

***Computer Science  
425/ECE 428/CSE  
424  
Distributed Systems  
(Fall 2009)***

**Lecture 19**  
**Distributed File Systems**  
**Reading: Chapter 8**

# ***Acknowledgement***

- **The slides during this semester are based on ideas and material from the following sources:**
  - Slides prepared by Professors M. Harandi, J. Hou, I. Gupta, N. Vaidya, Y-Ch. Hu, S. Mitra.
  - Slides from Professor S. Gosh's course at University of Iowa.

# ***Administrative***

- **MP2 posted October 5, 2009, on the course website,**
  - **Deadline November 6 (Friday)**
  - **Demonstration, 4-6pm, 11/6/2009**
  - **Tutorial for MP2 planned for October 28 evening if students send questions to TA by October 25. Send requests what you would like to hear in the tutorial.**

# ***Administrative***

- **MP3 proposal instructions**
  - **Deadline for MP3 proposal: October 25, 2009, email proposal to TA**
  - **At least one representative of each group meets with instructor or TA during October 26-28 during their office hours ) watch for extended office hours during these days.**
  - **Wednesday, October 28, 8:30-10am – instructor's office hours 3104 SC**
  - **No office hours, Thursday, 29, 9-10am**

# ***Administrative***

- **Homework 3 posted on Thursday, October 15**
  - Deadline: **Thursday, October 29, 2009** at the beginning of class
- **Midterm Re-grading Period by Instructor – additional office hours:**
  - October 27, 3:15-4pm – in 3104 SC
  - October 29, 3:15-4pm – in 3104 SC

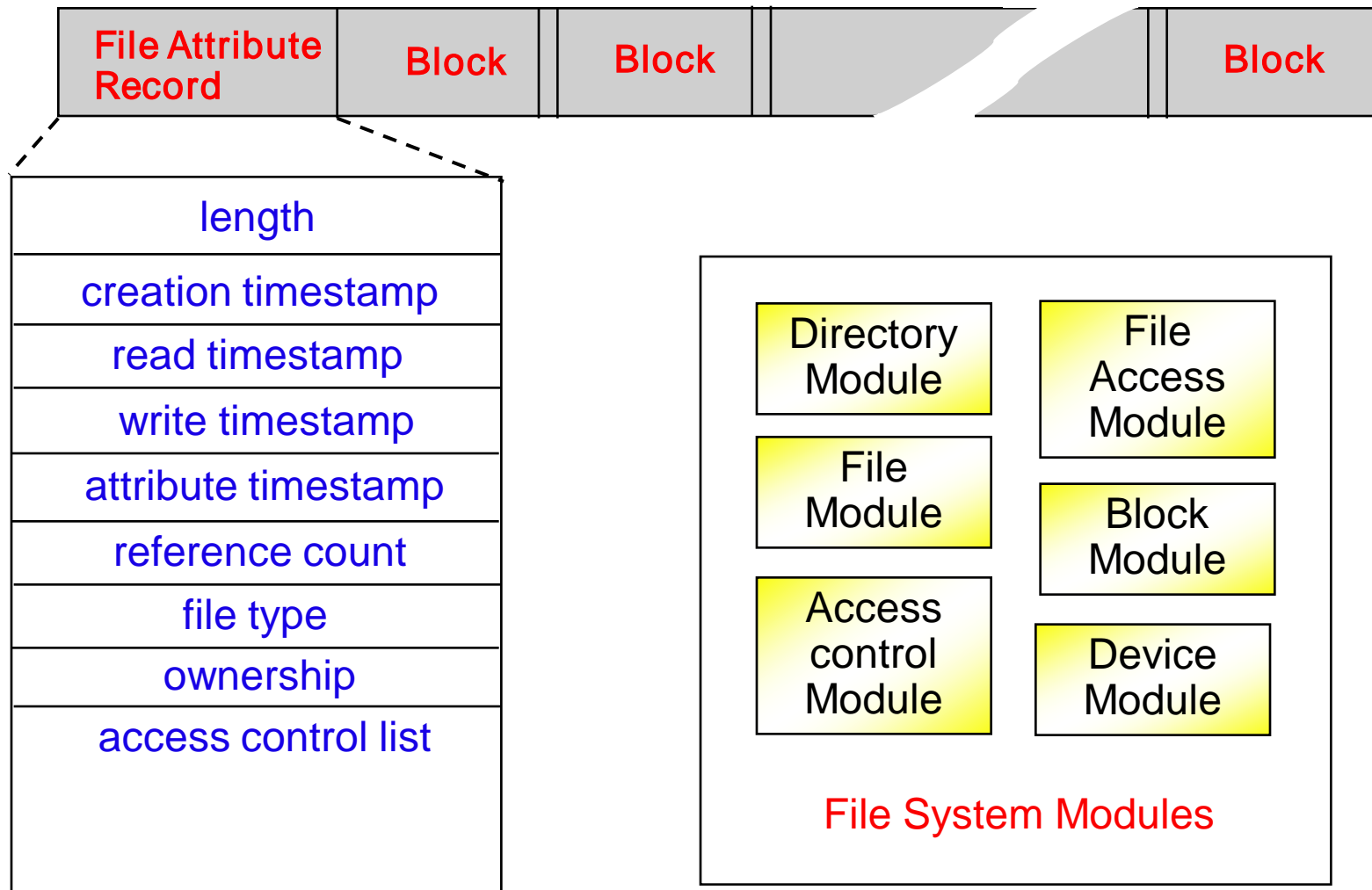
# ***Plan for Today***

- **File Systems – Review**
- **Distributed File Systems – Requirements**
- **File System Architecture**
- **Network File System (NFS)**
- **Andrew File System (AFS)**

# ***File Systems***

- ❖ A **file** is a collection of data with a user view (file structure) and a physical view (blocks).
- ❖ A directory is a file that provides a mapping from text names to internal file identifiers.
- ❖ **File systems** implement file management:
  - ❖ **Naming and locating** a file
  - ❖ **Accessing a file** – create, delete, open, close, read, write, append, truncate
  - ❖ **Physical allocation** of a file.
  - ❖ **Security and protection** of a file.
- ❖ A **distributed file system** (DFS) is a file system with distributed storage *and* distributed users. Files may be located remotely on servers, and accessed by multiple clients.
  - ❖ E.g., SUN NFS and AFS
- ❖ DFS provides **transparency** of location, access, and migration of files.
- ❖ DFS systems use **cache replicas** for efficiency and fault tolerance

# ***File Attributes & System Modules***





# ***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

(Single host File system. DFS may require additional components.)  
Layered architecture: each layer depends only on the layers below it.

# 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 <u>read-write pointer</u> .
<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 links to it, it is deleted from disk.
<i>status</i> = <i>link</i> ( <i>name1</i> , <i>name2</i> )	Creates a new link ( <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> .

---

# ***Distributed File System (DFS) Requirements***

- ❖ **Transparency** - server-side changes should be invisible to the client-side.
  - ❖ *Access transparency*: A single set of operations is provided for access to local/remote files.
  - ❖ *Location Transparency*: All client processes see a uniform file name space.
  - ❖ *Migration Transparency*: When files are moved from one server to another, users should not see it
  - ❖ *Performance Transparency*
  - ❖ *Scaling Transparency*
- ❖ **File Replication**
  - ❖ A file may be represented by several copies for service efficiency and fault tolerance.
- ❖ **Concurrent File Updates**
  - ❖ Changes to a file by one client should not interfere with the operation of other clients simultaneously accessing the same file.

# DFS Requirements (2)

## ❖ Concurrent File Updates

- ❖ **One-copy update** semantics: the file contents seen by all of the processes accessing or updating a given file are those they would see if only a single copy of the file existed.

## ❖ Fault Tolerance

- ❖ At most once invocation semantics.
- ❖ At least once semantics. OK for a server protocol designed for idempotent operations (i.e., duplicated requests do not result in invalid updates to files)

## ❖ Security

- ❖ **Access Control list** = per object, list of allowed users and access allowed to each
- ❖ **Capability list** = per user, list of objects allowed to access and type of access allowed (could be different for each (user,obj))
- ❖ User **Authentication**: need to authenticate requesting clients so that access control at the server is based on correct user identifiers.

## ❖ Efficiency

- ❖ Whole file v.s. block transfer

# ***Basic File Service Model***

## ❖ **An abstract model :**

### ❖ **Flat file service**

❖ implements create, delete, read, write, get attribute, set attribute and access control operations.

### ❖ **Directory service: is itself a client of (i.e., uses) flat file service.**

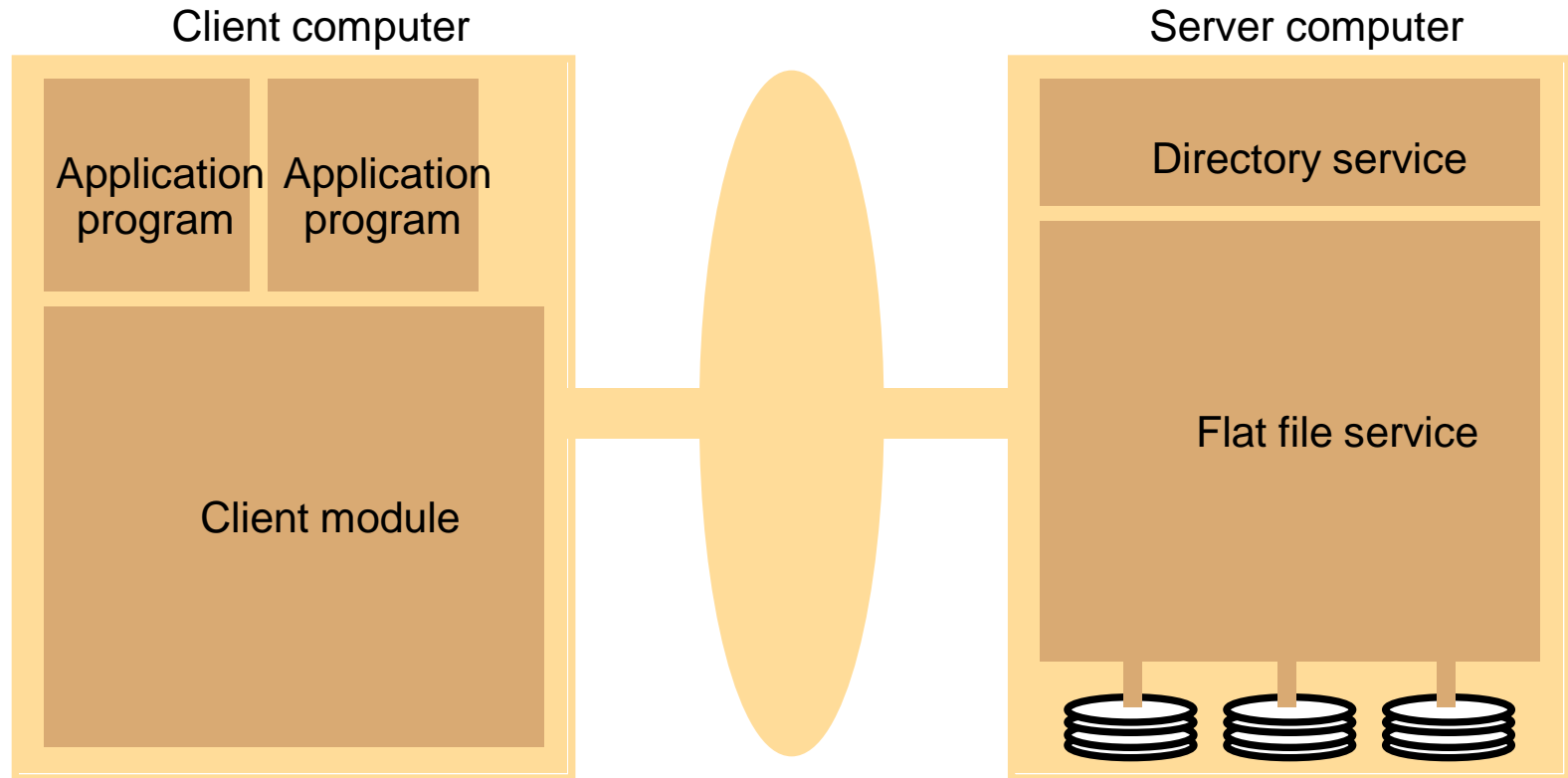
❖ Creates and updates directories (hierarchical file structures) and provides mappings between user names of files and the unique file ids in the flat file structure.

### ❖ **Client service: A client of directory and flat file services**

❖ Runs in each client's computer, integrating and expanding flat file and directory services to provide a unified API (e.g., the full set of UNIX file operations).

❖ Holds information about the locations of the flat file server and directory server processes.

# ***File Service Architecture***



# Flat File Service Operations

---

<i>Read(FileId, i, n) -&gt; 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() -&gt; 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) -&gt; 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 ).

---

- (1) Repeatable operation: No read-write pointer. Except for Create and delete, the operations are idempotent, allowing the use of at least once RPC semantics.
- (2) Stateless servers: No file descriptors. Stateless servers can be restarted after a failure and resume operation without the need to restore any state.

In contrast, the UNIX file operations are neither idempotent nor consistent, because

- (a) a read-write pointer is generated by the UNIX file system whenever a file is opened.
- (b) If an operation is accidentally repeated, the automatic advance of the read/write pointer results in access to different positions of the file.

# Access Control

- In UNIX, the user's access rights are checked against the access mode requested in the open call and the file is opened only if the user has the necessary rights.
- In DFS, a user identity has to be passed with requests – server first authenticates the user.
  - An access check is made whenever a file name is converted to a UFID (unique file id), and the results are encoded in the form of a **capability** which is returned to the client for future access.
    - » Capability = per user, list of objects allowed to access and type of access allowed (could be broken up per (user,obj))
  - A user identity is submitted with every client request, and an access check is performed for every file operation.



# 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*, *File*)  
— 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*. Like *grep*.

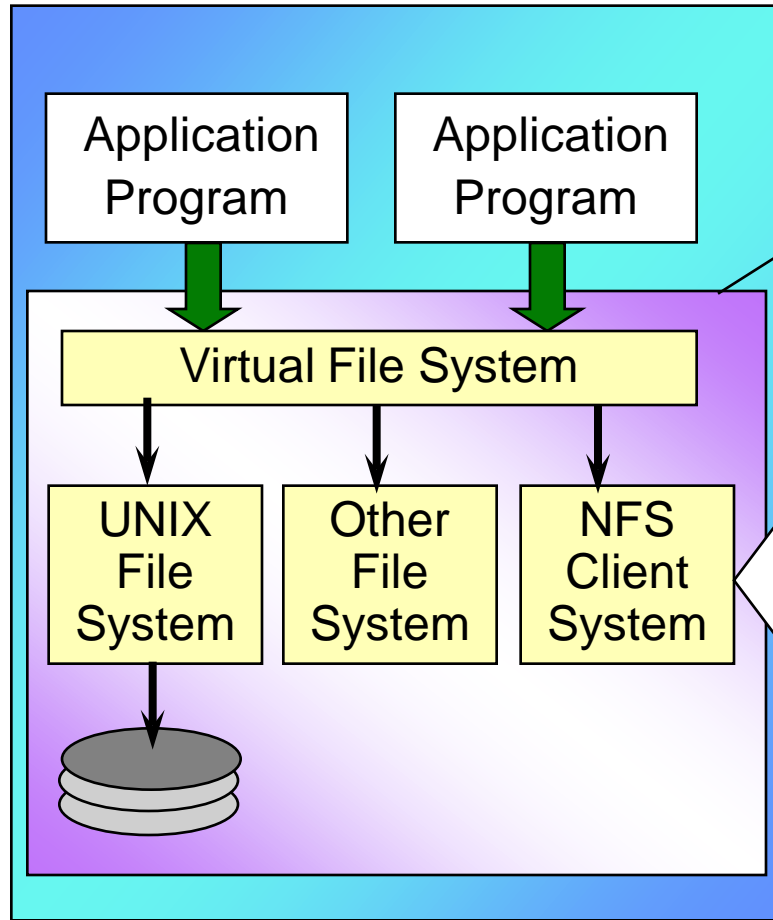
---

(1) Hierarchic file system: The client module provides a function that gets the UFID of a file given its pathname. The function interprets the pathname starting from the root, using *Lookup* to obtain the UFID of each directory in the path.

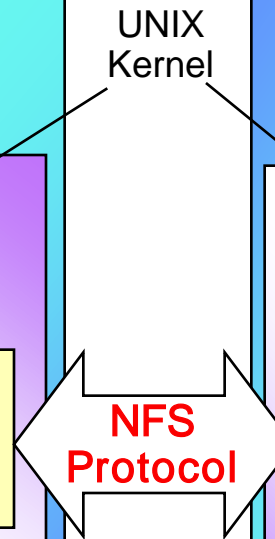
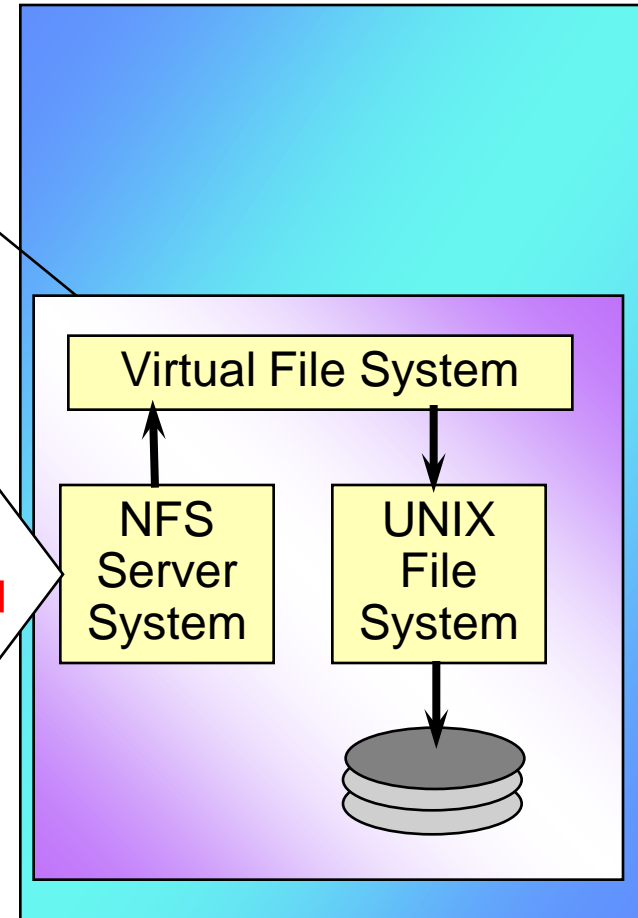
(2) Each server may hold several *file groups*, each of which is a collection of files located on the server. A file group identifier consists of IP address + date, and allows (i) file groups to migrate across servers, and (ii) clients to access file groups.

# Network File System (NFS)

Client Computer



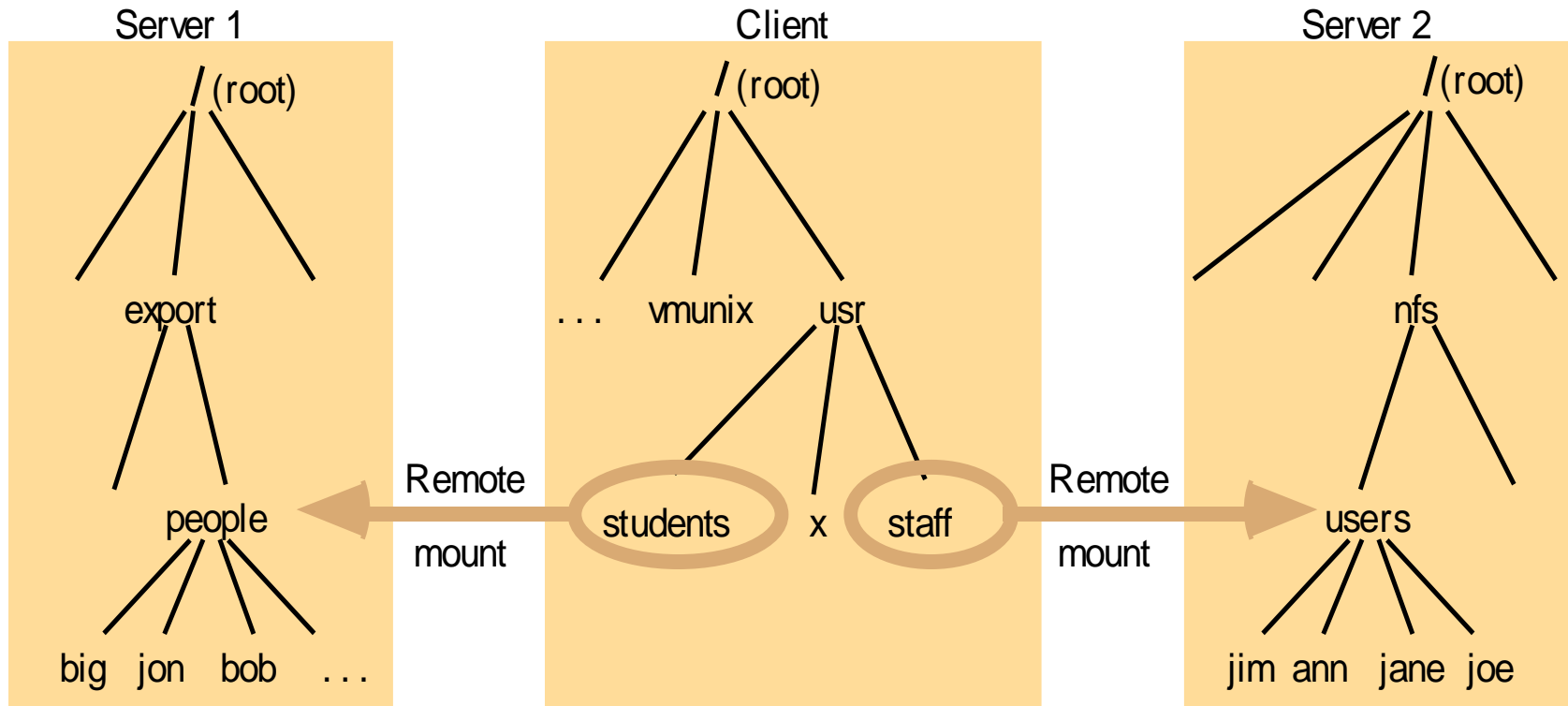
Server Computer



# NFS Architecture -- VFS

- Virtual file system module
  - Translates between NFS file identifiers and other file systems's (e.g., UNIX) identifiers.
    - » The NFS file identifiers are called *file handles*.
    - » File handle = *Filesystem/file group* identifier + *i-node* number of file + i-node generation number.
  - Keeps track of filesystems (i.e., NFS file groups, different from a “file system”) that are available locally and remotely.
    - » The client obtains the first file handle for a remote filesystem when it first *mounts* the filesystem. File handles are passed from server to client in the results of lookup, create, and mkdir operation.
  - Distinguishes between local and remote files.
    - » VFS keeps one VFS structure for each mounted filesystem and one v-node per open file.
      - A VFS structure relates a remote filesystem to the local directory on which it is mounted.
      - A v-node contains an indicator to show whether a file is local or remote. If the file is local, it contains a reference to the i-node; otherwise, it contains the file handle of the remote file if the file is remote.

# Local and Remote File Systems Accessible on an NFS client



**Hard mounting** (retry f.s. request on failure)

vs.

**Soft mounting** (return error on f.s. access failure) – Unix is more compatible with hard mounting

# ***NFS Client and Server***

- **Client**

- Plays the role of the client module from the basic/vanilla model.
- Integrated with the kernel, rather than being supplied as a library.
- Transfers blocks of files to and from server via RPC. Caches the blocks in the local memory.

- **Server**

- Provides a conventional RPC interface at a well-known port on each host.
- Plays the role of file and directory service modules in the architectural model.
- Mounting of sub-trees of remote filesystems by clients is supported by a separate mount service process on each NFS server.

# NFS Server Operations (simplified) – 1

---

<i>lookup(dirfh, name) -&gt; 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) -&gt; 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) -&gt; attr</i>	Returns file attributes of file <i>fh</i> . (Similar to the UNIX <i>stat</i> system call.)
<i>setattr(fh, attr) -&gt; 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) -&gt; 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) -&gt; 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) -&gt; 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) -&gt; 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> .

---

# **Server Caching**

- File pages, directories and file attributes that have been read from the disk retained in a *main memory buffer cache*.
- *Read-ahead* anticipates read accesses and fetches the pages following those that have most recently been read.
- In *delayed-write*, when a page has been altered, its new contents are written back to the disk only when the buffered page is required for another client.
  - In comparison, Unix *sync* operation writes pages to disk every 30 seconds
- In *write-through*, data in write operations is stored in the memory cache at the server immediately and written to disk before a reply is sent to the client.
  - *Better strategy to ensure data integrity even when server crashes occur. More expensive.*



# Client Caching

- A timestamp-based method is used to validate cached blocks before they are used.
- Each data item in the cache is tagged with
  - $T_c$ : the time when the cache entry was last validated.
  - $T_m$ : the time when the block was last modified at the server.
  - A cache entry at time  $T$  is valid if  $(T - T_c < t)$  or  $(Tm_{client} = Tm_{server})$ .
  - $t$ =freshness interval
    - » *Compromise between consistency and efficiency*
    - » *Sun Solaris:  $t$  is set adaptively between 3-30 seconds for files, 30-60 seconds for directories*

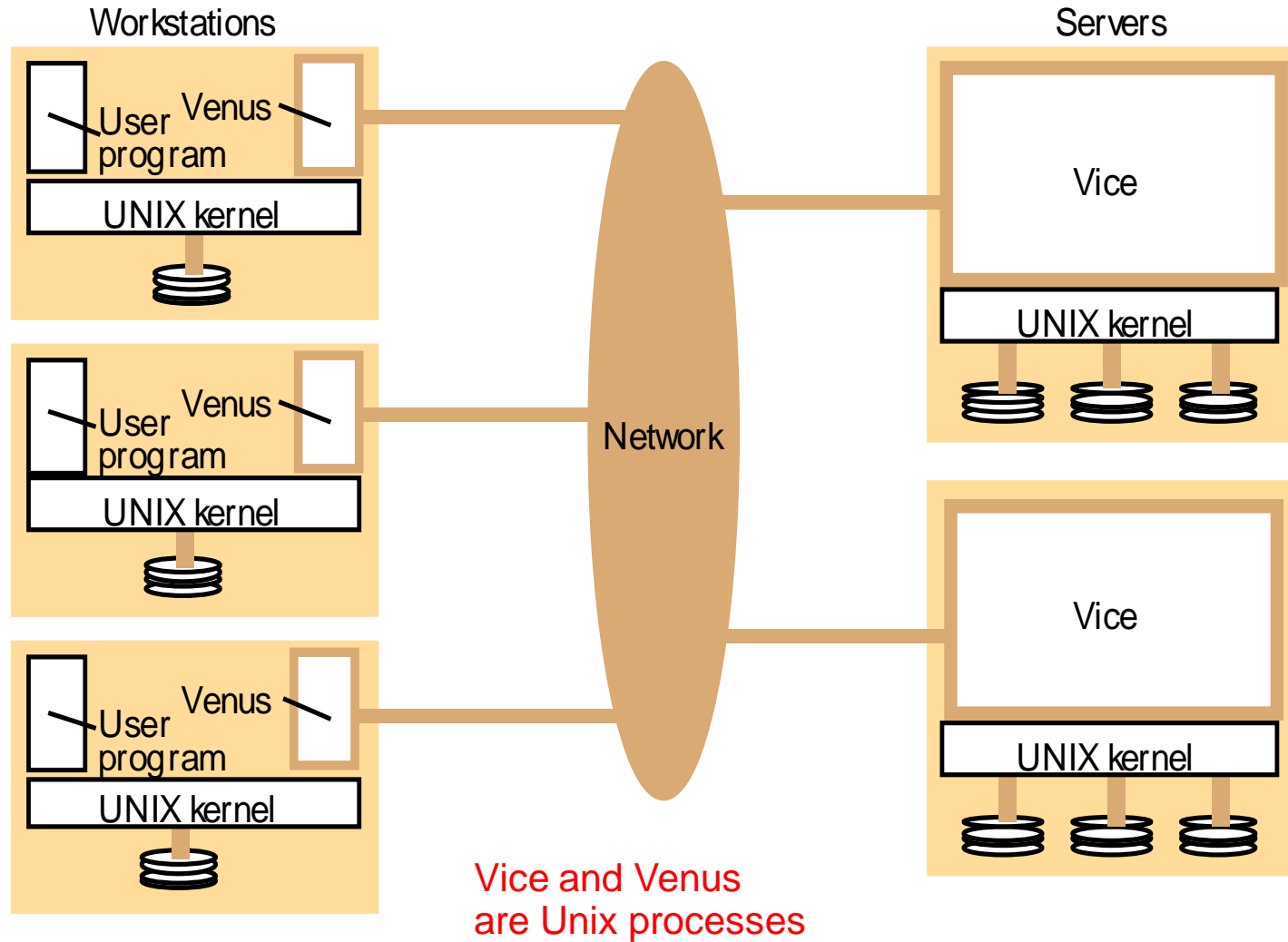
# ***Client Caching (Cont'd)***

- **When a cache entry is read, a validity check is performed.**
  - If the first half of validity condition (previous slide) is true, the second half need not be evaluated.
  - If the first half is not true,  $Tm_{server}$  is obtained (via *getattr()* to server) and compared against  $Tm_{client}$
- **When a cached page (not the whole file) is modified, it is marked as dirty and scheduled to be flushed to the server.**
  - Modified pages are flushed when the file is closed or a *sync* occurs at the client.
- **Does not guarantee one-copy update semantics.**
- **More details in textbook – please read up**

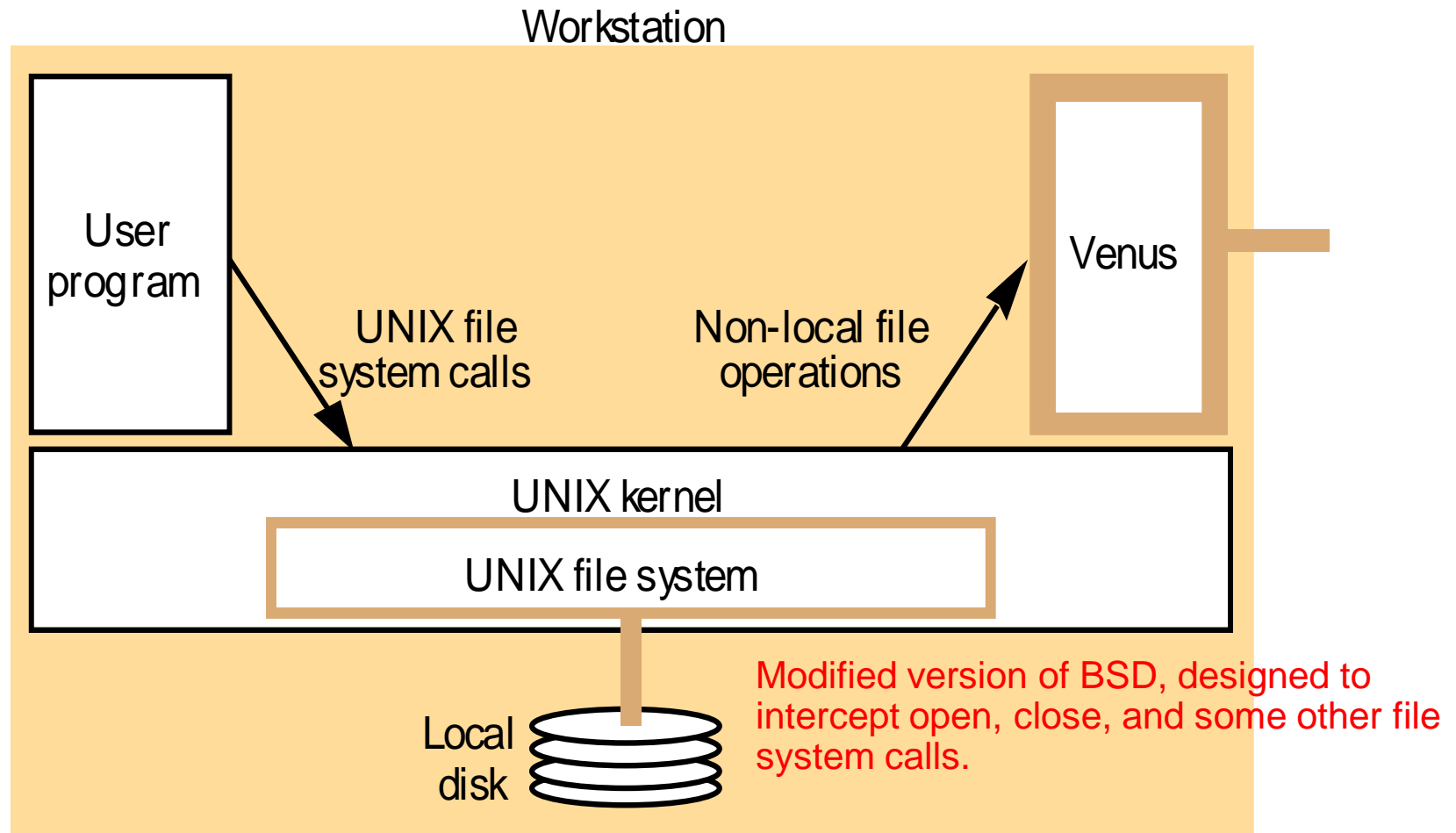
# ***Andrew File System (AFS)***

- **Two unusual design principles:**
  - Whole file serving
    - » Not in blocks
  - Whole file caching
    - » Permanent cache, survives reboots
- **Based on (validated) assumptions that**
  - Most file accesses are by a single user
  - Most files are small
  - Even a client cache as “large” as 100MB is supportable (e.g., in RAM)
  - File reads are much more often than file writes, and typically sequential
- **We'll see overview only**

# ***Distribution of Processes in the Andrew File System***



# System Call Interception in AFS



# Implementation of File System Calls in AFS

<i>User process</i>	<i>UNIX kernel</i>	<i>Venus</i>	<i>Net</i>	<i>Vice</i>
<i>open(FileName, mode)</i>	<p>If <i>FileName</i> refers to a file in shared file space, pass the request to Venus.</p> <p>Open the local file and return the file descriptor to the application.</p>	<p>Check list of files in local cache. If not present or there is no valid <i>callback promise</i>, send a request for the file to the Vice server that is custodian of the volume containing the file.</p> <p>Place the copy of the file in the local file system, enter its local name in the local cache list and return the local name to UNIX.</p>		<p>Transfer a copy of the file and a <i>callback promise</i> to the workstation. Log the callback promise.</p>
<i>read(FileDescriptor, Buffer, length)</i>	Perform a normal UNIX read operation on the local copy.			
<i>write(FileDescriptor, Buffer, length)</i>	Perform a normal UNIX write operation on the local copy.			
<i>close(FileDescriptor)</i>	Close the local copy and notify Venus that the file has been closed.	<p>If the local copy has been changed, send a copy to the Vice server that is the custodian of the file.</p>		<p>Replace the file contents and send a <i>callback</i> to all other clients holding <i>callback promises</i> on the file.</p>

# ***Summary***

- **Distributed File system requirements – transparency, etc.**
- **NFS and AFS**
- **Vnodes (NFS), mounting, caching, whole file caching (AFS)**