

UNIT 2 – Data Storage Technologies

Activity 2: MapReduce Optimization Experiment

Task 1: Program Development

1. Develop a MapReduce application for NOAA weather datasets analysis
2. Implement mapper to parse weather records into key-value pairs
3. Implement reducer to aggregate weather metrics
4. Inclusion of combiners when beneficial

Task 2: Optimization Experiments Trials on Hadoop

1. Execution of multiple trial runs, Vary the number of mappers and reducers.
2. Enable or disable combiners to observe differences.
3. Adjust block sizes or input split configurations.
4. Compare execution times and resource usage before and after tuning.

Task 3: Performance Monitoring and Analysis

1. Monitor job execution using Hadoop JobHistoryServer
2. Track time per phase (map, shuffle, reduce)
3. Record number of mappers and reducers allocated
4. Identify performance bottlenecks (CPU, I/O, network)
5. Analyze logs for optimization opportunities

When Mapper block size is 128 mb with and without Combiner

```
bytes written: 101
2025-12-16 10:03:18,996 INFO streaming.StreamJob: Output directory: /cloud/output/Pinky
hduser@master:~$ hdfs dfs -cat /cloud/output/Pinky
cat: '/cloud/output/Pinky': Is a directory
hduser@master:~$ hdfs dfs -ls /cloud/output/Pinky
Found 2 items
-rw-r--r-- 2 hduser supergroup          0 2025-12-16 10:03 /cloud/output/Pinky/_SUCCESS
-rw-r--r-- 2 hduser supergroup      101 2025-12-16 10:03 /cloud/output/Pinky/part-00000
hduser@master:~$ hdfs dfs -cat /cloud/output/Pinky/part-00000
clear    3432
cloudy   3797
drizzle  166
fog      154
haze     16
rain     664
snow     527
snow pellets 1
thunderstorms 27
hduser@master:~$ 
```

The screenshot shows the Hadoop JobHistory interface. On the left, there's a sidebar with 'Application' and 'Jobs' sections, and a 'Tools' section. The main area is titled 'JobHistory' and shows 'Retired jobs'. There are two rows of job information:

Submitted Time	Start Time	Finish Time	Job ID	Name	User	Queue	Date	Maps Total	Maps Completed	Reducers Total	Reducers Completed	Dropped Tasks
2025-12-16 10:49:39	2025-12-16 10:49:40	2025-12-16 10:50:19	hduser@10.10.10.124:12322_0002	streaming@10.10.10.124:12322_0002.jar	hduser	root/default	SUCCEEDED	87	87	1	1	00%
2025-12-16 10:48:13	2025-12-16 10:48:15	2025-12-16 10:47:29	hduser@10.10.10.124:12322_0003	streaming@10.10.10.124:12322_0003.jar	hduser	root/default	SUCCEEDED	87	87	1	1	00%, 91%, 100%

Without Combiner:-

```
File System Counters
  FILE: Number of bytes read=705351059
  FILE: Number of bytes written=1422610428
  FILE: Number of read operations=8
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=4919497435
  HDFS: Number of bytes written=116
  HDFS: Number of read operations=116
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=2
  HDFS: Number of bytes read erasure-coded=0

Job Counters
  Killed map tasks=1
  Launched map tasks=37
  Launched reduce tasks=1
  Data-local map tasks=37
  Total time spent by all maps in occupied slots (ms)=157412
  Total time spent by all reduces in occupied slots (ms)=54568
  Total time spent by all map tasks (ms)=157412
  Total time spent by all reduce tasks (ms)=54568
  Total vcore-milliseconds taken by all map tasks=157412
  Total vcore-milliseconds taken by all reduce tasks=54568
  Total megabyte-milliseconds taken by all map tasks=161189688
  Total megabyte-milliseconds taken by all reduce tasks=55869440

Map-Reduce Framework
  Map input records=68974923
  Map output records=68974923
  Map output bytes=567401801
  Map output materialized bytes=785351889
  Input split bytes=3626
  Combine input records=0
  Combine output records=0
  Reduce input groups=8
  Reduce shuffle bytes=705351869
  Reduce input records=68974923
  Reduce output records=8
  Spilled Records=137949846
  Shuffled Maps =37
  Failed Shuffles=0
  Merged Map outputs=37
  GC time elapsed (ms)=987
  CPU time spent (ms)=133528
  Physical memory (bytes) snapshot=13982107520
  Virtual memory (bytes) snapshot=107986483328
  Total committed heap usage (bytes)=8511291392
  Peak Map Physical memory (bytes)=368348982
  Peak Map Virtual memory (bytes)=2848714732
  Peak Reduce Physical memory (bytes)=874619688
  Peak Reduce Virtual memory (bytes)=2873985264

Shuffle Errors
  RAG ID=0
```

With Combiner:-

```
File System Counters
  FILE: Number of bytes read=4486
  FILE: Number of bytes written=11941418
  FILE: Number of read operations=8
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=4919497435
  HDFS: Number of bytes written=116
  HDFS: Number of read operations=116
  HDFS: Number of large read operations=8
  HDFS: Number of write operations=2
  HDFS: Number of bytes read erasure-coded=0

Job Counters
  Killed map tasks=1
  Launched map tasks=37
  Launched reduce tasks=1
  Data-local map tasks=37
  Total time spent by all maps in occupied slots (ms)=159813
  Total time spent by all reduces in occupied slots (ms)=21267
  Total time spent by all map tasks (ms)=159813
  Total time spent by all reduce tasks (ms)=21267
  Total vcore-milliseconds taken by all map tasks=159813
  Total vcore-milliseconds taken by all reduce tasks=21267
  Total megabyte-milliseconds taken by all map tasks=162829312
  Total megabyte-milliseconds taken by all reduce tasks=21777408

Map-Reduce Framework
  Map input records=68974924
  Map output records=68974923
  Map output bytes=567401081
  Map output materialized bytes=4622
  Input split bytes=3020
  Combine input records=68974923
  Combine output records=296
  Reduce input groups=8
  Reduce shuffle bytes=4622
  Reduce input records=296
  Reduce output records=8
  Spilled Records=592
  Shuffled Maps =37
  Failed Shuffles=0
  Merged Map outputs=37
  GC time elapsed (ms)=106870
  CPU time spent (ms)=106870
  Physical memory (bytes) snapshot=13448446464
  Virtual memory (bytes) snapshot=168622833152
  Total committed heap usage (bytes)=8018072864
  Peak Map Physical memory (bytes)=372338496
  Peak Map Virtual memory (bytes)=2648497664
  Peak Reduce Physical memory (bytes)=219353088
  Peak Reduce Virtual memory (bytes)=7846709152

Shuffle Errors
  BAD_ID=0
```

When Mapper block size is 64 mb with and without Combiner

```
hduser@master:~$ hadoop fs -cat /cloud/output/Pinky4/part-00000
clear    26949008
cloudy   29815349
drizzle  400466
fog      3345100
haze     125642
rain     4397305
snow     3926349
thunderstorms  15704
```

```

File System Counters
  FILE: Number of bytes read=4466
  FILE: Number of bytes written=119448392
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=4919497435
  HDFS: Number of bytes written=116
  HDFS: Number of read operations=116
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=2
  HDFS: Number of bytes read erasure-coded=0

Job Counters
  Killed map tasks=1
  Launched map tasks=37
  Launched reduce tasks=1
  Data-local map tasks=37
  Total time spent by all maps in occupied slots (ms)=177766
  Total time spent by all reduces in occupied slots (ms)=9884
  Total time spent by all map tasks (ms)=177766
  Total time spent by all reduce tasks (ms)=9884
  Total vcore-milliseconds taken by all map tasks=177766
  Total vcore-milliseconds taken by all reduce tasks=9884
  Total megabyte-milliseconds taken by all map tasks=382832304
  Total megabyte-milliseconds taken by all reduce tasks=18121216

Map-Reduce Framework
  Map input records=60974924
  Map output records=60974923
  Map output bytes=567461801
  Map output materialized bytes=4622
  Input split bytes=3626
  Combine input records=60974923
  Combine output records=296
  Reduce input groups=8
  Reduce shuffle bytes=4622
  Reduce input records=296
  Reduce output records=8
  Spilled Records=592
  Shuffled Maps =37
  Failed Shuffles=0
  Merged Map outputs=37
  GC time elapsed (ms)=1097
  CPU time spent (ms)=111728
  Physical memory (bytes) snapshot=13345013768
  Virtual memory (bytes) snapshot=107956510720
  Total committed heap usage (bytes)=8664829184
  Peak Map Physical memory (bytes)=370601984
  Peak Map Virtual memory (bytes)=2844717344
  Peak Reduce Physical memory (bytes)=221043456
  Peak Reduce Virtual memory (bytes)=2850045952

```

When Reducer is 1 with and without Combiner

```

2025-12-16 10:03:18,996 INFO streaming.StreamJob: Output directory: /cloud/output/Pinky
hduser@master:~$ hdfs dfs -cat /cloud/output/Pinky
cat: '/cloud/output/Pinky': Is a directory
hduser@master:~$ hdfs dfs -ls /cloud/output/Pinky
Found 2 items
-rw-r--r--  2 hduser supergroup      0 2025-12-16 10:03 /cloud/output/Pinky/_SUCCESS
-rw-r--r--  2 hduser supergroup    101 2025-12-16 10:03 /cloud/output/Pinky/part-00000
hduser@master:~$ hdfs dfs -cat /cloud/output/Pinky/part-00000
clear      3432
cloudy     3797
drizzle    166
fog        154
haze       16
rain       664
snow       527
snow pellets   1
thunderstorms 27
hduser@master:~$ []

```


JobHistory

Logged in as: admin

- Application	File	Jobs										
+ Tools												
Retired jobs												
Show 20 - 30779												
Submitted Time	Start Time	Finish Time	Job ID	Name	User	Queue	Status	Maps Total	Maps Completed	Reduces Total	Reduces Completed	Elapsed Time
2025-12-19 18:49:39	2025-12-19 18:49:40	2025-12-19 18:50:19	job_12008833101072_0002	UnescapeJob12008833101072_0002.jar	hduser	root/default	SUCCEEDED	37	37	1	1	00:05: 0000ms 30ms
2025-12-19 18:48:13	2025-12-19 18:48:15	2025-12-19 18:47:29	job_12008833101073_0003	UnescapeJob12008833101073_0003.jar	hduser	root/default	SUCCEEDED	37	37	1	1	00:05: 0100ms 100ms

Without Combiner:-

```

File System Counters
  FILE: Number of bytes read=705351059
  FILE: Number of bytes written=1422610428
  FILE: Number of read operations=8
  FILE: Number of large read operations=0
  FILE: Number of write operations=8
  HDFS: Number of bytes read=4919497435
  HDFS: Number of bytes written=116
  HDFS: Number of read operations=116
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=2
  HDFS: Number of bytes read erasure-coded=0

Job Counters
  Killed map tasks=1
  Launched map tasks=37
  Launched reduce tasks=1
  Data-local map tasks=37
  Total time spent by all maps in occupied slots (ms)=157412
  Total time spent by all reduces in occupied slots (ms)=54568
  Total time spent by all map tasks (ms)=157412
  Total time spent by all reduce tasks (ms)=54568
  Total vcore-milliseconds taken by all map tasks=54568
  Total vcore-milliseconds taken by all reduce tasks=54568
  Total megabyte-milliseconds taken by all map tasks=161189088
  Total megabyte-milliseconds taken by all reduce tasks=55869440

Map-Reduce Framework
  Map input records=68974924
  Map output records=68974923
  Map output bytes=567401801
  Map output materialized bytes=705351089
  Input split bytes=3626
  Combine input records=0
  Combine output records=0
  Reduce input groups=0
  Reduce shuffle bytes=705351069
  Reduce input records=68974923
  Reduce output records=0
  Spilled Records=137949846
  Shuffled Maps =37
  Failed Shuffles=0
  Merged Map outputs=37
  GC time elapsed (ms)=987
  CPU time spent (ms)=153528
  Physical memory (bytes) snapshot=13982117520
  Virtual memory (bytes) snapshot=107986403328
  Total committed heap usage (bytes)=8511291392
  Peak Map Physical memory (bytes)=366340592
  Peak Map Virtual memory (bytes)=2848714732
  Peak Reduce Physical memory (bytes)=874610688
  Peak Reduce Virtual memory (bytes)=2873995264
  shuffle Errors
    BAD_ID=0
  
```

With Combiner:-

```
File System Counters
  FILE: Number of bytes read=4486
  FILE: Number of bytes written=11941418
  FILE: Number of read operations=8
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=4919497435
  HDFS: Number of bytes written=116
  HDFS: Number of read operations=116
  HDFS: Number of large read operations=8
  HDFS: Number of write operations=2
  HDFS: Number of bytes read erasure-coded=0

Job Counters
  Killed map tasks=1
  Launched map tasks=37
  Launched reduce tasks=1
  Data-local map tasks=37
  Total time spent by all maps in occupied slots (ms)=159813
  Total time spent by all reduces in occupied slots (ms)=21267
  Total time spent by all map tasks (ms)=159813
  Total time spent by all reduce tasks (ms)=21267
  Total vcore-milliseconds taken by all map tasks=159813
  Total vcore-milliseconds taken by all reduce tasks=21267
  Total megabyte-milliseconds taken by all map tasks=162829312
  Total megabyte-milliseconds taken by all reduce tasks=21777408

Map-Reduce Framework
  Map input records=68974924
  Map output records=68974923
  Map output bytes=567401081
  Map output materialized bytes=4622
  Input split bytes=3020
  Combine input records=68974923
  Combine output records=296
  Reduce input groups=8
  Reduce shuffle bytes=4622
  Reduce input records=296
  Reduce output records=8
  Spilled Records=592
  Shuffled Maps =37
  Failed Shuffles=0
  Merged Map outputs=37
  GC time elapsed (ms)=106870
  CPU time spent (ms)=106870
  Physical memory (bytes) snapshot=13448446464
  Virtual memory (bytes) snapshot=168622833152
  Total committed heap usage (bytes)=8818072864
  Peak Map Physical memory (bytes)=372338496
  Peak Map Virtual memory (bytes)=2648497664
  Peak Reduce Physical memory (bytes)=219353088
  Peak Reduce Virtual memory (bytes)=7846769152

Shuffle Errors
  BND_ID=0
```

When Reducer is 2 Without Combiner and With Combiner

```
hduser@master:~$ hadoop fs -cat /cloud/output/Pinky3/part-00000
clear    26949008
cloudy   29815349
drizzle  400466
fog      3345100
haze     125642
rain     4397305
snow     3926349
thunderstorms 15704
```

Without Combiner:-

 **hadoop**

JobHistory

Report for job 00000000000000000000000000000000

Retired Jobs

Submit Time	Start Time	Finish Time	Job ID	Name	User	Queue	Status	Maps Total	Maps Completed	Reduces Total	Reduces Completed	Report Time
2025-12-16 17:14:08	2025-12-16 17:14:11	2025-12-16 17:15:25	job_00000000000000000000000000000000	mapred@ip-10-0-1-12:547254419678120948	hadoop	root/default	SUCCEEDED	37	37	2	2	00:00: 07m 59s 14ms
BT	BT	BT										

File System Counters

```

FILE: Number of bytes read=705331659
FILE: Number of bytes written=1422923957
FILE: Number of read operations=8
FILE: Number of large read operations=8
FILE: Number of write operations=0
HDFS: Number of bytes read=4919497435
HDFS: Number of bytes written=116
HDFS: Number of read operations=121
HDFS: Number of large read operations=8
HDFS: Number of write operations=4
HDFS: Number of bytes read erasure-coded=0

```

Job Counters

```

Killed map tasks=1
Launched map tasks=37
Launched reduce tasks=3
Data-local map tasks=37
Total time spent by all maps in occupied slots (ms)=170737
Total time spent by all reduces in occupied slots (ms)=128793
Total time spent by all map tasks (ms)=170737
Total time spent by all reduce tasks (ms)=128793
Total vcore-milliseconds taken by all map tasks=170737
Total vcore-milliseconds taken by all reduce tasks=128793
Total megabyte-milliseconds taken by all map tasks=174834608
Total megabyte-milliseconds taken by all reduce tasks=123692032

```

Map-Reduce Framework

```

Map input records=60974924
Map output records=60974923
Map output bytes=567401881
Map output materialized bytes=705352091
Input split bytes=3626
Combine input records=8
Combine output records=8
Reduce input groups=8
Reduce shuffle bytes=705352091
Reduce input records=60974923
Reduce output records=8
Spilled Records=137549946
Shuffled Maps =74
Failed Shuffles=8
Merged Map outputs=74
GC time elapsed (ms)=1200
CPU time spent (ms)=167980
Physical memory (bytes) snapshot=14081712128
Virtual memory (bytes) snapshot=110012504064
Total committed heap usage (bytes)=8624537688
Peak Map Physical memory (bytes)=370926664
Peak Map Virtual memory (bytes)=2046994432
Peak Reduce Physical memory (bytes)=619864544
Peak Reduce Virtual memory (bytes)=2876809216

```

Shuffle Errors

With Combiner:-



JobHistory

```
925-12-16 17:09:54,295 INFO mapreduce.Job: Job job_11055385073_0022 completed successfully  
925-12-16 17:09:54,295 INFO mapreduce.Job: Counters: 55
```

File System Counters

```
FILE: Number of bytes read=4412  
FILE: Number of bytes written=12255359  
FILE: Number of read operations=0  
FILE: Number of large read operations=0  
FILE: Number of write operations=0  
HDFS: Number of bytes read=4919497435  
HDFS: Number of bytes written=116  
HDFS: Number of read operations=121  
HDFS: Number of large read operations=0  
HDFS: Number of write operations=4  
HDFS: Number of bytes read erasure-coded=
```

Job Counters

```
Killed map tasks=1
Launched map tasks=37
Launched reduce tasks=2
Data-local map tasks=37
Total time spent by all maps in occupied slots (ms)=181442
Total time spent by all reduces-in occupied slots (ms)=13732
Total time spent by all map tasks (ms)=181442
Total time spent by all reduce tasks (ms)=13732
Total vcore-milliseconds taken by all map tasks=101442
Total vcore-milliseconds taken by all reduce tasks=13732
Total megabyte-milliseconds taken by all map tasks=185796608
Total megabyte-milliseconds taken by all reduce tasks=14861568
```

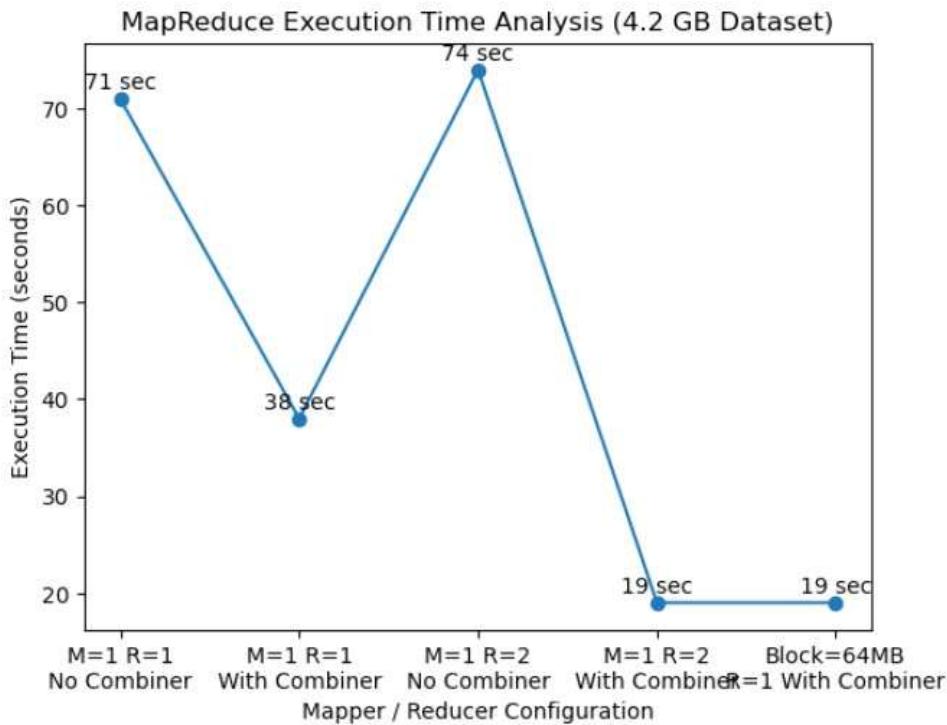
Map-Reduce Framework

```
Map input records=68974924
Map output records=68974923
Map output bytes=367481801
Map output materialized bytes=4844
Input split bytes=3626
Combine input records=68974923
Combine output records=290
Reduce input groups=11
Reduce shuffle bytes=8044
Reduce input records=296
Reduce output records=0
Spilled Records=592
Shuffled Maps=74
Failed Shuffles=6
Merged Map outputs=74
GC time elapsed (ms)=1896
CPU time spent (ms)=123870
Physical memory (bytes) snapshot=1
Virtual memory (bytes) snapshot=11
Total committed heap usage (bytes)
Peak Map Physical memory (bytes)=3
Peak Map Virtual memory (bytes)=28
Peak Reduce Physical memory (bytes)
Peak Reduce Virtual memory (bytes)
```

Shuffle Errors

Task 4: Visualization and Reporting

1. Graphs comparing execution times



2. Visual summaries of mapper/reducer variations

Mapper	Reducer	Combiner	Execution Time
1	1	No	1min 11 sec
1	1	Yes	38sec
1	2	No	1min 14 sec
1	2	Yes	19 sec
Block size = 128mb	1	No	1min 11 sec
Block size = 128mb	1	Yes	38sec
Block size = 64 mb	1	Yes	19 sec

Pipeline:-

Input --> Mapper --> Combiner --> Shuffle & Sort --> Reducer --> Output

3.A structured 300-word document stored on GitHub

4.Inclusion of experiment notes, visuals, and program code

Link= <https://github.com/rajsanodiya122/Fobd-Unit-2-Act2>

Short Summary :-

When Mapper = 1 :- The input dataset is processed as a single map task. Data processing is **sequential**, not parallel. CPU and cluster resources are **underutilized**. Execution time is **higher** for large datasets (like my 4.2 GB data). Network and reducer stages wait for one mapper to finish.

When Mapper = 2 :- The data is split into two parts, each part is processed in parallel. Improved **parallel processing**. Better **CPU and I/O utilization**. Reduced mapper execution time. Faster availability of intermediate data for reducers

When Reducer = 1:- All intermediate key-value pairs produced by mappers are sent to a **single reduce task**. Data aggregation happens **sequentially**. Shuffle phase sends **all intermediate data to one node**. Can become a **bottleneck** for large datasets. Simpler output handling (single output file)

In my experiment:

For a 4.2 GB dataset **without a combiner**, reducer = 1 completed in **71 seconds**, showing moderate performance but limited scalability.

When Reducer = 2:- Intermediate data is **partitioned** and processed in **parallel** by two reduce tasks. Improved **parallelism** in the reduce phase. Reduced load on a single reducer. Multiple output files are generated. Extra overhead due to partitioning and task coordination

In my experiment:

Without a combiner, reducer = 2 took **74 seconds**, slightly **more than reducer = 1**, because the overhead of managing an extra reducer outweighed the benefit for this dataset size.

Without Combiner:- All intermediate key–value pairs generated by the mappers are **directly sent to the reducer** through the shuffle phase. Large volume of intermediate data. High **network I/O during shuffle**. Reducers take longer to aggregate data. Execution time increases significantly. Cluster resources are less efficiently utilized

In my experiment:

Without a combiner, execution time was **71–74 seconds**, even with different reducer configurations, due to heavy data transfer between mapper and reducer.

With Combiner:- A **combiner** performs **local aggregation at the mapper level** before data is sent to reducers. Intermediate data size is **greatly reduced**. Shuffle phase becomes much faster. Network traffic is minimized. Reducers process fewer records. Overall execution time is drastically reduced

In my experiment:

With a combiner enabled, execution time dropped to **38 seconds** and further to **19 seconds**, showing a **major performance improvement**.