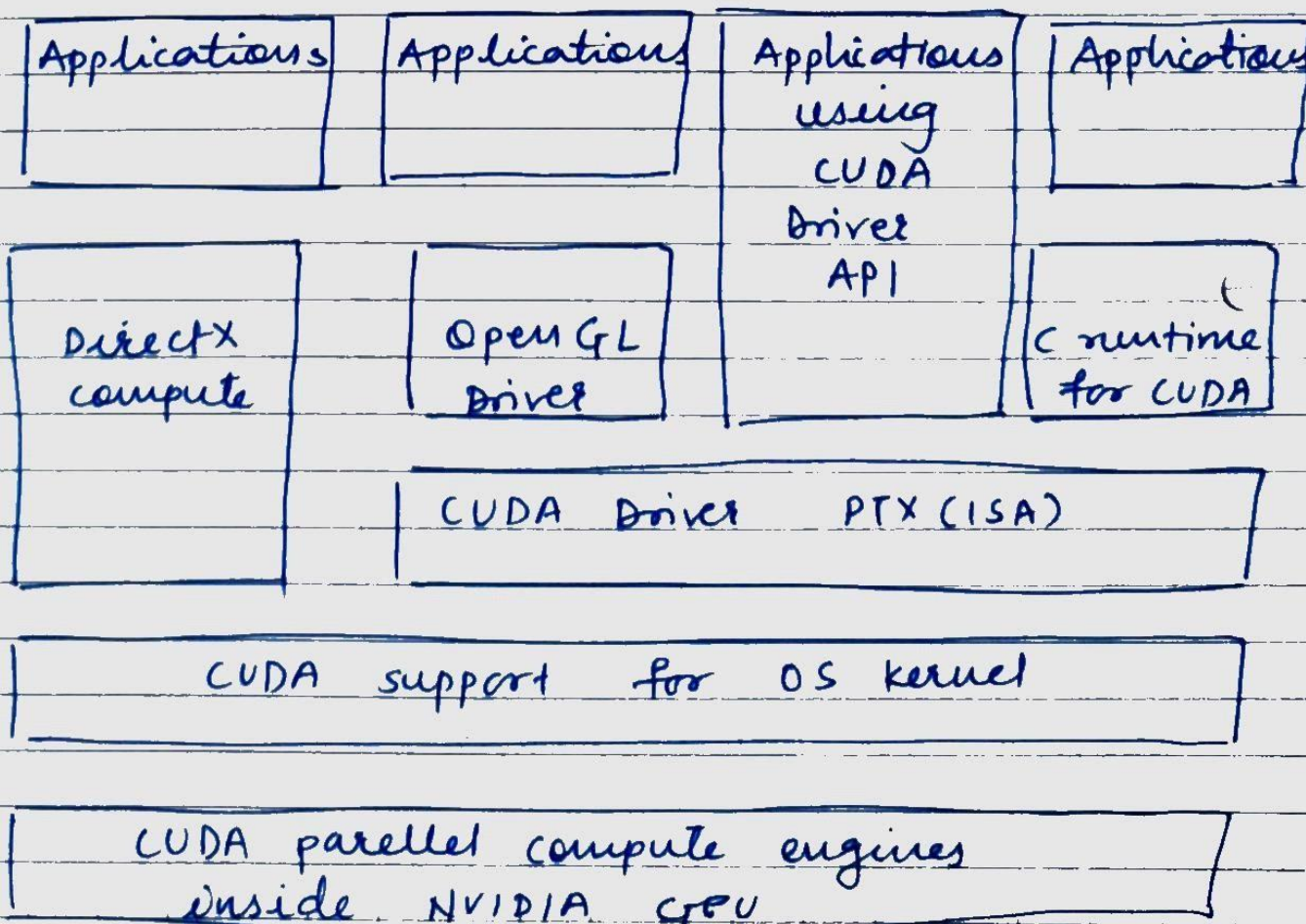
- OS: Fedora 20 / Ubuntu (64-bit)
- CUDA API
- nvcc compiler
- Nvidia GPU (Geforce 920M)
- Nsight Eclipse IDE
- RAM : 4GB
- HDD : 500GB

• Theory:

→ CUDA Architecture:

The architecture consists of several components as

- (1) Parallel compute engines in NVIDIA GPU
- (2) OS-kernel level support
- (3) user-mode driver providing device-level-API
- (4) PTX instruction set architecture for parallel computing kernel & functions.



→ Advantages of CUDA:

1. Unified Memory
2. scattered reads - code can be read from arbitrary addresses in memory
3. faster computational to & from GPU
4. full support for integers & bitwise operations
5. shared memory between threads.

→ Limitations of CUDA:

1. Newer versions of CUDA are now processed acc. to C++ syntax rules instead of C
2. Copying betⁿ host & device memory may incur a performance hit due to system bandwidth & latency.
3. CUDA enabled GPUs are available from NVIDIA only.
4. Threads should be running in groups of at least 32 for best performance.

→ CUDA programming:

Let M & N be the input vectors (or matrices) & P be the result obtained such that

$$P = M + N \quad (\text{vector addition})$$

$$P = M * N \quad (\text{vector-matrix multiplication})$$

$$P = M * N \quad (\text{matrix multiplication})$$

- With CUDA programming, each element in P can be obtained from one thread.
- As a result, the problem is decomposed into n (or $n \times n$) threads for parallel reduction.
- Thread IDs are used to differentiate the data with execution & storing the result.
- Test cases and Analysis:-

Operation	I/p size	Sequential time	Parallel time	Efficiency
Vector Addition	$n = 256$	0.01	0.02	0.5
	$n = 1024$	0.01	0.02	0.5
	$n = 2048$	0.02	0.01	2.0
Vector	$n = 256$	0.003	0.082	0.037
Matrix	$n = 1024$	0.013	0.135	0.689
Multi- plication	$n = 2048$	0.367	0.133	2.759
Matrix	$n = 256$	0.480	0.02	24.0
Multipl- ication	$n = 1024$	0.620	0.133	4.66
	$n = 2048$	0.754	0.136	5.544

$$\text{Efficiency} = \frac{HCSA}{WCPA}$$

Here, we observe that parallel algorithm is a major improvement for matrix-matrix multiplication while it is sufficiently significant for increasingly input size for vector-matrix multiplication. No significant improvement gained for vector addition of similar sizes but better for large values.

Input:

Vector 1 \rightarrow 5, 7, 9, 11, 4, 7, 6, 2, 0, 1

Vector 2 \rightarrow 6, 2, 7, 1, 0, 4, 13, 17, 12, 5

Output:

Sum of vectors \rightarrow 11, 9, 16, 12, 4, 11, 19, 19, 12, 6

- Conclusion:

Thus, we successfully implemented various vector & matrix operations parallelly using CUDA