

NETWORK DISTRIBUTED SYSTEMS

DISTRIBUTED FILE SYSTEMS

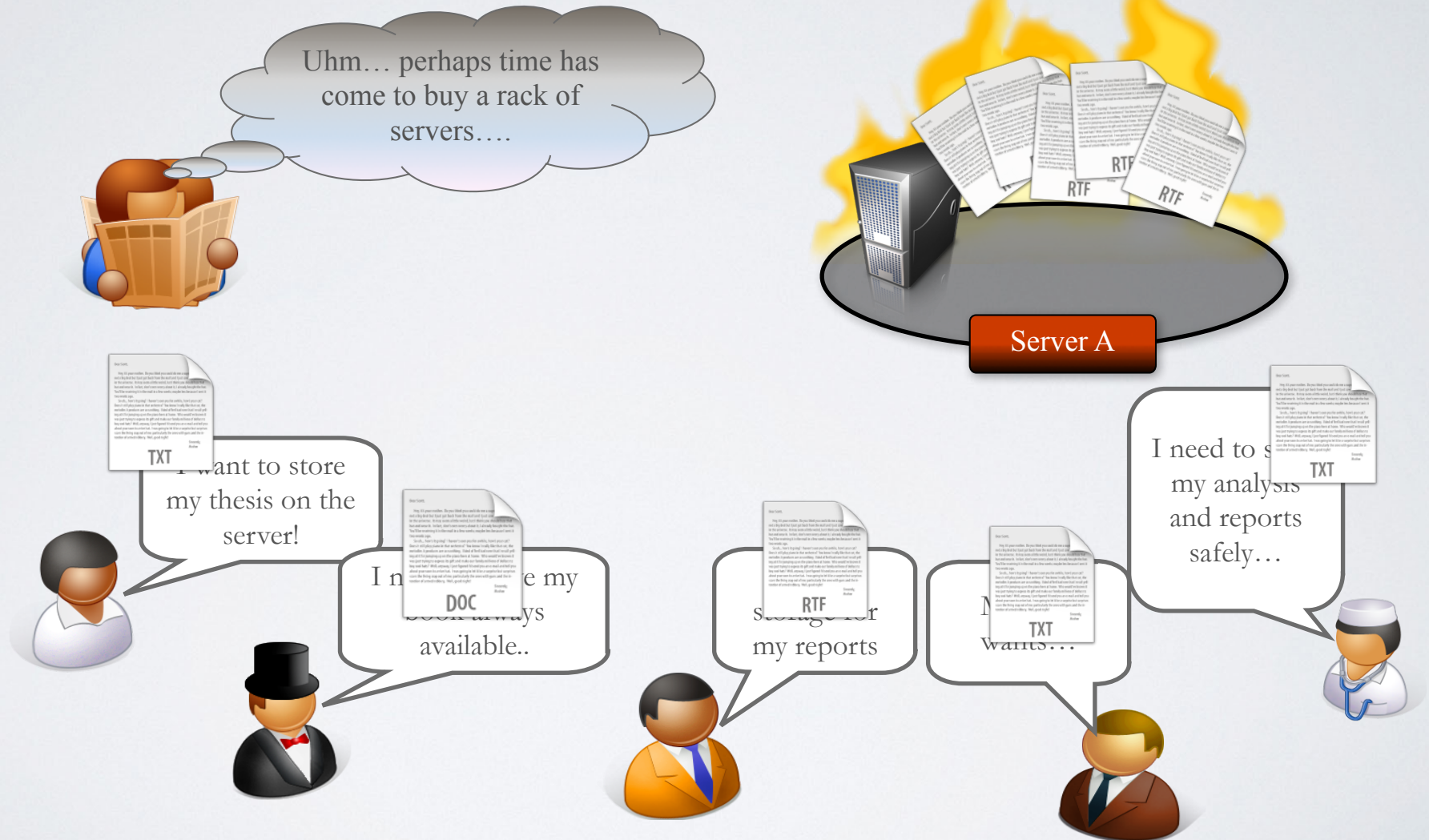
INTRODUCTION

- Why do we need a DFS?
 - Primary purpose of a Distributed System...

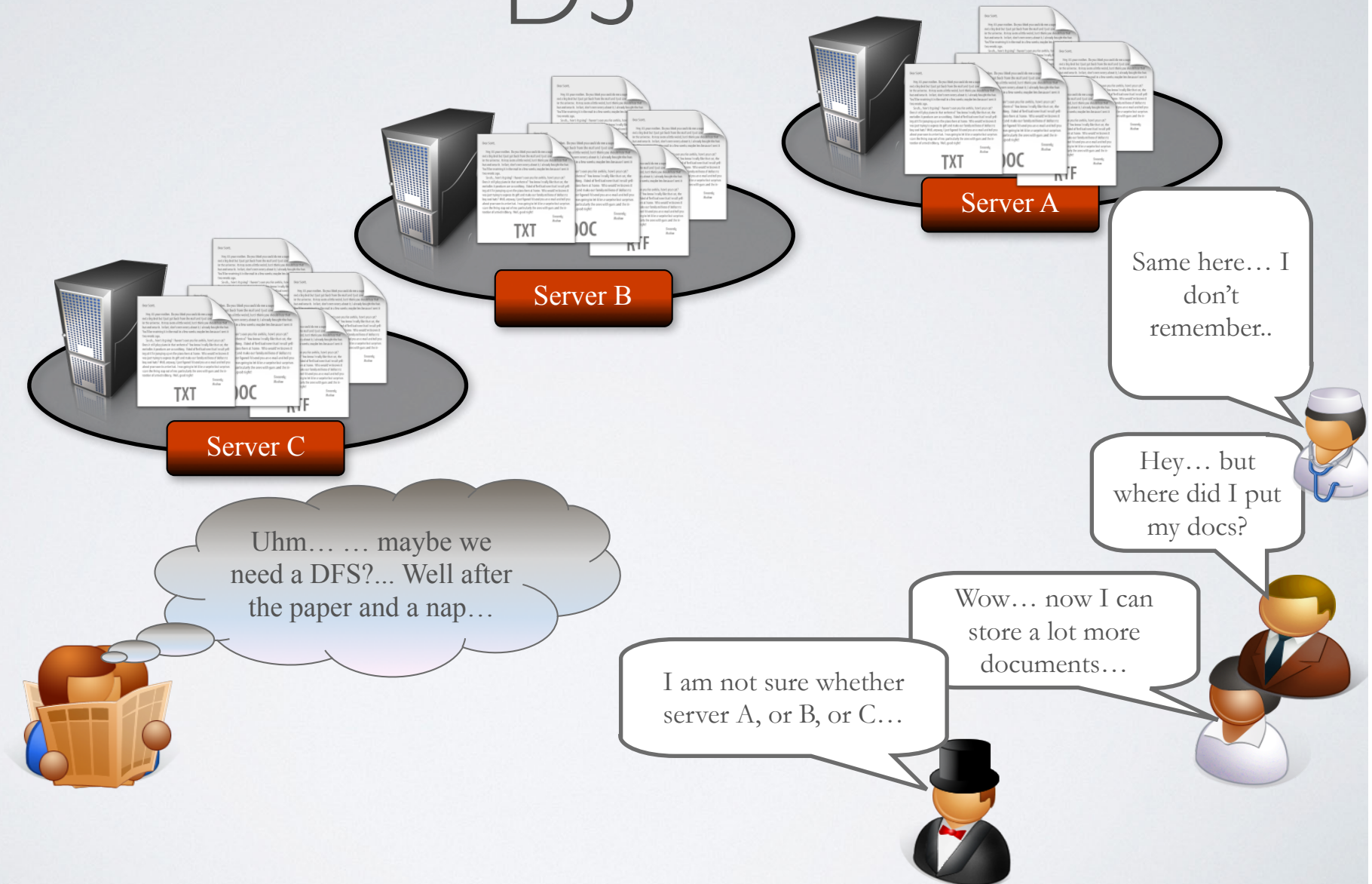
Connecting users and resources

- Resources...
 - ... can be inherently distributed
 - ... can actually be data (files, databases, ...) and...
 - ... their availability becomes a crucial issue for the performance of a Distributed System

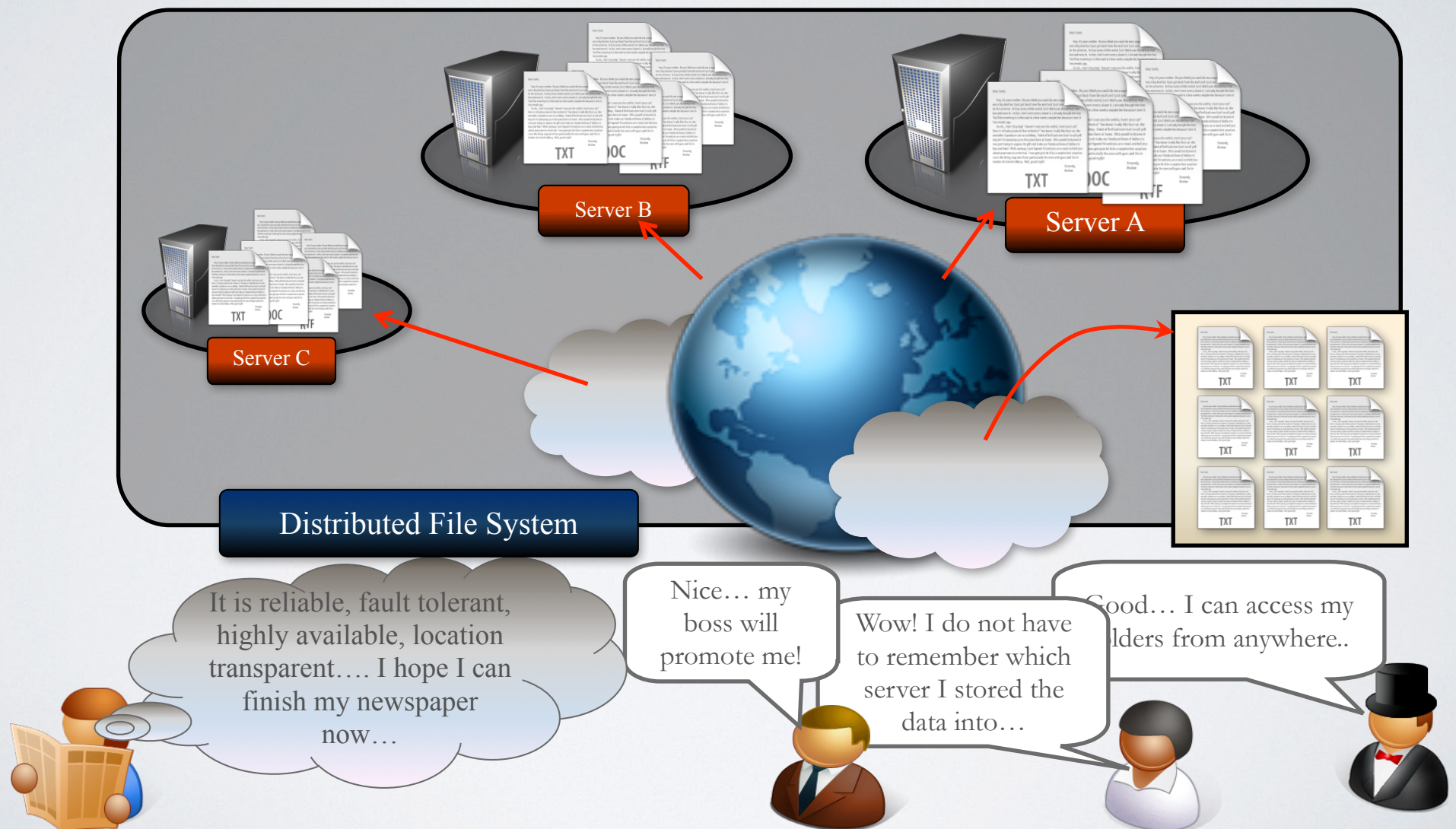
INTRODUCTION: A CASE FOR DFS



INTRODUCTION: A CASE FOR DS



INTRODUCTION: A CASE FOR DFS



STORAGE SYSTEMS AND THEIR PROPERTIES

- In first generation of distributed systems (1974-95), file systems (e.g. NFS) were the only networked storage systems.
- With the advent of distributed object systems (CORBA, Java) and the web, the picture has become more complex.
- Current focus is on large scale, scalable storage.
 - Google File System
 - Amazon S3
 - Cloud Storage (e.g., DropBox)

1974 - 1995

1995 - today

2007 - now



STORAGE SYSTEMS AND THEIR PROPERTIES

	<i>Sharing</i>	<i>Persistence</i>	<i>Distributed cache/replicas</i>	<i>Consistency maintenance</i>	<i>Example</i>
Main memory	×	×	×	1	RAM
File system	×	✓	×	1	UNIX file system
Distributed file system	✓	✓	✓	✓	Sun NFS
Web	✓	✓	✓	×	Web server
Distributed shared memory	✓	×	✓	✓	Ivy (DSM, Ch. 6)
Remote objects (RMI/ORB)	✓	×	×	1	CORBA
Persistent object store	✓	✓	×	1	CORBA Persistent State Service
Peer-to-peer storage system	✓	✓	✓	2	OceanStore (Ch. 10)

Types of consistency:

1: strict one-copy ✓: slightly weaker guarantees 2: considerably weaker guarantees

WHAT IS A FILE SYSTEM?

I

- Persistent stored data sets
- Hierarchic name space visible to all processes
- API with the following characteristics:
 - access and update operations on persistently stored data sets
 - Sequential access model (with additional random facilities)
- Sharing of data between users, with access control
- Concurrent access:
 - certainly for read-only access
 - what about updates?
- Other features:
 - mountable file stores
 - more? ...

WHAT IS A FILE SYSTEM?

2

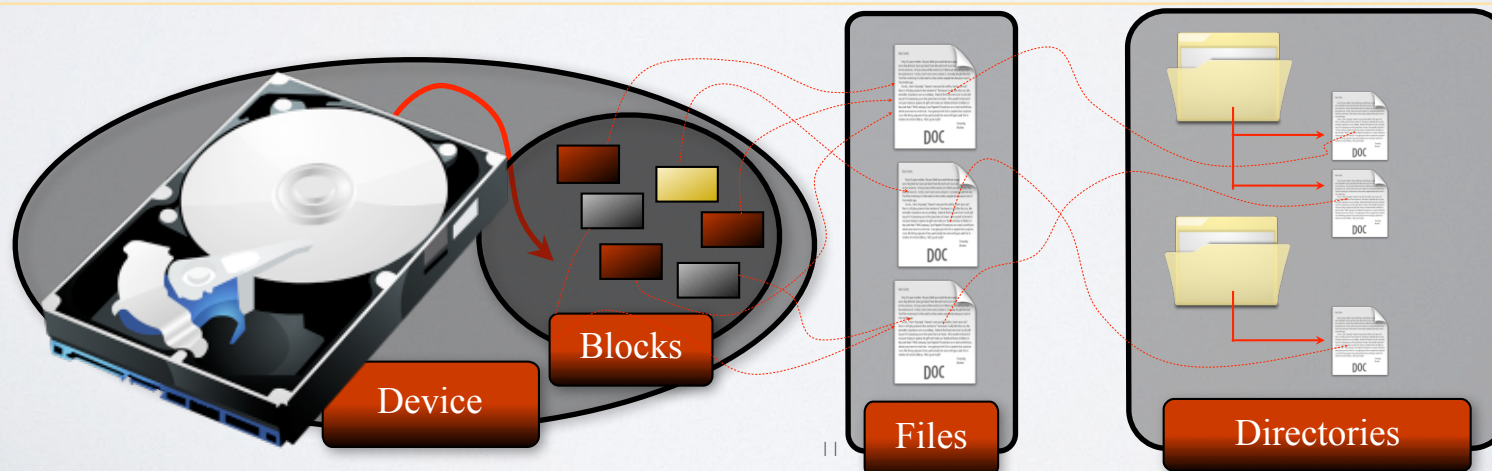
Figure 12.4 UNIX file system operations

<i>filedes</i> = <i>open</i> (<i>name</i> , <i>mode</i>)	Opens an existing file with the given <i>name</i> .
<i>filedes</i> = <i>creat</i> (<i>name</i> , <i>mode</i>)	Creates a new file with the given <i>name</i> .
	Both operations deliver a file descriptor referencing the open file. The <i>mode</i> is <i>read</i> , <i>write</i> or both.
<i>status</i> = <i>close</i> (<i>filedes</i>)	Closes the open file <i>filedes</i> .
<i>count</i> = <i>read</i> (<i>filedes</i> , <i>buffer</i> , <i>n</i>)	Transfers <i>n</i> bytes from the file referenced by <i>filedes</i> to <i>buffer</i> .
<i>count</i> = <i>write</i> (<i>filedes</i> , <i>buffer</i> , <i>n</i>)	Transfers <i>n</i> bytes to the file referenced by <i>filedes</i> from <i>buffer</i> .
	Both operations deliver the number of bytes actually transferred and advance the read-write pointer.
<i>pos</i> = <i>lseek</i> (<i>filedes</i> , <i>offset</i> , <i>whence</i>)	Moves the read-write pointer to <i>offset</i> (relative or absolute, depending on <i>whence</i>).
<i>status</i> = <i>unlink</i> (<i>name</i>)	Removes the file <i>name</i> from the directory structure. If the file has no other names, it is deleted.
<i>status</i> = <i>link</i> (<i>name1</i> , <i>name2</i>)	Adds a new name (<i>name2</i>) for a file (<i>name1</i>).
<i>status</i> = <i>stat</i> (<i>name</i> , <i>buffer</i>)	Puts the file attributes for file <i>name</i> into <i>buffer</i> .

WHAT IS A FILE SYSTEM?

(A TYPICAL MODULE STRUCTURE FOR IMPLEMENTATION OF NON-DFS)

Directory module:	relates file names to file IDs
File module:	relates file IDs to particular files
Access control module:	checks permission for operation requested
File access module:	reads or writes file data or attributes
Block module:	accesses and allocates disk blocks
Device module:	disk I/O and buffering



WHAT IS A FILE SYSTEM?

4

Figure 12.3 File attribute record structure

updated
by file
system:

File length

Creation timestamp

Read timestamp

Write timestamp

Attribute timestamp

Reference count

Owner

File type

Access control list

updated
by owner:

E.g. for UNIX: rw-rw-r--

DISTRIBUTED FILE SYSTEM/ SERVICE REQUIREMENTS



Transparency

Concurrency

Replication

Heterogeneity

Fault
Tolerance

Consistency

Security

Efficiency

Transparency

Concurrency

Replication

Heterogeneity

Fault
Tolerance

Consistency

Security

Efficiency

Transparency

Access: Same operations (client programs are unaware of distribution of files)

Location: Same name space after relocation of files or processes (client programs should see a uniform file name space)

Mobility: Automatic relocation of files is possible (neither client programs nor system admin tables in client nodes need to be changed when files are moved).

Performance: Satisfactory performance across a specified range of system loads

Scaling: Service can be expanded to meet additional loads or growth.

Transparency

Concurrency

Replication

Heterogeneity

Fault
Tolerance

Consistency

Security

Efficiency

Concurrency

Changes to a file by one client should not interfere with the operation of other clients simultaneously accessing or changing the same file.

Concurrency Properties

- Isolation
- File-level/record level locking
- Other forms of concurrency control to minimise contention

Transparency

Concurrency

Replication

Heterogeneity

Fault
Tolerance

Consistency

Security

Efficiency

Replication

Properties

File service maintains multiple identical copies of files

- Load-sharing between servers makes service more scalable
- Local access has better response (lower latency)
- Fault tolerance

Full replication is difficult to implement.

Caching (of all or part of a file) gives most of the benefits (except fault tolerance)

Transparency

Concurrency

Replication

Heterogeneity

Fault
Tolerance

Consistency

Security

Efficiency

Heterogeneity

Properties

Service can be accessed by clients running on (almost) any OS or hardware platform.

Design must be compatible with the file systems of different OSes

Service interfaces must be open - precise specifications of APIs are published.

Transparency

Concurrency

Replication

Heterogeneity

Fault
Tolerance

Consistency

Security

Efficiency

Fault Tolerance

Fault Tolerance

Service must continue to operate even when clients make errors or crash.

- at-most-once semantics
- at-least-once semantics
- requires idempotent operations

Service must resume after a server machine crashes.

If the service is replicated, it can continue to operate even during a server crash.

Transparency

Concurrency

Replication

Heterogeneity

Fault
Tolerance

Consistency

Security

Efficiency

Consistency

Unix offers one-copy update semantics for operations on local files

- caching is completely transparent.

Difficult to achieve the same for distributed file systems while maintaining good performance and scalability.

Transparency

Concurrency

Replication

Heterogeneity

Fault
Tolerance

Consistency

Security

Efficiency

Security

Must maintain access control and privacy as for local files.

- based on identity of user making request
- identities of remote users must be authenticated
- privacy requires secure communication

Service interfaces are open to all processes not excluded by a firewall.

- vulnerable to impersonation and other attacks

TRANSPARENT
SERVICE

Transparency

Concurrency

Replication

Heterogeneity

Fault
Tolerance

Consistency

Security

Efficiency

Efficiency

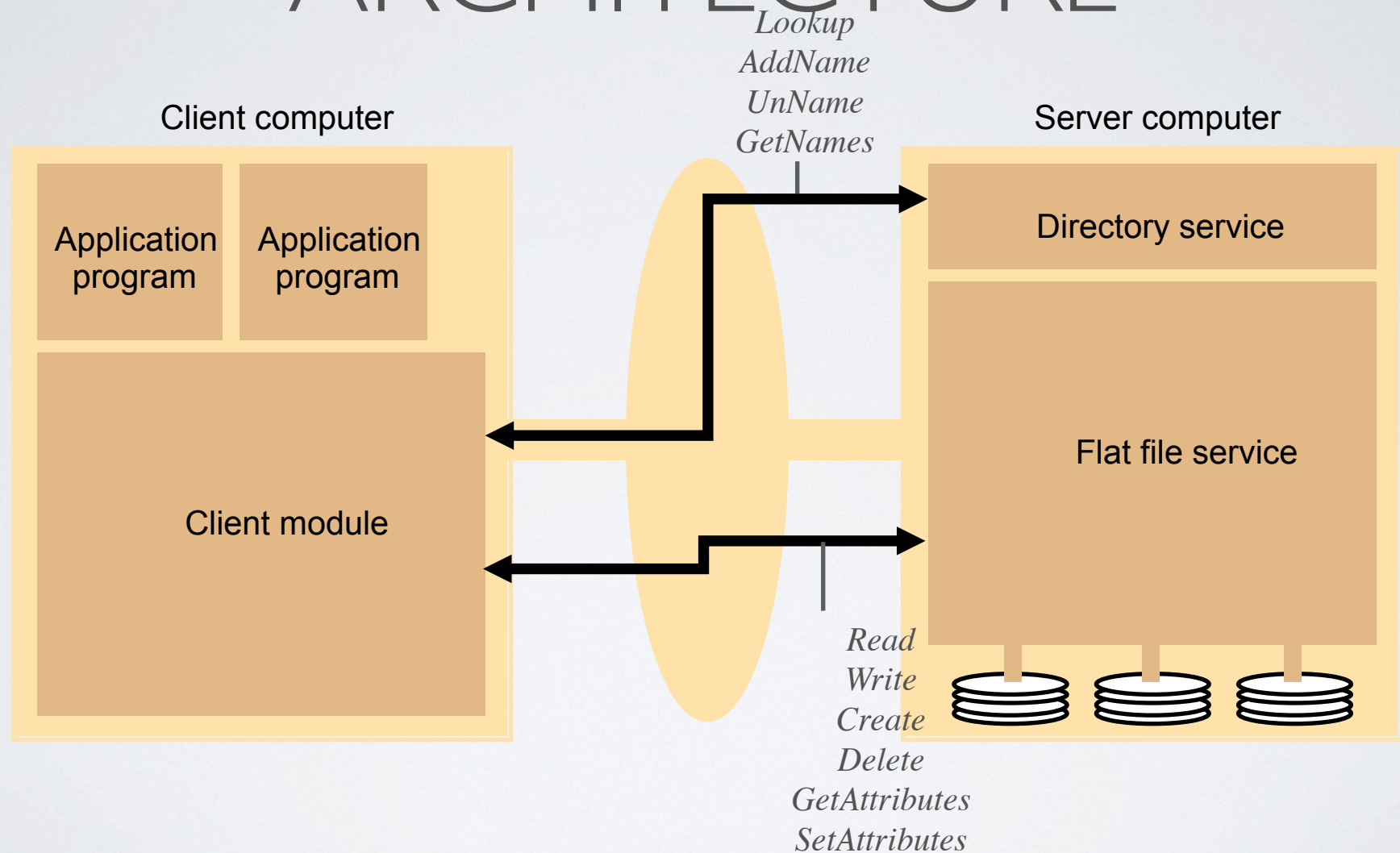
Goal for distributed file systems is usually performance comparable to local file system.

File service is most heavily loaded service in an intranet, so its functionality and performance are critical

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.
- The Client module implements exported interfaces by flat file and directory services on server side.

MODEL FILE SERVICE ARCHITECTURE



RESPONSIBILITIES OF VARIOUS MODULES

Client

Client Module

It runs on each computer and provides 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 achieve better performance through implementation of a cache of recently used file blocks at the client.

Server

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 generate directories, to add new files to directories.

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 among all of the files in a distributed system.

SERVER OPERATIONS/INTERFACES FOR THE MODEL FILE SERVICE

Figure 12.6 Flat file service operations

<i>Read</i> (<i>FileId</i> , <i>i</i> , <i>n</i>) → <i>Data</i> — throws <i>BadPosition</i>	If $1 \leq i \leq \text{Length}(\text{File})$: Reads a sequence of up to <i>n</i> items from a file starting at item <i>i</i> and returns it in <i>Data</i> .
<i>Write</i> (<i>FileId</i> , <i>i</i> , <i>Data</i>) — throws <i>BadPosition</i>	If $1 \leq i \leq \text{Length}(\text{File})+1$: Writes a sequence of <i>Data</i> to a file, starting at item <i>i</i> , extending the file if necessary.
<i>Create</i> () → <i>FileId</i>	Creates a new file of length 0 and delivers its <i>FileId</i> .
<i>Delete</i> (<i>FileId</i>)	Removes the file from the file store.
<i>GetAttributes</i> (<i>FileId</i>) → <i>Attr</i>	Returns the file attributes for the file.
<i>SetAttributes</i> (<i>FileId</i> , <i>Attr</i>)	Sets the file attributes (only those attributes that are shaded in Figure 12.3).

Pathname lookup

Pathnames such as '/usr/bin/tar' are resolved by iterative calls to `lookup()`, one call for each component of the path, starting with the ID of the root directory '/' which is known in every client.

Figure 12.7 Directory service operations

<i>Lookup</i> (<i>Dir</i> , <i>Name</i>) → <i>FileId</i> — throws <i>NotFound</i>	Locates the text name in the directory and returns the relevant UFID. If <i>Name</i> is not in the directory, throws an exception.
<i>AddName</i> (<i>Dir</i> , <i>Name</i> , <i>FileId</i>) — throws <i>NameDuplicate</i>	If <i>Name</i> is not in the directory, adds (<i>Name</i> , <i>File</i>) to the directory and updates the file's attribute record. If <i>Name</i> is already in the directory, throws an exception.
<i>UnName</i> (<i>Dir</i> , <i>Name</i>) — throws <i>NotFound</i>	If <i>Name</i> is in the directory, removes the entry of <i>Name</i> from the directory. If <i>Name</i> is not in the directory, throws an exception.
<i>GetNames</i> (<i>Dir</i> , <i>Pattern</i>) → <i>NameSeq</i>	Returns all the text names in the directory that match the regular expression <i>Pattern</i> .

Fileid

A unique identifier for files anywhere in the network. Similar to the remote object references described in Section 4.3.3.

HADOOP

- <https://www.youtube.com/watch?v=9S9H6OzJ9N4>

SUMMARY

- Distributed File systems provide illusion of a local file system and hide complexity from end users.
- Sun NFS is an excellent example of a distributed service designed to meet many important design requirements
- Effective client caching can produce file service performance equal to or better than local file systems
- Consistency versus update semantics versus fault tolerance remains an issue
- Most client and server failures can be masked