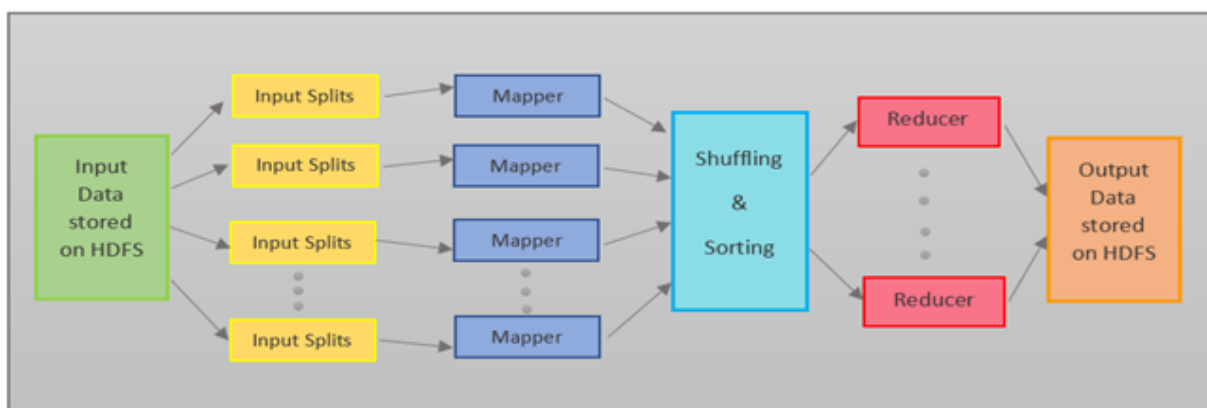


Map Reduce

In this assignment we have used map reduce to analyze the university students grades dataset using mrjob library. MapReduce enables parallel computation, improving efficiency, scalability, and fault tolerance when handling large datasets.

Each task follows the **MapReduce framework**, consisting of:

1. **Input Splitting:** The dataset is divided into smaller parts.
2. **Mapping:** Key-value pairs are generated from each data entry.
3. **Shuffling & Sorting:** Intermediate key-value pairs are grouped.
4. **Reducing:** Final computation is performed to summarize results.



This algorithm eliminates the bottleneck effect from the traditional low level operating system(OS) algorithms like: shortest job first(SJF), first come first serve(FCFS), Round

Robin(RR). As it computes all jobs concurrently with equal slices so it achieved the elimination of the bottleneck effect.

MRJOB: MRJob is a Python package that simplifies the process of writing MapReduce jobs. It provides a high-level API that abstracts away the low-level details of Hadoop MapReduce, making it easier for developers to write and run MapReduce jobs.

Task 1: Compute the Average Grade Per Course

Problem:

We need to calculate the **average grade** for each course and determine which course has the highest average.

Mapper Phase:

- Reads each line from the dataset.
- Extracts the Course Name and Grade.
- Returns (Course Name, Grade) pair.

Shuffling & Sorting:

- Groups all grades by Course Name.

Reducer Phase:

- Computes the average grade per course.
- Returns (Course Name, Average Grade).

Task 2: Compute the Average Grade Per University

Problem:

We need to calculate the **average grade** for each university and identify which university has the highest average.

Mapper Phase:

- Reads each line from the dataset.
- Extracts the **University Name** and **Grade**.
- returns (University Name, Grade) pair.

Shuffling & Sorting:

- Groups all grades by University Name.

Reducer Phase:

- Computes the average grade per university.
- Emits (University Name, Average Grade).

Task 3: Identify the Top 3 Highest Grades Per Year

Problem:**Mapper Phase:**

- Reads each line from the dataset.
- Extracts the Year and Grade.
- returns (Year, Grade) pair.

Shuffling & Sorting Phase:

- Groups all grades by Year.
- Sorts the grades in descending order.

Reducer Phase:

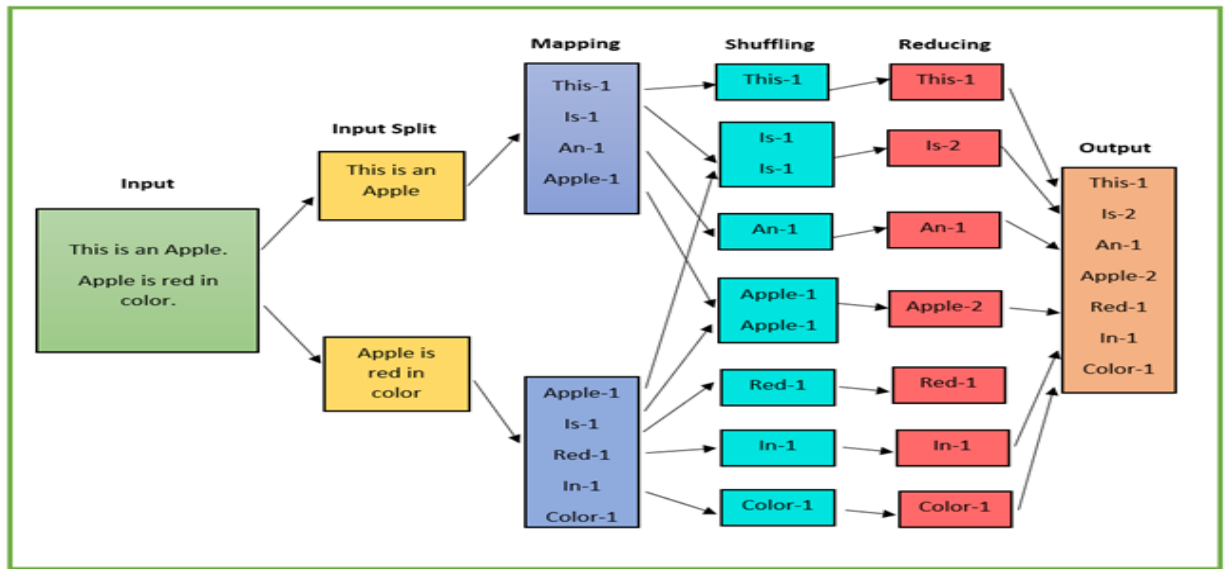
- Selects the top 3 grades for each year.
- Returns (Year, [Top 3 Grades]).

Code Execution:

```
python Task1.py coursegrades.txt > output_task1.txt
```

```
python Task2.py coursegrades.txt > output_task2.txt
```

```
python Task3.py coursegrades.txt > output_task3.txt
```



Project ▾ Distributed ▾ DB-Leader ▾

Task 1.py x coursegrades.txt Task 2.py Task 3.py

```
1 from mrjob.job import MRJob as mj
2
3 usage
4
5 class AvgGradePerCourse(mj):
6
7     def mapper(self, _, line):
8         _, course, grade, _ = line.split(',')
9         yield course.strip(), int(grade.strip())
10
11     def reducer(self, course, grades):
12         grade_list=list(grades)
13         yield course, sum(grade_list)/len(grade_list)
14
15 if __name__ == '__main__':
16     AvgGradePerCourse.run()
```

if __name__ == '__main__':

14:28 LF UTF-8 4 spaces Python 3.12

Project ▾ Distributed ▾ DB-Leader ▾

Task 1.py x coursegrades.txt Task 2.py Task 3.py

```
1 from mrjob.job import MRJob as mj
2
3 usage
4
5 class AvgGradePerCourse(mj):
6
7     def mapper(self, _, line):
8         _, course, grade, _ = line.split(',')
9         yield course.strip(), int(grade.strip())
10
11     def reducer(self, course, grades):
12         grade_list=list(grades)
13         yield course, sum(grade_list)/len(grade_list)
14
15 if __name__ == '__main__':
16     AvgGradePerCourse.run()
```

if __name__ == '__main__':

Terminal Local x + ▾

```
ali@Alis-MacBook-Air Distributed % source venv/bin/activate
(venv) ali@Alis-MacBook-Air Distributed %
```

14:28 LF UTF-8 4 spaces Python 3.12

Project

- Distributed ~/Desktop/Distributed
 - venv
 - coursegrades.txt
 - Task 1.py
 - Task 2.py
 - Task 3.py
 - External Libraries
 - Scratches and Consoles

Task 1.py

```
1 from mrjob.job import MRJob as mj
2
3 class AvgGradePerCourse(mj):
4
5     def mapper(self, _, line):
6         _, course, grade, _ = line.split(',')
7         yield course.strip(), int(grade.strip())
8
9     def reducer(self, course, grades):
10         grade_list = list(grades)
11         yield course, sum(grade_list)/len(grade_list)
12
13 if __name__ == '__main__':
14     AvgGradePerCourse.run()
```

Terminal

```
ali@Alis-MacBook-Air Distributed % source venv/bin/activate
(venv) ali@Alis-MacBook-Air Distributed % python "Task 1.py" coursegrades.txt > output_task1.txt
```

Project

- Distributed ~/Desktop/Distributed
 - venv
 - coursegrades.txt
 - output_task1.txt
 - Task 1.py
 - Task 2.py
 - Task 3.py
 - External Libraries
 - Scratches and Consoles

Task 1.py

```
1 Software Engineering 78.08695652173913
2 Cyber Security 78.58333333333333
3 Data Structures 81.45341614906832
4 Machine Learning 79.43125
5 Computer Vision 79.73958333333333
6 Artificial Intelligence 79.31012658227849
7
```

Terminal

```
ali@Alis-MacBook-Air Distributed % source venv/bin/activate
(venv) ali@Alis-MacBook-Air Distributed % python "Task 1.py" coursegrades.txt > output_task1.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 1.ali.20250312.082606.621532
Running step 1 of 1...
job output is in /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 1.ali.20250312.082606.621532/output
Streaming final output from /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 1.ali.20250312.082606.621532/output...
Removing temp directory /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 1.ali.20250312.082606.621532...
(venv) ali@Alis-MacBook-Air Distributed %
```

Project

- Distributed ~/Desktop/Distributed
 - venv
 - coursegrades.txt
 - output_task1.txt
 - Task 1.py
 - Task 2.py
 - Task 3.py
- External Libraries
- Scratches and Consoles

Task 1.py

```
1 from mrjob.job import MRJob as mj
2
3 usage
4
5 class AvgGradePerUniversity(mj):
6
7     def mapper(self, _, line):
8         _, _, grade, university = line.split(',')
9         yield university.strip(), int(grade.strip())
10
11     def reducer(self, university, grades):
12         grade_list = list(grades)
13         yield university, sum(grade_list) / len(grade_list)
14
15 if __name__ == '__main__':
16     AvgGradePerUniversity.run()
```

if __name__ == '__main__':

Project

- Distributed ~/Desktop/Distributed
 - venv
 - coursegrades.txt
 - output_task1.txt
 - Task 1.py
 - Task 2.py
 - Task 3.py
- External Libraries
- Scratches and Consoles

Task 1.py

```
1 from mrjob.job import MRJob as mj
2
3 usage
4
5 class AvgGradePerUniversity(mj):
6
7     def mapper(self, _, line):
8         _, _, grade, university = line.split(',')
9         yield university.strip(), int(grade.strip())
10
11     def reducer(self, university, grades):
12         grade_list = list(grades)
13         yield university, sum(grade_list) / len(grade_list)
14
15 if __name__ == '__main__':
16     AvgGradePerUniversity.run()
```

if __name__ == '__main__':

Terminal Local + -

```
ali@Alis-MacBook-Air Distributed % source venv/bin/activate
(venv) ali@Alis-MacBook-Air Distributed % python "Task 2.py" coursegrades.txt > output_task2.txt
```


Project

Distributed

~/Desktop/Distributed

venv

coursegrades.txt

output_task1.txt

output_task2.txt

Task 1.py

Task 2.py

Task 3.py

External Libraries

Scratches and Consoles

Task 1.py

output_task1.txt

coursegrades.txt

Task 2.py

output_task2.txt

Task 3.py

```
1 UC Berkeley 78.11242603550296
2 MIT 81.08904109589041
3 "Oxford University" 79.32558139534883
4 "Stanford University" 79.87647058823529
5 "Harvard University" 78.975
6 "Cambridge University" 79.44808743169399
7
```

Terminal

Local

```
ali@ALis-MacBook-Air Distributed % source venv/bin/activate
(venv) ali@ALis-MacBook-Air Distributed % python "Task 2.py" coursegrades.txt > output_task2.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 2.ali.20250312.082703.224554
Running step 1 of 1...
job output is in /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 2.ali.20250312.082703.224554/output
Streaming final output from /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 2.ali.20250312.082703.224554/output...
Removing temp directory /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 2.ali.20250312.082703.224554...
(venv) ali@ALis-MacBook-Air Distributed %
```

Distributed > output_task2.txt 1:1 LF UTF-8 4 spaces Python 3.12

Project

Distributed

~/Desktop/Distributed

venv

coursegrades.txt

output_task1.txt

output_task2.txt

Task 1.py

Task 2.py

Task 3.py

External Libraries

Scratches and Consoles

Task 1.py

output_task1.txt

coursegrades.txt

Task 2.py

output_task2.txt

Task 3.py

```
1 from mrjob.job import MRJob as mj
2
3 usage
4 class Top3GradesPerYear(mj):
5     def mapper(self, _, line):
6         year, _, grade, _ = line.split(',')
7         yield year.strip(), int(grade.strip())
8
9     def reducer(self, year, grades):
10         top_grades=sorted(grades, reverse=True)[:3]
11         yield year, top_grades
12
13 if __name__ == '__main__':
14     Top3GradesPerYear.run()
```

```
if __name__ == '__main__':
```

Distributed > Task 3.py 14:28 LF UTF-8 4 spaces Python 3.12

The screenshot shows a code editor with a project named "Distributed" located at ~/Desktop/Distributed. The project contains a "venv" directory with files "coursegrades.txt", "output_task1.txt", "output_task2.txt", "Task 1.py", "Task 2.py", and "Task 3.py". The "Task 3.py" file is open, showing a Python script that uses MRJob to process "coursegrades.txt" and output to "output_task3.txt". The script defines a "Top3GradesPerYear" class with a "mapper" and a "reducer" method. The "mapper" method splits each line into year and grade, and the "reducer" method sorts the grades for each year and outputs the top 3. The script is executed in a terminal window, showing the command "python 'Task 3.py' coursegrades.txt > output_task3.txt" and the output "output_task3.txt".

```
1 from mrjob.job import MRJob as mj
2
3 class Top3GradesPerYear(mj):
4
5     def mapper(self, _, line):
6         year, _, grade, _ = line.split(',')
7         yield year.strip(), int(grade.strip())
8
9     def reducer(self, year, grades):
10        top_grades=sorted(grades, reverse=True)[:3]
11        yield year, top_grades
12
13 if __name__ == '__main__':
14     Top3GradesPerYear.run()
```

Terminal Local x + v

```
ali@Alis-MacBook-Air Distributed % source venv/bin/activate
(venv) ali@Alis-MacBook-Air Distributed % python "Task 3.py" coursegrades.txt > output_task3.txt
```

The screenshot shows the same code editor with the "Task 3.py" file open. The terminal window shows the output of the script, which is a list of years and their top 3 grades. The output is "2024" [100, 100, 100], "2023" [100, 100, 100], and "2022" [100, 100, 100]. The terminal also shows the command "python 'Task 3.py' coursegrades.txt > output_task3.txt" and the output "output_task3.txt".

```
1 "2024" [100, 100, 100]
2 "2023" [100, 100, 100]
3 "2022" [100, 100, 100]
4
```

Terminal Local x

```
ali@Alis-MacBook-Air Distributed % source venv/bin/activate
(venv) ali@Alis-MacBook-Air Distributed % python "Task 3.py" coursegrades.txt > output_task3.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 3.ali.20250312.082805.466058
Running step 1 of 1...
job output is in /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 3.ali.20250312.082805.466058/output
Streaming final output from /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 3.ali.20250312.082805.466058/output...
Removing temp directory /var/folders/mz/chdvl3td5zvfq7_m8wLh6fsr0000gn/T/Task 3.ali.20250312.082805.466058...
(venv) ali@Alis-MacBook-Air Distributed %
```