

# Big Data Coursework Report

Name: Thanawat Thanaponpaiboon

Email: [Guy.Thanaponpaiboon@city.ac.uk](mailto:Guy.Thanaponpaiboon@city.ac.uk)

## 1d) Optimisation, experiments, and discussion (17%)

i) Improve parallelisation

```
firstRDD = sc.parallelize(filename, 16) ### TASK 1d ### (default setting = 2)
```

In order to improve parallelisation, we adjust the **second parameter** from the default (2) to (16) according to the example from section 1.2. This means we increase the amount of data into 16 partition and processes each partition in parallel across the nodes in the google cloud cluster.

- Two nodes (parallelize = 2)

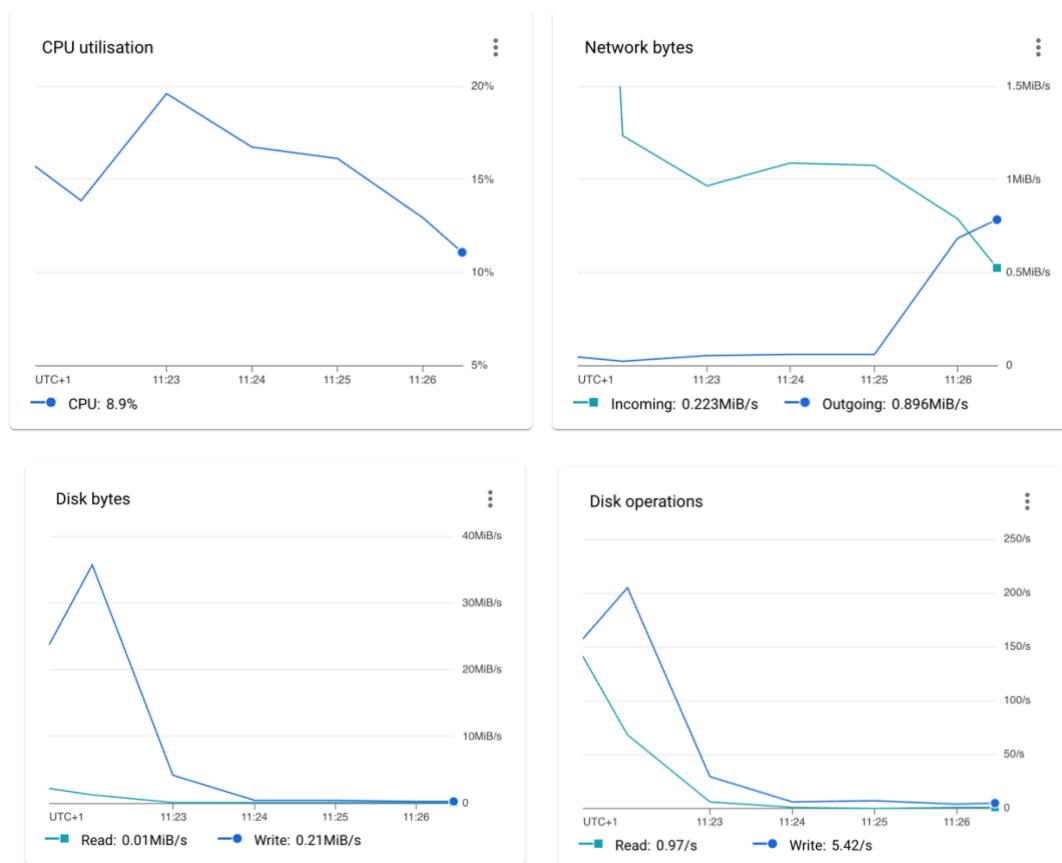
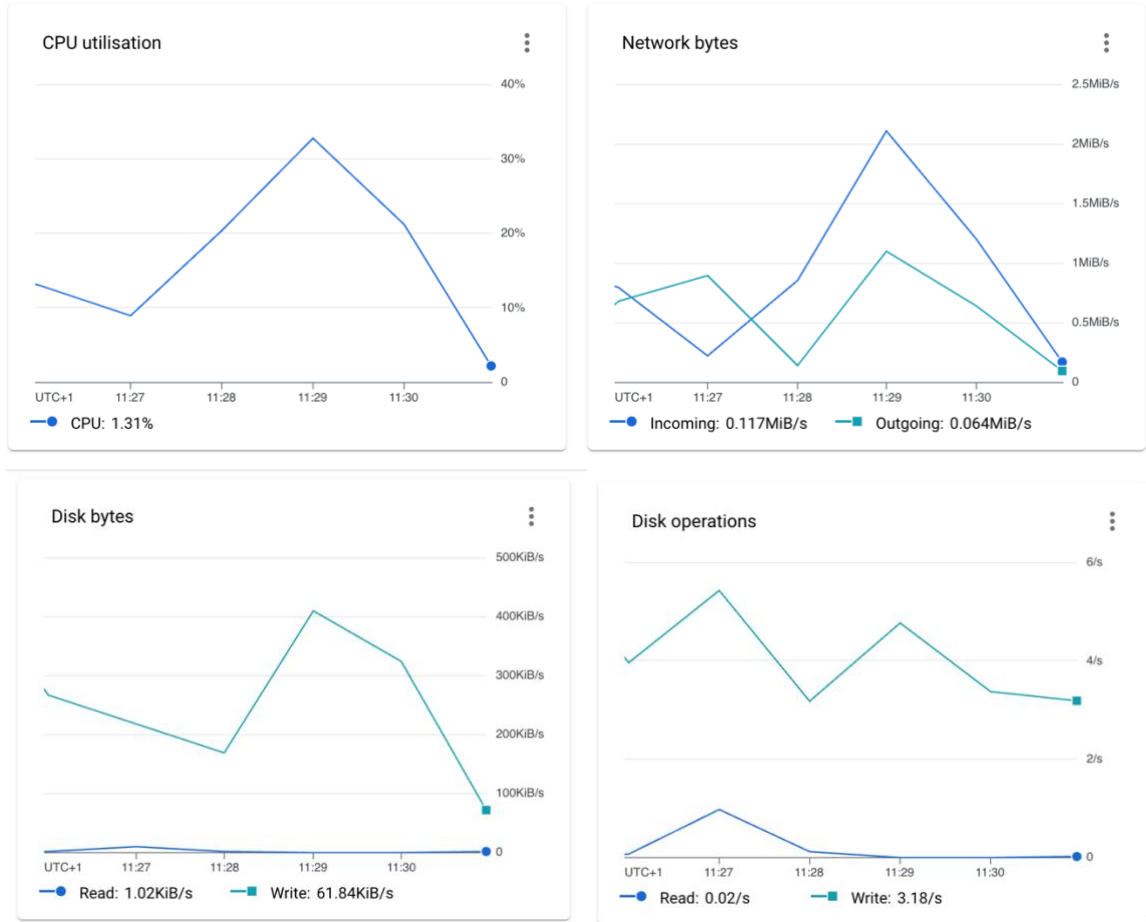


Figure 1: Single machine maximal SSD size (100) + 8 vCPUs. (2 partition)

- **Sixteen nodes (parallelize = 16)**



**Figure 2:** Single machine maximal SSD size (100) + 8 vCPUs. (16 partition)

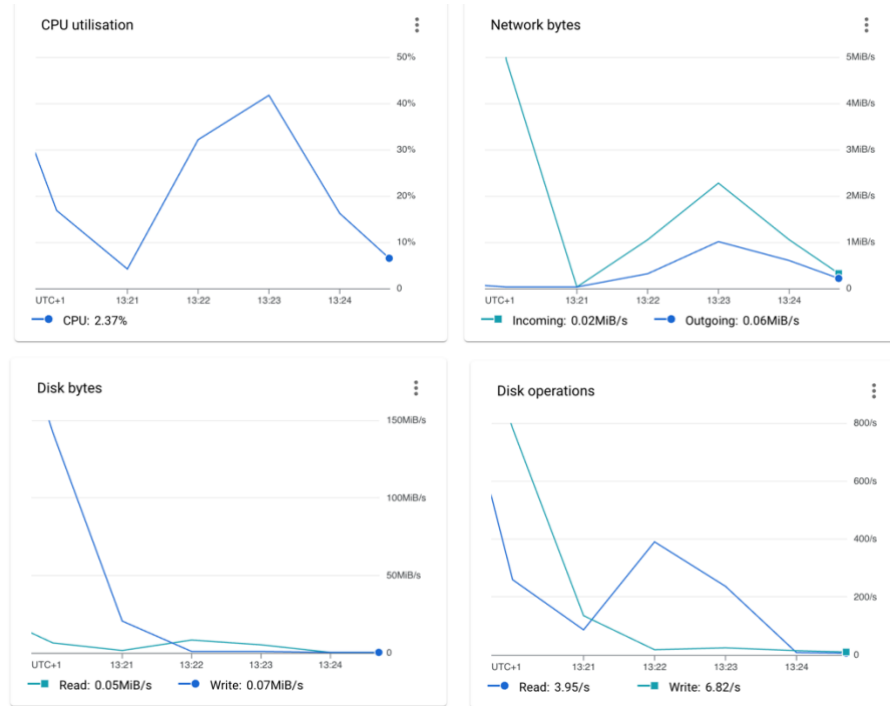
In general, when the second parameter was added from the default (2) to (16) the CPU utilisation went up from 19.63% to 32.71% at peak, increasing the number of multithreading may result in high CPU usages or increased CPU cycles. Although this costs some CPU time upfront, but network and disk bandwidth are saved. We can see from both figures, the lower number of disk bytes when we increase the number of partitions to 16 indicates that the disk is being used slightly compared to executing only 2 nodes.

Start time	Elapsed time	Labels
26 Apr 2023, 11:27:17	2 min 20 sec	nodes : 16
26 Apr 2023, 11:21:28	4 min 59 sec	nodes : 2

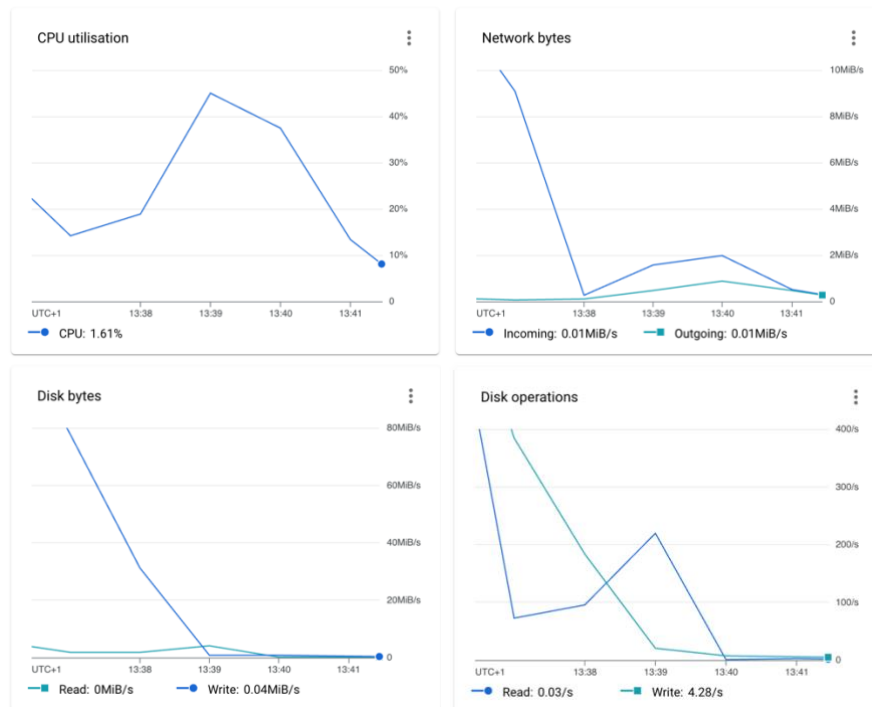
In terms of the **difference in the processing time**, increase the number of partitions from 2 to 16 increase the speed of computational time from around 4 minutes 59 seconds to 2 minute and 20 seconds. Hence, increasing the number of partitions improves processing speed by 2.1 times.

ii) Experiment with cluster configurations.

**The largest possible cluster** that we can create consists of 1 master and 7 worker nodes. Each of them with 1 (virtual) CPU. The master node gets a full 100GB capacity of SSD persistent disk, and the 7 worker nodes get equal shares of the standard disk capacity to maximize throughput which means 2000GB of standard persistent disk need to be divided by 7, 285GB will be the number of standard disk capacity per node approximately.



**Figure 3:** Largest possible cluster 8 VMs (1 master +7 worker nodes) (16 partition)



**Figure 4:** A machine with eightfold resources + 3 machines (2 vCPUs, memory, disk) (16 partition)

26 Apr 2023, 13:37:25	3 min 3 sec	machine : 4
26 Apr 2023, 13:20:57	2 min 29 sec	vm : 8

Creating the largest possible cluster can result in high CPU utilization. This can be the same scenario as when we created 4 machines with double virtual CPUs. Both clusters require a high CPU usage which lies around 41%-45% at the highest state. In the same way, a large cluster typically requires more disk I/O and network bandwidth than a smaller cluster, as you can see from the figure that after we create increase the worker machine network and disk also goes up to.

Large clusters can be beneficial for distributed computing such as Spark when we run across multiple nodes in parallel simultaneously. By distributing the workload across multiple nodes. This approach can reduce the computational time during the process and the risk of any hardware or software failures. However, creating a cluster that does not fit into the task can lead to poor performance, slow response times, and also end up wasting resources such as storage, and networking bandwidth. As can be seen from the experiment increasing the worker node in the cluster seem not to have a profound effect in terms of increasing the speed of processing time compared to the single master node. And reducing the number of worker nodes to 3 along with using a couple of Virtual CPUs in the working node produces a similar result as the experimentation with 8 virtual machines.

iii) Explain the difference between this use of Spark and most standard applications

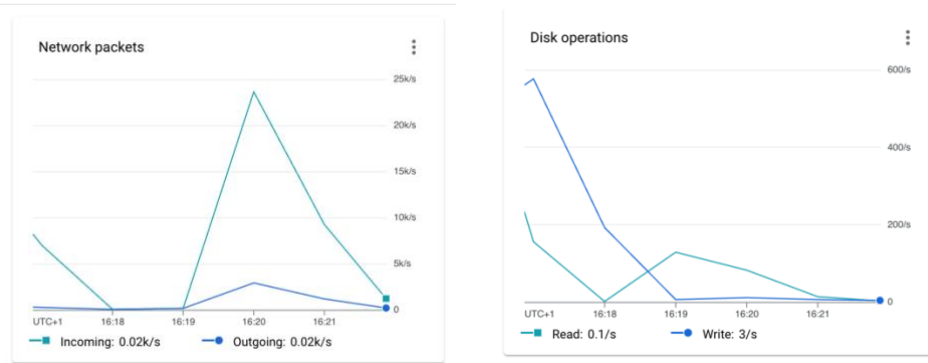
In the lab session, data is usually stored in a single machine like pandas dataframe array, excel database or CSV file etc. It always executes locally in a platform such as Jupiter. In contrast, Spark use the technique called **parallelize** for distribute the data over partitions, which mean that the data is distributed across multiple machines or nodes in a cluster can be use simultaneously. The result is drastically increasing the speed of processing time when perform on large datasets. However, spark distribution needs to rely on extra settings and more external environments like google cloud which might affect the processing time as we mentioned. Poor settings can result in poor performance and wasteful resources. While most of the time standard applications proceed locally without any reinforcement.

## 2b) Testing the code and collecting results (4%)

I) First, test locally with %run: it takes a minute to run through all step locally

```

1m 16:18:00 16:18:01 16:18:02 16:18:03 16:18:04 16:18:05 16:18:06 16:18:07 16:18:08 16:18:09 16:18:10 16:18:11 16:18:12 16:18:13 16:18:14 16:18:15 16:18:16 16:18:17 16:18:18 16:18:19 16:18:20 16:18:21 16:18:22 16:18:23 16:18:24 16:18:25 16:18:26 16:18:27 16:18:28 16:18:29 16:18:30 16:18:31 16:18:32 16:18:33 16:18:34 16:18:35 16:18:36 16:18:37 16:18:38 16:18:39 16:18:40 16:18:41 16:18:42 16:18:43 16:18:44 16:18:45 16:18:46 16:18:47 16:18:48 16:18:49 16:18:50 16:18:51 16:18:52 16:18:53 16:18:54 16:18:55 16:18:56 16:18:57 16:18:58 16:18:59 16:19:00 16:19:01 16:19:02 16:19:03 16:19:04 16:19:05 16:19:06 16:19:07 16:19:08 16:19:09 16:19:10 16:19:11 16:19:12 16:19:13 16:19:14 16:19:15 16:19:16 16:19:17 16:19:18 16:19:19 16:19:20 16:19:21 16:19:22 16:19:23 16:19:24 16:19:25 16:19:26 16:19:27 16:19:28 16:19:29 16:19:30 16:19:31 16:19:32 16:19:33 16:19:34 16:19:35 16:19:36 16:19:37 16:19:38 16:19:39 16:19:40 16:19:41 16:19:42 16:19:43 16:19:44 16:19:45 16:19:46 16:19:47 16:19:48 16:19:49 16:19:50 16:19:51 16:19:52 16:19:53 16:19:54 16:19:55 16:19:56 16:19:57 16:19:58 16:19:59 16:20:00 16:20:01 16:20:02 16:20:03 16:20:04 16:20:05 16:20:06 16:20:07 16:20:08 16:20:09 16:20:10 16:20:11 16:20:12 16:20:13 16:20:14 16:20:15 16:20:16 16:20:17 16:20:18 16:20:19 16:20:20 16:20:21 16:20:22 16:20:23 16:20:24 16:20:25 16:20:26 16:20:27 16:20:28 16:20:29 16:20:30 16:20:31 16:20:32 16:20:33 16:20:34 16:20:35 16:20:36 16:20:37 16:20:38 16:20:39 16:20:40 16:20:41 16:20:42 16:20:43 16:20:44 16:20:45 16:20:46 16:20:47 16:20:48 16:20:49 16:20:50 16:20:51 16:20:52 16:20:53 16:20:54 16:20:55 16:20:56 16:20:57 16:20:58 16:20:59 16:21:00 16:21:01 16:21:02 16:21:03 16:21:04 16:21:05 16:21:06 16:21:07 16:21:08 16:21:09 16:21:10 16:21:11 16:21:12 16:21:13 16:21:14 16:21:15 16:21:16 16:21:17 16:21:18 16:21:19 16:21:20 16:21:21 16:21:22 16:21:23 16:21:24 16:21:25 16:21:26 16:21:27 16:21:28 16:21:29 16:21:30 16:21:31 16:21:32 16:21:33 16:21:34 16:21:35 16:21:36 16:21:37 16:21:38 16:21:39 16:21:40 16:21:41 16:21:42 16:21:43 16:21:44 16:21:45 16:21:46 16:21:47 16:21:48 16:21:49 16:21:50 16:21:51 16:21:52 16:21:53 16:21:54 16:21:55 16:21:56 16:21:57 16:21:58 16:21:59 16:22:00 16:22:01 16:22:02 16:22:03 16:22:04 16:22:05 16:22:06 16:22:07 16:22:08 16:22:09 16:22:10 16:22:11 16:22:12 16:22:13 16:22:14 16:22:15 16:22:16 16:22:17 16:22:18 16:22:19 16:22:20 16:22:21 16:22:22 16:22:23 16:22:24 16:22:25 16:22:26 16:22:27 16:22:28 16:22:29 16:22:30 16:22:31 16:22:32 16:22:33 16:22:34 16:22:35 16:22:36 16:22:37 16:22:38 16:22:39 16:22:40 16:22:41 16:22:42 16:22:43 16:22:44 16:22:45 16:22:46 16:22:47 16:22:48 16:22:49 16:22:50 16:22:51 16:22:52 16:22:53 16:22:54 16:22:55 16:22:56 16:22:57 16:22:58 16:22:59 16:23:00 16:23:01 16:23:02 16:23:03 16:23:04 16:23:05 16:23:06 16:23:07 16:23:08 16:23:09 16:23:10 16:23:11 16:23:12 16:23:13 16:23:14 16:23:15 16:23:16 16:23:17 16:23:18 16:23:19 16:23:20 16:23:21 16:23:22 16:23:23 16:23:24 16:23:25 16:23:26 16:23:27 16:23:28 16:23:29 16:23:30 16:23:31 16:23:32 16:23:33 16:23:34 16:23:35 16:23:36 16:23:37 16:23:38 16:23:39 16:23:40 16:23:41 16:23:42 16:23:43 16:23:44 16:23:45 16:23:46 16:23:47 16:23:48 16:23:49 16:23:50 16:23:51 16:23:52 16:23:53 16:23:54 16:23:55 16:23:56 16:23:57 16:23:58 16:23:59 16:24:00 16:24:01 16:24:02 16:24:03 16:24:04 16:24:05 16:24:06 16:24:07 16:24:08 16:24:09 16:24:10 16:24:11 16:24:12 16:24:13 16:24:14 16:24:15 16:24:16 16:24:17 16:24:18 16:24:19 16:24:20 16:24:21 16:24:22 16:24:23 16:24:24 16:24:25 16:24:26 16:24:27 16:24:28 16:24:29 16:24:30 16:24:31 16:24:32 16:24:33 16:24:34 16:24:35 16:24:36 16:24:37 16:24:38 16:24:39 16:24:40 16:24:41 16:24:42 16:24:43 16:24:44 16:24:45 16:24:46 16:24:47 16:24:48 16:24:49 16:24:50 16:24:51 16:24:52 16:24:53 16:24:54 16:24:55 16:24:56 16:24:57 16:24:58 16:24:59 16:25:00 16:25:01 16:25:02 16:25:03 16:25:04 16:25:05 16:25:06 16:25:07 16:25:08 16:25:09 16:25:10 16:25:11 16:25:12 16:25:13 16:25:14 16:25:15 16:25:16 16:25:17 16:25:18 16:25:19 16:25:20 16:25:21 16:25:22 16:25:23 16:25:24 16:25:25 16:25:26 16:25:27 16:25:28 16:25:29 16:25:30 16:25:31 16:25:32 16:25:33 16:25:34 16:25:35 16:25:36 16:25:37 16:25:38 16:25:39 16:25:40 16:25:41 16:25:42 16:25:43 16:25:44 16:25:45 16:25:46 16:25:47 16:25:48 16:25:49 16:25:50 16:25:51 16:25:52 16:25:53 16:25:54 16:25:55 16:25:56 16:25:57 16:25:58 16:25:59 16:26:00 16:26:01 16:26:02 16:26:03 16:26:04 16:26:05 16:26:06 16:26:07 16:26:08 16:26:09 16:26:10 16:26:11 16:26:12 16:26:13 16:26:14 16:26:15 16:26:16 16:26:17 16:26:18 16:26:19 16:26:20 16:26:21 16:26:22 16:26:23 16:26:24 16:26:25 16:26:26 16:26:27 16:26:28 16:26:29 16:26:30 16:26:31 16:26:32 16:26:33 16:26:34 16:26:35 16:26:36 16:26:37 16:26:38 16:26:39 16:26:40 16:26:41 16:26:42 16:26:43 16:26:44 16:26:45 16:26:46 16:26:47 16:26:48 16:26:49 16:26:50 16:26:51 16:26:52 16:26:53 16:26:54 16:26:55 16:26:56 16:26:57 16:26:58 16:26:59 16:27:00 16:27:01 16:27:02 16:27:03 16:27:04 16:27:05 16:27:06 16:27:07 16:27:08 16:27:09 16:27:10 16:27:11 16:27:12 16:27:13 16:27:14 16:27:15 16:27:16 16:27:17 16:27:18 16:27:19 16:27:20 16:27:21 16:27:22 16:27:23 16:27:24 16:27:25 16:27:26 16:27:27 16:27:28 16:27:29 16:27:30 16:27:31 16:27:32 16:27:33 16:27:34 16:27:35 16:27:36 16:27:37 16:27:38 16:27:39 16:27:40 16:27:41 16:27:42 16:27:43 16:27:44 16:27:45 16:27:46 16:27:47 16:27:48 16:27:49 16:27:50 16:27:51 16:27:52 16:27:53 16:27:54 16:27:55 16:27:56 16:27:57 16:27:58 16:27:59 16:28:00 16:28:01 16:28:02 16:28:03 16:28:04 16:28:05 16:28:06 16:28:07 16:28:08 16:28:09 16:28:10 16:28:11 16:28:12 16:28:13 16:28:14 16:28:15 16:28:16 16:28:17 16:28:18 16:28:19 16:28:20 16:28:21 16:28:22 16:28:23 16:28:24 16:28:25 16:28:26 16:28:27 16:28:28 16:28:29 16:28:30 16:28:31 16:28:32 16:28:33 16:28:34 16:28:35 16:28:36 16:28:37 16:28:38 16:28:39 16:28:40 16:28:41 16:28:42 16:28:43 16:28:44 16:28:45 16:28:46 16:28:47 16:28:48 16:28:49 16:28:50 16:28:51 16:28:52 16:28:53 16:28:54 16:28:55 16:28:56 16:28:57 16:28:58 16:28:59 16:29:00 16:29:01 16:29:02 16:29:03 16:29:04 16:29:05 16:29:06 16:29:07 16:29:08 16:29:09 16:29:10 16:29:11 16:29:12 16:29:13 16:29:14 16:29:15 16:29:16 16:29:17 16:29:18 16:29:19 16:29:20 16:29:21 16:29:22 16:29:23 16:29:24 16:29:25 16:29:26 16:29:27 16:29:28 16:29:29 16:29:30 16:29:31 16:29:32 16:29:33 16:29:34 16:29:35 16:29:36 16:29:37 16:29:38 16:29:39 16:29:40 16:29:41 16:29:42 16:29:43 16:29:44 16:29:45 16:29:46 16:29:47 16:29:48 16:29:49 16:29:50 16:29:51 16:29:52 16:29:53 16:29:54 16:29:55 16:29:56 16:29:57 16:29:58 16:29:59 16:30:00 16:30:01 16:30:02 16:30:03 16:30:04 16:30:05 16:30:06 16:30:07 16:30:08 16:30:09 16:30:10 16:30:11 16:30:12 16:30:13 16:30:14 16:30:15 16:30:16 16:30:17 16:30:18 16:30:19 16:30:20 16:30:21 16:30:22 16:30:23 16:30:24 16:30:25 16:30:26 16:30:27 16:30:28 16:30:29 16:30:30 16:30:31 16:30:32 16:30:33 16:30:34 16:30:35 16:30:36 16:30:37 16:30:38 16:30:39 16:30:40 16:30:41 16:30:42 16:30:43 16:30:44 16:30:45 16:30:46 16:30:47 16:30:48 16:30:49 16:30:50 16:30:51 16:30:52 16:30:53 16:30:54 16:30:55 16:30:56 16:30:57 16:30:58 16:30:59 16:31:00 16:31:01 16:31:02 16:31:03 16:31:04 16:31:05 16:31:06 16:31:07 16:31:08 16:31:09 16:31:10 16:31:11 16:31:12 16:31:13 16:31:14 16:31:15 16:31:16 16:31:17 16:31:18 16:31:19 16:31:20 16:31:21 16:31:22 16:31:23 16:31:24 16:31:25 16:31:26 16:31:27 16:31:28 16:31:29 16:31:30 16:31:31 16:31:32 16:31:33 16:31:34 16:31:35 16:31:36 16:31:37 16:31:38 16:31:39 16:31:40 16:31:41 16:31:42 16:31:43 16:31:44 16:31:45 16:31:46 16:31:47 16:31:48 16:31:49 16:31:50 16:31:51 16:31:52 16:31:53 16:31:54 16:31:55 16:31:56 16:31:57 16:31:58 16:31:59 16:32:00 16:32:01 16:32:02 16:32:03 16:32:04 16:32:05 16:32:06 16:32:07 16:32:08 16:32:09 16:32:10 16:32:11 16:32:12 16:32:13 16:32:14 16:32:15 16:32:16 16:32:17 16:32:18 16:32:19 16:32:20 16:32:21 16:32:22 16:32:23 16:32:24 16:32:25 16:32:26 16:32:27 16:32:28 16:32:29 16:32:30 16:32:31 16:32:32 16:32:33 16:32:34 16:32:35 16:32:36 16:32:37 16:32:38 16:32:39 16:32:40 16:32:41 16:32:42 16:32:43 16:32:44 16:32:45 16:32:46 16:32:47 16:32:48 16:32:49 16:32:50 16:32:51 16:32:52 16:32:53 16:32:54 16:32:55 16:32:56 16:32:57 16:32:58 16:32:59 16:33:00 16:33:01 16:33:02 16:33:03 16:33:04 16:33:05 16:33:06 16:33:07 16:33:08 16:33:09 16:33:10 16:33:11 16:33:12 16:33:13 16:33:14 16:33:15 16:33:16 16:33:17 16:33:18 16:33:19 16:33:20 16:33:21 16:33:22 16:33:23 16:33:24 16:33:25 16:33:26 16:33:27 16:33:28 16:33:29 16:33:30 16:33:31 16:33:32 16:33:33 16:33:34 16:33:35 16:33:36 16:33:37 16:33:38 16:33:39 16:33:40 16:33:41 16:33:42 16:33:43 16:33:44 16:33:45 16:33:46 16:33:47 16:33:48 16:33:49 16:33:50 16:33:51 16:33:52 16:33:53 16:33:54 16:33:55 16:33:56 16:33:57 16:33:58 16:33:59 16:34:00 16:34:01 16:34:02 16:34:03 16:34:04 16:34:05 16:34:06 16:34:07 16:34:08 16:34:09 16:34:10 16:34:11 16:34:12 16:34:13 16:34:14 16:34:15 16:34:16 16:34:17 16:34:18 16:34:19 16:34:20 16:34:21 16:34:22 16:34:23 16:34:24 16:34:25 16:34:26 16:34:27 16:34:28 16:34:29 16:34:30 16:34:31 16:34:32 16:34:33 16:34:34 16:34:35 16:34:36 16:34:37 16:34:38 16:34:39 16:34:40 16:34:41 16:34:42 16:34:43 16:34:44 16:34:45 16:34:46 16:34:47 16:34:48 16:34:49 16:34:50 16:34:51 16:34:52 16:34:53 16:34:54 16:34:55 16:34:56 16:34:57 16:34:58 16:34:59 16:35:00 16:35:01 16:35:02 16:35:03 16:35:04 16:35:05 16:35:06 16:35:07 16:35:08 16:35:09 16:35:10 16:35:11 16:35:12 16:35:13 16:35:14 16:35:15 16:35:16 16:35:17 16:35:18 16:35:19 16:35:20 16:35:21 16:35:22 16:35:23 16:35:24 16:35:25 16:35:26 16:35:27 16:35:28 16:35:29 16:35:30 16:35:31 16:35:32 16:35:33 16:35:34 16:35:35 16:35:36 16:35:37 16:35:38 16:35:39 16:35:40 16:35:41 16:35:42 16:35:43 16:35:44 16:35:45 16:35:46 16:35:47 16:35:48 16:35:49 16:35:50 16:35:51 16:35:52 16:35:53 16:35:54 16:35:55 16:35:56 16:35:57 16:35:58 16:35:59 16:36:00 16:36:01 16:36:02 16:36:03 16:36:04 16:36:05 16:36:06 16:36:07 16:36:08 16:36:09 16:36:10 16:36:11 16:36:12 16:36:13 16:36:14 16:36:15 16:36:16 16:36:17 16:36:18 16:36:19 16:36:20 16:36:21 16:36:22 16:36:23 16:36:24 16:36:25 16:36:26 16:36:27 16:36:28 16:36:29 16:36:30 16:36:31 16:36:32 16:36:33 16:36:34 16:36:35 16:36:36 16:36:37 16:36:38 16:36:39 16:36:40 16:36:41 16:36:42 16:36:43 16:36:44 16:36:45 16:36:46 16:36:47 16:36:48 16:36:49 16:36:50 16:36:51 16:36:52 16:36:53 16:36:54 16:36:55 16:36:56 16:36:57 16:36:58 16:36:59 16:37:00 16:37:01 16:37:02 16:37:03 16:37:04 16:37:05 16:37:06 16:37:07 16:37:08 16:37:09 16:37:10 16:37:11 16:37:12 16:37:13 16:37:14 16:37:15 16:37:16 16:37:17 16:37:18 16:37:19 16:37:20 16:37:21 16:37:22 16:37:23 16:37:24 16:37:25 16:37:26 16:37:27 16:37:28 16:37:29 16:37:30 16:37:31 16:37:32 16:37:33 16:37:34 16:37:35 16:37:36 16:37:37 16:37:38 16:37:39 16:37:40 16:37:41 16:37:42 16:37:43 16:37:44 16:37:45 16:37:46 16:37:47 16:37:48 16:37:49 16:37:50 16:37:51 16:37:52 16:37:53 16:37:54 16:37:55 16:37:56 16:37:57 16:37:58 16:37:59 16:38:00 16:38:01 16:38:02 16:38:03 16:38:04 16:38:05 16:38:06 16:38:07 16:38:08 16:38:09 16:38:10 16:38:11 16:38:12 16:38:13 16:38:14 16:38:15 16:38:16 16:38:17 16:38:18 16:38:19 16:38:20 16:38:21 16:38:22 16:38:23 16:38:24 16:38:25 16:38:26 16:38:27 16:38:28 16:38:29 16:38:30 16:38:31 16:38:32 16:38:33 16:38:34 16:38:35 16:38:36 16:38:37 16:38:38 16:38:39 16:38:40 16:38:41 16:38:42 16:38:43 16:38:44 16:38:45 16:38:46 16:38:47 16:38:48 16:38:49 16:38:50 16:38:51 16:38:52 16:38:53 16:38:54 16:38:55 16:38:56 16:38:57 16:38:58 16:38:59 16:39:00 16:39:01 16:39:02 16:39:03 16:39:04 16:39:05 16:39:06 16:39:07 16:39:08 16:39:09 16:39:10 16:39:11 16:39:12 16:39:13 16:39:14 16:39:15 16:39:16 16:39:17 16:39:18 16:39:19 16:39:20 16:39:21 16:39:22 
```



**Figure 5:** A single machine using the maximal SSD size (100) and 8 vCPUs (Speed test for Task 2b)

## 2c) Improve efficiency (6%)

We implemented a straightforward version into Cloud and found **inefficiency in your code**. Running on Cloud takes quite a long time similar to running locally. Because we cannot fully utilize the power of the spark.

```
ds_comb = parameters_rdd.flatMap(new_time_configs_ds)
tf_comb = parameters_rdd.flatMap(new_time_configs_tf)

### TASK 2c ###
ds_comb.cache()
tf_comb.cache()

# create rdd_results for dataset : batch_size & images_per_second
rdd_dataset_bs_speed = ds_comb.map(lambda x: (x[0],x[5]))
dataset_bs_speed_result = rdd_dataset_bs_speed.collect()

# create rdd_results for dataset : batch_number & images_per_second
rdd_dataset_bn_speed = ds_comb.map(lambda x: (x[1],x[5]))
dataset_bn_speed_result = rdd_dataset_bn_speed.collect()

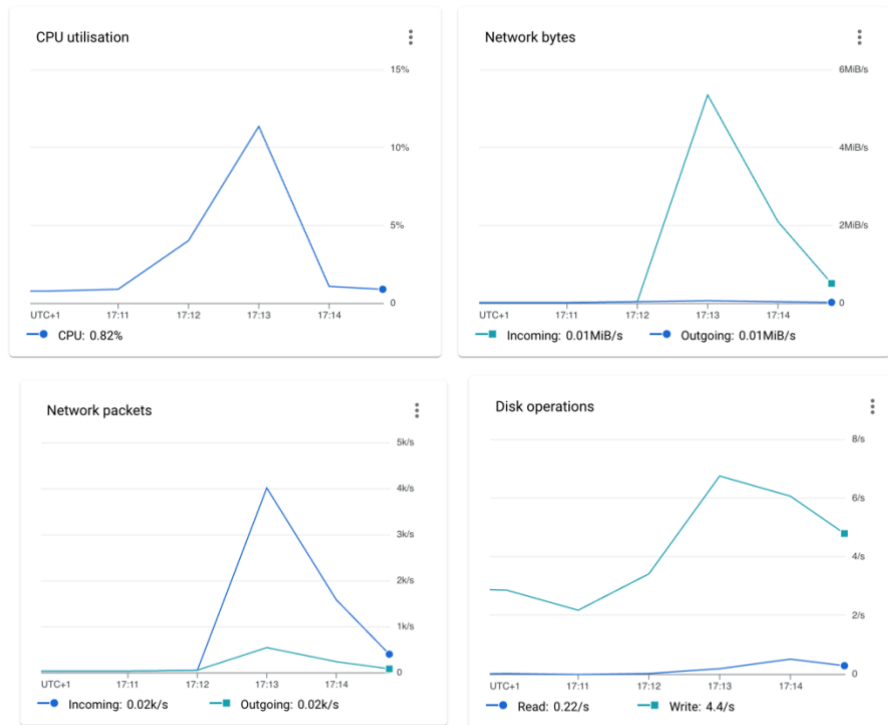
# create rdd_results for dataset : repetition & images_per_second
rdd_dataset_rep_speed = ds_comb.map(lambda x: (x[2],x[5]))
dataset_rep_speed_result = rdd_dataset_rep_speed.collect()

# create rdd_results for tf_dataset : batch_size & images_per_second
rdd_tfdataset_bs_speed = tf_comb.map(lambda x: (x[0],x[5]))
```

We apply **RDD.cache()** function to our RDD to improve efficiency. The first noticeable thing is that the computational time for compiling this code on cloud is faster. Due to this function catch the RDD in the memory or disk instead of remaking it every time you call that RDD. In task 2a, we often reuse the same RDD (ds\_comb, tf\_comb) to generate the new result RDD for each combination of parameters compared to speed. The result after applying cache function significantly increases the performance because caching the RDD reduce the amount of time spent on recomputing the RDD from its source data.

Start time	Elapsed time	Labels
1 May 2023, 17:11:51	47 sec	cache : true
1 May 2023, 16:18:37	1 min 30 sec	cache : non

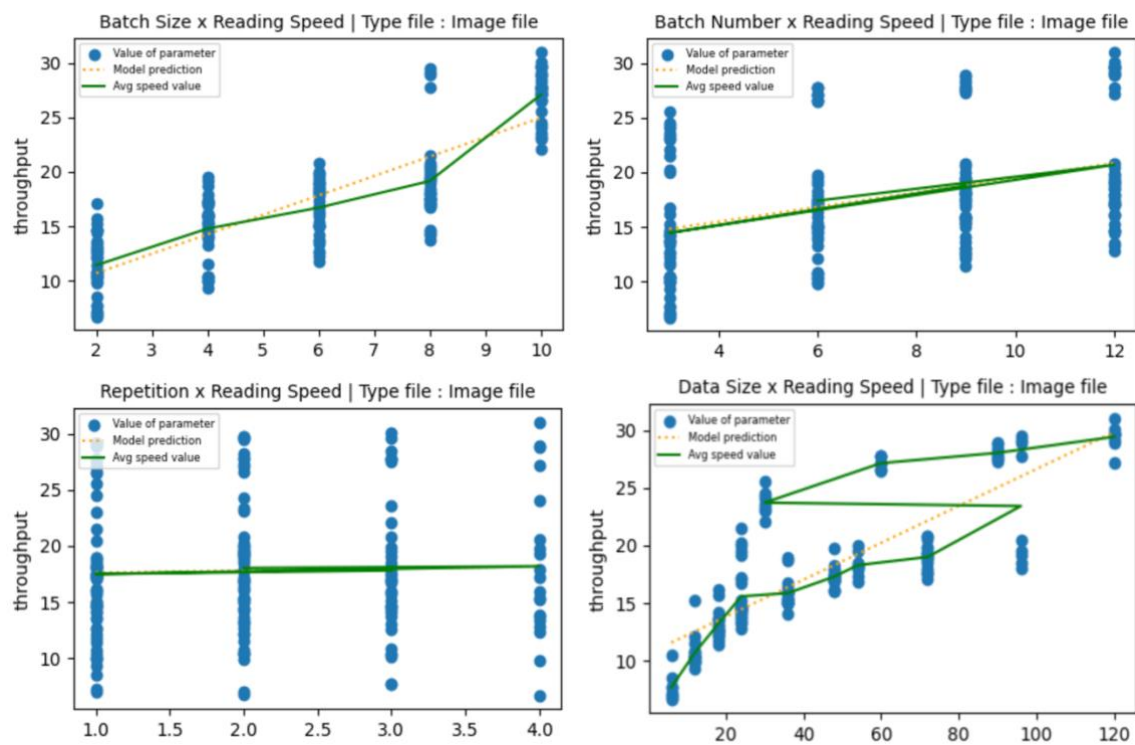
However, cache function consumes a significant amount of memory or disk space. Hence, you should make sure when applying cache RDDs that the task will reuse the same RDD more than once.



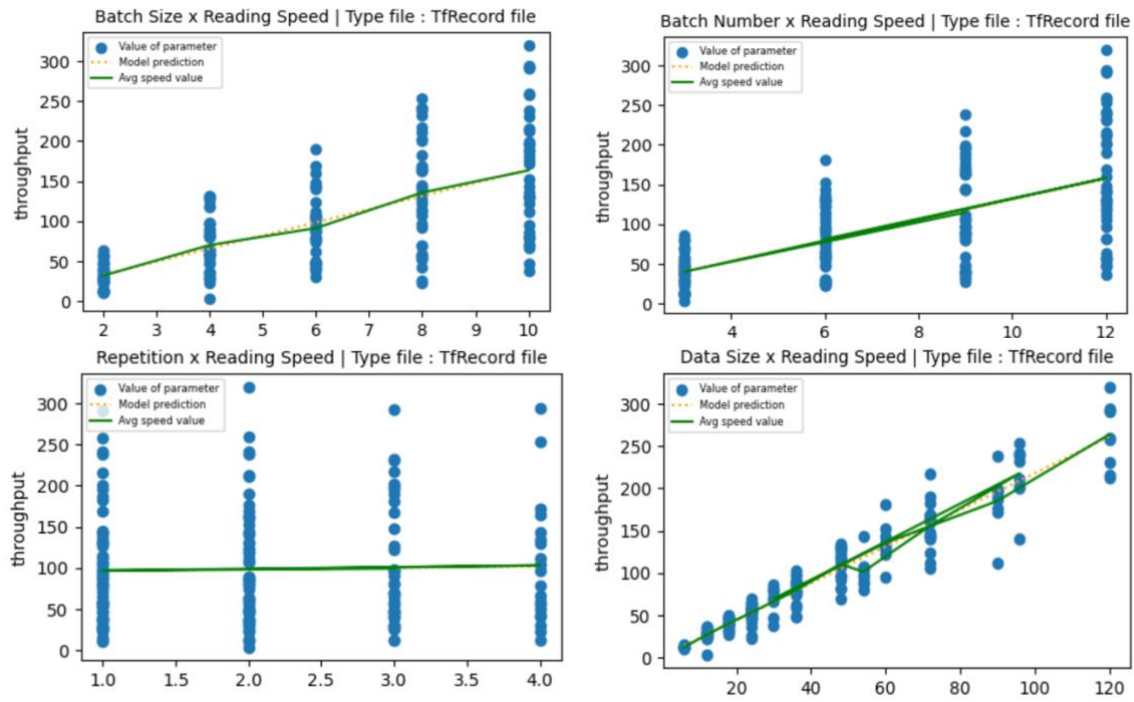
**Figure 6:** A single machine using the maximal SSD size (100) and 8 vCPUs after apply cache() (Speed test for Task 2c)

## 2d) Retrieve, analyse and discuss the output (12%)

After we retrieve the result of each combination from pickle file in GCP bucket, we perform linear regression to find a global trend and observe the correlation between each parameter compared with throughput (Image per second) in the time data points. Below are visualizations of the relationship between throughput across different parameters in each specific case which are image file and TFRecord file.



**Figure 7:** Regression plots of each parameter for image file case



**Figure 8:** Regression plots of each parameter for TFRecord file case

Scale for two cases obviously is seen that TFRecord file has a larger amount of throughput around ten times than image file with the same parameter indicating that at the same given time, TFRecord files are more effective and faster at reading flower images.

Figures 7, 8 and the table below indicate that there is a strong correlation of throughput with batch size and batch number in both TFRecord file and image file case. We can see from the high value of coefficient and low value of p-value shows that the results from the model are replicable. Similar results were found in data size, which is the product of batch size and number, which revealed a linear positive trend. The result suggests that increasing batch parameters such as batch number and batch size lead to a significantly larger throughput per second as well. There is no relationship between repetition and throughput in both cases. It can be seen that the lines of the average and prediction from the model are flat on the horizontal line no matter how much the number of repetitions increases. In addition, high p-value demonstrate the result cannot conclude that the explanatory repetition affects the throughput. Therefore, repetition does not affect the amount of throughput.

Case Type	Parameter	Coefficient	Intercept	P-value
Image file	Batch Size	1.784	7.106	7.9e-49
	Batch Number	0.670	12.787	2.4e-7
	Repetition	0.188	17.417	0.679
	Data Size	0.160	10.607	1.3e-44
TFRecord file	Batch Size	16.406	0.124	1.7e-24
	Batch Number	12.902	1.799	1.2e-20
	Repetition	1.750	94.869	0.739
	Data Size	2.168	0.989	7.6e-98

**Table 1:** Result of the parameters

It is much more beneficial to run testing in local single machine before running on cloud to avoid technical errors. However, for **large-scale machine learning**, running on cloud with parallelization might be more effective in both performance and training time. Providing a better environment to run the model simultaneously when dealing with larger datasets. As a result, it will lead to higher accuracy of the linear model as well as reduce the amount of time for training model.

However, there are some circumstances that we might encounter with bottleneck problem with limits of the bucket reading speed even though we parallelize the data and model and run simultaneously we have to consider the amount of retrieved input that can feed into the model. Setting a higher number of parallels or maximum cluster without considering the retrieved data can result in wasteful of resources and cost.

### Task 3: Discussion in context. (24%)

#### 3a) Contextualise

Throughout the experiment on cloud with various setting of clusters from this coursework, we can observe the diversity of the results in both good and bad performance in terms of processing time, CPU, network, and disk bandwidth according to the parameter and type that we set to our cloud cluster machine. Setting the right variables leads to significant improvement in processing time as well as performance. In contrast, poor setting may result in a slow process or wasteful resources and costs or even worse than running locally.

Since there are dozens of possible VM instance types and even more cluster sizes to choose from, choosing the right cloud configuration may take time and resources. As we can see from the previous tasks setting cluster and execute on google cloud take an amount of time especially with the large cluster can take up to 5-10 minutes just only to create single cluster.

Cherry-pick can be applied to mitigate this problem and find optimal configurations in few times instead of running through all possible sets of clusters manually. Cherry-pick take advantage of Bayesian Optimization method to build accurate models to distinguish the near-best cloud configuration in a few runs. As a result, it can save time and searching costs for optimized cloud configuration and also guarantee that we will end up with the near-best likely optimal configuration.

#### 3b) Strategise

The flower images dataset in this coursework uses the benefit of **batch** processing to break the workload into a smaller chunk of data and process it simultaneously. In spark programming we use method called parallelization to create Resilient Distributed Datasets (RDD) to perform parallel processing across a cluster or machines. This leverages the use of resources such as CPU, memory, and disk space significantly improving overall performance and processing time. Batch processing usually perform on large scale dataset which required a lot of processing power and storage to be processed efficiently. This method typically executes on existing datasets without the need for real-time processing such as in the lab datasets are normally maintained in the storage and ready to use which we focus on quality and performance in terms of accuracy and time processing.

On the other hand, **stream** processing might involve the real-time dataset a great example would it be continuous flow data that come from cameras or sensors, or sales transactions. This kind of data plays an important part in business nowadays providing tremendous insights and being able to process in real time. Giving the instant result or immediate feedback to the user. However, similarly to batch, stream processing also relies on scalable real-time processing services such as GCP to distribute the workload across multiple machines to increase the processing time and performance which is easier to control and limited the scale of resources. The result may benefit the end-user experience by giving instant responses and improving the computational time.



These two scenarios required different resources and configurations in order to reach the highest performance. As we discussed in section 3a to maximize the full potential of both resources and configurations. Cherry-Pick is a decent method to be applied for searching the best cloud configurations by improving the estimation for configurations that are closer to the optimal. Using Bayesian optimization technique, searching across different cloud configurations allows this algorithm to work well even with limited given information. Consequently, performance, reliability, and cost for these two scenarios can be improved by applying Cherry-Pick. We can ensure that we will end up with at least near-best optimize configuration which lead to an increase in the performance of the task and minimize the spending cost.

**Word Count:** 1898 words (not including Question, figures, and table)

**Colab Link:** <https://colab.research.google.com/drive/13JmvSEmvPOmuRUwCd0gJsJFQvdLz8V9?usp=sharing>