

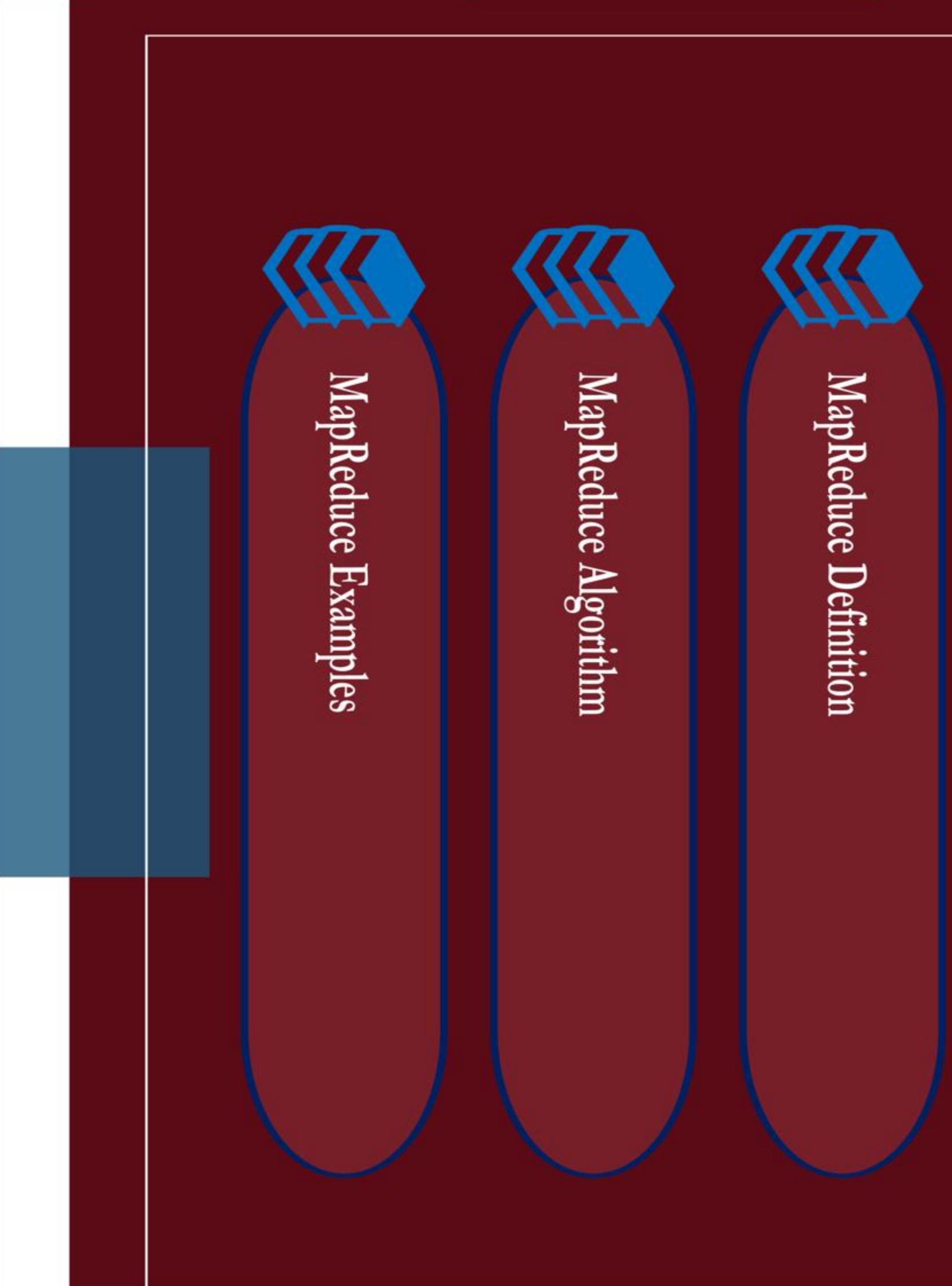
Agenda

Lecture 2

Big Data Processing Techniques

(MapReduce)

Dr. Lydia Wahid



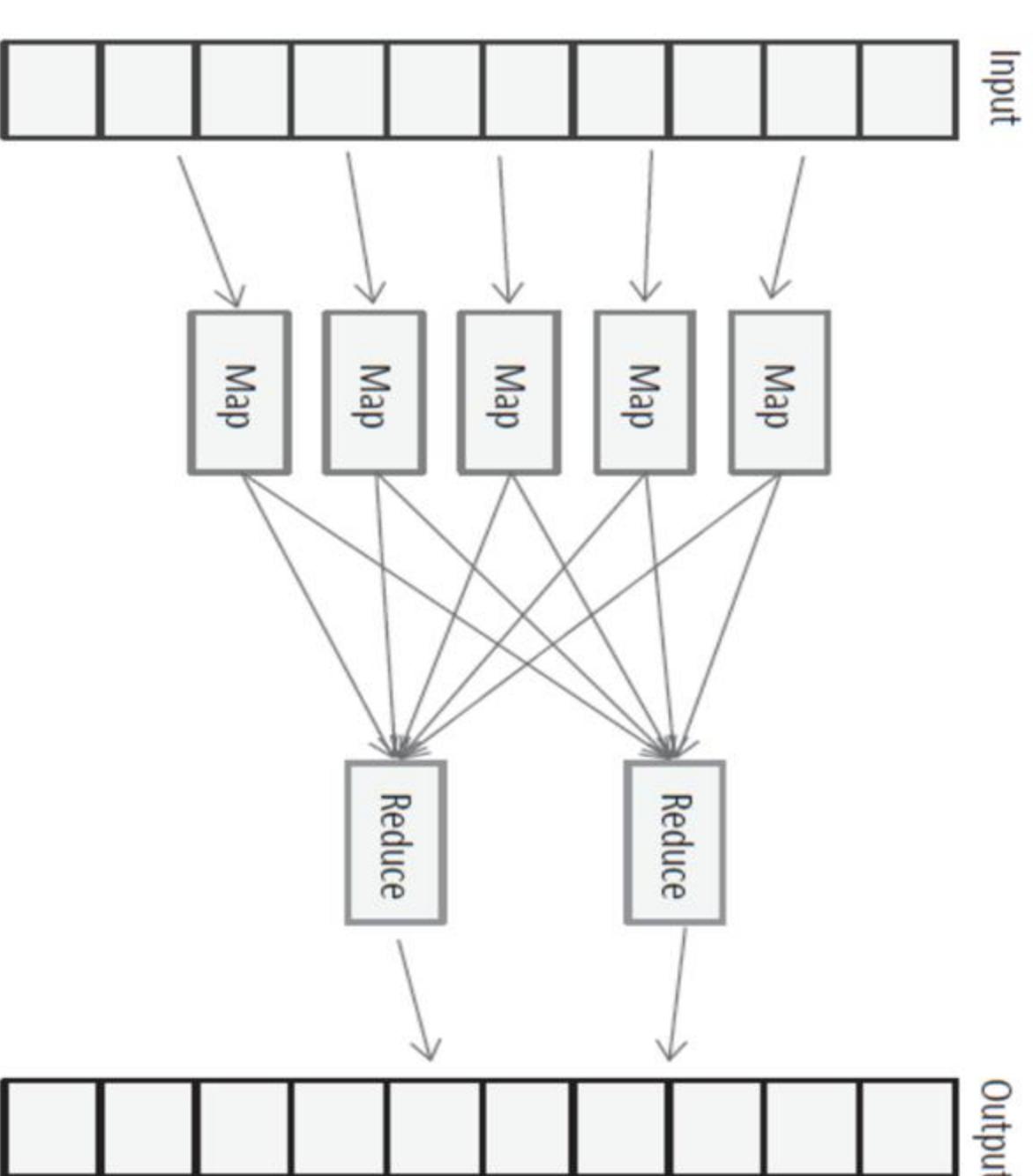
MapReduce Definition

- MapReduce is a widely used Big data processing technique.
 - It processes large datasets using parallel processing deployed over clusters of hardware.
 - It is based on the principle of **divide-and-conquer**. It divides a big problem into a collection of smaller problems that can each be solved quickly.
 - A **dataset is broken down** into multiple smaller parts, and operations are performed on each part independently and in parallel.
 - The results from all operations are then **combined** to arrive at the result of the whole dataset.

MapReduce Definition

- Each MapReduce job is composed of a **map phase** and a **reduce phase** and each phase consists of multiple stages.
 - The Map and Reduce phases run **sequentially** in a cluster.
 - The Map phase is executed first then the Reduce phase.
 - The output of the Map phase becomes the input of the Reduce phase.
 - MapReduce does not require that the input data conform to any particular data model.

MapReduce Definition



MapReduce Definition

- In MapReduce, all map and reduce tasks run in parallel.

- First of all, all map tasks are independently run.

- Meanwhile, reduce tasks wait until their respective maps are finished.

- Then, reduce tasks process their data concurrently and independently.



MapReduce Algorithm

- We will now apply and explain each stage on the following example:

- Problem Statement:**

Count the number of occurrences of each word available in a DataSet.

Input Dataset



Required Output

1 Black = 3
2 Blue = 6
3 Green = 5
4 Orange = 1
5 Red = 7
6 White = 3
7

- By applying this stage on our example, we get the following:

Input Dataset



Split stage

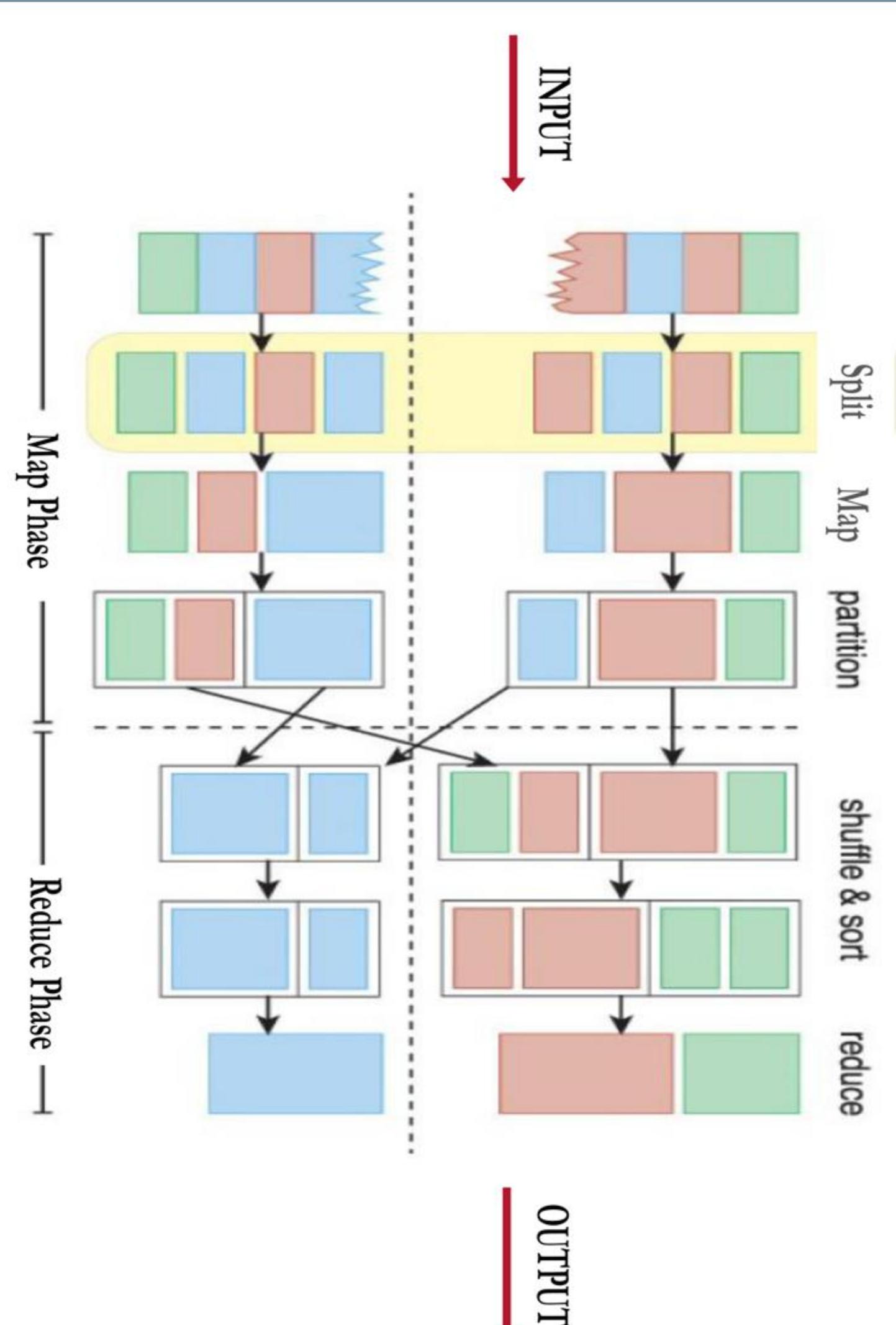
1 Red Blue Red Blue Green Red Blue Green
2 White Black
3 Red White Black
4 Orange Green
5 Red Blue Red
6 Blue Green Red Blue
7 Green White Black

MapReduce Algorithm

- Split stage:

- Takes input DataSet and divides it into smaller Sub-DataSets called splits.
- Each split is parsed into its constituent records as a key-value pair. The key is usually the ordinal position of the record, and the value is the actual record.
- A common example will read a directory full of text files and return each line as a record.
- The key-value pairs for each split are then sent to a map function (or mapper).

By applying this stage on our example, we get the following:



MapReduce Algorithm

- We will now apply and explain each stage on the following example:

- Problem Statement:**

Count the number of occurrences of each word available in a DataSet.

Input Dataset

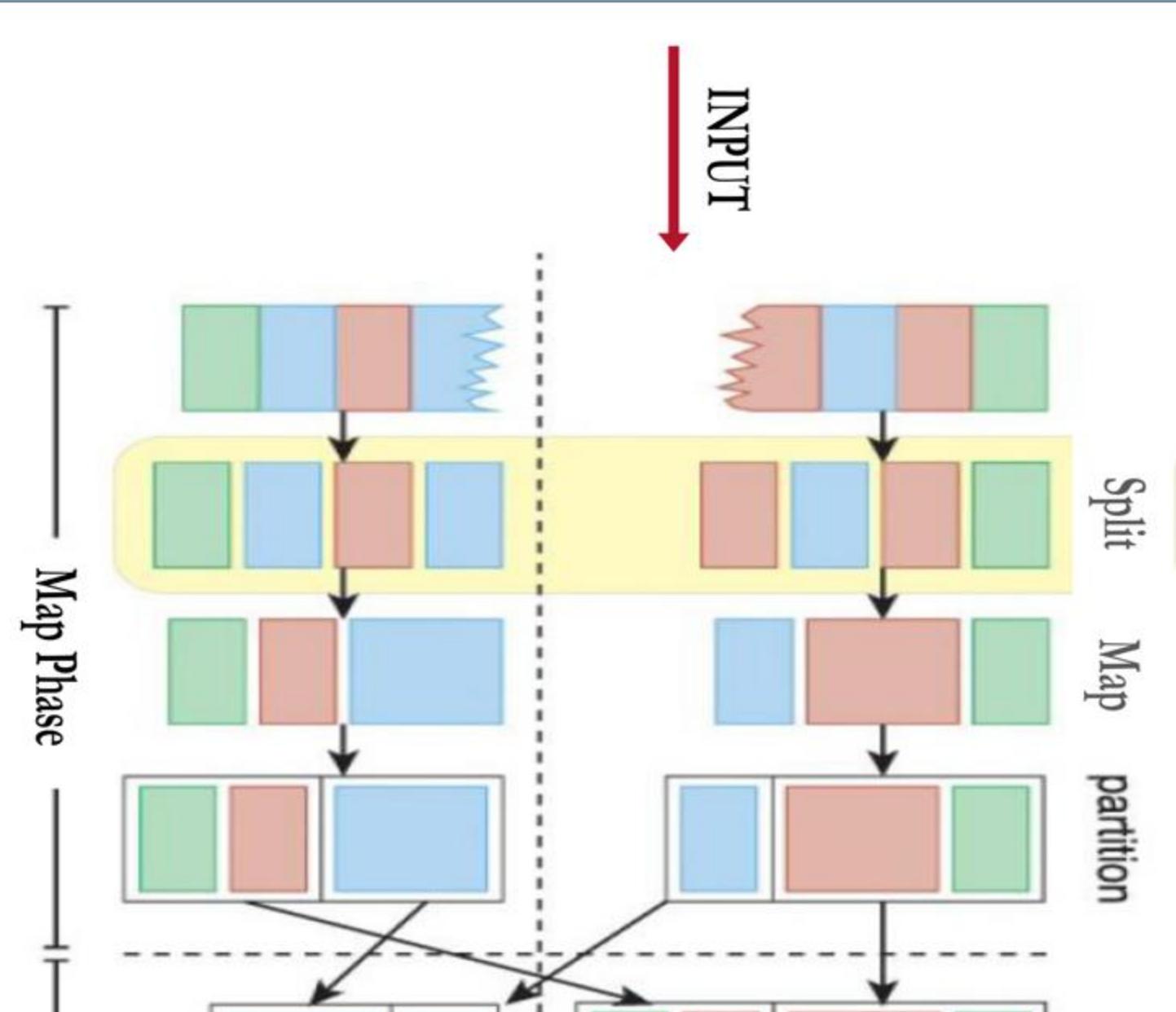


1 Sub-DataSet-1
2
3 Red Blue Red Blue Green Red Blue Green
4 White Black
5 Red White Black

By applying this stage on our example, we get the following:

1 Sub-DataSet-2
2
3 Orange Green
4 Red Blue Red
5 Blue Green Red Blue

1 Red Blue Red Blue Green Red Blue Green
2 White Black
3 Red White Black
4 Orange Green
5 Red Blue Red
6 Blue Green Red Blue
7 Green White Black

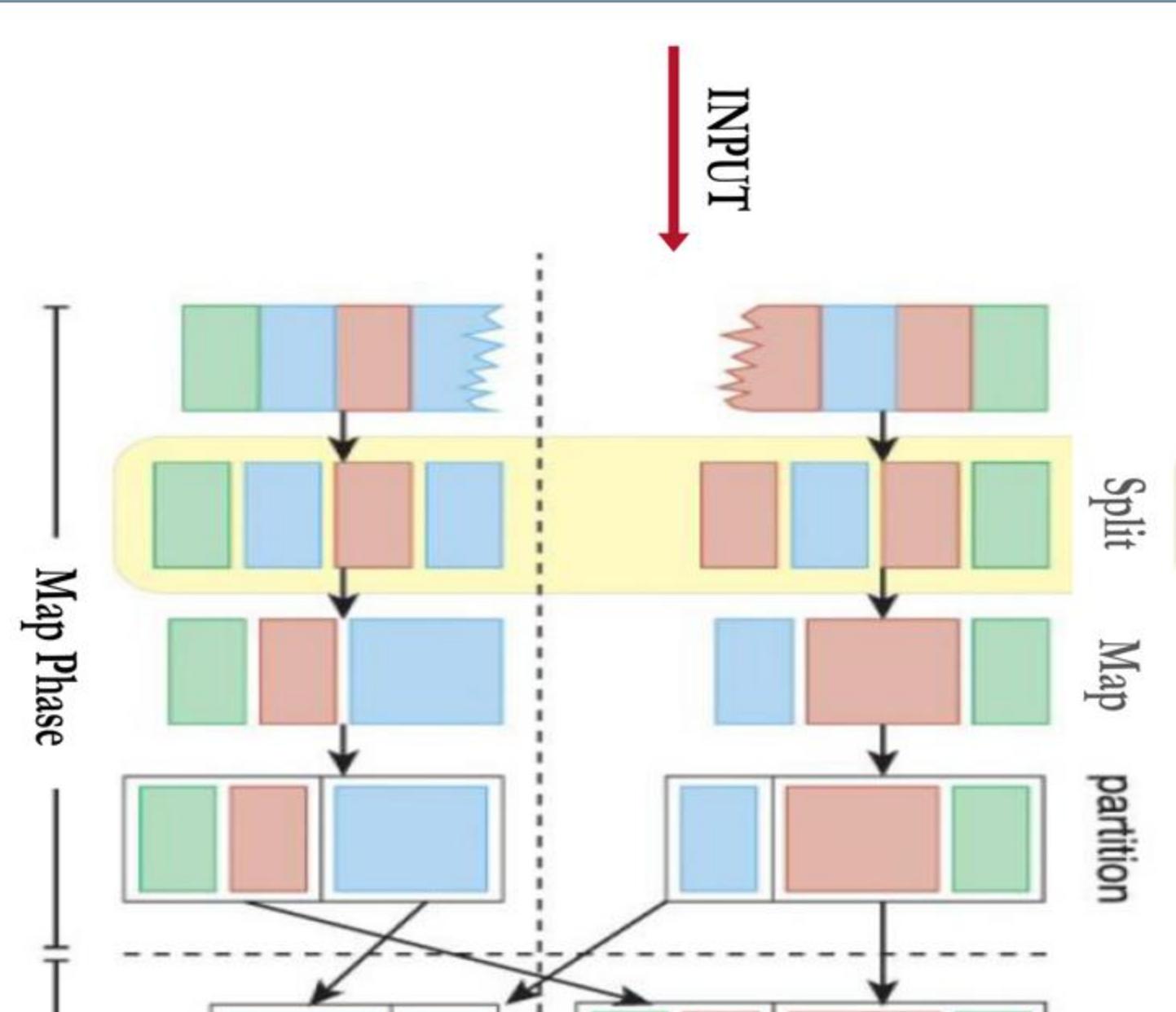


1 Sub-DataSet-1
2
3 Red Blue Red Blue Green Red Blue Green
4 White Black
5 Red White Black

By applying this stage on our example, we get the following:

1 Sub-DataSet-2
2
3 Orange Green
4 Red Blue Red
5 Blue Green Red Blue

1 Red Blue Red Blue Green Red Blue Green
2 White Black
3 Red White Black
4 Orange Green
5 Red Blue Red
6 Blue Green Red Blue
7 Green White Black

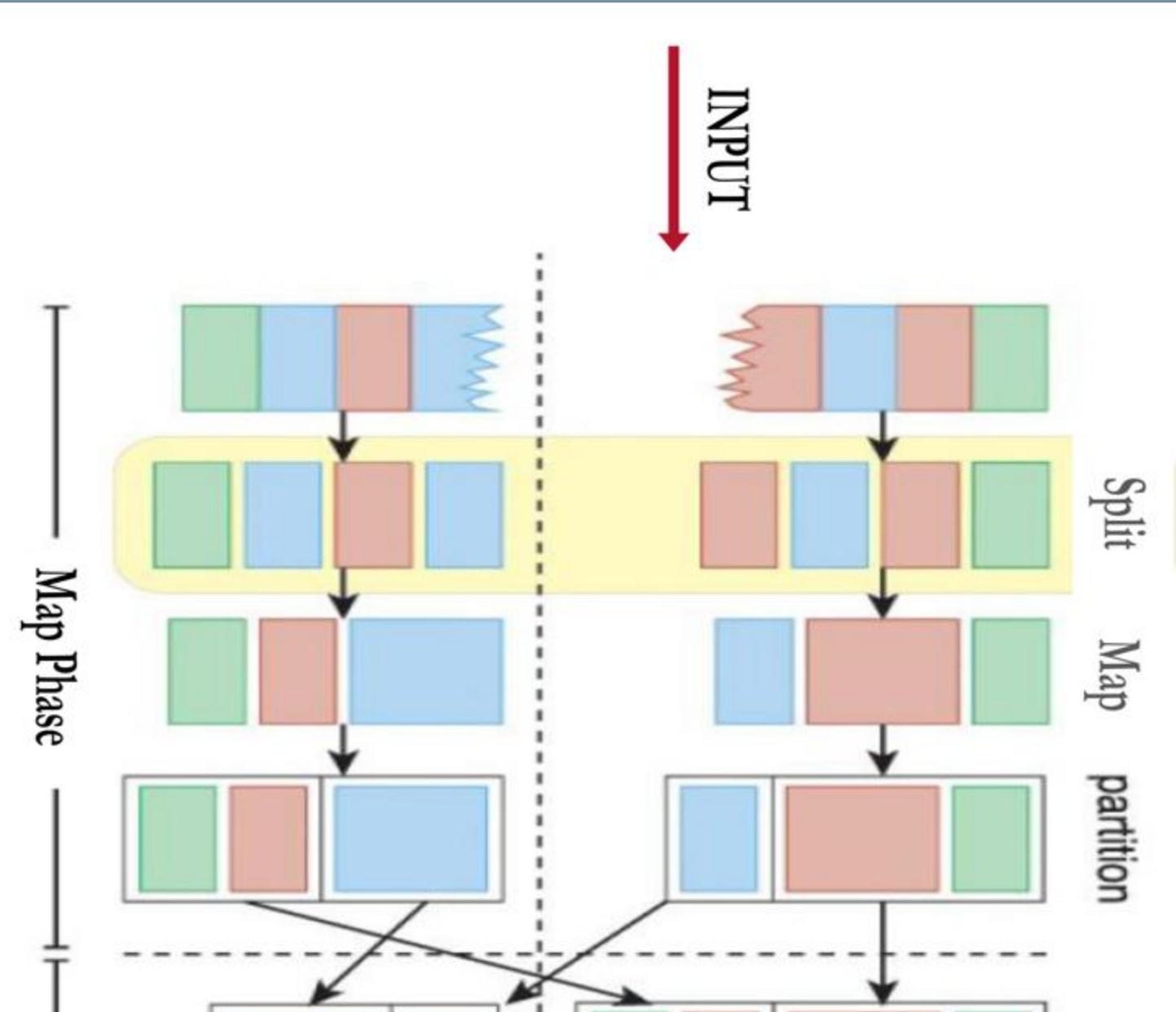


1 Sub-DataSet-1
2
3 Red Blue Red Blue Green Red Blue Green
4 White Black
5 Red White Black

By applying this stage on our example, we get the following:

1 Sub-DataSet-2
2
3 Orange Green
4 Red Blue Red
5 Blue Green Red Blue

1 Red Blue Red Blue Green Red Blue Green
2 White Black
3 Red White Black
4 Orange Green
5 Red Blue Red
6 Blue Green Red Blue
7 Green White Black

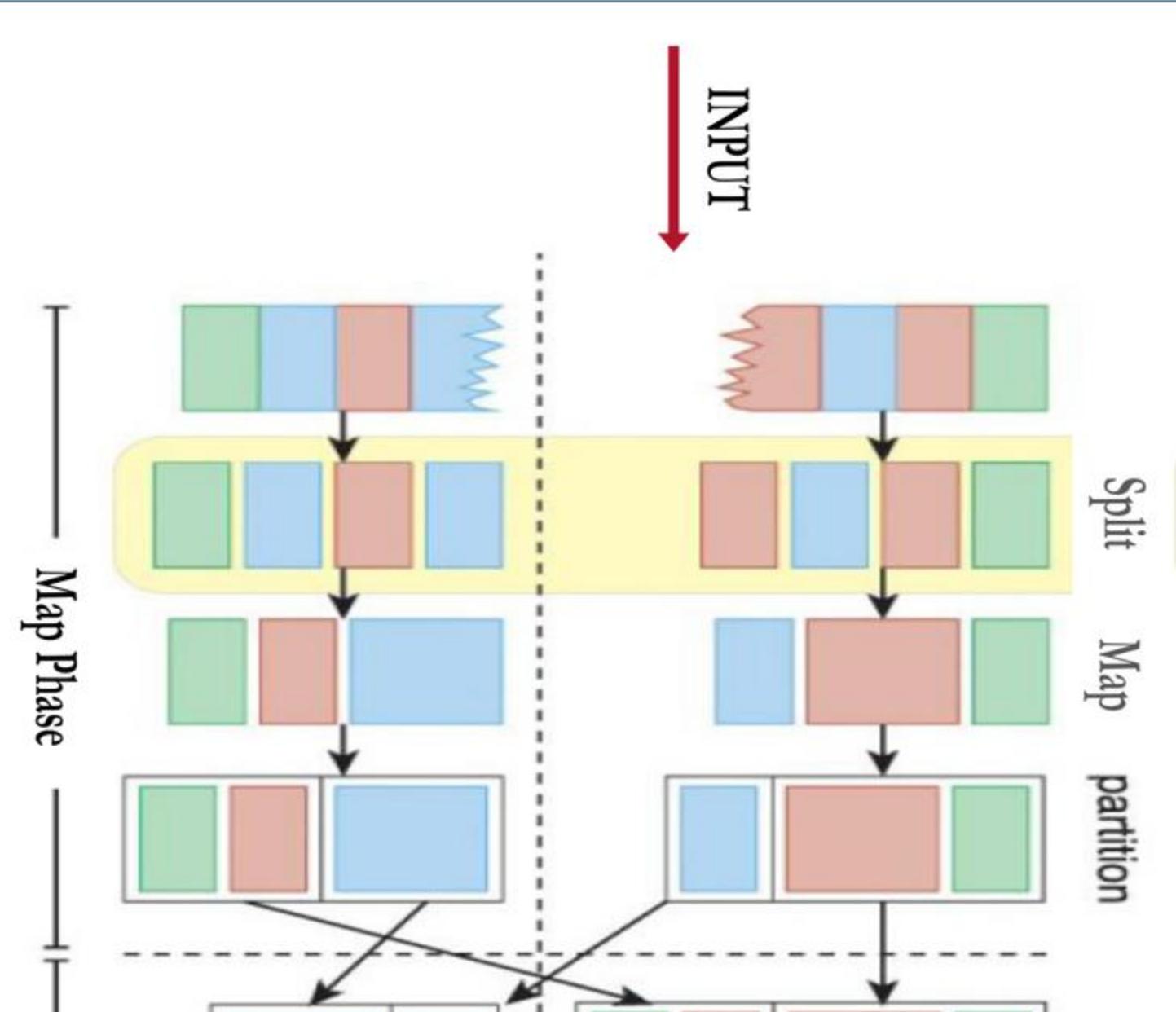


1 Sub-DataSet-1
2
3 Red Blue Red Blue Green Red Blue Green
4 White Black
5 Red White Black

By applying this stage on our example, we get the following:

1 Sub-DataSet-2
2
3 Orange Green
4 Red Blue Red
5 Blue Green Red Blue

1 Red Blue Red Blue Green Red Blue Green
2 White Black
3 Red White Black
4 Orange Green
5 Red Blue Red
6 Blue Green Red Blue
7 Green White Black

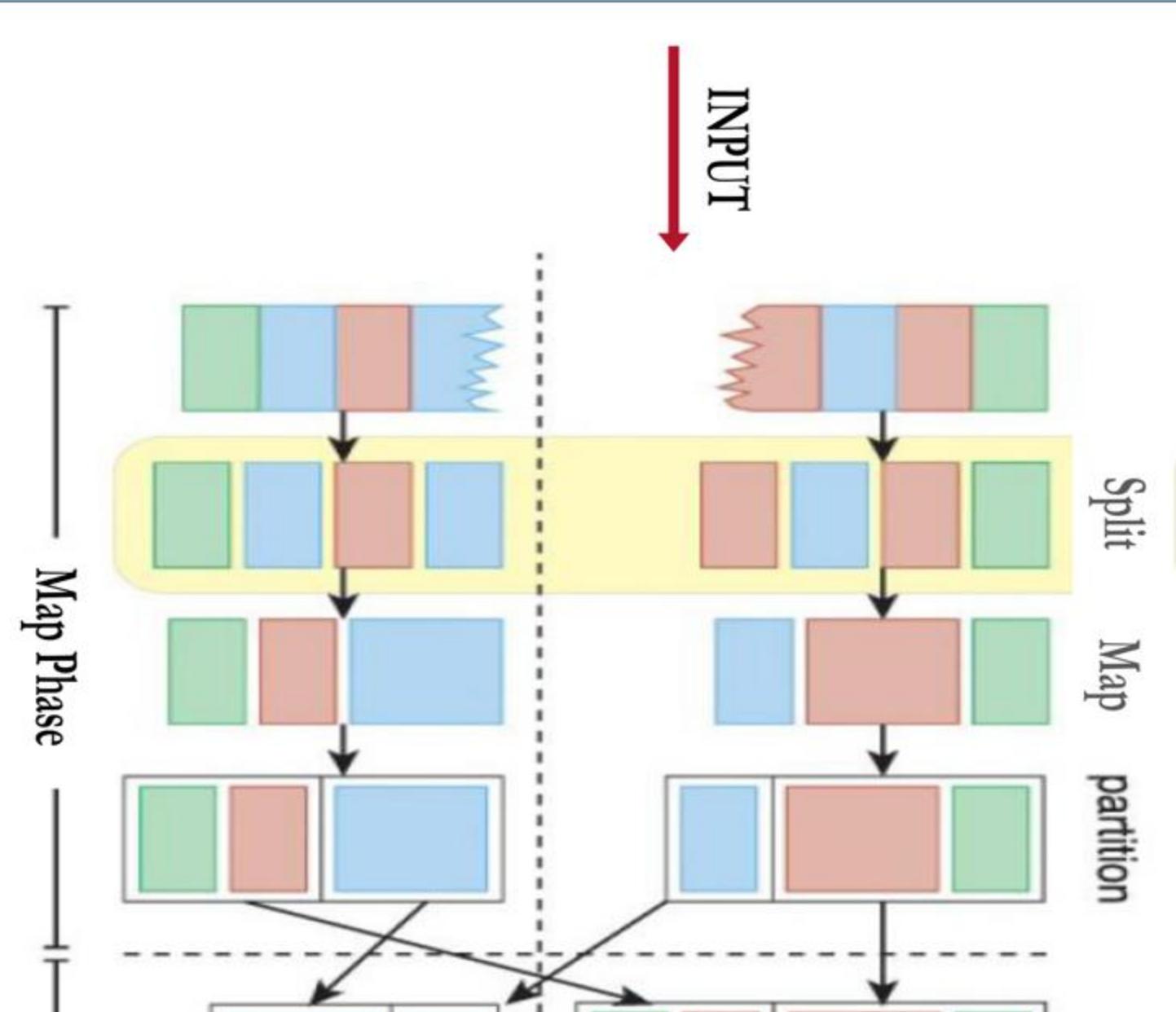


1 Sub-DataSet-1
2
3 Red Blue Red Blue Green Red Blue Green
4 White Black
5 Red White Black

By applying this stage on our example, we get the following:

1 Sub-DataSet-2
2
3 Orange Green
4 Red Blue Red
5 Blue Green Red Blue

1 Red Blue Red Blue Green Red Blue Green
2 White Black
3 Red White Black
4 Orange Green
5 Red Blue Red
6 Blue Green Red Blue
7 Green White Black



1 Sub-DataSet-1

<tbl_r cells="1" ix="4" maxcspan="1" maxrspan="1" usedcols

MapReduce Algorithm

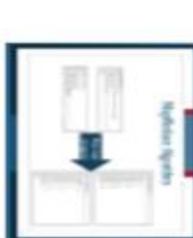
Map stage:

- This is the **map function** or **mapper** that executes **user-defined logic**.
- The mapper processes each key-value pair as per the user-defined logic and further generates a key-value pair as its output.
- The output key can either be the same as the input key or a substring value from the input value, or another user-defined object.

- Similarly, the output value can either be the same as the input value or a substring value from the input value, or another user-defined object.

- When all records of the split have been processed, the output is a list of key-value pairs where multiple key-value pairs can exist for the same key.

► By applying this stage on our example, we get the following:



13

MapReduce Algorithm

Shuffle and Sort stage:

- During the first stage of the reduce task, output from all partitioners is copied across the network to the nodes running the reduce task. This is known as **shuffling**.

- The output list of key-value pairs from each partitioner can contain the same key multiple times, so **sorting and merging** of the key-value pairs is done according to the keys so that the output contains a sorted list of all input keys and their values with the same keys appearing together.

- This merge creates a single key-value pair per group, where key is the group key and the value is the list of all group values.
- The way in which keys are sorted and merged can be *customized*.

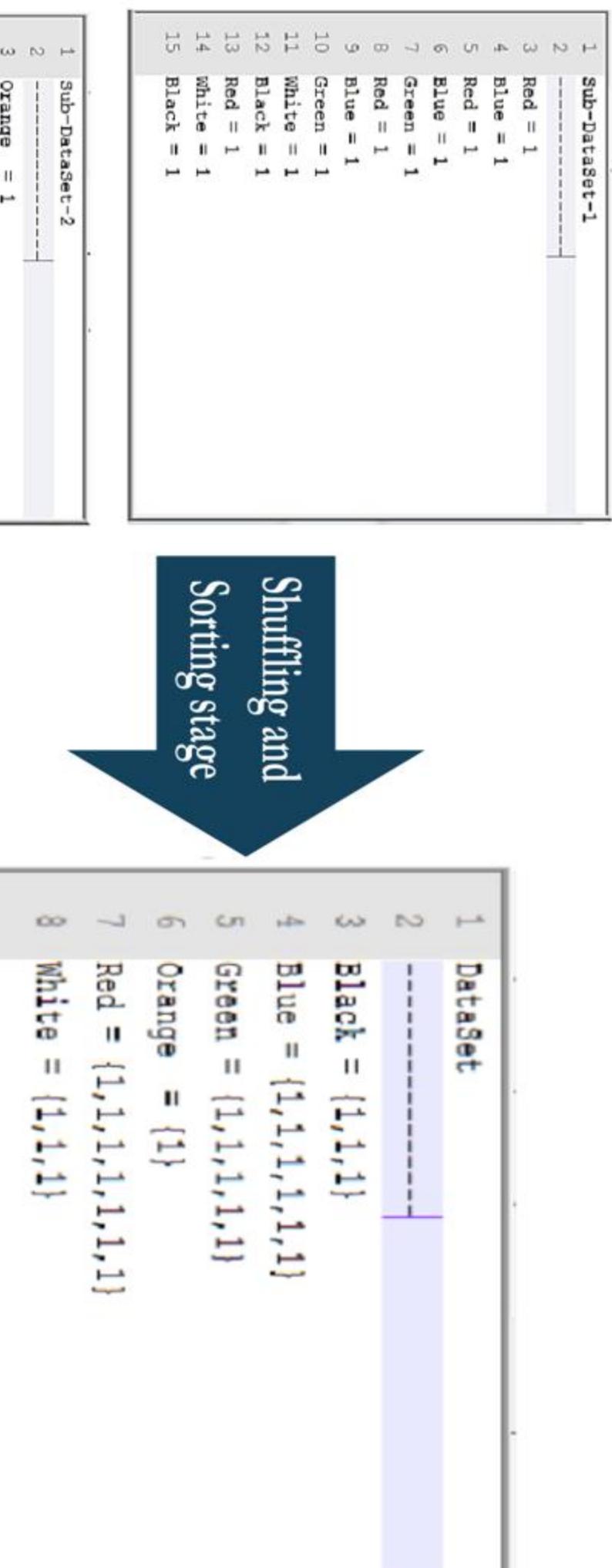
► By applying this stage on our example, we get the following:



14

MapReduce Algorithm

Shuffling and
Sorting stage



14

MapReduce Algorithm

Reduce stage:

- Reduce is the final stage of the reduce phase.
- Depending on the **user-defined logic** specified in the **reduce function** or **reducer**, the reducer will either further summarize its input or will emit the output without making any changes.

- The output key can either be the same as the input key or a substring value from the input value, or another user-defined object.
- The output value can either be the same as the input value or a substring value from the input value, or another user-defined object.

► By applying this stage on our example, we get the following:



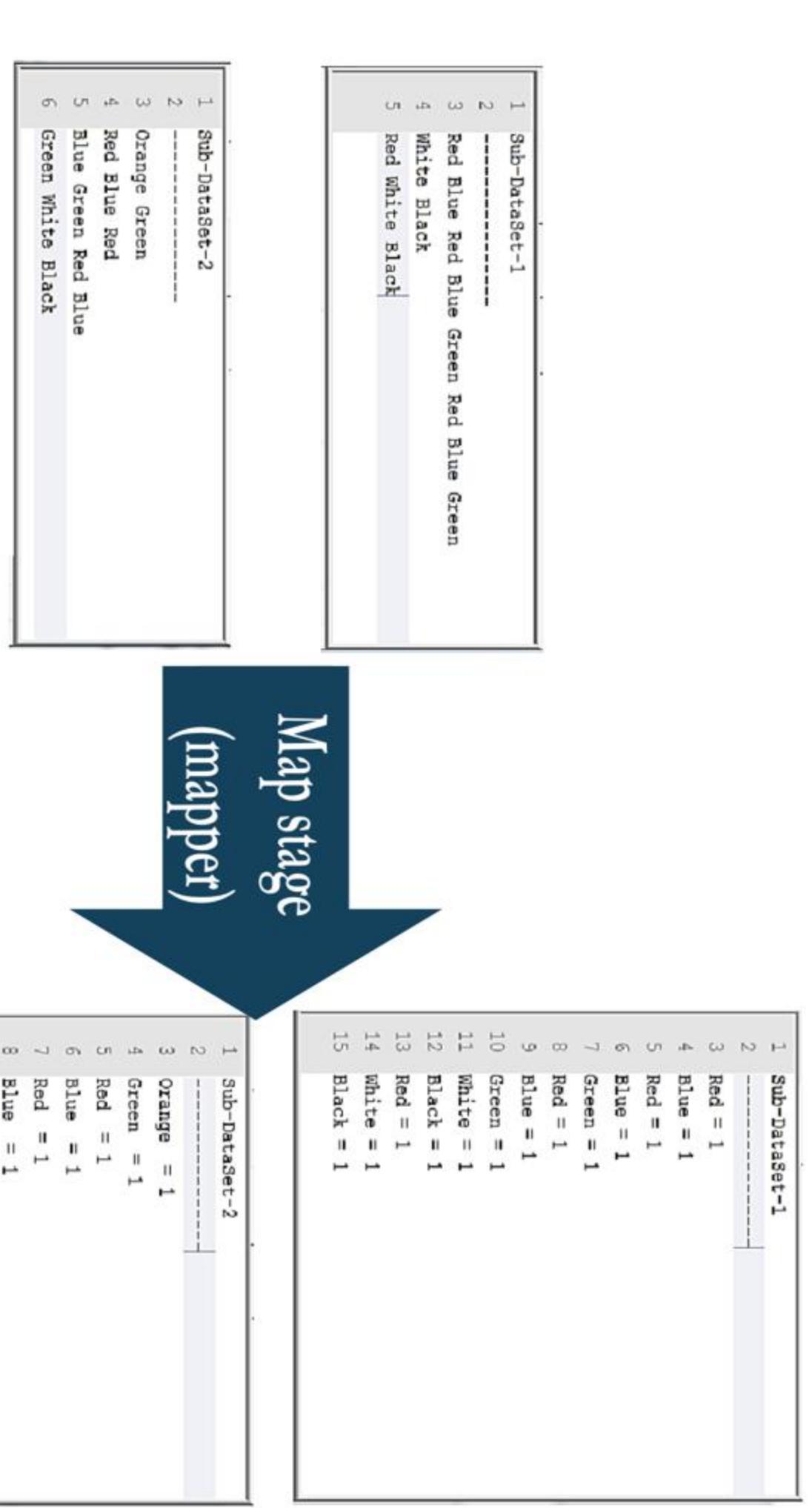
15

MapReduce Algorithm

Partition stage:

- During the partition stage, if more than one reducer is involved, a partitioner divides the output from the mapper into partitions between reducer instances.
- **All records for a particular key are assigned to the same reducer**.
- The MapReduce algorithm guarantees a random and fair distribution between reducers while making sure that all of the same keys across multiple mappers end up with the same reducer.

► Assume here in our example, that we have only one reducer.



16

MapReduce Algorithm

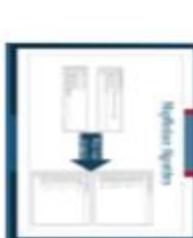
Map stage:

- This is the **map function** or **mapper** that executes **user-defined logic**.
- The mapper processes each key-value pair as per the user-defined logic and further generates a key-value pair as its output.
- The output key can either be the same as the input key or a substring value from the input value, or another user-defined object.

- Similarly, the output value can either be the same as the input value or a substring value from the input value, or another user-defined object.

- When all records of the split have been processed, the output is a list of key-value pairs where multiple key-value pairs can exist for the same key.

► By applying this stage on our example, we get the following:



13

MapReduce Algorithm

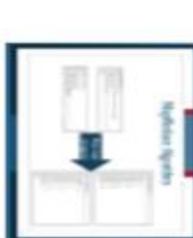
Map stage (mapper):

- This is the **map function** or **mapper** that executes **user-defined logic**.
- The mapper processes each key-value pair as per the user-defined logic and further generates a key-value pair as its output.
- The output key can either be the same as the input key or a substring value from the input value, or another user-defined object.

- Similarly, the output value can either be the same as the input value or a substring value from the input value, or another user-defined object.

- When all records of the split have been processed, the output is a list of key-value pairs where multiple key-value pairs can exist for the same key.

► By applying this stage on our example, we get the following:



14

MapReduce Algorithm

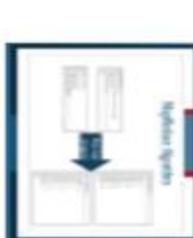
Map stage (mapper):

- This is the **map function** or **mapper** that executes **user-defined logic**.
- The mapper processes each key-value pair as per the user-defined logic and further generates a key-value pair as its output.
- The output key can either be the same as the input key or a substring value from the input value, or another user-defined object.

- Similarly, the output value can either be the same as the input value or a substring value from the input value, or another user-defined object.

- When all records of the split have been processed, the output is a list of key-value pairs where multiple key-value pairs can exist for the same key.

► By applying this stage on our example, we get the following:



15

MapReduce Algorithm

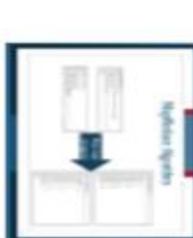
Map stage (mapper):

- This is the **map function** or **mapper** that executes **user-defined logic**.
- The mapper processes each key-value pair as per the user-defined logic and further generates a key-value pair as its output.
- The output key can either be the same as the input key or a substring value from the input value, or another user-defined object.

- Similarly, the output value can either be the same as the input value or a substring value from the input value, or another user-defined object.

- When all records of the split have been processed, the output is a list of key-value pairs where multiple key-value pairs can exist for the same key.

► By applying this stage on our example, we get the following:



16

MapReduce Algorithm

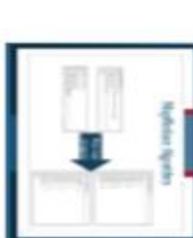
Map stage (mapper):

- This is the **map function** or **mapper** that executes **user-defined logic**.
- The mapper processes each key-value pair as per the user-defined logic and further generates a key-value pair as its output.
- The output key can either be the same as the input key or a substring value from the input value, or another user-defined object.

- Similarly, the output value can either be the same as the input value or a substring value from the input value, or another user-defined object.

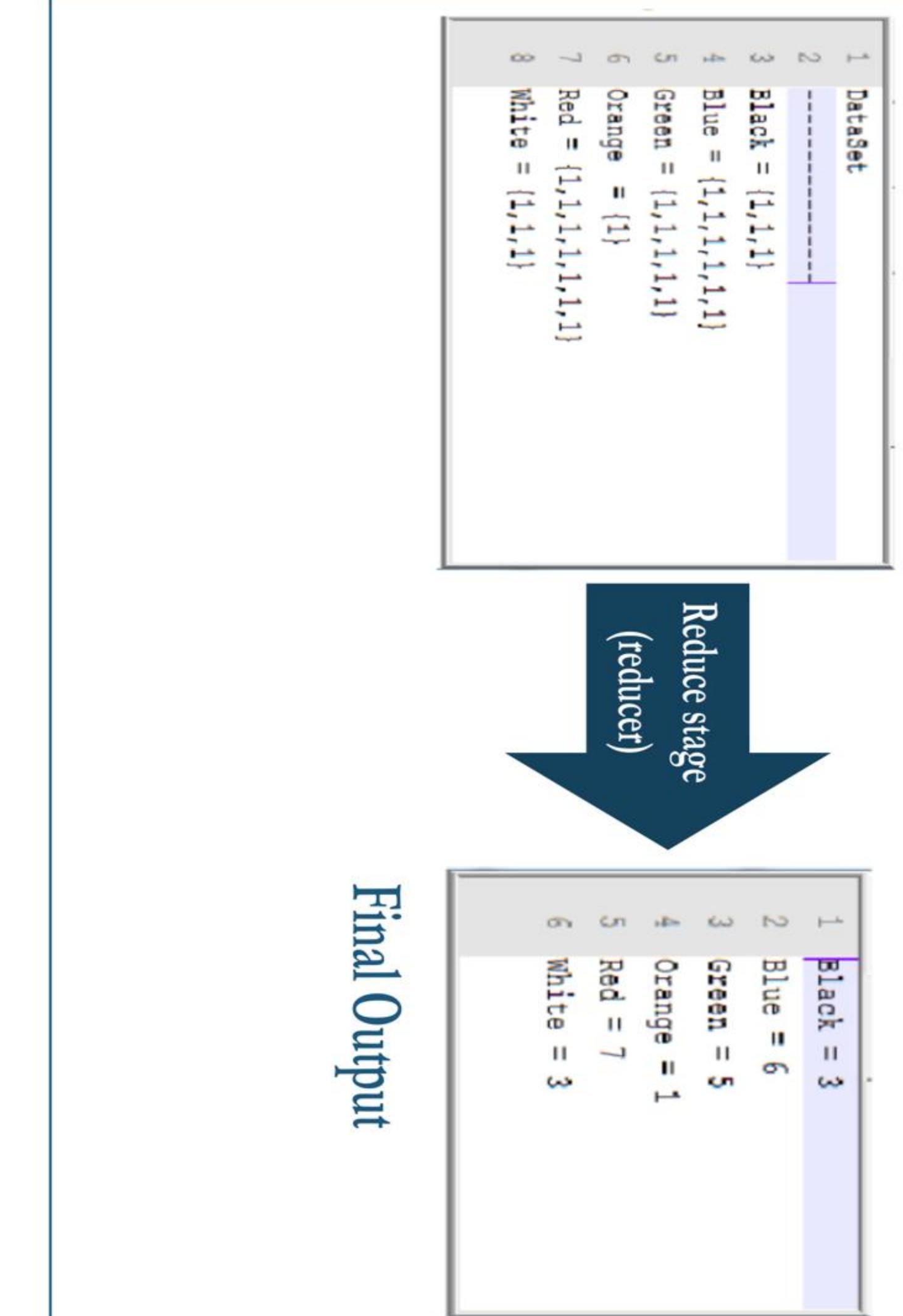
- When all records of the split have been processed, the output is a list of key-value pairs where multiple key-value pairs can exist for the same key.

► By applying this stage on our example, we get the following:



17

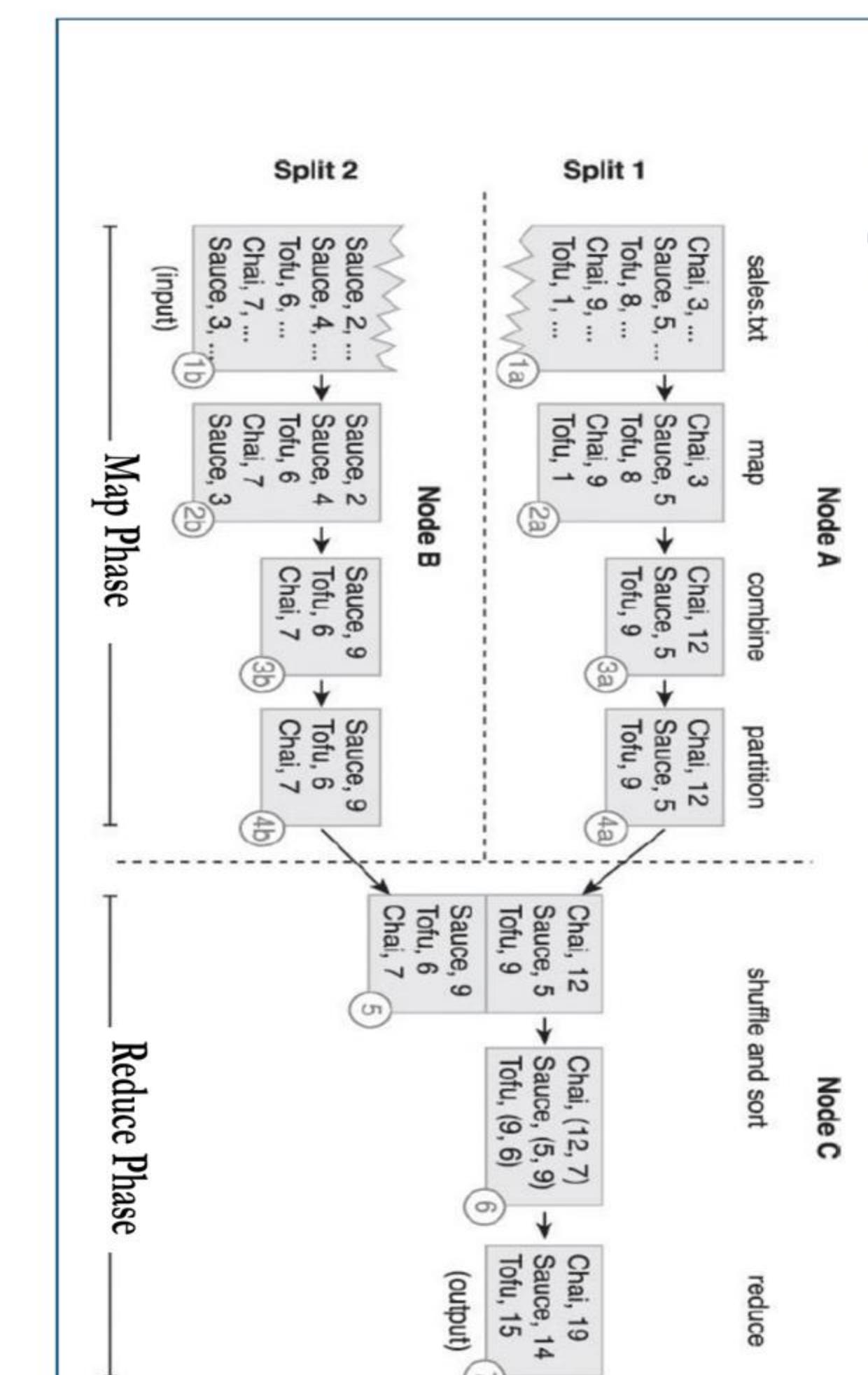
MapReduce Algorithm



MapReduce Algorithm

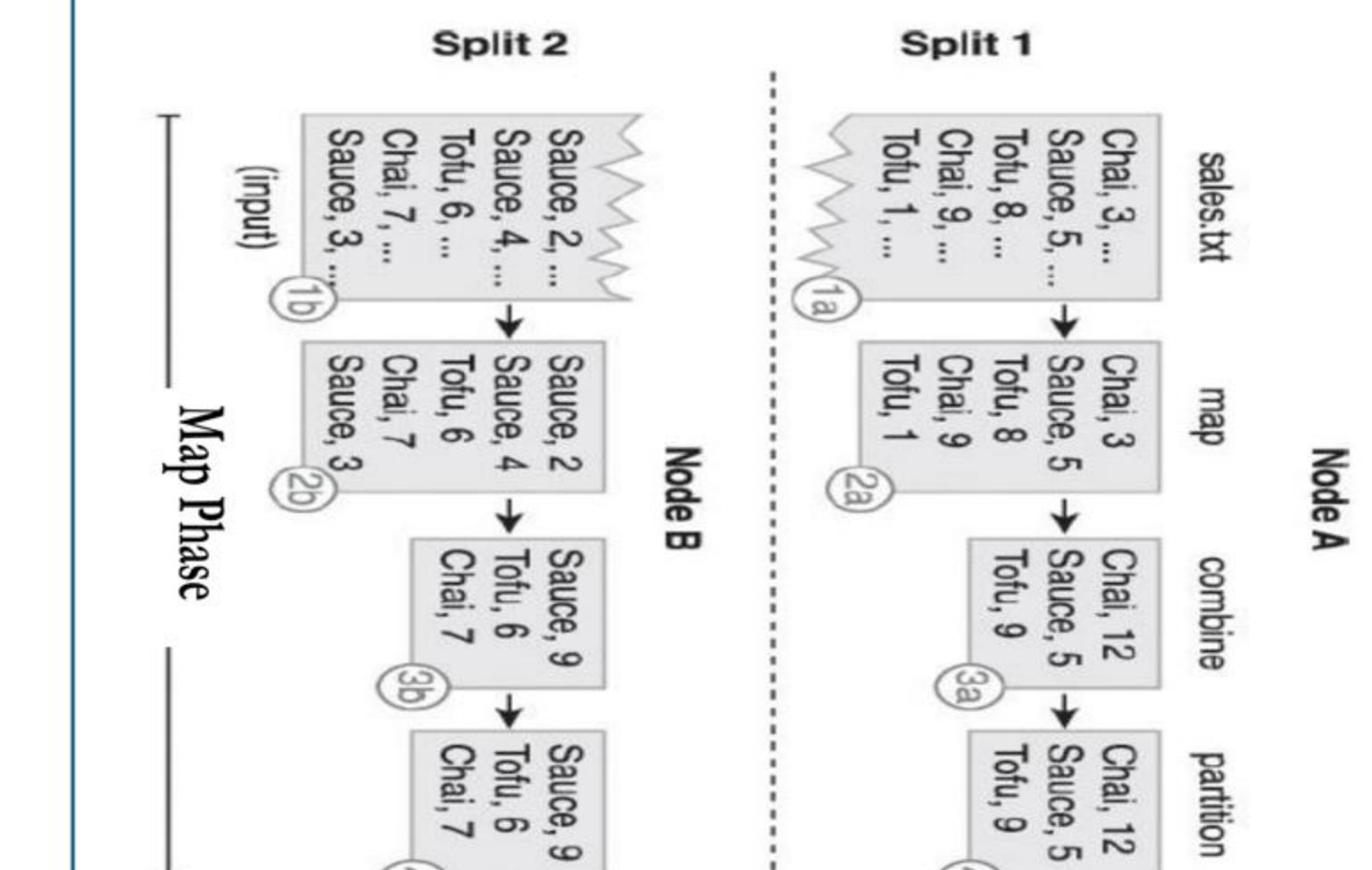
- Consider another example as follows:

- We have products information as input and we need as output the quantity of each product.



MapReduce Algorithm

1. The input (`sales.txt`) is divided into two splits.
 2. Two map tasks running on two different nodes, Node A and Node B, extract product and quantity from the respective split's records in parallel. The output from each map function is a key-value pair where product is the key while quantity is the value.
 3. The combiner then performs local summation of product quantities. (A combiner is essentially a reducer function that locally groups a mapper's output on the same node as the mapper.)
 4. As there is only one reduce task, no partitioning is performed.



MapReduce Algorithm

5. The output from the two map

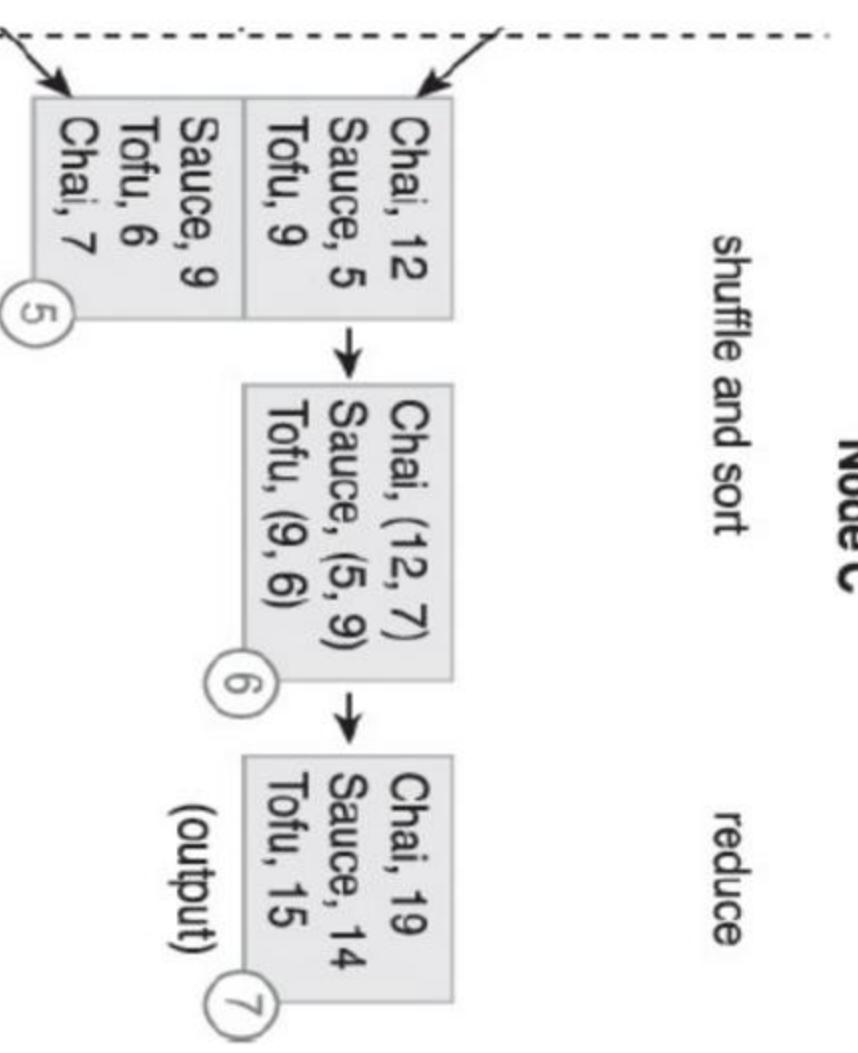
Node

The output from the two map tasks is then copied to a third node, Node C, that runs the shuffle stage as part of the reduce task.

6. The sort stage then groups all quantities of the same product together as a list.

7. The reduce function then sums up the quantities of each unique product in order to create the output.

MapReduce Examples



MapReduce Examples

- For the examples in this section, we will use data similar to the data collected by a web analytics service that shows various statistics for page visits for a website.

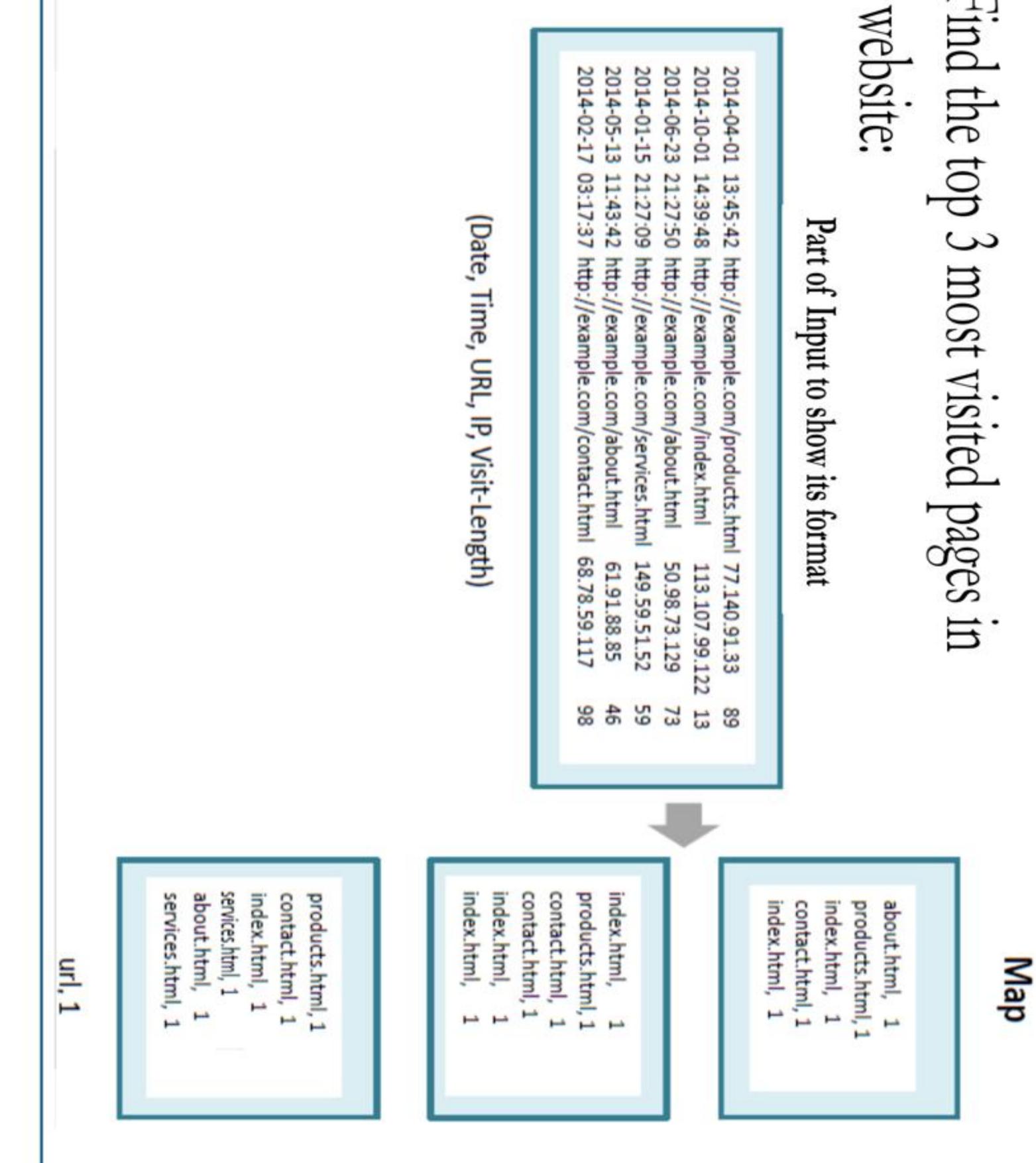
- Each page has some tracking code which sends the visitor's IP address along with a timestamp to the web analytics service. The web analytics service keeps a record of all page visits and the visitor IP addresses and uses MapReduce programs for computing various statistics.

- Each visit to a page is logged as one row in the log. The log file contains the following columns:

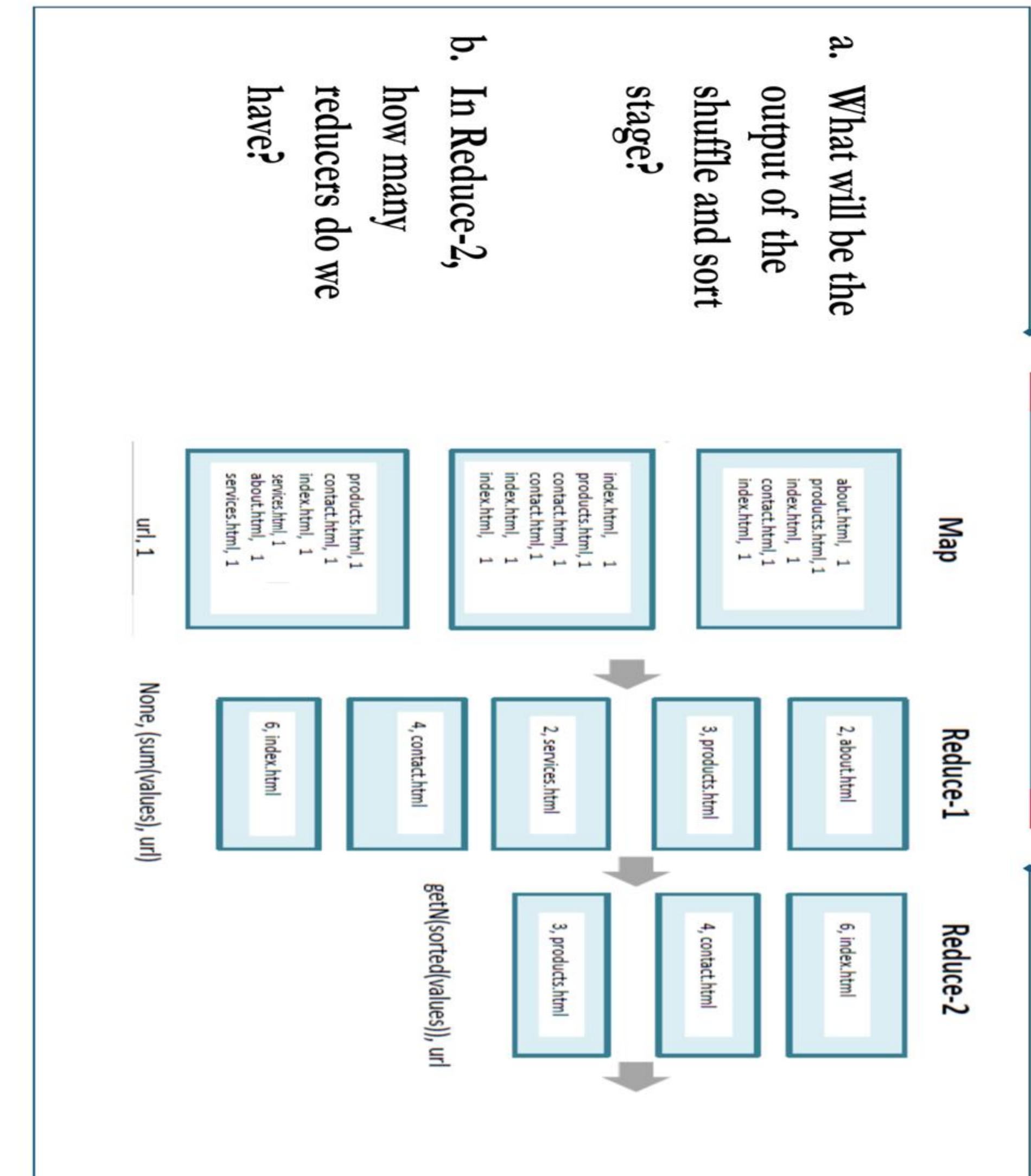
Date (YYYY-MM-DD), Time (HH:MM:SS), URL, IP, Visit-Length.

MapReduce Examples

3. Top-N: Find the top 3 most visited pages in the given website:



a. What will be the output of the shuffle and sort stage?



Top-N computation Explanation:

MapReduce Examples

4. Filtering:

- Filter out a subset of the records based on a filtering criteria.
 - For example: filtering all page visits for the page ‘contact.html’ in the month of Dec 2014.

(Date, Time, URL, IP, Visit-Length)

List of Input to show its initial					
2014-04-01 13:45:42	http://example.com/products.html	77.140.91.33	89		
2014-10-01 14:39:40	http://example.com/index.html	113.107.99.122	13		
2014-06-23 21:27:50	http://example.com/about.html	50.98.73.129	73		
2014-01-15 21:27:09	http://example.com/services.html	149.59.51.52	59		
2014-05-13 11:43:42	http://example.com/about.html	61.91.88.85	46		
2014-02-17 08:17:37	http://example.com/contact.html	68.78.59.117	98		

MapReduce Examples

1

1

四

•

1

1

2

- Filtering is useful when you want to get a subset of data for processing

MapReduce Examples

MapReduce Examples

- Each mapper filters out its local records based on the filtering criteria in the map function.
- The mapper function in this example parses each line of the input, extracts the month, year and page URL and emits key-value pairs if the month and year are Dec 2014 and the page URL is '<http://example.com/contact.html>'.
- The key is the URL, and the value is a tuple containing the rest of the parsed fields.

21

URL, Date, Time, IP, Visit-Length

27

3