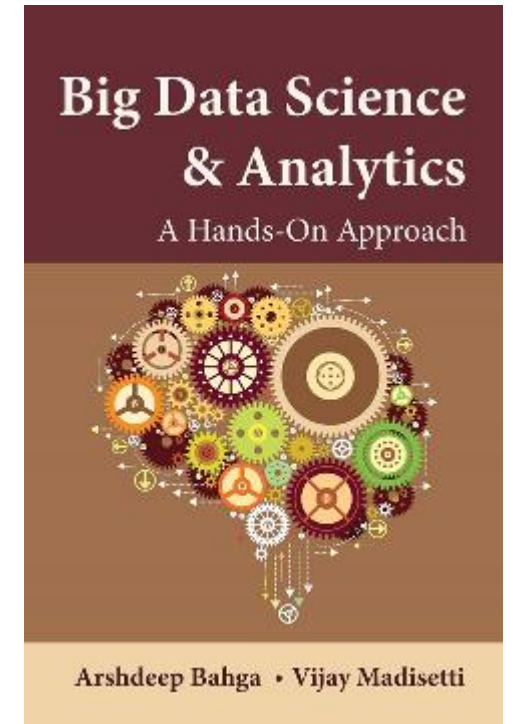# Big Data Storage

## Prof. Gheith Abandah

# Reference

- Chapter 6: **Big Data Storage**

- Arshdeep Bahga and Vijay Madisetti, **Big Data Science and Analytics: A Hands-On Approach**, 2019.
  - Web site: http://www.hands-on-books-series.com/

# Outline

- HDFS
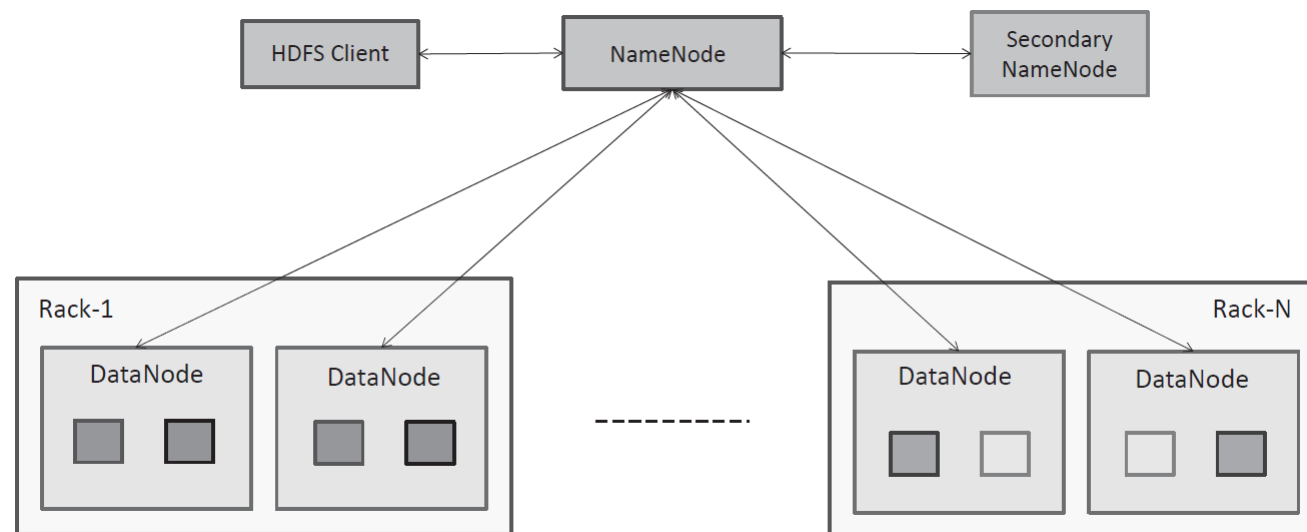- HDFS Architecture
- HDFS Usage Examples

# HDFS

- **Distributed file system** (DFS) from **Hadoop** framework that runs on **large clusters** and provides **high-throughput** access to data.

- Highly **fault-tolerant** system and is designed to work with **commodity hardware**.

- Stores each file as a sequence of **blocks**.

- The blocks of each file are **replicated** on multiple machines in a cluster to provide fault tolerance.

- MapReduce programs take advantage of **locality of data** and the data processing takes place on the nodes where the data resides.

# HDFS Characteristics

- **Scalable Storage for Large Files**: Large files are broken into chunks and each chunk is replicated across multiple machines in the cluster.

- **Replication**: The default block size used is 64MB and the default replication factor is 3.

- **Streaming Data Access**: is not suited for applications that require low-latency access to data; it provides high throughput data access.

- **File Appends**: HDFS was originally designed to have immutable files. Recent versions of HDFS have introduced the append capability.
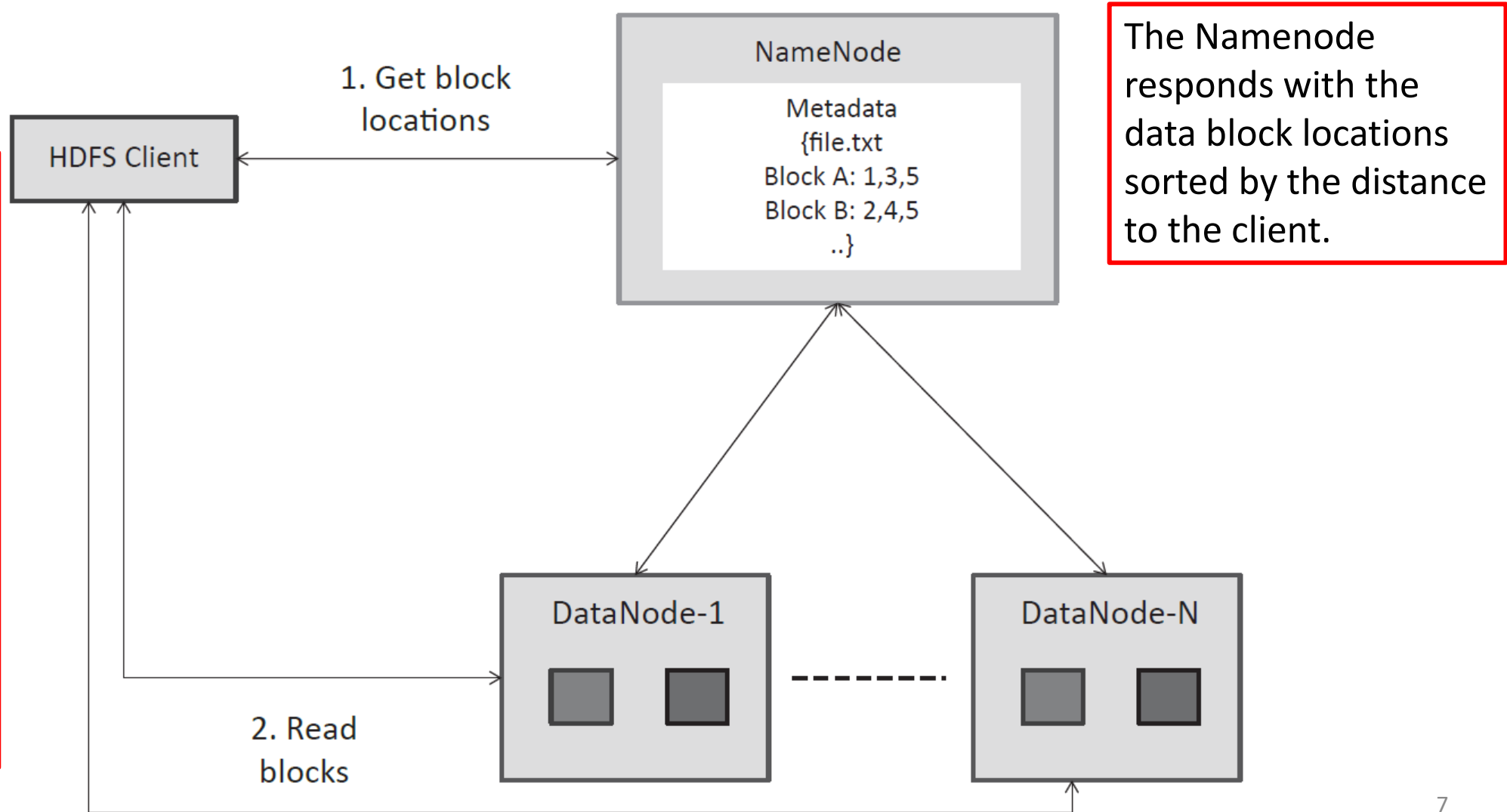
# HDFS Architecture

- **Namenode**: manages and stores the meta data and mappings of the blocks.

- **Secondary Namenode**: delegated to apply the mappings updates.

- **Datanodes** organized in racks and send heartbeats.

- **Replication**: one replica in the local rack and 2 in a remote rack.

# HDFS Read Path

The Datanodes stream the data to the client. During the read process, if a replica becomes unavailable, the client can read another replica on a different Datanode.

1. Get block locations

**NameNode**

Metadata
{file.txt
Block A: 1,3,5
Block B: 2,4,5
..}

The Namenode responds with the data block locations sorted by the distance to the client.

HDFS Client

2. Read blocks

**DataNode-1**

- - - - - - - - -

**DataNode-N**

# HDFS Write Path

**HDFS Client**

1. Create file →

8. Complete

**NameNode**

Metadata
{file.txt
Blocks: [],
...}

The Namenode checks permissions and responds with an output stream object.

The data packets consumed from the data queue are written to the first Datanode on the replication pipeline, which writes data to the second Datanode in the pipeline and so on.

7. Ack

**DataNode-1**

6. Ack

**DataNode-N**

5. Ack

**DataNode-N**

2. Write block

3. Write block

4. Write block

# HDFS Usage Examples

```
# Copy file to HDFS format
hdfs dfs -put <local src> <HDFS dest>
# Example:
hdfs dfs -put file /user/hadoop/file


# Get file from HDFS format
hdfs dfs -get <src on hdfs> <local dest>
# Example:
hdfs dfs -get /user/hadoop/file file
```

```
# List files on HDFS format
hdfs dfs -ls <args>
# Example:
hdfs dfs -ls /user/hadoop/


# Remove a file on HDFS format
hdfs dfs -rm <HDFS Path>
# Example:
hdfs dfs -rm /user/hadoop/file
```

# Accessing HDFS with Python

```python
from snakebite.client import Client
client = Client("localhost", 8020, use_trash=False)

# Listing files on HDFS with Python
list(client.ls(["/"]))

# Reading a file from HDFS with Python
list(client.text(["/user/hadoop/input.txt"]))

# Copying a file from HDFS with Python
list(client.copyToLocal(["/user/hadoop/input.txt"], '/home/ubuntu/'))
```

# Summary

- HDFS
- HDFS Architecture
- HDFS Usage Examples