

CS259-L: Performance Comparison of MapReduce and CUDA on the HG19 Human Genome Dataset

Chung-Hsuan Wu
Computer Science Department
San Jose State University
San Jose, CA 95192
408-924-1000

chung-hsuan.wu@sjsu.edu

ABSTRACT

Although both MapReduce and CUDA are efficient techniques for parallel processing, they differ fundamentally in design and application. MapReduce emphasizes distributed data processing through Mapper and Reducer operations, while CUDA exploits massive thread parallelism on GPUs. This project compares CPU-based MapReduce and GPU-based CUDA by implementing a k-mer counting task on a human gene dataset. Performance will be evaluated in terms of execution time, scalability, and resource utilization. This project will also focus on analyzing the factors driving the observed differences, and discussing the types of bioinformatics and data-intensive workloads best suited for each framework.

1. INTRODUCTION

Human genome analysis has been one of the most important areas in bioinformatic research. K-mer counting is one fundamental task in genome assembly, sequence comparison, taxonomic classification, and error estimation. However, as the size of human genomic datasets continues to grow, traditional methods based on a hash map will lead to memory exhaustion. Single-core processors are insufficient for such tasks, and even multi-core systems on a single machine cannot meet the computational demand. Therefore, distributed and parallel computing systems such as Hadoop or CUDA are becoming more and more popular in this field. Due to the essential characteristic in K-mer counting, MapReduce becomes extremely suitable for this task. This project will first use Hadoop to establish a cluster system for MapReduce to evaluate its efficiency and scalability. The same algorithm is also tested on a single CPU as a baseline. Finally, NVIDIA CUDA will be utilized to perform the same task, allowing for a direct comparison of performance and hardware utilization.

2. DATASET

The *hg19* human genome assembly is used in this project. *Hg19* is an older version of the human genome reference assembly with a total size of approximately 3.8 GB. There are only five characters in this dataset, which are A, C, T, G, and N. Although it is not widely used in recent bioinformatic research, it is still a foundational reference genome used for identifying genes, genetic variations, and understanding diseases.

3. METHODOLOGY

This project will use the MapReduce framework to perform the K-mer counting task. There are two types of phases in the

MapReduce framework: Mapper and Reducer. The Mapper phase is for generating intermediate key-value pairs. The result for Mapper is, for example, ("GAT: 1") when K is 3. Notably, the value from Mapper will always be one, even if the keys are the same. The next step is shuffling, where all records with the same keys are grouped together. The final phase is Reducer, which is to group key-value pairs by summing the count for each key. This inherent property makes MapReduce suitable for execution on different machines or threads. In this project, two different architectures will be implemented and tested to analyze the efficiencies.

- Hadoop

Hadoop is a cluster system framework that is highly suitable for distributed processing. Hadoop also supports fault tolerance and is commonly used in big data analysis.

- CUDA

CUDA is a parallel computing platform and programming model developed by NVIDIA. It allows the use of an NVIDIA GPU for general-purpose processing. Although CUDA supports C and C++ programming, some advanced data structures, such as sets and maps, are not supported in GPU computation; only basic data structures, including pointers and arrays, can be used in CUDA.

Since a map cannot be used in this project, an array will be used to simulate a map to store the key-value pairs. The size of the array relies on the number of k. In this study, the value of k is set to 3, which means only three bits will be considered each time. There are only five choices for each digit, which are A, C, T, G, and N. The size of the array becomes 5 to the power of k, meaning the size becomes $5^3 = 125$ for this project.

In CUDA, a thread is the basic unit for parallel processing. Combining multiple threads becomes a block. Similarly, the blocks of parallel threads make up a grid. The following figure shows the relations between thread, block, and grid.

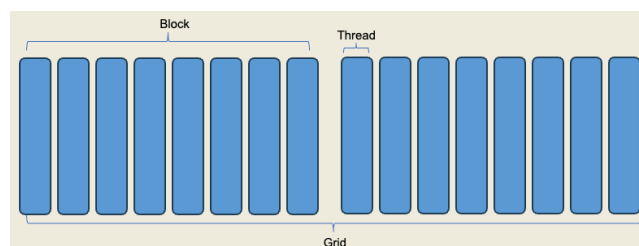


Figure 1. The relations between thread, block, and grid.

CUDA introduces a special feature known as unified memory. Unified memory provides a single virtual address space accessible by both CPU and GPU, meaning no consecutive transferring needs to be done during the process. Figure 2 shows how unified memory works.

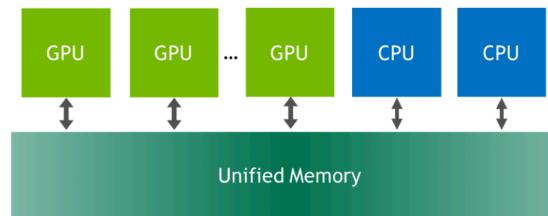


Figure 2. Unified Memory

4. HARDWARE CONFIGURATION

The following table describe the hardware configurations in this study.

Table 1. hardware configurations

CPU	i7-14700F
GPU	4060 8GB
RAM	16GB
System	Ubuntu 24.04

5. RESULTS

5.1 Hadoop

Figure 3 shows the results for Hadoop.

Show 23 of 23 rows																	
ID	User	Name	Application Type	Application Tags	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU Vcores	Allocated Memory MB	Allocated GPUs	Reserved CPU Vcores	Reserved Memory MB
application_17058429984_0003	hadoop	elasticsearch-20240708000434.jar	MAPREDUCE		not default	0	Fri Jul 21 16:37:57 16:37:57	Fri Jul 21 16:37:57 16:37:57	Fri Jul 21 17:22:58 17:22:58	FINISHED	SUCCEEDED	N/A	N/A	N/A	N/A	N/A	N/A

Figure 3. Results for Hadoop

Performance analysis focuses on the start time and finish time in this study. It took approximately 30 minutes to complete the MapReduce task. There are multiple reasons for it. First, both Mapper and Reducer are using standard output (STDOUT). However, it takes a long time to perform the printing command. Secondly, for Mapper, the intermediate results will be written to disk. As a result, Reducer must get the key-value pairs from disk. Disk reading and writing take a huge amount of time compared to RAM or cache. This will slow down the whole process. Besides, the RAM size in this project is set to 8GB. Page faults might occur and therefore hinder the process. Notably, in general, those different machines for Mappers and Reducers are usually connected through the internet, making the bandwidth of the internet a concern. In this study, only one machine is used in the Hadoop system, making this issue able to be ignored.

5.2 Sequential C++

Figure 4 shows the compiling and executing process, and the analysis for a normal C++ program.

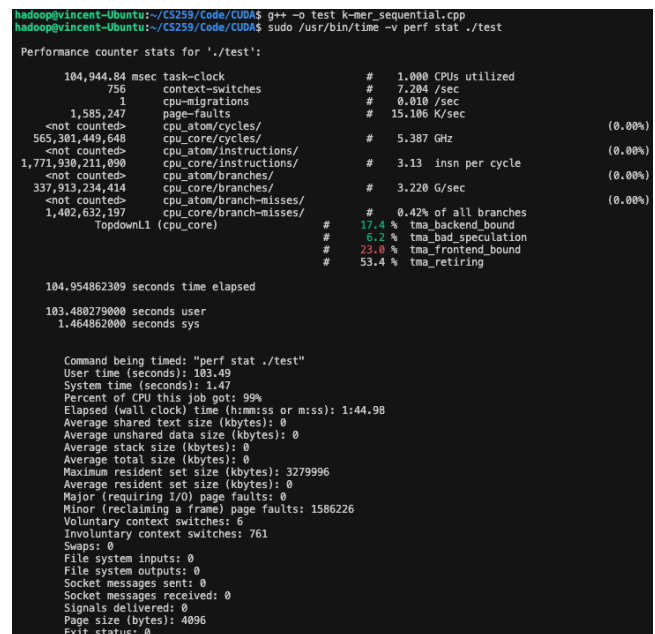


Figure 4. Analysis for normal C++ program.

For CPU, it shows the program took approximately 104 seconds to complete the task, with only 1 CPU, and merely 1 second is used for system calls. This program is running 3 instructions per cycle for pipelining, and almost achieves 99.5% of successful branch predictions. For RAM, there are 1,585,247 page faults in this program. There are two types of page faults: Major and Minor. A major page fault means the data is not in RAM. Minor page fault represents that the data is not in the process page table. Since in this program, the array is initialized to empty, minor page faults are unavoidable.

5.3 CUDA

5.3.1 Single Block

Figure 5 shows the analysis for single block with 256 threads.

Time (ns)	Total Time (ns)	Num Calls	Avg (ns)	Med (ns)	Min (ns)	Max (ns)	Stddev (ns)	Name
76.1	268,979,872	61	4,397,345.4	4,347,544.0	9,976	4,455,427	63,875.1	cudaMemcpy
21.5	7,787,685	36	36,853,352.5	36,853,352.5	29,957	73,681,688	52,800,212.6	cudaMalloc
0.0	6,189,080	60	114,784.0	116,851.3	2,745	124,719	14,825.6	cudaMemcpySynchronize
0.0	1,842,779,176	32	57,586,867.5	57,586,867.5	92,968	114,754,784	92,968,000.0	cudaMemcpy
0.1	323,655	60	5,394.3	2,848.5	2,256	15,120	19,080.1	cudaLaunchCntrl
10.0	7,787,685	36	214,784.0	214,784.0	172,968	1,132,128	1,132,128.0	cudaMemcpy
0.0	7,738	1	7,738.0	7,738.0	7,738	7,738	0.000.1	cudaMemset
0.0	4,448	1	4,448.0	4,448.0	4,448	4,448	2,722.2	cudaEventWait
0.0	1,480	1	1,480.0	1,480.0	1,480	1,480	0.000.1	cudaEventWait
0.0	966	2	483.0	483.0	176	790	434.2	cudaEventRecord
0.0	966	2	483.0	483.0	176	790	434.2	cudaEventRecordSynchronize

Figure 5. Analysis for single block in CUDA.

The longest time in a single block with CUDA is the memory copying, which is to transfer data into unified memory. “cudaDeviceSynchronize” is where the CPU waits for GPUs to finish their tasks. In a single block, the total amount of time for GPU work is 6,887,085 ns, and is approximately 2% of the total amount of time.

5.3.2 Multiple Blocks

Figure 6 shows the analysis for multiple blocks in CUDA.

Time (ns)	Total Time (ns)	Num Calls	Avg (ns)	Med (ns)	Min (ns)	Max (ns)	StdDev (ns)	Name
41.4	318,714,948	61	5,313,915.5	5,259,795.3	5,454,854.6	7,648,677	617,868.6	cudaDeviceSynchronize
39.7	768,133,388	61	4,766,293.2	4,723,658.5	4,766,293.2	5,454,854.6	545,454.5	cudaDeviceSynchronize
118.4	74,371,866	2	37,185,935.3	37,185,935.3	33,274	73,339,592	52,147,994.2	cudaMalloc
0.0	47,473,658	1	47,473,658.5	47,473,658.5	47,473,658.5	54,545,454.5	54,545,454.5	cudaMalloc
0.0	319,214	2	5,178.2	5,178.2	2,358	16,376	18,544.0	cudaLaunchKernel
0.0	1,517.4	2	5,178.2	5,178.2	0	9,356	1,668.0	cudaEventRecord
0.0	7,859	2	3,929.5	3,929.5	7,859	7,859	0	cudaEventRecord
0.0	4,939	2	2,469.5	2,469.5	5	4,924	3,846.0	cudaEventRecord
0.0	1,673	2	836.5	836.5	1,673	1,673	0	cudaDeviceSynchronize
0.0	1,094	2	547.0	547.0	284	890	485.1	cudaEventDestroy
0.0	1,094	2	547.0	547.0	284	890	485.1	cudaEventDestroy

Figure 6. Multiples blocks in CUDA

For multiple blocks, the time GPUs actually work is 318,714,948 ns, and is 48.7% of the total amount of time. Surprisingly, the time increases when multiple blocks are used. This is due to a race condition. Although more threads are used to perform the task, they might intend to write the same key-value pair in the array, which means one of them should idle while the other is writing data. The following figure shows an example of when a race condition occurs.

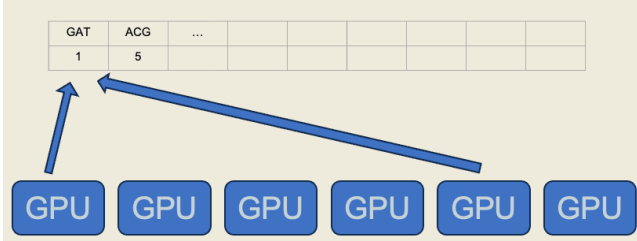


Figure 7. Race Condition

Since more workers will not necessarily speed up the task, this project aims to solve the latency for memory copying.

5.3.3 Pipeline

Two different approaches can solve the memory copying issue: prefetching and pipelining. Prefetching can be used when the size that needs to be allocated is known in advance. In this study, the HG-19 dataset is too huge to be entirely put in the unified memory, as well as the array with size 125. The figure below shows an example of prefetching.

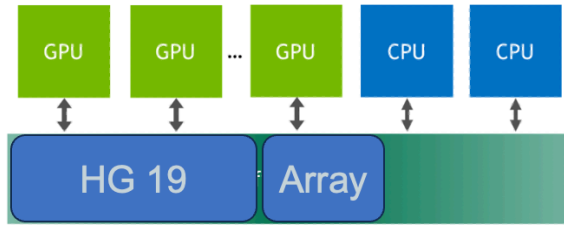


Figure 8. Prefetching

As a result, splitting data into different pieces is necessary. In this project, the data is split into 50 MB. Therefore, repetitively moving those chunks of data needs to be performed. This process is completed by the CPU. Figure 9 shows this process.

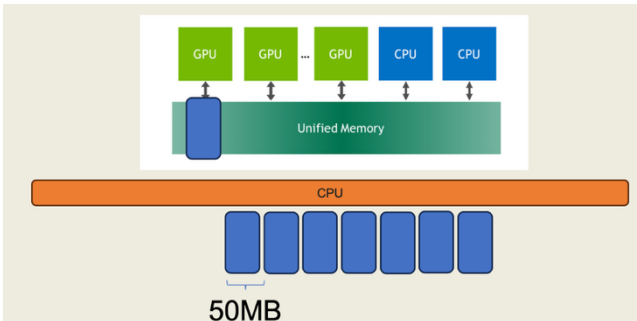


Figure 9. Splitting data

Since the transfer is executed by the CPU, GPUs can compute those chunks of data simultaneously. The figures below show the workflow with and without pipelining.

CPU	MOVE DATA	MOVE DATA	MOVE DATA
GPU	COMPUTE	COMPUTE	COMPUTE

Figure 10. workflow without pipelining.

CPU	MOVE DATA	MOVE DATA	MOVE DATA	MOVE DATA	MOVE DATA	MOVE DATA
GPU	COMPUTE	COMPUTE	COMPUTE	COMPUTE	COMPUTE	COMPUTE

Figure 11. workflow with pipelining.

5.3.4 Pipeline with single block

Figure 12 shows the analysis with single block utilizing pipelining.

Time (%)	Total Time (ns)	Num Calls	Avg (ns)	Med (ns)	Min (ns)	Max (ns)	StdDev (ns)	Name
55.1	125,476,237	1	125,476,237.0	125,476,237.0	125,476,237	125,476,237	0.0	cudaDeviceSynchronize
32.1	73,135,584	3	24,378,528.0	44,577.0	34,577	73,856,438	42,156,306.0	cudaMalloc
7.7	17,433,520	2	8,716,760.0	8,716,760.0	8,631,470	8,882,158	128,680.0	cudaMalloc
2.7	6,047,320	2	3,023,660.0	3,023,660.0	2,688,385	3,358,435	473,415.1	cudaFreeHost
2.0	4,479,745	0	74,554.1	2,906.5	2,320	4,236,495	525,222.2	cudaLaunchKernel
0.4	865,185	3	321,781.7	51,648.0	38,762	875,195	479,368.0	cudaFree
0.1	215,559	60	3,582.7	3,394.0	2,830	12,138	1,213.8	cudaMemcpyAsync
0.0	15,089	2	9,384.5	9,384.5	2,862	15,107	9,130.4	cudaStreamCreate
0.0	13,781	2	6,898.5	6,898.5	5,919	7,862	1,373.9	cudaEventRecord
0.0	12,371	1	12,371.0	12,371.0	12,371	12,371	0.0	cudaMemcpy
0.0	8,822	1	8,822.0	8,822.0	8,822	8,822	0.0	cudaMemset
0.0	5,318	2	2,659.0	2,659.0	1,548	3,778	1,571.2	cudaStreamDestroy
0.0	3,989	2	1,994.5	1,994.5	235	3,474	2,431.7	cudaEventCreate
0.0	1,972	1	1,972.0	1,972.0	1,972	1,972	0.0	cudaEventSynchronize
0.0	1,056	2	528.0	528.0	288	888	357.7	cudaEventDestroy
0.0	548	1	548.0	548.0	548	548	0.0	cudaModuleLoadMode

Figure 12. Single block with pipelining.

The longest time in this approach is “cudaDeviceSynchronize”, which means GPUs are working. Compared to the non-pipeline version of single block, the total amount of time decreased, and the percentage of time increased as well. Pipelining significantly makes the whole process much more efficient.

5.3.5 Pipeline with multiple blocks

The figure below shows the total time needed for each process using multiple blocks with pipelining.

Time (%)	Total Time (ns)	Num Calls	Avg (ns)	Med (ns)	Min (ns)	Max (ns)	StdDev (ns)	Name
66.6	219,588,414	1	219,588,414.0	219,588,414	219,588,414	219,588,414	0.0	cudaDeviceSynchronize
24.5	88,889,178	3	29,629,392.7	48,497.0	34,752	80,793,929	46,624,674.2	cudaMalloc
5.5	17,522,426	2	8,761,212.5	8,761,212.5	8,687,273	8,834,852	89,909.4	cudaMalloc
1.0	5,953,828	2	2,976,918.0	2,976,918.0	2,588,316	3,365,584	549,554.9	cudaFreeHost
1.5	4,465,882	0	74,438.0	2,845.0	2,245	4,234,153	553,997.4	cudaLaunchKernel
0.3	943,518	3	314,506.0	39,999.0	30,489	972,438	465,160.1	cudaFree
0.1	215,470	60	3,657.8	3,487.5	2,888	11,369	1,164.3	cudaMemcpyAsync
0.0	15,369	2	9,784.5	9,784.5	2,633	16,726	9,338.8	cudaStreamCreate
0.0	13,930	2	6,965.0	6,965.0	4,783	9,147	3,885.8	cudaEventRecord
0.0	12,316	1	12,316.0	12,316.0	12,316	12,316	0.0	cudaMemcpy
0.0	8,377	1	8,377.0	8,377.0	8,377	8,377	0.0	cudaMemset
0.0	4,321	2	2,160.5	2,160.5	1,422	2,980	1,821.5	cudaStreamDestroy
0.0	3,914	2	1,957.0	1,957.0	283	3,631	2,387.4	cudaEventCreate
0.0	2,377	1	2,377.0	2,377.0	2,377	2,377	0.0	cudaEventSynchronize
0.0	815	2	407.5	407.5	283	612	289.2	cudaEventDestroy
0.0	581	1	581.0	581.0	581	581	0.0	cudaModuleLoadMode

Figure X. Multiple blocks with pipelining.

With pipelining, multiple blocks can still achieve an efficient computing process. The total time needed decreased approximately 1/3 while the percentage of executing time increased.

These evaluations both show that the pipeline can obtain a more efficient computing approach.

6. CONCLUSION

In this study, Hadoop and CUDA are examined to determine which approach can outperform the other in the K-mer counting task using the HG19 human genome assembly dataset. Hadoop can achieve a more efficient computation with multiple machines when it comes to big data analysis, while on a single machine, CUDA can outperform Hadoop by fully utilizing GPUs. However, due to the limited RAM size in CUDA, data splitting may be necessary, meaning the transfer between RAM and VRAM in GPUs needs to be implemented. This process can be performed by the CPU, while GPUs can compute the results simultaneously. Pipelining can get a higher throughput. Nonetheless, more threads that can perform more parallel processing do not necessarily speed up the process due to a race condition issue. In future studies, the divide and conquer method can be examined to figure out if there is a way to avoid a race condition.

7. REFERENCES

- [1] Mark Harris, “An Even Easier Introduction to CUDA, ” 2025. [Online] Available: <https://developer.nvidia.com/blog/even-easier-introduction-cuda/>
- [2] Mark Harris, “Unified Memory, ” 2013 [Online] Available: <https://developer.nvidia.com/blog/unified-memory-in-cuda-6/>
- [3] “MapReduce: Simplified Data Processing on Large Clusters,” research.google. <https://research.google/pubs/mapreduce-simplified-data-processing-on-large-clusters/>
- [4] Apache Hadoop, “MapReduce Tutorial, ” 2025 [Online] Available: <https://hadoop.apache.org/docs/stable/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html>