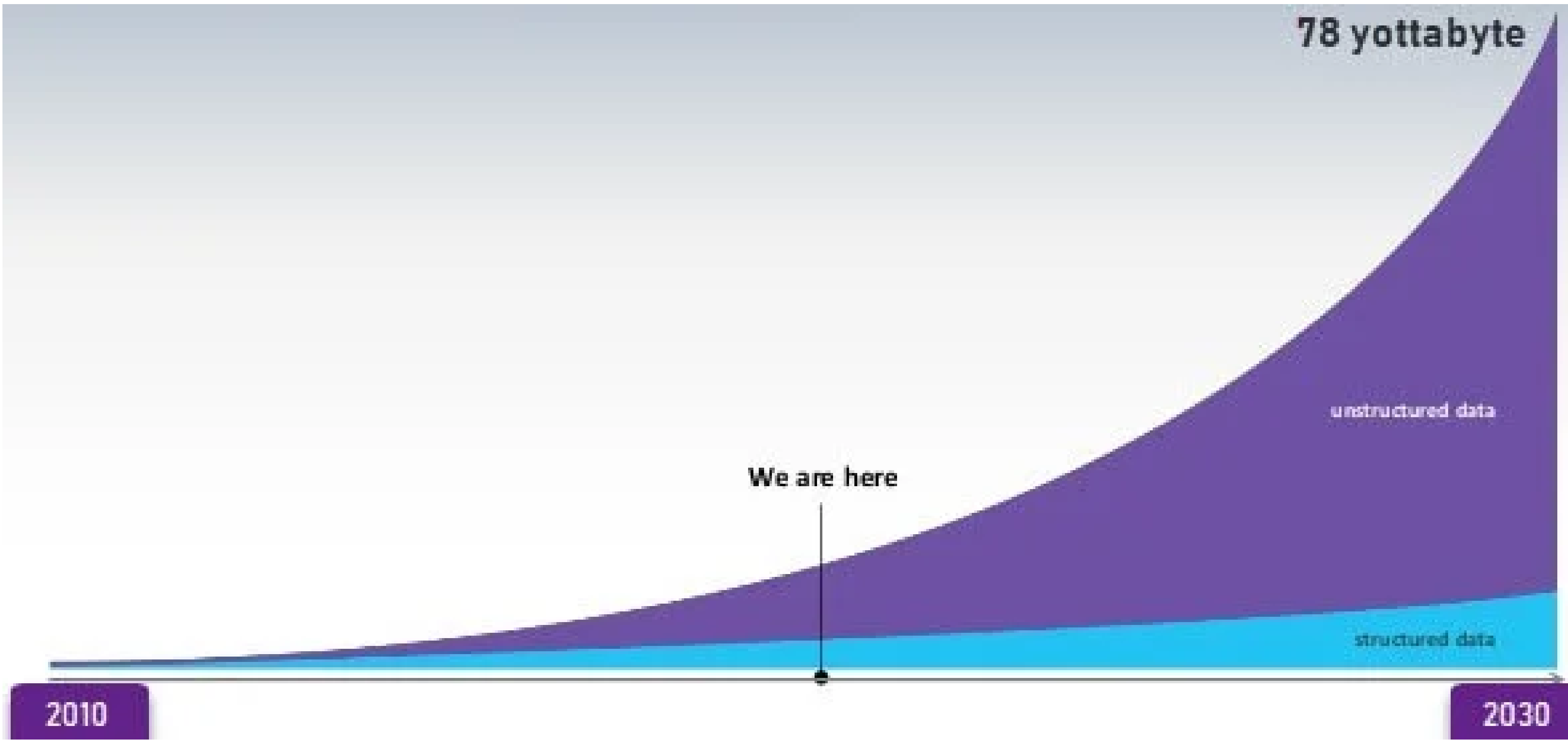# introduction to big data

## DR. AMRITPAL SINGH

Data that is always increasing and cannot be processed and stored on a single machine is termed as Big Data.

# Data Growth over the years

# Big Data Examples

The New York Stock Exchange is an example of Big Data that generates about one terabyte of new trade data per day.

# Big Data Examples

statistic shows that 500+terabytes of new data get ingested into the databases of social media site Facebook, every day.

# Big Data Examples

A single Jet engine can generate 10+terabytes of data in 30 minutes of flight time. With many thousand flights per day, generation of data reaches up to many Petabytes.
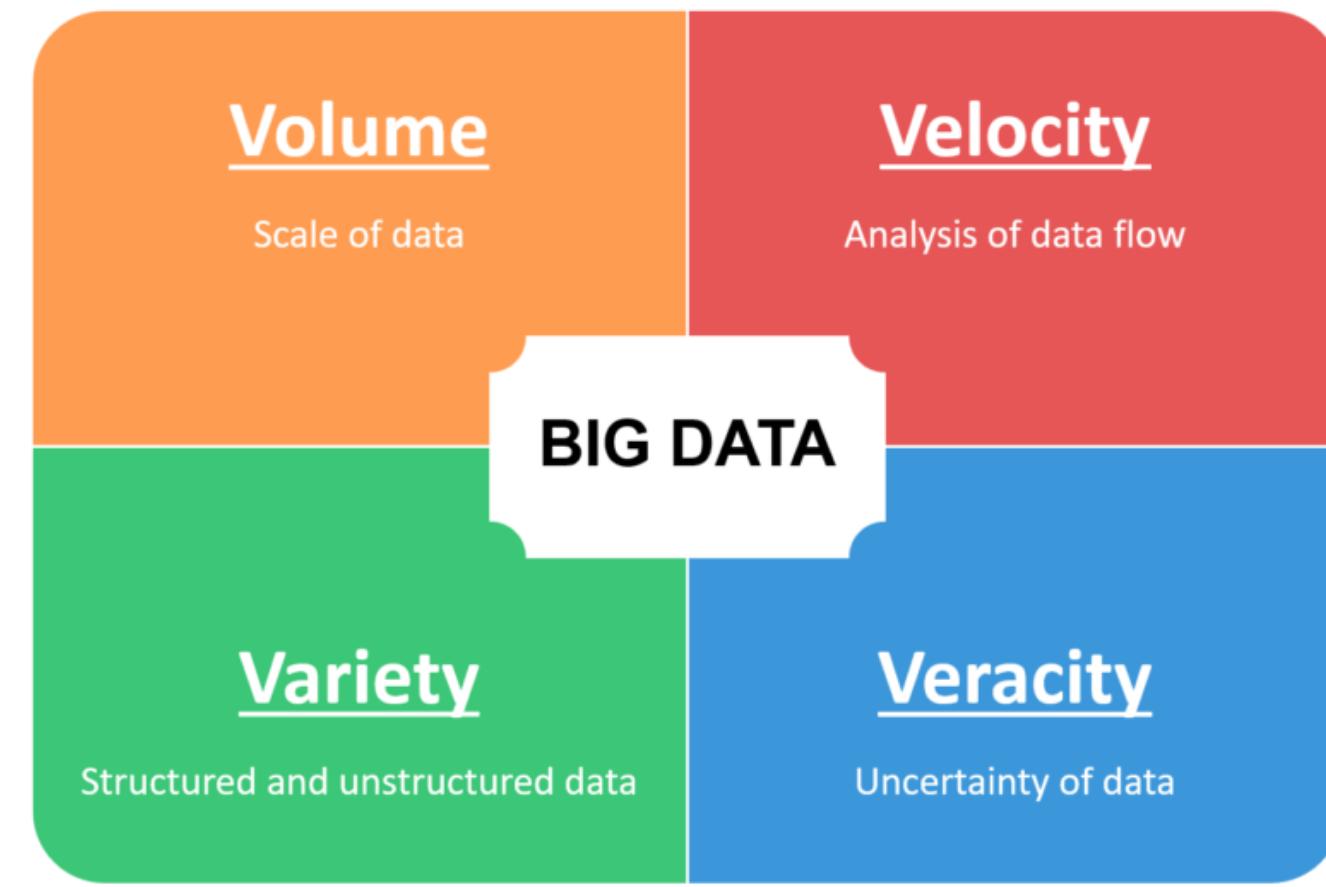
# Types of Data

Structured Data
Unstructured Data
Semi-Structured Data

# V's of Big Data

Volume
Variety
Velocity
Variability
Veracity
Visualization
Value

# Big Data Tools

| | |
|---|---|
| Data Storage & Processing | No sequel Database (cassandra) |
| Data Analysis | Big Data Tools | Realtime Processing (Spark) |
| Data Warehouse (HIVE) | Messaging System |

# Big Data Pipeline



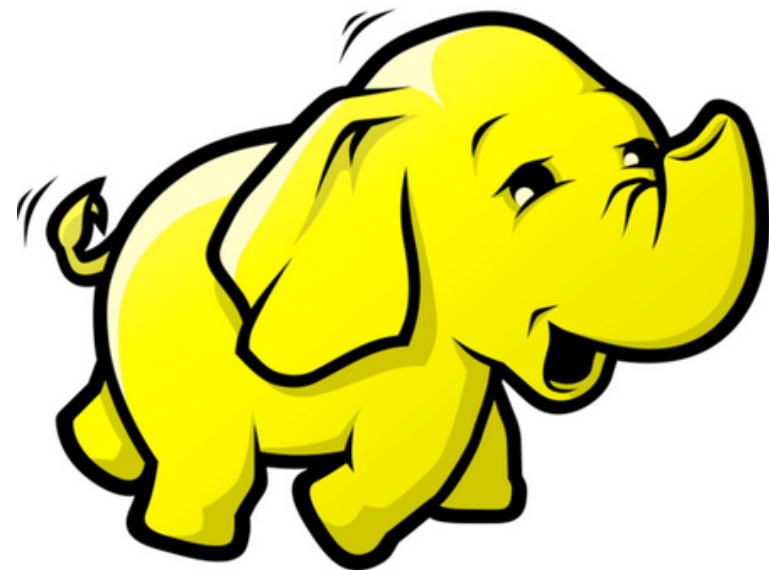**Big Data Ingestion** → **Data Validation, Cleanup & Processing** → **Data Analysis** → **Visualization**

# Intro to Hadoop

Apache Hadoop is an open source framework that is used to efficiently store and process large datasets ranging from GBs to PBs.
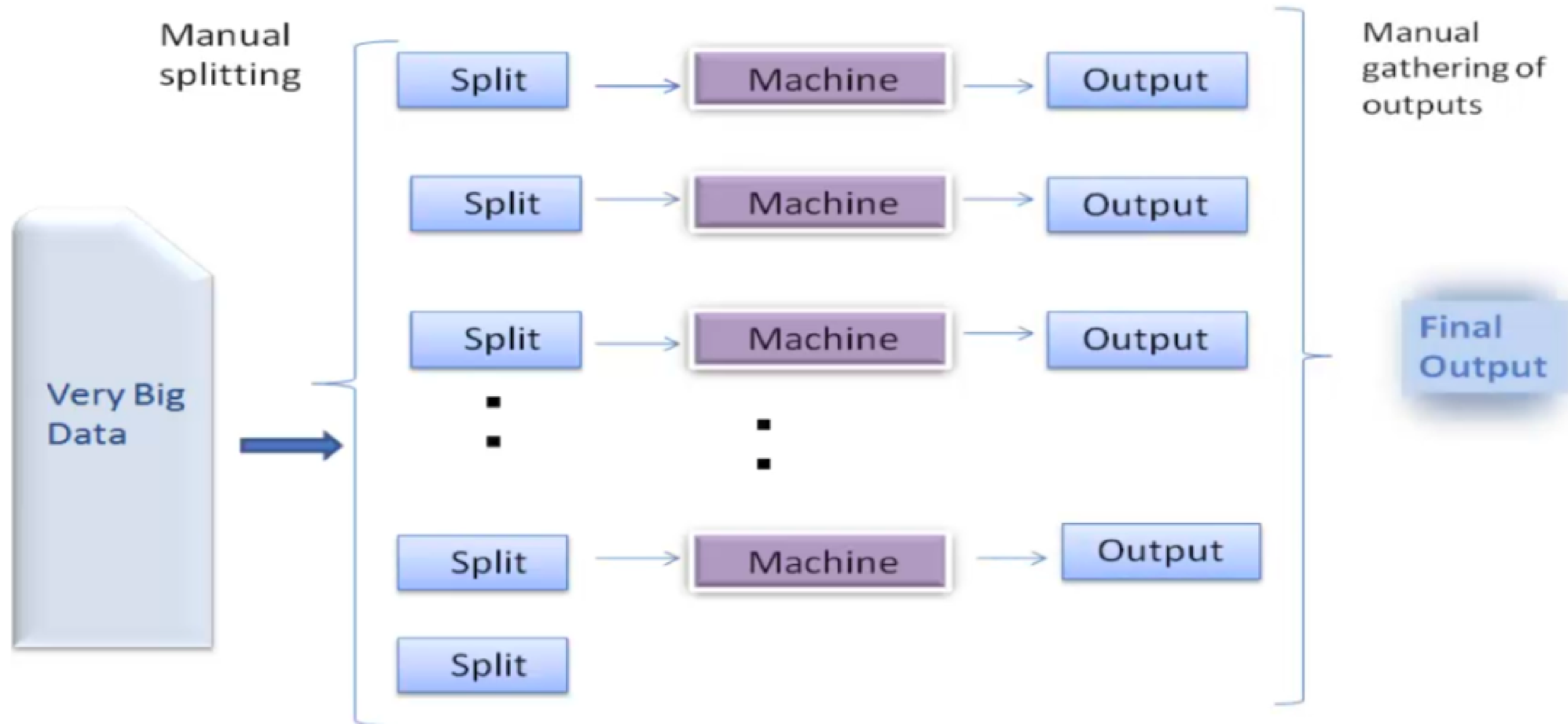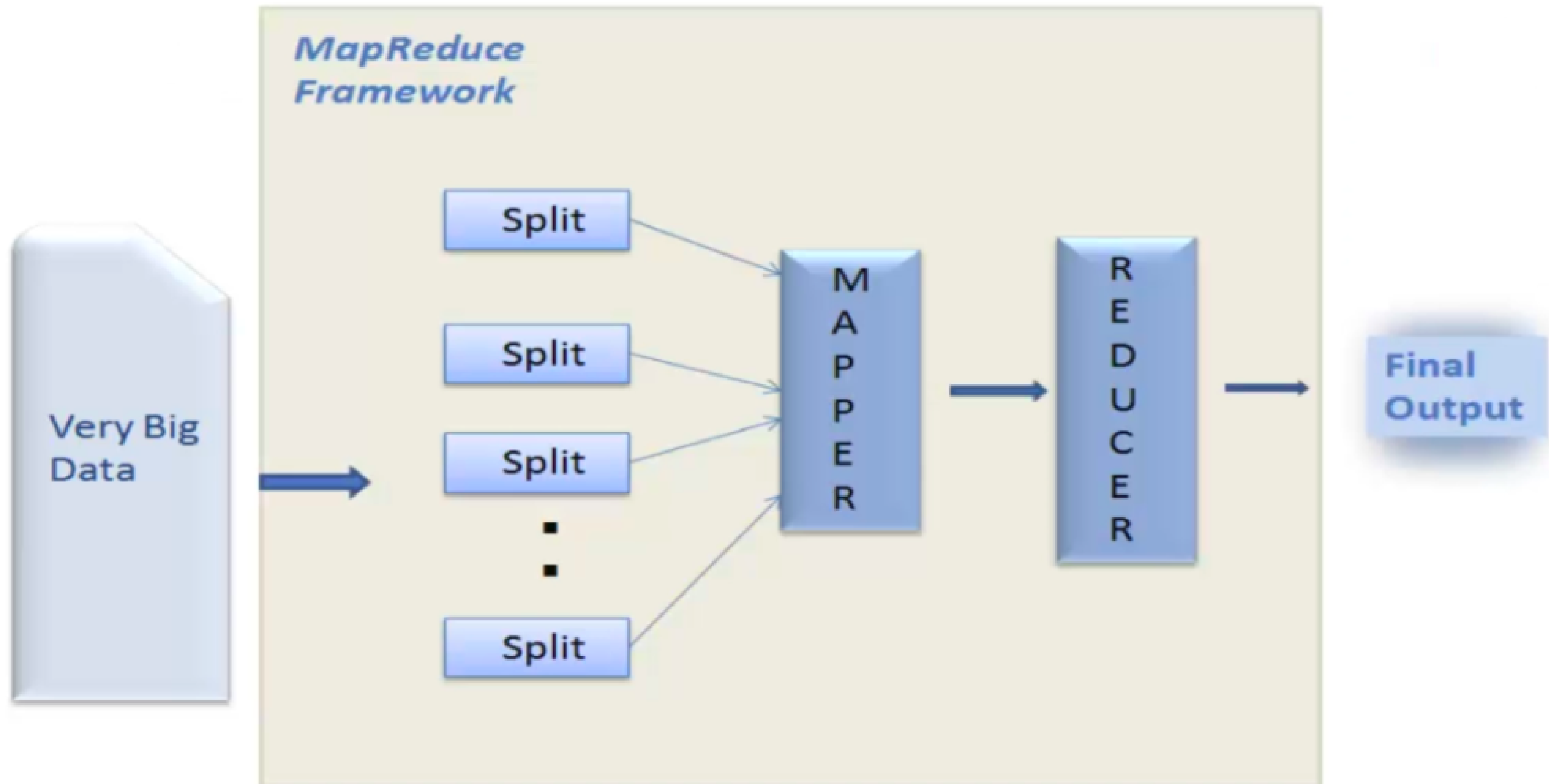
# WHAT IS MAP REDUCE?

MapReduce is the core component for data processing in Hadoop framework.

In layman's term Mapreduce helps to split the input data set into a number of parts and run a program on all data parts parallel at once. The term MapReduce refers to two separate and distinct tasks.

# TRADITIONAL APPROACH

# MAP REDUCE APPROACH

# MAP REDUCE APPROACH