# Apache Hadoop - A Herd of Elephants

Navigating the Hadoop Ecosystem

**Vinayak Hegde**

**Inmobi**

# The Data Revolution

- Social networks
- System generated data
- Financial Data
- Open Data – Government Data
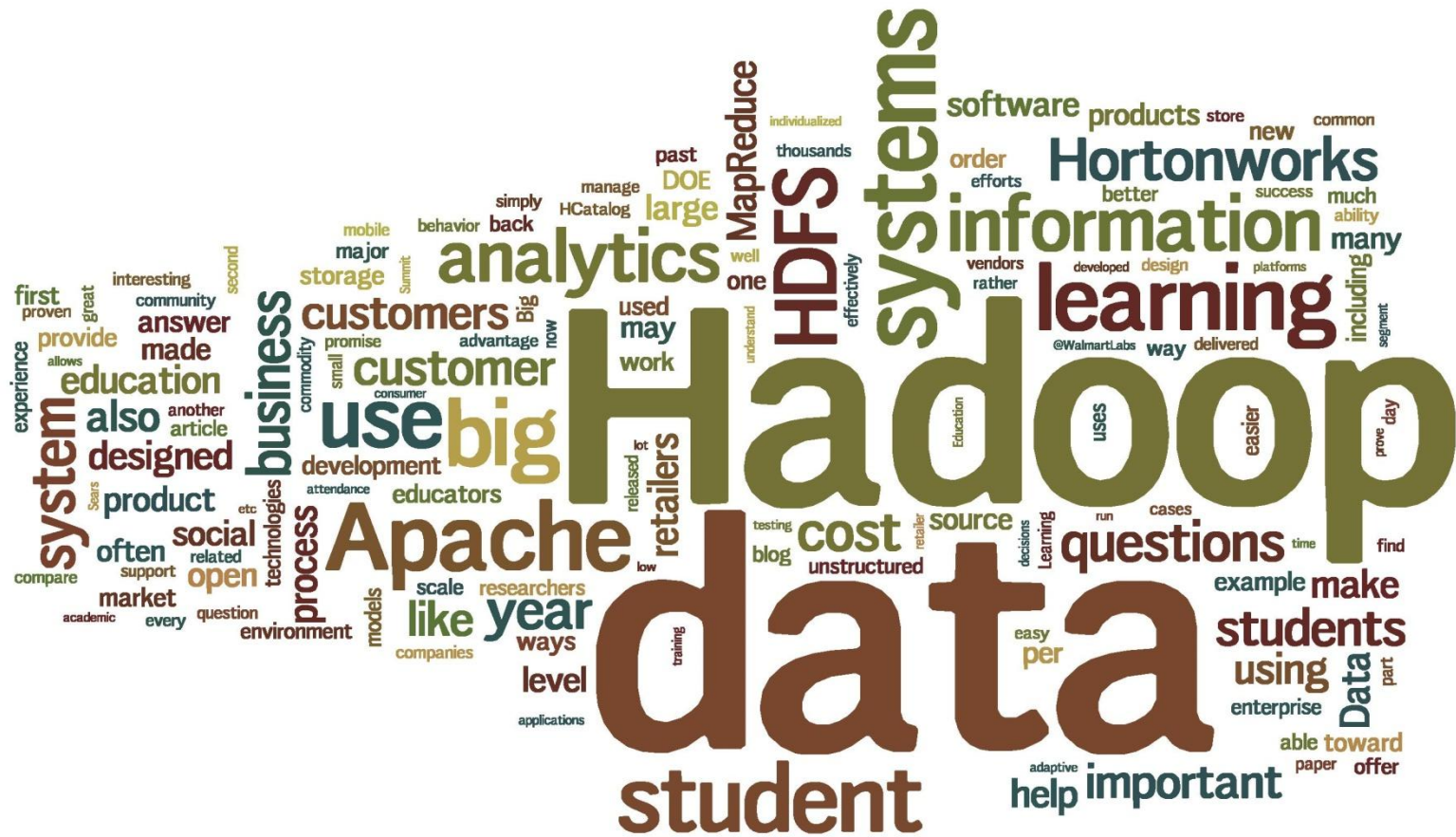- Audio data
- High resolution videos

# The Data Stack

Data Visualization

Data Analysis

Data Processing

Data Storage

# Hadoop – The Basics

- **HDFS**

  HDFS is a filesystem designed for storing very large files with streaming data access patterns, running on clusters of commodity hardware.

A very short introduction to HDFS architecture

- Namenode
- Datanode
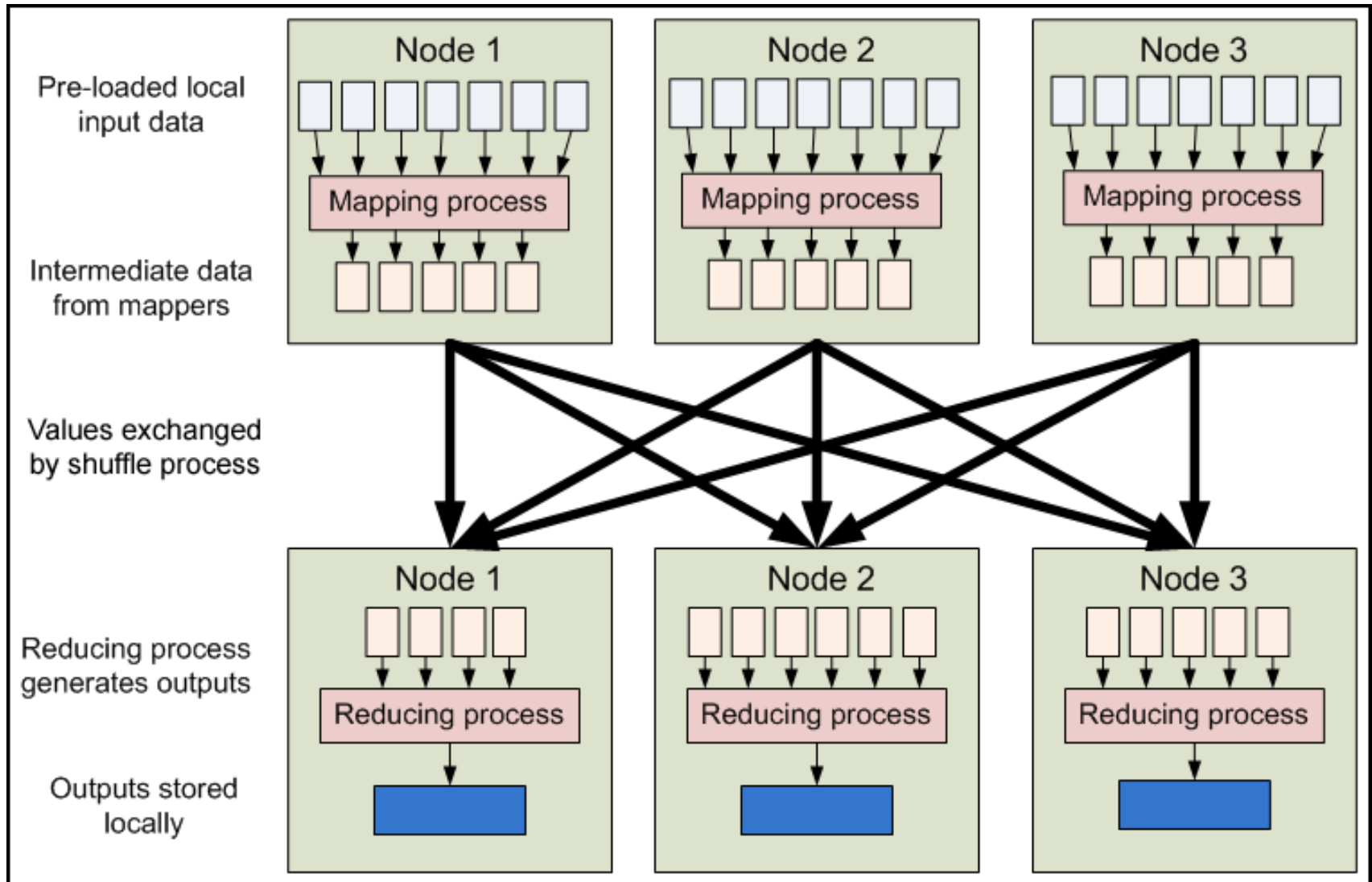- Replication factor
- Rackaware-ness

# Hadoop – The Basics

- **MapReduce**

  Mapreduce is a paradigm based on two functions – Map and Reduce. This is the basis of the programming model of Hadoop.

  Demo : A example in Python

# Map Reduce Flow

# Data Storage

- **Data Storage**
  - HDFS
  - S3 (native, Block based)
  - HSFTP (HDFS over HTTP)
  - KFS (Kosmos File System), Ceph
- **Data Serialization / Archiving**
  - Avro
  - Protobufs
  - Thrift
  - RCFile , HAr ,HCatalog

# Data Storage

- **Data Movement /Transport**
  - Sqoop
  - Flume
  - Scribe
  - Kafka
  - Message queues ?
- **Columnar Storage**
  - Zebra
- **KV Store**
  - HBase

# Data Processing

- **Schedulers**
  - Azkhaban
  - Oozie
  - Ivory (A data processing framework)
- **Cluster monitoring**
  - Ganglia + Nagios
  - Chukwa
  - Zookeeper (Distributed Coordination service)

# Data Analysis

- Mahout

- Piggybank

- Hive

- Pegasus

- Giraph

- Pig

- AllReduce (Vowpal Rabbit)

- Mallet

# Data Visualization

- Ambrose

- Tableau

- GWT

- D3 / infovis

- R / Python – matplotlib / Octave

- Gephi

# Pig – A dataflow language

- Pig is made up of two components:
  - The dataflow language called *Pig Latin*.
  - The execution environment
- Good for fast prototyping & ease of use
- UDFs (User defined Functions)

# Pig – A dataflow language

**input_lines = LOAD '/tmp/my-copy-of-all-pages-on-internet' AS (line:chararray);**

--Extract words from each line and put them into a pig bag datatype,

-- then flatten the bag to get one word on eachrow

**words = FOREACH input_lines GENERATE FLATTEN(TOKENIZE(line)) AS word;**

-- filter out any words that are just white spaces

**filtered_words = FILTER words BY word MATCHES '\\w+';**

-- create a group for each word

**word_groups = GROUP filtered_words BY word;**

-- count the entries in each group

**word_count = FOREACH word_groups GENERATE COUNT(filtered_words) AS count, group AS word;**

-- order the records by count

**ordered_word_count = ORDER word_count BY count DESC;**

**STORE ordered_word_count INTO '/tmp/number-of-words-on-internet';**

# Hive – A Data Warehousing Framework

- Hive consists of 3 components
  - Hive Execution environment and Hive Shell
  - HQL – the SQL like structured language
  - Metastore
- Hive Query Language
  - Good for rapid prototyping structured temporal data
  - Gentle learning curve
  - Good aggregation functions built-in
  - Good support for sampling

# Hive – A Data Warehousing Framework

```
CREATE TABLE page_view
        (viewTime INT, userid BIGINT, page_url STRING,
        referrer_url STRING, ip STRING
COMMENT 'IP Address of the User')
COMMENT 'This is the page view table'
PARTITIONED BY(dt STRING, country STRING)
CLUSTERED BY(userid) SORTED BY(viewTime) INTO 32 BUCKETS
ROW FORMAT DELIMITED
        FIELDS TERMINATED BY '1'
        COLLECTION ITEMS TERMINATED BY '2' MAP KEYS
        TERMINATED BY '3'
STORED AS SEQUENCEFILE;
```

# HBase – A Key/Value Store

- HBase is a distributed column-oriented database built on top of HDFS.

- Based on BigTable

- Key Concepts
  - Column families
  - Regions
  - Locking

# HBase – A Key/Value Store

- No real indexes
- Automatic partitioning
- Scale linearly and automatically with new nodes
- Commodity hardware
- Fault tolerance
- Batch processing

# Commercial vendors

- Cloudera
- Hortonworks
- Mapr
- IBM

# Thank you

- Email : vinayakh@gmail.com

- Twitter : @vinayakh

- Linkedin : www.linkedin.com/in/vinayakh

- Github : https://github.com/vinayakh