



Walchand Linux Users' Group



Presents Club Service On



Topics Covered :

- Big data
- Hadoop intro
- Architecture
- Working
- Advantages

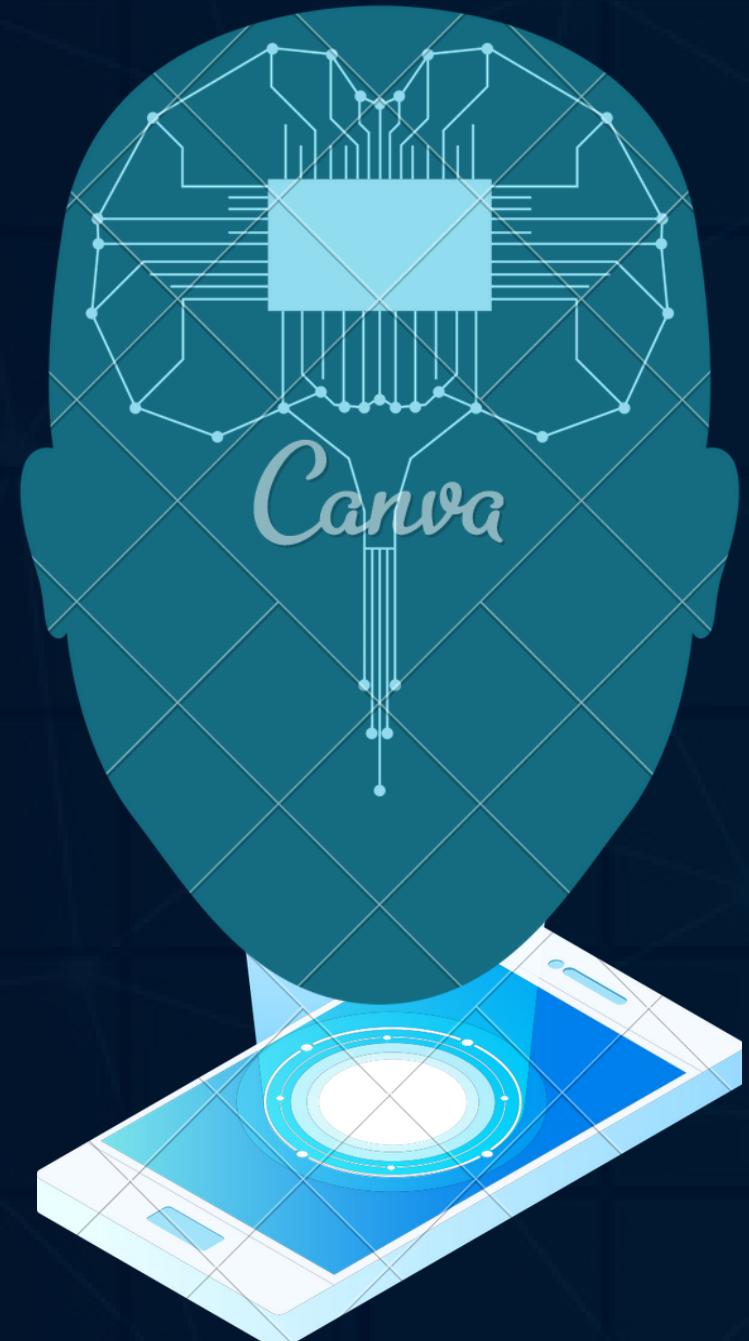


CONNECT WITH US



www.wcewlug.org

15th Sept 6:00 PM
N2 Classroom



HADOOP TUTORIAL



TABLE OF CONTENT

01 Big Data

The world is one big data problem.

02 Hadoop intro

Open Source

03 Architecture

Hadoop Architecture

04 Working

Hadoop Working.

Update package:
sudo apt get update

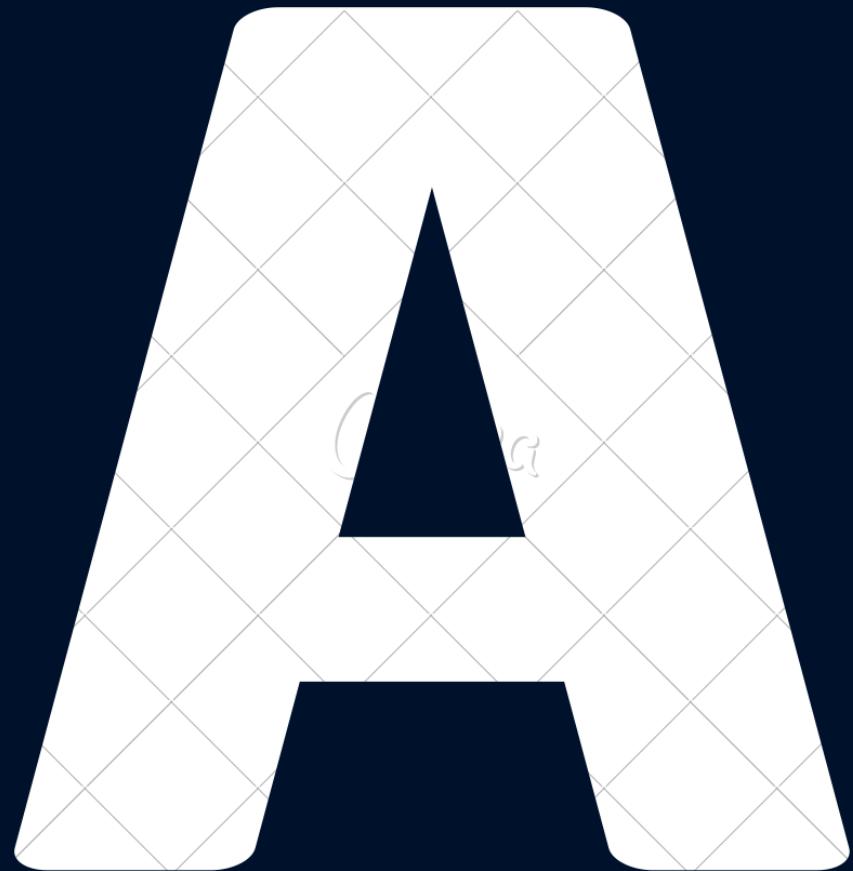
Accept required certificates:
sudo apt-get install apt-transport-https ca-certificates

Install docker
sudo apt install docker.io

Pull the hadoop image:
docker pull sequenceiq/hadoop-docker:2.7.0

BIG DATA

- Since the Dawn of time to 2005 - Humans have created 130 exabytes of data.



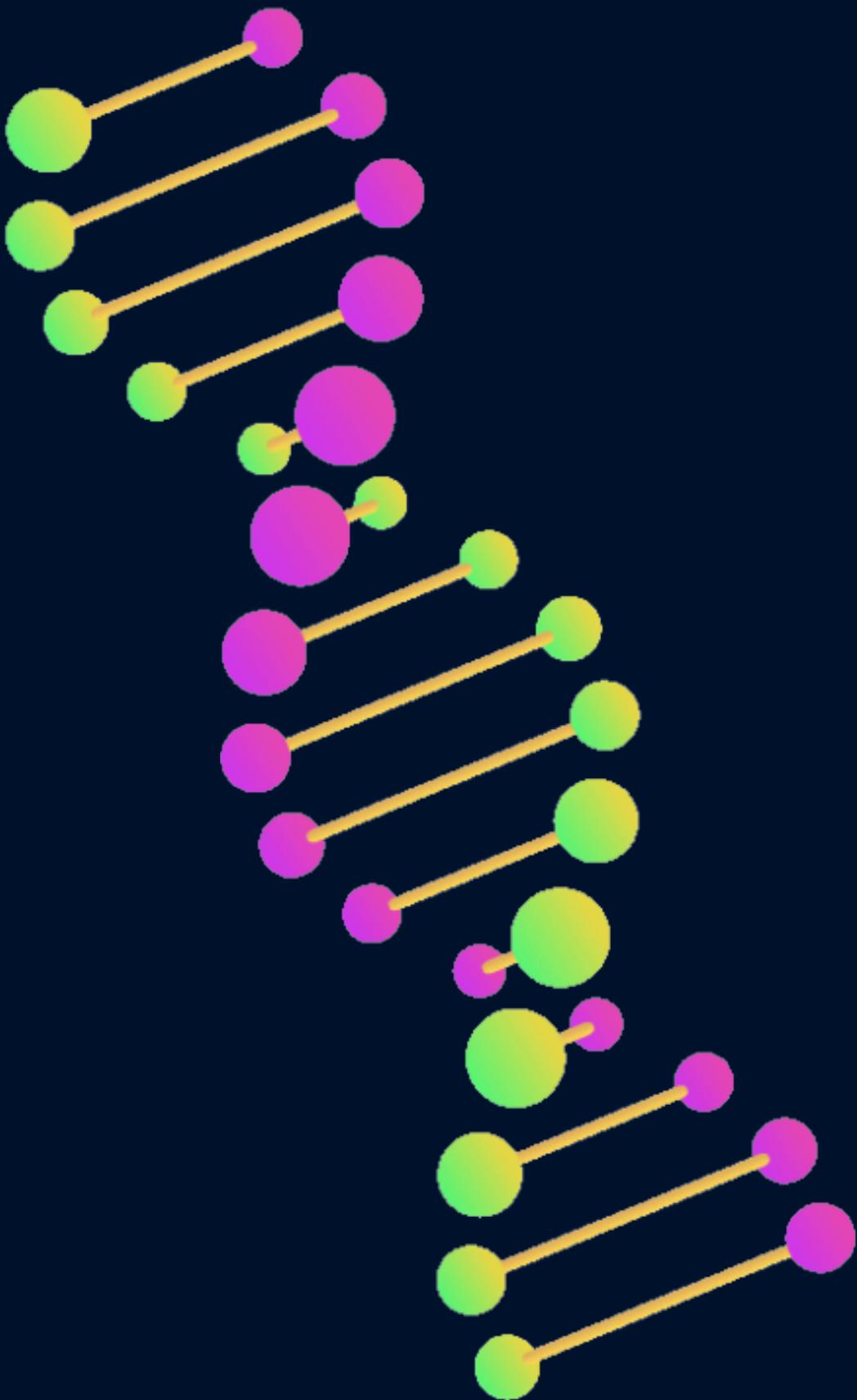
1 Byte



1 kiloBytes = 1000 Bytes



1 MegaBytes ~ 1000 KiloBytes



1 GigaBytes ~ 1000 MegaBytes

Human GNOME can
be Encoded in Digital
memory of around
1 GB size.



1 TeraBytes ~ 1000 GigaBytes

1 PetaBytes ~ 1000 TeraBytes



390 billion trees

7.9 billion people on earth



1 ExaBytes ~ 1000 PetaBytes

BIG DATA

- Since the Dawn of time to 2005, Humans have created 130 exabytes of data.

BIG DATA

- Since the Dawn of time, Humans have created 130 exabytes of data.
- Till 2010 - 1,200 ExaBytes

BIG DATA

- Since the Dawn of time, Humans have created 130 exabytes of data.
- Till 2010 - 1,200 ExaBytes
- Till 2015 - 7,900 ExaBytes

BIG DATA

- Since the Dawn of time, Humans have created 130 exabytes of data.
- Till 2010 - 1,200 ExaBytes
- Till 2015 - 7,900 ExaBytes
- Till 2020 - 49,900 ExaBytes

SOLUTIONS ?





TRADITIONAL SOLUTIONS

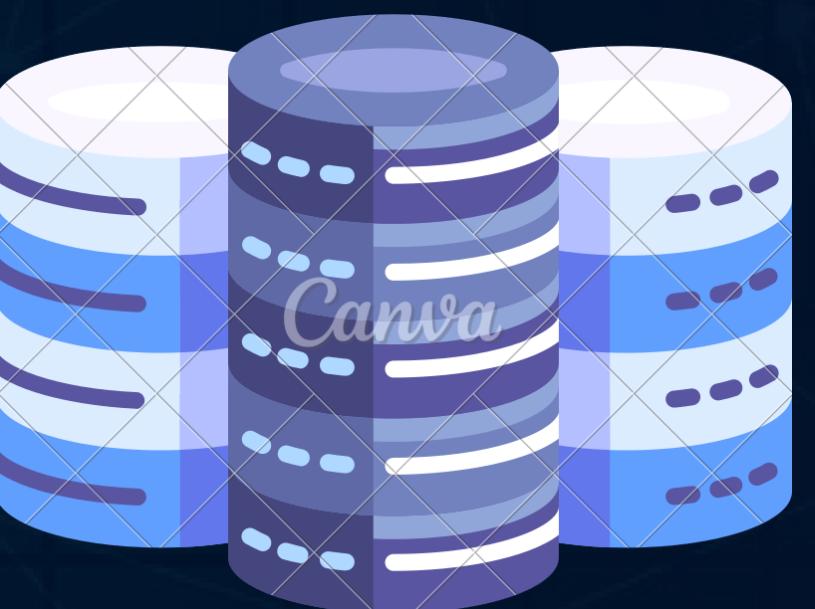
IBM AND
ORACLE
SQL
MongoDB



User

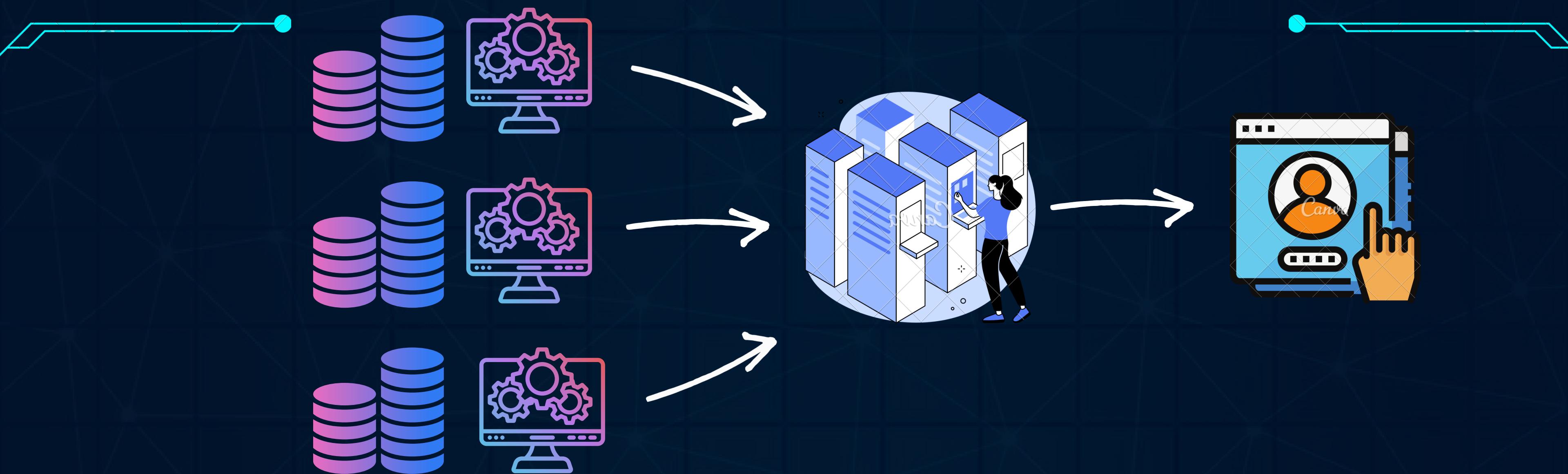


Centralized
system



DataBases

GOOGLE SOLUTIONS





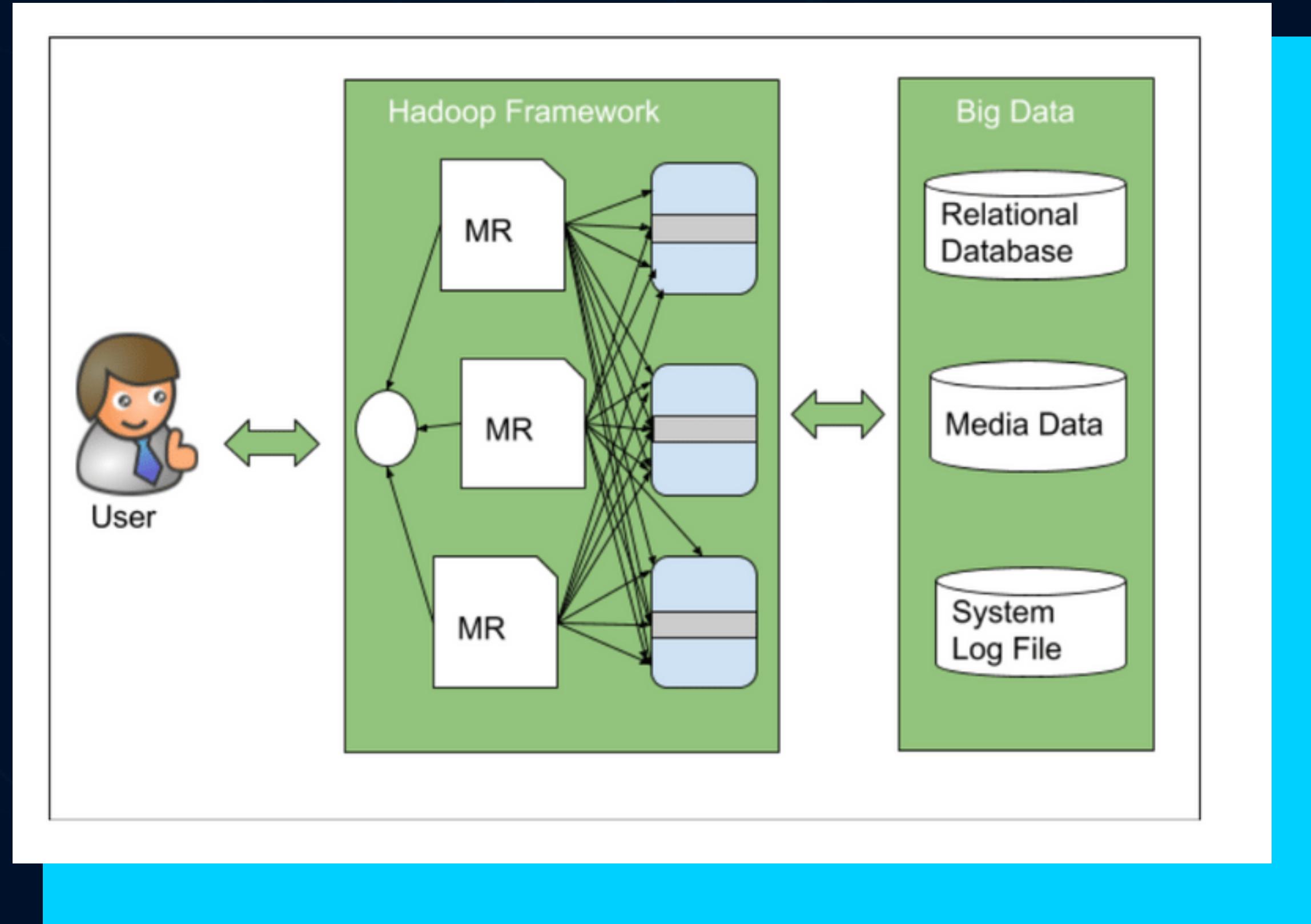
HADOOP

Using the solution provided by Google, Doug Cutting and his team developed an Open Source Project called HADOOP.

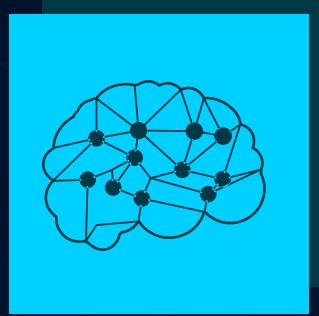
INTRODUCTION

Hadoop is an open-source framework that allows to store and process big data in a distributed environment across clusters of computers using simple programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage.

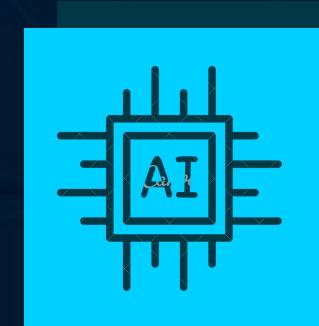
HADOOP



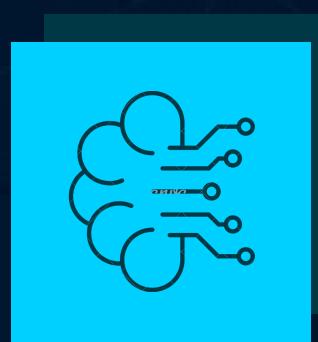
HADOOP



APACHE OPEN
SOURCE
FRAMEWORK



WRITTEN
IN
JAVA



MAP
REDUCE
ALGORITHM

HADOOP ARCHITECTURE

HADOOP ARCHITECTURE

HDFS****
Hadoop
Distributed
File System

**MAP
REDUCE**
Processing/
Computation
layer



HADOOP DISTRIBUTED FILE SYSTEM

HDFS

- Based on Google File System.
- Some similarities with Existing Distributed File System
- Highly fault tolerant
- Can be deployed on low cost hardware
- High throughput access to application data
- Suitable for application having large datasets
- Provides a command line interface to interact with HDFS

HDFS CONCEPTS

Name Node: Stores Meta Data

Meta Data:
`/data/pristine/catalina.log.>1, 2, 4`
`/data/pristine/myfile. >3,5`

Data Node 1

1 2 4
5

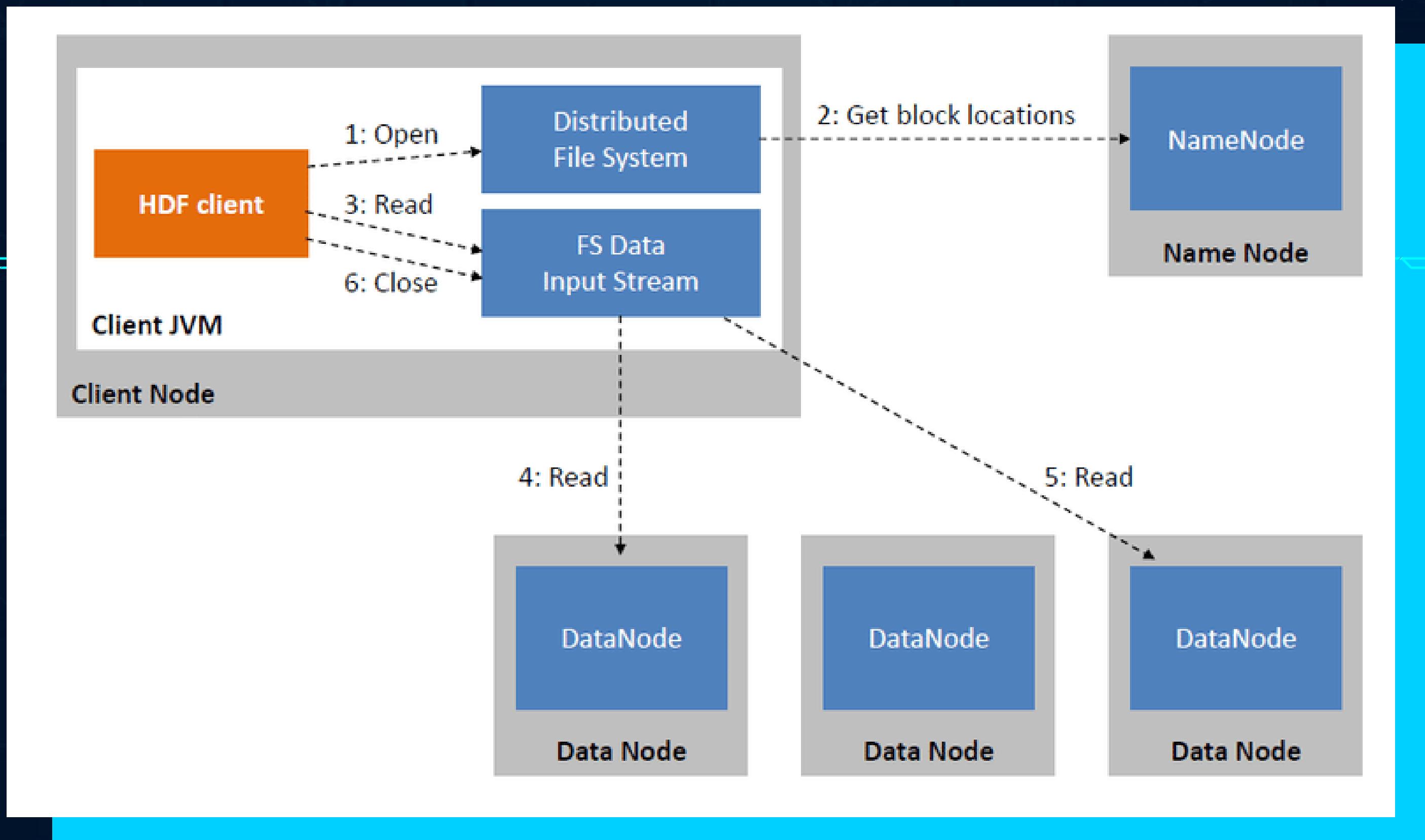
Data Node 2

5 2 3

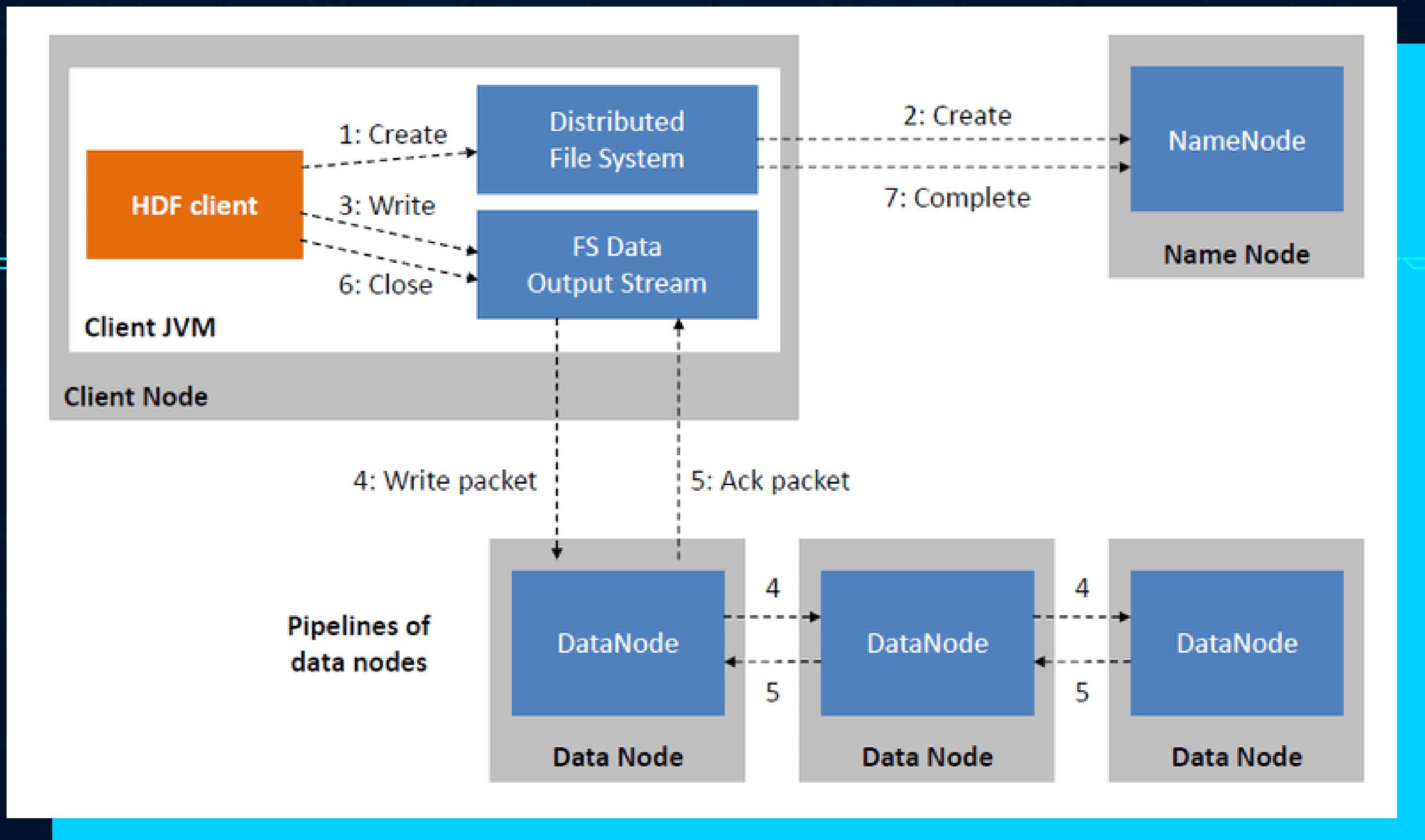
Data Node 3

4 1 3

HDFS READ



HDFS WRITE





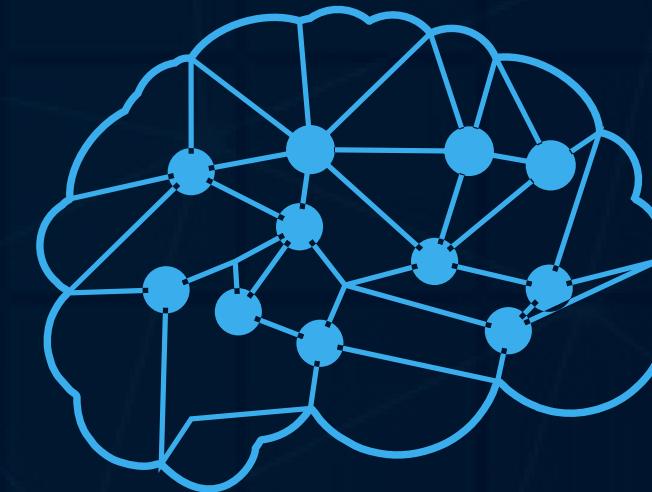
HADOOP MAPREDUCE

HADOOP MAPREDUCE



MapReduce is used for parallel processing of the Big Data, which is stored in HDFS

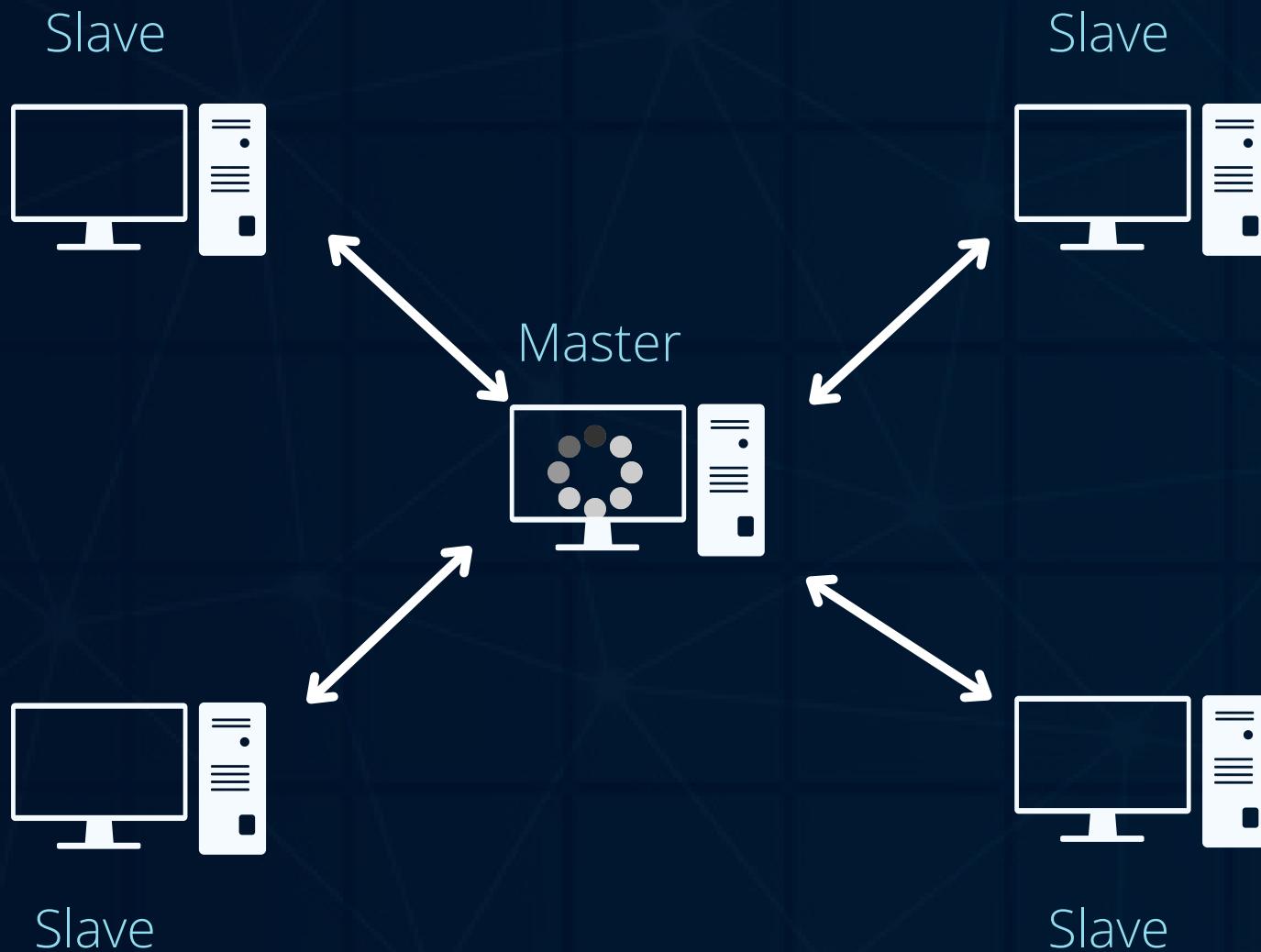
HADOOP MAPREDUCE



- Hadoop Mapreduce is a processing component for Hadoop
- It Processes Data in Distributed Environment

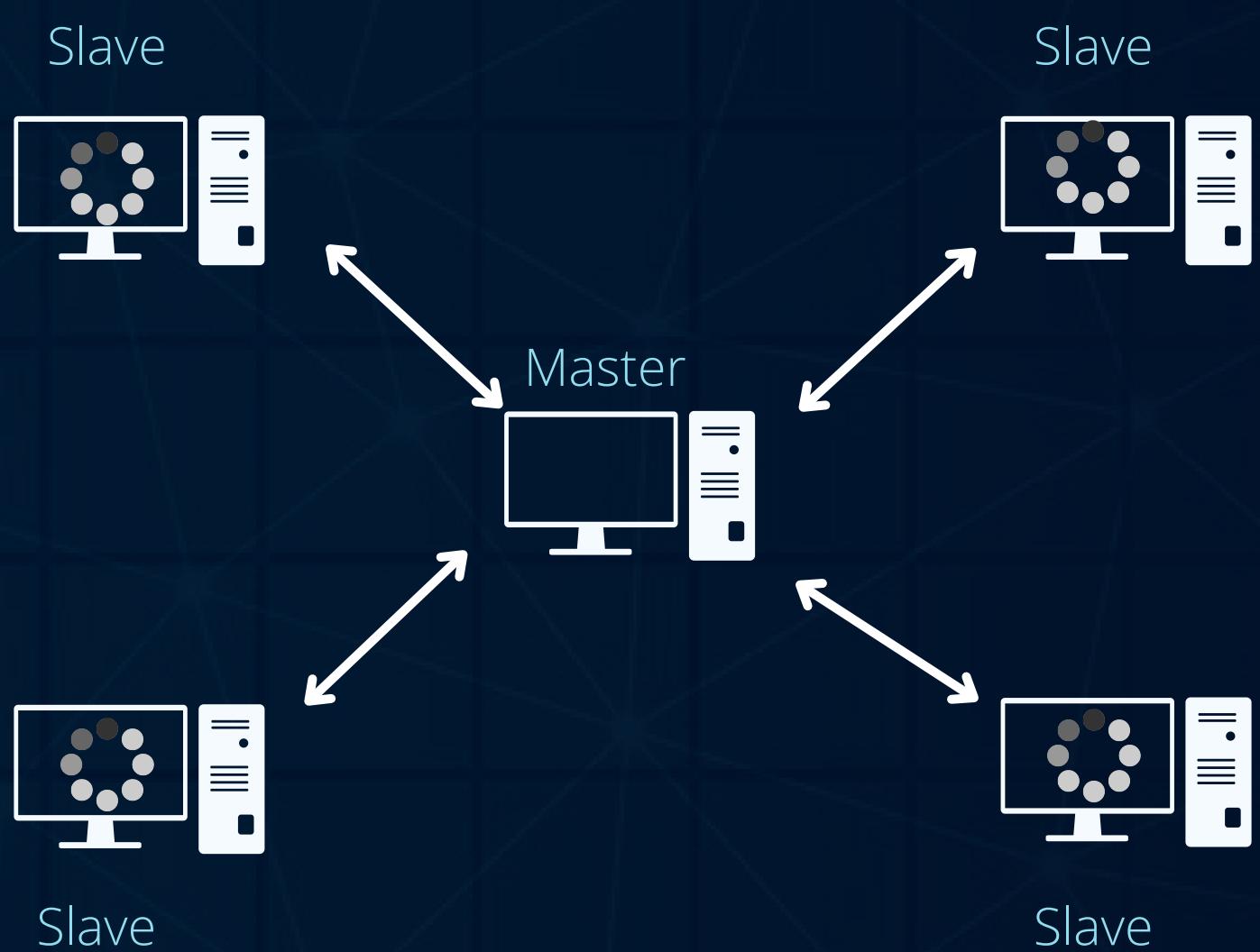
Traditional Approach

Data is processed at the Master



MapReduce Approach

Data is processed at the Slave Node

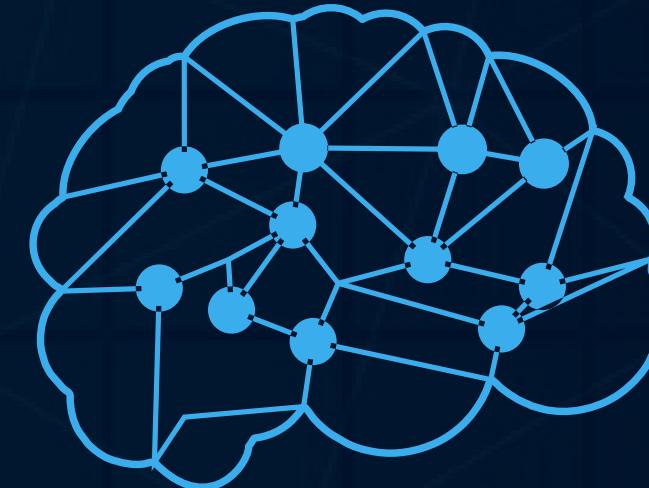


ADVANTAGE 1: PARALLEL PROCESSING



- Data is Processed in Parallel
- Processing becomes Faster

ADVANTAGE 2: DATA LOCALITY- PROCESSING TO STORAGE



- Moving Data to Processing is Very Costly
- In MapReduce, we move Processing to Data

ELECTION VOTES COUNTING

Election Votes Casting-

- Votes are stored in Different Booths
- Result Center has result of all the booths.

Data



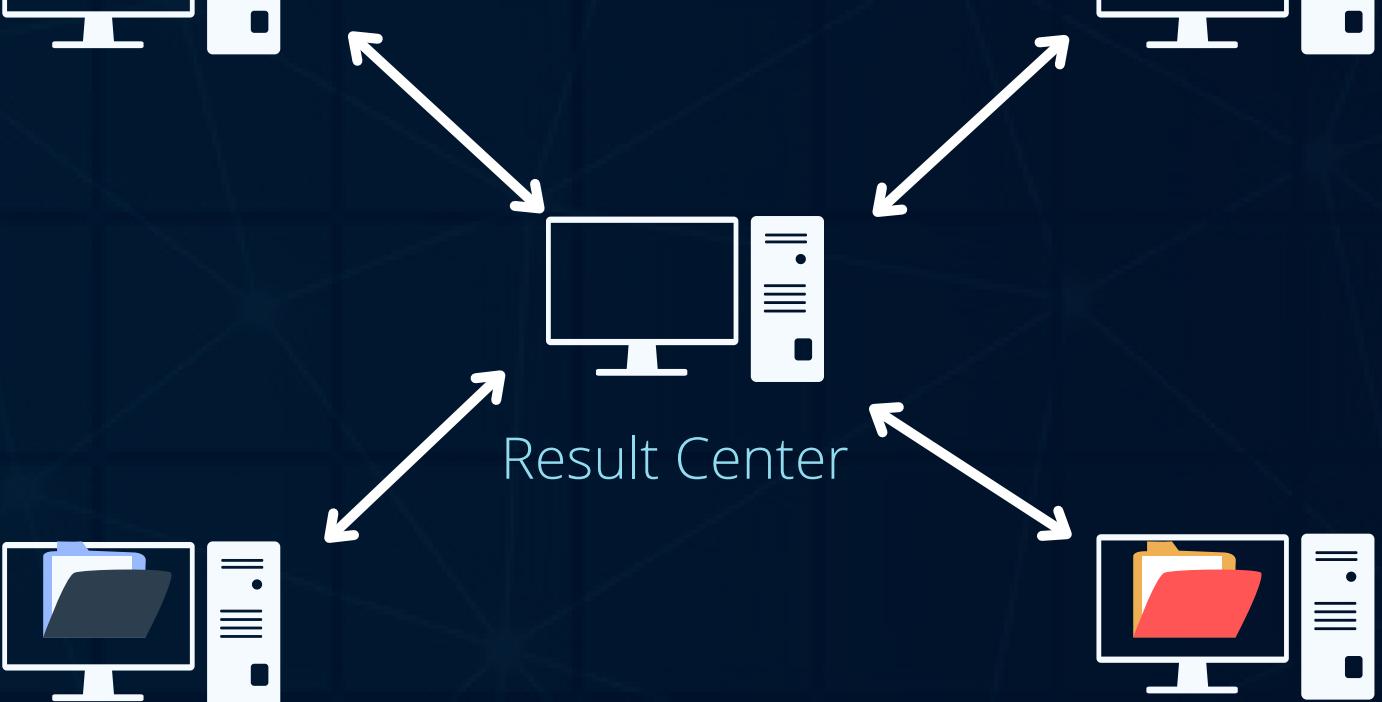
Booth 1



Booth 2



Result Center



Booth 3



Booth 4

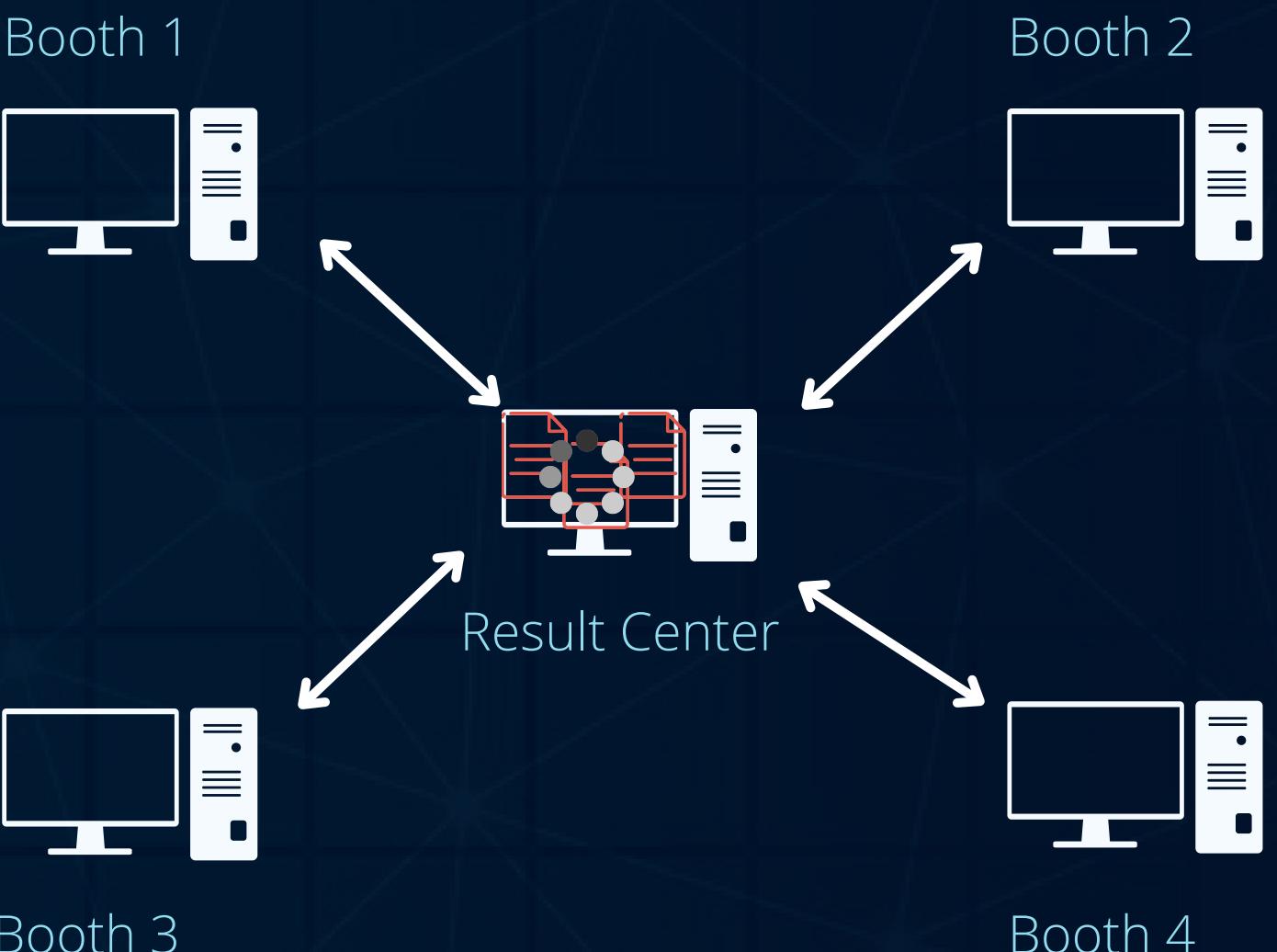


ELECTION VOTES COUNTING TRADITIONAL WAY

Counting - Traditional Approach

- Votes are moved to Result Center for counting
- Moving all the votes to Center is costly
- Result Center is over-burdened
- Counting takes time

Data

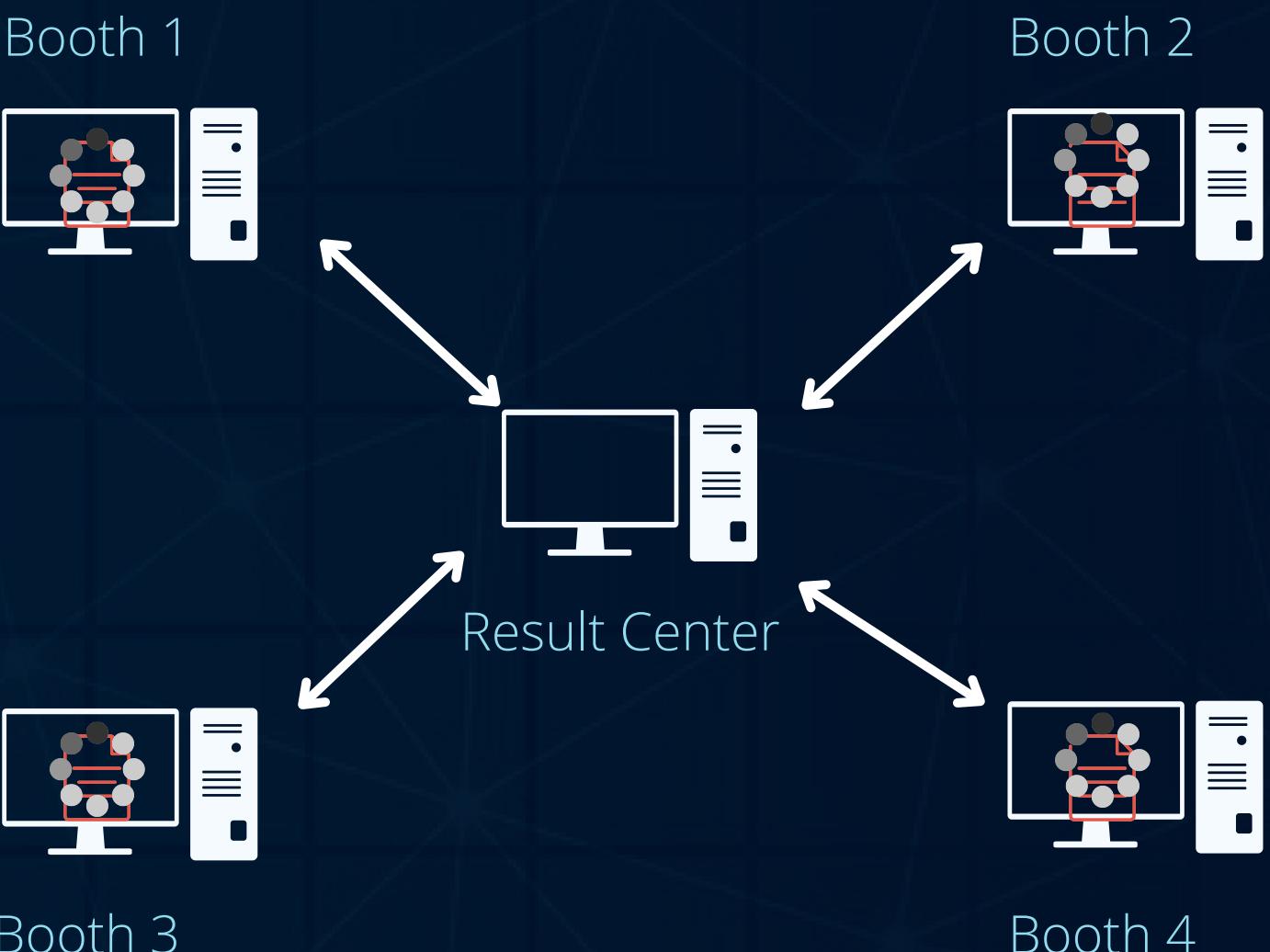


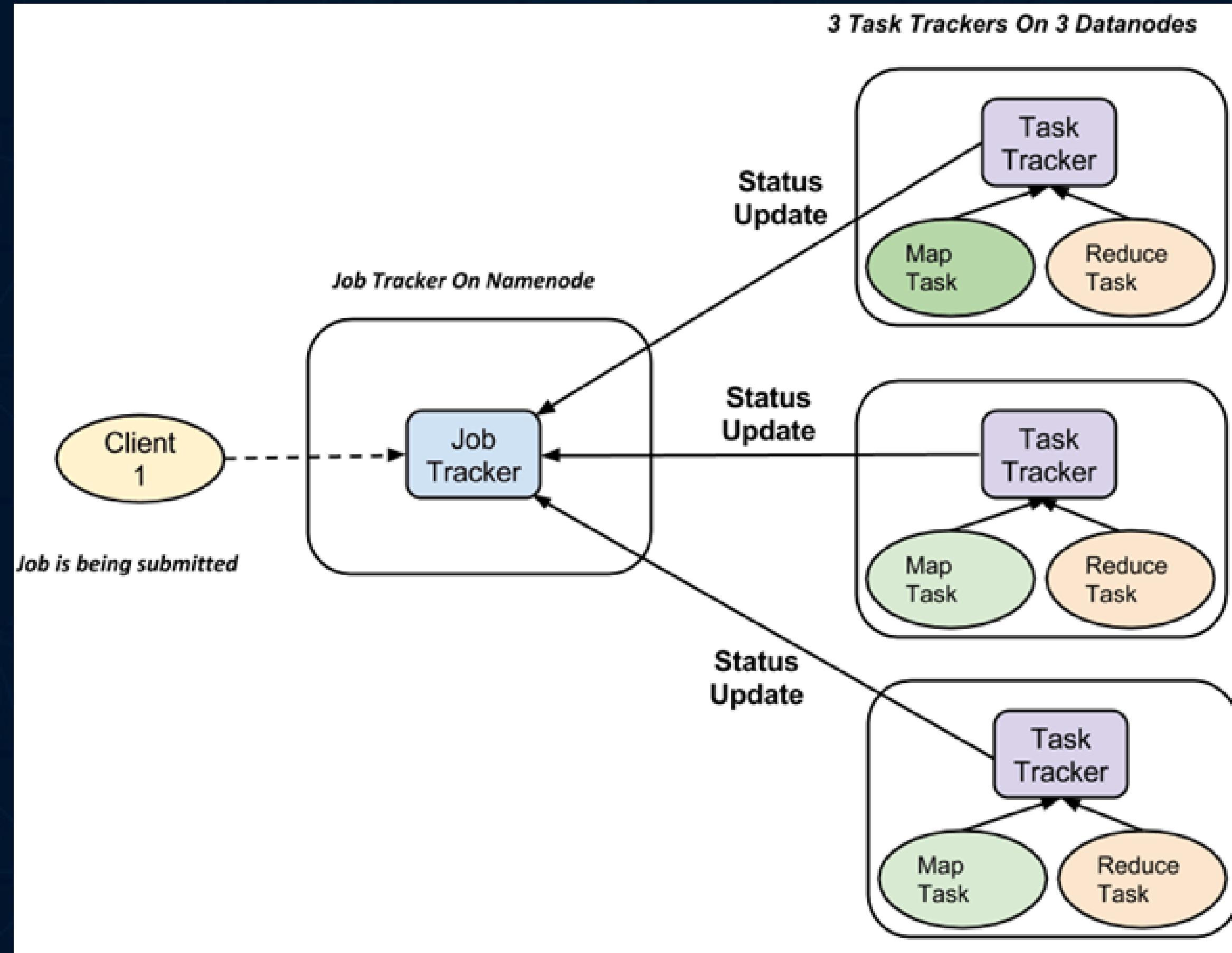
ELECTION VOTES COUNTING MAPREDUCE WAY

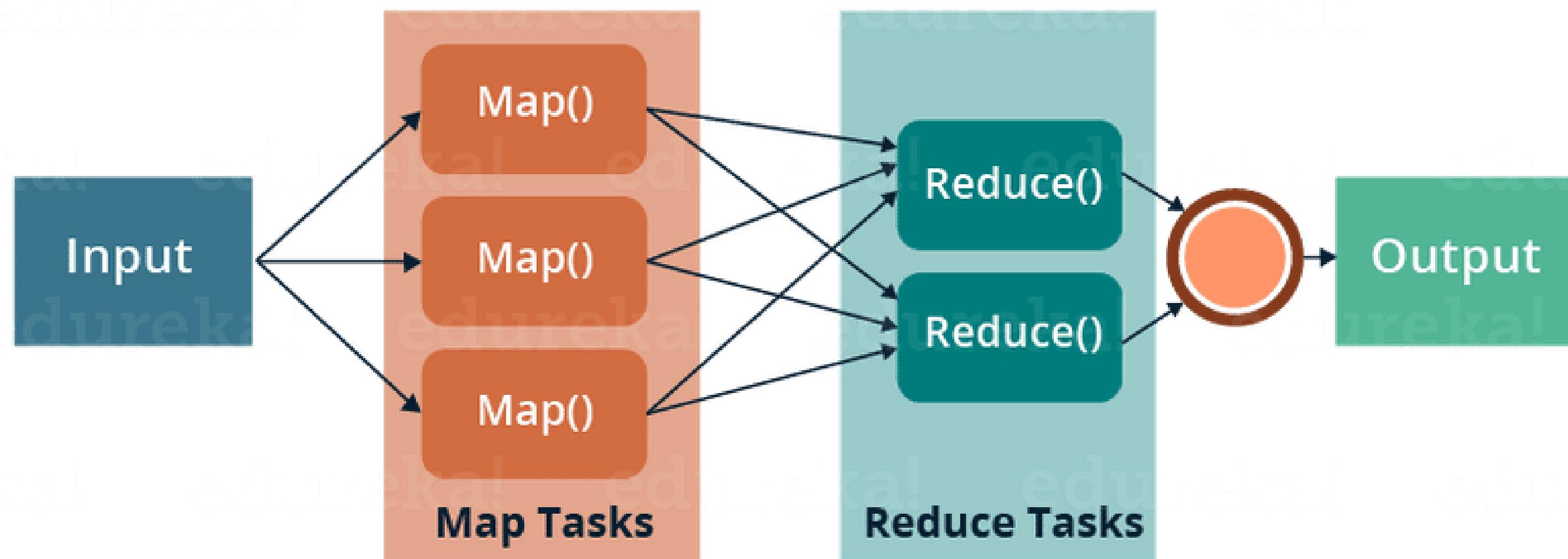
Counting - MapReduce Approach

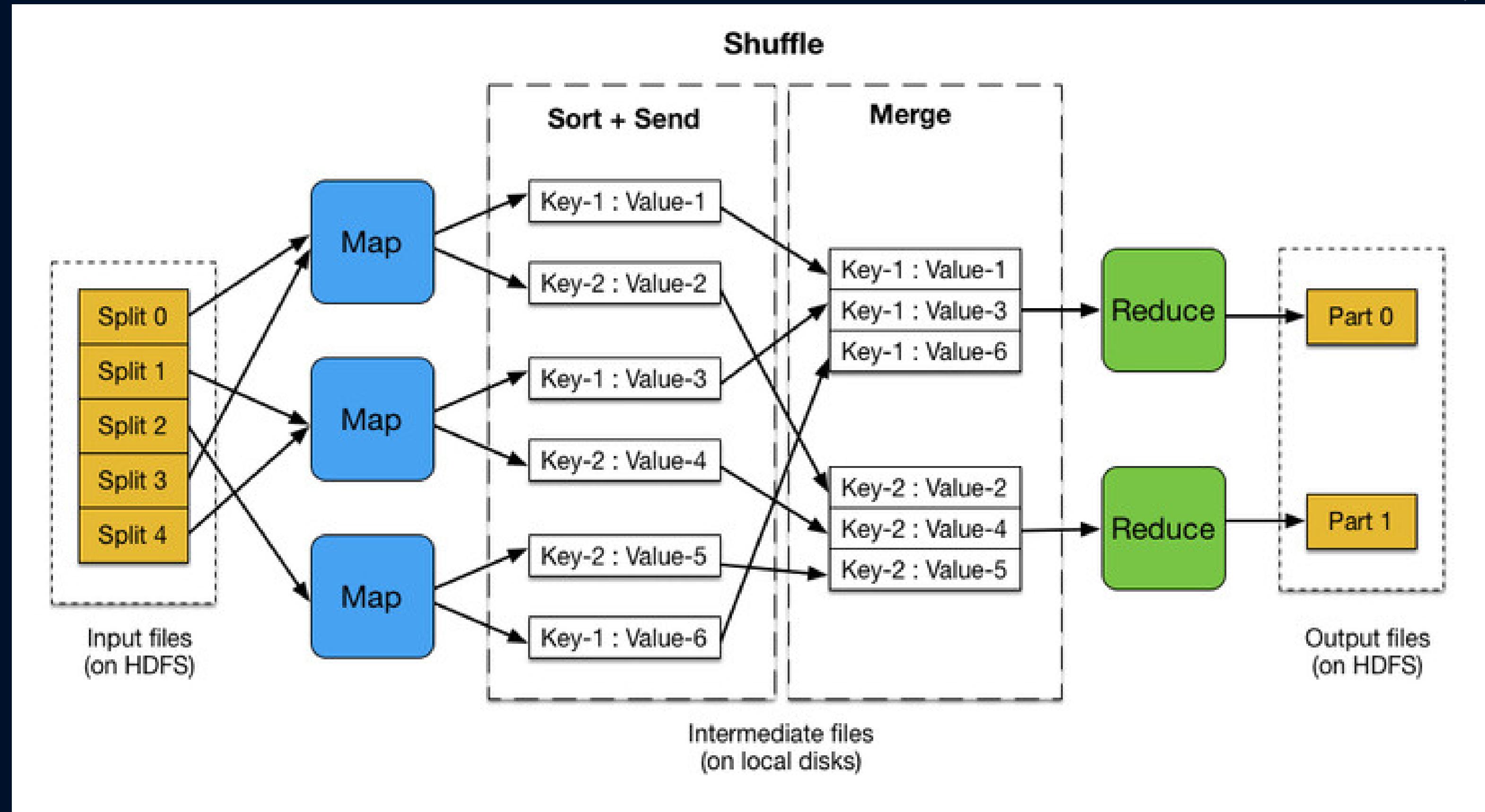
- Votes are counted at individual booths
- Booth-wise results are sent back to the result center
- Final Result is declared easily and quickly using this way

Data

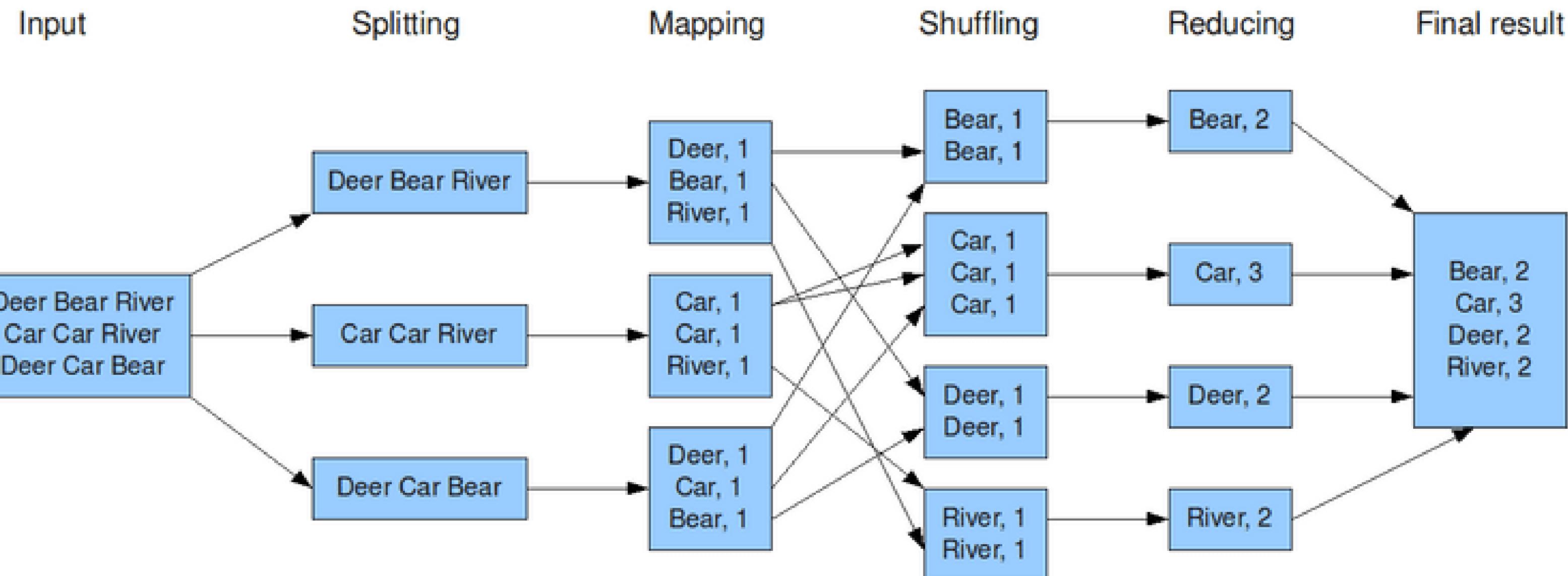








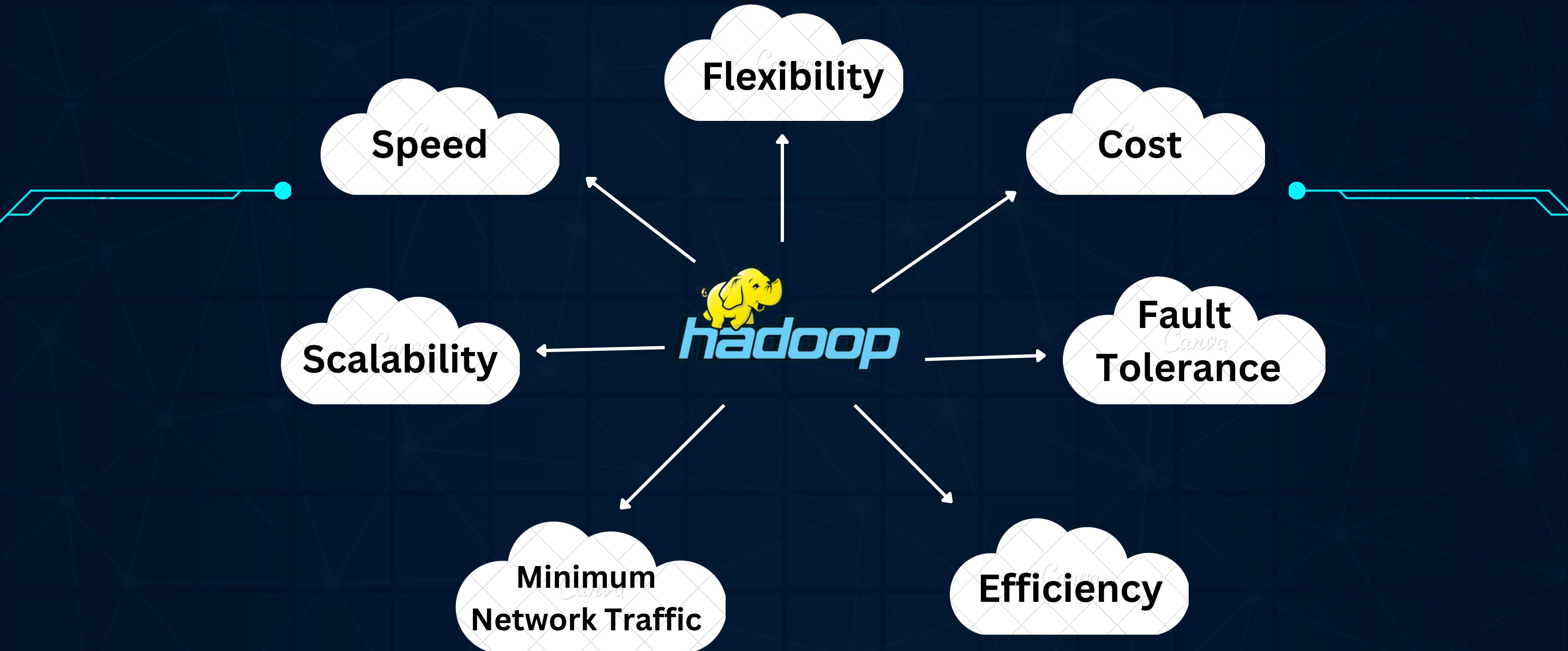
The overall MapReduce word count process





HADOOP ADVANTAGE

ADVANTAGES



THANK
YOU

QUIZZ
TIME!

