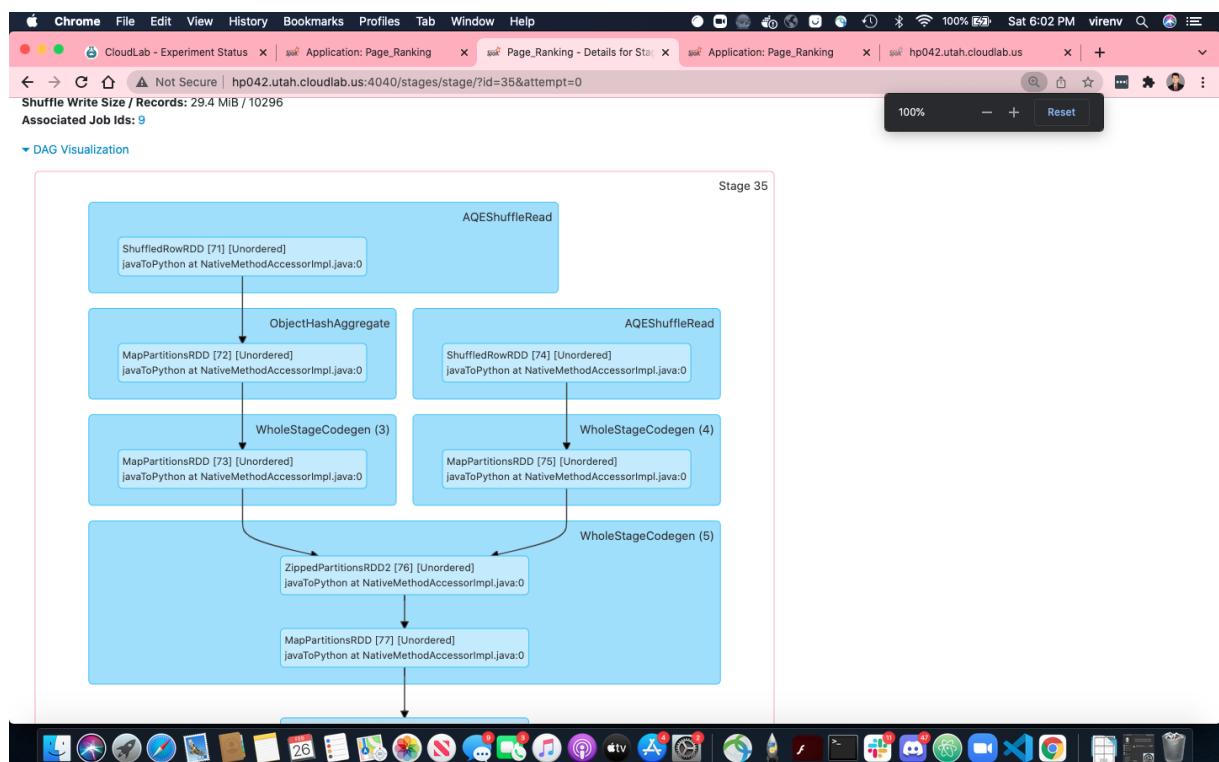


Pranav Akinepalli
Sameer Haniyur
Viren Velacheri
February 27th,2022

Assignment 2 Report

Task 1 Analysis

Small Dataset Analysis



(JavaToPython Job Stage Details)

Sat 6:05 PM virenv

Workers (3)

Worker Id	Address	State	Cores	Memory	Resources
worker-20220225111033-128.110.218.113-40187	128.110.218.113:40187	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-20220225111034-128.110.218.81-34071	128.110.218.81:34071	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-20220225111034-128.110.218.87-36219	128.110.218.87:36219	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	

Running Applications (0)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
Completed Applications (11)								
Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
app-20220226170213-0010	Page_Ranking	30	29.0 GiB		2022/02/26 17:02:13	pa8789	FINISHED	1.3 min
app-20220226165902-0009	Page_Ranking	30	29.0 GiB		2022/02/26 16:59:02	pa8789	FINISHED	1.3 min
app-20220225161116-0008	Page_Ranking	30	29.0 GiB		2022/02/25 16:11:16	pa8789	FINISHED	3.4 min
app-20220225121650-0007	Page_Ranking	30	29.0 GiB		2022/02/25 12:16:50	pa8789	FINISHED	20 min
app-20220225121040-0006	Page_Ranking	30	29.0 GiB		2022/02/25 12:10:40	pa8789	FINISHED	1.3 min
app-20220225120630-0005	Page_Ranking	30	29.0 GiB		2022/02/25 12:06:30	pa8789	FINISHED	1.3 min
app-20220225120005-0004	Page_Ranking	30	29.0 GiB		2022/02/25 12:00:05	pa8789	FINISHED	1.2 min
app-20220225115633-0003	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:33	pa8789	FINISHED	1.3 min
app-20220225115600-0002	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:00	pa8789	FINISHED	2 s
app-20220225114951-0001	Page_Ranking	30	29.0 GiB		2022/02/25 11:49:51	pa8789	FINISHED	10 s
app-20220225111717-0000	Page_Ranking	30	29.0 GiB		2022/02/25 11:17:17	pa8789	FINISHED	20 min

(1.3 minutes to completion)

Sat 6:03 PM virenv

CloudLab - Experiment Status | Application: Page_Ranking | Page_Ranking - Spark Jobs | Application: Page_Ranking | hp042.utah.cloudlab.us | + | Not Secure | hp042.utah.cloudlab.us:4040/jobs/

Spark Jobs (?)

User: pa8789
Total Uptime: 1.2 min
Scheduling Mode: FIFO
Active Jobs: 1
Completed Jobs: 27

Event Timeline

Enable zooming

Executors

- Added
- Removed

Jobs

- Succeeded
- Failed
- Running

Timeline: 15 20 25 30 35 40 45 50 55 0 5 10 15 20

26 February 17:02 26 February 17:03

Active Jobs (1)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

(Application Timeline)

Takeaway: The algorithm runs pretty fast. With that said, as shown by the Spark Jobs timeline, many jobs are taking place. In this case, about 30 jobs were completed by the end of the application's lifetime. Further, it was interesting seeing how the jobs were broken up further into stages. A good example is the DAG visualization for one of the common javaToPython jobs where we can see that it is being broken down into different operations like shuffle read and so forth. The timeline above also shows that there is a pattern of repetition in the jobs being run which is probably representative of the 10 iterations that the algorithm runs for. Overall, the page ranking algorithm for the small dataset took around 1.3 minutes (78 seconds) to complete.

Large Dataset Analysis

The screenshot shows a web browser window with the following tabs:

- CloudLab - Experiment Status
- Spark Master at spark://hp042...
- Task 1 Big Screenshots - Google

The main content area displays the following information:

Resources in use:

- Applications: 0 Running, 12 Completed
- Drivers: 0 Running, 0 Completed
- Status: ALIVE

Workers (3)

Worker Id	Address	State	Cores	Memory	Resources
worker-20220225111033-128.110.218.113-40187	128.110.218.113:40187	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-20220225111034-128.110.218.81-34071	128.110.218.81:34071	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-20220225111034-128.110.218.87-36219	128.110.218.87:36219	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	

Running Applications (0)

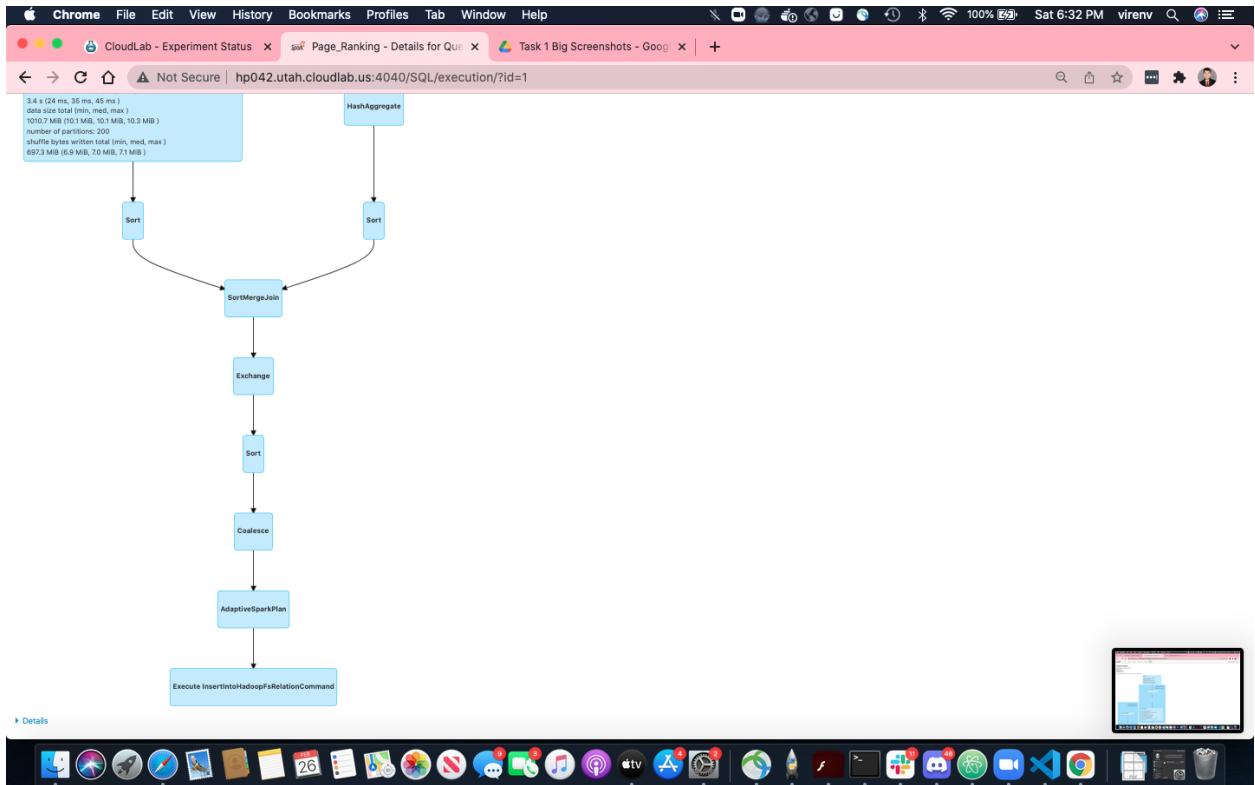
Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
No applications are currently running.								

Completed Applications (12)

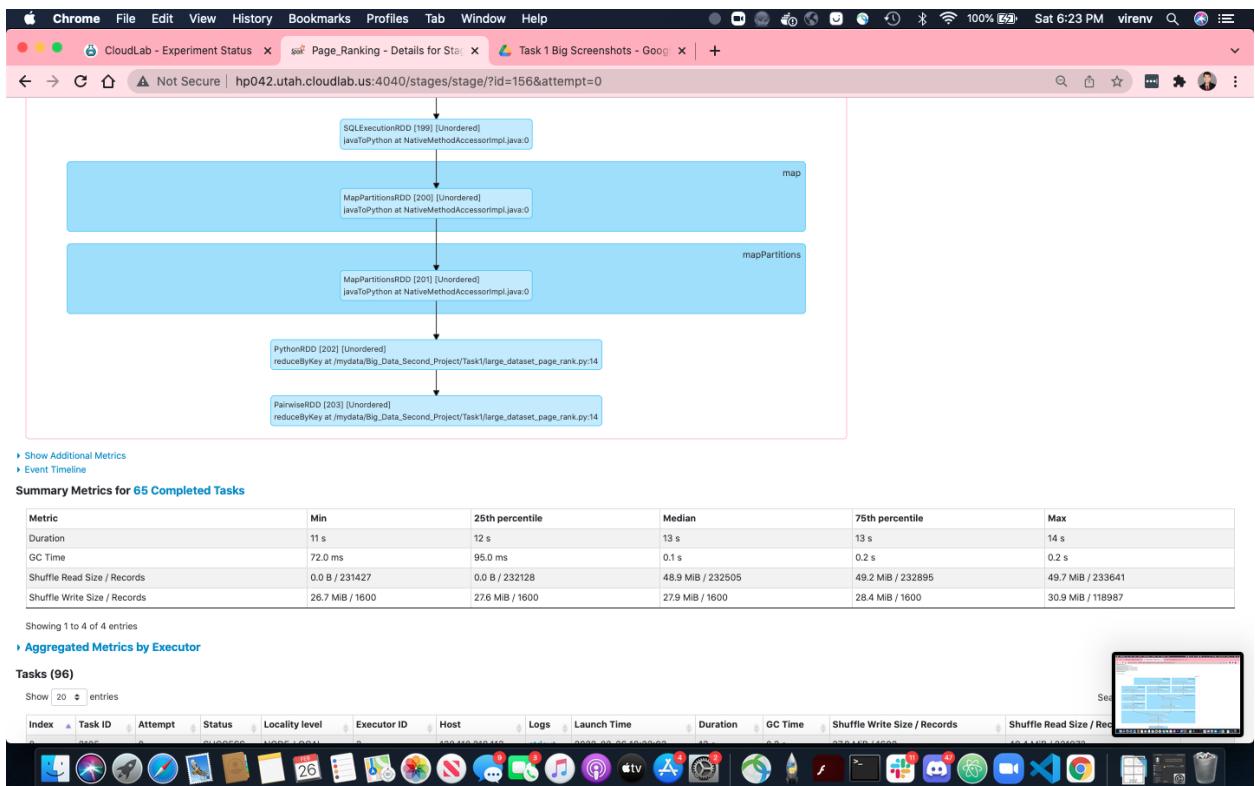
Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
app-20220226171211-0011	Page_Ranking	30	29.0 GiB		2022/02/26 17:12:11	pa8789	FINISHED	20 min
app-20220226170213-0010	Page_Ranking	30	29.0 GiB		2022/02/26 17:02:13	pa8789	FINISHED	1.3 min
app-20220226165902-0009	Page_Ranking	30	29.0 GiB		2022/02/26 16:59:02	pa8789	FINISHED	1.3 min
app-20220225161116-0008	Page_Ranking	30	29.0 GiB		2022/02/25 16:11:16	pa8789	FINISHED	3.4 min
app-20220225121650-0007	Page_Ranking	30	29.0 GiB		2022/02/25 12:16:50	pa8789	FINISHED	20 min
app-20220225121040-0006	Page_Ranking	30	29.0 GiB		2022/02/25 12:10:40	pa8789	FINISHED	1.3 min
app-20220225120630-0005	Page_Ranking	30	29.0 GiB		2022/02/25 12:06:30	pa8789	FINISHED	1.3 min
app-20220225120005-0004	Page_Ranking	30	29.0 GiB		2022/02/25 12:00:05	pa8789	FINISHED	1.2 min
app-20220225115633-0003	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:33	pa8789	FINISHED	1.3 min
app-20220225115600-0002	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:00	pa8789	FINISHED	2 s
app-20220225114951-0001	Page_Ranking	30	29.0 GiB		2022/02/25 11:49:51	pa8789	FINISHED	10 s
app-20220225111717-0000	Page_Ranking	30	29.0 GiB		2022/02/25 11:17:17	pa8789	FINISHED	20 min

At the bottom of the browser window, there is a dock with various application icons.

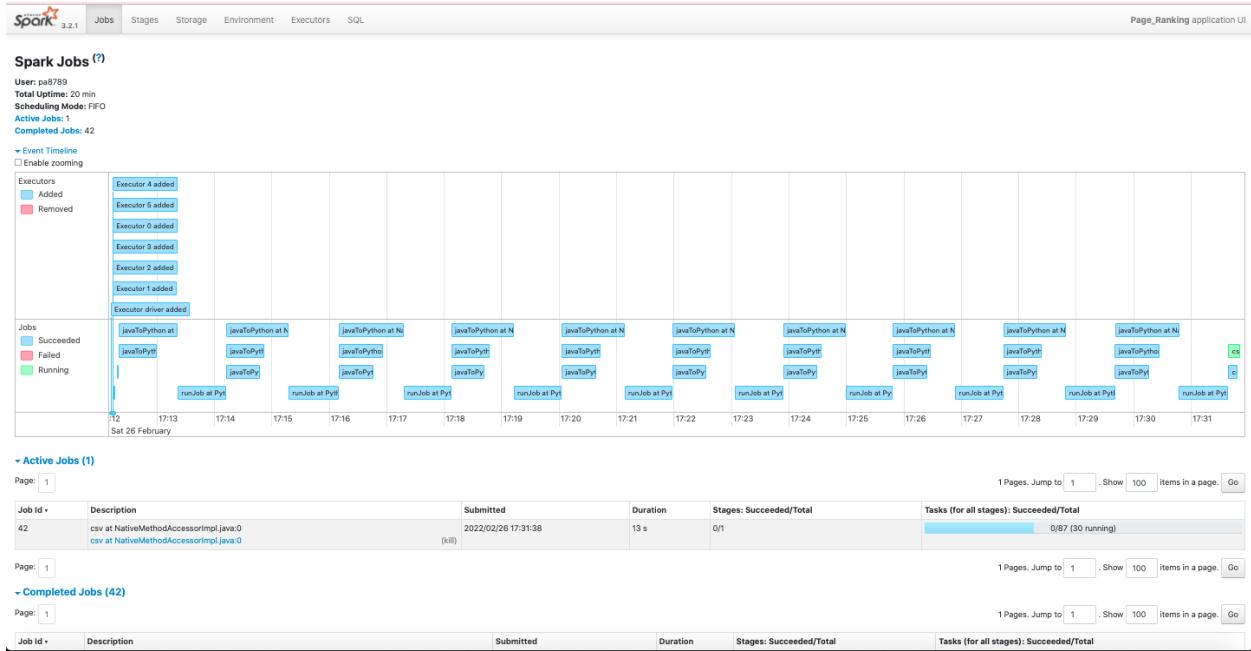
(20 min completion time)



(SQL Execution plan)



(Summary Statistics, ReduceByKey in lineage graph)



(Timeline)

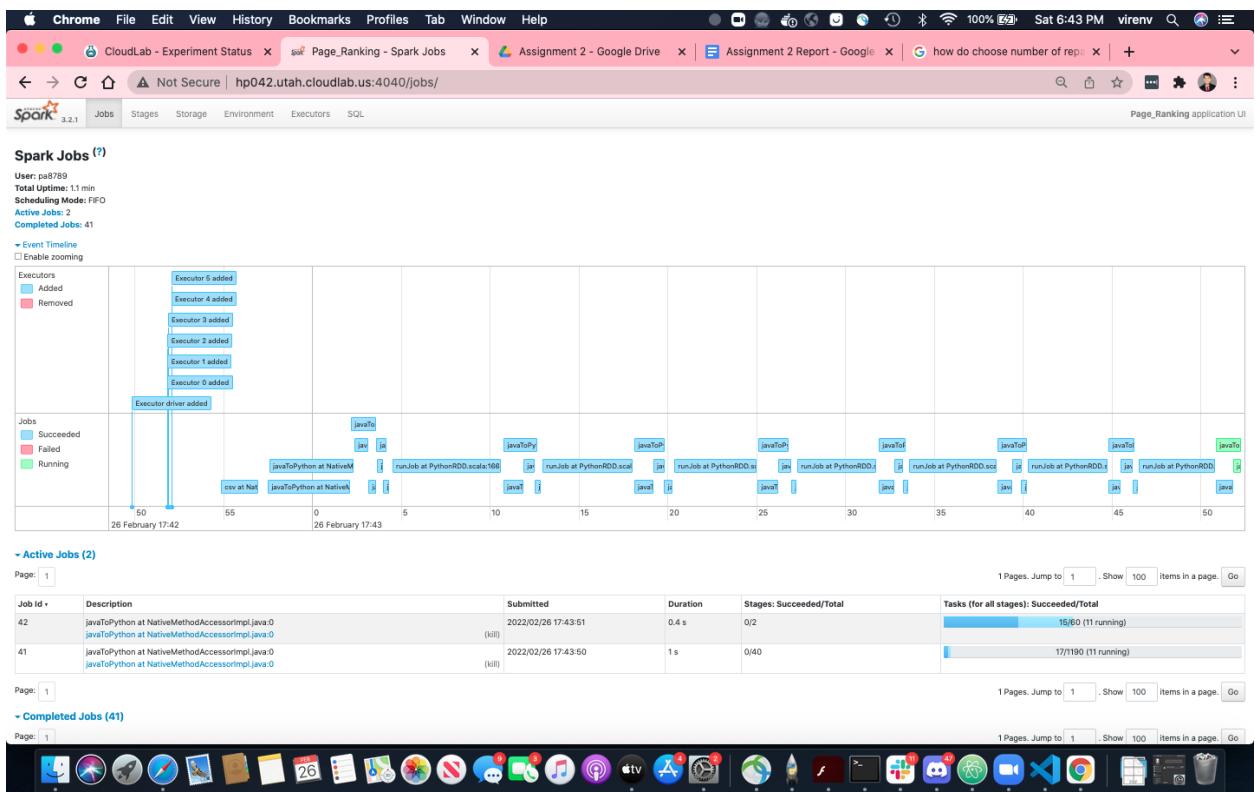
Takeaway: As expected given the large dataset size, the algorithm completed in a much longer period of time, approximately 20 minutes total. We can see that running through this dataset resulted in over 65 completed tasks at one instance in a given job. It is also interesting to visualize the execution plan for what appears to be a single job in the application. We can clearly see various stages in the physical execution plan generated by the query planner, such as the SortMergeJoin algorithm for joining, likely a result of our join query with the links and ranks datasets.

Task 2 Analysis

Small Dataset - 30 Repartitions

URL: spark://hp04t2.utah.cloudlab.us:7077	Alive Workers: 3	Cores in use: 60 Total: 0 Used	Memory in use: 195.4 GB Total: 0.0 B Used	Applications: 0 Running, 13 Completed	Drivers: 0 Running, 0 Completed	Status: ALIVE		
- Workers (3)								
Worker ID	Address	State	Cores	Memory	Resources			
worker-20220226111033-128.110.218.113-40187	128.110.218.113-40187	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)				
worker-20220226111034-128.110.218.81-34071	128.110.218.81-34071	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)				
worker-20220226111034-128.110.218.87-36219	128.110.218.87-36219	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)				
- Running Applications (0)								
Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
- Completed Applications (13)								
Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
app-20220226174249-0012	Page_Ranking	30	29.0 GiB		2022/02/26 17:42:49	pa8789	FINISHED	1.4 min
app-20220226171211-0011	Page_Ranking	30	29.0 GiB		2022/02/26 17:12:11	pa8789	FINISHED	20 min
app-20220226170213-0010	Page_Ranking	30	29.0 GiB		2022/02/26 17:02:13	pa8789	FINISHED	1.3 min

(Completion time 1.4 min, a bit slower than Task 1 small)



(Timeline)

The above timeline shows the numbers of jobs run through at 41 respectively. This is quite a bit more than the previous run of the program with no repartitions.

Summary Metrics for 33 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	0.8 s	1 s	1 s	1 s	2 s
GC Time	0.0 ms	0.0 ms	0.0 ms	30.0 ms	0.3 s
Shuffle Read Size / Records	2.3 MB / 40360	2.3 MB / 40756	2.3 MB / 41021	2.3 MB / 41205	3.1 MB / 54863
Shuffle Write Size / Records	1.2 MB / 429	1.2 MB / 429	1.2 MB / 429	1.2 MB / 429	1.5 MB / 462

Showing 1 to 4 of 4 entries

Aggregated Metrics by Executor

Tasks (33)

Index	Task ID	Attempt	Status	Locality Level	Executor ID	Host	Logs	Launch Time	Duration	GC Time	Shuffle Write Size / Records	Shuffle Read Size / Records	Errors
0	346	0	SUCCESS	NODE_LOCAL	4	128.110.218.81	stdout stderr	2022-02-26 18:43:12	1 s		1.2 MB / 429	2.3 MB / 4109	
1	347	0	SUCCESS	NODE_LOCAL	5	128.110.218.81	stdout stderr	2022-02-26 18:43:12	2 s	0.3 s	1.2 MB / 429	2.3 MB / 40947	
2	348	0	SUCCESS	NODE_LOCAL	4	128.110.218.81	stdout stderr	2022-02-26 18:43:12	2 s		1.2 MB / 429	2.3 MB / 40645	
3	349	0	SUCCESS	NODE_LOCAL	5	128.110.218.81	stdout stderr	2022-02-26 18:43:12	2 s	0.3 s	1.2 MB / 429	2.3 MB / 40819	
4	350	0	SUCCESS	NODE_LOCAL	4	128.110.218.81	stdout stderr	2022-02-26 18:43:12	1 s		1.2 MB / 429	2.3 MB / 40724	
5	351	0	SUCCESS	NODE_LOCAL	5	128.110.218.81	stdout stderr	2022-02-26 18:43:12	2 s	0.3 s	1.2 MB / 429	2.3 MB / 4076	
6	352	0	SUCCESS	NODE_LOCAL	4	128.110.218.81	stdout stderr	2022-02-26 18:43:12	1 s		1.2 MB / 429	2.3 MB / 40900	
7	353	0	SUCCESS	NODE_LOCAL	5	128.110.218.81	stdout stderr	2022-02-26 18:43:12	2 s	0.3 s	1.2 MB / 429	2.3 MB / 40360	
8	354	0	SUCCESS	NODE_LOCAL	4	128.110.218.81	stdout stderr	2022-02-26 18:43:12	1 s		1.2 MB / 429	2.3 MB / 41132	
9	355	0	SUCCESS	NODE_LOCAL	5	128.110.218.81	stdout stderr	2022-02-26 18:43:12	2 s	0.3 s	1.2 MB / 429	2.3 MB / 41112	
10	356	0	SUCCESS	NODE_LOCAL	4	128.110.218.81	stdout stderr	2022-02-26 18:43:14	1 s		1.2 MB / 429	2.3 MB / 41260	
11	357	0	SUCCESS	NODE_LOCAL	4	128.110.218.81	stdout stderr	2022-02-26 18:43:14	1 s		1.2 MB / 429	2.3 MB / 41121	
12	358	0	SUCCESS	NODE_LOCAL	4	128.110.218.81	stdout	2022-02-26 18:43:14	1 s	24.0 ms	1.2 MB / 429	2.3 MB / 40757	

(Summary statistics for 30 partition small dataset)

Takeaway: Compared to Task 1 where there was no repartitioning, the page ranking algorithm on the small dataset with 30 partitions took longer (1.4 minutes instead of 1.3 or around 6 seconds longer). Looking at the timeline, the overall pattern is still the same as with no repartitioning in terms of having a repeated structure of jobs probably representing the 10 iterations being run. However, the main difference between this run and having no repartitioning is the number of jobs that had to be run. Around the same time, this run had 41 jobs completed whereas the run with no repartitioning had 27 jobs completed. Thus, there were significantly more jobs that ran which we suspect is due to repartitioning and having more partitions meant more jobs and tasks that had to be delegated to handle the increased number of partitions. We also think that the algorithm took a bit longer to run due to the overhead of repartitioning and all the shuffling that it is required to perform.

Small Dataset - 60 Repartitions:

Screenshot of the CloudLab interface showing cluster status and completed applications.

Workers (3)

Worker Id	Address	State	Cores	Memory	Resources
worker-2022022511033-128.110.218.113-40187	128.110.218.113:40187	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-2022022511034-128.110.218.81-34071	128.110.218.81:34071	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-2022022511034-128.110.218.87-36219	128.110.218.87:36219	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	

Running Applications (0)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
----------------	------	-------	---------------------	------------------------	----------------	------	-------	----------

Completed Applications (16)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
app-20220226183235-0015	Page_Ranking	30	29.0 GiB		2022/02/26 18:32:35	pa8789	FINISHED	1.5 min
app-20220226174912-0014	Page_Ranking	30	29.0 GiB		2022/02/26 17:49:12	pa8789	FINISHED	21 min
app-20220226174813-0013	Page_Ranking	30	29.0 GiB		2022/02/26 17:48:13	pa8789	FINISHED	8 s
app-20220226174249-0012	Page_Ranking	30	29.0 GiB		2022/02/26 17:42:49	pa8789	FINISHED	1.4 min
app-20220226171211-0011	Page_Ranking	30	29.0 GiB		2022/02/26 17:12:11	pa8789	FINISHED	20 min
app-20220226170213-0010	Page_Ranking	30	29.0 GiB		2022/02/26 17:02:13	pa8789	FINISHED	1.3 min
app-20220226165902-0009	Page_Ranking	30	29.0 GiB		2022/02/26 16:59:02	pa8789	FINISHED	1.3 min
app-20220225161116-0008	Page_Ranking	30	29.0 GiB		2022/02/25 16:11:16	pa8789	FINISHED	3.4 min
app-20220225121650-0007	Page_Ranking	30	29.0 GiB		2022/02/25 12:16:50	pa8789	FINISHED	20 min
app-20220225121040-0006	Page_Ranking	30	29.0 GiB		2022/02/25 12:10:40	pa8789	FINISHED	1.3 min
app-20220225120630-0005	Page_Ranking	30	29.0 GiB		2022/02/25 12:06:30	pa8789	FINISHED	1.3 min
app-20220225120005-0004	Page_Ranking	30	29.0 GiB		2022/02/25 12:00:05	pa8789	FINISHED	1.2 min
app-20220225115633-0003	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:33	pa8789	FINISHED	1.3 min

(Completion time for small dataset at 60 repartition, slightly slower still than 30 repartition)

Screenshot of the CloudLab interface showing summary metrics and task details for 33 completed tasks.

Summary Metrics for 33 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	1 s	2 s	2 s	2 s	3 s
GC Time	0.0 ms	0.0 ms	40.0 ms	69.0 ms	72.0 ms
Shuffle Read Size / Records	2.3 MiB / 40216	2.4 MiB / 40606	2.4 MiB / 40900	2.4 MiB / 41036	3.2 MiB / 54688
Shuffle Write Size / Records	908.3 Kib / 429	916.6 Kib / 429	926.3 Kib / 429	930.9 Kib / 429	1.1 MiB / 462

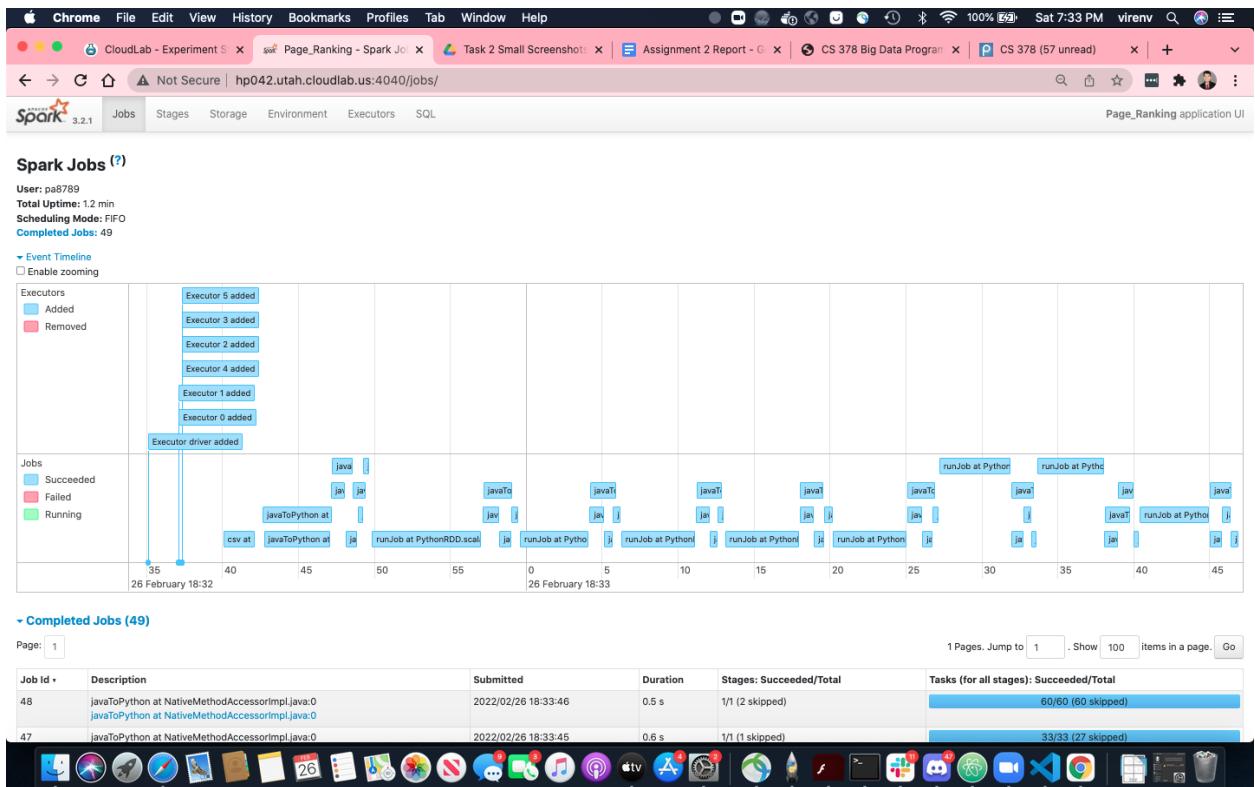
Showing 1 to 4 of 4 entries

Aggregated Metrics by Executor

Tasks (33)

Index	Task ID	Attempt	Status	Locality level	Executor ID	Host	Logs	Launch Time	Duration	GC Time	Shuffle Write Size / Records	Shuffle Read Size / Records	Errors
0	280	0	SUCCESS	NODE_LOCAL	0	128.110.218.87	stdout stderr	2022-02-26 19:32:49	2 s	72.0 ms	931.6 Kib / 429	2.4 MiB / 40948	
1	281	0	SUCCESS	NODE_LOCAL	1	128.110.218.87	stdout stderr	2022-02-26 19:32:49	2 s	41.0 ms	908.3 Kib / 429	2.4 MiB / 40798	
2	282	0	SUCCESS	NODE_LOCAL	0	128.110.218.87	stdout stderr	2022-02-26 19:32:49	2 s	72.0 ms	913.6 Kib / 429	2.3 MiB / 40514	
3	283	0	SUCCESS	NODE_LOCAL	1	128.110.218.87	stdout stderr	2022-02-26 19:32:49	2 s	41.0 ms	912.7 Kib / 429	2.4 MiB / 40670	
4	284	0	SUCCESS	NODE_LOCAL	0	128.110.218.87	stdout stderr	2022-02-26 19:32:49	2 s	72.0 ms	926.4 Kib / 429	2.4 MiB / 40586	
5	285	0	SUCCESS	NODE_LOCAL	1	128.110.218.87	stdout stderr	2022-02-26 19:32:49	3 s	41.0 ms	921.8 Kib / 429	2.4 MiB / 40604	
6	286	0	SUCCESS	NODE_LOCAL	0	128.110.218.87	stdout stderr	2022-02-26 19:32:49	2 s	72.0 ms	922.7 Kib / 429	2.4 MiB / 40756	
7	287	0	SUCCESS	NODE_LOCAL	1	128.110.218.87	stdout stderr	2022-02-26 19:32:49	2 s	41.0 ms	921.3 Kib / 429	2.3 MiB / 40216	
8	288	0	SUCCESS	NODE_LOCAL	0	128.110.218.87	stdout stderr	2022-02-26 19:32:49	2 s	72.0 ms	927.4 Kib / 429	2.4 MiB / 41008	
9	289	0	SUCCESS	NODE_LOCAL	1	128.110.218.87	stdout stderr	2022-02-26 19:32:49	2 s	41.0 ms	930.9 Kib / 429	2.4 MiB / 40978	
10	290	0	SUCCESS	NODE_LOCAL	0	128.110.218.87	stdout	2022-02-26 19:32:52	2 s		915.9 Kib / 429	2.4 MiB / 41112	

(Summary stats for 60 repartition small dataset)



(Timeline)

Takeaway: Doubling the number of repartitions made the application even slower as it finished in 1.5 minutes instead of 1.4 minutes. It appears that adding more and more repartitions, at least in the context of this Page Ranking application, adds overhead that leads to the application running slower than before. Based on the summary metrics, it is apparent that things like garbage collection took more time. Adding more repartitions led to more jobs having to be done as about 50 or so jobs were probably completed by the end of application based on the above timeline.

Large Dataset 30 Repartitions

Not Secure | hp042.utah.cloudlab.us:4040/stages/stage/?id=967&attempt=0

PythonRDD [465] [Unordered]
reduceByKey at /mydata/Big_Data_Second_Project/Task2/large_dataset_page_rank.py:15

PairwiseRDD [466] [Unordered]
reduceByKey at /mydata/Big_Data_Second_Project/Task2/large_dataset_page_rank.py:15

Show Additional Metrics
Event Timeline

Summary Metrics for 21 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	12 s	12 s	12 s	12 s	13 s
GC Time	19.0 ms	58.0 ms	0.1 s	0.1 s	0.2 s
Shuffle Read Size / Records	0.0 B / 230827	0.0 B / 231394	0.0 B / 231995	0.0 B / 232194	49 MiB / 232981
Shuffle Write Size / Records	26.8 MiB / 1600	27.6 MiB / 1600	27.8 MiB / 1600	27.9 MiB / 1600	29.6 MiB / 1600

Showing 1 to 4 of 4 entries

Aggregated Metrics by Executor

Tasks [50]

Index	Task ID	Attempt	Status	Locality level	Executor ID	Host	Logs	Launch Time	Duration	GC Time	Shuffle Write Size / Records	Shuffle Read Size / Records	Errors
0	5410	0	SUCCESS	NODE_LOCAL	0	128.110.218.87	stdout stderr	2022-02-26 19:08:19	12 s	0.2 s	27.8 MiB / 1600	49.3 MiB / 231388	
1	5411	0	SUCCESS	NODE_LOCAL	2	128.110.218.113	stdout stderr	2022-02-26 19:08:19	13 s	0.1 s	28.6 MiB / 1600	50.1 MiB / 232315	
2	5412	0	SUCCESS	NODE_LOCAL	3	128.110.218.113	stdout stderr	2022-02-26 19:08:19	12 s	0.1 s	27.8 MiB / 1600	48.8 MiB / 232517	
3	5413	0	SUCCESS	NODE_LOCAL	1	128.110.218.87	stdout stderr	2022-02-26 19:08:19	12 s	58.0 ms	27.6 MiB / 1600	48.8 MiB / 231995	
4	5415	0	SUCCESS	NODE_LOCAL	2	128.110.218.113	stdout stderr	2022-02-26 19:08:19	12 s	0.1 s	27.9 MiB / 1600	49.7 MiB / 232383	
5	5414	0	SUCCESS	NODE_LOCAL	0	128.110.218.87	stdout stderr	2022-02-26 19:08:19	12 s	0.2 s	27.6 MiB / 1600	48.7 MiB / 232194	

(Summary metrics for large dataset, 30 repartitions: low shuffle read size but large shuffle write size compared to small dataset at 30 partitions. Garbage collection time > small dataset GC time)

↓ Active Jobs (3)

CloudLab - Experiment S | Spark Master at spark:// | Task 2 Big Screenshots | Assignment 2 Report - G | CS 378 Big Data Program | CS 378 (57 unread) | + | Sat 9:43 PM virenv | Not Secure | hp042.utah.cloudlab.us:8080

Applications: 0 Running, 18 Completed
Drivers: 0 Running, 0 Completed
Status: ALIVE

Workers (3)

Worker Id	Address	State	Cores	Memory	Resources
worker-2022022511033-128.110.218.113-40187	128.110.218.113:40187	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-2022022511034-128.110.218.81-34071	128.110.218.81:34071	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-2022022511034-128.110.218.87-36219	128.110.218.87:36219	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	

Running Applications (0)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
Completed Applications (18)								
Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
app-20220226204158-0017	Page_Ranking	30	29.0 GiB		2022/02/26 20:41:58	pa8789	FINISHED	1.6 min
app-20220226183836-0016	Page_Ranking	30	29.0 GiB		2022/02/26 18:36:36	pa8789	FINISHED	20 min
app-20220226183238-0015	Page_Ranking	30	29.0 GiB		2022/02/26 18:32:38	pa8789	FINISHED	1.5 min
app-20220226174912-0014	Page_Ranking	30	29.0 GiB		2022/02/26 17:49:12	pa8789	FINISHED	21 min
app-20220226174819-0013	Page_Ranking	30	29.0 GiB		2022/02/26 17:48:19	pa8789	FINISHED	8 s
app-20220226174249-0012	Page_Ranking	30	29.0 GiB		2022/02/26 17:42:49	pa8789	FINISHED	1.4 min
app-20220226171211-0001	Page_Ranking	30	29.0 GiB		2022/02/26 17:12:11	pa8789	FINISHED	20 min
app-20220226170213-0010	Page_Ranking	30	29.0 GiB		2022/02/26 17:02:13	pa8789	FINISHED	1.3 min
app-20220226165902-0009	Page_Ranking	30	29.0 GiB		2022/02/26 16:59:02	pa8789	FINISHED	1.3 min
app-20220225161116-0008	Page_Ranking	30	29.0 GiB		2022/02/25 16:11:16	pa8789	FINISHED	3.4 min
app-20220225121650-0007	Page_Ranking	30	29.0 GiB		2022/02/25 12:16:50	pa8789	FINISHED	20 min
app-20220225121040-0006	Page_Ranking	30	29.0 GiB		2022/02/25 12:10:40	pa8789	FINISHED	1.3 min
app-20220225120630-0005	Page_Ranking	30	29.0 GiB		2022/02/25 12:06:30	pa8789	FINISHED	1.3 min
app-20220225120005-0004	Page_Ranking	30	29.0 GiB		2022/02/25 12:00:05	pa8789	FINISHED	1.2 min
app-20220225115633-0003	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:33	pa8789	FINISHED	1.3 min
app-20220225115600-0002	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:00	pa8789	FINISHED	2 s
app-20220225114404-0001	Page_Ranking	30	29.0 GiB		2022/02/25 11:44:04	pa8789	FINISHED	10 s

(Total time for application was 21 minutes, so just a minute slower than without any repartitions, but still around the same time)

Takeaway: It seems like the page ranking algorithm repartitioned for 30 partitions run on the large dataset mostly shares the same qualities of the one run on the small dataset with amplified effects. For example, it still shares the quality of taking slightly longer than the original run on the dataset (21 minutes vs 20 minutes so a minute longer). This again is most likely due to the overhead of repartitioning with the shuffling involved. It also has the same pattern in the timeline representing the 10 iterations being run. We also see that it completes significantly more jobs than the original run on the large dataset (around 80 jobs vs 40 jobs so double the amount). Compared to the 30 repartitioning run on the small dataset, we can see that the tasks on average take much longer to run (around 12 seconds each versus 1 second each) and the GC time is also much longer on average (100s of ms versus 10s of ms). However, the shuffle read size / records is much lower (close to 0 B vs a consistent 2.3 MiB) while the shuffle write size / records is much higher (28 MiB versus 1.2 MiB).

Large Dataset 60 Repartitions

Applications: 0 Running, 18 Completed
Drivers: 0 Running, 0 Completed
Status: ALIVE

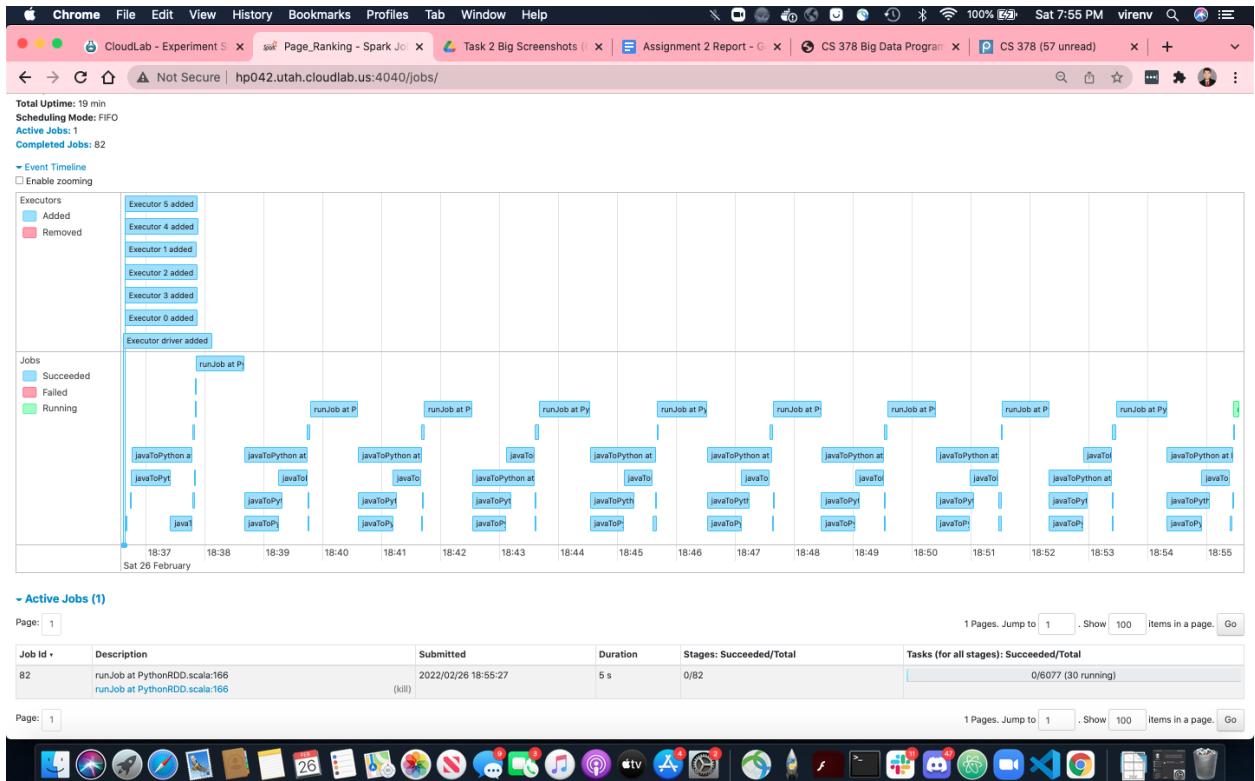
▼ Workers (3)

Worker Id	Address	State	Cores	Memory	Resources
worker-2022022511033-128.110.218.113-40187	128.110.218.113:40187	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-2022022511034-128.110.218.81-34071	128.110.218.81:34071	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-2022022511034-128.110.218.87-36219	128.110.218.87:36219	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	

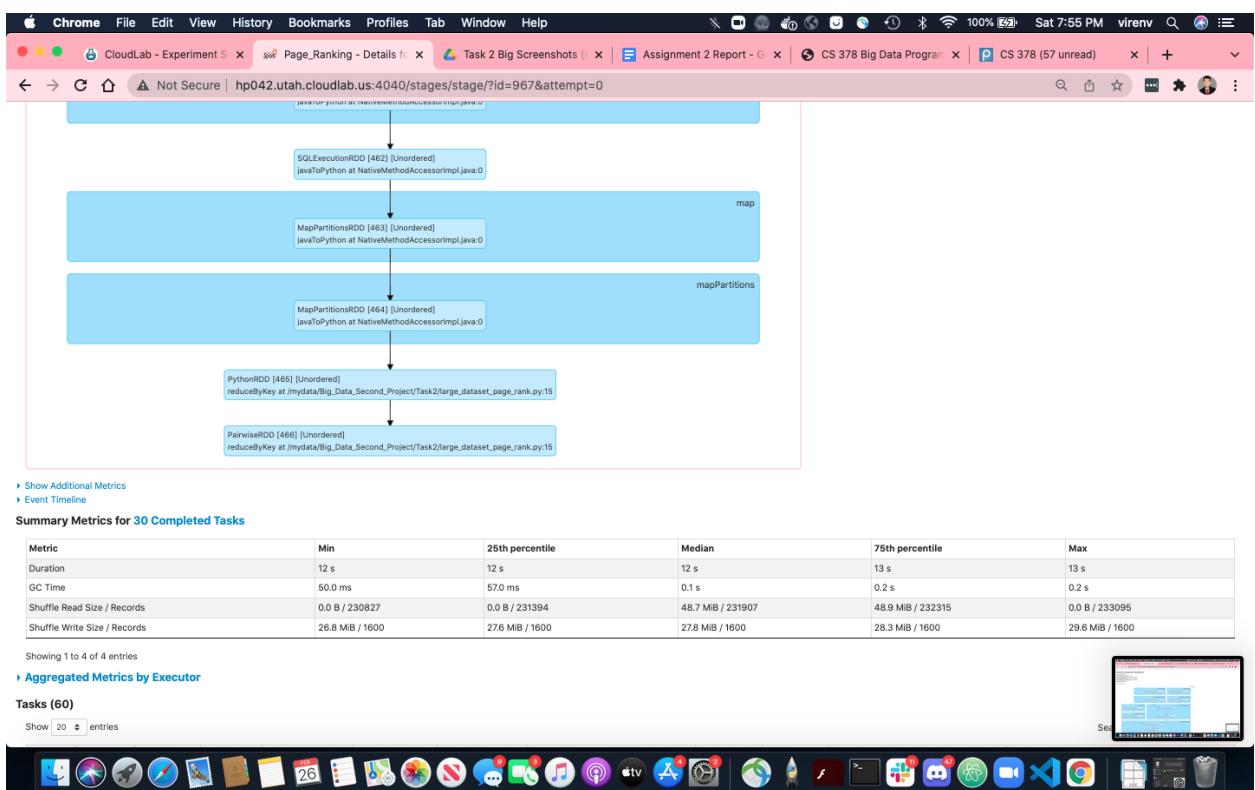
▼ Running Applications (0)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
▼ Completed Applications (18)								
Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
app-20220226204158-0017	Page_Ranking	30	29.0 GiB		2022/02/26 20:41:58	pa8789	FINISHED	1.6 min
app-20220226183636-0016	Page_Ranking	30	29.0 GiB		2022/02/26 18:36:36	pa8789	FINISHED	20 min
app-20220226183235-0015	Page_Ranking	30	29.0 GiB		2022/02/26 18:32:35	pa8789	FINISHED	1.5 min
app-20220226174912-0014	Page_Ranking	30	29.0 GiB		2022/02/26 17:49:12	pa8789	FINISHED	21 min
app-20220226174813-0013	Page_Ranking	30	29.0 GiB		2022/02/26 17:48:13	pa8789	FINISHED	8 s
app-20220226174249-0012	Page_Ranking	30	29.0 GiB		2022/02/26 17:42:49	pa8789	FINISHED	1.4 min
app-2022022617211-0011	Page_Ranking	30	29.0 GiB		2022/02/26 17:21:11	pa8789	FINISHED	20 min
app-20220226170213-0010	Page_Ranking	30	29.0 GiB		2022/02/26 17:02:13	pa8789	FINISHED	1.3 min
app-20220226165902-0009	Page_Ranking	30	29.0 GiB		2022/02/26 16:59:02	pa8789	FINISHED	1.3 min
app-20220226161116-0008	Page_Ranking	30	29.0 GiB		2022/02/26 16:11:16	pa8789	FINISHED	3.4 min
app-20220225121650-0007	Page_Ranking	30	29.0 GiB		2022/02/25 12:16:50	pa8789	FINISHED	20 min
app-20220225120404-0006	Page_Ranking	30	29.0 GiB		2022/02/25 12:10:40	pa8789	FINISHED	1.3 min
app-20220225120630-0005	Page_Ranking	30	29.0 GiB		2022/02/25 12:06:30	pa8789	FINISHED	1.3 min
app-20220225120005-0004	Page_Ranking	30	29.0 GiB		2022/02/25 12:00:05	pa8789	FINISHED	1.2 min
app-20220225115633-0003	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:33	pa8789	FINISHED	1.3 min
app-20220225115600-0002	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:00	pa8789	FINISHED	2 s

(Total time for application was 20 minutes, so a little better than one with 30 repartitions, but essentially the same as application without any repartitions. Goes to show that repartitions for this kind of application give back pretty much the same or slightly worse results)



(Total number of jobs for Application)



(Lineage graph and summary stats)

Takeaway: We observed that despite doubling the number of partitions, the completion time for PageRank with the large dataset was not significantly affected. As a matter of fact, it was rather similar to the time from Task 1 at around 20 minutes, about 5% faster than the completion time with 30 partitions. The network shuffling overhead that results from partitioning could have been slightly negated by simply having more partitions across more nodes at work, which could subsequently lead to a greater leveraging of parallelism to compensate for network latency. We see that 82 jobs were completed with this particular run, slightly more than the 76 jobs with 30 partitions and almost double the job count for the large dataset in Task 1. The average completion time for each job is relatively similar to running the large dataset with 30 partitions at 12 seconds on average.

Task 3:

Small Dataset Analysis

The screenshot shows a Chrome browser window with several tabs open. The tabs include "CloudLab - Experiment S", "Spark Master at spark://", "Task 2 Big Screenshots", "Assignment 2 Report", "CS 378 Big Data Program", "CS 378 (57 unread)", and a new tab. The main content area displays experimental statistics and application logs.

Applications: 0 Running, 18 Completed
Drivers: 0 Running, 0 Completed
Status: ALIVE

Workers (3)

Worker Id	Address	State	Cores	Memory	Resources
worker-2022022511033-128.110.218.113-40187	128.110.218.113:40187	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-2022022511034-128.110.218.81-34071	128.110.218.81:34071	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	
worker-2022022511034-128.110.218.87-36219	128.110.218.87:36219	ALIVE	20 (0 Used)	61.8 GiB (0.0 B Used)	

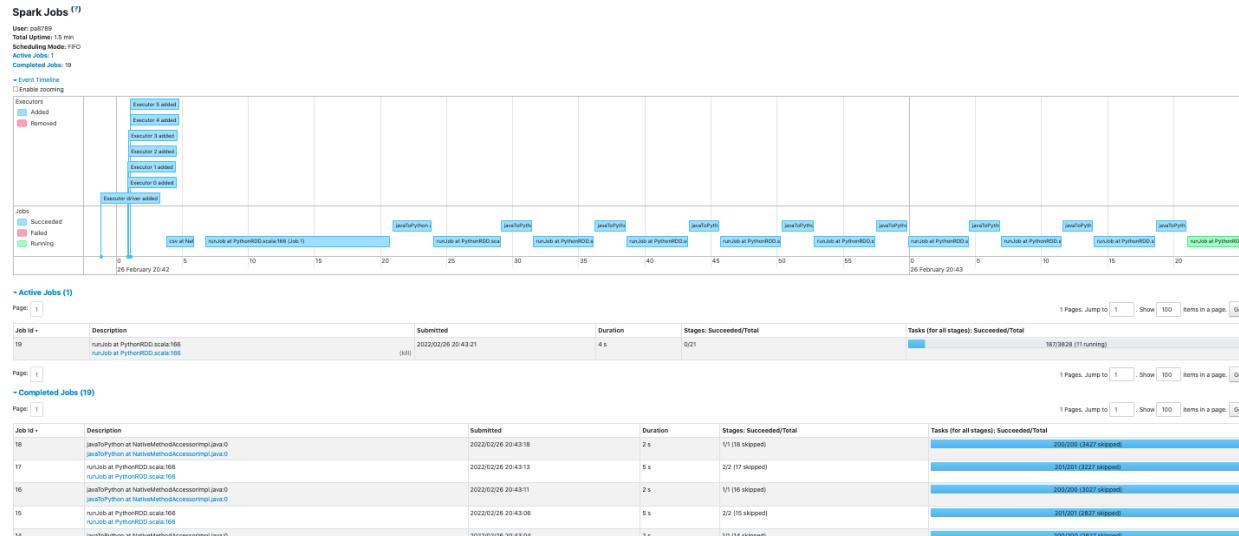
Running Applications (0)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration

Completed Applications (18)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
app-20220226040158-0017	Page_Ranking	30	29.0 GiB		2022/02/26 20:41:58	pa8789	FINISHED	1.6 min
app-20220226183636-0016	Page_Ranking	30	29.0 GiB		2022/02/26 18:36:36	pa8789	FINISHED	20 min
app-20220226183235-0015	Page_Ranking	30	29.0 GiB		2022/02/26 18:32:35	pa8789	FINISHED	1.5 min
app-20220226174912-0014	Page_Ranking	30	29.0 GiB		2022/02/26 17:49:12	pa8789	FINISHED	21 min
app-20220226174813-0013	Page_Ranking	30	29.0 GiB		2022/02/26 17:48:13	pa8789	FINISHED	8 s
app-20220226174249-0012	Page_Ranking	30	29.0 GiB		2022/02/26 17:42:49	pa8789	FINISHED	1.4 min
app-20220226171211-0001	Page_Ranking	30	29.0 GiB		2022/02/26 17:12:11	pa8789	FINISHED	20 min
app-20220226170213-0010	Page_Ranking	30	29.0 GiB		2022/02/26 17:02:13	pa8789	FINISHED	1.3 min
app-20220226165902-0009	Page_Ranking	30	29.0 GiB		2022/02/26 16:59:02	pa8789	FINISHED	1.3 min
app-20220225161116-0008	Page_Ranking	30	29.0 GiB		2022/02/25 16:11:16	pa8789	FINISHED	3.4 min
app-20220225121650-0007	Page_Ranking	30	29.0 GiB		2022/02/25 12:16:50	pa8789	FINISHED	20 min
app-20220225121040-0006	Page_Ranking	30	29.0 GiB		2022/02/25 12:10:40	pa8789	FINISHED	1.3 min
app-20220225120630-0005	Page_Ranking	30	29.0 GiB		2022/02/25 12:06:30	pa8789	FINISHED	1.3 min
app-20220225120005-0004	Page_Ranking	30	29.0 GiB		2022/02/25 12:00:05	pa8789	FINISHED	1.2 min
app-20220225115633-0003	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:33	pa8789	FINISHED	1.3 min
app-20220225115600-0002	Page_Ranking	30	29.0 GiB		2022/02/25 11:56:00	pa8789	FINISHED	2 s

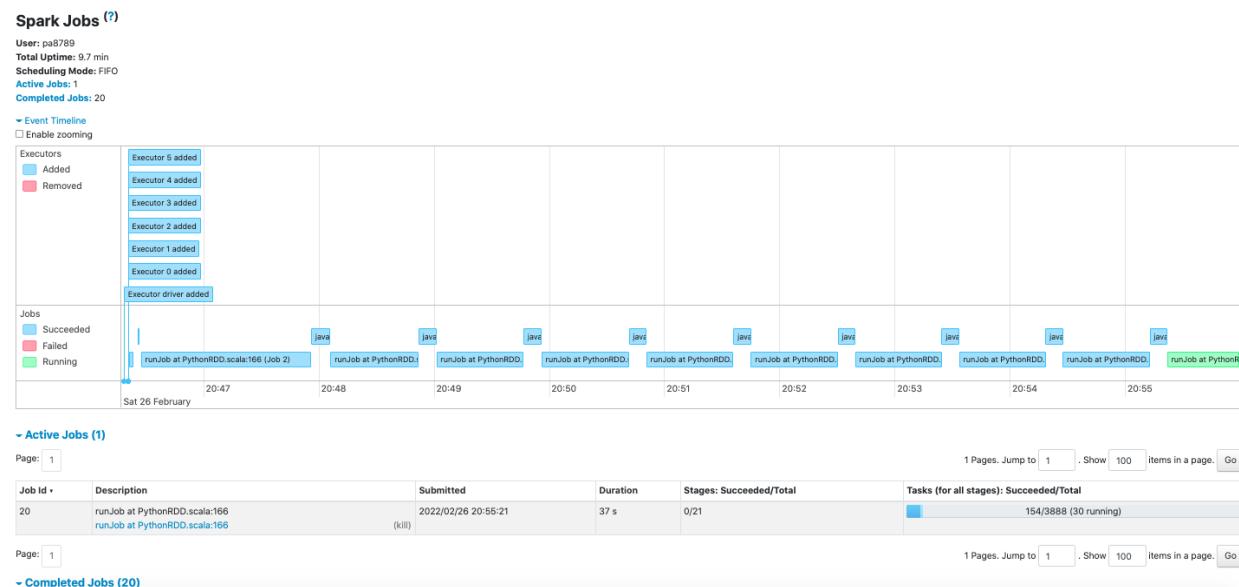
(Completion Time)



(Job/Application Timeline)

Takeaway: It appears that even though we were caching/persisting, the dataset is so small that it isn't worth doing so as it makes the application run slower at 1.6 minutes as opposed to 1.3. The number of jobs run though was half of before as only 15 jobs were needed to be completed as opposed to around 30 with no caching/persisting. The degradation in completion time could be due to some additional overhead caused by creating and maintaining the links dataset in cached copies.

Big Dataset Analysis



(Half the jobs as without persisting. Forty originally, now only 20 completed jobs, job completion time was recorded as 10 minutes, half of the time as without the cache)

Not Secure | hp042.utah.cloudlab.us:4040/stages/stage/?id=3&attempt=0

[Show Additional Metrics](#)

[Event Timeline](#)

Summary Metrics for 30 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	7 s	9 s	9 s	10 s	10 s
GC Time	0.0 ms	59.0 ms	77.0 ms	95.0 ms	0.2 s
Input Size / Records	37.7 MiB / 10	39.3 MiB / 10	40.3 MiB / 11	41.8 MiB / 11	46.7 MiB / 12
Shuffle Read Size / Records	0.0 B / 5504	0.0 B / 5606	0.0 B / 5667	0.0 B / 5692	0.0 B / 5771
Shuffle Write Size / Records	14.9 MiB / 2995	15.5 MiB / 2998	15.7 MiB / 2999	16.1 MiB / 2999	17.1 MiB / 3000

Showing 1 to 5 of 5 entries

[Aggregated Metrics by Executor](#)

Tasks (61)

Show 20 entries

Search:

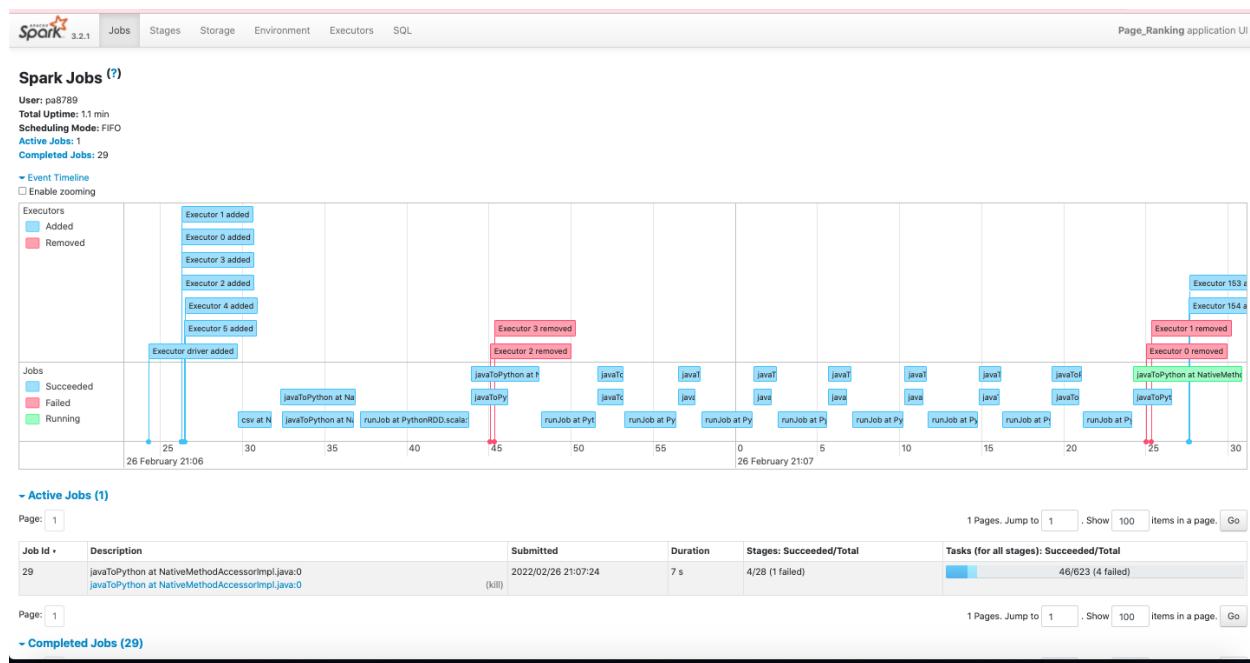
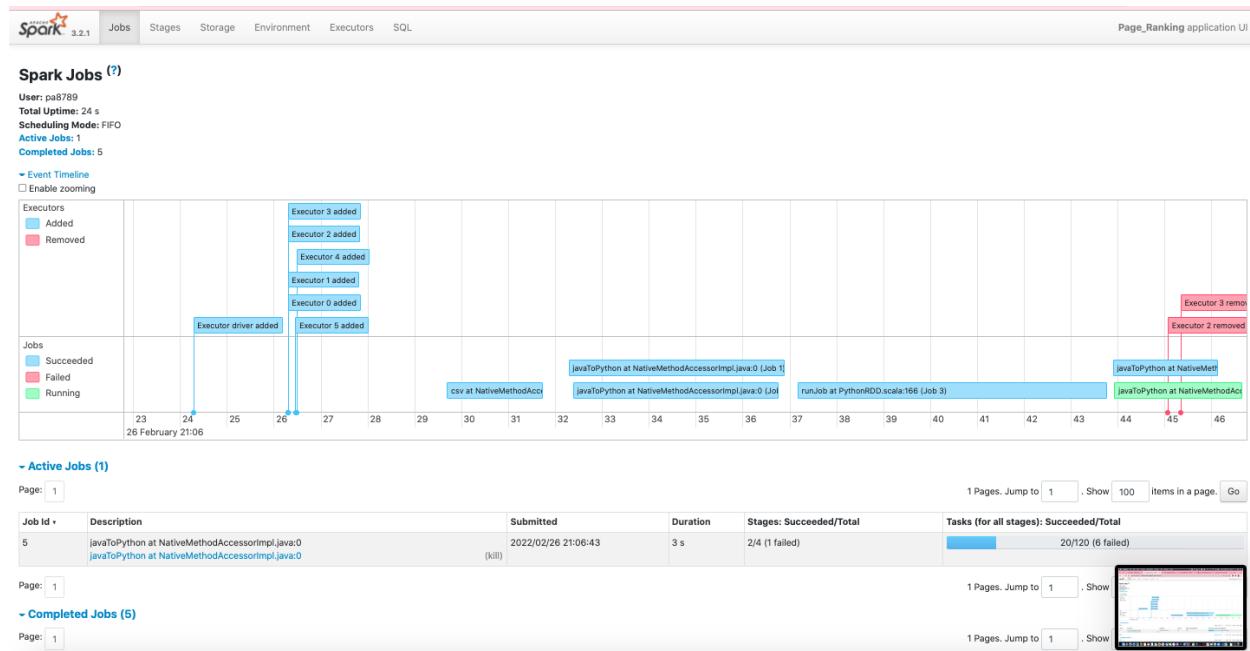
Index	Task ID	Attempt	Status	Locality level	Executor ID	Host	Logs	Launch Time	Duration	GC Time	Input Size / Records	Shuffle Write Size / Records	Shuffle Read Size / Records	Errors
0	160	0	SUCCESS	PROCESS_LOCAL	0	128.110.218.87	stdout stderr	2022-02-26 21:47:05	9 s	80.0 ms	40.3 MiB / 11	15.7 MiB / 2999	21.1 MiB / 5756	
1	161	0	SUCCESS	PROCESS_LOCAL	1	128.110.218.87	stdout stderr	2022-02-26 21:47:05	9 s	0.2 s	40.5 MiB / 11	15.7 MiB / 2998	21.2 MiB / 5571	
2	162	0	SUCCESS	PROCESS_LOCAL	5	128.110.218.81	stdout stderr	2022-02-26 21:47:05	9 s	77.0 ms	40.8 MiB / 11	16.1 MiB / 2998	21.5 MiB / 5692	
3	163	0	SUCCESS	PROCESS_LOCAL	4	128.110.218.81	stdout stderr	2022-02-26 21:47:05	10 s	65.0 ms	41.8 MiB / 11	16.1 MiB / 2997	21.6 MiB / 5801	

(summary metrics for big dataset after persisting links)

Takeaway: We see the major difference caching/persisting has on a large dataset such as this one. It only took 10 minutes for the application to run to completion as opposed to 20 minutes, a 50% reduction. The number of jobs also run was about half as before as only 20 jobs were needed to be completed. This could have been the result of some computations or HDFS reads that were eliminated along the way as a result of caching, so any part of the lineage that was originally cached could have been quickly accessed if any recomputation was needed.

Task 4

Small Dataset Analysis



Spark Jobs (?)

User: pa8789
Total Uptime: 1.7 min
Scheduling Mode: FIFO
Active Jobs: 1
Completed Jobs: 29

Event Timeline
Enable zooming

Executors:
Added (Blue)
Removed (Red)

Jobs:
Succeeded (Blue)
Failed (Red)
Running (Green)

Timeline: 26 February 21:06 to 26 February 21:07

Active Jobs (1)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
29	javaToPython at NativeMethodAccessImpl.java:0 javaToPython at NativeMethodAccessImpl.java:0 (kill)	2022/02/26 21:07:24	38 s	26/28 (1 failed)	308/623 (4 failed)

Completed Jobs (29)

(Can clearly see the killed workers highlighted in red, note the increased total time. Overall completion time ended up being 1.8 minutes)

CloudLab - Experiment S | Page_Ranking - Executor

Show Additional Metrics
Select All
On Heap Memory
Off Heap Memory
Peak JVM Memory OnHeap / OffHeap
Peak Execution Memory OnHeap / Off-heap
Peak Storage Memory OnHeap / OffHeap
Peak Pool Memory Direct / Mapped
Resources
Resource Profile Id
Exec Loss Reason

Summary

	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Excluded
Active(5)	0	2 MB / 76.5 GiB	0.0 B	20	5	1	553	559	10 min (19 s)	1.1 GiB	609.8 MiB	760.9 MiB	0
Dead(4)	0	1.4 MiB / 61.2 GiB	0.0 B	20	-6	22	354	370	7.5 min (16 s)	610 MiB	432.3 MiB	474.5 MiB	0
Total(9)	0	3.4 MiB / 137.6 GiB	0.0 B	40	-1	23	907	929	18 min (34 s)	1.7 GiB	1 GiB	1.2 GiB	0

Executors

Executor ID	Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Logs	Thread Dump
0	128.110.218.87-43467	Dead	0	554 KiB / 15.3 GiB	0.0 B	5	-1	9	153	161	3.1 min (6 s)	254.9 MiB	196.9 MiB	209.3 MiB	stdout stderr	Thread Dump
driver	hp042.utah.cloudlab.us:33315	Active	0	656.5 KiB / 15.3 GiB	0.0 B	0	0	0	0	0	0.0 ms (0.0 ms)	0.0 B	0.0 B	0.0 B		Thread Dump
1	128.110.218.87-43309	Dead	0	593.9 KiB / 15.3 GiB	0.0 B	5	-2	8	176	182	3.3 min (6 s)	273.9 MiB	228 MiB	231.5 MiB	stdout stderr	Thread Dump
2	128.110.218.113-42625	Dead	0	140.6 KiB / 15.3 GiB	0.0 B	5	0	0	13	13	33 s (3 s)	40.6 MiB	4.4 MiB	17.6 MiB	stdout stderr	Thread Dump
3	128.110.218.113-44011	Dead	0	140.6 KiB / 15.3 GiB	0.0 B	5	-3	5	12	14	32 s (2 s)	40.6 MiB	2.9 MiB	16.2 MiB	stdout stderr	Thread Dump
4	128.110.218.81-41957	Active	0	561.7 KiB / 15.3 GiB	0.0 B	5	0	0	253	253	4.3 min (6 s)	477.2 MiB	302.3 MiB	351.6 MiB	stdout stderr	Thread Dump
5	128.110.218.81-33659	Active	0	623.9 KiB / 15.3 GiB	0.0 B	5	5	1	270	276	4.4 min (5 s)	594.8 MiB	292.7 MiB	365.4 MiB	stdout stderr	Thread Dump
153	128.110.218.87-38917	Active	0	95.8 KiB / 15.3 GiB	0.0 B	5	0	0	15	15	46 s (4 s)	47.8 MiB	7.4 MiB	23.7 MiB	stdout stderr	Thread Dump
154	128.110.218.87-44403	Active	0	161 KiB / 15.3 GiB	0.0 B	5	0	0	15	15	43 s (3 s)	47.8 MiB	7.4 MiB	20.1 MiB	stdout stderr	Thread Dump

(Executor Statistics)

Summary Metrics for 33 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	0.2 s	0.2 s	0.3 s	0.3 s	0.3 s
GC Time	0.0 ms	0.0 ms	0.0 ms	0.0 ms	0.0 ms
Shuffle Read Size / Records	910.7 KIB / 430	925.4 KIB / 430	932.2 KIB / 430	937.5 KIB / 430	958.1 KIB / 430
Shuffle Write Size / Records	333.6 KIB / 20514	335.8 KIB / 20698	336.9 KIB / 20760	338 KIB / 20835	340.2 KIB / 21000

Showing 1 to 4 of 4 entries

Aggregated Metrics by Executor

Tasks (33)

Index	Task ID	Attempt	Status	Locality level	Executor ID	Host	Logs	Launch Time	Duration	GC Time	Shuffle Write Size / Records	Shuffle Read Size / Records
0	803	0	SUCCESS	NODE_LOCAL	4	128.110.218.81	stdout	2022-02-26 22:07:36	0.3 s	339.3 KIB / 20909	939.3 KIB / 430	

(Summary statistics for small dataset)

Takeaway: It is not a surprise at all that killing a worker at both 25% and 75% lifetime mark of application would lead to a much slower runtime. We can see the effect in the two timelines when a worker was killed as shown by the removal of a couple of the executors in red. This led to some failed tasks which can be seen in the Executor statistics screenshot. It took the application 1.8 minutes to reach completion which is about a 38% increase as the baseline application from Task 1 only took 1.3 minutes. The number of jobs was the same as the baseline application, with just a longer completion time as expected.

Big Dataset Analysis

Spark Master at spark://hp042.utah.cloudlab.us:7077

URL: spark://hp042.utah.cloudlab.us:7077

Alive Workers: 1

Cores in use: 20 Total: 0 Used

Memory in use: 61.8 Gib Total: 0.0 8 Used

Resources in use:

Applications: 0 Running, 1 Completed

Drivers: 0 Running, 0 Completed

Status: ALIVE

Workers (1)

Worker Id	Address	State	Cores	Memory	Resources
worker-2022022613729-128.110.218.81-34699	128.110.218.81-34699	ALIVE	20 (0 Used)	61.8 Gib (0.0 8 Used)	

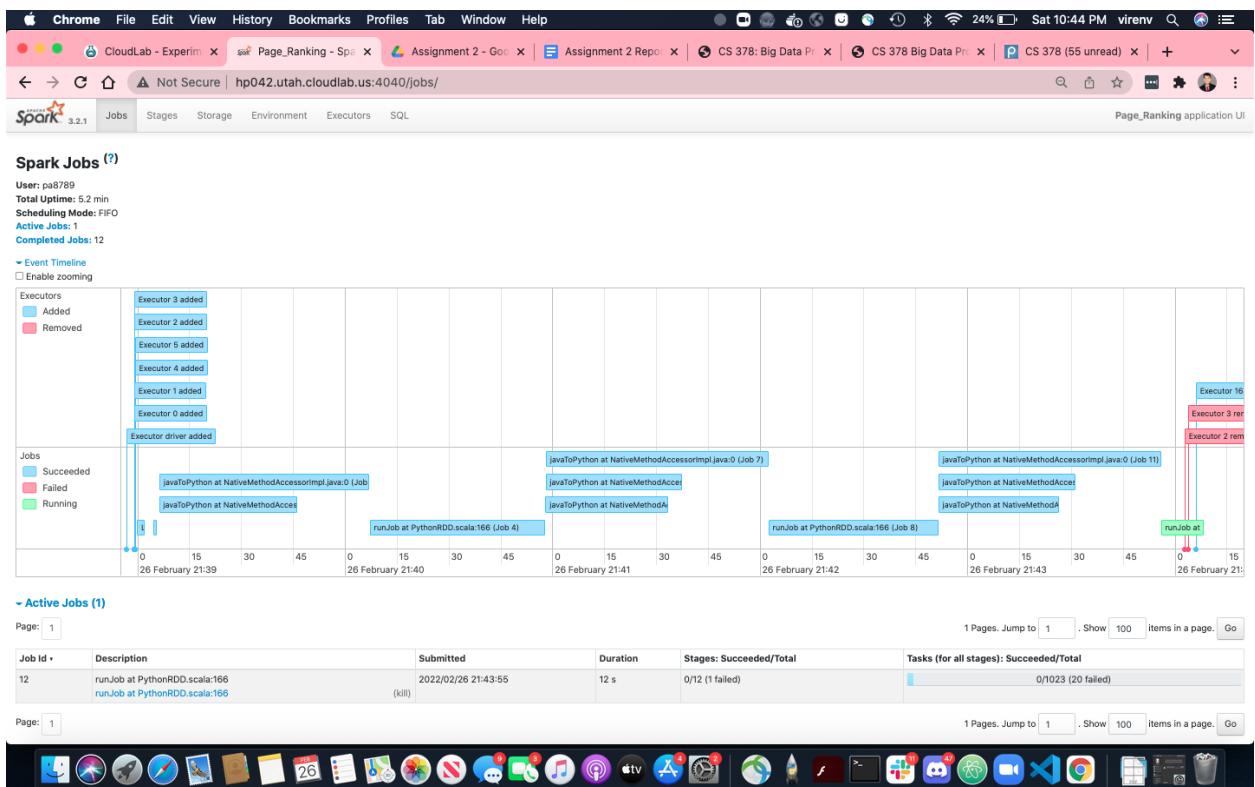
Running Applications (0)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration

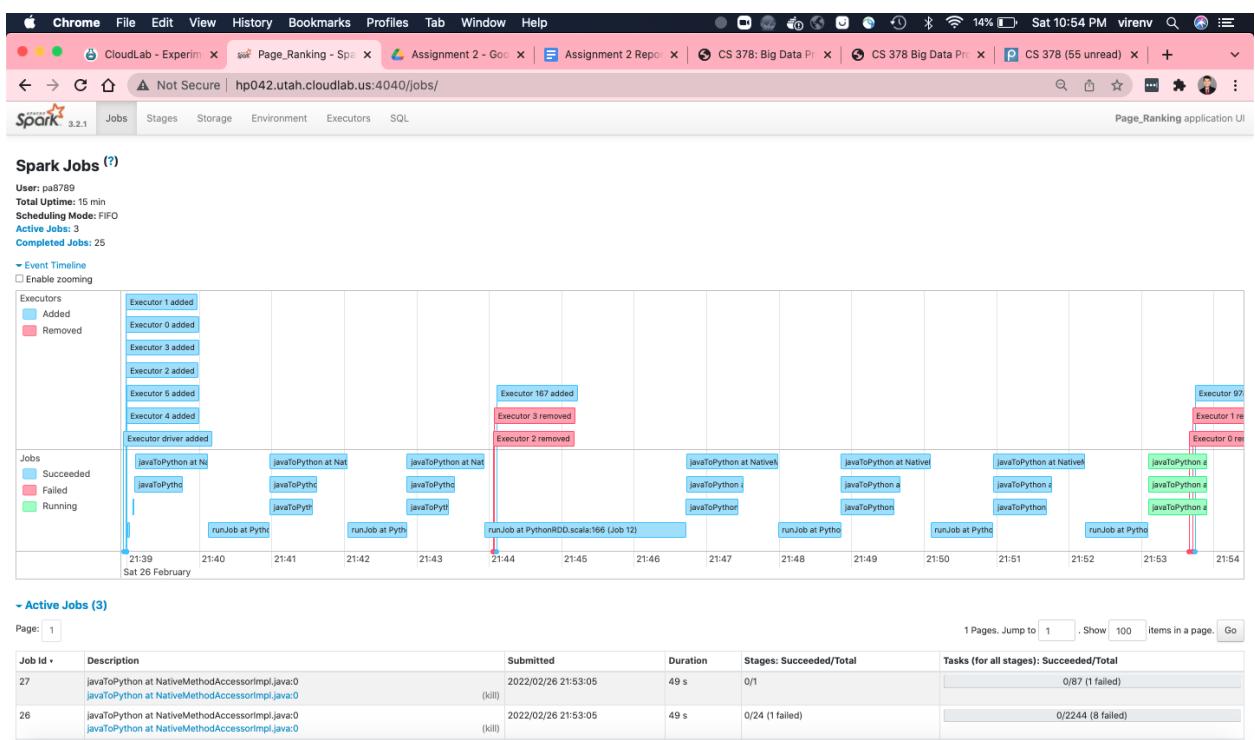
Completed Applications (1)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
app-2022022613856-0000	page.Ranking	10	29.0 Gib		2022/02/26 21:38:56	pa8789	FINISHED	31 min

(Completion time of 31 minutes)



(Timeline at 5 min, killed worker shown)



(Timeline at 15 minutes when second worker is killed)

CloudLab - Experiment | Page_Ranking - Spark | Task 4 Big Screens! | Assignment 2 Repo | CS 378: Big Data Pr | CS 378 Big Data Pr | CS 378 (55 unread) | + | Sat 11:09 PM virenv

Not Secure | hp042.utah.cloudlab.us:4040/jobs/

Spark Jobs (2)

User: pa8799
Total Uptime: 31 min
Scheduling Mode: FIFO
Active Jobs: 2
Completed Jobs: 41

Event Timeline (Enable zooming)

Executors: Added (Blue), Removed (Red)

Jobs: Succeeded (Blue), Failed (Red), Running (Green)

Sat 26 February 21:39 21:40 21:41 21:42 21:43 21:44 21:45 21:46 21:47 21:48 21:49 21:50 21:51 21:52 21:53 21:54 21:55 21:56 21:57 21:58 21:59 22:00 22:01 22:02 22:03 22:04 22:05 22:06 22:07 22:08 22:09

Active Jobs (2)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total	
42	csv at NativeMethodAccessorImpl.java:0 csv at NativeMethodAccessorImpl.java:0	(kill)	2022/02/26 22:09:22	9 s	0/1	0/87
41	csv at NativeMethodAccessorImpl.java:0 csv at NativeMethodAccessorImpl.java:0	(kill)	2022/02/26 22:09:22	9 s	0/40	69/3740 (21 running)

Completed Jobs (41)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
40	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
39	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
38	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
37	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
36	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
35	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
34	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
33	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
32	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
31	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
30	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
29	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
28	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
27	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
26	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
25	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
24	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
23	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
22	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
21	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
20	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
19	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
18	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
17	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
16	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
15	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
14	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
13	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
12	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
11	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
10	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
9	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
8	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
7	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
6	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
5	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
4	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
3	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
2	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
1	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1
0	runJob at PythonRDD.scala:16	2022/02/26 22:09:22	9 s	0/1	0/1

(Overall Timeline)

CloudLab - Experiment | Page_Ranking - Spark | Assignment 2 - Go | Assignment 2 Repo | CS 378: Big Data Pr | CS 378 Big Data Pr | CS 378 (55 unread) | + | Sat 11:06 PM virenv

Not Secure | hp042.utah.cloudlab.us:4040/executors/

Executors

Show Additional Metrics

- On Heap
- Off Heap Memory
- Off Off Heap
- Peak JVM Memory OnHeap / OffHeap
- Peak Execution Memory OnHeap / OffHeap
- Peak Storage Memory OnHeap / OffHeap
- Peak Pool Memory Direct / Mapped
- Resources
- Resource Profile Id
- Error Loss Reason

Summary

Active	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Excluded
Active(9)	0	3.7 MB / 76.5 GB	0.0 B	20	23	19	3255	3291	7.1 h (3.3 min)	173.5 GB	62.7 GB	66.1 GB	0
Dead(4)	0	2.1 MB / 61.2 GB	0.0 B	20	-12	34	1250	1272	3.1 h (1.7 min)	72.6 GB	21.4 GB	26.3 GB	0
Total(9)	0	5.8 MB / 137.6 GB	0.0 B	40	11	47	4505	4563	10.1 h (5.0 min)	2461.0 GB	84.1 GB	92.4 GB	0

Executors

Show 20 entries

Executor ID	Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Logs	Thread Dump
0	128.110.218.87:45107	Dead	0	871.5 KIB / 15.3 GB	0.0 B	5	-1	6	463	468	1.1 h (36 s)	26.1 GB	8.5 GB	9.7 GB	stdout	Thread Dump
driver	hp042.utah.cloudlab.us:43423	Active	0	759.7 KIB / 15.3 GB	0.0 B	0	0	0	0	0	0.0 ms (0.0 ms)	0.0 B	0.0 B	0.0 B		Thread Dump
1	128.110.218.87:38397	Dead	0	525.3 KIB / 15.3 GB	0.0 B	5	-11	18	458	465	1.2 h (43 s)	26.9 GB	8.2 GB	10.2 GB	stdout	Thread Dump
2	128.110.218.113:39855	Dead	0	388.6 KIB / 15.3 GB	0.0 B	5	0	5	169	174	23 min (14 s)	9.9 GB	2.4 GB	3.2 GB	stdout	Thread Dump
3	128.110.218.113:39823	Dead	0	388.6 KIB / 15.3 GB	0.0 B	5	0	6	160	165	23 min (11 s)	9.6 GB	2.3 GB	3.2 GB	stdout	Thread Dump
4	128.110.218.81:42491	Active	0	747.2 KIB / 15.3 GB	0.0 B	5	6	6	914	922	2.2 h (1.2 min)	48 GB	17.8 GB	19.0 GB	stdout	Thread Dump
5	128.110.218.81:43345	Active	0	753.4 KIB / 15.3 GB	0.0 B	5	6	6	977	984	2.2 h (43 s)	47.5 GB	19.3 GB	19.2 GB	stdout	Thread Dump
167	128.110.218.113:42801	Active	0	747.2 KIB / 15.3 GB	0.0 B	5	6	10	894	910	1.8 h (48 s)	48.6 GB	17.5 GB	18.5 GB	stdout	Thread Dump
978	128.110.218.87:36471	Active	0	753.4 KIB / 15.3 GB	0.0 B	5	5	0	470	475	59 min (32 s)	29.5 GB	8.1 GB	9.4 GB	stdout	Thread Dump

(Executor Statistics)

CloudLab - Experiment | Page_Ranking - Det | Assignment 2 - Goo | Assignment 2 Repo | CS 378: Big Data Pr | CS 378 Big Data Pr | CS 378 (55 unread) | + | Sat 11:06 PM virenv

Not Secure | hp042.utah.cloudlab.us:4040/stages/stage/?id=342&attempt=0

Summary Metrics for 40 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	9 s	10 s	12 s	12 s	13 s
GC Time	16.0 ms	19.0 ms	30.0 ms	54.0 ms	56.0 ms
Shuffle Read Size / Records	0.0 B / 231427	48.6 MB / 231972	48.9 MB / 232409	49.2 MB / 232754	49.6 MB / 233641
Shuffle Write Size / Records	26.7 MB / 1600	27.6 MB / 1600	27.8 MB / 1600	28.1 MB / 1600	28.6 MB / 1600

Showing 1 to 4 of 4 entries

Aggregated Metrics by Executor

Executor ID	Logs	Address	Task Time	Total Tasks	Failed Tasks	Killed Tasks	Succeeded Tasks	Excluded	Shuffle Read Size / Records	Shuffle Write Size / Records
4	stdout stderr	128.110.218.81:42491	2.1 min	10	0	0	10	false	490.9 MB / 2324406	279 MB / 16000
5	stdout stderr	128.110.218.81:43345	2.1 min	10	0	0	10	false	734.6 MB / 3473051	276.6 MB / 16000
167	stdout stderr	128.110.218.113:42801	1.7 min	10	0	0	10	false	740.4 MB / 3487169	278.7 MB / 16000
978	stdout stderr	128.110.218.87:36471	1.7 min	10	0	0	10	false	742 MB / 3487467	280.9 MB / 16000

Search:

(Summary Statistics)

CloudLab - Experiment | Page_Ranking - Det | Assignment 2 - Goo | Assignment 2 Repo | CS 378: Big Data Pr | CS 378 Big Data Pr | CS 378 (55 unread) | + | Sat 11:06 PM virenv

Not Secure | hp042.utah.cloudlab.us:4040/stages/stage/?id=342&attempt=0

Total Time Across All Tasks: 7.5 min
Locality Level Summary: Node local: 60
Shuffle Read Size / Records: 2.5 GB / 11852932
Shuffle Write Size / Records: 1115.2 MB / 64000
Associated Job Ids: 36

DAG Visualization

Stage 342

(Interesting Lineage graph right after 25% mark)

Takeaway: Like for the small dataset, the loss of workers led to a larger application runtime. However, for the large dataset it was even worse as the application completion time was 31 minutes, a 55% increase compared to the baseline. As seen from the executor statistics, there were even more failed tasks that took place. The jobs themselves seemed larger and more complex as seen by the above DAG visualization for one of the jobs. Like for the small dataset, the number of jobs completed was the same as the baseline application at a little over 40 jobs, but it just took longer.

Specific Contributions of each Group Member

Note: In general, all of the members worked together on all aspects and would really only work on the project when everyone was present. However, each member specialized in driving some aspect of the project.

Pranav: Pranav was mainly in charge of implementing the actual page ranking algorithm and determining the logic of what needed to be done. There was pseudo code given in a prior lecture written in Scala and not handling preprocessing. Since this project was done in Python, some effort needed to be made in terms of looking at the PySpark APIs and figuring out how the code should be written to meet the requirements of the algorithm as detailed on the assignment page. Pranav had the most experience with Scala and Spark so it was easier for him to lead the coding aspect due to his prior experience. However, coding was still largely a group effort.

Viren: Viren was the one primarily in charge of setting up Hadoop and Spark and running the experiment (s). Most of the time when we met up, his screen was shared as he drove most of the time. He also set up the repository on Github and helped in organizing a lot of the things, such as the run scripts for the different tasks as well as the README files.

Sameer: Sameer was responsible primarily for gathering and compiling the necessary evidence (including screenshots) and organizing them within the report. He assisted with communicating various trends and interesting features from each of the tasks to the rest of the team and helped drive some of the discussions surrounding the causes of these events and how to move forward with any further tests.