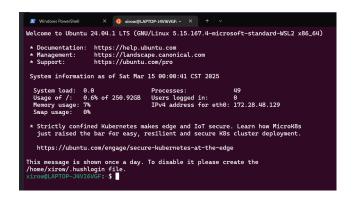# A report about assessment of GOODLAB

Liwei Huang

## I. TASK ONE:

**Index Terms**—Config & Kernel Development

### A. Environment config:

*1) :* Use WSL2 to deploy Ubuntu

One of the advantages of using WSL2 is The Ubuntu could <u>connect directly the graphic card</u> of ur Windows host!



*2) :* Deploy CUDA Toolkit and Pytorch

search the method of installing CUDA and PyTorch in Ubuntu and configure the environment.

Firstly, employing the CUDA:



Secondly, configuring the conda environment. I decide to use it to carry my python and pytorch environment for its advantages: 1.isolating dependencies 2.simplifying manage of dependencies 3.Supporting Cross-platform



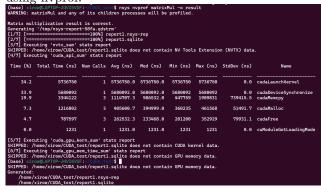Then installing the PyTorch Frame:

Detecting GPU and the computing of cuda!



*3) :* GPU computing power

Verify the availability of the GPU computing power by using nvprof.



### B. Exploit CUDA kernel:

*1) :* To pave the way for exploiting kernel

Firstly
CUDA program:
I had a general understanding of cuda programming through study a course of it, which made me have a moderate command of cuda, such as: kernel function,memory model of cuda,thread-block-warp, resource allocation,error check...

Secondly
LayerNorm
LayerNorm are always broken apart into two parts which is normalization and then by layers:
When training for neural net,we often use LayNorm to avoid gradient disappearance or gradient explosion. LayerNorm could encapsulate those values produced by neurons within a much small range,typically centered around zero. What this allow for is much more stable training as we actually perform a gradient step!
And every neuron in every layer is normalized such that all the activation will have a center like zero or standard deviation of one!
The following is the formula to achieve LayerNorm:

$$\hat{x_i} = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$\epsilon$ is a small value to avoid zero appearing in denominator simplifying:

$$\hat{x_i} = \frac{x_i - \mu}{\sigma}$$

$$\mu = \frac{1}{n}\sum_{i=1}^{n} x_i \qquad \sigma = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(x_i - \mu)^2}$$

*2) :* process of exploiting kernel

Firstly
Using vector float4 to access memory:
Through reading four floats one times to reduce times of accessing memory so that it can improve efficiency. It also takes advantage of the power of parallel computing better!

Secondly
Using Warp Reduce:
It is such an efficient optimizing technology which can reduce the computing complexity.Because a cuda core computes in parallel in warps which contain 32 threads.So using some instructions to make the thread execute computing tasks simultaneously.In that case,We can apply Warp Reduce to every adding calculation!

Thirdly
Supporting Dynamic shapes automatic adaptation:

Calculating the size of grid and block dynamically. Parameterize the kernel using templates or macros. Dynamic Parallelism.

Fourthly
Realizing Shared-memory Double Buffering system to avoid Bank Conflict:
The core idea of it is using two shared memory areas,one of which is used to read another is used to write. So that when a buffer is being used, another can load and store to improve efficiency!

# The whole code is on my github repository!

## II. TASK TWO:

*Index Terms*—Compare &
**Analyze the performance of operator**

### A. Calculating TFLOPS:

*1) :* definition:
"Tera Floating Point Operations Per Second"
TFLOPS is a vital index of assessing the performance of computing device when disposing complex calculating tasks, especially in deep learning field where vast calculations happen.

### B. Global Load/Store Efficiency:

*1) :* definition:
Global Load/Store Efficiency measures how efficiently a CUDA kernel accesses global memory. It is an important metric in performance tuning because global memory accesses are significantly slower than shared memory or registers

### C. Register/shared memory usage optimization:

*1) :* definition:
Registers are the fastest memory on a GPU, located inside each Streaming Multiprocessor (SM).
Shared memory is faster than global memory and shared among threads in a block.

### D. Acceleration Ratio of different scales:

*1) :* definition:
The Acceleration Ratio depends on various factors, including the specific application, workload characteristics, and how effectively the software utilizes the GPU's architecture.

# III. PROBLEMS & SOLUTIONS

*A.*

**F**IRST of all,
I had no idea about all of the terms in the Task.pdf,which really troubled me! However,with my determination and ambition for new knowledge,I decided to get a general understanding of these new terms firstly to have a overview of my following tasks. Then I used search engine to accomplish above-mentioned things.

*B. TaskOne*

PB1: Before the kernel development, I prefered to have a rough command of CUDA and programming on it through bilibili and csdn.

PB2: However,when I decided to use the "nvprof" instruction,I finded out that it was replaced by some new tools NVIDIA Nsight Systems and NVIDIA Nsight Compute. When I struggled to install and configure NVIDIA Nsight Systems and Compute, its all-English UI really troubled me. So I have to learn how to use this tool to verify the availability of GPU computing power.
It will take some time,but I am interested in it!
Solution to PB2: When I read the help document of Nsys ,I was glad to find out that the instruction "nvprof" could still use by using nsys nvprof!

PB3: when I used the Nsight System to analyze the cuda program,it will always report an error when I use nsys-ui which confused me a lot!!
Solution to PB3: Through reading documents and referring information,I realized that error could be ignored which would not impact my analysis on Nsight!

PB4 and Solution: optimizing the cuda kernel was so difficult which made me have to turn to AI and numerous videos to understand them step by step statement by statement.And I could not take a good command of them within a short days!

*C. TaskTwo*

PB1 and Solution: I haven't learn the proper nouns of operator performance, so I need to take a general comprehension of them.

PB2: I couldn't find the TFLOPS data and some relevant data on Nsight Compute although I had used Nsight to analyze my CUDA kernel program.

PB3: I tried to use python to write test code that contains native PyTorch LayerNorm and custom CUDA kernel But some unknown and difficult errors really confused me, and I think all of those confusion and errors couldn't be resolved in a short period.