# Foundations of Big Data

An analysis of the milestones, technological shifts, and socio-economic cycles defining the history of big data.

Ghassem Tofighi

# Contents

## Overview

Big Data refers to extremely large, diverse, and rapidly growing collections of information that traditional tools cannot process efficiently. It has become essential across industries for deriving insights, improving decisions, reducing costs, and creating new opportunities.

This document explores the explosive growth of data, core characteristics (including the widely recognized 5 Vs), sources, benefits, challenges, and technologies. It focuses on Apache Hadoop as a foundational solution for distributed storage and processing, covering its architecture, HDFS, MapReduce, Hive, and Spark. These concepts remain key to understanding scalable data systems, even as modern cloud and lakehouse approaches build upon them.

The material draws from multiple educational and official sources to provide a complete, accurate overview suitable for learning distributed computing and big data principles.

## Big Data – Rapid Increase in Data Creation

The amount of information produced worldwide has expanded enormously thanks to new devices, communication channels, and platforms such as social networks.

From the start of human record-keeping until 2003, approximately 5 billion gigabytes (5 exabytes) of data existed. By 2011, this volume was being generated every two days. By 2013, the same quantity appeared roughly every ten minutes. This pace of growth has continued to rise significantly.

A large portion of this information holds potential usefulness when examined properly, yet much of it remains unprocessed.

## The 5 Vs of Big Data

Big Data is commonly characterized by five key attributes, often referred to as the 5 Vs. These help define why special approaches are needed.

| V | Description | Key Implications |
| --- | --- | --- |
| Volume | The sheer scale or quantity of data, often measured in terabytes, petabytes, or beyond | Requires scalable storage and processing infrastructure |
| Velocity | The speed at which data is generated, collected, and must be processed (real-time or near real-time in many cases) | Demands high-throughput ingestion and low-latency handling |
| Variety | The diversity of data types and formats (structured, semi-structured, unstructured) | Necessitates flexible tools that handle multiple forms without rigid schemas |
| Veracity | The quality, accuracy, trustworthiness, and reliability of the data (including potential noise, inconsistencies, or uncertainty) | Involves cleaning, validation, and governance to ensure dependable insights |
| Value | The usefulness and actionable insights that can be extracted from the data | The ultimate goal—turning raw information into business or operational benefits |

These five dimensions highlight the challenges and opportunities in managing modern datasets.

## Understanding Big Data

Big Data consists of extremely large collections of information that exceed the processing abilities of standard computing methods and tools. It represents an entire field rather than one specific method, incorporating a range of approaches, software, and platforms.

## Categories and Sources Within Big Data

Big Data includes information created by many different systems and applications. Key areas include:

| Type of Source | Description | Typical Examples |
|---|---|---|
| Black Box Records | Captures audio from crew, device performance metrics | Airplanes, helicopters, jets |
| Social Networking Content | Opinions, posts, and interactions shared publicly | Facebook, Twitter/X |
| Stock Trading Information | Buy and sell orders for company shares | Stock exchanges worldwide |
| Electricity Grid Monitoring | Usage details at specific points relative to stations | Power distribution networks |
| Transportation Details | Vehicle specifications, routes, availability | Logistics systems, public transit |
| Search System Logs | Queries and retrieved content from various repositories | Major internet search engines |

The information falls into three categories:

- **Structured** — Organized in relational tables
- **Semi-structured** — Tagged formats such as XML
- **Unstructured** — Documents, PDFs, text files, media, logs

## Advantages Gained from Big Data

Analyzing stored information provides several practical benefits:

- Marketing groups study responses to advertisements and campaigns using social platform data.
- Companies review consumer likes and opinions to guide manufacturing plans.
- Medical facilities examine past patient records to deliver faster, more effective treatment.

Big Data tools deliver precise evaluations, supporting better choices, improved efficiency, lower expenses, and decreased risks.

## Big Data Technologies – Operational and Analytical

Two main groups exist:

| Group | Focus | Examples | Key Traits |
|-------|-------|----------|------------|
| Operational | Real-time data handling and access | MongoDB, other NoSQL systems | Interactive workloads, cloud-friendly |
| Analytical | Deep, historical examination | MapReduce-based systems, MPP databases | Large-scale batch and complex queries |

These groups complement each other and are often deployed side by side.

## Difficulties Associated with Big Data

Organizations face several hurdles:

• Collecting incoming streams
• Maintaining data quality and organization
• Storing vast amounts economically
• Enabling fast searches
• Allowing secure sharing
• Transferring large volumes
• Conducting detailed analysis
• Presenting results clearly

Enterprise-grade servers typically help address these issues.

## Traditional Methods Versus Modern Solutions

### Conventional Single-System Method

An organization uses one powerful computer for storage and computation, often with a vendor database. Users interact through applications that manage storage and processing.

This setup functions adequately for modest data sizes fitting within standard server limits. However, when dealing with rapidly expanding, massive volumes, single-system constraints create severe bottlenecks.

### Google's Innovative Approach

Google developed MapReduce, an algorithm that breaks tasks into smaller pieces, distributes them across many machines, and combines outputs to form complete results.

### Development of Hadoop

Inspired by Google's work, Doug Cutting and collaborators created Hadoop, an open-source framework. Hadoop applies MapReduce principles to enable parallel processing of enormous datasets.

Hadoop supports applications performing full statistical evaluations on very large information collections.

## Introduction to Hadoop

Hadoop is an Apache open-source Java-based framework enabling distributed computation over large datasets on computer clusters. It operates in environments providing distributed storage and processing, scaling from one machine to thousands, each contributing local resources.

# Hadoop Architecture

Hadoop consists of two primary layers:

- **Processing Layer** — MapReduce for computation
- **Storage Layer** — Hadoop Distributed File System (HDFS)

Additional modules include:

- **Hadoop Common** — Shared Java libraries and utilities
- **YARN** — Job scheduling and cluster resource management
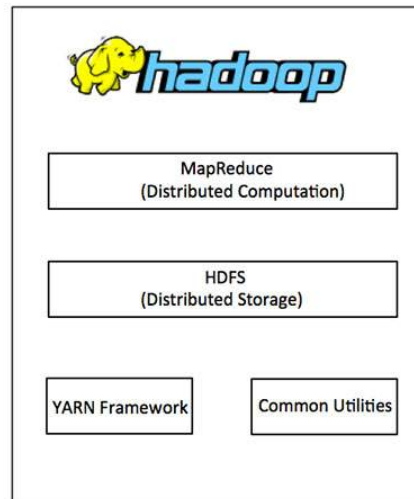
Figure 1 provides an overview of the ecosystem.



Figure 1:  Hadoop Ecosystem Architecture

**Figure 1: Overview of Hadoop components and layers**

# Operation of Hadoop

Building large high-end servers is costly. Instead, Hadoop connects many inexpensive computers into a unified distributed system. This approach reads data in parallel, achieving superior throughput at lower cost.

Core operations:

- Files split into directories and uniform blocks (typically 128 MB or 64 MB, often 128 MB preferred)
- Blocks distributed to cluster nodes
- HDFS oversees processing atop local filesystems
- Blocks replicated to protect against hardware issues
- Ensures correct task completion
- Manages sorting between map and reduce phases
- Routes sorted data appropriately
- Records debugging information per job

# Benefits of Hadoop

- Enables rapid development and testing of distributed applications

- Automatically spreads data and workload, leveraging CPU parallelism
- Provides fault tolerance and availability through software, not hardware
- Supports dynamic addition or removal of servers without stopping operations
- Open-source and compatible across platforms due to Java foundation

# Overview of HDFS

HDFS follows distributed filesystem principles and operates on commodity hardware. Unlike many similar systems, it emphasizes high fault tolerance and suitability for low-cost setups.

HDFS manages very large datasets with easy access. Files spread across machines with redundancy to prevent loss from failures. It supports parallel application access.

## Features of HDFS

- Designed for distributed storage and processing
- Offers command-line interaction
- Includes built-in NameNode and DataNode servers for cluster status checks
- Provides streaming data access
- Implements file permissions and authentication

## HDFS Architecture

HDFS uses a master-slave model with these elements:

- **NameNode** — Master server running on commodity hardware with GNU/Linux and NameNode software. Manages namespace, client access, and operations like renaming, opening, closing files/directories.
- **DataNode** — One per cluster node, running on commodity hardware with GNU/Linux and DataNode software. Manages local storage, performs read/write as requested, creates/deletes/replicates blocks per NameNode instructions.
- **Block** — User data divided into segments stored on DataNodes. Minimum read/write unit (default 64 MB, configurable, often increased to 128 MB).
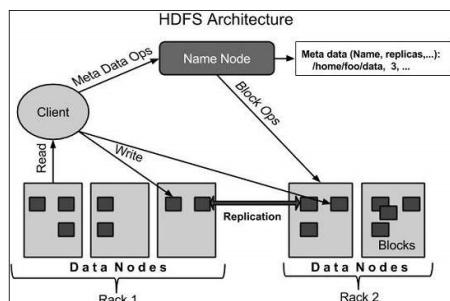
Figure 2 shows this structure.



Figure 2:  HDFS Architecture – NameNode and DataNodes

**Figure 2: Master-slave design of HDFS**

## Goals of HDFS

- Quick, automatic fault detection and recovery due to frequent component issues on commodity hardware

- Support for hundreds of nodes per cluster handling huge datasets
- Perform computation near data to reduce network traffic and boost throughput

# MapReduce – Processing Framework

MapReduce enables reliable parallel processing of large datasets on commodity clusters.

It is a Java-based distributed computing model with two main phases:

- **Map** — Breaks input into key/value tuples
- **Reduce** — Aggregates tuples into summarized output

The name reflects the sequence: map first, then reduce.

## Advantages

- Scales easily across many machines with minimal configuration changes
- Simplifies writing parallel applications

## Execution Stages

- **Map Stage** — Processes input line-by-line, produces intermediate chunks
- **Shuffle Stage** — Groups and sorts intermediate data
- **Reduce Stage** — Combines grouped data into final output

Hadoop distributes map and reduce tasks, manages data movement, verifies completion, minimizes network usage via local processing, and collects results.

## Input/Output Perspective (Java)

Jobs operate on `<key, value>` pairs:

(Input) `<k1, v1>` → map → `<k2, v2>` → reduce → `<k3, v3>` (Output)

Keys implement sorting; values and keys are serializable.

## Terminology

- **Payload** — Map and Reduce implementations
- **Mapper** — Produces intermediate pairs
- **NamedNode** — Manages HDFS
- **DataNode** — Holds data
- **MasterNode** — Runs JobTracker
- **SlaveNode** — Executes Map/Reduce
- **JobTracker** — Schedules and tracks
- **TaskTracker** — Monitors tasks
- **Job** — MapReduce execution on dataset
- **Task** — Map or Reduce on data slice
- **TaskAttempt** — Single execution attempt

## Example Scenario

Consider monthly electricity usage records across years or industries. Finding maximum/minimum usage years is simple with few records but challenging with statewide industrial data due to time, network load, etc.

MapReduce solves this by parallelizing across clusters.
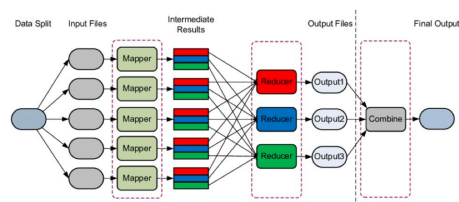
Figure 3 illustrates the process.



Figure 3: MapReduce Data Flow

**Figure 3: MapReduce workflow stages**

# Hive – Data Warehouse Tool on Hadoop

Big Data challenges traditional management, leading to Hadoop with modules like MapReduce and HDFS.

The ecosystem includes tools such as Sqoop (data transfer), Pig (procedural scripting), and Hive (SQL-like querying).

Hive processes structured data on Hadoop, summarizing large collections and simplifying queries.

Developed initially by Facebook, now Apache Hive, used by companies including Amazon Elastic MapReduce.

Hive is **not** a relational database, OLTP system, or real-time query tool.

## Features
- Stores schema in database, data in HDFS
- Designed for OLAP workloads
- Uses HiveQL (SQL-like) for queries
- Familiar, scalable, extensible

# Apache Spark – High-Speed Cluster Computing

Hadoop's MapReduce is scalable but slow for iterative or interactive tasks due to disk I/O.

Spark, from Apache, accelerates processing, extending MapReduce for interactive queries, streaming, etc.

Spark is not a Hadoop modification but can use Hadoop for storage while managing its own clusters.

It achieves speed through in-memory computing.

## Features
- Runs applications up to 100× faster in memory, 10× on disk
- Supports Java, Scala, Python APIs
- Includes 80+ high-level operators
- Handles batch, iterative, interactive, streaming workloads
- Supports SQL, machine learning, graph algorithms

## Built on Hadoop Options

- Standalone
- Hadoop YARN
- Spark in MapReduce (SIMR)
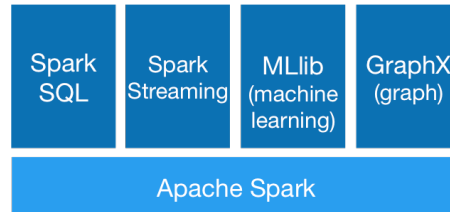
Figure 4 shows the component stack.



Figure 4:  Apache Spark Ecosystem

**Figure 4: Spark components and layered design**

## Spark Core and RDD

Spark Core provides the execution engine, supporting in-memory datasets (RDDs).

**RDD** — Immutable, partitioned, fault-tolerant collection computed in parallel. Created by parallelizing collections or referencing external data (HDFS, etc.).

RDDs enable faster operations than MapReduce by avoiding repeated disk writes.

## Spark SQL

Introduces DataFrame abstraction for structured processing, acting as distributed SQL engine.

Features:

- Mix SQL with Spark code
- Unified access to various sources
- Hive compatibility
- JDBC/ODBC connectivity
- Scalability with fault tolerance

Architecture layers: Language API, Schema RDD (DataFrame), Data Sources (Parquet, JSON, Hive, etc.).

# References

This document was compiled and rewritten based on multiple educational and official sources, including:

- Apache Hadoop Project Documentation (hadoop.apache.org)
- Apache Spark Official Documentation (spark.apache.org)
- TutorialsPoint Hadoop and Spark Tutorials (tutorialspoint.com)
- GeeksforGeeks Big Data and Hadoop Articles
- IBM and Oracle Big Data Overviews
- Wikipedia Entries on Big Data and Hadoop
- Edureka and ProjectPro Educational Guides
- TechTarget, Simplilearn, and other resources on the 5 Vs of Big Data

All content has been rephrased originally while faithfully covering every point from the provided source material.