

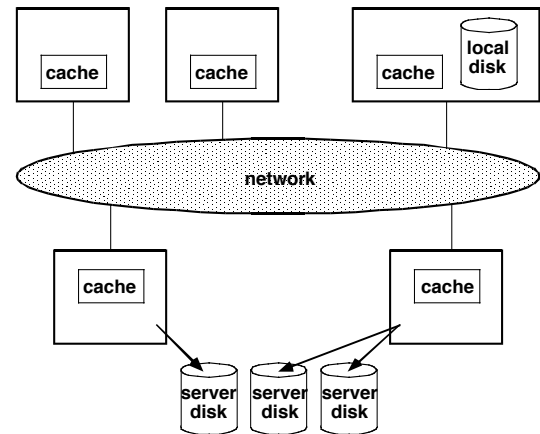
## Distributed File Systems

- *Distributed file system* — a distributed implementation of a file system
  - *File service* — specification of the file system interface as seen by the clients
  - *File server* — a process running on some machine which helps implement the file service by supplying files
- Goals of a distributed file system
  - *Network transparency*
    - Provide same operations for accessing remote and local files
    - Ideally, clients should not have to know the location of files to access them
  - *Availability / robustness* — file service should be maintained even in the presence of partial system failures
  - *Performance* — should overcome bottlenecks of a centralized file system

1

Fall 2005, Lecture 14

## Distributed File Systems (cont.)



- In principle, files in a distributed file system can be stored at any machine
  - For increased performance and reliability, a typical distributed environment has a few dedicated machines called *file servers* that store most of the files

2

Fall 2005, Lecture 14

## Review of Some File Concepts

- Logical components of a file
  - File name, file attributes, data blocks
  - Directory maps file name to file descriptor (inode in Unix terms)
  - File descriptor contains file attributes and pointers to data blocks
- Basic operations
  - Create / delete, open / close, read / write
- Types of file access
  - Sequential, direct / random, keyed
  - File pointer keeps track of location in file on a per-process basis
- Two separate concepts:
  - File lookup / naming (directory service)
  - File access (file service)

3

Fall 2005, Lecture 14

## Distributed File System Services — File Service Interface

- Need operations for creating and deleting, opening and closing, and reading and writing, files
- Upload / download model
  - File service provides:
    - Read — transfer entire file to client
    - Write — transfer entire file to server
  - Client works on file locally (in memory or on disk)
  - ✓ Simple, efficient if working on entire file
  - ✗ Must move entire file
  - ✗ Needs local disk space
- Remote access model
  - File service provides usual file operations
  - File stays on server

4

Fall 2005, Lecture 14

## Distributed Naming Structures

- Need operations for name translation, support for multilevel directories and links
  - *Location transparency* — the name of the file does not reveal the physical storage location
    - True for many naming schemes
  - *Location independence* — the name of the file need not change if the file's storage location changes
    - False for most naming schemes
- Absolute names
  - Names of form: *machine : pathname*
  - Used by:
    - Old UNIX distributed file systems
    - Current web browsers (e.g., Netscape)
  - ✓ User can use same tools and file operations for local and remote access
  - ✗ Not location transparent or independent

5

Fall 2005, Lecture 14

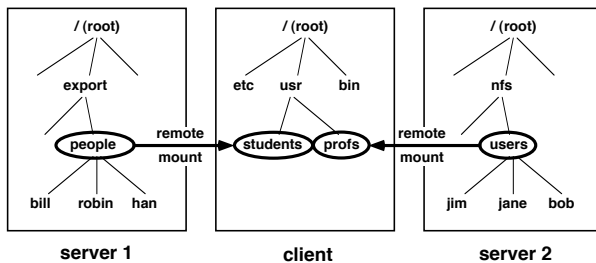
## Sun's Network File System

- Designed by Sun Microsystems
  - First distributed file service designed as a project, introduced in 1985
  - To encourage its adoption as a standard
    - Definitions of the key interfaces were placed in the public domain in 1989
    - Source code for a reference implementation was made available to other computer vendors under license
    - Currently the *de facto* standard for LANs
- Provides transparent access to remote files on a LAN, for clients running on UNIX and other operating systems
  - A UNIX computer typically has a NFS client and server module in its OS kernel
    - Available for almost any UNIX and MACH
  - Client modules are available for Macintosh OS and Windows

6

Fall 2005, Lecture 14

## Mounting Remote File Systems



- NFS supports mounting of remote file systems onto a specific location (mount point) in a client machine's file system
- NFS does not enforce a single network-wide name space
  - Name space seen by each client may be different; same file on server may have different path names on different clients
  - A uniform name space (and location transparency) can be set up if desired

7

Fall 2005, Lecture 14

## Mounting Remote File Systems (cont.)

- On each server
  - There is a file (usually */etc/exports*) containing the names of local file systems that are available for remote mounting
  - An access list is associated with each name, and indicates which hosts are permitted to mount that file system
- On each client
  - A modified version of the UNIX *mount* command mounts a remote file system
    - Based on RPC — specifies remote host name, pathname of a directory in the remote file system, and local name where it is to be mounted
    - Mount requests are usually performed when the system is initialized (booted)
      - Usually specified in */etc/fstab*
    - User may also be able to mount other remote file systems

8

Fall 2005, Lecture 14

## Mounting Remote File Systems (cont.)

### ■ Automounting

- The *automounter* dynamically mounts a file system whenever an “empty” mount point is referenced by a client
  - Further accesses do not result in further requests to the automounter...
  - Unless there are no references to the remote file system for several minutes, in which case the automounter unmounts it

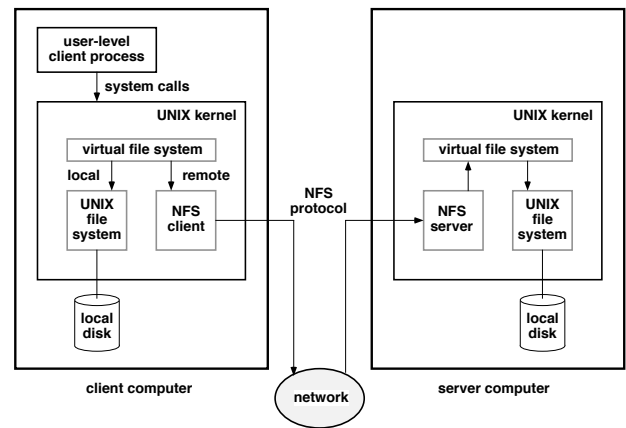
### ■ Remote file systems may be

- *Hard mounted* — when a user-level process accesses a file, it is suspended until the request can be completed
  - If a server crashes, the user-level process will be suspended until recovers
- *Soft mounted* — after a small number of retries, the NFS client returns a failure code to the user process
  - Most UNIX utilities don't check this code...

9

Fall 2005, Lecture 14

## NFS Software Architecture



### ■ Virtual file system:

- Separates generic file-system operations from their implementation (can have different types of local file systems)
- Based on a file descriptor called a vnode that is unique networkwide (UNIX inodes are only unique on a single file system)

10

Fall 2005, Lecture 14

## NFS Protocol

### ■ NFS protocol provides a set of RPCs for remote file operations

- Looking up a file within a directory
- Manipulating links and directories
- Creating, renaming, and removing files
- Getting and setting file attributes
- Reading and writing files

### ■ NFS is stateless

- Servers do not maintain information about their clients from one access to the next
  - There are no open-file tables on the server
- There are no open and close operations
  - Each request must provide a unique file identifier, and an offset within the file
- Easy to recover from a crash, but file operations must be idempotent

11

Fall 2005, Lecture 14

## NFS Protocol (cont.)

- Because NFS is stateless, all modified data must be written to the server's disk before results are returned to the client
  - Server crash and recovery should be invisible to client — data should be intact
  - ✗ Lose benefits of caching
    - Solution — RAM disks with battery backup (un-interruptable power supply), written to disk periodically
- A single NFS write is guaranteed to be atomic, and not intermixed with other writes to the same file
  - However, NFS does not provide concurrency control
    - A write system call may be decomposed into several NFS writes, which may be interleaved
    - Since NFS is stateless, this is not considered to be an NFS problem

12

Fall 2005, Lecture 14

## Distributed Naming Structures (cont.)

- Single name space for remote and local directories
  - Names of form: */.../machinelfs/pathname*
  - Used by:
    - CMU's Andrew, now in OSF's Distributed Computing Environment (DCE)
    - Berkeley's Sprite
  - File names are always the same, whether file is remote or local
  - As clients access a file, the server sends a copy to the client's workstation, and the workstation caches the file
    - In Andrew, local disks are used
    - In Sprite, large memories are used, and workstations are diskless
    - More details on these two next time...
  - Location independent, not location transparent

## CMU's Andrew File System

- Designed by Carnegie Mellon University
  - Developed during mid-1980s as part of the Andrew distributed computing environment
  - Designed to support a WAN of more than 5000 workstations
  - Much of the core technology is now part of the Open Software Foundation (OSF) Distributed Computing Environment (DCE), available for most UNIX and some other operating systems
- Provides transparent access to remote files on a WAN, for clients running on UNIX and other operating systems
  - Access to all files is via the usual UNIX file primitives
  - Compatible with NFS — servers can mount NFS file systems