
Chapter 12: Distributed File Systems

Distributed File Systems

❑ Objectives

- To be reminded on storage systems and their properties
- To understand the basics of file system modules and attributes
- To compare differences of file systems and their distributed sisters as well as understand major characteristics of examples

❑ Topics

- Storage, file systems, and file attribute records
- Reference counting
- Distributed file systems: Operations and requirements
- File service architecture
 - Network File System (NFS)
 - Inter Planetary File System (IPFS)
 - Google File System

Storage Systems and Their Properties

<i>System</i>	<i>Sharing across network</i>	<i>Persis- tence</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	✓	✓	✓	2 ✓	Sun NFS
Web	✓	✓	✓	×	Web server
Distributed shared memory	✓	×	✓	2	Ivy
Remote objects (RMI)	✓	×	×	1	CORBA
Persistent object store	✓	✓	×	1	CORBA Persistent Object Service
Peer-to-peer storage system	✓	✓	✓	3	OceanStore

Consistency types: 1: Strictly one-copy
 ✓ 2: Slightly weaker guarantees
 3: Considerably weaker guarantees

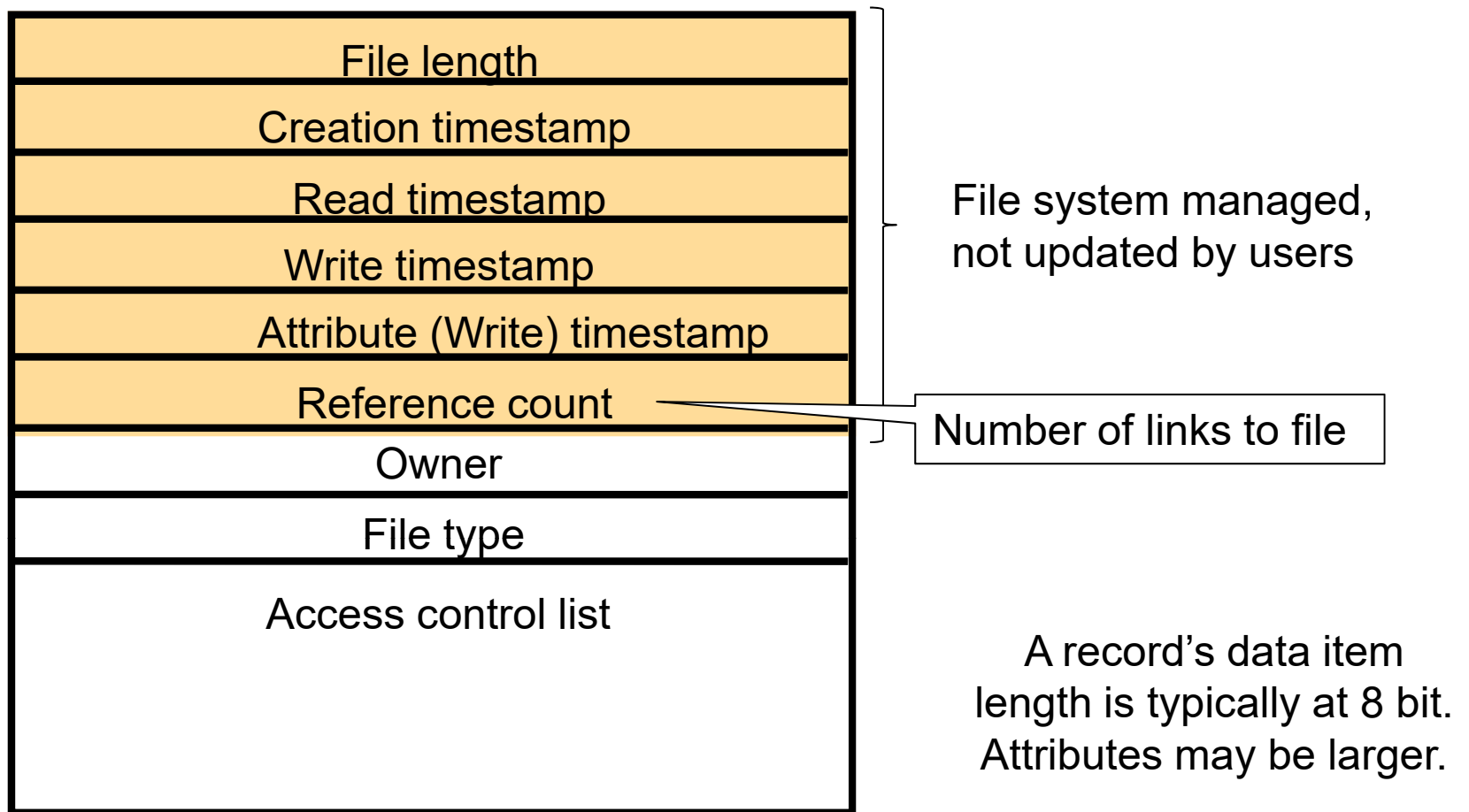
File System Modules

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

File Attribute Record Structure

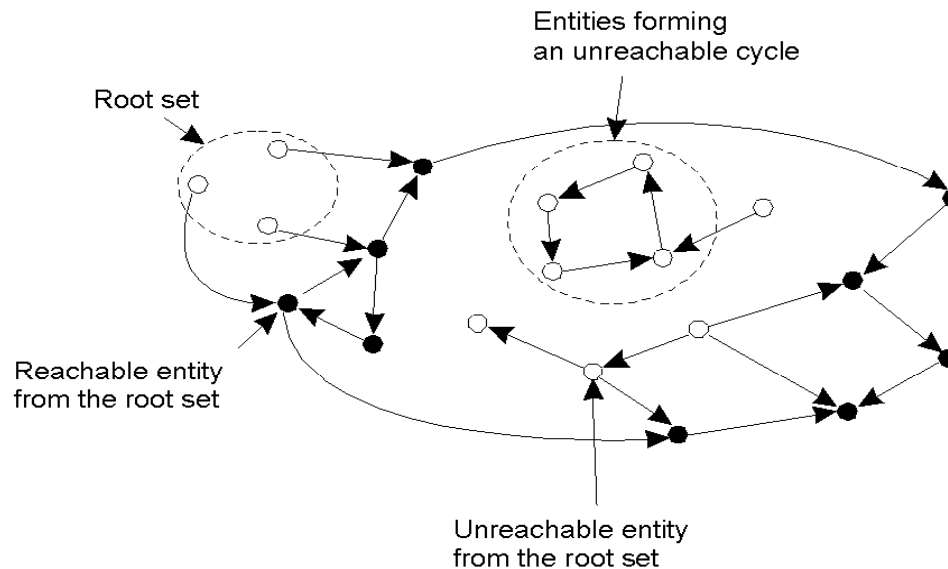
File systems: organization, storage, retrieval, naming, sharing, protection of files

Files: contain attributes (maintained by file system) and data (updated by user)



Reference Counting

- ❑ Maintenance task for distributed systems, OSes
 - Applied to files and objects
- ❑ *E.g.*, **garbage collection**
 - Delete object/file, if not reachable from root

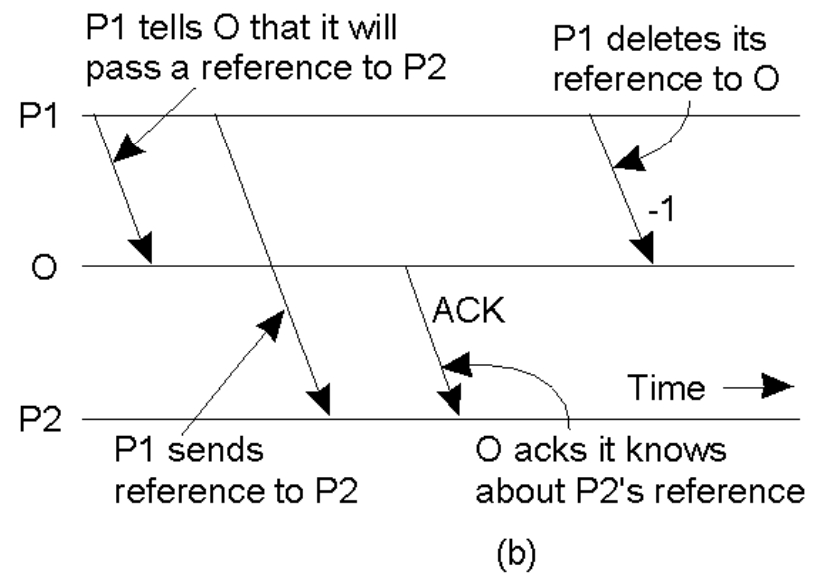
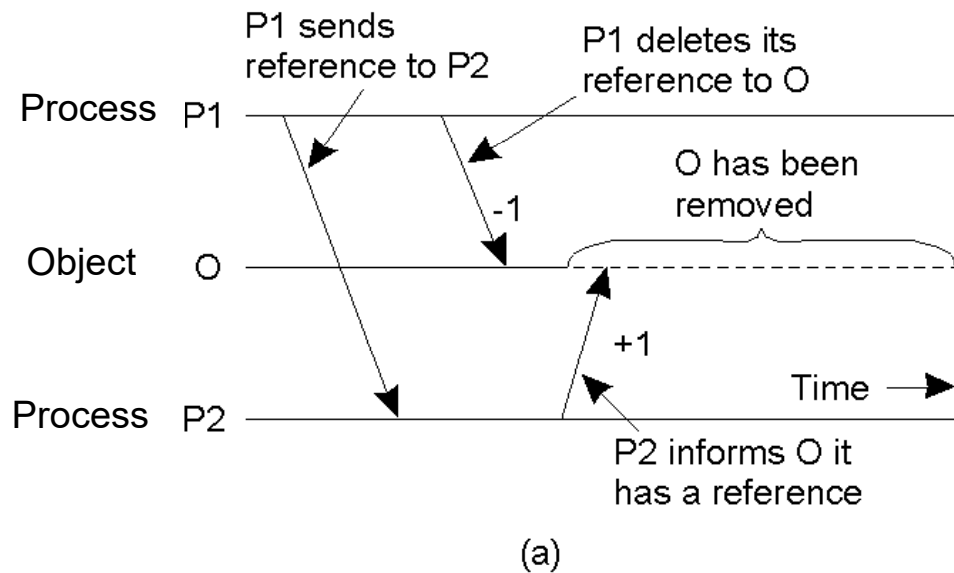


- ❑ How to run such a scheme in a distributed manner?

Distributed Reference Counting

- ❑ Copying a reference to another process and incrementing the counter can be too late

- ❑ A solution



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 offset (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>)	Gets the file attributes for file <i>name</i> into <i>buffer</i> .

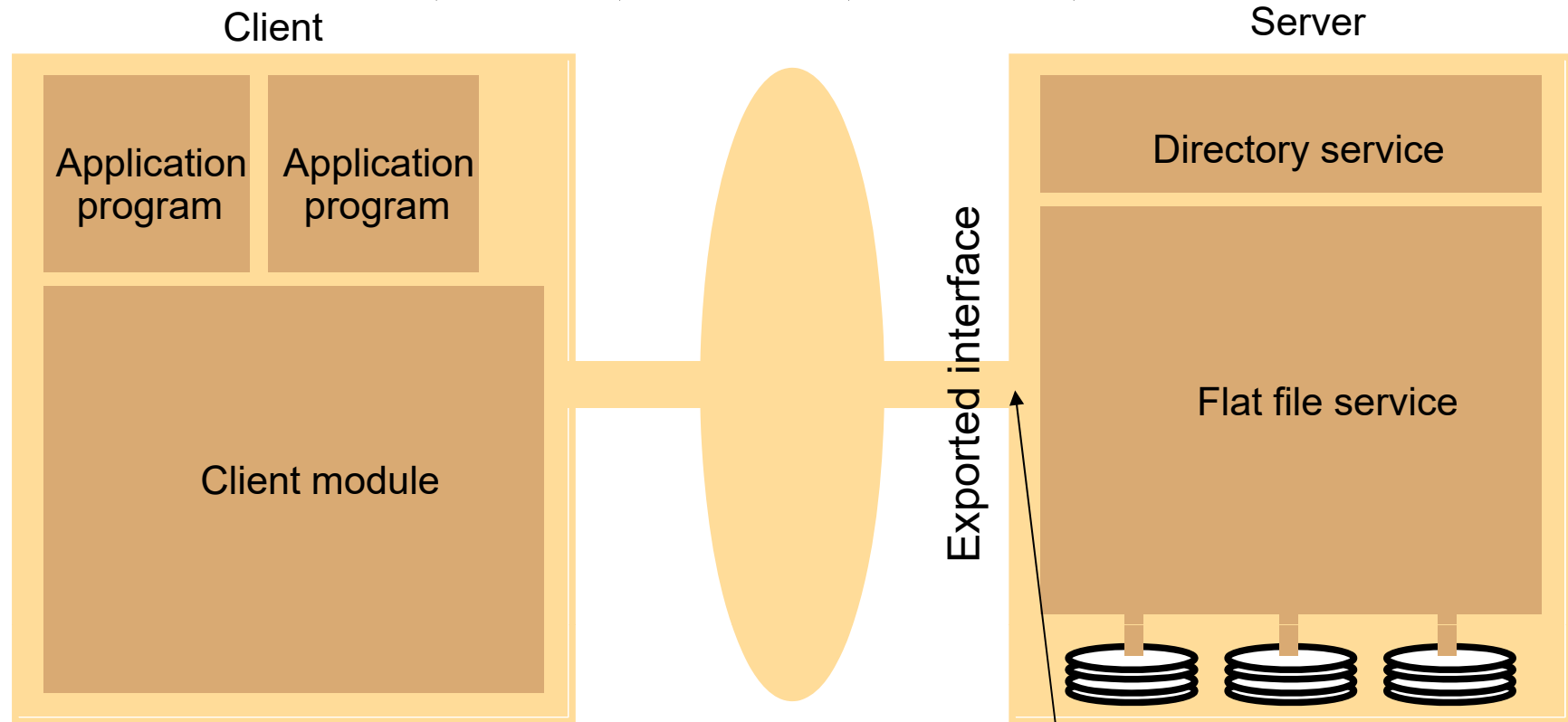
Figure 12.4, CDK

Distributed File System Requirements

- ❑ Transparency
 - Access, Location, Scaling, and all the others
- ❑ Concurrent file updates
- ❑ File replication
- ❑ Hardware and software heterogeneity
- ❑ Fault tolerance
- ❑ Consistency (one-copy update semantics)
- ❑ Security
- ❑ Efficiency

File Service Architecture

- ❑ Services: read, write, create, delete, ...



- ❑ **Unique File Identifiers (UFID)** used
 - Referring to files in all requests for service operations

Flat File Service Operations

<i>Read(FileId, i, n) -> Data</i> — throws <i>BadPosition</i>	If $1 \leq i \leq \text{Length}(\text{File})$: Reads a sequence of up to n items from a file starting at item i and returns it in <i>Data</i> .
<i>Write(FileId, 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 , extending the file if necessary.
<i>Create() -> FileId</i>	Creates a new file of length 0 and delivers a UFID for it.
<i>Delete(FileId)</i>	Removes the file from the file store.
<i>GetAttributes(FileId) -> Attr</i>	Returns the file attributes for the file.
<i>SetAttributes(FileId, Attr)</i>	Sets the file attributes (only those attributes that are not shaded in Figure 8.3).

Note: This interface is repeatable (*i.e.*, idempotent) and stateless!

Directory Service Operations

Lookup(Dir, Name) -> FileId
— throws *NotFound*

Locates the text name in the directory and returns the relevant UFID. If *Name* is not in the directory, throws an exception.

AddName(Dir, Name, FileId)
— throws *NameDuplicate*

If *Name* is not in the directory, adds (*Name, File*) to the directory and updates the file's attribute record.
If *Name* is already in the directory: throws an exception.

UnName(Dir, Name)
— throws *NotFound*

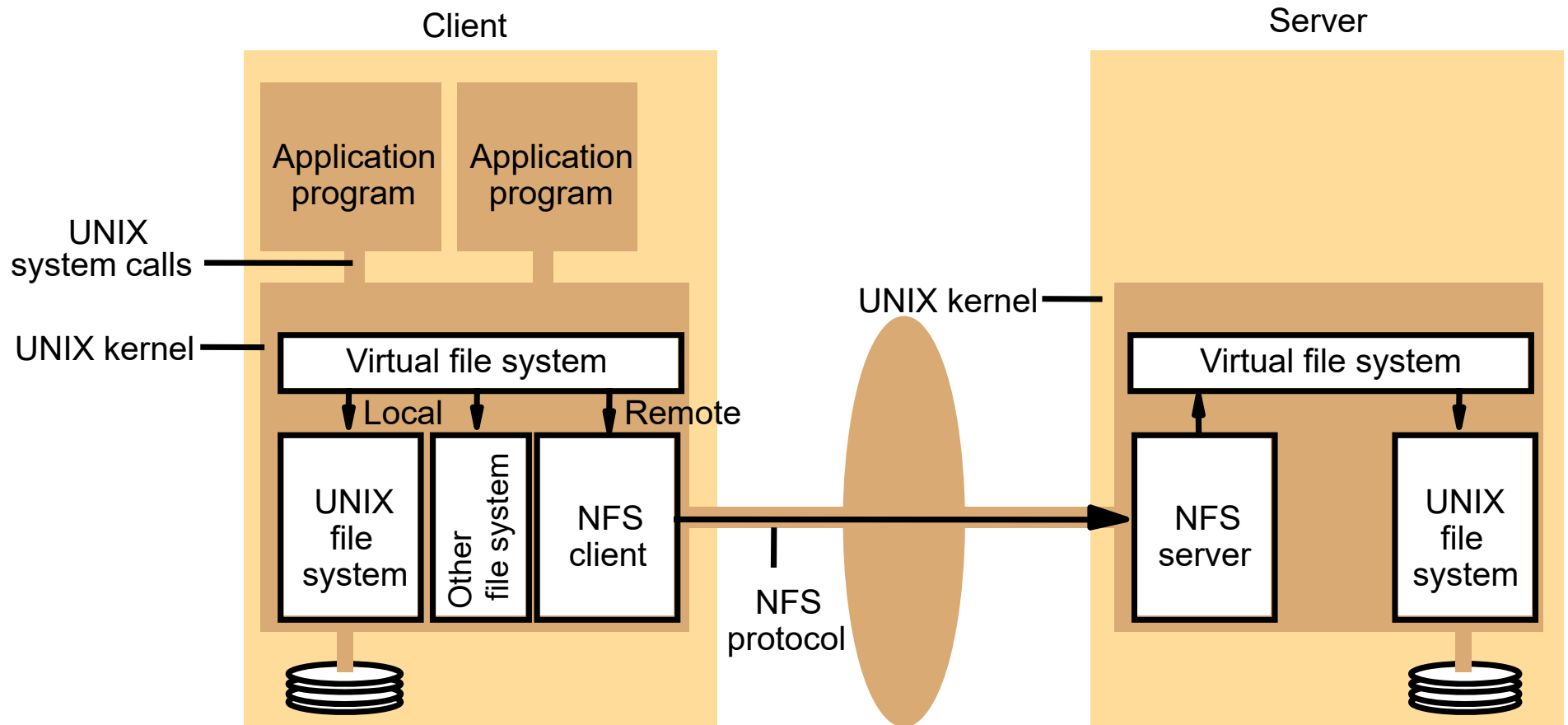
If *Name* is in the directory: the entry containing *Name* is removed from the directory.
If *Name* is not in the directory: throws an exception.

GetNames(Dir, Pattern) -> NameSeq

Returns all the text names in the directory that match the regular expression *Pattern*.

Case Study 1: NFS Architecture

❑ Network File System (NFS)



NFS Server Operations Simplified (1)

<i>lookup(dirfh, name) -> fh, attr</i>	Returns file handle and attributes for the file <i>name</i> in the directory <i>dirfh</i> .
<i>create(dirfh, name, attr) -> newfh, attr</i>	Creates a new file name in directory <i>dirfh</i> with attributes <i>attr</i> and returns the new file handle and attributes.
<i>remove(dirfh, name) status</i>	Removes file name from directory <i>dirfh</i> .
<i>getattr(fh) -> attr</i>	Returns file attributes of file <i>fh</i> . (Similar to the UNIX <i>stat</i> system call.)
<i>setattr(fh, attr) -> attr</i>	Sets the attributes (mode, user id, group id, size, access time and modify time of a file). Setting the size to 0 truncates the file.
<i>read(fh, offset, count) -> attr, data</i>	Returns up to <i>count</i> bytes of data from a file starting at <i>offset</i> . Also returns the latest attributes of the file.
<i>write(fh, offset, count, data) -> attr</i>	Writes <i>count</i> bytes of data to a file starting at <i>offset</i> . Returns the attributes of the file after the write has taken place.
<i>rename(dirfh, name, todirfh, toname) -> status</i>	Changes the name of file <i>name</i> in directory <i>dirfh</i> to <i>toname</i> in directory to <i>todirfh</i>
<i>link(newdirfh, newname, dirfh, name) -> status</i>	Creates an entry <i>newname</i> in the directory <i>newdirfh</i> which refers to file <i>name</i> in the directory <i>dirfh</i> .

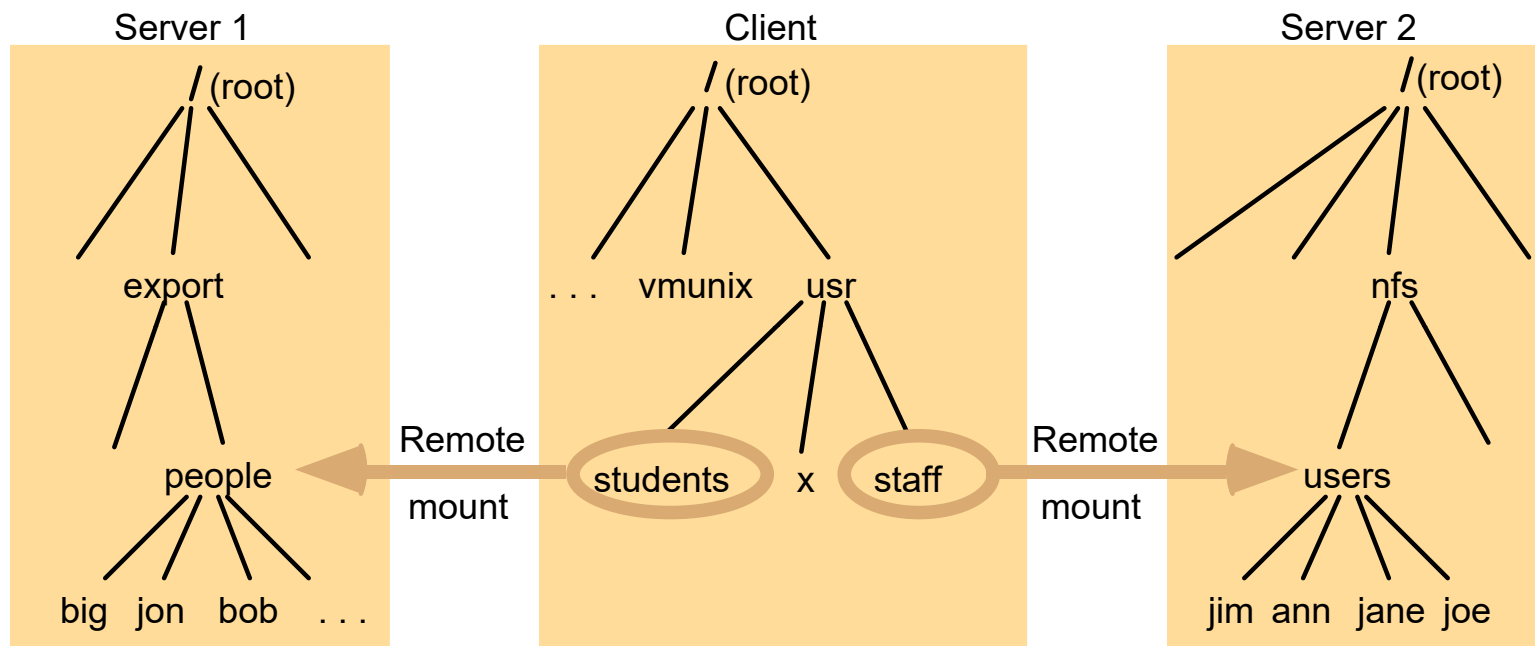
NFS Server Operations Simplified (2)

<i>symlink(newdirfh, newname, string)</i> -> <i>status</i>	Creates an entry <i>newname</i> in the directory <i>newdirfh</i> of type symbolic link with the value <i>string</i> . The server does not interpret the <i>string</i> but makes a symbolic link file to hold it.
<i>readlink(fh)</i> -> <i>string</i>	Returns the string that is associated with the symbolic link file identified by <i>fh</i> .
<i>mkdir(dirfh, name, attr)</i> -> <i>newfh, attr</i>	Creates a new directory <i>name</i> with attributes <i>attr</i> and returns the new file handle and attributes.
<i>rmdir(dirfh, name)</i> -> <i>status</i>	Removes the empty directory <i>name</i> from the parent directory <i>dirfh</i> . Fails if the directory is not empty.
<i>readdir(dirfh, cookie, count)</i> -> <i>entries</i>	Returns up to <i>count</i> bytes of directory entries from the directory <i>dirfh</i> . Each entry contains a file name, a file handle, and an opaque pointer to the next directory entry, called a <i>cookie</i> . The <i>cookie</i> is used in subsequent <i>readdir</i> calls to start reading from the following entry. If the value of <i>cookie</i> is 0, reads from the first entry in the directory.
<i>statfs(fh)</i> -> <i>fsstats</i>	Returns file system information (such as block size, number of free blocks and so on) for the file system containing a file <i>fh</i> .

Local and Remote File Systems Access

❑ Accessible from an NFS client

- 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



NFS Concerns

- ❑ Auto-mounter
 - Mount on access of any non-mounted file
 - Valid for a set of registered directories
- ❑ Server caching
 - Conventional UNIX: read-ahead and delayed write
 - NFS: write-through
- ❑ Client caching
 - Timeouts
 - Cache valid iff $(T - T_c < c)$ or $(\text{FileAttrclient} = \text{FileAttrserver})$
- ❑ Security
 - Quite flaky → Better with access control systems (Kerberos)

Case Study 1: Summary

- ❑ Transparency

- Access
- Location
- Mobility
- Scaling

- ❑ Concurrent file updates

- ❑ File replication

- ❑ Hardware and software heterogeneity

- ❑ Fault tolerance

- ❑ Consistency

- ❑ Security

- ❑ Efficiency

NFS client process

File name space identical

Remounts possible

Only if no “hot spot files”

Not supported

Read-only files

Given

Stateless and idempotent

Close to one-copy semantics

Medium

Seems ok (still in larger use)

Case Study 2: IPFS



❑ Inter Planetary File System

- Defined as a protocol for the file system defining a **decentralized storage** and **delivery network** build upon peer-to-peer (p2p) and content-addressed file-system principles

❑ Goal

- Replacing centralized systems and protocols, which serve today as the backbone of the Web providing a more open, available, and faster Internet

• **Open**: decentralization of services and content

• **Available**: content is decentralized in the network

• **Faster**: pages are loaded from multiple peers instead of a single host

Service Unavailable

The server is temporarily unable to service your request due to maintenance or capacity problems. Please try again later.

Apache/2.4.7 (Ubuntu) Server at coostord.az.10010.com

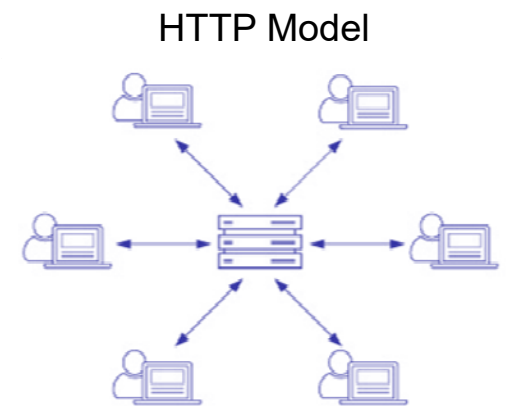
Not Found

The requested URL /oldpage.html was not found on this server.
Apache/2.2.3 (CentOS) Server at www.example.com Port 80

Contrasting Juxtaposition: HTTP and IPFS

❑ Hyper Text Transfer Protocol (HTTP)

- Backbone protocol of the Internet
- Essentially **centralized** implemented
 - Relying on hosts to provide content/services
 - Widely deployed



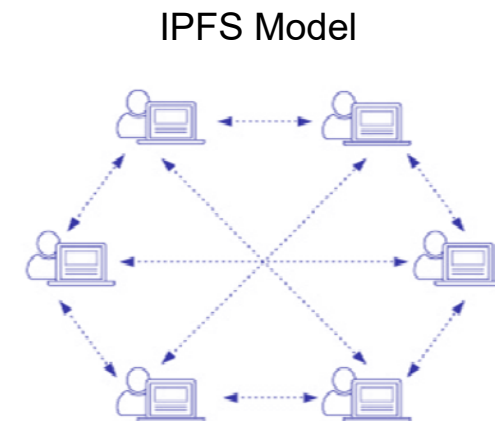
❑ IPFS

- Essentially **decentralized**
 - Reliably persisting characteristics
 - Secured against DDoS attacks
 - Reduced time to load webpages

– Drawbacks

DDoS: Distributed Denial-of-Service Attacks

- Hard maintenance and privacy concerns



Feature-based Comparison: HTTP and IPFS

Feature	IPFS	HTTP
Architecture	Decentralized	Centralized
Addressing	Content-based whereas files are identified by a hash	Location-based whereas hosts are identified by an IP address
Bandwidth	Higher as data is retrieved from multiple peers	Lower since data retrieval' relies on host's upload bandwidth
Availability	Higher since data is decentralized across multiple peers	Lower since data is stored in a single host
Privacy	Lower since data is decentralized across multiple unknown peers	Higher since data is stored at known hosts
Adoption	Slower since it is a relatively new protocol	Established and industry-standard protocol
Costs	Lower since every peer can host content	Higher as a single host should provide the entire content/service

IPFS Building Blocks

- ❑ IPFS as a **distributed file system** seeks to connect all computing devices with the same system of files
- ❑ IPFS built upon major components
 - A. **Merkle DAG**: uses a Merkle Tree or a Merkle DAG similar to the one used in the Git version control system
 - B. **Block Exchange and Network**: used in the BitTorrent Protocol (also known as BitSwap) to exchange data between nodes, it is a p2p file sharing protocol, which coordinates data exchange between untrusted swarms
 - C. **Distributed Hash Tables (DHT)**: Kademlia is used to store and retrieve data across nodes in the network (*i.e.*, content-based location)

DAG: Directed Acyclic Graph

A. IPFS Merkle DAG (1)

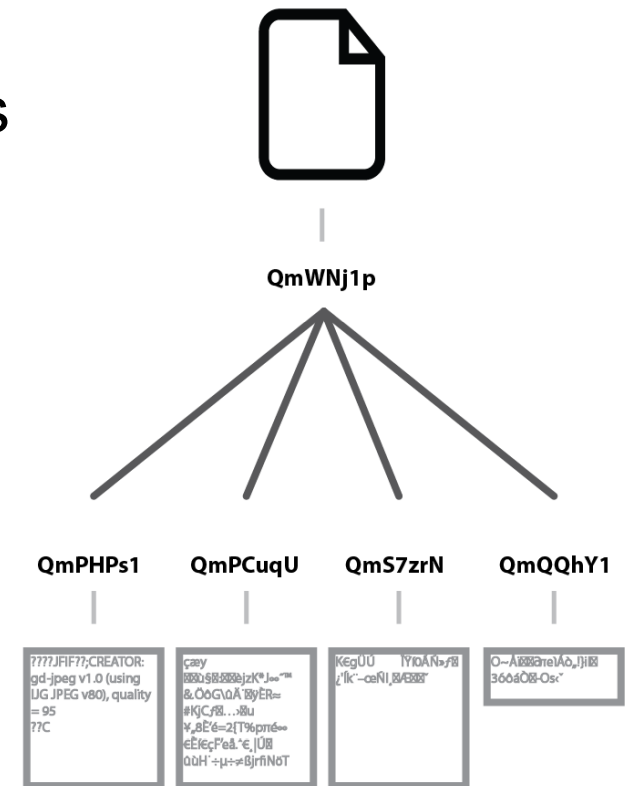
❑ “Heart” of the IPFS

– Directed Acyclic Graph (DAG) with links

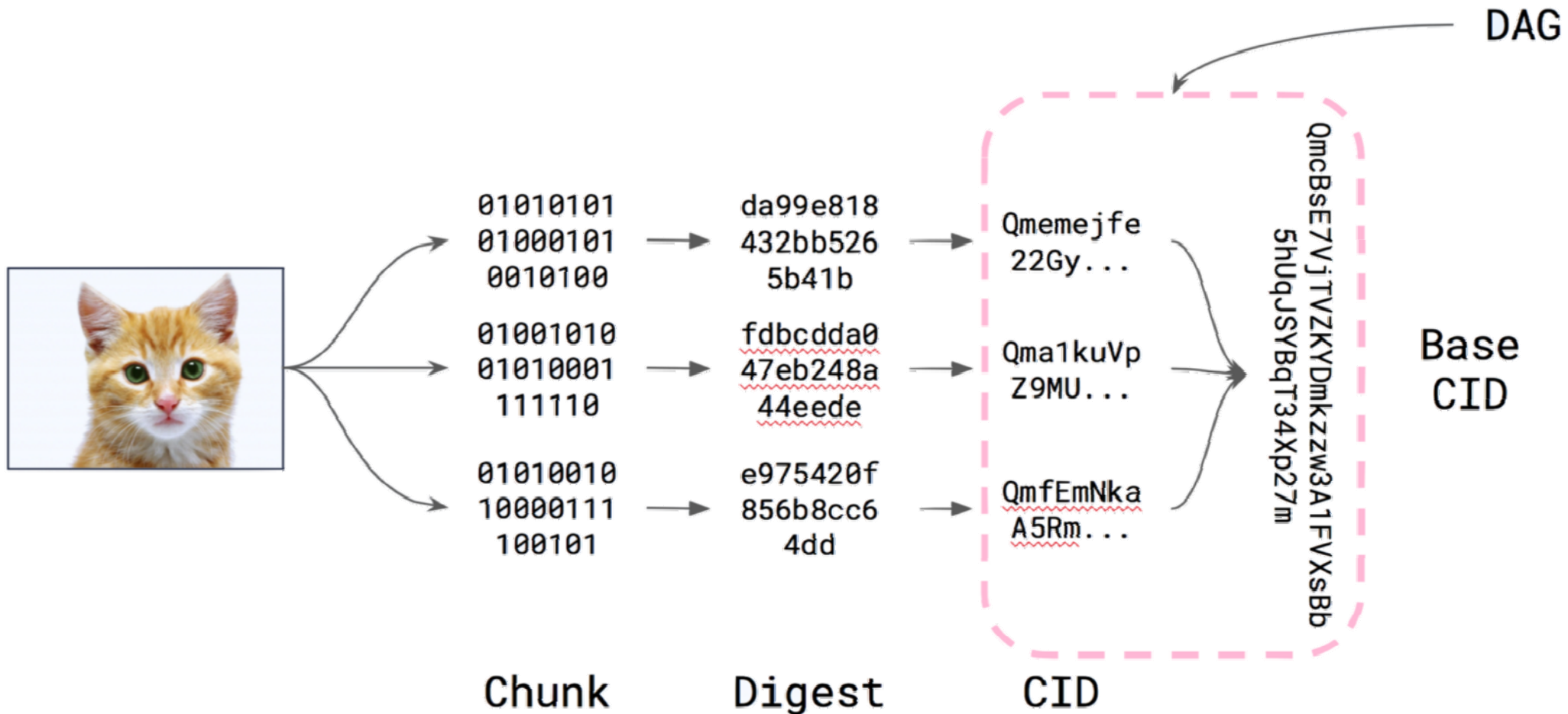
- Links are implemented as hashes

– Properties

- Authenticated
 - Content can be hashed and verified against the link
- Permanent
 - Once fetched, objects can be cached forever
- Universal
 - Any data structure can be represented as a Merkle DAG
- Decentralized
 - Objects can be created by anyone, without centralized writers

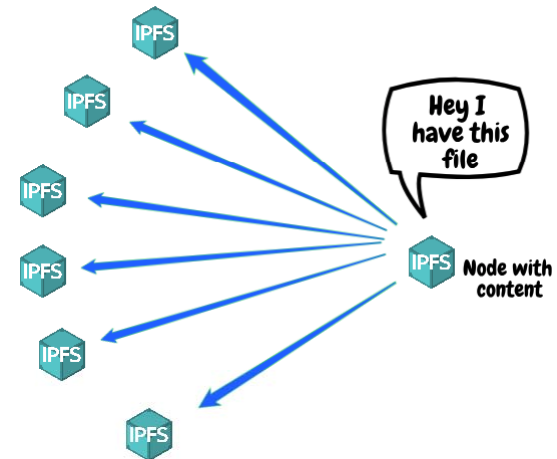


A. IPFS Merkle DAG (2)



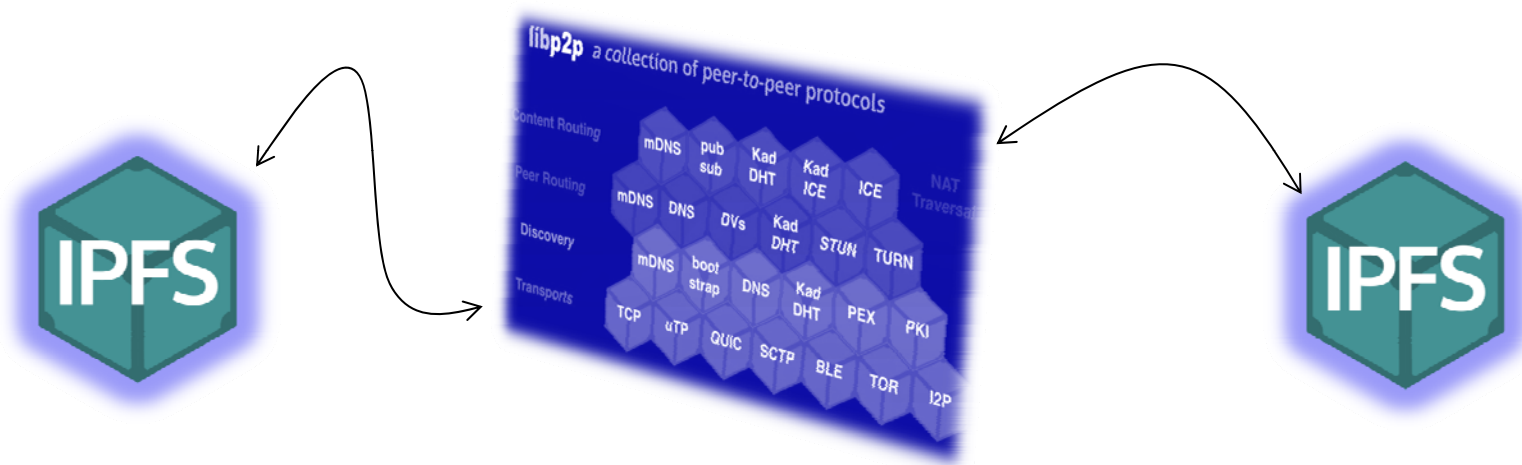
B. Block and Node

- ❑ **Blocks** represent data to be stored
- ❑ **IPFS node** is a component to find, publish, and replicate “merkledag” objects
 1. “IPFS Block Exchange” negotiates bulk data transfers
 - Bitswap is the main protocol for exchanging data
 - Generalization of BitTorrent to work with arbitrary and not known a priori DAGs
 2. HTTP as a simple exchange implemented with HTTP clients and servers
 - Compatibility reasons
- ❑ Protocol to **handle object’s operations b/w IPFS nodes**
 - File sharing protocol coordinates data exchange



B. Block Exchange and Network

- ❑ IPFS' underlying **network** provides
 - Point-to-point transport connections (reliable and unreliable) between any two IPFS nodes
- ❑ Protocol needed to **exchange data (blocks) between IPFS nodes**
 - Based on **Bitswap** (block exchange) and on **libp2p** (routing)



B. Routing

❑ IPFS Routing

- Protocol defined within a separate layer
 - [Peer Routing](#) to find other nodes
 - [Content Routing](#) to find data published to IPFS

❑ IFPS Routing System

- Defined as an interface
- Satisfied by different kind of implementations
 - [DHT \(Distributed Hash Table\)](#): used to create a semi-persistent routing record for a distributed cache in the network
 - [MDNS \(Multicast DNS\)](#): used to find services advertised locally
 - [DNS \(Domain Name System\)](#): used to find a node for the storing and sharing of data, can utilize the IPNS (Inter Planetary Name Space)

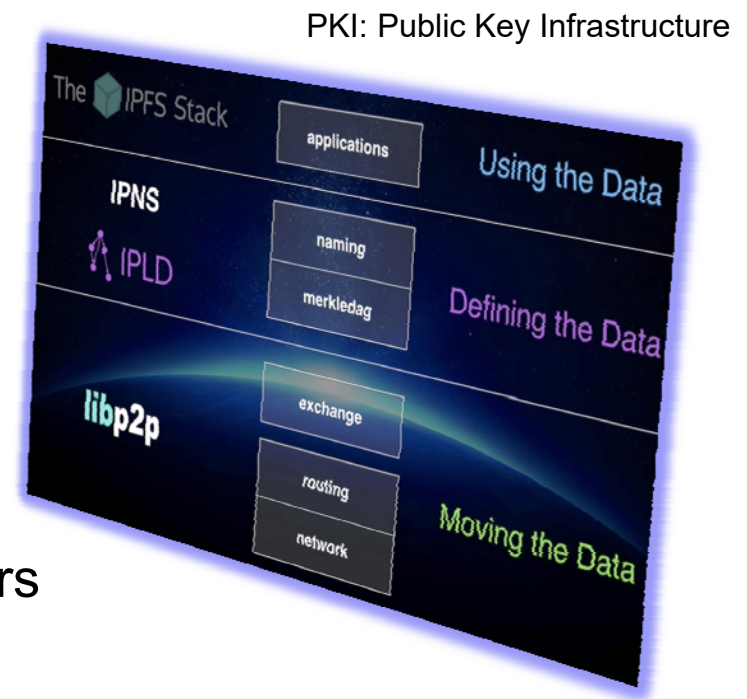
C. Distributed Hash Table (DHT)

- ❑ DHT used as fundamental IPFS component
 - For the routing system
 - Acts like a cross between a catalog and a navigation system

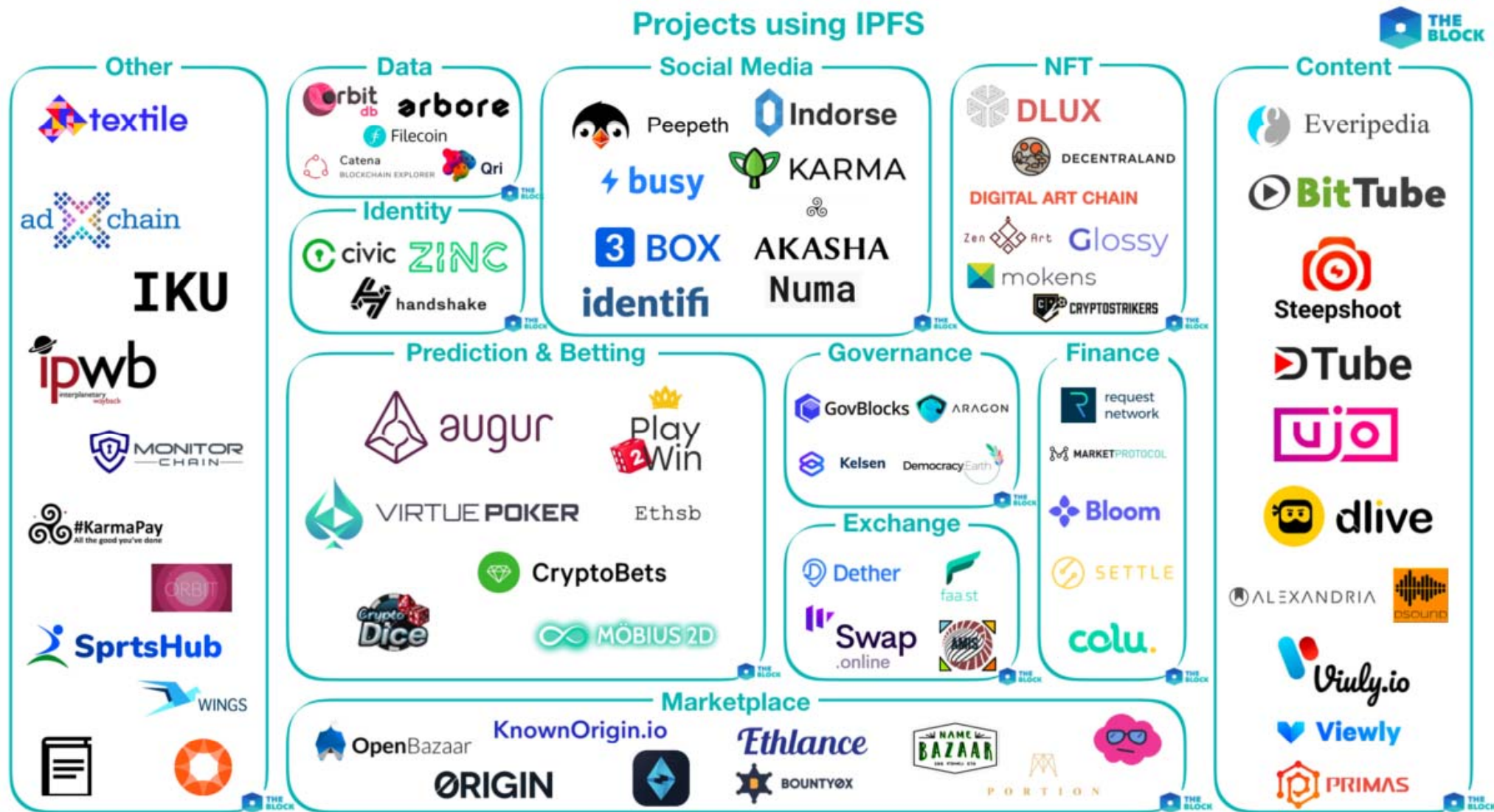
Type	Purpose	Used by
Provider records	Map a data identifier (<i>i.e.</i> , a multi-hash) to a peer that has advertised that they have that content and are willing to provide it to you.	<ul style="list-style-type: none">- IPFS to find content- IPNS over PubSub to find other members of the pubsub <i>topic</i>.
IPNS records	Map an IPNS key (<i>i.e.</i> , the hash of a public key) to an IPNS record (<i>i.e.</i> , a signed and versioned pointer to a path like <code>/ipfs/bafyxyz...</code>)	<ul style="list-style-type: none">- IPNS
Peer records	Map a peerID to a set of multi-addresses at which the peer may be reached	<ul style="list-style-type: none">- IPFS when we know of a peer with content, but do not know its address.- Manual connections (<i>e.g.</i>, <code>ipfs swarm connect /p2p/Qmxyz...</code>)

IPFS Protocol Stack

- ❑ IPFS protocol stack offers **modular protocols**
 - Each layer may have multiple implementations
 - All may use **different modules**
 - Applications: Web apps, Blockchain
 - Naming: self-certifying PKI for IPNS
 - MerkleDAG:
data structure format (thin waist)
 - Exchange:
block transport and replication
 - Routing: locating peers and objects
 - Network:
establishing connections between peers



IFPS Applications and Projects



IPFS' Privacy

- ❑ IPFS utilizes a **public network** – the Internet
- ❑ Methods to ensure privacy are (partially) essential
 - **Private IPFS Network**
 - Behaves similarly to the public network except for the fact that participants are only able to communicate with other nodes inside that same private IPFS network
 - **Content Encryption**
 - Considering encryption of content uploaded to the IPFS network,
 - It still can be tracked, but content remains “unreadable”
 - **Gateway Usage**
 - Requesting content through a public gateway allows a user to retrieve content from the IPFS network without running their own node
 - Hiding identities behind gateways

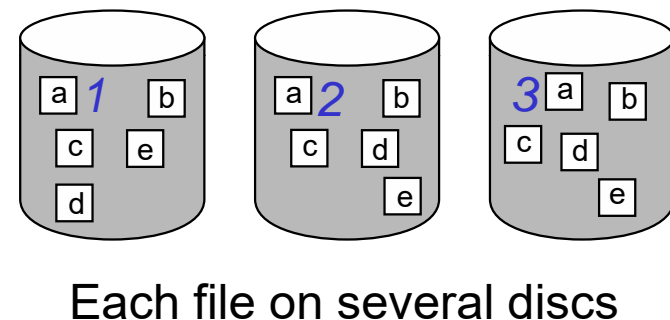
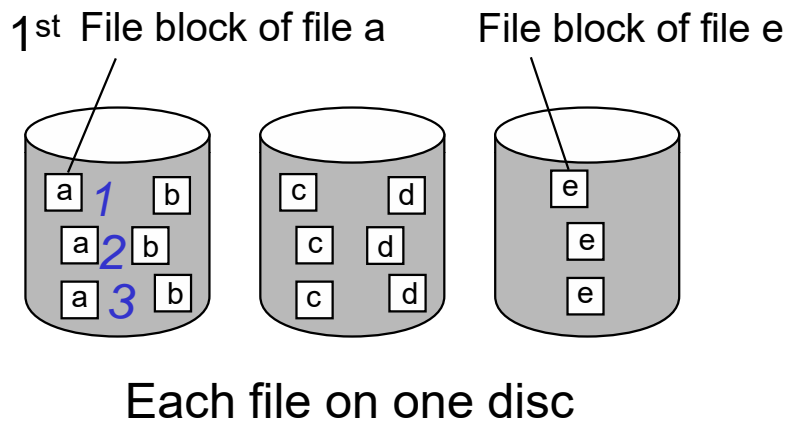
Case Study 2: Summary

- ❑ IPFS defines a system and related protocols to **decentralize services and content offered** on the Internet
 - Increased performance and availability due to its decentralized nature, but not “private”
 - On-going project with “complex” configurations for daily users

- ❑ Growing list of **use cases and applications**
 - Standalone and decentralized file system
 - Complementary decentralized storage to existing HTTP-based and centralized systems
 - Important role as decentralized storage for Blockchains

Case Study 3: Google File System

- ❑ Designed for large files and appends to files
- ❑ Files are distributed over multiple servers in chunks of 64 MB (like in a **RAID**) RAID: Redundant Array of Independent Disks
→ File striping
- ❑ Chunks are replicated (not shown below)



Google File System Architecture

- ❑ Master serves acts as **directory** only
- ❑ All chunks are distributed over chunk server
→ better scalability

