# Agenda
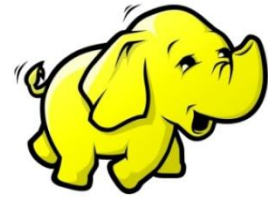
- ✓ Revision
- ✓ Where Map-Reduce is used?
- ✓ Weather Data( Where MapReduce is Used
- ✓ Usecases in Healthcare & Architecture
- ✓ Traditional Way
- ✓ Map Reduce Process
- ✓ Anatomy of a MapReduce program
- ✓ MapReduce Way
- ✓ Revisit de-identification architecture
- ✓ Why MapReduce
- ✓ Input Splits
- ✓ Relation between Input Splits and HDFS blocks
- ✓ MapReduce Flow
- ✓ Overview of MapReduce
- ✓ Combiner & Partitioner
- ✓ Input Formats

**Conf/ files**

- ✓ Core-site.xml
- ✓ Mapred-site.xml
- ✓ Hdfs-site.xml
- ✓ Slaves
- ✓ Masters

# Where MapReduce Is Used?

✓ Weather Forecasting

✓ Health Care

✓ Oil & Gas Industry

# Weather Data

## edureka!

ftp://ftp.ncdc.noaa.gov/pub/data/uscrn/products/daily01/

✓ **Problem Statement:**
  ✓ De-identify personal health information.


✓ **Challenges:**
  ✓ Huge amount of data flows into the systems daily and there are multiple data sources that we need to aggregate data from.
  ✓ Crunching this huge data and deidentifying it in a traditional way had problems.

Taking DB dump in CSV format and ingest into HDFS

*Sqoop*

matches

Store Deidentified CSV file into HDFS

HDFS

Read CSV file from HDFS

Map Task 1
Map Task 2
.
.

Deidentify columns based on configurations

Reduce Task 1
Reduce Task 2
.
.

**edureka!**



Very Big Data → Split Data → grep → matches

Split Data → grep → matches

Split Data → grep → matches

⋮

Split Data → grep → matches

→ cat → All matches

edureka!

Key

Value

MapReduce

Map:

(K1, V1)

List (K2, V2)

Reduce:

(K2, list (V2))

List (K3, V3)

Taking DB dump in CSV format
and ingest into HDFS

*Sqoop*

matches

Store Deidentified CSV
file into HDFS

HDFS

Read CSV
file from
HDFS

Map Task 1

Map Task 2

.

.

Deidentify
columns based on
configurations

Reduce Task 1

Reduce Task 2

.

.

✓ **Two biggest Advantages:**

 ✓ Taking processing to the data
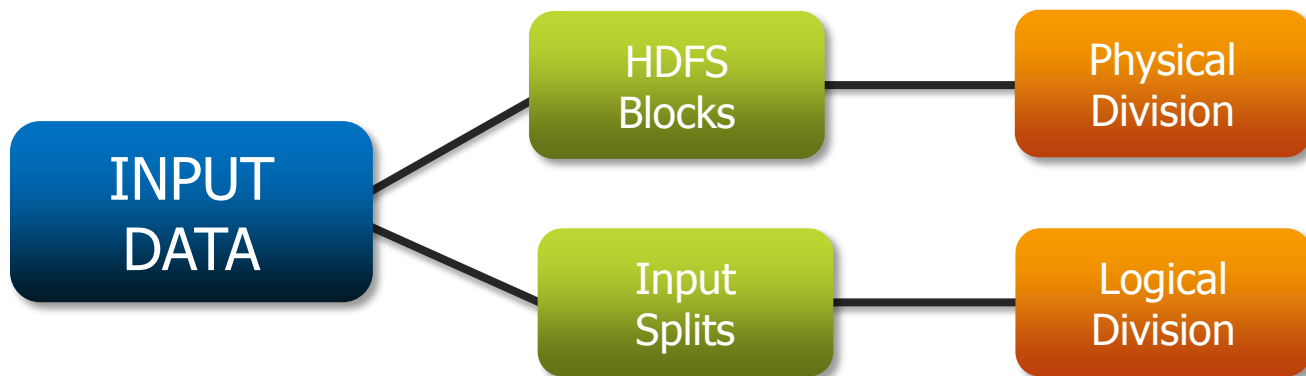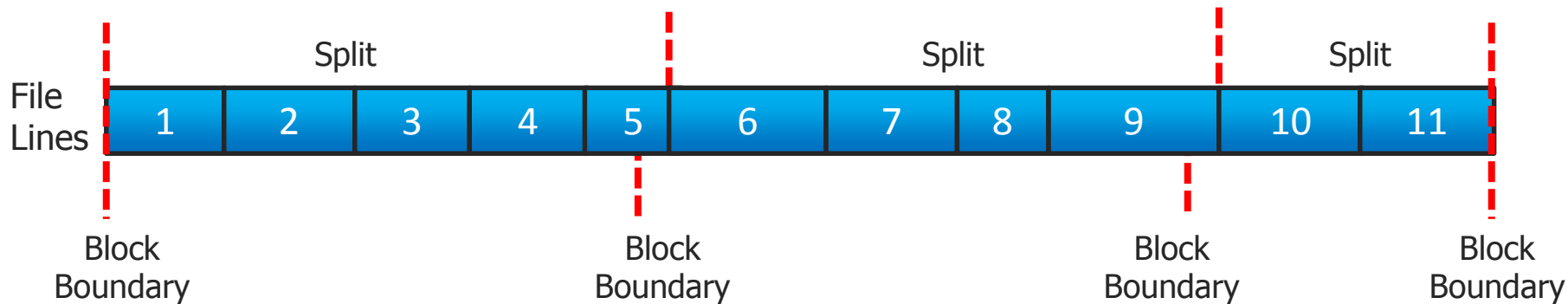
 ✓ Processing data in parallel

# Relation between Input Splits and HDFS Blocks

✓ Logical records do not fit neatly into the HDFS blocks.

✓ Logical records are lines that cross the boundary of the blocks.

✓ First split contains line 5 although it spans across blocks.

**INPUT DATA**

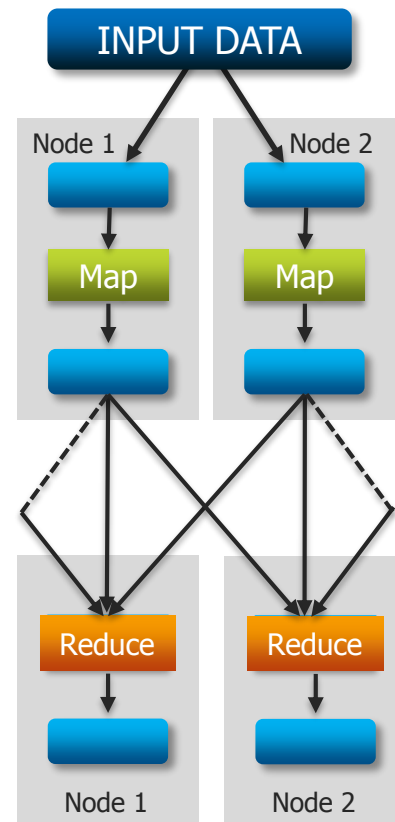Input data is distributed to nodes

Each map task works on a "split" of data
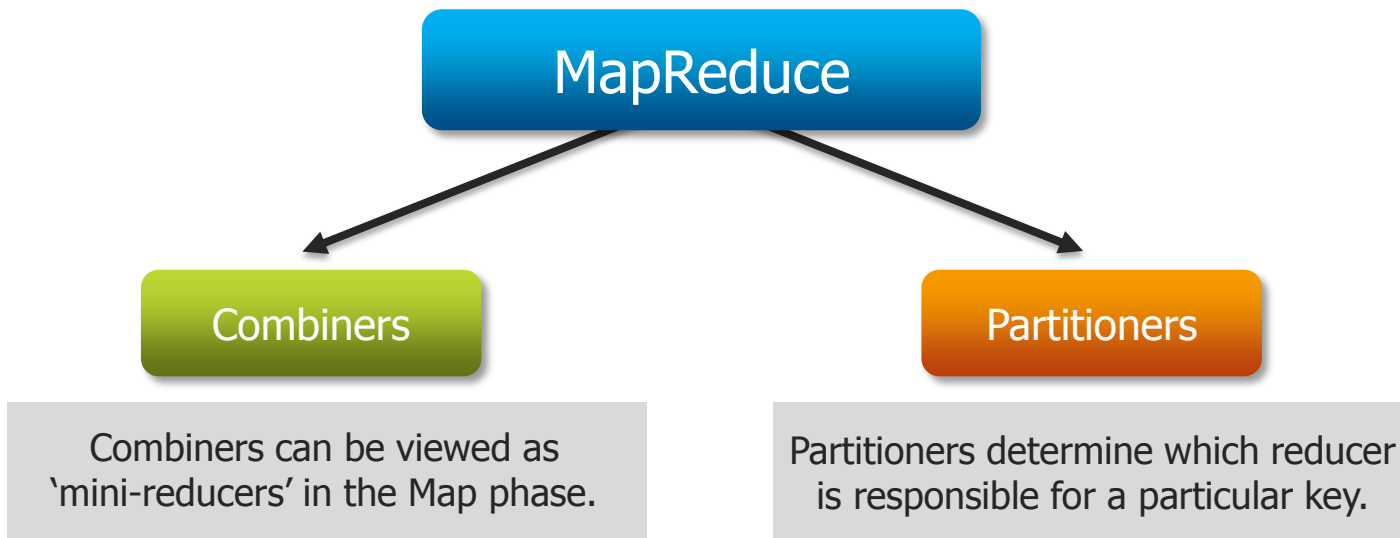
Mapper outputs intermediate data

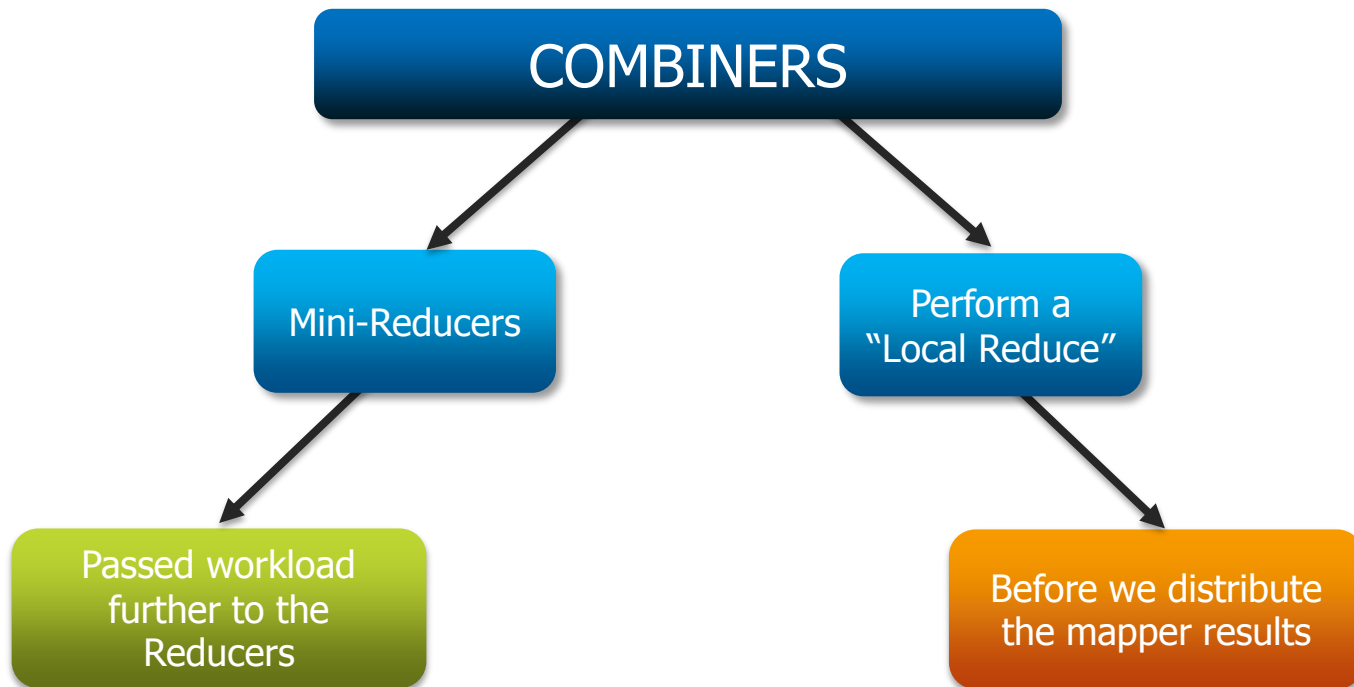Data exchange between nodes in a "shuffle" process

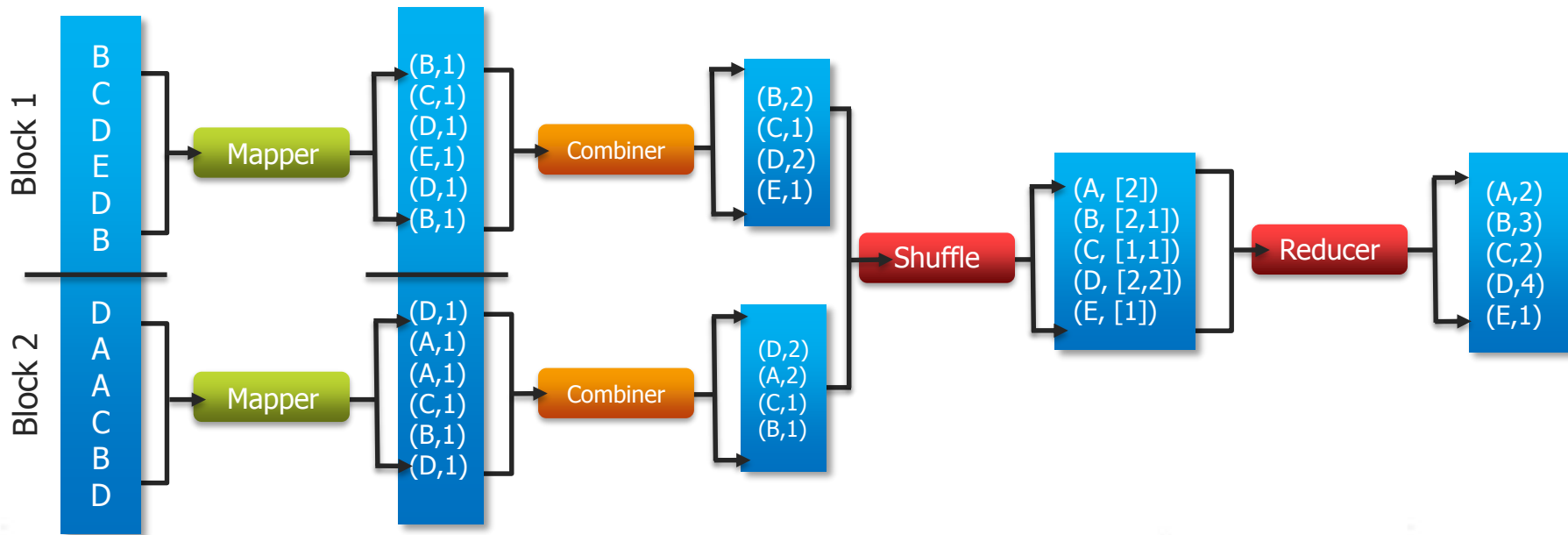Intermediate data of the same key goes to the same reducer

Reducer output is stored

Node 1    Node 2

Map    Map

Reduce    Reduce

Node 1    Node 2

**edureka!**

Complete view of MapReduce, illustrating combiners and partitioners
in  addition to Mappers and Reducers

MapReduce

Combiners

Partitioners

Combiners can be viewed as
'mini-reducers' in the Map phase.

Partitioners determine which reducer
is responsible for a particular key.

Each line in the text files is a **record.**

| KeyValueTextInputFormat | → | Text Files |

| SequenceFileInputFormat | → | Sequence Files |

# Assignment

Attempt the following assignment using the items present in the LMS under the tab Week 3:

- ✓ Watch video "Running MapReduce Program".

- ✓ Execute Partitioner Code.

- ✓ Attempt Assignment Week 3 – Wordcount , Patents, Temperature ,Temperature2 & Alphabets

# edureka!
# Thank You
## See You in Class Next Week