



# MAP-REDUCE

**Submitted By**

Maimoona Khilji

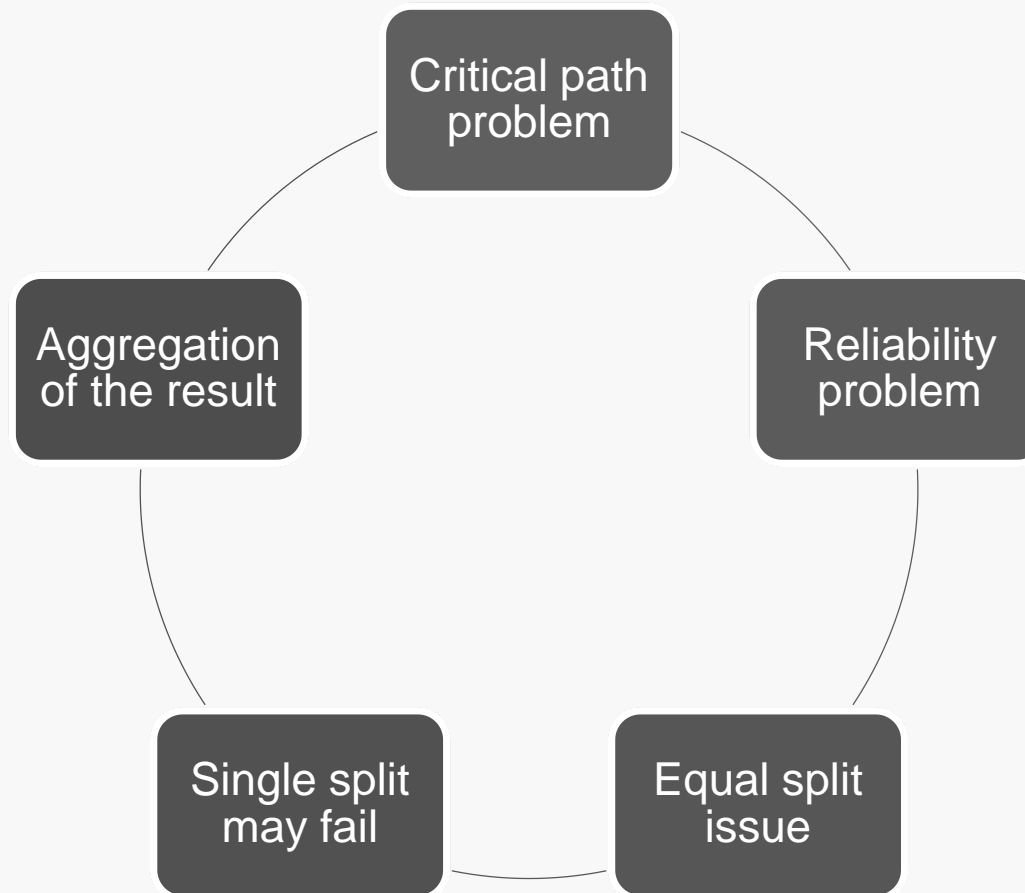
Registration no. 195300273

BS-Data Science

Semester - V



# Challenges in Traditional approach



# History

- MapReduce was **first popularized as a programming model in 2004** by Jeffery Dean and Sanjay Ghemawat of Google (Dean & Ghemawat, 2004).
- In their paper, “MAPREDUCE: SIMPLIFIED DATA PROCESSING ON LARGE CLUSTERS,” they discussed Google's approach to collecting and analyzing website data for search optimizations

# What is MapReduce?

- Processing component of apache Hadoop
- Processes data parallel in distributed environment.
- MapReduce is a programming model or pattern within the Hadoop framework
- It is used to access big data stored in the Hadoop File System (HDFS).

# Comparison between Traditional and Hadoop mapReduce



Human brings food toward the mouth

**VS.**



Tiger brings its mouth toward the food

# Comparison between Traditional and Hadoop mapReduce

1. Parallel Processing
  - Divide and Conquer process
2. Data Locality
  - Move the processing unit to the data in the MapReduce Framework

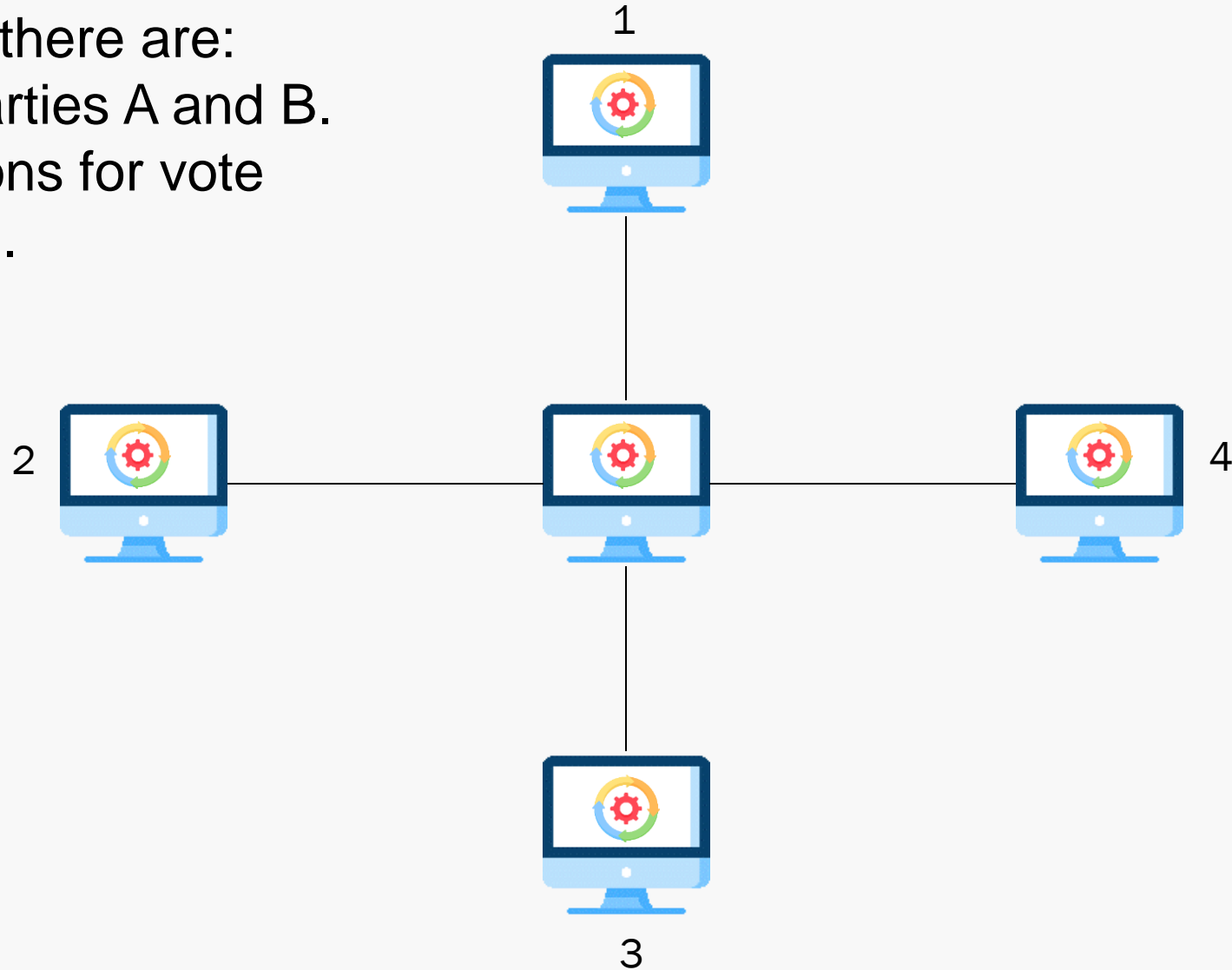
# Mechanism

- MapReduce facilitates concurrent processing
  - *by splitting petabytes of data into smaller chunks*
  - *processing them in parallel on Hadoop commodity servers.*
  - *In the end, it aggregates all the data from multiple servers to return an output back to the application.*

# Votes Counting

Suppose there are:

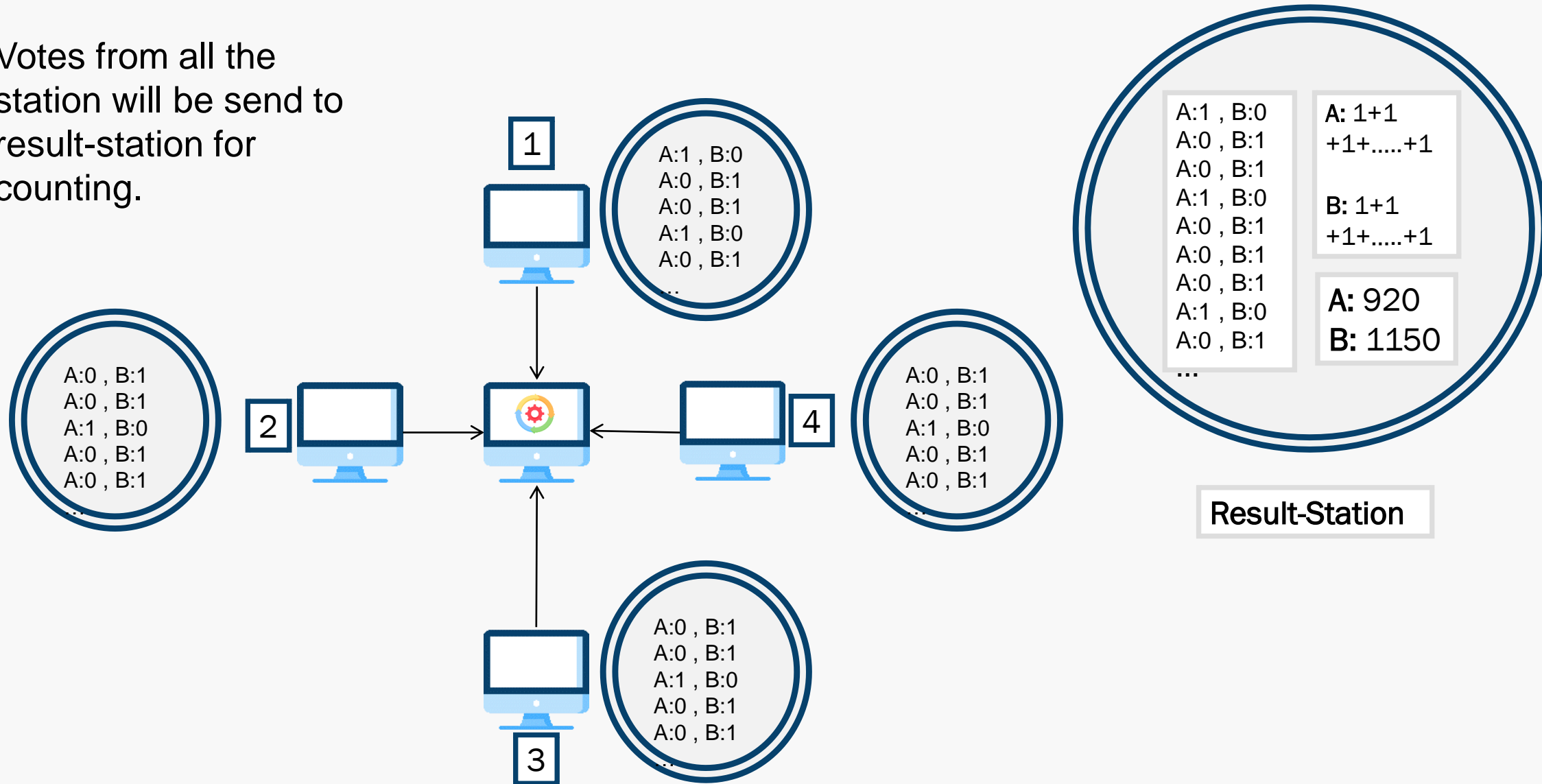
- Two parties A and B.
- 4 stations for vote casting.





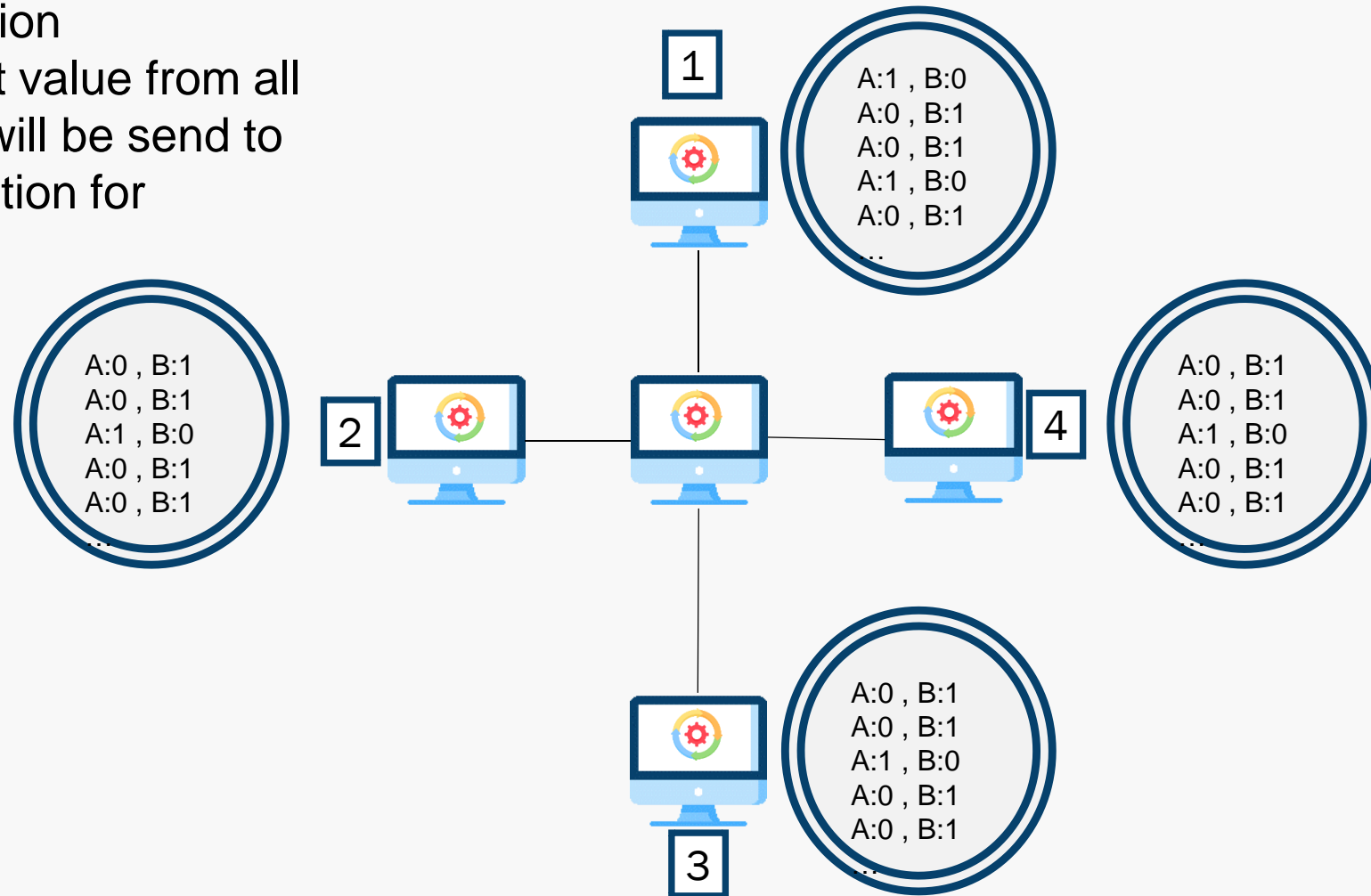
# Traditional way of vote counting

- Votes from all the station will be send to result-station for counting.



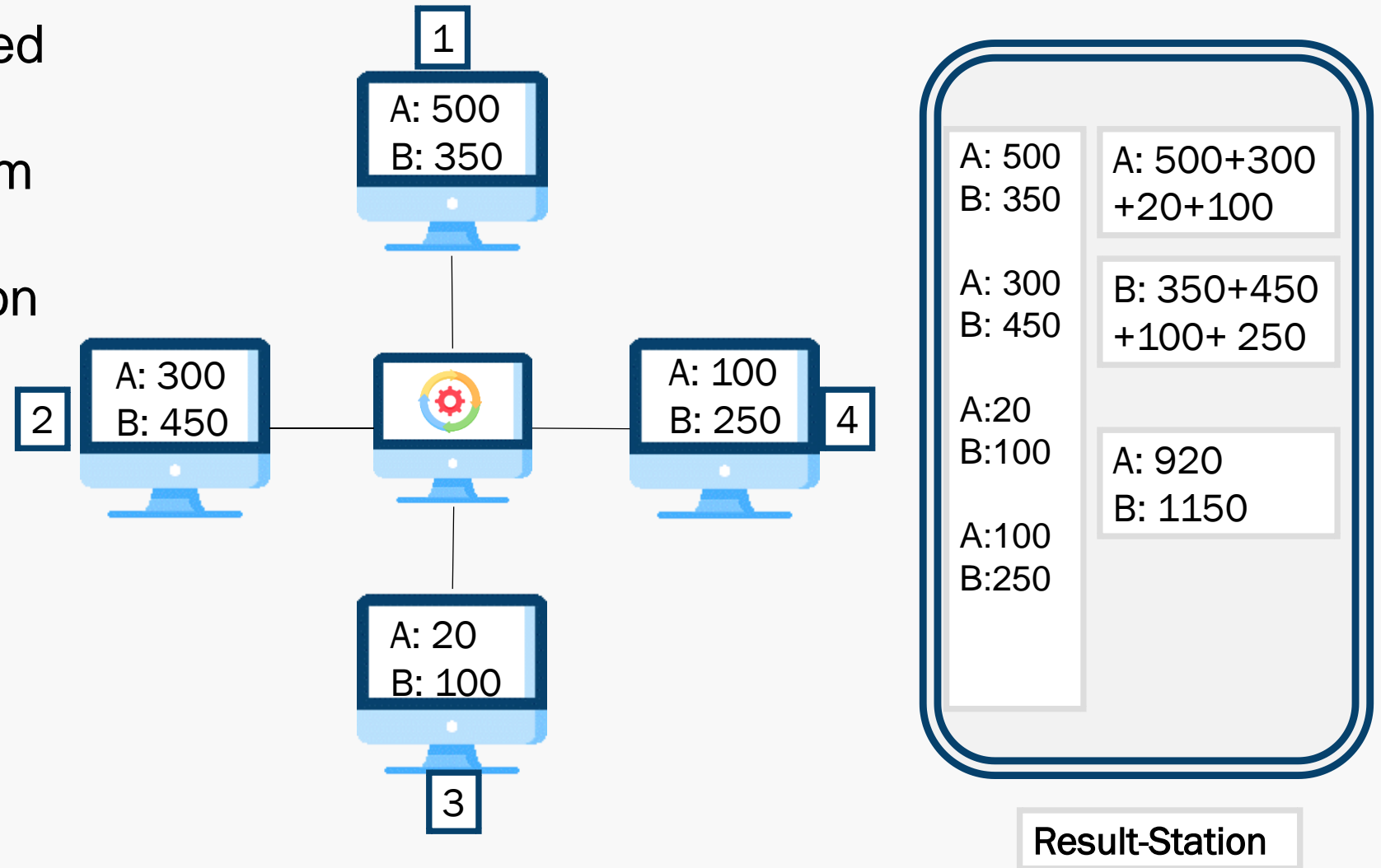
# MapReduce way of vote counting

- Votes will be counted in each station
- Resultant value from all stations will be send to result-station for finalizing.

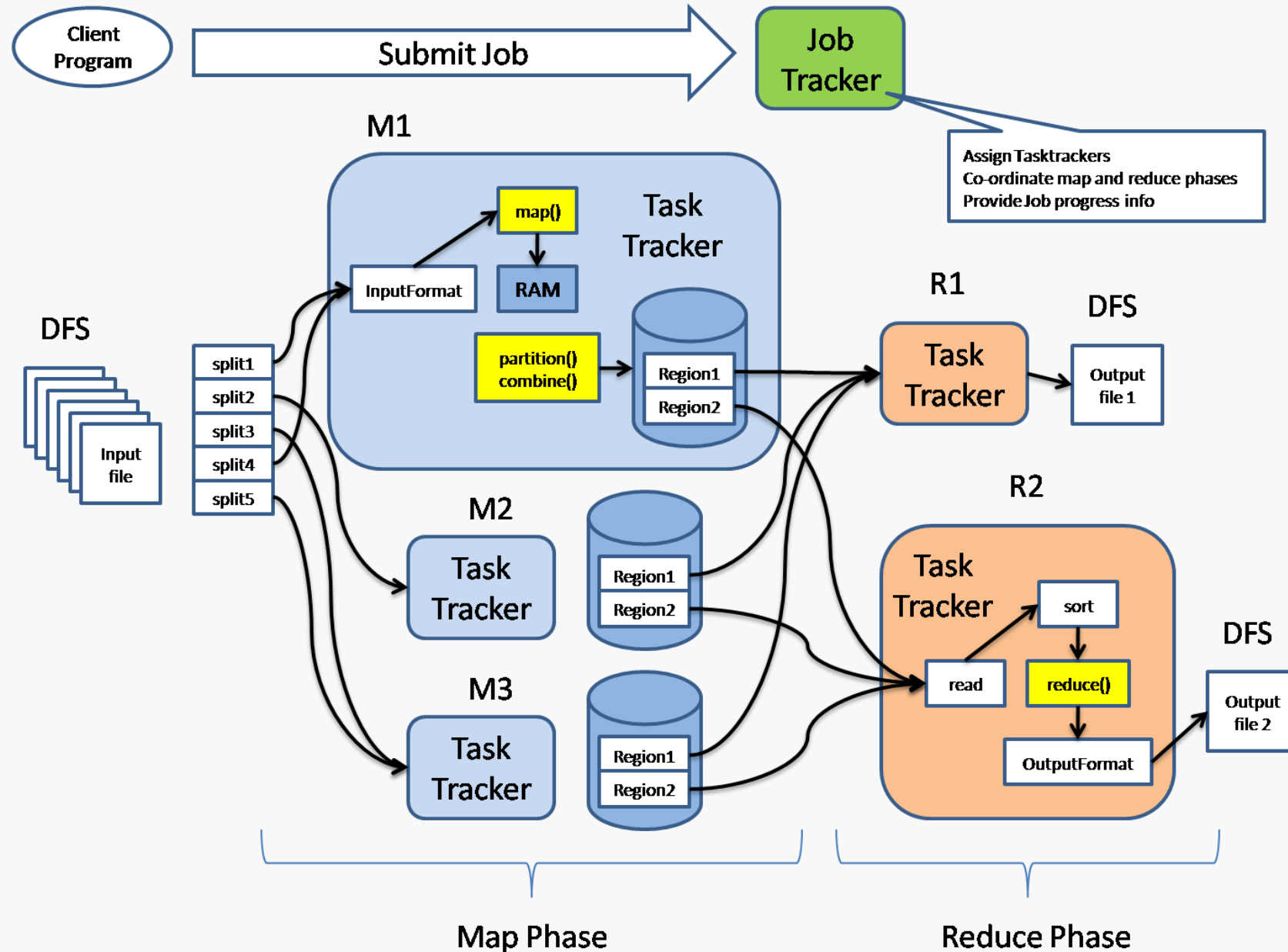


# MapReduce way of vote counting

- Votes will be counted in each station
- Resultant value from all stations will be send to result-station for finalizing.



# How MapReduce in Hadoop works



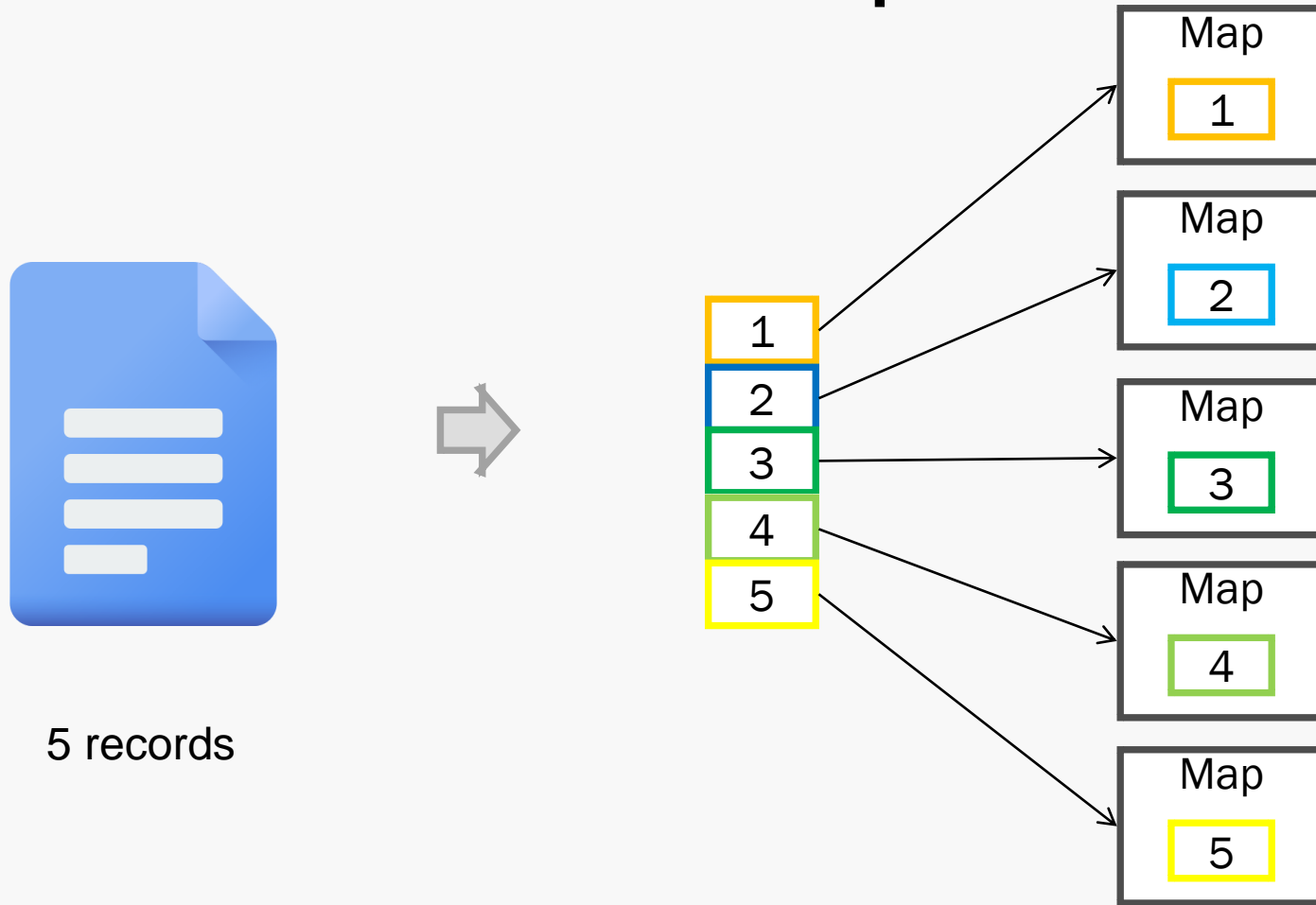
# MapReduce Functions

- At the crux of MapReduce are two functions: Map and Reduce. They are sequenced one after the other.
  - The **Map** function takes input from the disk as *<key,value> pairs*, processes them, and produces another set of intermediate *<key,value> pairs* as output.
  - The **Reduce** function also takes inputs as *<key,value> pairs*, and produces *<key,value> pairs* as output.

# Map

- The input data is first split into smaller blocks. Each block is then assigned to a mapper for processing.
- For example, if a file has 5 records to be processed, 5 mappers can run together to process single records each.

# Map



The Hadoop framework decides how many mappers to use, based on the size of the data to be processed and the memory block available on each mapper server.

# Reduce

- After all the mappers complete processing, the framework shuffles and sorts the results before passing them on to the reducers.
- A reducer cannot start while a mapper is still in progress.
- All the map output values that have the same key are assigned to a single reducer, which then aggregates the values for that key.



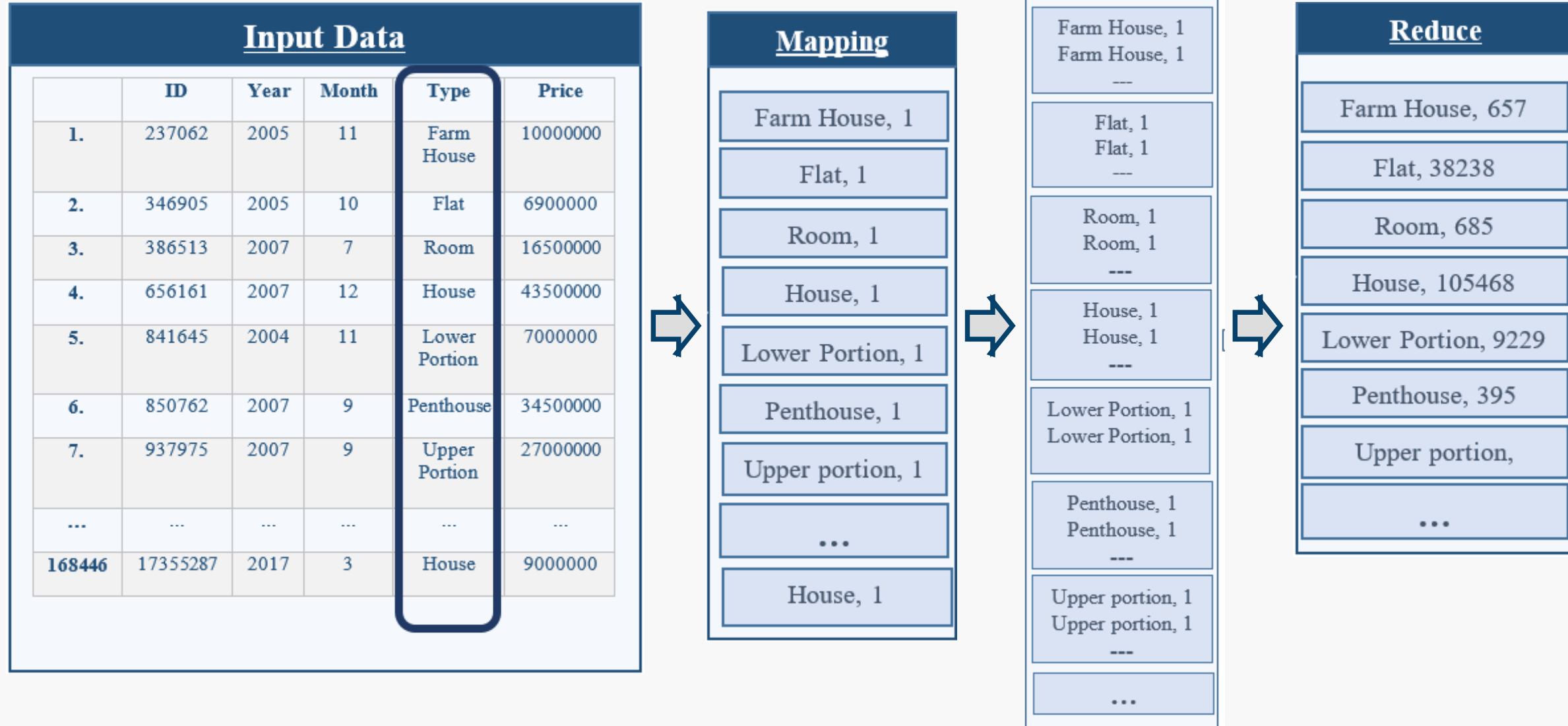
Let's have a practical view of  
MapReduce

# Property Dataset

	A	B	C	D	E
1	ID	Year	Month	Type	Price
2	237062	2005	11	Flat	10000000
3	346905	2005	10	Flat	6900000
4	386513	2007	7	House	16500000
5	656161	2007	12	House	43500000
6	841645	2004	11	House	7000000
168448	..	..	..	..	..
168449	17355287	2017	3	House	9000000

**Problem:** How many times has a particular type of property been sold?

**Problem:** How many times has a particular type of property been sold?



# Implementation in Spyder

The screenshot displays the Spyder IDE interface with a dark theme. The main editor window shows a Python script named `property.py` with the following code:

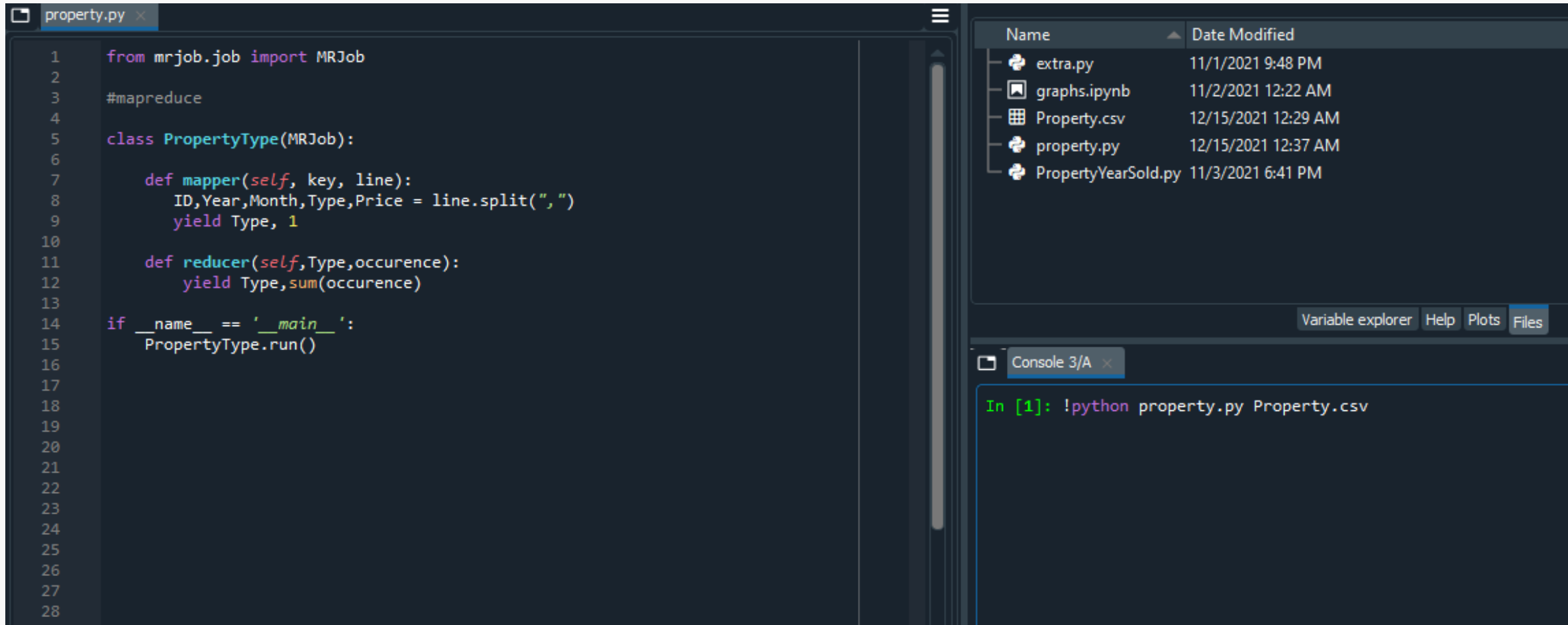
```
1  from mrjob.job import MRJob
2
3  #mapreduce
4
5  class PropertyType(MRJob):
6
7      def mapper(self, key, line):
8          ID,Year,Month,Type,Price = line.split(",")
9          yield Type, 1
10
11     def reducer(self,Type,occurence):
12         yield Type,sum(occurence)
13
14 if __name__ == '__main__':
15     PropertyType.run()
16
17
18
19
20
21
22
23
24
25
26
27
28
```

On the right side, the 'Files' panel shows a list of files in the current directory:

Name	Date Modified
extra.py	11/1/2021 9:48 PM
graphs.ipynb	11/2/2021 12:22 AM
Property.csv	12/15/2021 12:29 AM
property.py	11/1/2021 2:03 AM
PropertyYearSold.py	11/3/2021 6:41 PM

Below the file list, there are tabs for 'Variable explorer', 'Help', 'Plots', and 'Files'. At the bottom, the 'Console' panel shows the prompt `In [1]:`.

# Implementation in Spyder



The screenshot displays the Spyder IDE interface. The main editor window shows a Python script named `property.py` with the following code:

```
1 from mrjob.job import MRJob
2
3 #mapreduce
4
5 class PropertyType(MRJob):
6
7     def mapper(self, key, line):
8         ID,Year,Month,Type,Price = line.split(",")
9         yield Type, 1
10
11     def reducer(self,Type,occurence):
12         yield Type,sum(occurence)
13
14 if __name__ == '__main__':
15     PropertyType.run()
16
17
18
19
20
21
22
23
24
25
26
27
28
```

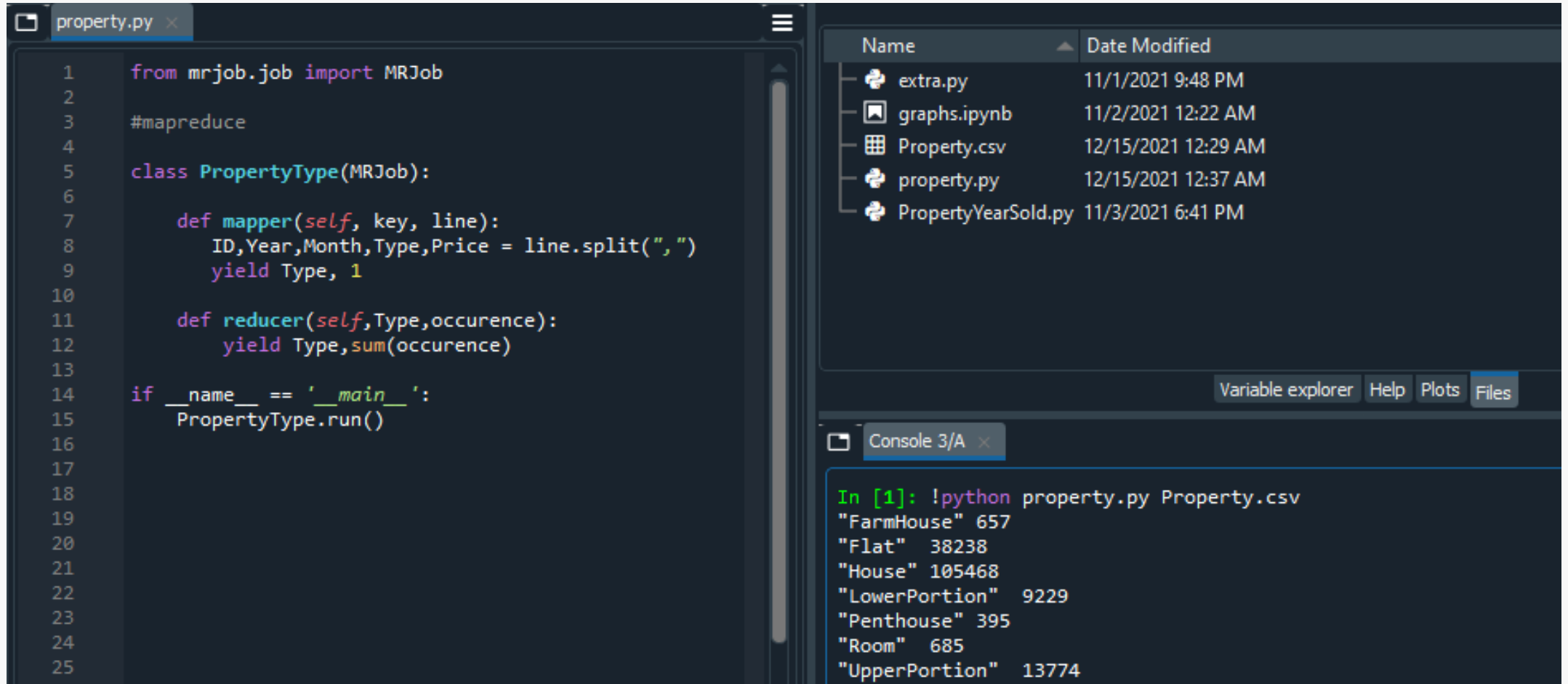
On the right side, the File explorer panel shows a list of files in the current project:

Name	Date Modified
extra.py	11/1/2021 9:48 PM
graphs.ipynb	11/2/2021 12:22 AM
Property.csv	12/15/2021 12:29 AM
property.py	12/15/2021 12:37 AM
PropertyYearSold.py	11/3/2021 6:41 PM

Below the File explorer, the Console panel shows the execution of the script:

```
In [1]: !python property.py Property.csv
```

# Implementation in Spyder



The screenshot displays the Spyder IDE interface. The main editor window shows a Python script named `property.py` with the following code:

```
1 from mrjob.job import MRJob
2
3 #mapreduce
4
5 class PropertyType(MRJob):
6
7     def mapper(self, key, line):
8         ID,Year,Month,Type,Price = line.split(",")
9         yield Type, 1
10
11     def reducer(self,Type,occurence):
12         yield Type,sum(occurence)
13
14 if __name__ == '__main__':
15     PropertyType.run()
16
17
18
19
20
21
22
23
24
25
```

The right sidebar contains a file explorer showing a list of files and their modification dates:

Name	Date Modified
extra.py	11/1/2021 9:48 PM
graphs.ipynb	11/2/2021 12:22 AM
Property.csv	12/15/2021 12:29 AM
property.py	12/15/2021 12:37 AM
PropertyYearSold.py	11/3/2021 6:41 PM

Below the file explorer, there are tabs for `Variable explorer`, `Help`, `Plots`, and `Files`. The `Console 3/A` tab is active, showing the output of the command `!python property.py Property.csv`:

```
In [1]: !python property.py Property.csv
"FarmHouse" 657
"Flat" 38238
"House" 105468
"LowerPortion" 9229
"Penthouse" 395
"Room" 685
"UpperPortion" 13774
```

# Result

**Problem:** How many times has a particular type of property been sold?

**Result:**

A number of particular type of property been sold is:

- **Farmhouse** has been sold **657** times
- **Flat** has been sold **38238** times
- **House** has been sold **105468** times
- **Lower Portion** has been sold **9229** times
- **Penthouse** has been sold **395** times
- **Room** has been sold **685** times
- **Upper Portion** has been sold **13774** times

# Property Dataset

	A	B	C	D	E
1	ID	Year	Month	Type	Price
2	237062	2005	11	Flat	10000000
3	346905	2005	10	Flat	6900000
4	386513	2007	7	House	16500000
5	656161	2007	12	House	43500000
6	841645	2004	11	House	7000000
168448	..	..	..	..	..
168449	17355287	2017	3	House	9000000

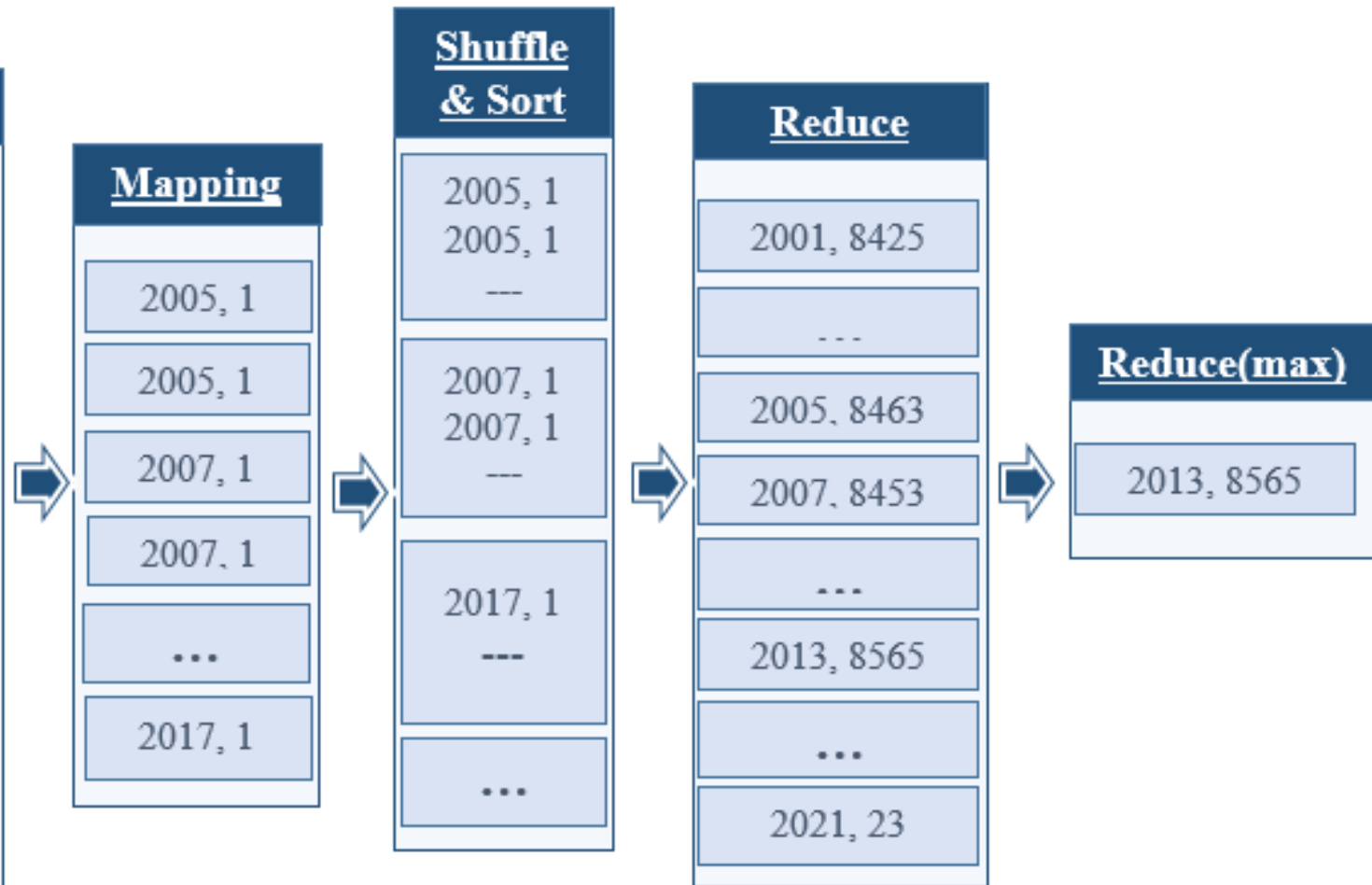
**Problem:** In what year were most of the properties sold?

1. Count of properties Sold per year?
2. Maximum of that count?

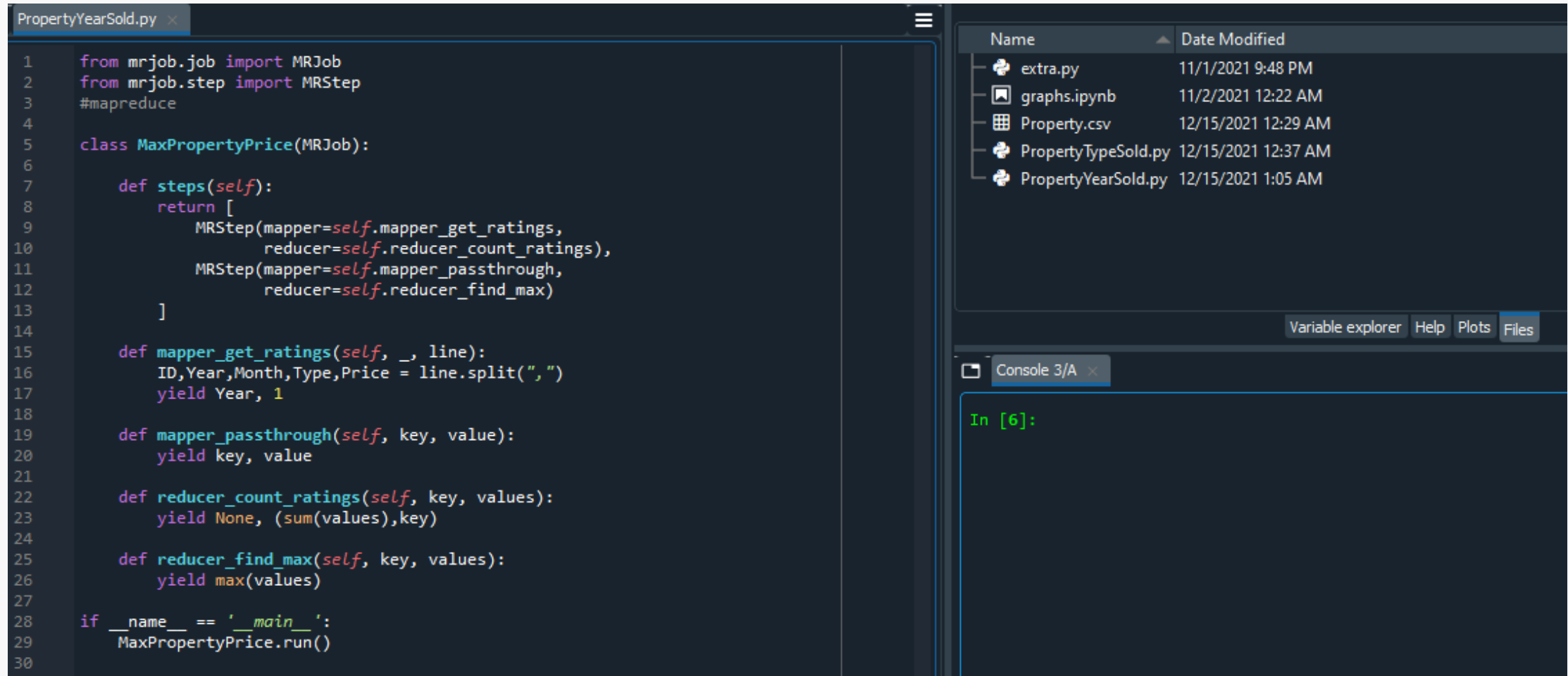


**Problem:** In what year were most of the properties sold?

Input Data					
	ID	Year	Month	Type	Price
1.	237062	2005	11	Farm House	10000000
2.	346905	2005	10	Flat	6900000
3.	386513	2007	7	Room	16500000
4.	656161	2007	12	House	43500000
...	...	...	...	...	...
168446	17355287	2017	3	House	9000000



# Implementation in Spyder



The screenshot displays the Spyder IDE interface. The main editor window shows a Python script named `PropertyYearSold.py` with the following code:

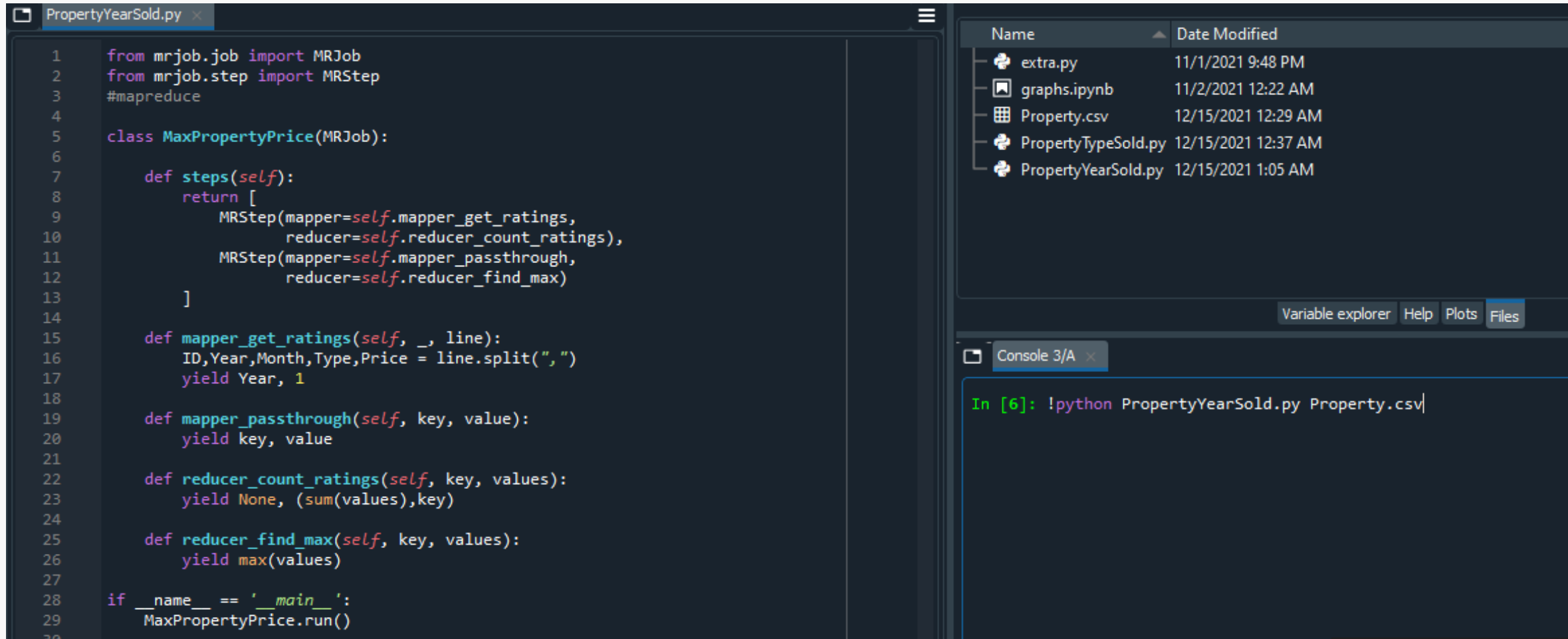
```
1  from mrjob.job import MRJob
2  from mrjob.step import MRStep
3  #mapreduce
4
5  class MaxPropertyPrice(MRJob):
6
7      def steps(self):
8          return [
9              MRStep(mapper=self.mapper_get_ratings,
10                     reducer=self.reducer_count_ratings),
11              MRStep(mapper=self.mapper_passthrough,
12                     reducer=self.reducer_find_max)
13          ]
14
15      def mapper_get_ratings(self, _, line):
16          ID,Year,Month,Type,Price = line.split(",")
17          yield Year, 1
18
19      def mapper_passthrough(self, key, value):
20          yield key, value
21
22      def reducer_count_ratings(self, key, values):
23          yield None, (sum(values),key)
24
25      def reducer_find_max(self, key, values):
26          yield max(values)
27
28  if __name__ == '__main__':
29      MaxPropertyPrice.run()
30
```

The right sidebar contains a file explorer with the following files and their modification dates:

Name	Date Modified
extra.py	11/1/2021 9:48 PM
graphs.ipynb	11/2/2021 12:22 AM
Property.csv	12/15/2021 12:29 AM
PropertyTypeSold.py	12/15/2021 12:37 AM
PropertyYearSold.py	12/15/2021 1:05 AM

Below the file explorer, there are tabs for `Variable explorer`, `Help`, `Plots`, and `Files`. The `Console 3/A` tab is active, showing the prompt `In [6]:`.

# Implementation in Spyder



The screenshot displays the Spyder IDE interface. The main editor window shows a Python script named `PropertyYearSold.py`. The script implements a MapReduce job to find the maximum property price. It imports `MRJob` and `MRStep` from `mrjob.job` and `mrjob.step`. A class `MaxPropertyPrice` inherits from `MRJob` and defines four methods: `steps`, `mapper_get_ratings`, `mapper_passthrough`, and `reducer_count_ratings` and `reducer_find_max`. The `steps` method returns a list of `MRStep` objects. The `mapper_get_ratings` method splits each line by commas and yields the year and price. The `mapper_passthrough` method yields the key and value. The `reducer_count_ratings` method yields the sum of values for each key. The `reducer_find_max` method yields the maximum value for each key. The script is executed in the console using the command `!python PropertyYearSold.py Property.csv`.

```
1 from mrjob.job import MRJob
2 from mrjob.step import MRStep
3 #mapreduce
4
5 class MaxPropertyPrice(MRJob):
6
7     def steps(self):
8         return [
9             MRStep(mapper=self.mapper_get_ratings,
10                    reducer=self.reducer_count_ratings),
11             MRStep(mapper=self.mapper_passthrough,
12                    reducer=self.reducer_find_max)
13         ]
14
15     def mapper_get_ratings(self, _, line):
16         ID,Year,Month,Type,Price = line.split(",")
17         yield Year, 1
18
19     def mapper_passthrough(self, key, value):
20         yield key, value
21
22     def reducer_count_ratings(self, key, values):
23         yield None, (sum(values),key)
24
25     def reducer_find_max(self, key, values):
26         yield max(values)
27
28 if __name__ == '__main__':
29     MaxPropertyPrice.run()
30
```

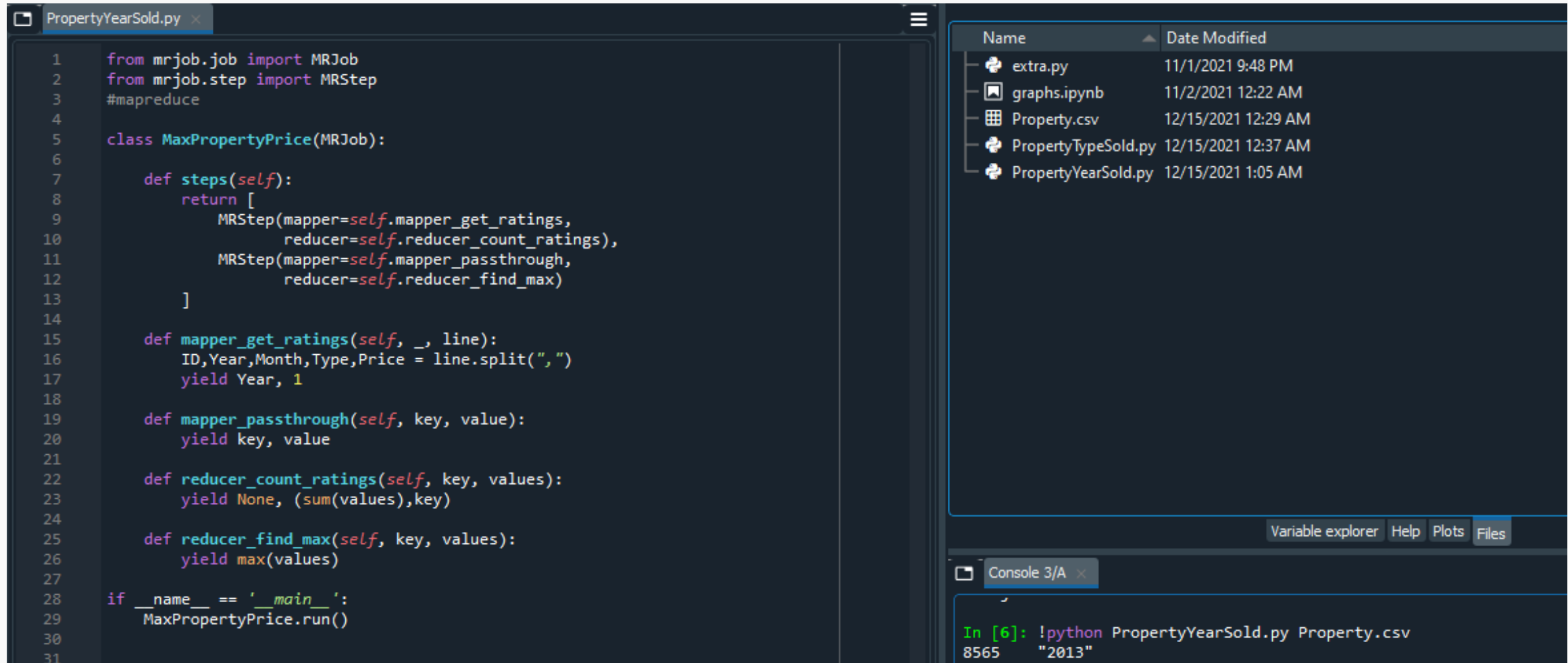
Name	Date Modified
extra.py	11/1/2021 9:48 PM
graphs.ipynb	11/2/2021 12:22 AM
Property.csv	12/15/2021 12:29 AM
PropertyTypeSold.py	12/15/2021 12:37 AM
PropertyYearSold.py	12/15/2021 1:05 AM

Variable explorer Help Plots Files

Console 3/A

```
In [6]: !python PropertyYearSold.py Property.csv
```

# Implementation in Spyder



The screenshot displays the Spyder IDE interface. The main editor window shows a Python script named `PropertyYearSold.py` with the following code:

```
1  from mrjob.job import MRJob
2  from mrjob.step import MRStep
3  #mapreduce
4
5  class MaxPropertyPrice(MRJob):
6
7      def steps(self):
8          return [
9              MRStep(mapper=self.mapper_get_ratings,
10                     reducer=self.reducer_count_ratings),
11              MRStep(mapper=self.mapper_passthrough,
12                     reducer=self.reducer_find_max)
13          ]
14
15      def mapper_get_ratings(self, _, line):
16          ID,Year,Month,Type,Price = line.split(",")
17          yield Year, 1
18
19      def mapper_passthrough(self, key, value):
20          yield key, value
21
22      def reducer_count_ratings(self, key, values):
23          yield None, (sum(values),key)
24
25      def reducer_find_max(self, key, values):
26          yield max(values)
27
28  if __name__ == '__main__':
29      MaxPropertyPrice.run()
30
31
```

The right sidebar contains a file explorer showing the project structure:

Name	Date Modified
extra.py	11/1/2021 9:48 PM
graphs.ipynb	11/2/2021 12:22 AM
Property.csv	12/15/2021 12:29 AM
PropertyTypeSold.py	12/15/2021 12:37 AM
PropertyYearSold.py	12/15/2021 1:05 AM

Below the file explorer, the console window shows the execution of the script:

```
In [6]: !python PropertyYearSold.py Property.csv
8565    "2013"
```

# Result

**Problem:** In what year were most of the properties sold?

## Result

Most of the properties were sold in **2013**.

- **As** 8565 Properties were sold in 2013 which is the maximum count.)

# Store Dataset

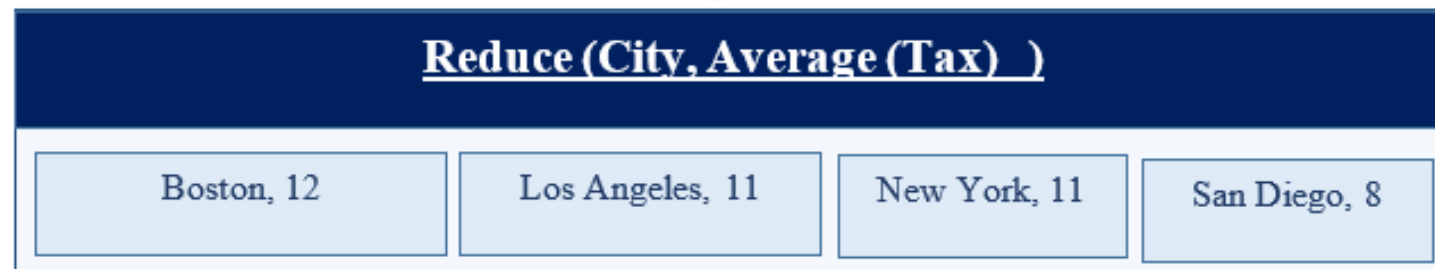
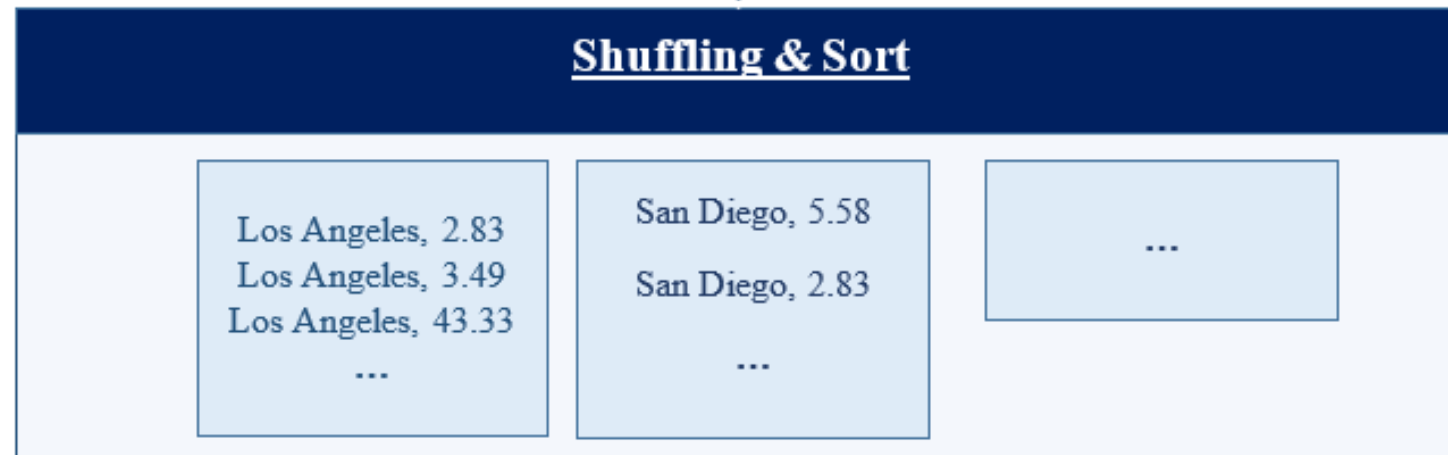
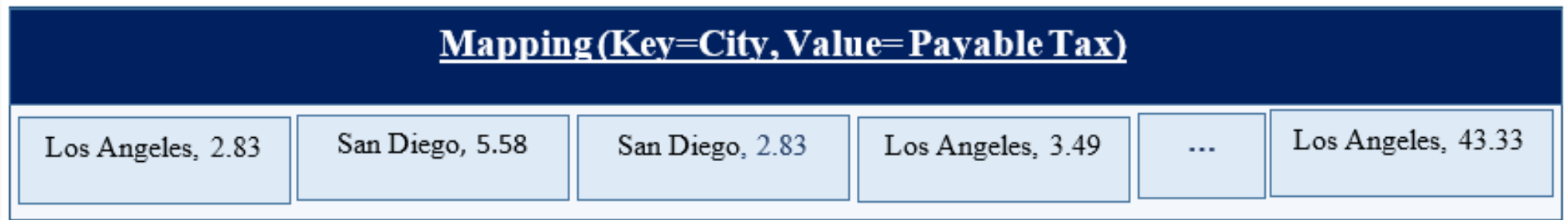
	A	B	C	D	E	F	G	H	I	J	K
1	OrderDate	Region	City	Category	Product	Quantity	UnitPrice	TotalPrice	Payable_tax	Payable_amount	
2	43577	West	Los Angeles	Bars	Carrot	20	1.77	Rs 35.40	2.83	38.23	
3	43628	West	San Diego	Crackers	Whole Wheat	20	3.49	Rs 69.80	5.58	75.38	
4	43760	West	San Diego	Bars	Carrot	20	1.77	Rs 35.40	2.83	38.23	
5	44039	West	Los Angeles	Cookies	Arrowroot	20	2.18	Rs 43.60	3.49	47.09	
6	44090	East	New York	Snacks	Potato Chips	20	1.68	Rs 33.60	2.69	36.29	
7	44156	West	San Diego	Bars	Carrot	20	1.77	Rs 35.40	2.83	38.23	
8	44165	East	Boston	Crackers	Whole Wheat	20	3.49	Rs 69.80	5.58	75.38	
9	43598	West	Los Angeles	Crackers	Whole Wheat	21	3.49	Rs 73.29	5.86	79.15	
10	43700	East	Boston	Crackers	Whole Wheat	21	3.49	Rs 73.29	5.86	79.15	

**Problem:** What is the average tax paid by each City?

**Problem:** What is the average tax paid by each City?

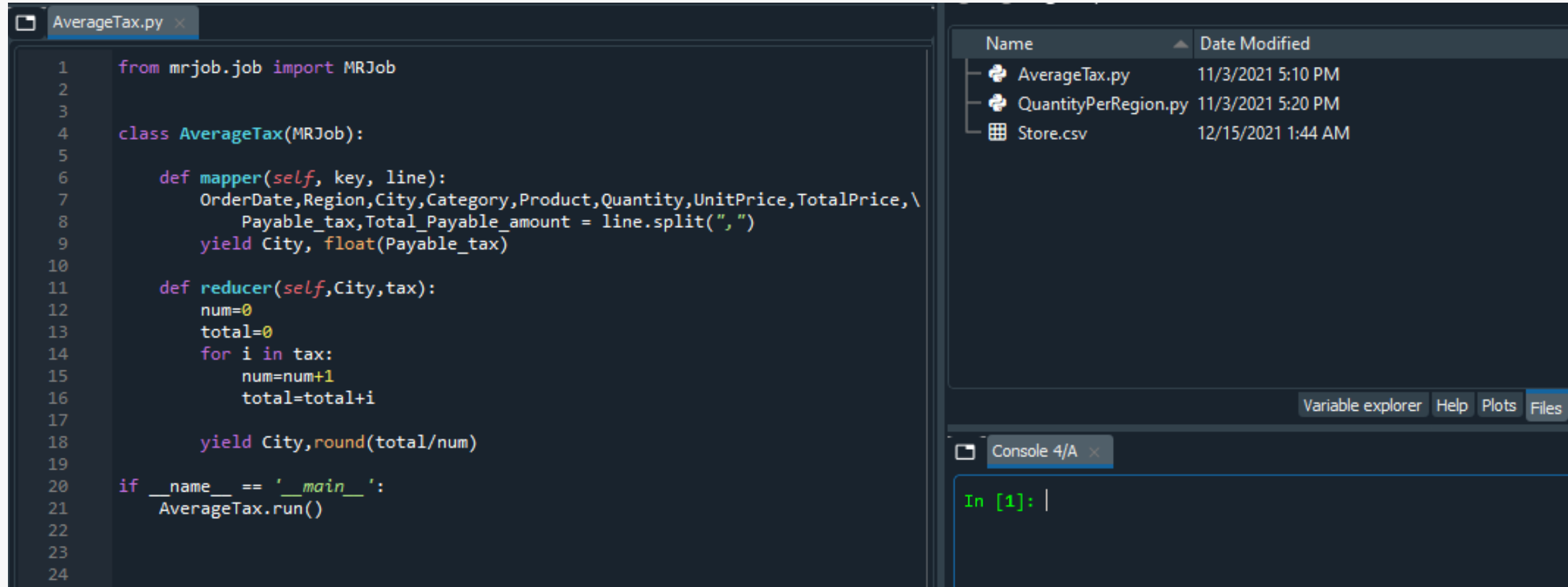
<u>Input Data</u>										
	OrderDate	Region	City	Category	Product	Quantity	UnitPrice	TotalPrice	Payable_tax	Total_Payable_amount
1.	43577	West	Los Angeles	Bars	Carrot	20	1.77	Rs 35.40	2.83	38.23
2.	43628	West	San Diego	Crackers	Whole Wheat	20	3.49	Rs 69.80	5.58	75.38
3.	43760	West	San Diego	Bars	Carrot	20	1.77	Rs 35.40	2.83	38.23
4.	44039	West	Los Angeles	Cookies	Arrowroot	20	2.18	Rs 43.60	3.49	47.09
...	...									
244	43637	West	Los Angeles	Bars	Carrot	306	1.77	Rs 541.62	43.33	584.95

**Problem:** What is the average tax paid by each City?





# Implementation in Spyder



The screenshot displays the Spyder IDE interface. The main editor window shows a Python script named `AverageTax.py`. The script imports `MRJob` from `mrjob.job` and defines a class `AverageTax` that inherits from `MRJob`. The class contains two methods: `mapper` and `reducer`. The `mapper` method takes `key` and `line` as arguments, splits the `line` by commas, and yields the `City` and the `Payable_tax` as a float. The `reducer` method takes `City` and `tax` as arguments, calculates the sum of `tax` values, and yields the `City` and the average tax value. The script also includes a standard `if __name__ == '__main__':` block to run the `AverageTax` class.

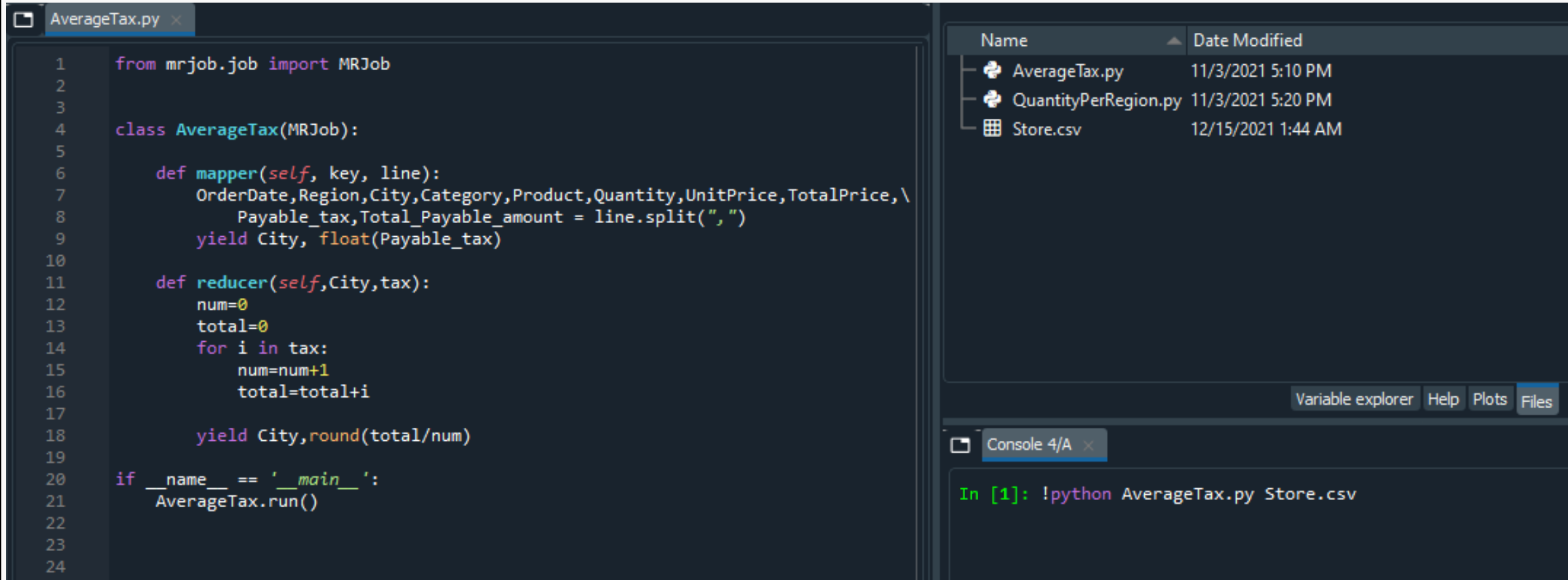
```
1 from mrjob.job import MRJob
2
3
4 class AverageTax(MRJob):
5
6     def mapper(self, key, line):
7         OrderDate,Region,City,Category,Product,Quantity,UnitPrice,TotalPrice,\
8             Payable_tax,Total_Payable_amount = line.split(",")
9         yield City, float(Payable_tax)
10
11    def reducer(self, City, tax):
12        num=0
13        total=0
14        for i in tax:
15            num=num+1
16            total=total+i
17
18        yield City, round(total/num)
19
20 if __name__ == '__main__':
21     AverageTax.run()
22
23
24
```

The right sidebar shows the file explorer with the following files and their modification dates:

Name	Date Modified
AverageTax.py	11/3/2021 5:10 PM
QuantityPerRegion.py	11/3/2021 5:20 PM
Store.csv	12/15/2021 1:44 AM

Below the file explorer, there are tabs for `Variable explorer`, `Help`, `Plots`, and `Files`. The `Console 4/A` tab is active, showing the prompt `In [1]:`.

# Implementation in Spyder



```
1 from mrjob.job import MRJob
2
3
4 class AverageTax(MRJob):
5
6     def mapper(self, key, line):
7         OrderDate,Region,City,Category,Product,Quantity,UnitPrice,TotalPrice,\
8             Payable_tax,Total_Payable_amount = line.split(",")
9         yield City, float(Payable_tax)
10
11     def reducer(self, City, tax):
12         num=0
13         total=0
14         for i in tax:
15             num=num+1
16             total=total+i
17
18         yield City, round(total/num)
19
20 if __name__ == '__main__':
21     AverageTax.run()
22
23
24
```

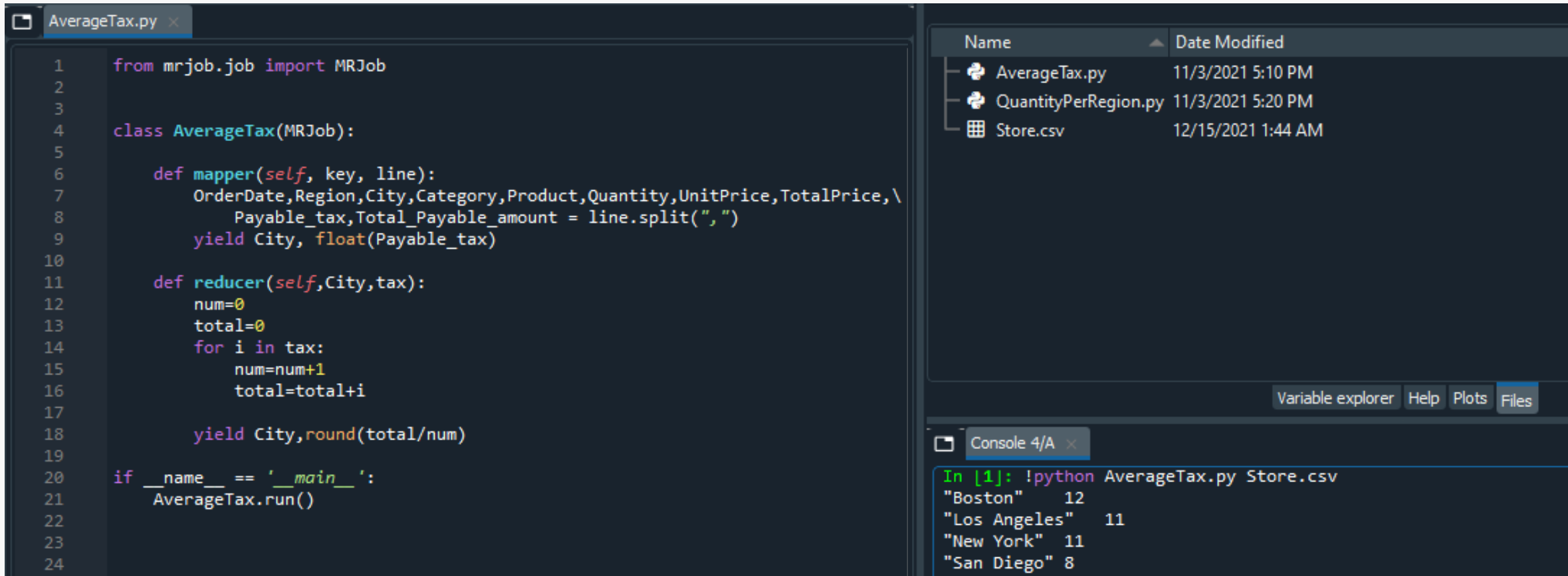
Name	Date Modified
AverageTax.py	11/3/2021 5:10 PM
QuantityPerRegion.py	11/3/2021 5:20 PM
Store.csv	12/15/2021 1:44 AM

Variable explorer Help Plots Files

Console 4/A

```
In [1]: !python AverageTax.py Store.csv
```

# Implementation in Spyder



The screenshot displays the Spyder IDE interface. The main editor window shows the file `AverageTax.py` with the following Python code:

```
1 from mrjob.job import MRJob
2
3
4 class AverageTax(MRJob):
5
6     def mapper(self, key, line):
7         OrderDate,Region,City,Category,Product,Quantity,UnitPrice,TotalPrice,\
8             Payable_tax,Total_Payable_amount = line.split(",")
9         yield City, float(Payable_tax)
10
11     def reducer(self, City, tax):
12         num=0
13         total=0
14         for i in tax:
15             num=num+1
16             total=total+i
17
18         yield City, round(total/num)
19
20 if __name__ == '__main__':
21     AverageTax.run()
22
23
24
```

On the right side, the 'Files' pane shows a directory structure with the following files and their modification dates:

Name	Date Modified
AverageTax.py	11/3/2021 5:10 PM
QuantityPerRegion.py	11/3/2021 5:20 PM
Store.csv	12/15/2021 1:44 AM

Below the 'Files' pane, the 'Console' pane shows the output of running the script:

```
In [1]: !python AverageTax.py Store.csv
Boston"    12
Los Angeles"  11
New York"   11
San Diego"   8
```

# Result

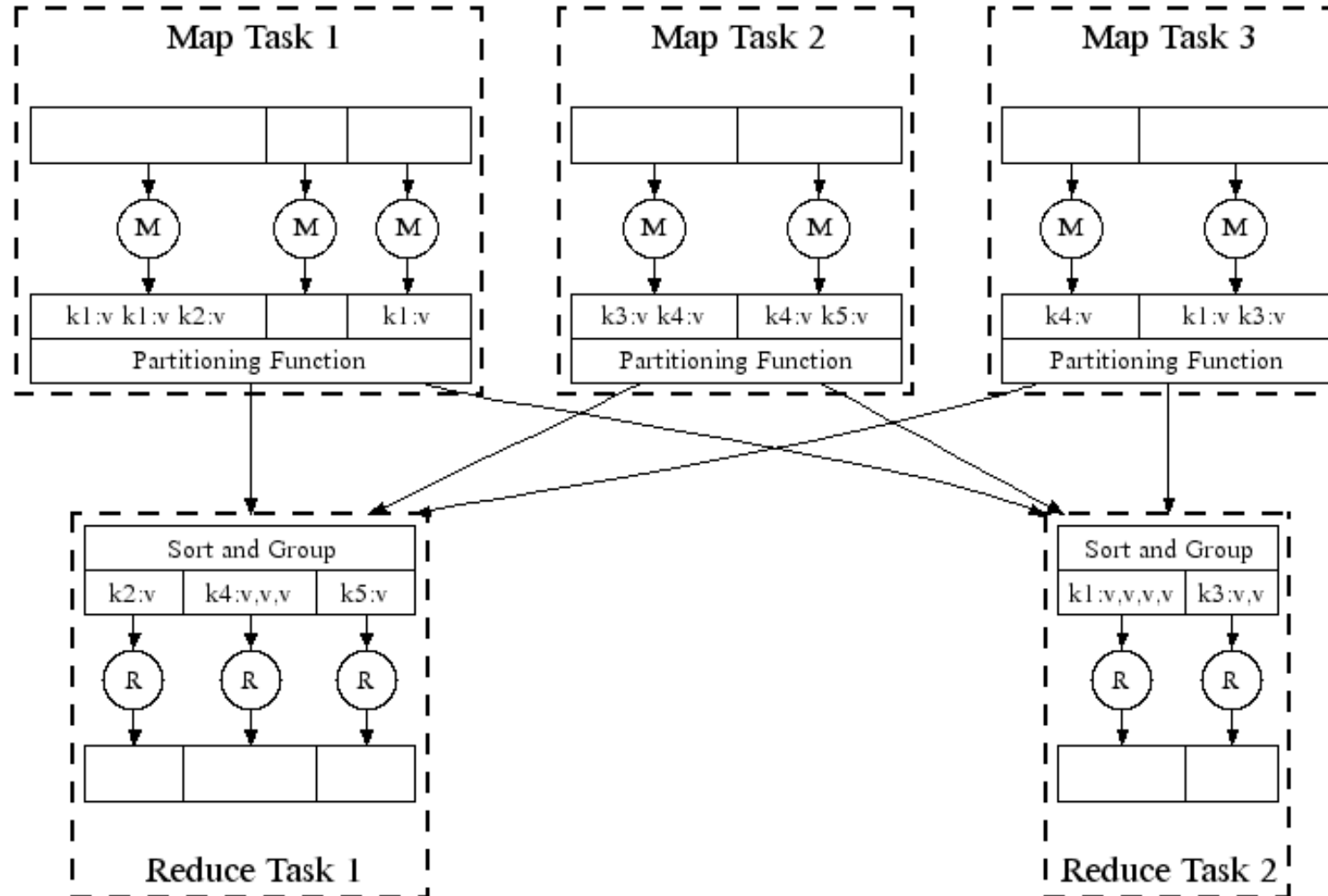
**Problem:** What is the average tax paid by each City?

## Result

The average tax paid by each City is given below:

- **Boston** has paid average tax of **12**.
- **Los Angeles** has paid average tax of **11**.
- **New York** has paid average tax of **11**.
- **San Diego** has paid average tax of **8**.

# How it is related to Big Data?



Let's Discuss a Problem

# Problem: Word Count



- **Statement:** Different observers saw different animals in forest. Count the occurrence of each animal, use Map-Reduce to the count ( occurrence) of words.

Deer, Bear, Dog  
Cat, Cat, Dog  
Deer, Cat, Bear



**Question:** What are the Mapper and Reducer Functions?



# Mapping



Deer



Bear



Dog

Deer: 1  
Bear: 1  
Dog: 1



Cat



Cat



Dog

Cat: 1  
Cat: 1  
Dog: 1



Deer



Cat



Bear

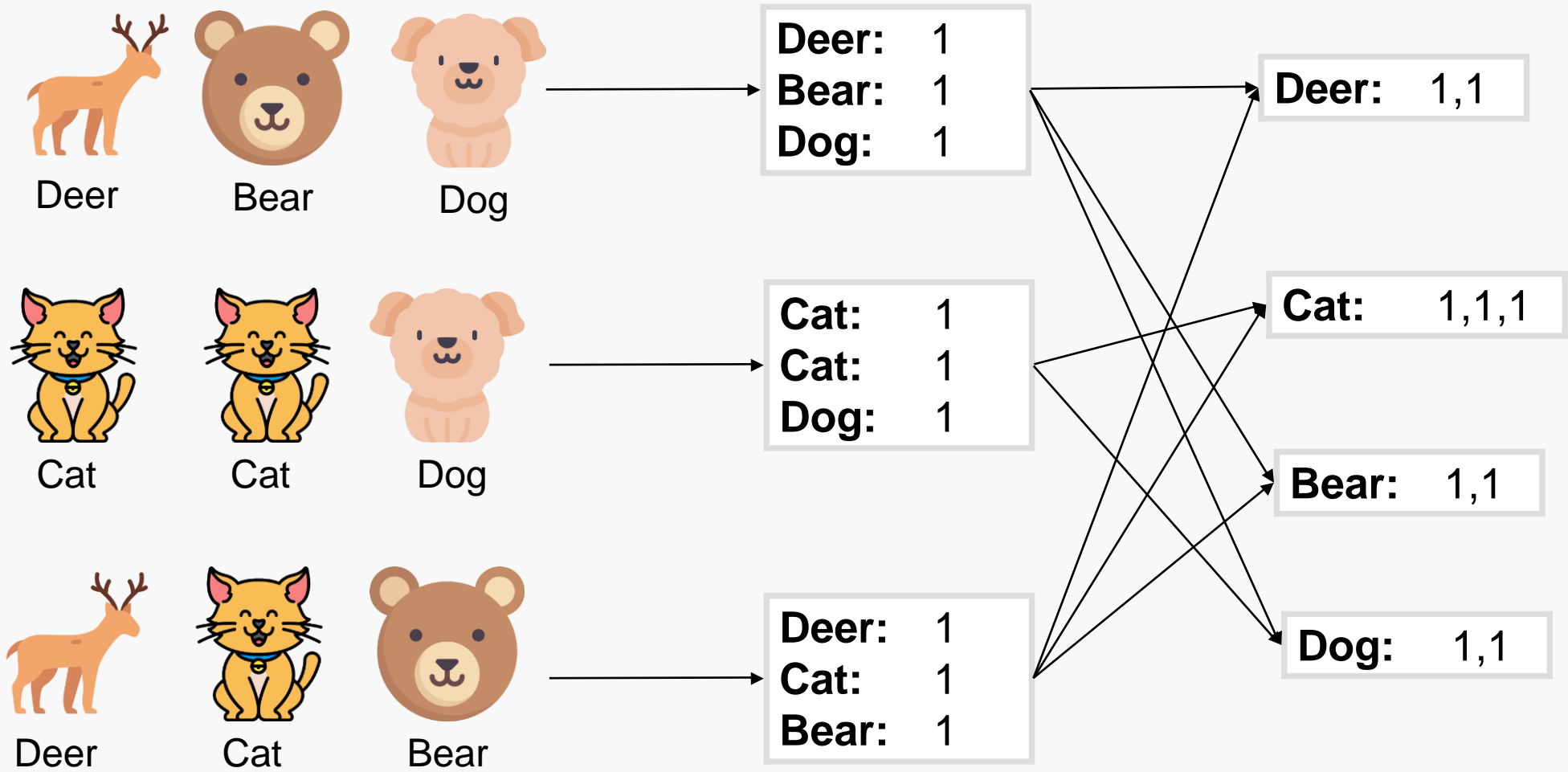
Deer: 1  
Cat: 1  
Bear: 1

## MAP

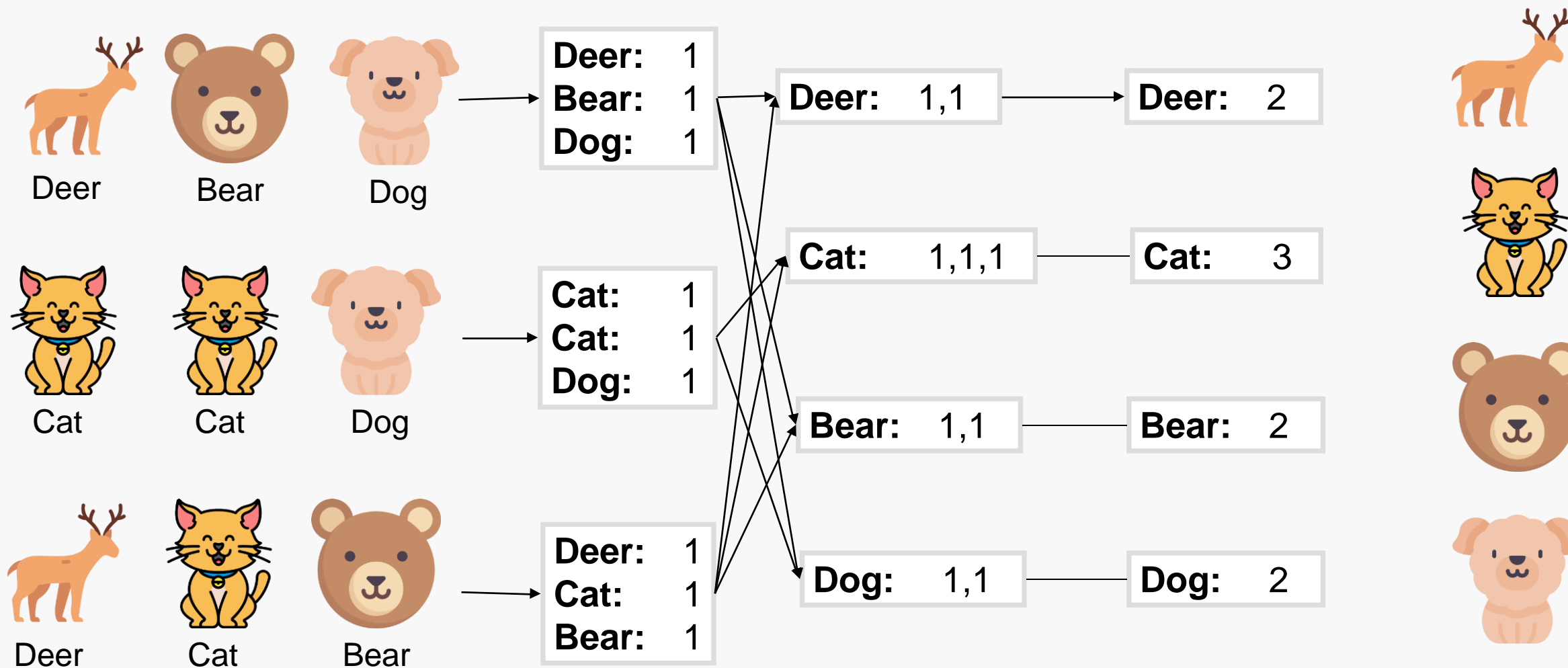
(key, value)  
(animal, occurrence)



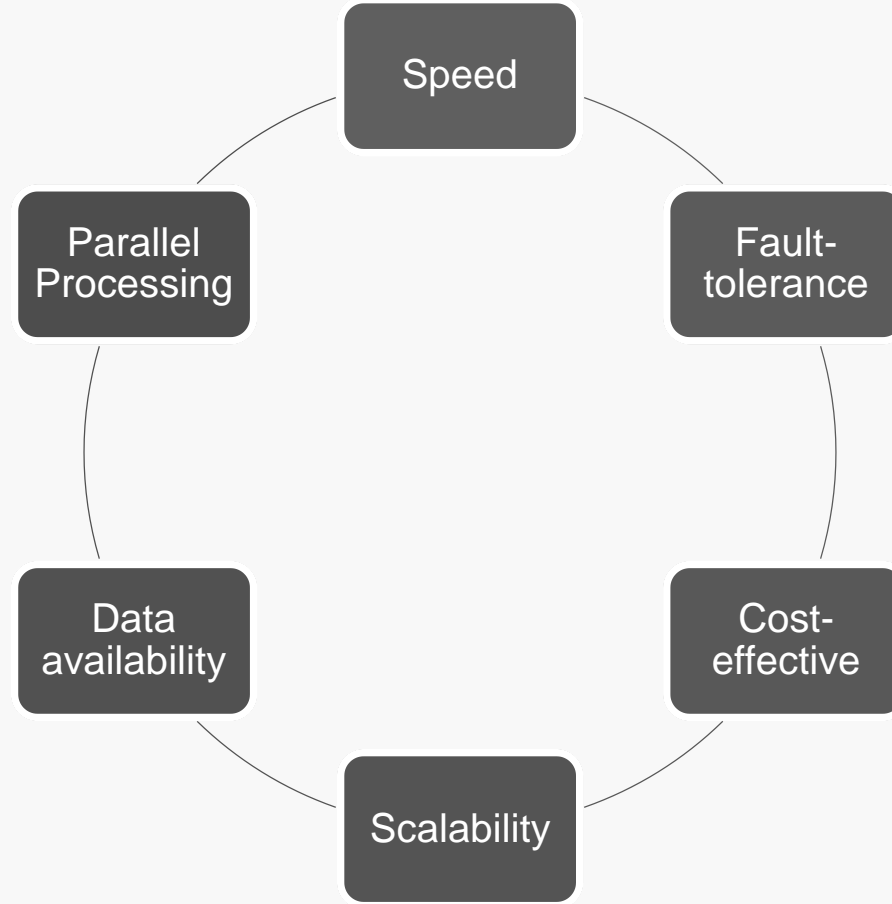
# Shuffling and Sort



# Reduce



# Benefits of Hadoop MapReduce



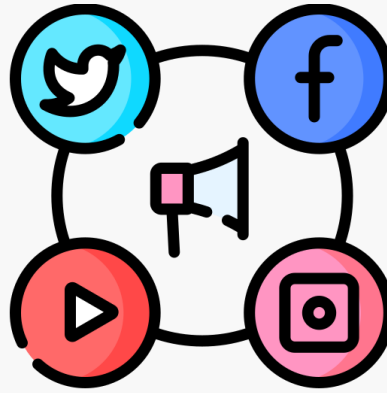
# Used-in

- Index and Search
- Classification
- Recommendation system
- Analytics

# Applications of MapReduce



E-Commerce



Social Media



Entertainment

# Summary

- MapReduce is a crucial processing component of the Hadoop framework. It's a quick, scalable, and cost-effective program that can help data analysts and developers process huge data.
- This programming model is a suitable tool for analyzing usage patterns on websites and e-commerce platforms. Companies providing online services can utilize this framework to improve their marketing strategies.



# Limitation

Map Reduce is bad for:

- Frequently changing data
- Dependent
- Interactive Analytics



# References

- **Fundamentals of MapReduce with MapReduce Example**

<https://medium.com/edureka/mapreduce-tutorial-3d9535ddbe7c>

- **History and Advantages of Hadoop MapReduce Programming**

<https://mindmajix.com/mapreduce/history-and-advantages-of-hadoop-mapreduce-programming>

- **What is MapReduce?**

<https://www.talend.com/resources/what-is-mapreduce/>

- **MapReduce:  
Simplified Data Processing on Large Clusters**

<https://research.google.com/archive/mapreduce-osdi04-slides/index.html>

Thank You!