

Parallel and Distributed Computing

CS3006 (BCS-6C/6D)

Lecture 22

Instructor: Dr. Syed Mohammad Irteza

Assistant Professor, Department of Computer Science, FAST

18 April, 2023

Previous Lecture

- MPI
 - Odd-Even Sort (BubbleSort variant)
- Fault Handling
 - Fault Tolerance
 - Redundancy
 - HW
 - Passive (TMR)
 - Active (hot/cold sparing), eviction types
 - Hybrid
 - Information
 - Time

File Systems

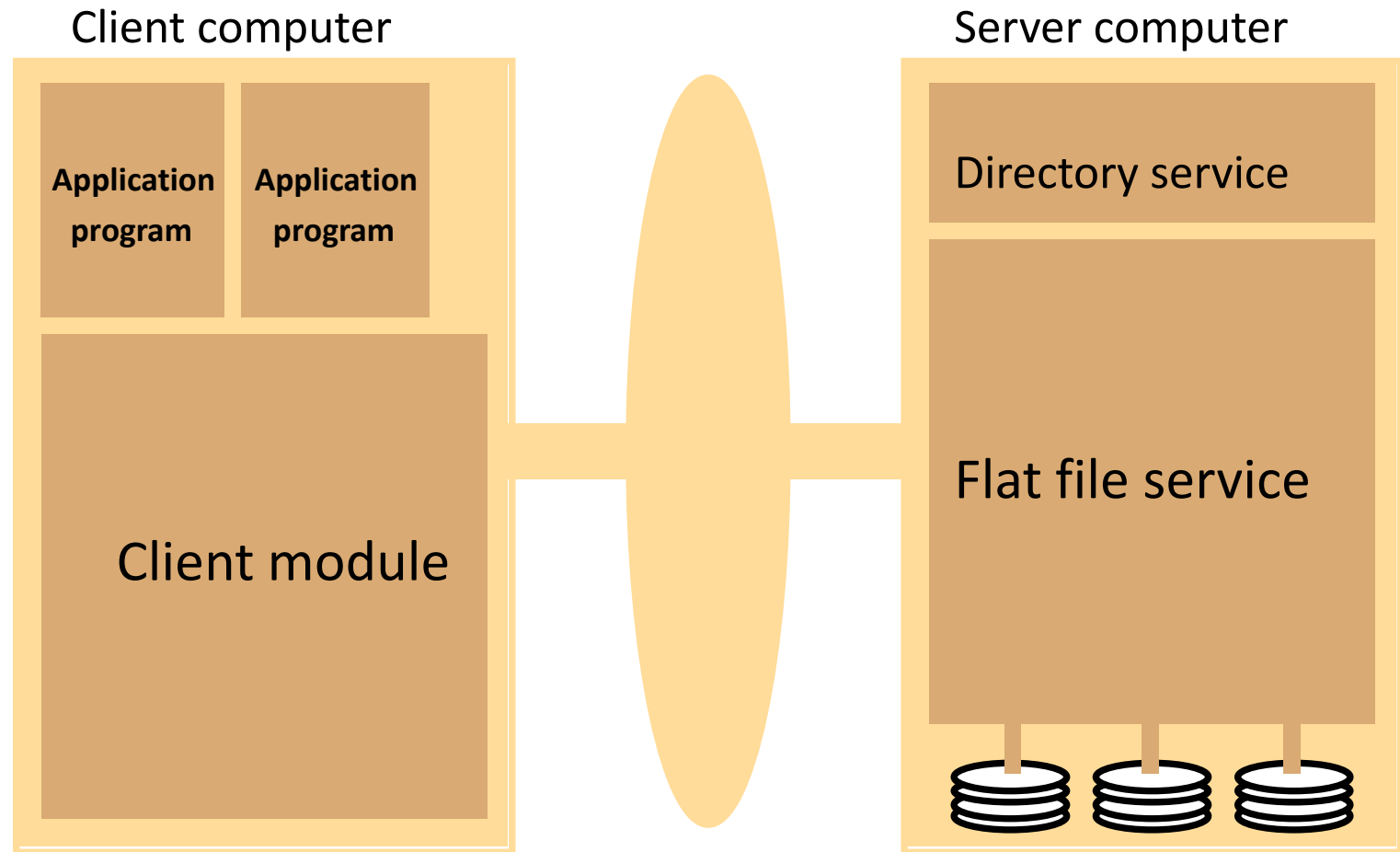
- File systems *were originally developed for centralized computer systems* and desktop computers.
- The file system was an operating system facility providing *a convenient programming interface to disk storage*.
- File systems are responsible for the organization, storage, retrieval, naming, sharing and protection of files.
- Files contain both *data* and *attributes*.

Distributed File Systems

- Distributed file systems *support the sharing of information* in the form of files and hardware resources.
- A DFS enables programs to *store and access remote files/storage* exactly as local ones.
- The *performance and reliability* of such access should be comparable to that for files stored locally.
- Recent advances in higher bandwidth connectivity of switched local networks and disk organization have led to higher performance and *highly scalable file systems*.
- *Functional requirements*: open, close, read, write, access control, directory organization, etc.
- *Non-functional requirements*: scalable, fault-tolerant, secure

General Distributed File Service Architecture

- An architecture that offers a *clear separation* of the main concerns in providing access to files is obtained by structuring the file service as three components:
 - A *flat file service*
 - A *directory service*
 - A *client module*



Distributed File Service Architecture

- The *client module* implements exported interfaces of *flat file* and *directory services* available on the server side.
- The responsibilities of the various *modules* can be defined as follows:
- *Flat file service:*
 - Concerned with the implementation of *operations on the contents of file*. Unique File Identifiers (**UFIDs**) are used to refer to files in all requests for flat file service operations. **UFIDs** are *long sequences of bits* chosen so that each file has a unique ID among all of the files in a distributed system.

Distributed File Service Architecture

- *Directory service*

- Provides mapping between *text names for the files* and their **UFIDs**. Clients may obtain the **UFID** of a file by quoting its text name to directory service. Directory service *supports functions needed to generate directories and to add new files to directories.*

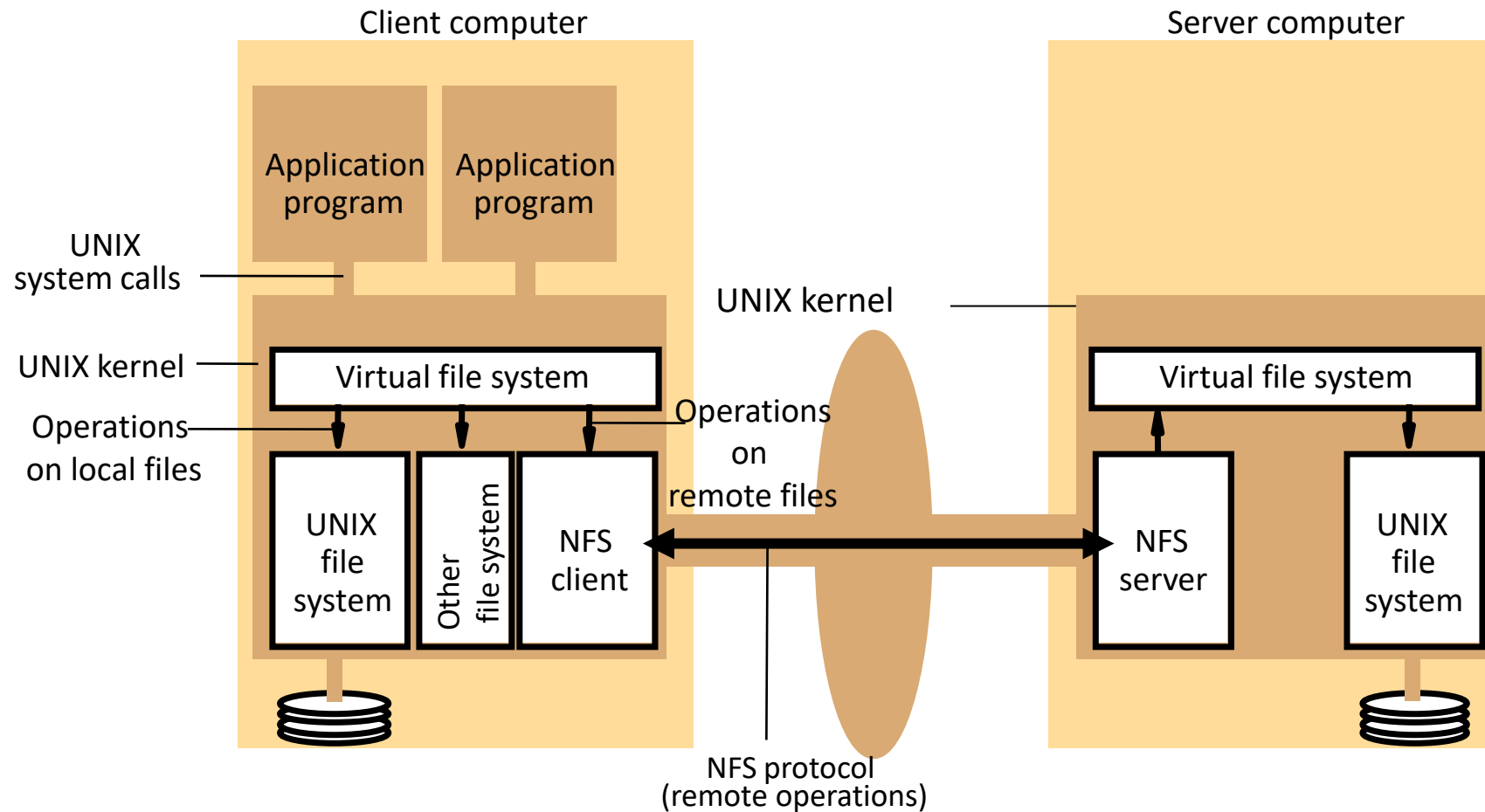
- *Client module*

- It runs on each computer and provides an integrated service (flat file and directory) as a *single API to application programs*. For example, in UNIX hosts, a client module emulates the full set of Unix file operations.
- It holds information about the *network locations of flat-file and directory server processes*; and achieves better performance through the implementation of a *cache of recently used file blocks at the client.*

Distributed File System Examples

- Network File System (NFS) ~ *Sun Microsystems*
- Andrew File System (AFS) ~ *Carnegie Mellon University*
- Sprite File System ~ Berkeley
- Google File System (GFS)
- Hadoop Distributed File System (HDFS)

NFS Architecture



NFS Architecture

- The NFS client and server modules communicate using *Remote Procedure Calls (RPCs)*.
- Sun's RPC system
 - *Sun RPC* was developed for use in *NFS*. It can be configured to use either *UDP* or *TCP*, and the *NFS protocol* is compatible with both

Virtual file system (VFS)

- **NFS provides access transparency through VFS**
 - user programs *can issue file operations for local or remote files without distinction*
 - *Other distributed file systems* may be present that support UNIX system calls, and if so, they *could be integrated via VFS*
 - Distinguishes between local and remote file identifiers
 - *Keeps track of the available file systems* – both local and remote
- The virtual file system layer has *one VFS structure for each mounted file system and one v-node per open file*
 - *A VFS structure relates a remote file system to the local* directory on which it is mounted
 - The *v-node contains an indicator to show whether a file is local or remote*. If the file is local, the v-node contains a reference to the index of the local file (an *i-node* in a UNIX implementation). If the *file is remote*, it contains the file handle of the remote file.

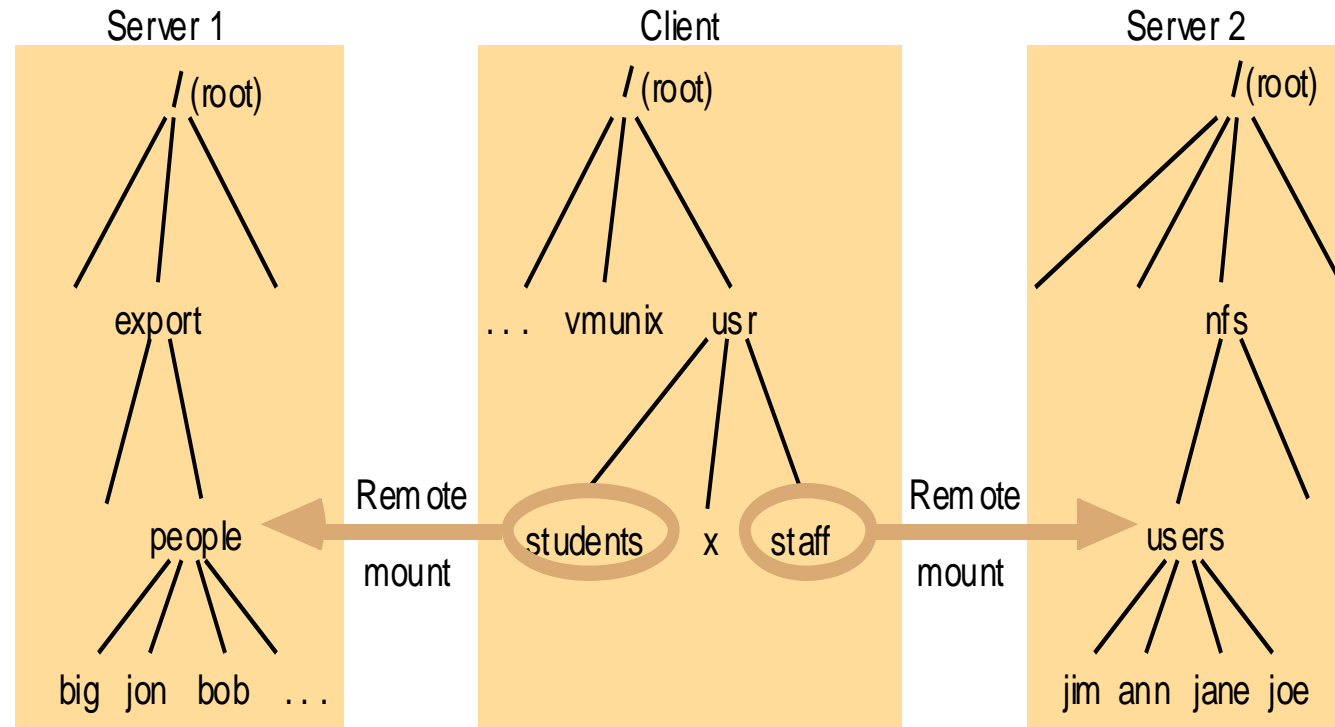
Hierarchic file system

- A hierarchic file system *consists of a number of directories arranged in a tree structure.*
- Any file or directory can be *referenced using a pathname* – a multi-part name
- A UNIX-like file-naming system *can be implemented by the client module using the flat file and directory services* that we have defined.
- A tree-structured network of directories is constructed with *files at the leaves and directories at the other nodes* of the tree. The root of the tree is a directory with a ‘well-known’ **UFID**.
- A function can be provided in the client module that gets the **UFID** of a file given its pathname. The *function interprets the pathname*.
 - Pathname starting from the root, using Lookup to obtain the UFID of each directory in the path.

Mount service

- The *mounting of subtrees of remote file systems* by clients is supported by a separate *mount service* process that runs at user level on each NFS server computer.
- On each server, there is a file with a well-known name (*/etc/exports*) containing the names of *local file systems* that are available for remote mounting.
- An *access list* is associated with each file system name indicating which hosts are permitted to mount the file system
- **Mount operation:**
 `mount(remotehost, remotedirectory, localdirectory)`
- Each client maintains *a table of mounted file systems* holding:
 < IP address, port number, file handle >

Mount Service

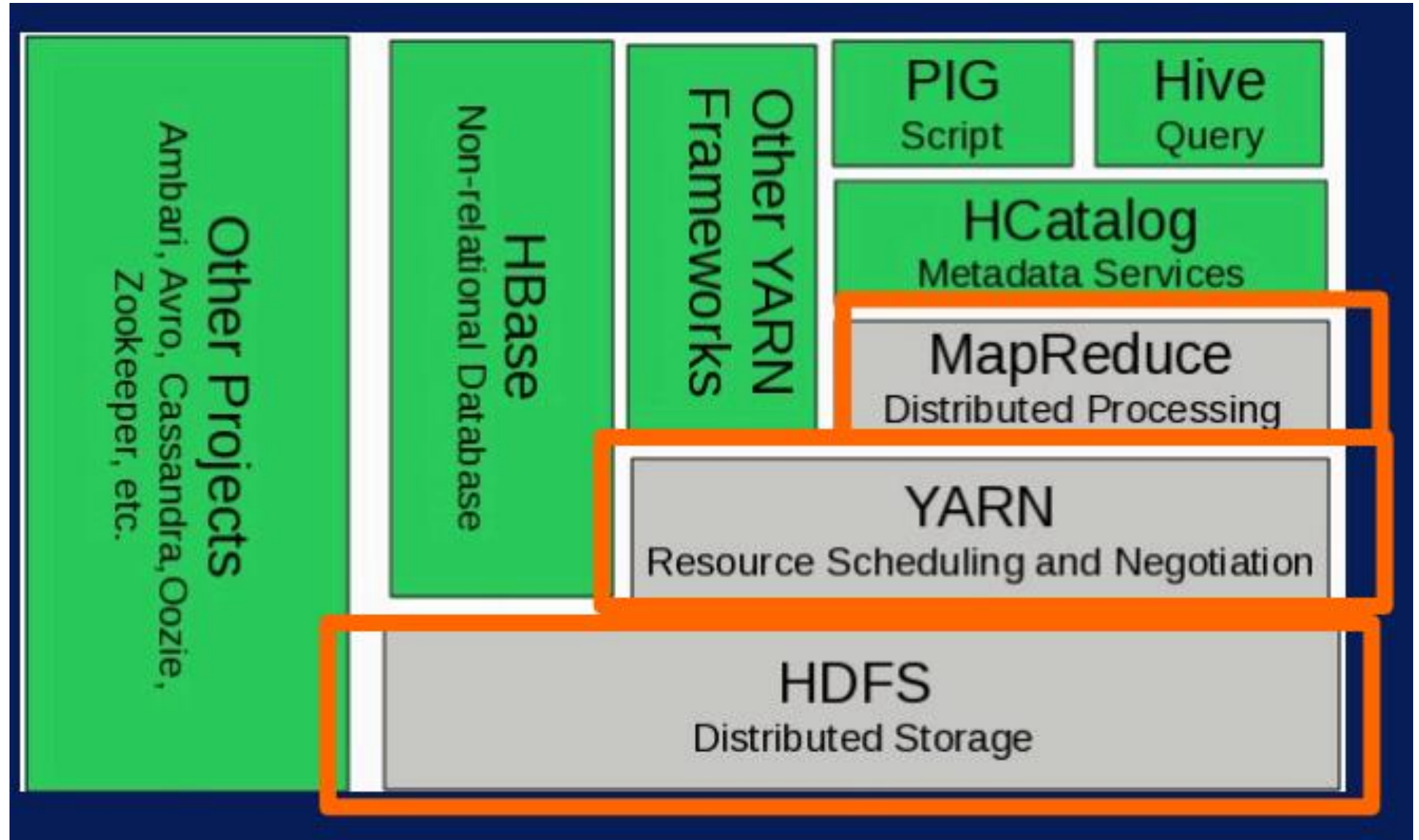


- The file system mounted at `/usr/students` in the client is actually the sub-tree located at `/export/people` in Server 1;
- the file system mounted at `/usr/staff` in the client is actually the sub-tree located at `/nfs/users` in Server 2.

Hadoop

- ***Apache Hadoop*** is an open source software framework for storage and large scale processing of data-sets on *clusters of commodity hardware*.
- It consists of the following basic modules:
 - *Hadoop Distributed File System (HDFS)*
 - *Hadoop YARN*
 - *Hadoop MapReduce*

Hadoop Module



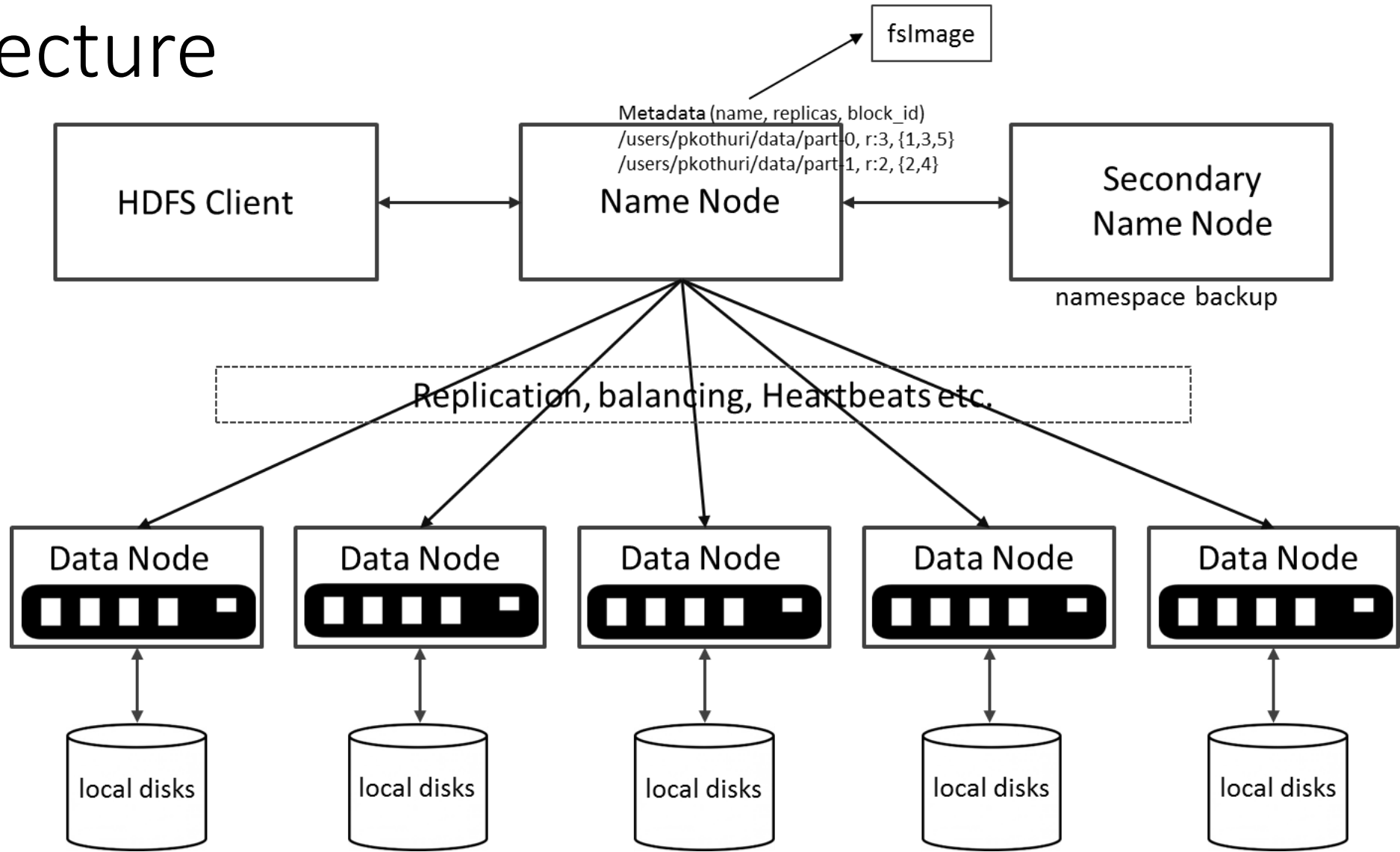
Hadoop Distributed File System

- **HDFS** is a distributed file system written in Java that is fault tolerant, scalable and extremely easy to expand.
- **HDFS** is the primary distributed storage for Hadoop applications.
- **HDFS** provides interfaces for applications to move themselves closer to data.

There are two types of machines in a **HDFS** cluster.

- NameNode is the heart of an **HDFS** filesystem, it maintains and manages the file system metadata. For example, what blocks make up a file, and on which DataNodes those blocks are stored.
- DataNode where **HDFS** stores the actual data, there are usually quite a few of these.

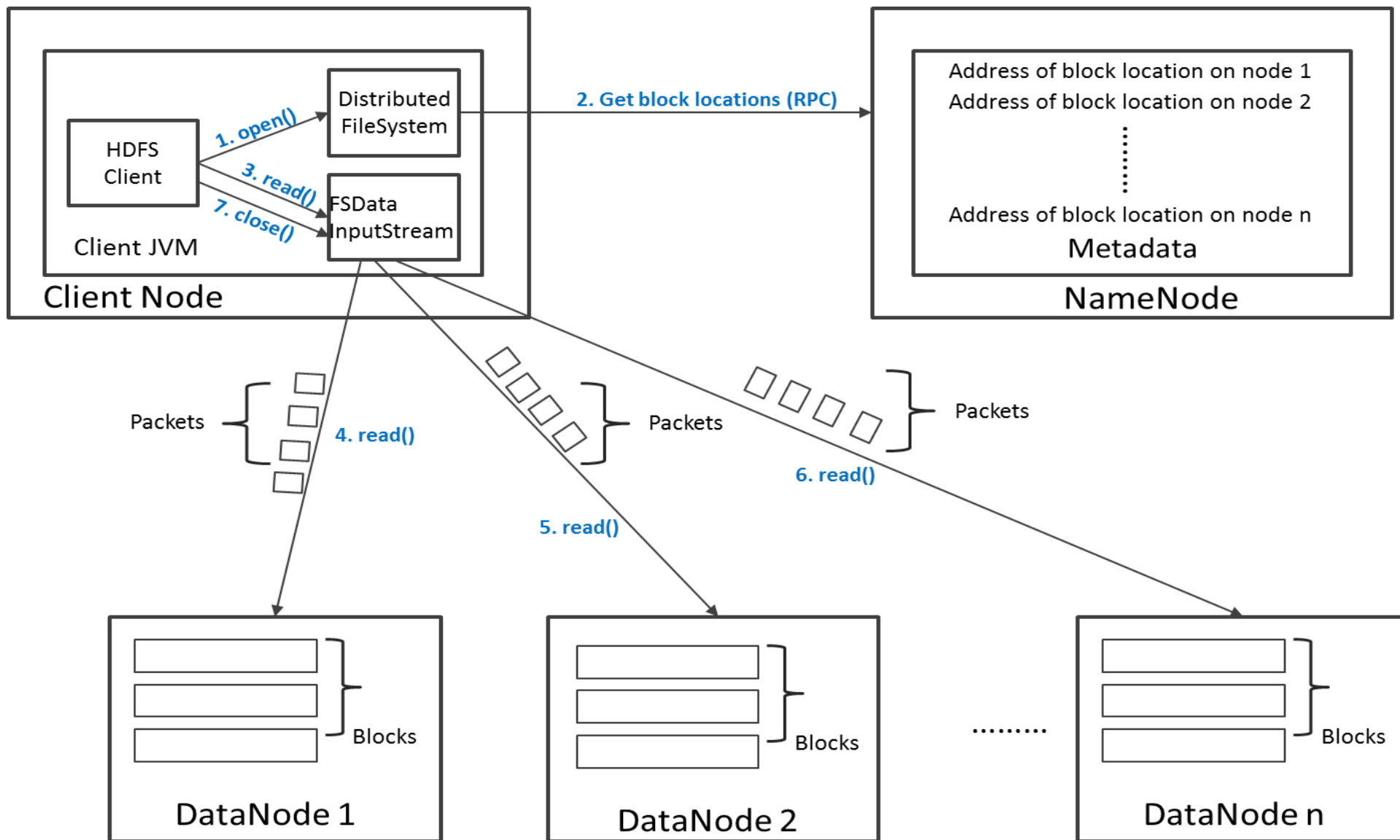
HDFS Architecture



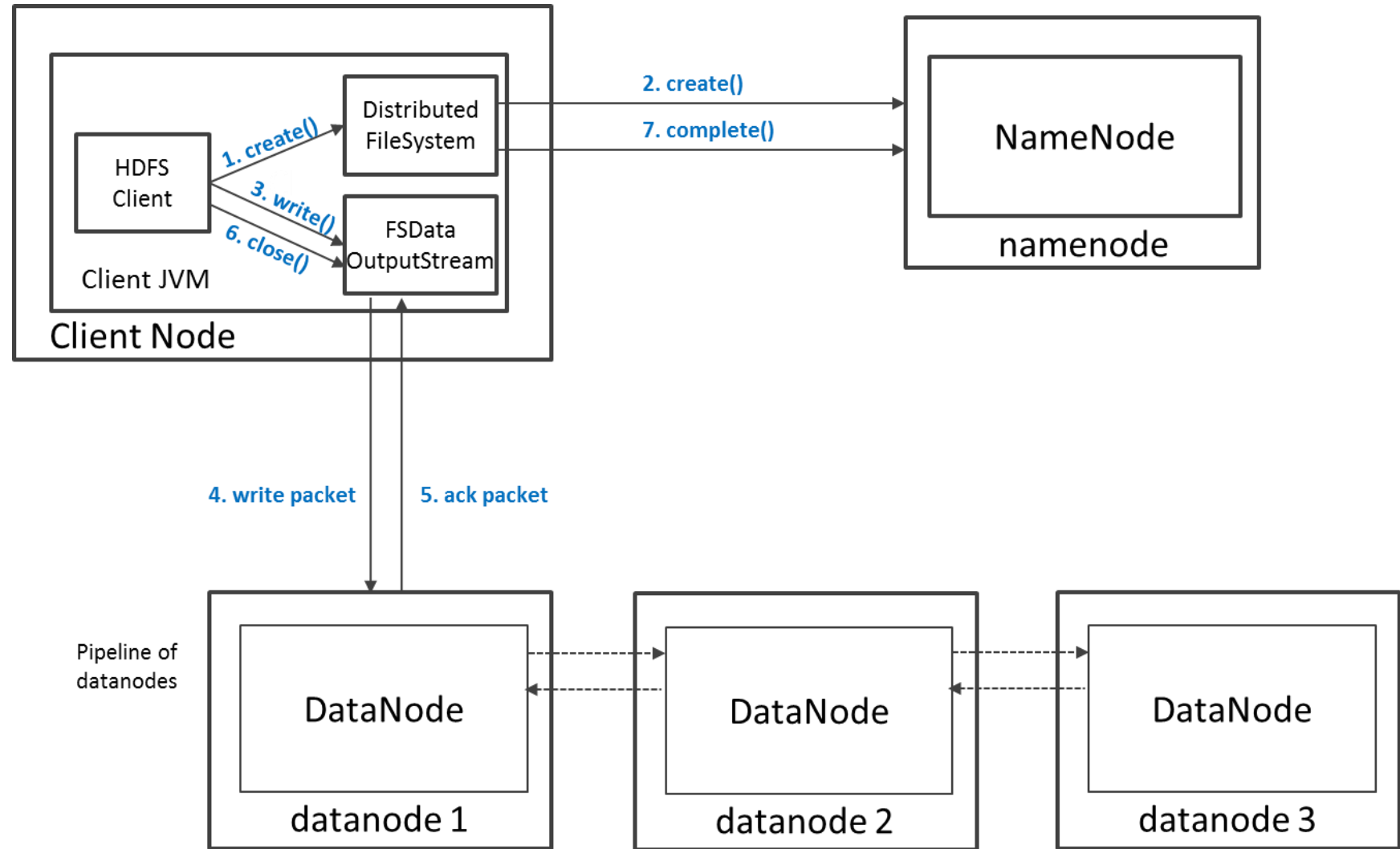
HDFS features

- Failure tolerant - data is duplicated across multiple DataNodes to protect against machine failures. The default is a replication factor of 3 (every block is stored on three machines).
- Scalability - data transfers happen directly with the DataNodes so your read/write capacity scales fairly well with the number of DataNodes
- Space - need more disk space? Just add more DataNodes and re-balance
- Industry standard - Other distributed applications are built on top of HDFS (HBase, Map-Reduce)

Read Operation in HDFS



Write Operation in HDFS



References

1. Slides of Dr. Rana Asif Rehman & Dr. Haroon Mahmood

Helpful Links:

1. <https://hadoop.apache.org/>
2. <https://www.cloudera.com/products/open-source/apache-hadoop.html>
3. <https://www.techtarget.com/searchenterprisedesktop/definition/Network-File-System>