Assignment 1 Report: Parallel Vector Addition with PyCUDA and PyOpenCL

Weijie Wang
ww2739@columbia.edu

https://forums.developer.nvidia.com/t/info-finished-with-code-256-error-install-of-driver-component-failed/107661
https://medium.com/@abhishekjainindore24/gpus-vs-cpus-threads-core-speed-75b5732ce389

In this assignment, the introduction to PyCUDA and PyOpenCL was introduced, and the theory of the Parallel Vector Addition was learned. Compare and contrast different methods of host-to-device memory allocation for simple elementwise operations on vectors. Also, profiling is introduced.

(As I mentioned in the Ed discussion board, I suddenly lost access to the CUDA toolkit and thus lost the plots or profiling files. I only found a PyOpenCL output screenshot in my folder. https://edstem.org/us/courses/66062/discussion/5372617 )

The output of the PyCuda.py is attached below.

Figure 1 PyCUDA code output

2.
The reason why no execution time without memory transfer is that we can only measure the total execution time because the data transfers and kernel execution are inseparable in this method.

4.
Compare the results with CPU_numpy_Add.

Plot of the average execution times, including memory transfer for GPU operations, against vector size
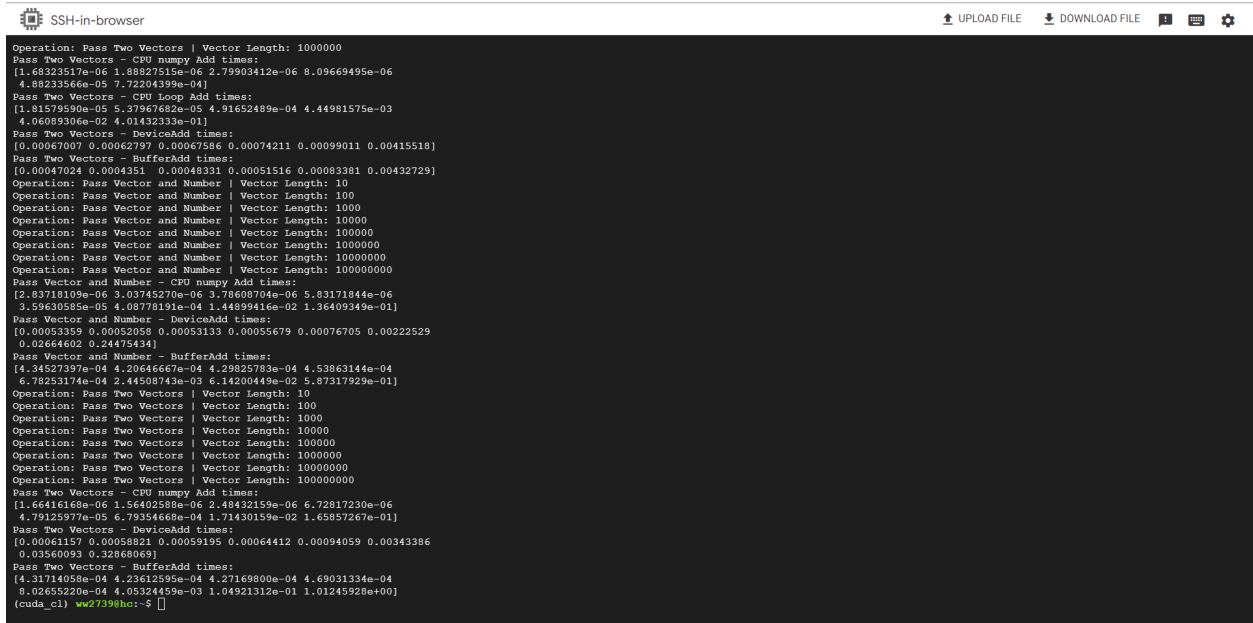
Plot of the average execution times, excluding memory transfer for GPU operations, against vector size

9. Profiling and explain.
CLI

Compute UI

The output of the PyOpenCL.py is attached below.



Figure 2 PyOpenCL code output

The average execution times for the four GPU cases and the faster CPU case against vector size

Theory

In a CPU, a thread is the smallest unit that shares memory and resources of the process which enable parallelism in multi-core CPUs. In GPUs, thread is a register per work item which is ideally low and a nice divisor of the number of hardware registers per computer unit. A task is a unit of work that is scheduled for execution. In a CPU, each core is capable of handling complex tasks independently, with a focus on general-purpose computing. In a GPU, it comprises thousands of smaller cores designed to handle multiple tasks simultaneously. A process is an independent execution unit that has its own memory space. In a CPU, processes rely on context switching for parallelism while in a GPU, a process can spawn multiple threads to run different computations in parallel.

The CPU_numpy_Add method is proved to be faster for CPU methods. It comes with a smaller processing time because it takes advantage of NumPy's optimized, low-level implementations for array operations which maximize the performance of operations. Meanwhile, CPU_loop_Add is slower due to the inefficiency inherent in executing explicit loops.

These parallel approaches in PyCUDA and PyOpenCL can perform thousands of operations simultaneously while the CPU methods require multi-threading and multi-processing, which limit the number of operations that happen simultaneously. All in all, parallel approaches in PyCUDA and PyOpenCL bring more computations at the same time and are more efficient in dealing with large workloads.

In PyOpenCL, the deviceAdd method is faster as indicated in the previous results output. It takes advantage by making all operations performed on the GPU, while host-to-device and device-to-host memory transfers are much slower than directly on the device.

In PyCUDA, the add_gpuarrary_no_kernel is the fastest as indicated in the previous results. The reason it is faster than the other methods is that it avoids the overhead of custom kernel compilation and launches, utilizes efficient memory management provided by gpuarray, and has fewer layers of abstraction.

I would prefer PyCUDA. Since PyCUDA targets NVIDIA GPUs and we're using NVIDIA T4 GPUs, it may be more compatible with the platform while we're coding. PyCUDA is optimized for NVIDIA GPUs which brings high performance for specific tasks. PyCUDA also seems to be easier to ease than PyOpenCL with easier access to highly optimized libraries.