# 實作課程 Lesson 02

蔡秉文

Cookie Tsai

# Big Data 的起源 Google 三篇論文

- The Google File System
  http://research.google.com/archive/gfs.html

- MapReduce: Simplified Data Processing on Large Clusters
  http://research.google.com/archive/mapreduce.html

- Bigtable: A Distributed Storage System for Structured Data
  http://research.google.com/archive/bigtable.html

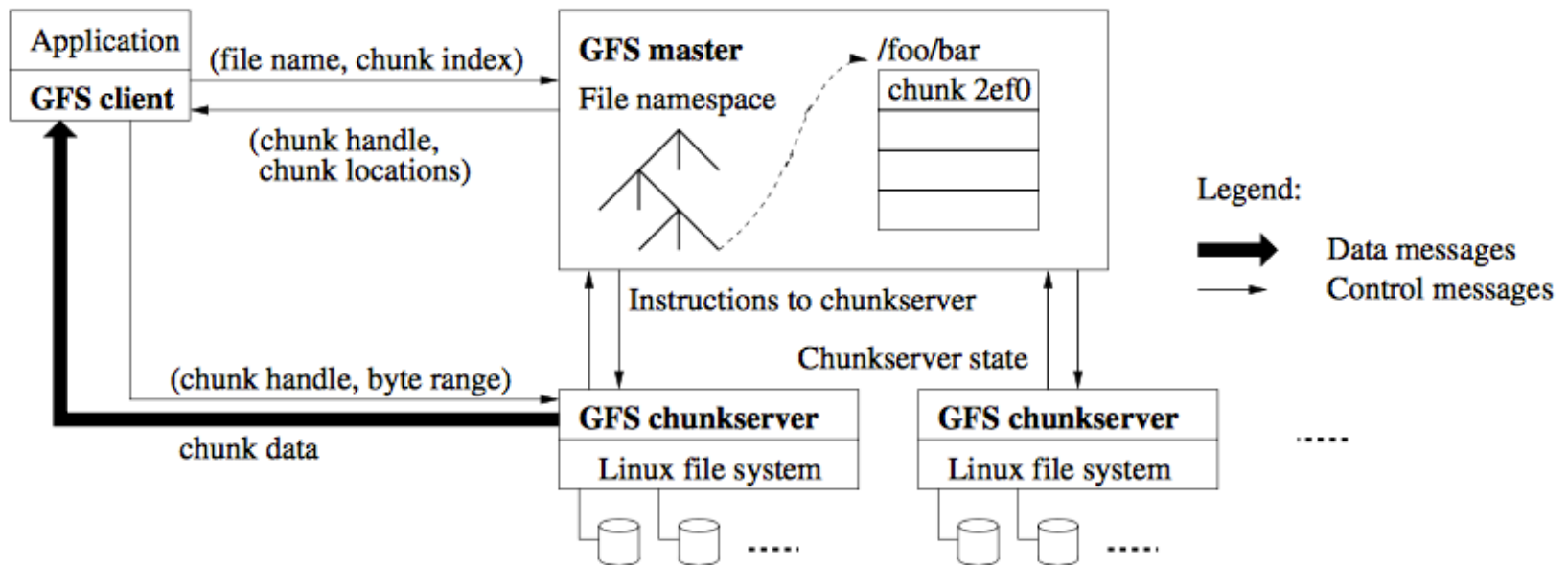# The Google File System Architecture



Figure 1: GFS Architecture

# What is Apache Hadoop

- ## Hadoop Common

  The common utilities that support the other Hadoop modules.

- ## Hadoop Distributed File System (HDFS™)

  A distributed file system that provides high-throughput access to application data.

- ## Hadoop YARN

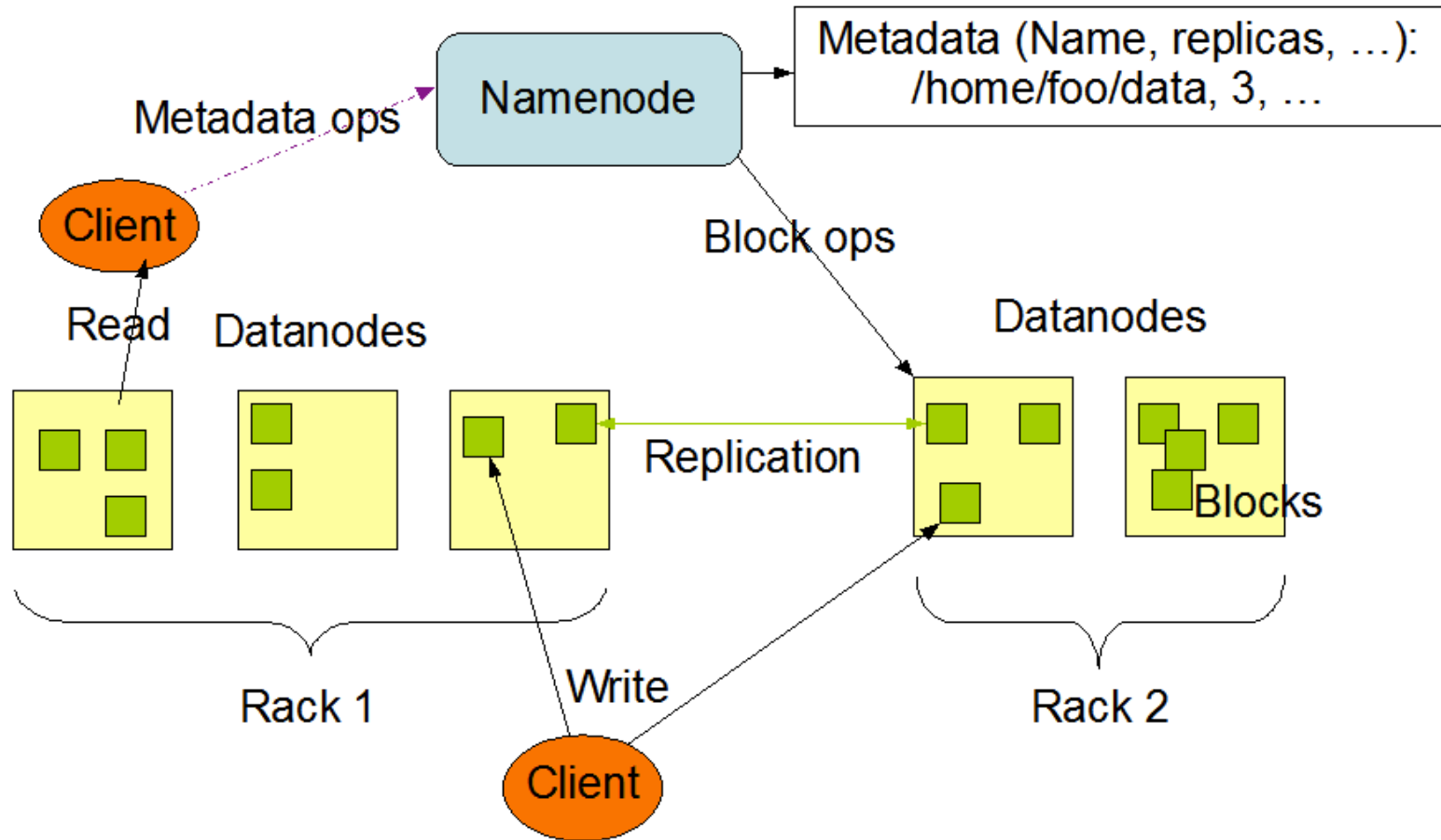  A framework for job scheduling and cluster resource management.

- ## Hadoop MapReduce

  A YARN-based system for parallel processing of large data sets.

# HDFS Overview

- HDFS is the primary distributed storage used by Hadoop applications.

- A HDFS cluster primarily consists of a NameNode that manages the file system metadata and DataNodes that store the actual data. The HDFS Architecture Guide describes HDFS in detail.

- This user guide primarily deals with the interaction of users and administrators with HDFS clusters.

- The HDFS architecture diagram depicts basic interactions among NameNode, the DataNodes, and the clients.

- Clients contact NameNode for file metadata or file modifications and perform actual file I/O directly with the DataNodes.
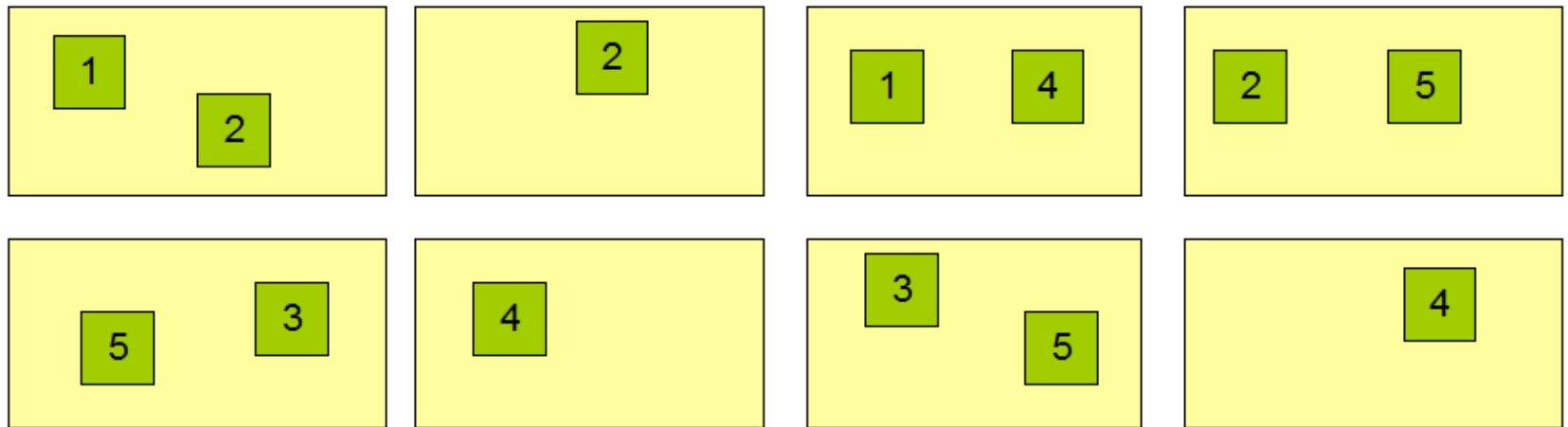
# HDFS Architecture

# Block Replication

Namenode (Filename, numReplicas, block-ids, …)
/users/sameerp/data/part-0, r:2, {1,3}, …
/users/sameerp/data/part-1, r:3, {2,4,5}, …

Datanodes

# HDFS Shell Commands

## Overview

The File System (FS) shell includes various shell-like commands that directly interact with the Hadoop Distributed File System (HDFS) as well as other file systems that Hadoop supports, such as Local FS, HFTP FS, S3 FS, and others. The FS shell is invoked by:

```
bin/hadoop fs <args>
```

All FS shell commands take path URIs as arguments. The URI format is `scheme://authority/path`. For HDFS the scheme is `hdfs`, and for the Local FS the scheme is `file`. The scheme and authority are optional. If not specified, the default scheme specified in the configuration is used. An HDFS file or directory such as /parent/child can be specified as `hdfs://namenodehost/parent/child` or simply as `/parent/child` (given that your configuration is set to point to `hdfs://namenodehost`).

Most of the commands in FS shell behave like corresponding Unix commands. Differences are described with each of the commands. Error information is sent to stderr and the output is sent to stdout.

# Reference

- HDFS Users Guide
  - [http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsUserGuide.html](http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsUserGuide.html)

- File System Shell Guide
  - [http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/FileSystemShell.html](http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/FileSystemShell.html)

- HDFS Architecture
  - [http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsDesign.html](http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsDesign.html)

Thank you for your listening