

# COEN242 PA2 - Top K words in a big dataset

Hui-yu Liu, Ching-Yueh Huang

## I. Introduction

In the previous work, we used python to find top K words in a big dataset, including 300MB, 2.5GB and 16GB. This report is another way using Hadoop MapReduce to implement it. The goal is to find the Top 100 frequent words and the words having more than 6 characters using MapReduce. In addition, we will be utilizing the `nlargest` function from Python's `heapq` module to efficiently analyze large datasets and gain insights into the most commonly used words in a given text corpus.

This report aims to introduce the method we utilized in our function and setting files. We will first provide a detailed explanation of the method and how it was implemented. In addition, we will discuss some observations and optimization techniques that we adopted during the testing phases to improve the performance of our method. Furthermore, we will summarize the results obtained from our experiments and provide insights into the effectiveness of our approach. Finally, we will conclude our work by discussing the implications of our findings and potential future improvements.

## II. Methods

To perform word count using Hadoop MapReduce, we began by creating a mapper to split each line of text into individual words and convert them to lowercase. This was done because all the stopwords provided were in lowercase. By converting all the words to lowercase, we made it easier to identify the stopwords. The output of the mapper was a set of key-value pairs where the word was the key and the value was set to 1. To optimize performance, we also implemented a combiner to aggregate the intermediate results for each key before they were sent to the reducer.

To meet the requirement of retrieving the top 100 words with each word having more than 6 characters, we wrote another mapper program. Here, we also converted each word to lowercase and counted the number of occurrences of each word using counters. After aggregating the results, we passed them to the reducer

for further processing of the data. In the reducer program, we created a dictionary to store the key-value pairs. Moreover, we implemented 'heapq' module with 'nlargest' function, which is used to return the n largest elements from the key-value dataset. The algorithm of heap queue works by maintaining the the heap size, and iterating over the collection, adding each element to the heap and then popping the smallest element off if the heap size exceeds the given number. This algorithm has the time complexity of  $O(\log n)$ , proving that it is an efficient way to maintain a heap in memory. In terms of space efficiency, since only a small portion fo the input data needs to be stored in memory at any given time, the heap queue is space efficient. To get the top k element from the input data, we customized our heapq algorithm with 'nlargest' function, which allowed us to sort the data based on a specified key and handle any ties between elements. This approach proved to be the effective solution for sorting and filtering data in Python.

Next, we used Hadoop Streaming to run the MapReduce job on a Hadoop cluster. We specified the number of reducers and partitioned the data based on the first two characters of each word. And end up writing a reducer function to sum up the value of each key to get the total count for each word. In addition, we changed some variables in mapred-site.xml to optimize the whole process. For example, we increase the 'mapreduce.task.io.sort.mb' parameter which results in improved performance with more memory available for the map output buffer and reduced amount of data that needs to be spilled to disk during the sort and shuffle phase. It also reduced disk I/O. With more memory available for the map output buffer, the amount of data that needs to be spilled to disk during the sort and shuffle phase is reduced. This resulted in faster job completion times. Plus, the setting parameter increased parallelism with more memory available for sorting. This improved overall job throughput by allowing more tasks executing in parallel. However, we were aware of the tradeoffs involved in increasing this parameter, such as requiring more memory, which could be problematic if the cluster had limited resources., which can be a problem if the cluster has limited resources. So we carefully monitored the test result and adjusted the the parameter by calculating the most suitable requirements based on our resources and the efficiency of MapReduce tasks accordingly.

### III. Analysis

#### A. MapReduce Optimization

When Map starts generating output, it does not simply write the data to disk, as frequent disk operations can cause severe performance degradation. It is more complex in that the data is first written to a buffer in memory and some pre-sorting is done to improve efficiency.

Each Map task has a circular memory buffer used to write the output data. The default size of this buffer is 100MB, and the exact size can be set via the `io.sort.mb` property. When the amount of data in the buffer reaches a specific threshold (`io.sort.mb * io.sort.spill.percentage`, where `io.sort.spill.percentage` is 0.80 by default), a background thread is started to spill the contents of the buffer to disk. During the spill process, the output of Map will continue to be written to the buffer, but if the buffer is full, Map will be blocked until the spill is complete. spill thread will do a second quick sort on the buffer before writing it to disk, first sorting by the partition the data belongs to, and then sorting by Key in each partition. in each partition. The output consists of an index file and a data file.

If a Combiner is set up, it will run on top of the sorted output, which is a Mini Reducer that runs on the node itself performing the Map task, first doing a simple Reduce on the Map output, making the Map output more compact, with less data being written to disk and transferred to the Reducer.

The spill file is stored in the directory specified by `mapred.local.dir` and deleted after the Map task is completed.

A new spill file is created every time the data in memory reaches the spill threshold, so there may be multiple spill files by the time the Map task finishes writing its last output record. Before the Map task finishes, all spill files will be merged and sorted into one index file and one data file. This is a multi-way merge process, and the maximum number of merges is controlled by `io.sort.factor` (default is 10). If Combiner is set and the number of spill files is at least 3 (controlled by the `min.num.spills.for.combine` property), then Combiner will be run to compress the data before the output files are written to disk.

Compressing the data written to disk is usually a good way to do this, as it makes the data write to disk faster, saves disk space, and reduces the amount of data that needs to be transferred to the Reducer. The default output is not compressed, but it can be enabled by simply setting `mapred.compress.map.output` to `true`. The library used for compression is set by `mapred.map.output.compression.codec`.

When the spill files are consolidated, Map will delete all temporary spill files and tell TaskTracker that the task is complete. The number of threads used to transfer the partitions data is controlled by `tasktracker.http.threads`, which is set for each TaskTracker, not a single Map, and defaults to 40.

For Intermediate data, this is always a good idea to use LZO compression. Every Hadoop job that generates a non-negligible amount of map output will benefit from intermediate data compression with LZO. Although LZO adds a little bit of overhead to the CPU, it saves time by reducing the amount of disk IO during the shuffle.

Name	Default	Changed
mapreduce.task.io.sort.mb	100	220
mapreduce.task.io.sort.factor	10	100
mapreduce.map.sort.spill.percent	0.8	0.9
mapred.map.child.java.opts	-Xmx200m	-Xmx1024m
mapred.reduce.child.java.opts	-Xmx1024m	-Xmx1024m
mapreduce.reduce.input.buffer.percent	0.0	0.2
mapred.job.shuffle.merge.percent	0.66	0.8
min.num.spills.for.combine	3	5

Data changed in `mapred-site.xml`

## B. Cases of experiment

In Hadoop, the input data is divided into blocks, typically 64MB or 128MB in size. For a 16GB input, assuming a block size of 128MB, you

would have approximately 110 blocks. Each block will be processed by a mapper. We will give some cases below to check how different cases work.

### 1) Case 1

When setting only one reducer and allowing the number of mappers to be defined by the data size, the system will automatically determine the optimal number of mappers based on the input data.

`hadoop jar`

`$HADOOP_HOME/share/hadoop/tools/lib/hadoop-streaming-3.3.5.jar`

`r \`

`-file <path_to_mapper.py> \`

`-mapper <path_to_mapper.py> \`

`-file <path_to_reducer.py> \`

`-reducer <path_to_reducer.py> \`

`-file <path_to_stopword.txt> \`

`-input <hadoop_path_to_data_16GB.txt> \`

`-output /output/16gb/case1`

Job Overview	
<b>Job Name:</b>	streamjob5080891716127857877.jar
<b>User Name:</b>	sandy
<b>Queue:</b>	default
<b>State:</b>	SUCCEEDED
<b>Uberized:</b>	false
<b>Submitted:</b>	Fri May 12 21:02:59 PDT 2023
<b>Started:</b>	Fri May 12 21:03:05 PDT 2023
<b>Finished:</b>	Fri May 12 21:50:25 PDT 2023
<b>Elapsed:</b>	47mins, 19sec
<b>Diagnostics:</b>	
<b>Average Map Time</b>	1mins, 19sec
<b>Average Shuffle Time</b>	24mins, 4sec
<b>Average Merge Time</b>	2mins, 12sec
<b>Average Reduce Time</b>	15mins, 8sec

ApplicationMaster			
Attempt Number	Start Time	Node	Logs
1	Fri May 12 21:03:01 PDT 2023	172.31.94.170:8042	logs

Task Type	Total		Complete
<b>Map</b>	115		115
<b>Reduce</b>	1		1
Attempt Type	Failed	Killed	Successful
<b>Maps</b>	0	1	115
<b>Reduces</b>	0	0	1

As you can see, for Case 1, there are 115 mappers and one reducer, we spend most of the time in shuffle and reduce. The total time of it is about 47min. We can also check the MapReduce Framework in Counter, the spilled records of mapper is more than 2,900,000,000. Each map spill for about 28MB, so that's why we tried to change `mapreduce.task.io.sort.mb` from default 100 to 200.

## 2) Case 2

To enhance the efficiency of the data processing workflow, we conducted an analysis to identify areas for improvement. In Case1, where no combiner was employed, the execution time for the Reduce phase was significant. In an effort to mitigate this, we introduced a combiner, which operates as an intermediate step between the mapper and reducer stages.

```
hadoop jar
$HADOOP_HOME/share/hadoop/tools/lib/hadoop-streaming-3.3.5.jar \
r \
-file <path_to_mapper.py> \
-mapper <path_to_mapper.py> \
-file <path_to_reducer.py> \
-reducer <path_to_reducer.py> \
-file <path_to_combiner.py> \
-combiner <path_to_combiner.py> \
-file <path_to_stopword.txt> \
-input <hadoop_path_to_data_16GB.txt> \
-output /output/16gb/case2
```

Job Overview	
Job Name:	streamjob2613114076023928026.jar
User Name:	sandy
Queue:	default
State:	SUCCEEDED
Uberized:	false
Submitted:	Sat May 13 09:52:06 PDT 2023
Started:	Sat May 13 09:52:13 PDT 2023
Finished:	Sat May 13 10:20:42 PDT 2023
Elapsed:	28mins, 29sec
Diagnostics:	
Average Map Time	1mins, 14sec
Average Shuffle Time	22mins, 18sec
Average Merge Time	5sec
Average Reduce Time	12sec

ApplicationMaster			
Attempt Number	Start Time	Node	Logs
1	Sat May 13 09:52:08 PDT 2023	<a href="#">172.20.205.169:8042</a>	<a href="#">logs</a>

Task Type	Total		Complete
Map	115		115
Reduce	1		1
Attempt Type	Failed	Killed	Successful
Maps	0	1	115
Reduces	0	0	1

By incorporating the combiner, we observed a remarkable reduction in overall execution time, with a decrease of approximately 40% compared to Case1. This improvement can be

attributed to the combiner's ability to locally aggregate and compress the intermediate key-value pairs generated by the mapper, reducing the data volume transferred over the network and optimizing the reducer's workload. The utilization of a combiner not only minimizes network congestion and data transfer, but also optimizes the reducer's workload by performing an initial aggregation step locally.

### 3) Case 3

For Case3, our objective was to evaluate the impact of employing multiple reducers on reducing the execution time of the data processing task. By distributing the workload across multiple reducers, we aimed to parallelize the computation and potentially achieve a faster overall processing time. In this experiment, we increased the number of reducers and monitored the execution time to assess the effectiveness of this approach. By dividing the data into multiple partitions and assigning each partition to a separate reducer, we expected to leverage parallel processing capabilities and potentially achieve a significant reduction in execution time.

```
hadoop jar
$HADOOP_HOME/share/hadoop/tools/lib/hadoop-streaming-3.3.5.jar \
-D mapred.reduce.tasks=50 \
-file mapper.py \
-mapper mapper.py \
-file reducer.py \
-reducer reducer.py \
-file combiner.py \
-combiner combiner.py \
-file stopword.txt \
-input /wordcount/input/data_16GB.txt \
-output /wordcount/output/16gb/case3-temp \
&&\
hadoop jar
$HADOOP_HOME/share/hadoop/tools/lib/hadoop-streaming-3.3.5.jar \
-file sort_map.py \
-mapper sort_map.py \
-file sort_reduce.py \
-reducer sort_reduce.py \
```

```

-file combiner.py \
-combiner combiner.py \
-input /wordcount/output/16gb/case3b-temp \
-output /wordcount/output/16gb/case3b

```

Job Overview			
<b>Job Name:</b>	streamjob9167437224214791307.jar		
<b>User Name:</b>	sandy		
<b>Queue:</b>	default		
<b>State:</b>	SUCCEEDED		
<b>Uberized:</b>	false		
<b>Submitted:</b>	Sun May 14 12:08:21 PDT 2023		
<b>Started:</b>	Sun May 14 12:08:30 PDT 2023		
<b>Finished:</b>	Sun May 14 12:58:08 PDT 2023		
<b>Elapsed:</b>	49mins, 37sec		
<b>Diagnostics:</b>			
<b>Average Map Time</b>	55sec		
<b>Average Shuffle Time</b>	45sec		
<b>Average Merge Time</b>	0sec		
<b>Average Reduce Time</b>	1sec		

  

ApplicationMaster			
Attempt Number	Start Time	Node	Logs
1	Sun May 14 12:08:26 PDT 2023	<a href="#">100.76.111.243:8042</a>	<a href="#">logs</a>

  

Task Type	Total	Complete
<b>Map</b>	115	115
<b>Reduce</b>	50	50

  

Attempt Type	Failed	Killed	Successful
<b>Maps</b>	0	0	115
<b>Reduces</b>	0	0	50

However, upon analyzing the results, we found that the expected improvement in execution time was not realized. Despite using multiple reducers, the overall processing time did not demonstrate a substantial reduction compared to previous cases.

This outcome suggests that the specific characteristics of the data, the nature of the computation, or other factors may have limited the benefits of employing multiple reducers in this particular scenario. Further analysis is required to understand the underlying reasons for the lack of performance improvement and to identify potential optimizations that could enhance the efficiency of the data processing task.

Based on the results of our experimentation, it can be concluded that Case 2, which involved the utilization of a combiner in the data processing pipeline, exhibited better performance compared to Case 3, which employed multiple reducers.



In Case 2, the introduction of the combiner allowed for partial aggregation of intermediate key-value pairs within the mapper phase itself. This reduced the amount of data transferred between the mappers and reducers, leading to a more efficient data processing workflow. As a result, the overall execution time was significantly reduced by approximately 40% compared to Case 1.

## IV. Conclusion

The primary objective of this project was to efficiently process a large dataset using the Hadoop Data Processing Framework - MapReduce, and extract the top 100 words in parallel across clusters. To achieve this, we implemented nlargest and heapq algorithms, along with the sort\_mapper and sort\_reducer files, which helped streamline the sorting process and identify the most frequently occurring words.

After running several test cases, we identified areas for improvement in our program, including adding a combiner function, making parameter changes in the .xml files, and setting up the Hadoop LZO environment to enhance performance further. We were able to achieve significant reductions in processing time by optimizing hardware and software performance. For instance, in case one from part three, we increased the amount of data that can be spilled to disk during the sort and shuffle phase, reducing processing time from over an hour to just 47 minutes. In case two, by adding a combiner function and strictly coding it, we were able to achieve a 40% reduction in processing time.

This project demonstrates the power of MapReduce in processing large datasets and the various techniques that can be employed to optimize performance. Throughout this project, we have gained insights into how MapReduce processes large volumes of data in a parallel and scalable manner, while also simplifying complex computations through effective management of Mappers and Reducers. By constantly evaluating and improving our methods, we can continue to refine our approach and tackle even more complex big data challenges in the future.

## V. Reference

[Apache Hadoop 3.3.5 – MapReduce Tutorial](#)

[Hadoop Streaming](#)

Pusukuri, Kishore. COEN 242 Class Slides. Santa Clara University

[Hadoop LZO](#)

[Hadoop Mapreduce Client](#)

## Top100 Words for 16GB.txt

said 16983038	much 1878812	right 1336449
would 5829109	say 1866770	court 1330434
one 5827854	day 1831556	team 1329924
new 5619251	week 1829329	united 1310698
also 4618231	home 1827482	need 1307171
us 4602724	take 1785403	report 1299343
people 4302078	per 1779050	country 1295488
last 4096906	work 1770896	help 1289802
year 4011803	going 1744159	according 1281617
two 3964144	think 1670557	business 1275145
first 3839807	company 1665814	market 1267854
mr 3746747	good 1635409	life 1256423
years 3637652	dont 1629480	long 1242001
time 3635760	next 1605312	set 1232190
could 3374880	including 1594887	months 1227463
like 3058690	see 1561688	man 1226580
government 2520240	states 1543254	best 1211261
says 2417351	another 1536907	come 1208192
may 2393105	group 1521302	cent 1203960
get 2346949	go 1508945	officials 1195387
many 2344542	public 1506218	family 1186525
back 2295703	de 1499434	left 1178734
million 2255486	house 1494732	high 1177794
three 2250200	around 1482330	show 1172154
president 2193692	city 1474866	health 1163668
even 2175400	former 1461398	york 1161016
made 2130243	part 1446451	
make 2104945	second 1433375	
told 2092903	national 1415826	
since 2041708	obama 1405174	
world 2039952	know 1386991	
percent 2022481	want 1379621	
police 2014250	game 1371927	
well 1965773	four 1371032	
still 1956831	news 1368259	
state 1948101	found 1366151	
way 1936600	end 1353926	

### Top100 Words(>6) for 16GB.txt

government	2520240	something	896500	despite	652919
million	2255486	military	893360	official	643514
president	2193692	several	877379	council	639609
percent	2022481	companies	869826	election	635236
company	1665814	washington	866772	announced	632756
including	1594887	members	861953	looking	630482
another	1536907	service	853089	possible	629503
national	1415826	federal	852384	thought	628703
country	1295488	campaign	820329	control	628021
according	1281617	statement	817540	current	624576
business	1275145	economy	816287	program	624140
officials	1195387	british	813789	getting	623774
american	1129214	university	806441	interest	612373
international	1116052	services	799299	outside	605728
financial	1097729	director	788859	england	599716
support	1080825	important	773502	minutes	599437
children	1072801	working	766645	spokesman	596504
however	1072387	players	755612	continue	594807
security	1058065	countries	740881	increase	588451
billion	1049368	decision	732901	private	586730
without	1039688	earlier	732872	authorities	584908
political	1018915	different	722896	meeting	578137
european	1018667	general	722268	history	576197
whether	1014215	research	716089	together	575985
yearold	981085	capital	705038	results	575269
tuesday	969122	saturday	702016	community	571257
already	960627	believe	701746	reports	565968
minister	942536	industry	700059	leaders	565284
information	927619	department	699764	problem	562464
wednesday	920376	executive	696821	hospital	559119
reported	917579	quarter	696623	process	556804
expected	916793	following	665860	markets	555948
economic	910616	although	662413		
thursday	910329	foreign	654706		

## Top100 Words for 2.5GB.txt

said 2616266	home 304861	life 218625
one 949555	week 299482	national 216146
would 917212	much 296828	family 214071
new 852822	cent 294331	best 207276
also 727934	work 294328	news 206894
last 700735	le 284535	need 205295
people 688828	take 284214	help 204715
us 670794	going 283188	set 202012
mr 659562	good 278834	come 201761
year 654244	say 278437	got 201136
de 643122	dont 277780	long 200086
two 637405	think 274485	according 199469
first 616295	next 269513	business 199119
years 608882	state 264274	left 197811
time 602016	around 259091	les 196404
could 548856	see 258524	et 195339
like 484097	including 250803	months 193875
says 419576	president 250617	really 193662
government 405470	go 249882	used 193435
back 398012	another 249601	market 192509
get 397140	found 245247	percent 191929
may 372870	company 244010	im 191567
police 371212	former 239811	big 189219
per 367298	public 236849	place 188953
la 360237	part 235478	
three 360087	court 234684	
told 359229	group 234072	
many 358774	team 232873	
made 344154	man 231083	
even 335035	know 230360	
make 333551	game 230229	
well 326187	city 227040	
world 319405	house 225066	
million 318478	want 224799	
since 318416	second 223861	
way 316957	end 222602	
still 312846	right 221258	
day 305766	four 220644	

## Top100 Words(>6) for 2.5GB.txt

government	405470	tuesday	126216	european	96292
million	318478	members	126102	general	96204
including	250803	federal	125761	campaign	94970
president	250617	director	125292	started	94334
another	249601	companies	123237	foreign	93788
company	244010	statement	123219	together	93721
national	216146	military	122988	continue	93269
according	199469	thursday	122559	morning	92841
business	199119	different	122406	interest	92423
percent	191929	research	121958	michael	92310
children	186405	wednesday	121808	announced	92247
country	185323	services	120953	history	91906
yearold	174923	earlier	120110	playing	90794
however	171875	decision	120034	council	90489
minister	171512	following	118936	evidence	90275
support	160054	important	116119	forward	89655
without	159819	reported	113628	community	89601
whether	157119	executive	112042	official	89275
international	155251	despite	112028	current	89202
something	150041	hospital	111788	spokesman	87907
already	147866	thought	111391	building	87419
players	147136	economic	108564	private	87286
officials	145338	looking	108465	released	86908
security	143445	industry	108024	countries	86402
expected	142353	economy	107246	parents	86343
american	140006	department	105927	election	85852
billion	139672	england	105272		
financial	135819	saturday	104744		
british	134480	getting	103814		
political	133593	believe	103005		
information	133546	although	102923		
australia	133494	minutes	102270		
australian	133022	outside	99852		
university	132830	washington	99784		
service	132019	possible	97903		
several	127788	capital	97281		
working	126900	control	96887		

## Top100 Words for 300MB.txt

european	318743	many	53998	international	40222
mr	210690	social	53696	country	39653
would	181912	way	53401	directive	39575
also	180118	believe	52845	said	39528
commission	172783	development	52466	state	39025
must	156856	commissioner	51201	within	38778
president	152134	say	50620	already	38700
union	130344	proposal	50250	much	38680
states	130270	market	49889	cooperation	38527
member	126301	fact	48978	good	38321
parliament	122059	debate	48595	common	38316
report	119173	human	48400	vote	38217
like	111176	group	46352	world	38038
council	107720	question	45847	part	37862
one	102551	think	45818	clear	37847
europe	95645	agreement	45493	members	37768
countries	93490	years	45489	mrs	37471
us	91379	even	45344	may	37147
eu	86935	issue	44464	still	36947
need	84622	citizens	44331	year	36773
policy	82662	point	44227	particular	36433
important	81684	future	43162	use	36111
new	81139	order	43091	know	36086
people	78712	situation	42969	see	35543
time	78348	public	42716	system	34983
rights	69717	financial	42361	energy	34781
therefore	68313	well	42255		
support	68006	right	42231		
however	65871	measures	42147		
take	64840	two	42045		
make	64650	cannot	42028		
work	62135	could	41419		
economic	60152	possible	41337		
committee	59456	community	41307		
political	56396	want	41292		
made	55424	today	40822		
first	54666	national	40257		

## Top100 Words(>6) for 300MB.txt

european	318743	protection	31263	concerned	21513
commission	172783	problem	30630	another	21287
president	152134	example	30564	towards	21233
parliament	122059	gentlemen	30454	welcome	21101
council	107720	amendments	30308	research	20969
countries	93490	problems	29849	something	20907
important	81684	government	29573	approach	20826
therefore	68313	rapporteur	29433	environment	20733
support	68006	particularly	29247	commissions	20637
however	65871	programme	29047	current	20384
economic	60152	necessary	28937	forward	20348
committee	59456	resolution	28790	resources	20189
political	56396	position	27628	conditions	20158
believe	52845	respect	27523	clearly	20086
development	52466	whether	26696	finally	20040
commissioner	51201	framework	26478	together	20040
proposal	50250	presidency	26427	including	20030
question	45847	services	25850	industry	20024
agreement	45493	regulation	25611	principle	20000
citizens	44331	proposals	25403	procedure	19964
situation	42969	amendment	25115	fundamental	19823
financial	42361	certain	24890	authorities	19651
measures	42147	adopted	24128	importance	19645
possible	41337	legislation	23702	proposed	19548
community	41307	different	23433	specific	19452
national	40257	strategy	22929	policies	19429
international	40222	institutions	22850		
country	39653	decision	22848		
directive	39575	present	22830		
already	38700	progress	22575		
cooperation	38527	subject	22475		
members	37768	continue	22407		
particular	36433	opinion	22225		
security	33259	employment	22066		
information	32899	transport	21805		
without	32726	working	21633		
process	32647	general	21518		