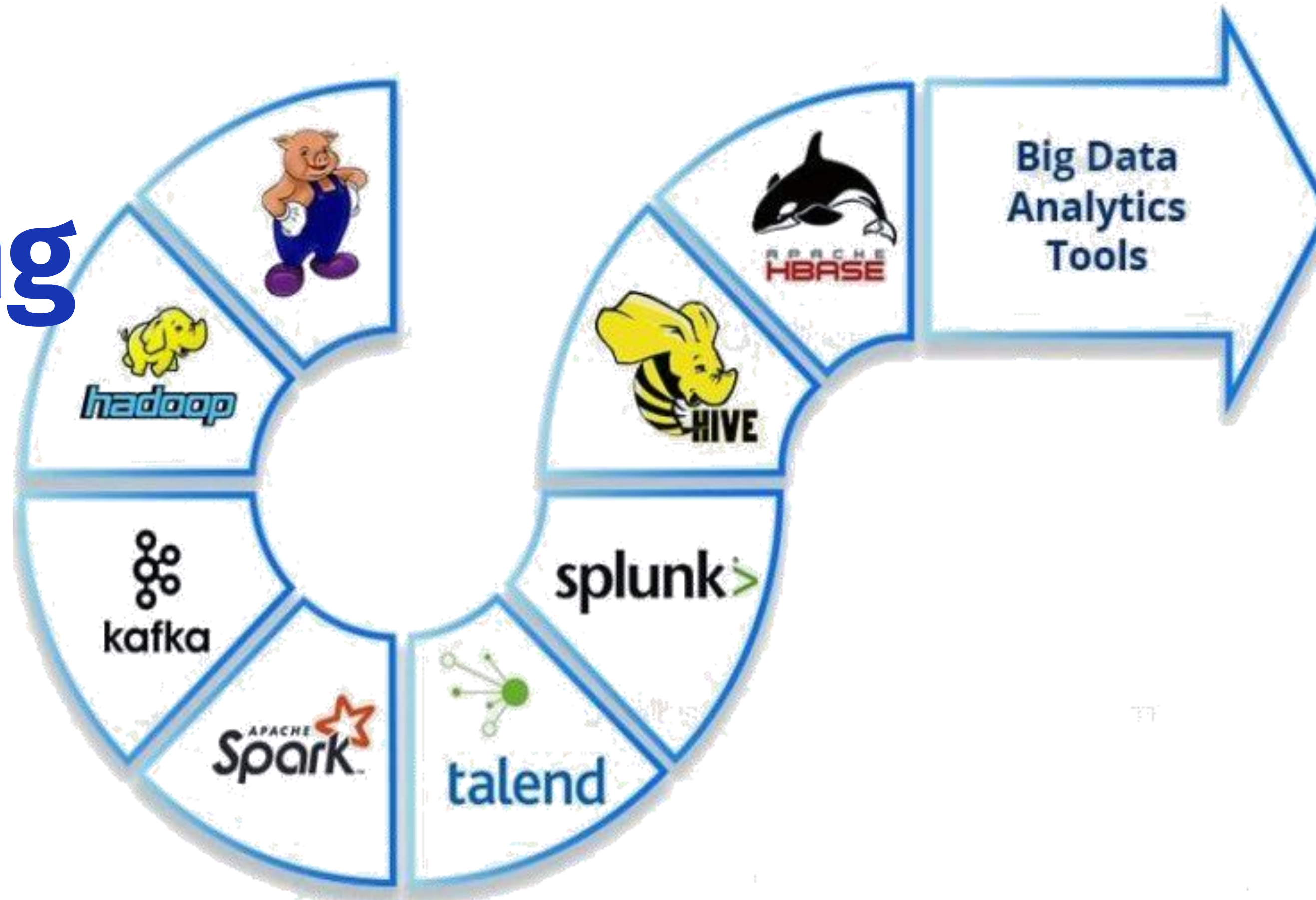




# Parallel Processing in Big Data

Professor : Dr.Tadayon

Presented by : Mobarake Aftabi







# Parallel DBMS

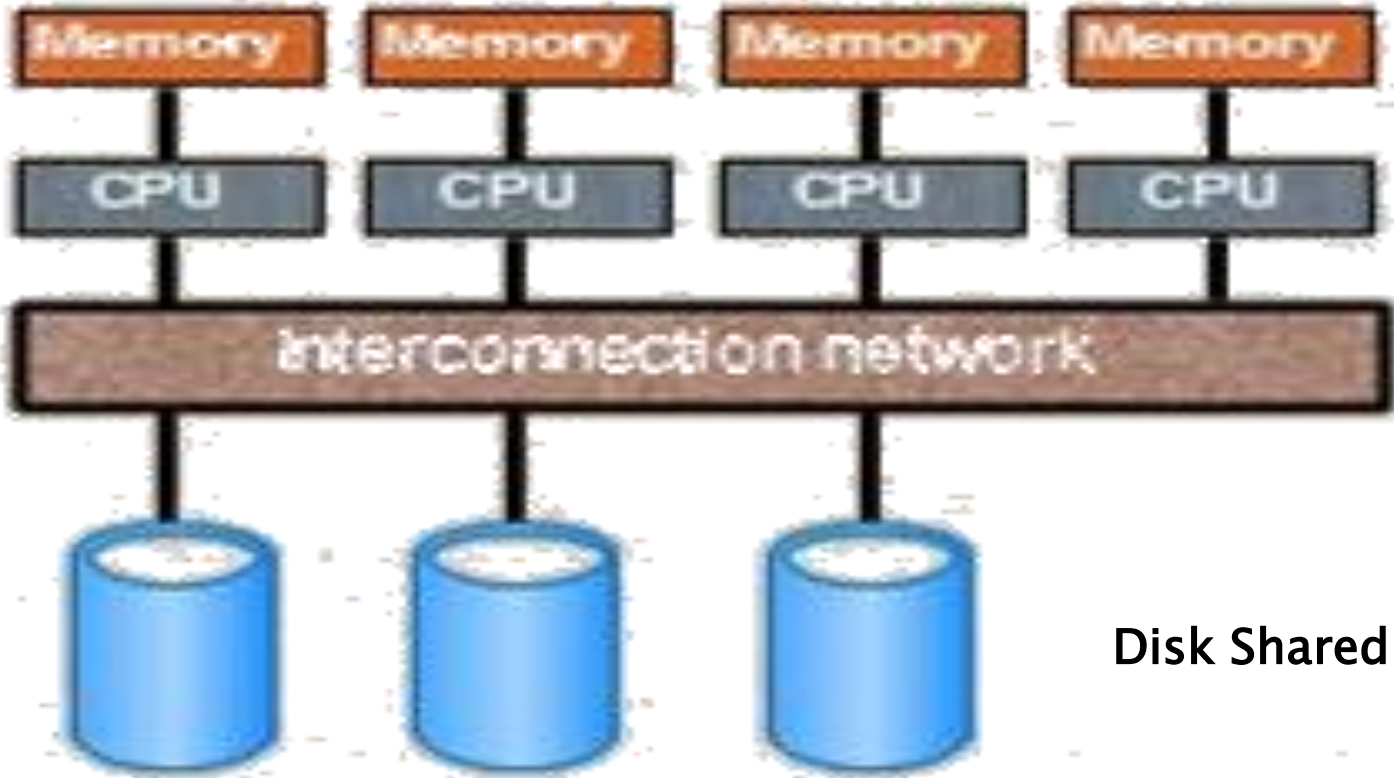
A parallel database management system seeks to improve performance through operational parallelization such as data loading, indexing, and query evaluation.

Has been able to run database systems on clusters of SN nodes since the 1980s. These systems support standard SQL spreadsheets, and as a result, the fact that data is distributed across multiple systems is obscured from the end user.

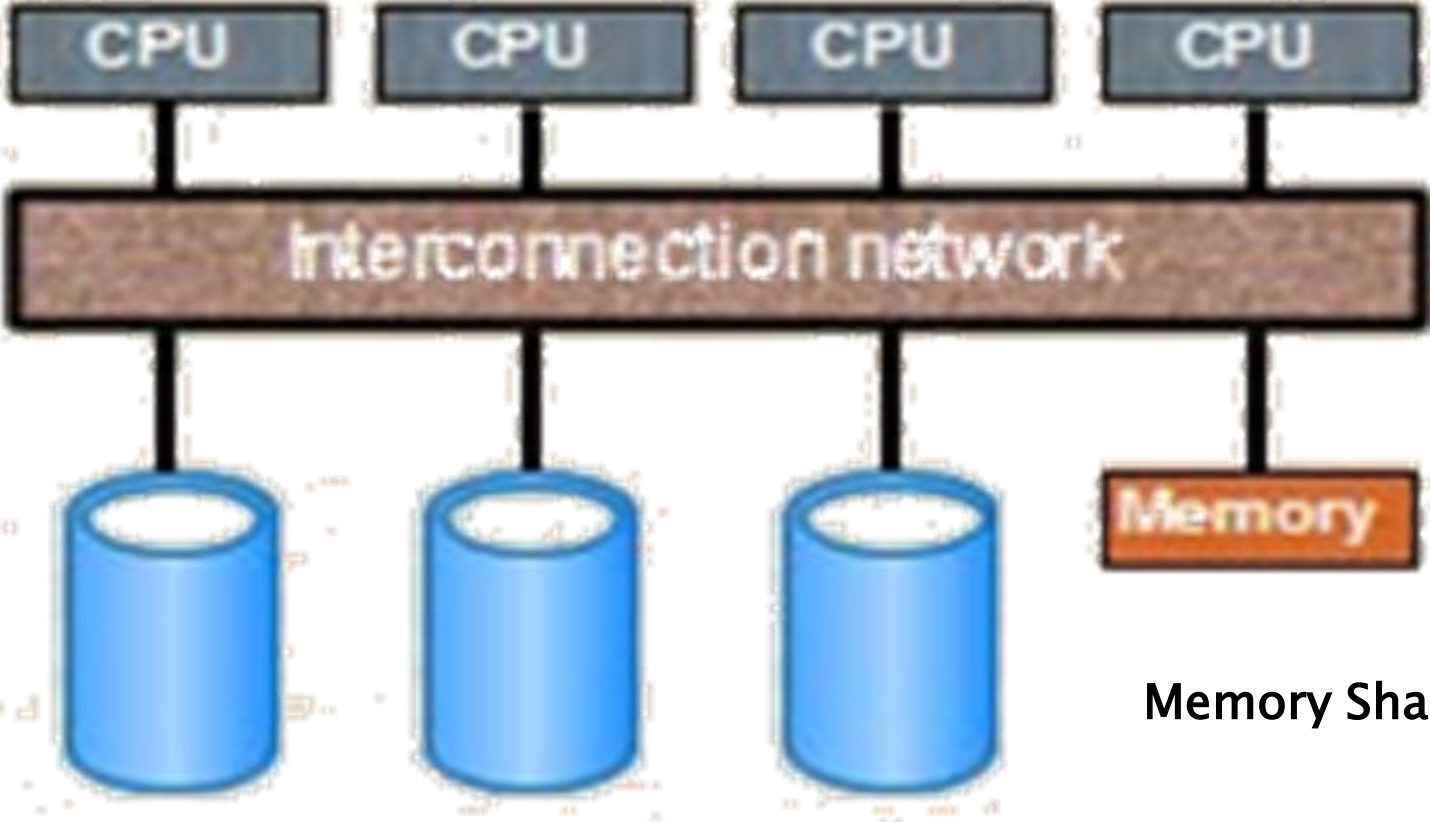
All or most of the tables are divided into nodes in a cluster

The system uses an optimizer to convert SQL statements to a query map whose execution is distributed across several nodes.

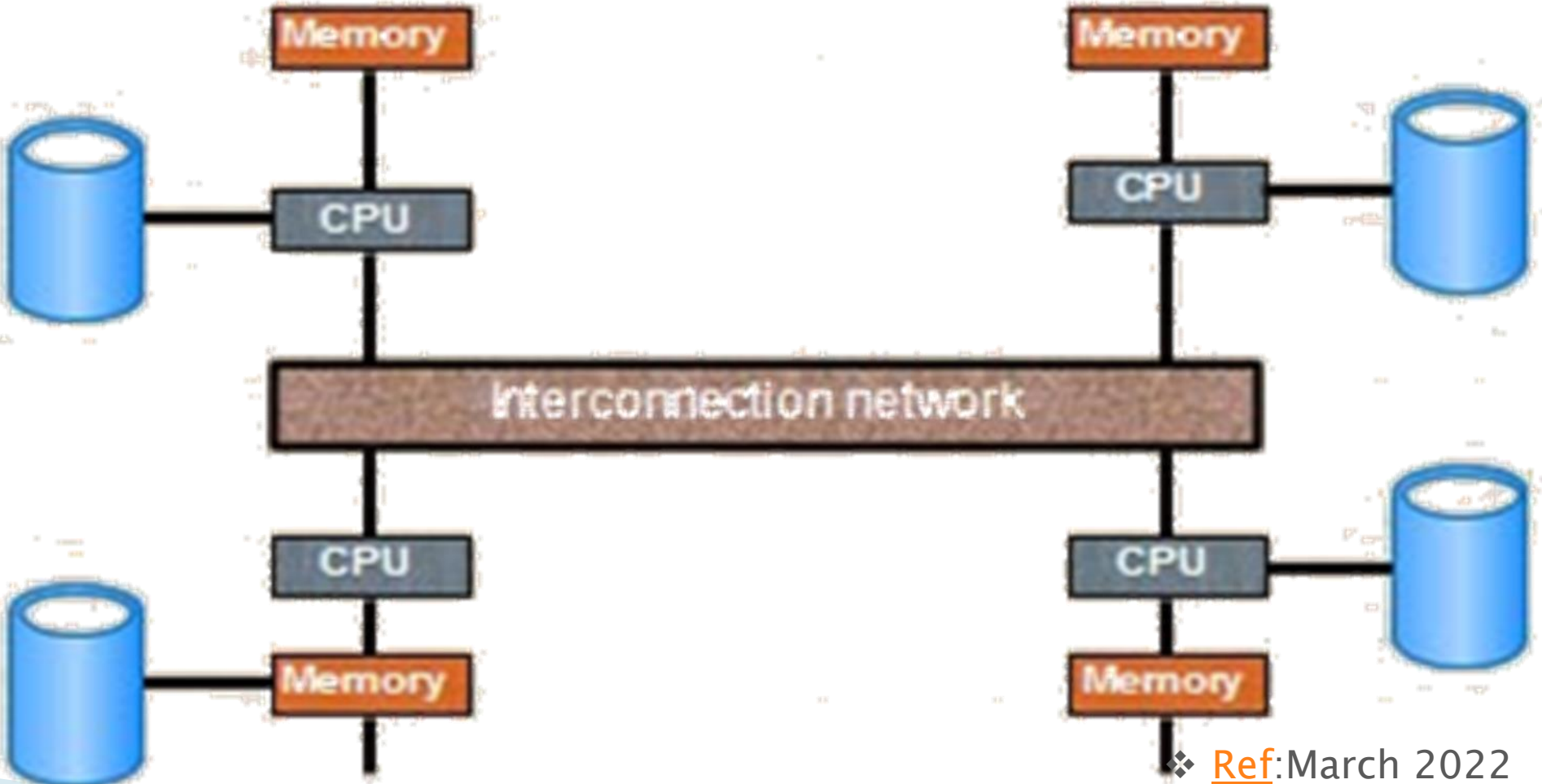
# Parallel DBMS



Disk Shared



Memory Shared



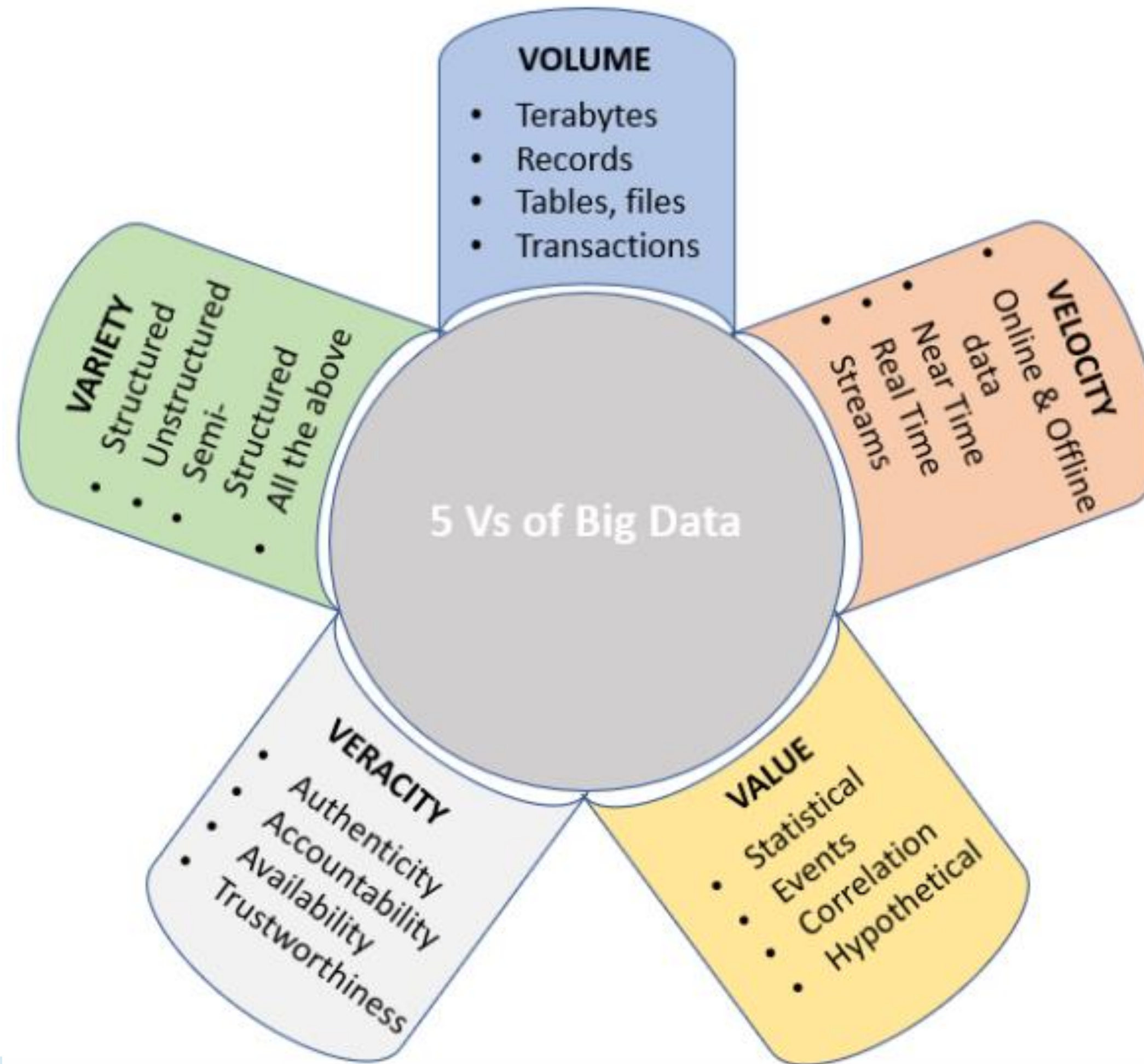
Share Nothing

## What is Big Data?(Cont.)



Author's name	Definition
Batty	Big data are massive in size and cannot fit into Excel spreadsheets comprising approximately 16 000 columns and 1 million rows.
Havens et al.	Big data cannot be loaded into local storage devices (computer memory).
Fisher et al.	Big data cannot be easily processed and managed in a straightforward manner.
The State Council of People's Republic of China	Big data have several characteristics, such as high application value, fast access speed, large volume, and multiple types.
Bayer and Laney	Big data have large volume, variety, and velocity that demand cost effectiveness and are helpful in decision making.

# Characteristics of big data





Data stores	Column	Cassandra, HBase, Hypertable
	Document	MongoDB, JSON ODM, CouchDB
	Graph	Neo4j, ArangoDB, Bigdata
	Key-Value	DynamoDB, Azure Table Storage, Cassandra, PNUTS, Berkeley DB

(a)

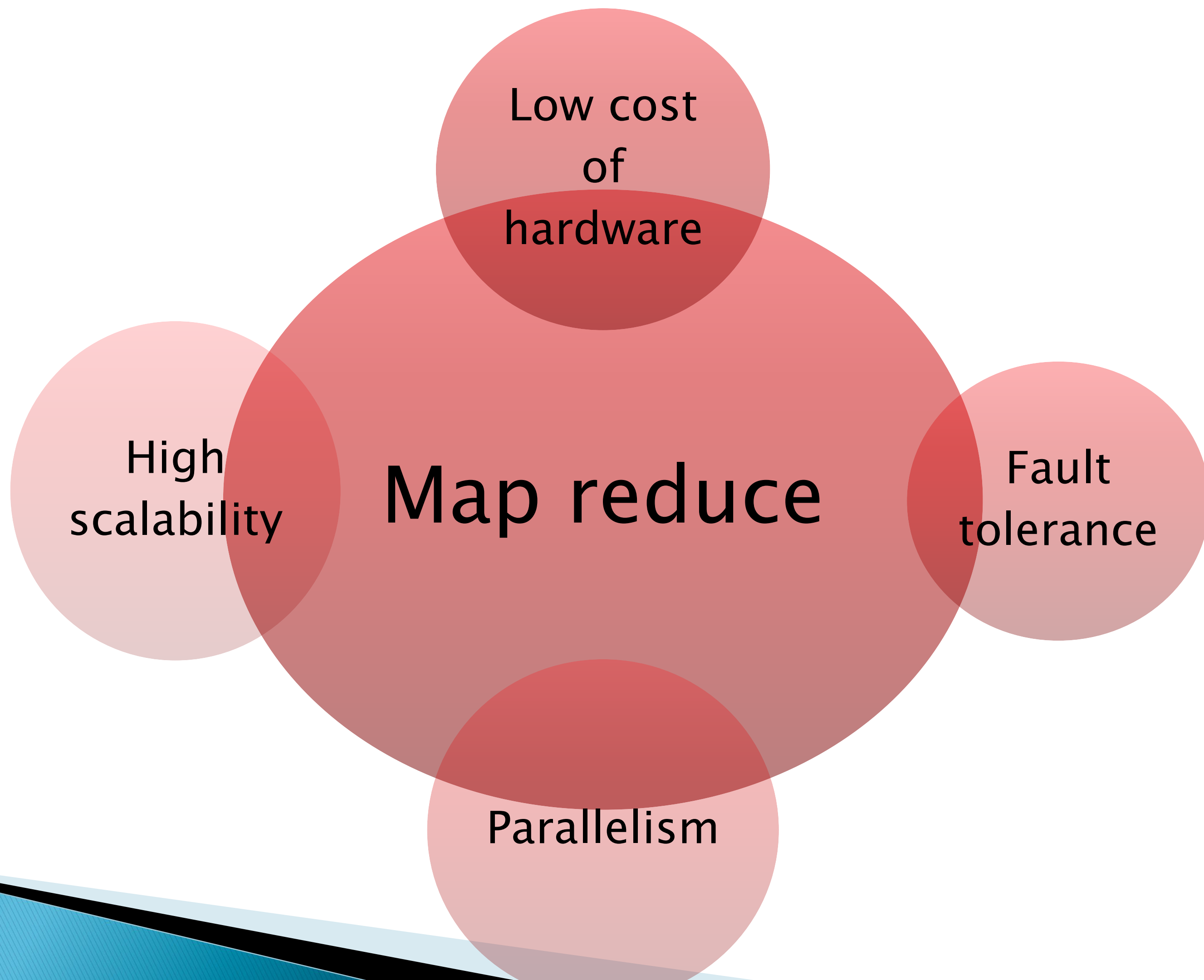
Data processing	Programming	SQL-like	Hive, Shark, Scope, Pig, Dremel, AsterixDB
		User-defined method	Hadoop, Spark, Flink, Dryad, GraphLab, Pregel, Storm, S4
	Input	Batch	Hadoop, Spark, Dryad, Flink, Nephele, Dremel, AsterixDB
		Stream	Storm, S4
		Graph	GraphLab, PowerGraph, Pregel, Hama
		ML	Petuum, MLBase, Mahout

(b)

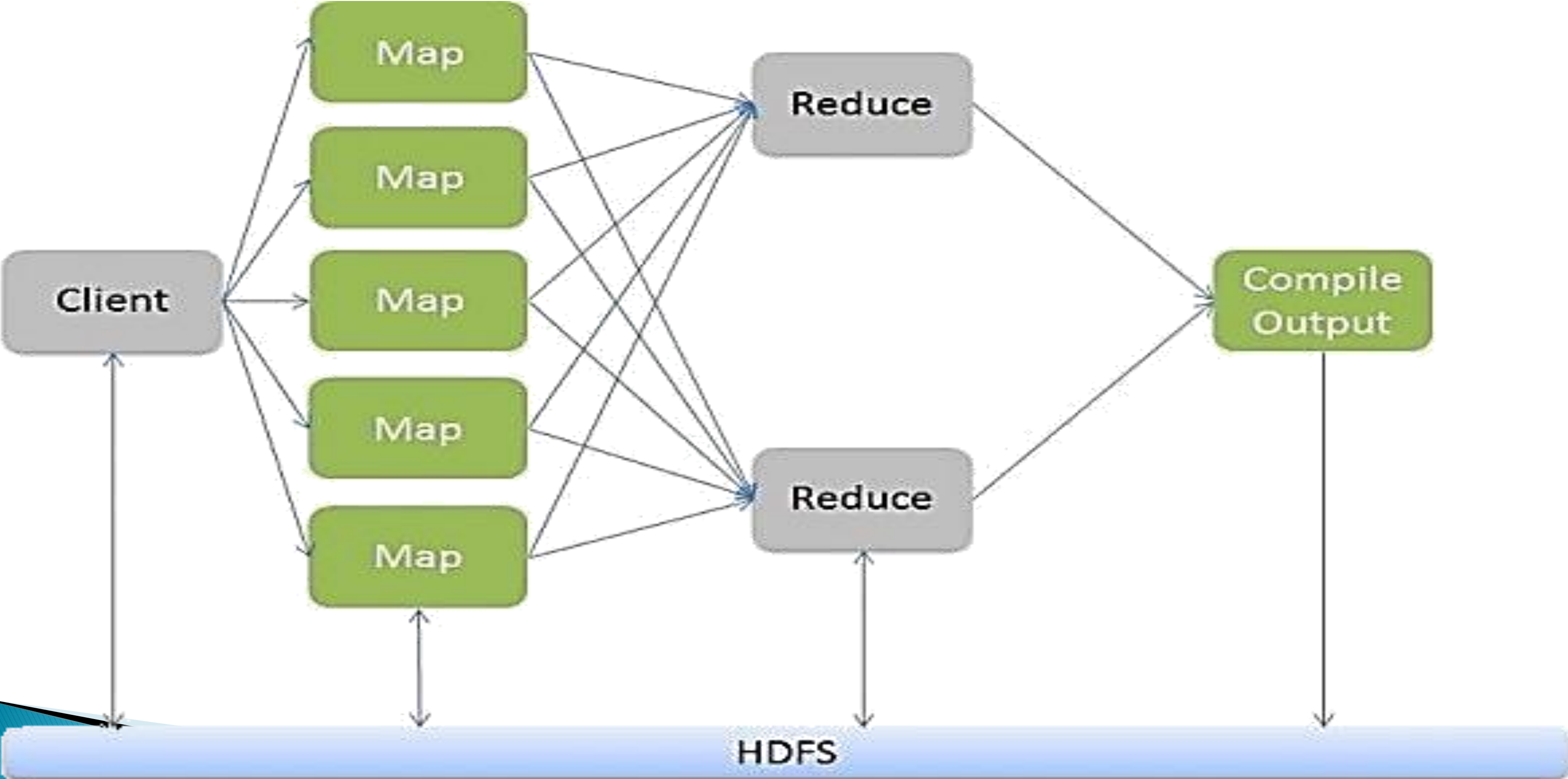
# Map Reduce

MapReduce is a parallel programming model for processing data on clusters that consists of two main phases including the mapping phase (Map) and the reduction phase (Reduce).

This framework divides the big data into subcategories, then divides them into different machines, then combines the separate processes into the final result.

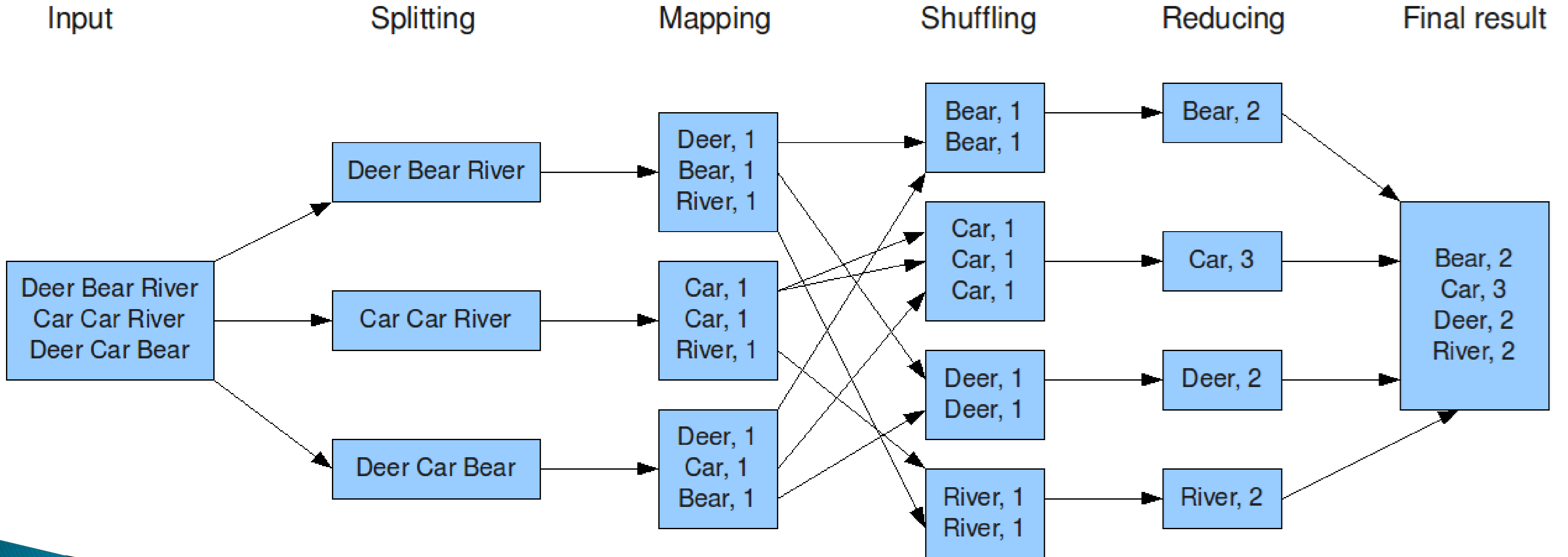


# Map Reduce



# Example

The overall MapReduce word count process



# Map Radius implementations

Apache Hadoop  
MapReduce

Google  
MapReduce

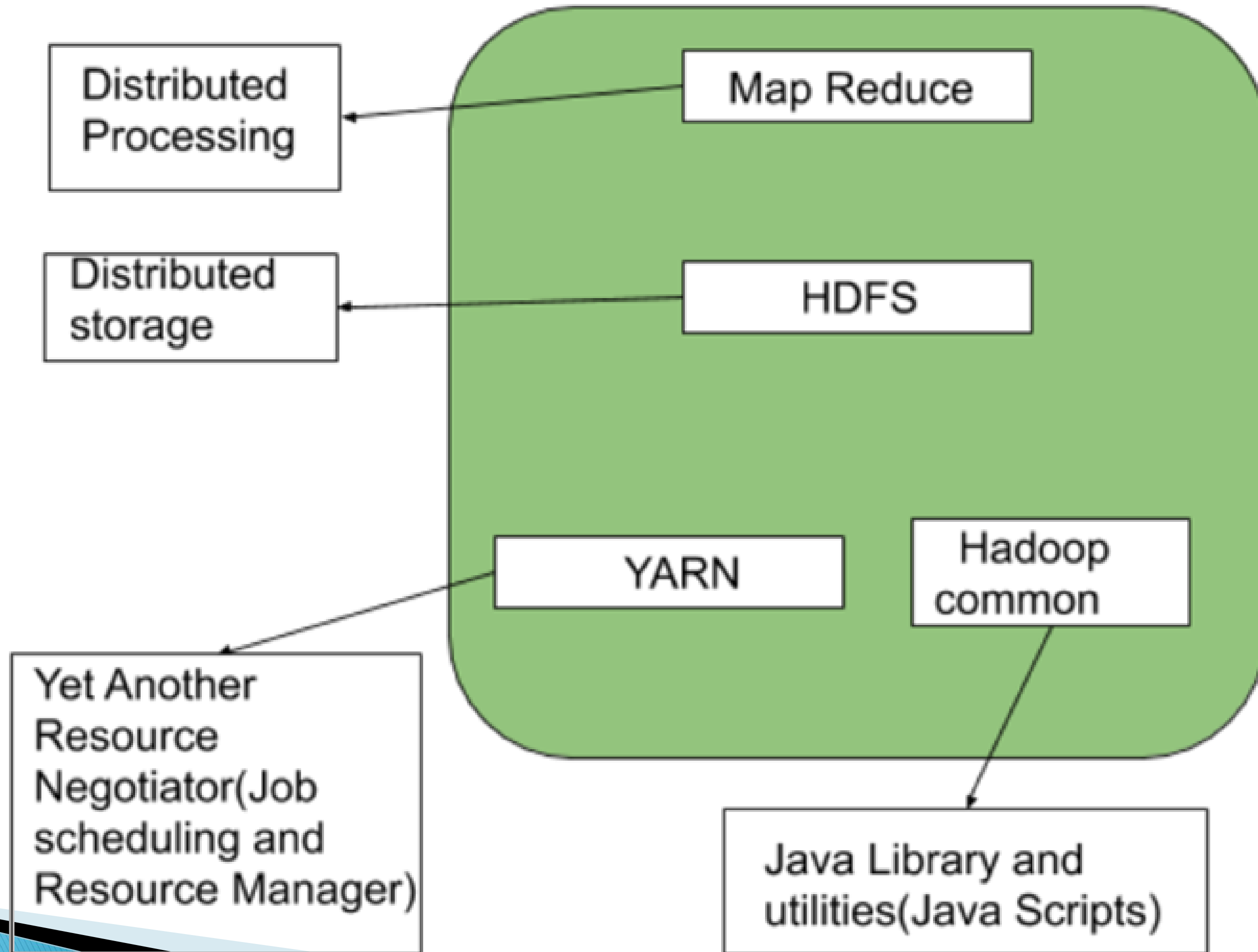
Amazon Elastic  
Mapreduce

**Table 1** Comparison of Parallel DBMSs and MapReduce

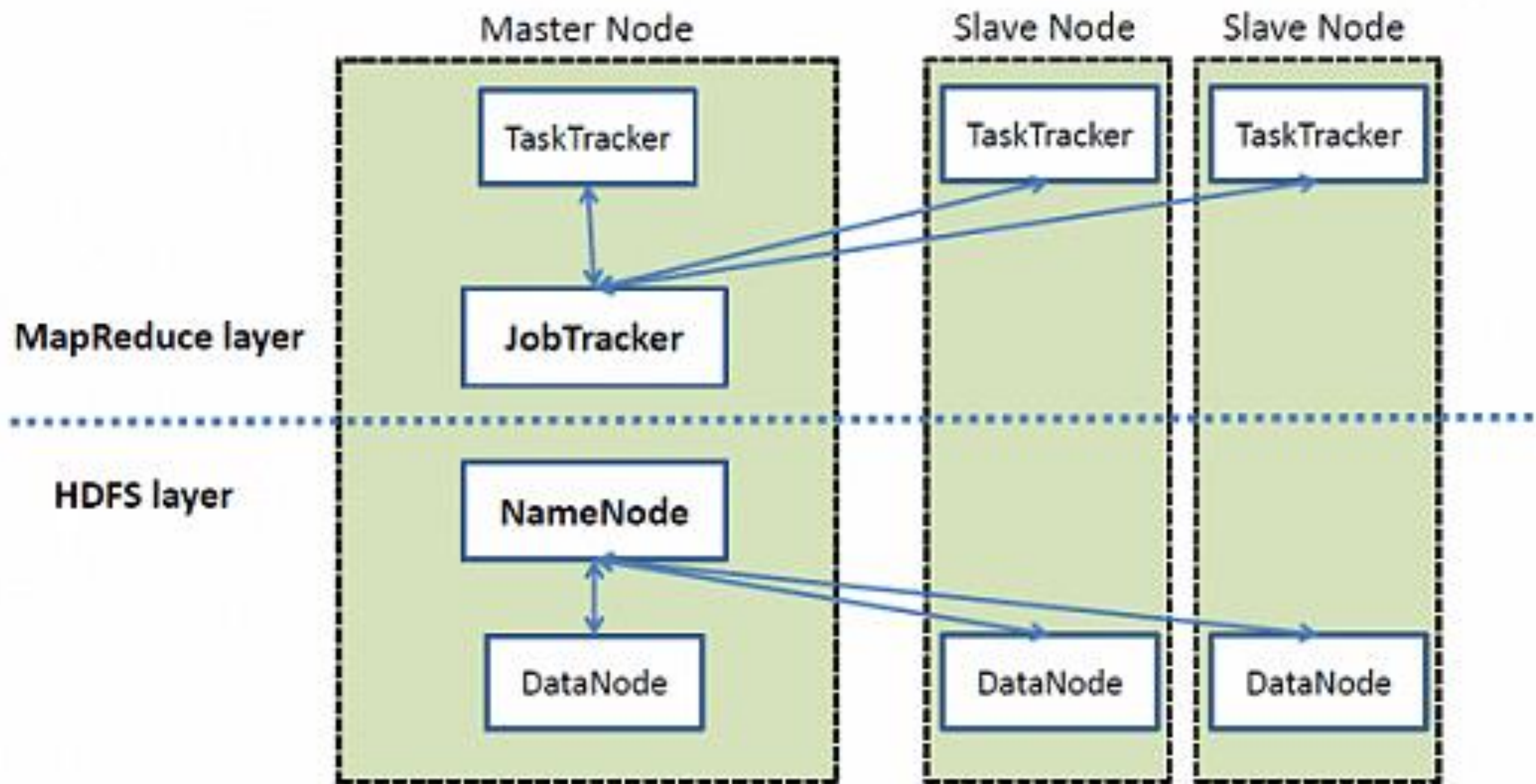
	<b>Parallel DBMS</b>	<b>MapReduce</b>
Schema	✓	Not naturally
Index	✓	Not naturally
Programming	Declarative (SQL)	Imperative (C++, Java, ...)
Optimization	Compression, Column storage, ...	Not naturally
Pre-parsing	✓	Not naturally
Flexibility	Not naturally	✓
Fault tolerance	Transaction-level	Task-level



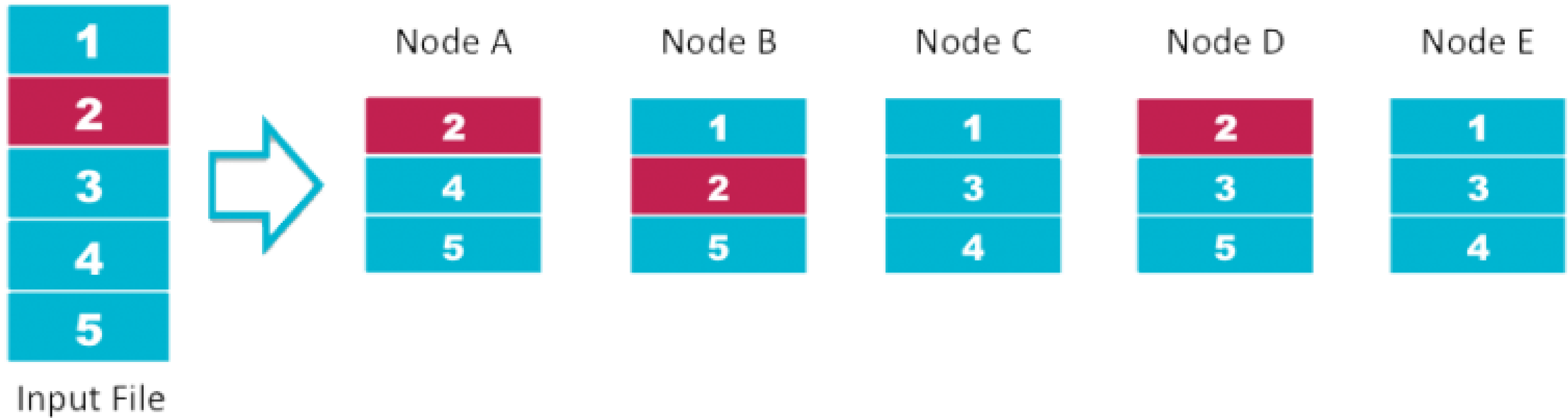




# High Level Architecture of Hadoop



# HDFS Data Distribution





**Hadoop**  
MapReduce

VS

**Spark**

# Infographic



Speed

Spark

100x times than MapReduce

Hadoop

Faster than traditional system

Compact & easier than Hadoop

Spark



Easy of Use

Complex & lengthy process

Hadoop



Data Storage

Spark

Store data in-memory

Hadoop

Store data on disk

Scala

Spark



Written in

Java

Hadoop



Data Processing

Spark

Batch / real-time / iterative /  
interactive / graph

Hadoop

Batch processing

RDD provide the fault tolerant

Spark



Fault  
Tolerance

Hadoop is also fault tolerant

Hadoop



Machine  
Learning

Spark

MLlib on top of Spark provides ML algo  
with lightening fast speed

Hadoop

Mahout on top of Hadoop provides  
ML algo

Cluster of 8000 nodes operational

Spark



Scalability

Hadoop has been tested on 15000 nodes

Hadoop



Latency

Spark

Provides low-latency computing

Hadoop

A high latency computing framework

Process real time data stream

Spark



Streaming

Process data only in batch mode

Hadoop



Caching

Spark

Caches the data in-memory &  
enhances the system performance

Hadoop

Doesn't support caching of data

Higher than Hadoop

Spark



Cost

Low cost

Hadoop

# Big data benchmark

Measuring  
and  
comparing big  
data systems  
and  
architecture

Being data-  
centric

Diverse and  
representative  
workloads

Covering  
representative  
software  
stacks

State-of-the-  
art techniques

Usability

**Table 3** The Summary of BigDataBench. Adaptation With the Permission of Wang *et al.* [78]

Application Scenarios	Application Type	Workloads	Data Types	Data Source	Software Stacks
Micro Benchmarks	Offline Analytics	Sort	Unstructured	Text	Hadoop, Spark, MPI
		Grep			
		WordCount		Graph	
		BFS			
Basic Datastore Operations	Online Service	Read	Semi-structured	Table	Hbase, Cassandra, MongoDB, MySQL
		Write			
		Scan			
Relational Query	Realtime Analytics	Select Query	Structured	Table	Impala, MySQL, Hive, Shark
		Aggregate Query			
		Join Query			
Search Engine	Online Services	Nutch Server	Un-structured	Text	Hadoop
	Offline Analytics	Index		Graph	Hadoop, Spark, MPI
		PageRank			
Social Network	Online Services	Olio Server	Un-structured	Graph	Apache+MySQL
	Offline Analytics	Kmeans			Hadoop, Spark, MPI
		Connected Components			
E-commerce	Online Services	Rubis Server	Structured	Table	Apache+JBoss+MySQL
	Offline Analytics	Collaborative Filtering	Semi-structured	Text	Hadoop, Spark, MPI
		Naive Bayes			

## Conclusion

Big data-processing systems have been widely researched by academia and industry. Based on the processing paradigm, we categorize those systems into batch, stream, graph, and machine learning processing. The paper first introduced the basic framework of MapReduce and its outstanding features as well as deficiencies compared to DBMSs. According to the deficiencies, we discussed the extensions and optimizations for MapReduce platform, including support for flexible dataflows, efficient data access and communication, parameter tuning, as well as energy. We then surveyed other batch-processing systems, including general-purpose systems Dryad, Nephelē/PACT and Spark. SQL-like systems involved in this paper are Hive, Shark, SCOPE, AsterixDB, and Dremel. For stream processing systems, Storm and S4 are introduced as representatives. Scalability is one of the ML algorithms bottlenecks. We then discussed how graph-centric systems like Pregel and GraphLab, and ML-centric systems like Petuum, parallelize the graph and ML model, as well as their distinctive characteristics.



1. <https://ieeexplore.ieee.org/document/7547948>
2. <https://ieeexplore.ieee.org/document/9378386>
3. [https://doi.org/10.1007/978-3-319-63962-8\\_165-1](https://doi.org/10.1007/978-3-319-63962-8_165-1)

# REFERENCE



Thank  
you

