

Assignment HAC-2

Title: Vector and operations using CUDA

Problem Statement

Design parallel algorithm to

- 1.) Add two large vectors
- 2.) Multiply vector and matrix
- 3.) Multiply two $N \times N$ arrays n^2 processes

Objectives

To learn CUDA architecture and programming language

Outcomes

To be able to learn CUDA architecture and programming concepts

Software and hardware requirements

OS = Fedora / Windows 10 (64-bit)

CUDA API, NVCC compiler, NVIDIA GPU / Google Colab

RAM: 4 GB, HDD: 500 GB

Concepts related theory

CUDA architecture

The architecture consists of several components like

- 1) Parallel compute engines inside NVIDIA GPUs
- 2) OS - kernel level support
- 3) user - mode driver providing device level - API
- 4) Pix instruction set architecture for parallel computing kernels and functions

CUDA memory hierarchy

- Each thread has private local memory
- Each Thread block has shared memory visible to all threads in block
- All threads have access to global memory.

Compilation with NVCC - compiler driver to simply process
embedding of C++ program.

Algorithms

1. Addition

In addition of vectors, add i^{th} element from first array to i^{th} of second. Each array addition can be done in different threads. Cases -

i) n blocks and 1 thread per block

$$\rightarrow id = block \mid dx.x$$

$$cur \ll n, 1 \gg (1);$$

ii) 1 block and n threads per block

$$\rightarrow id = thread \mid dx.x$$

$$cur \ll \ll 1, n \gg \gg (1);$$

iii) n blocks and n threads per block

$$id = block \mid dx.x \oplus block \mid dx.x + thread \mid dx.x$$

2. Multiplication

i) 2D blocks and one thread per block

$$\text{Here, } x = block \mid dx.x$$

$$y = block \mid dx.y$$

grid dimensions \Rightarrow dim 3 grid (col 2, row 1)

Function kernel \Rightarrow matproduct $\ll grid, 1 \gg (1, n, n);$

$$\text{Speedup} = \frac{\text{time serial}}{\text{time parallel}}$$

Matrix linearisation in CUDA

row = block₁ dx.y * block Dim.y + thread₁ dx.y

col = block₁ dx.x * block Dim.x + thread₁ dx.x

offset = row * N + column

Test Cases

	Function	Size	Serial time (ms)	Parallel time (ms)	Speed up
1.	Vector addition	a) 10000	0.0475	0.059	0.8
		b) 50000	0.0195	0.145	1.34
2.	Vector-matrix multiplication	a) 1x70 and 100x100	0.059	0.053	1.11
3.	Matrix-matrix multiply	a) N=10	0.006	0.022	0.267
		b) N=50	0.057	0.074	0.82

Conclusion

Thus we successfully implemented vector and matrix operations using parallel computing and CUDA. Analysed the speedup for various operations. For small arrays speedup is less than 1, but as the size increases parallel computing becomes less costly.