

---

# Distributed File Systems



# Distributed File Systems

---

- Distributing file access across multiple nodes
  - single homogeneous large file system

**NFS: Network File System**

**AFS: Andrew File System**



# Distributed Disk Systems

---

**-Distributing data over multiple disks**

**- higher disk access bandwidth**

**- higher reliability**

**RAID: Reliable Array of Inexpensive Disks**

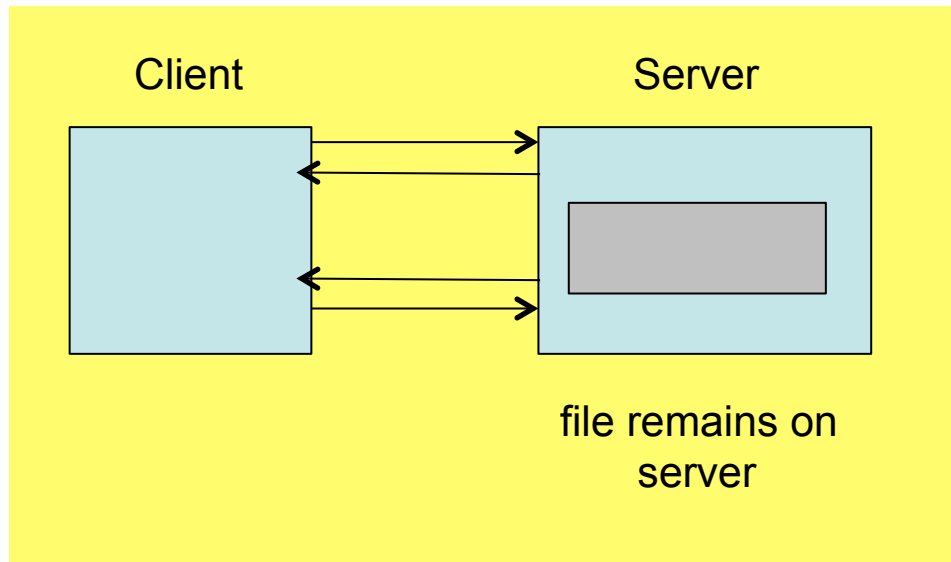


# Distributed File Systems

---

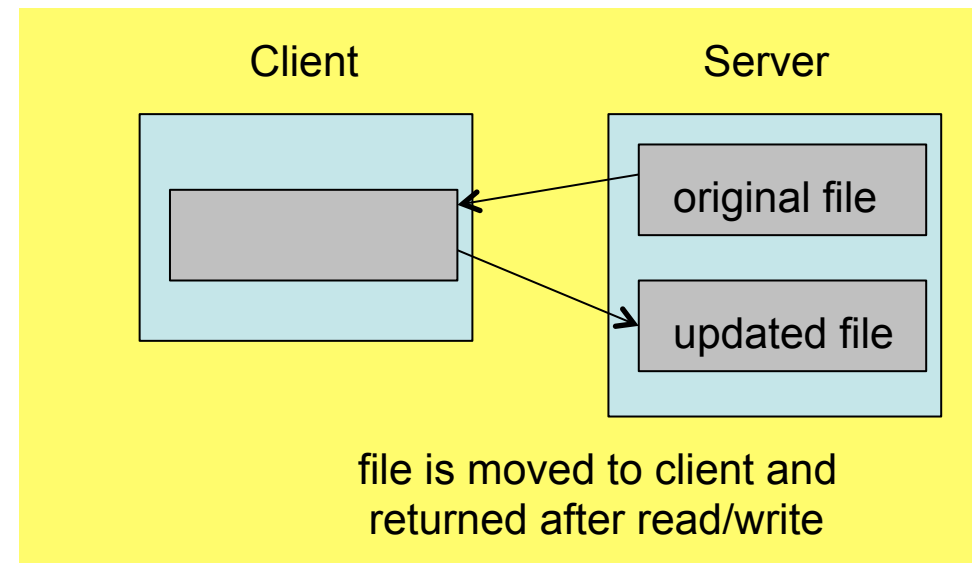
## Models of Remote Access

remote access model



remote execution  
of file operations

upload/download model



local execution  
of file operations



# Distributed File Systems

---

Naming and Mounting

↑  
policy

↑  
mechanism

Local Interface to Distributed File System  
**Issues?**

Synchronization and Caching



# Requirements for Distributed File Systems

---

- ➔ **Transparencies (access, location, mobility, performance, scalability)**
- ➔ **Concurrent File Update**
- ➔ **Replication of Files**
- ➔ **Openess (Heterogeneity of OS and Hardware)**
- ➔ **Fault-Tolerance**
- ➔ **Consistency**
- ➔ **Security**
- ➔ **Efficiency**



# Early milestones in distributed file systems

- ➔ D.R. Brownsbridge, L.F. Marshall, B. Randell: "The Newcastle Connection or UNIXes of the World Unite!", Software-Practice and Experience, Vol.12, 1147-1162, 1982
- ➔ B. Walker, G. Propek, R. English, C. Kline, and G. Thiel (UCLA)  
The LOCUS Distributed Operating System  
Proceedings of the Ninth ACM Symposium on Operating Systems Principles, October 10-13, 1983, pages. 49-70
- ➔ R. Sandberg, D. Goldberg, S. Kleinman, D. Walsh  
The Design and Implementation of the SUN Network File System  
Proceedings Usenix Conference, Portland, Oregon 1985
- ➔ J. Morris, M. Satyanarayanan, M.H. Conner, J.H. Howard, D.S. Rosenthal, F.D. Smith  
Andrew: A distributed personal computing environment  
Comm. of the ACM, Vol.29, No. 3, 1986

} **first  
commercial  
system**

**AFS inspired the development of the "Distributed Computing Environment (DCE)"**



# First Approaches: The Newcastle Connection

---

SOFTWARE-PRACTICE AND EXPERIENCE. VOL. 12. 1147-1162 (1982)

## The Newcastle Connection

or

### UNIXes of the World Unite!

D. R. BROWNBRIDGE, L. F. MARSHALL AND B. RANDELL

*Computing Laboratory, The University, Newcastle upon Tyne NE1 7RU, England*

#### SUMMARY

In this paper we describe a software subsystem that can be added to each of a set of physically interconnected UNIX or UNIX look-alike systems, so as to construct a distributed system which is functionally indistinguishable at both the user and the program level from a conventional single-processor UNIX system. The techniques used are applicable to a variety and multiplicity of both local and wide area networks, and enable all issues of inter-processor communication, network protocols, etc., to be hidden. A brief account is given of experience with such a distributed system, which is currently operational on a set of PDP11s connected by a Cambridge Ring. The final sections compare our scheme to various precursor schemes and discuss its potential relevance to other operating systems.



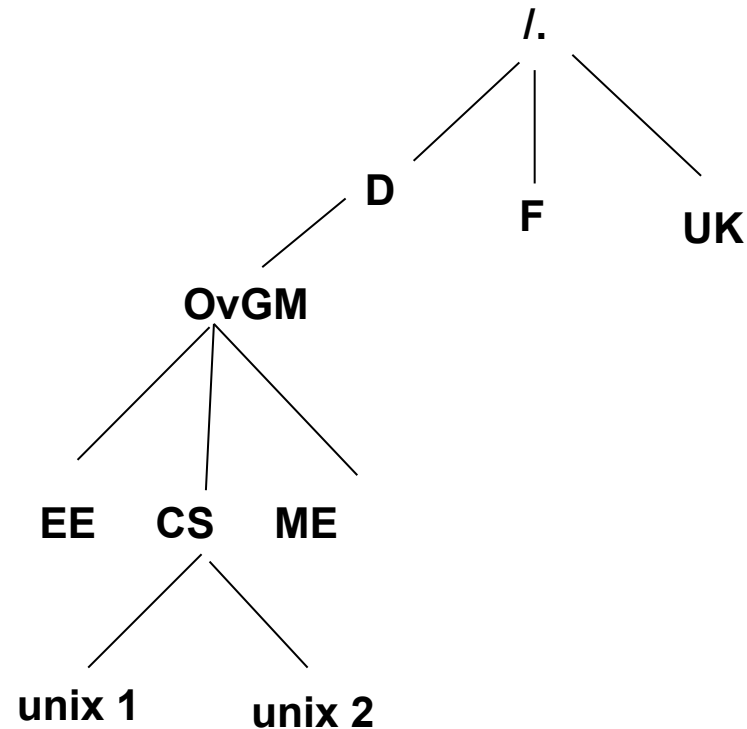
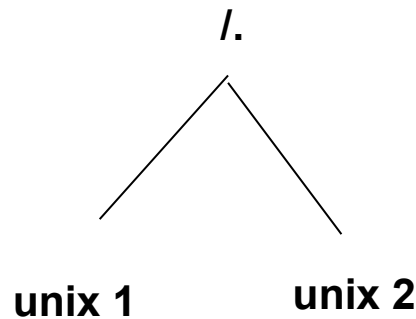


# First Approaches: The Newcastle Connection

---

## Principles:

- Extending the hierarchical Unix Naming Scheme by a "Super Root",
- Using RPC to perform remote file access



# Network File Systems

---

**Newcastle connection provides a single name space for files.**

**Problems with the Newcastle Connection:**

**No Location transparency**

**No Replication or Chaching**

**No Mobility Transparency**



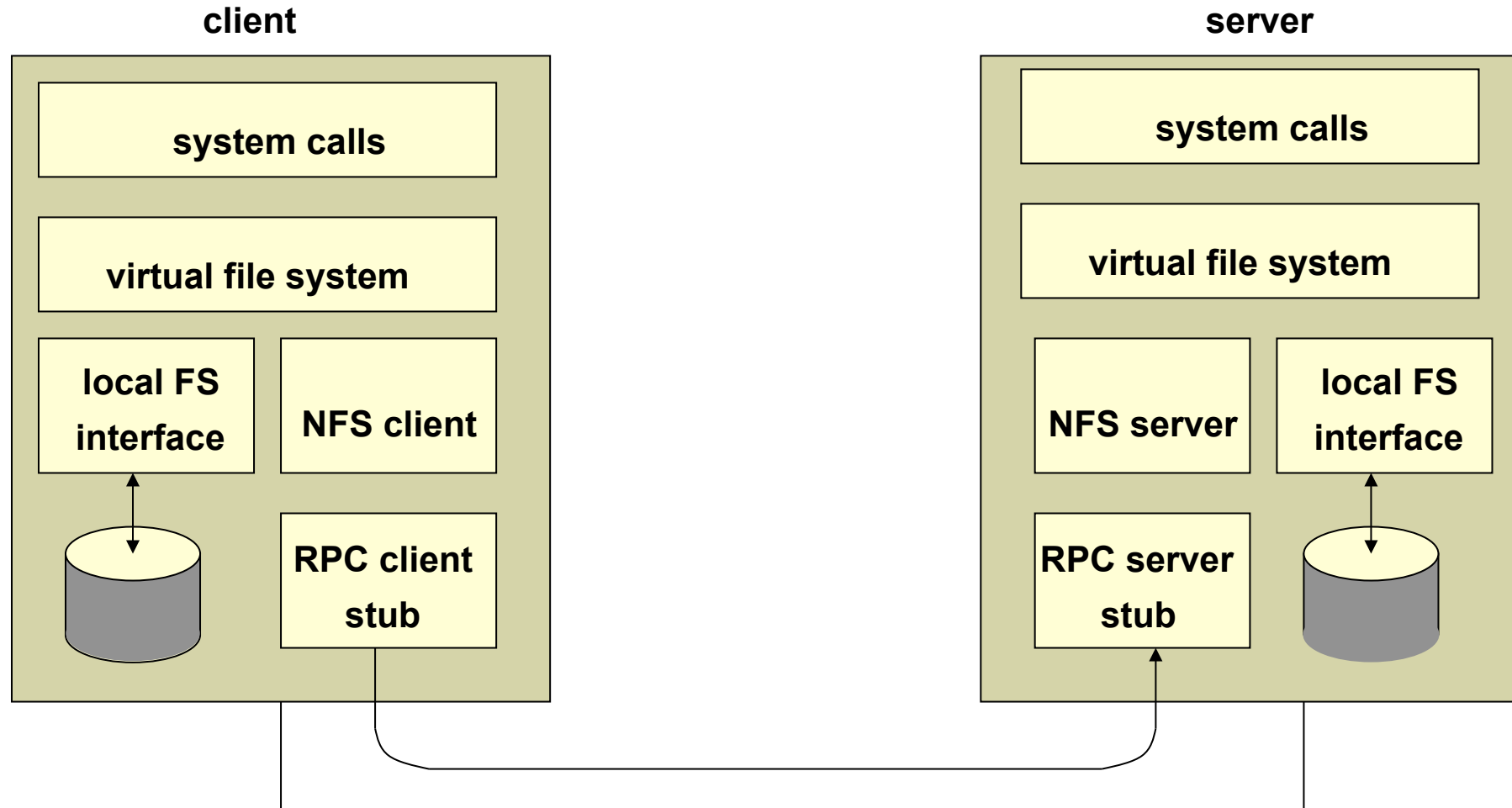
# Network File Service Architecture

---

- **location transparency**
- **migration transparency**
- **robustness against client and server faults**

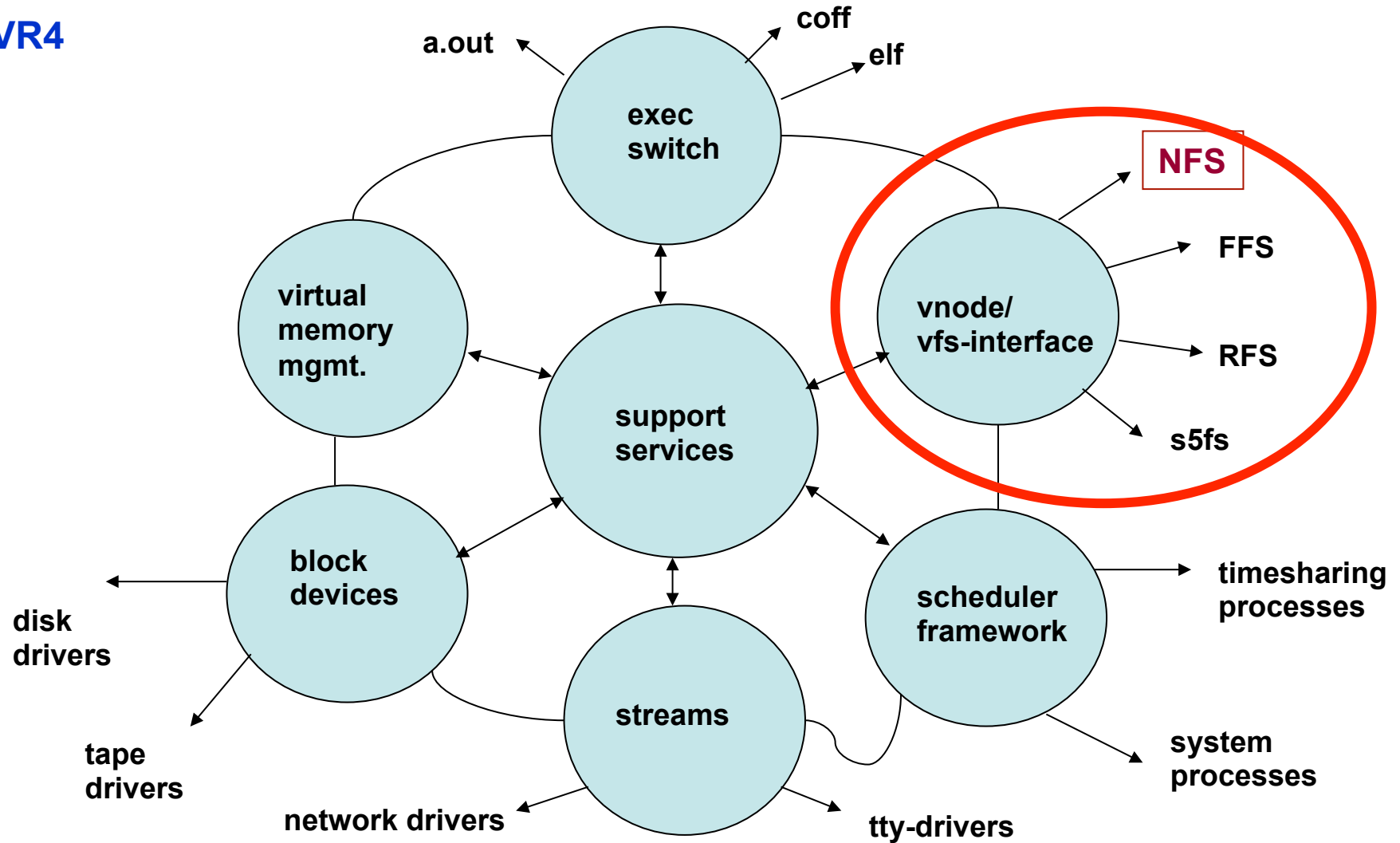


# NFS Client-Server Architectures



# Recall (BS I): Modern Unix-Kernel (Vahalia 1996)

SVR4



# Differences to the Unix File System API

---

## Stateless File Server:

- no state information about open file
- no information about the number and state of clients
  - ➔ every request must be self-contained.

**Benefit:** A client or a server crash does not require extensive recovery activities.

- no open or close
- operations are **idempotent** except "create"



# Flat File Service Operations

---

Used by the client module not used by programs at user level file access!

***Read (FileId, i, n) → Data***  
- throws *BadPosition*

**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 *Data***

***Write (FileId, i, n) → Data***  
- throws *BadPosition*

**If  $1 \leq i \leq \text{Length}(\text{File})+1$ : Writes a sequence of *Data* to a file starting at item  $i$ , extending the file if necessary**

***Create() → FileId***

**Creates a new file of length 0 and delivers a UFID for it.**

***Delete(FileId)***

**Removes a file from the file store.**

***GetAttributes(FileId) → Attr***

**Returns the file attributes for the file.**

***SetAttributes(FileId, Attr)***

**Sets the file attributes for the file (except owner, type and ACL).**



# Directory Service Operations

---

***Lookup (Dir, Name) → FileId***

**- throws *NotFound***

**Locates the text name in the directory and returns the respective UFID. If *Name* is not found, an exception is raised.**

***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. Throws an exception if *Name* is already in the directory.**

***UnName (Dir, Name)***

**- throws *NotFound***

**If *Name* is in the directory it is removed.**

**If *Name* is not in the directory an exception is raised.**

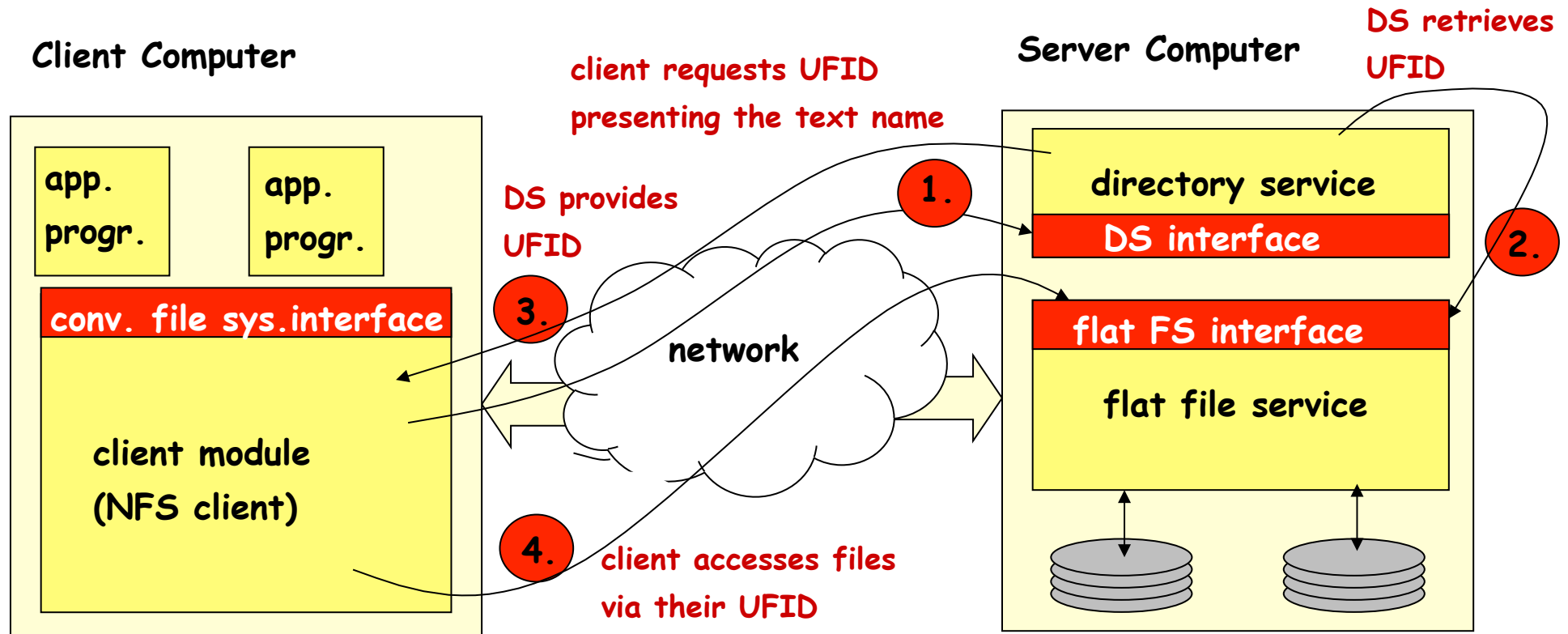
***GetNames (Dir, Pattern) → NameSeq***

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





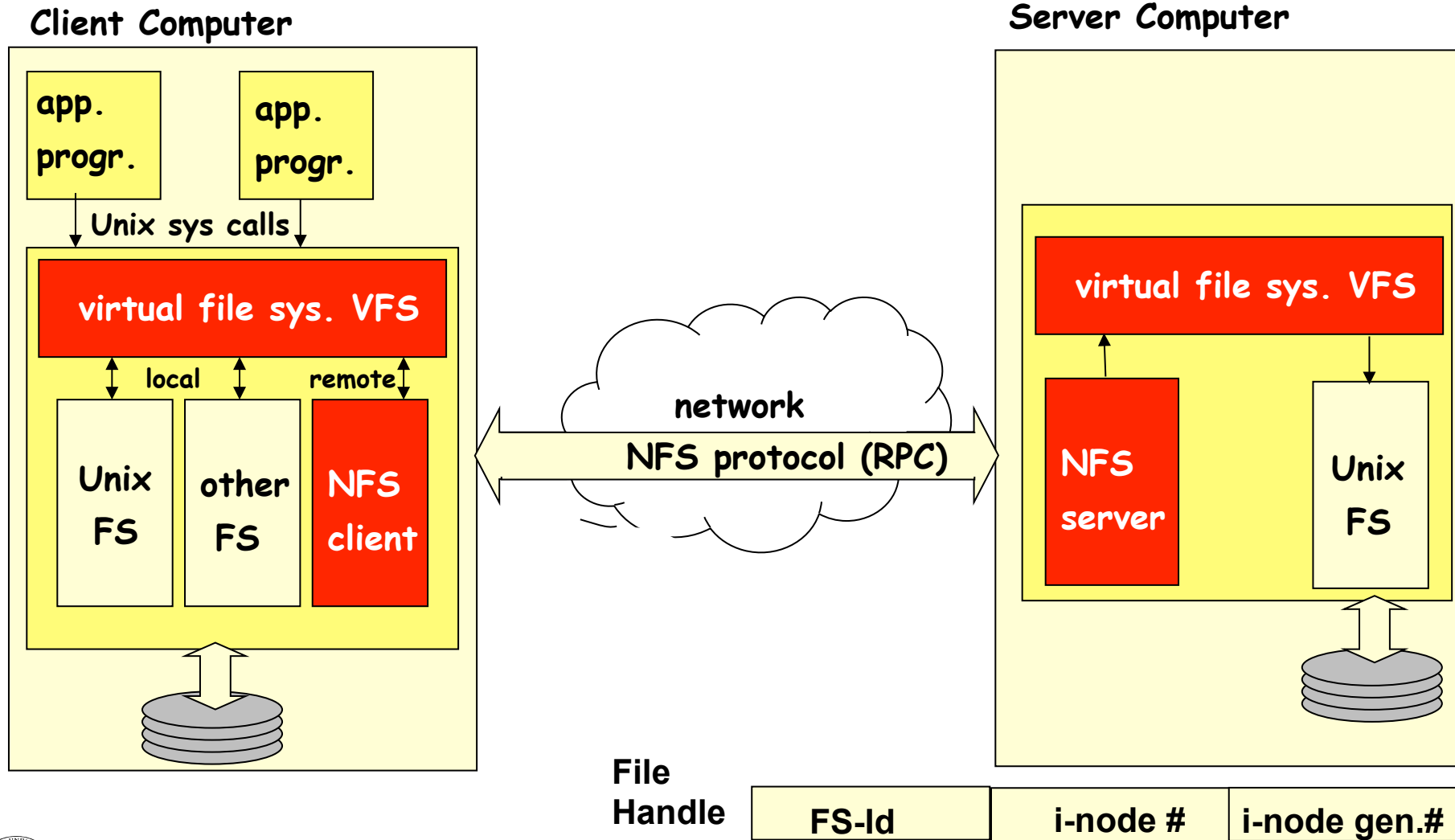
# Network File System Architecture



- ➔ Client-Server architecture using SUN RPC
- ➔ Flat FS uses Unique File IDs (UFIDs) instead of hierarchical path names
- ➔ DS associates file text names with Unique File IDs (UFID)

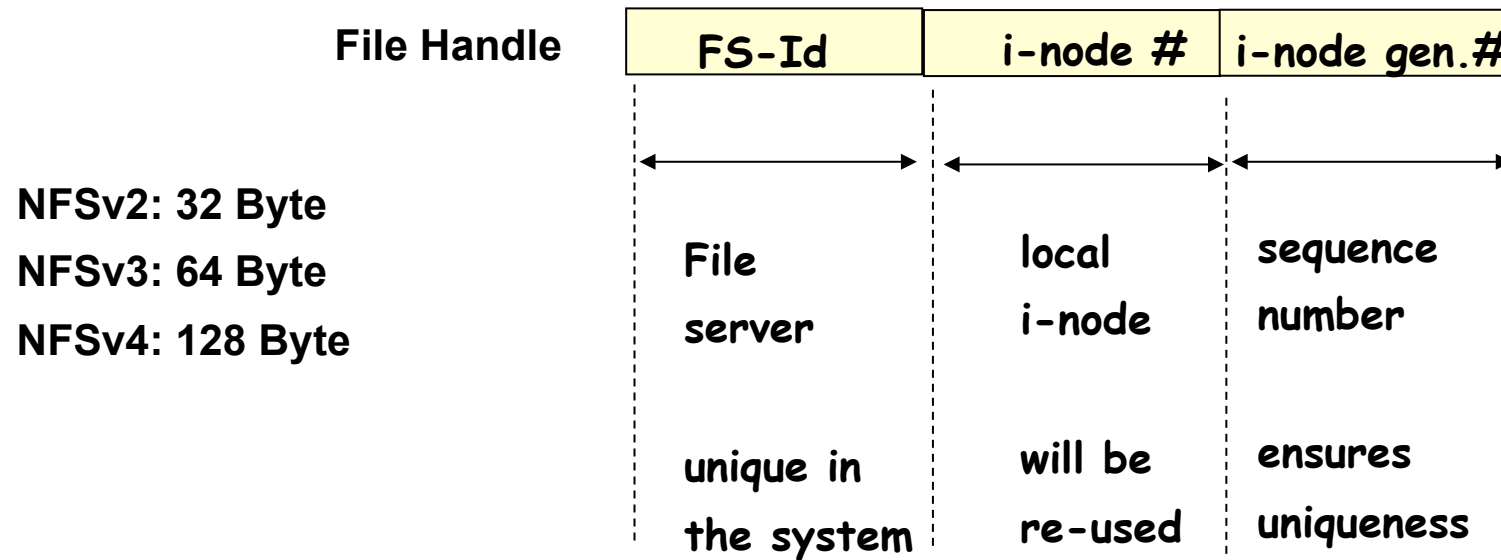


# SUN NFS Architecture



# NFS File Handle

---



The File Handle enables file access to any file in the distributed file system without looking it up in the name server.

How to obtain a file handle in a remote file system subtree?



# NFS server operations (simplified)

---

<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 <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> .

Continues on next slide ...

# NFS server operations (simplified)

---

<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> .

---

# Naming in Network File Systems

---

**Naming distinguishes between:**

- **User-Level Names e.g. UNIX path names (structured ns)**
- **Unique File Identifiers (UFID) System-wide unambiguous number (flat ns)**
  
- **Hierarchical naming system is established using (flat) file system UIDs (UFID), and a directory service.**
  
- **UFIDs support location transparency.**



# NFS mount service

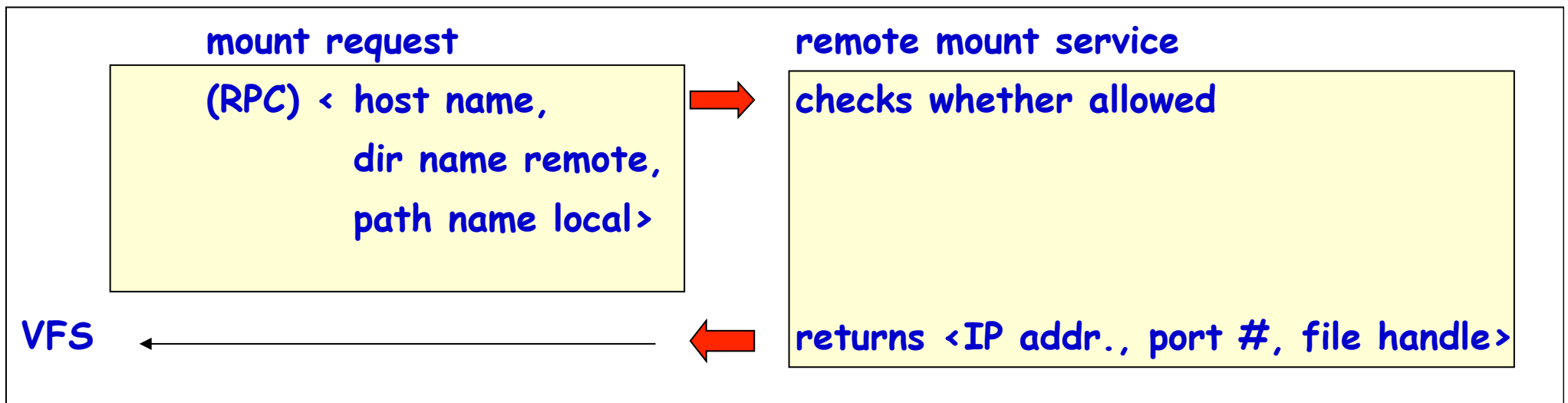
---

**Mount Service Process: executed on every server**

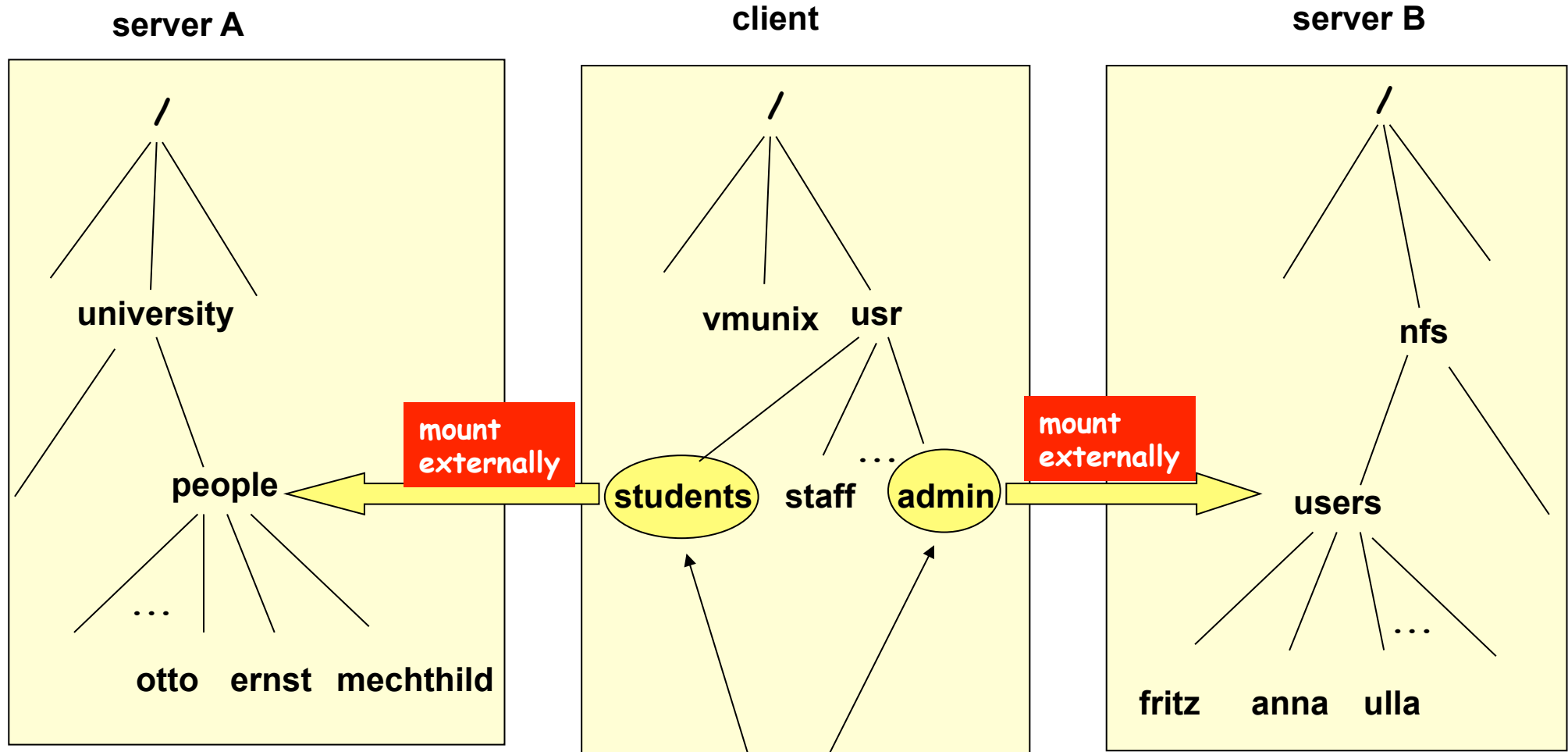
**Data Structures:**

**Server: etc/exports**

contains names of local FS which may be mounted ext.  
For every file system a list of names of (client) hosts is associated which are allowed to mount the FS.



# NFS mount service





# NFS mount service

---

**Hard-Mounted:** requesting application-level service blocks until the request is serviced. Server crashes and subsequent recovery is transparent for the application process.

**Soft-Mounted:** if the request cannot be serviced, the NFS client module signals an error condition to the application.

**Soft-Mounting** needs a meaningful reaction of the application process. In most cases the transparency of the hard-mounting is preferred.



# NFS Server Caching

---

## Standard Unix FS mechanisms

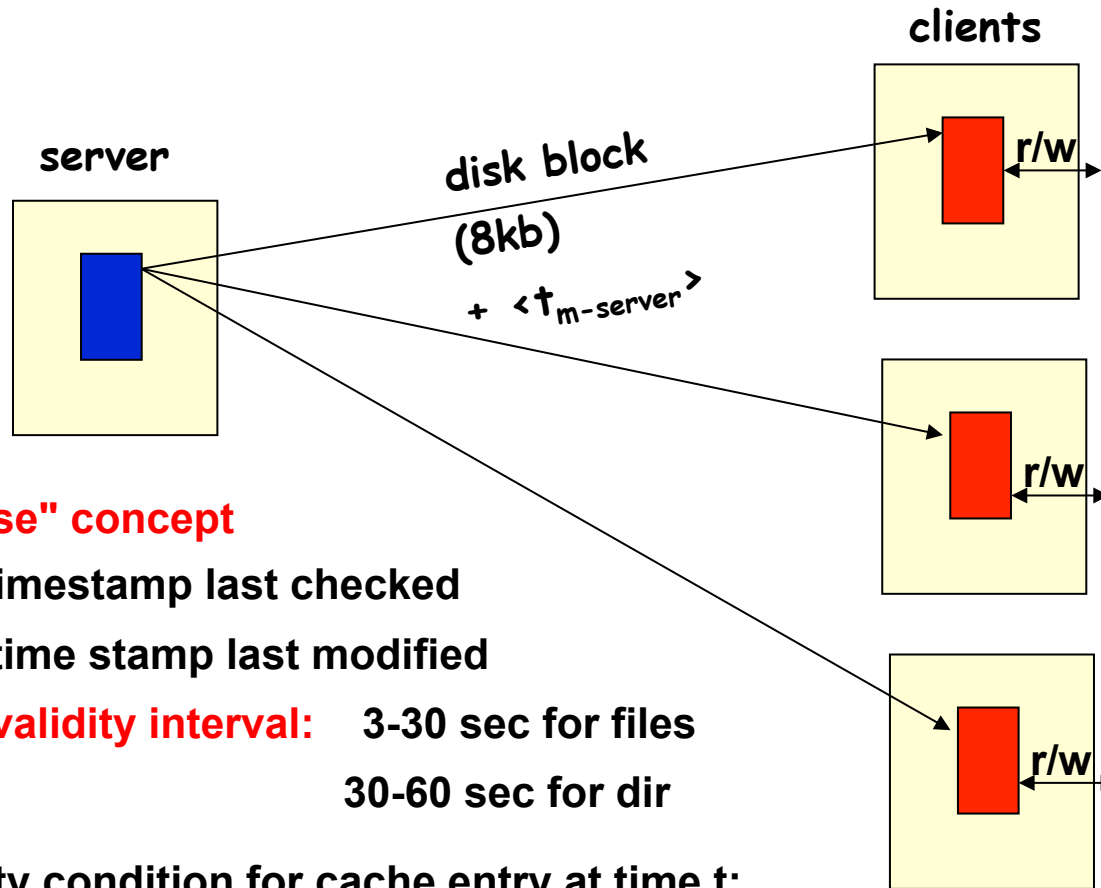
- buffer cache
- read ahead
- delayed write
- sync (periods of 30 sec)

## Additionally: Two options for write (NFS version 3)

- 1.) Data from clients is written to the buffer cache AND the disk (write through). ⇒ Data is persistent when RPC returns.
- 2.) Data will be held in the cache only. Explicit **commit**-operation makes data persistent. Default mode for Standard NFS clients. Commit is issued when closing a file.



# NFS Client Caching



## "lease" concept

$t_c$  : timestamp last checked

$t_m$  : time stamp last modified

$\Delta t$  : **validity interval:** 3-30 sec for files  
30-60 sec for dir

Validity condition for cache entry at time  $t$ :

$$(t - t_c < \Delta t) \vee (t_{m-client} = t_{m-server})$$

## READ:

all **reads** in an interval of  $\Delta t$  after caching only go to the cache. **Reads** occurring after that time check the validity of the copy with the server. If still valid they may use it another  $\Delta t$ .

## WRITE:

cached locally until a sync of the client or if file is closed.

**Mechanism only approximates  
1-Copy-Consistency !**



# Dealing with shared Files

---

**Unix Semantics: Every operation is instantaneously visible to all processes.**

**Session Semantics: No changes are visible to other processes until the file is closed.**

**Immutable files: No updates possible. On update a new file is created.**

**Transactions: All changes are atomic**



# Locking Files

---

<b>Operation</b>	<b>Description</b>
<b>Lock</b>	<b>Create a lock for a range of bytes</b>
<b>Lockt</b>	<b>Test whether a conflicting lock has been created</b>
<b>LockU</b>	<b>Remove a lock from a range of bytes</b>
<b>Renew</b>	<b>Renew the lease on a specified block</b>



# "Share reservations" (NFS 4)

weak form of type-specific access request

requested access	Current file denial state			
	none	read	write	both
read	succeed	fail	succeed	fail
write	succeed	succeed	fail	fail
both	succeed	fail	fail	fail

a client tries to open a file that has a certain denial status under the a certain access status

current access stat	Requested file denial state			
	none	read	write	both
read	succeed	fail	succeed	fail
write	succeed	succeed	fail	fail
both	succeed	fail	fail	fail

a client want to change the access status dynamically when other clients already have access under a certain denial state.



# NFS Properties

---

<b>Access Transparency</b>	<b>++</b>	
<b>Location Transparency</b>	<b>++</b>	
<b>Migration Transparency</b>	<b>+ -</b>	
<b>Scalability</b>	<b>+</b>	
<b>File Replication</b>	<b>+ -</b>	<b>only read replication</b>
<b>Heterogeneity</b>	<b>++</b>	<b>available for many platforms</b>
<b>Fault-Tolerance</b>	<b>+</b>	<b>stateless, restricted fault model</b>
<b>Consistency</b>	<b>+ -</b>	<b>"almost" one copy</b>
<b>Security</b>	<b>-</b>	<b>needs additions (e.g. Cerberos)</b>
<b>Efficiency</b>	<b>++</b>	



---

# Network File System (NFS) version 4 Protocol

<http://www.ietf.org/rfc/rfc3530.txt>





# New features of NFSv4

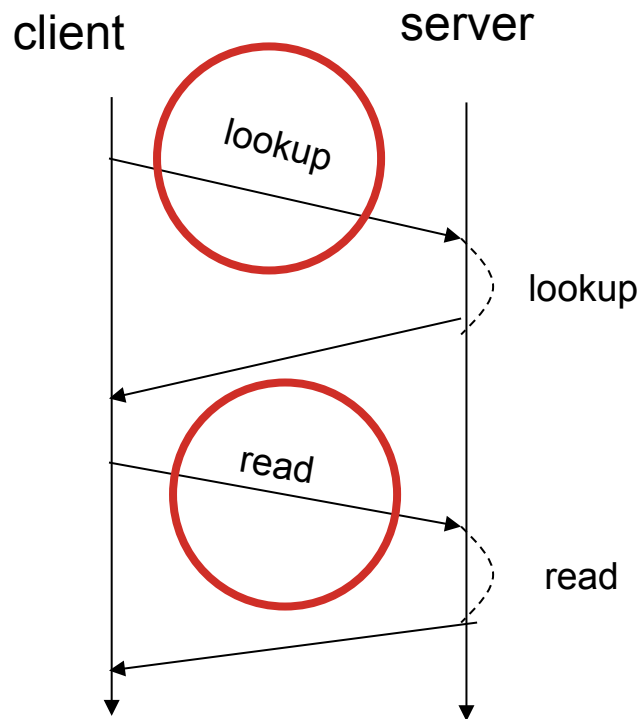
---

- NFSv4 introduces state. NFSv4 is a stateful protocol unlike NFSv2 or NFSv3.
- NFSv4 introduces file delegation. An NFSv4 server can enable an NFSv4 client to access and modify a file in its cache without sending any network requests to the server.
- NFSv4 uses compound remote procedure calls(RPCs) to reduce network traffic. An NFSv4 client can combine several traditional NFS operations (LOOKUP, OPEN, and READ) into a single compound RPC request to carry out a complex operation in one network round trip.
- NFSv4 specifies a number of sophisticated security mechanisms including Kerberos5 and Access Control Lists.
- NFSv4 can seamlessly coexist with NFSv3 and NFSv2 clients and servers.

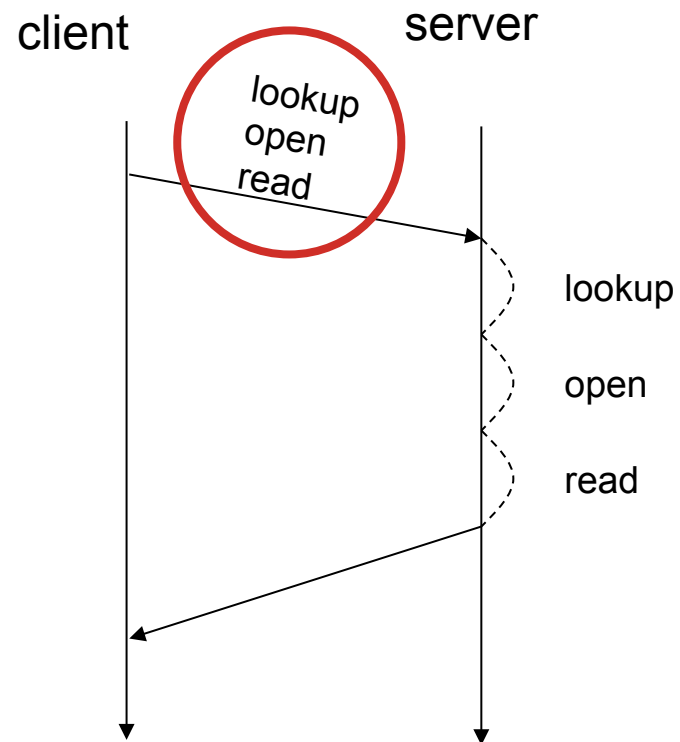


# Compound RPCs in NFS

## NFS V3



## NFS V4



# NFS V4 Compound (mount) Request

```
nfernand@nf734153:/pdfs026/simple2
1 Network File System
2   Program Version: 4
3   V4 Procedure: COMPOUND (1)
4   Tag: mount ←
5     length: 12
6     contents: mount
7   minorversion: 0
8   Operations (count: 5)
9     Opcode: PUTROOTFH (24)
10    Opcode: GETFH (10)
11    Opcode: LOOKUP (15)
12    ...
13    ...
14    Opcode: GETFH (10)
15    Opcode: GETATTR (9)
16      attrmask
17        mand_attr: FATTR4_SUPPORTED_ATTRS (0)
18        mand_attr: FATTR4_TYPE (1)
19        ...
20        ...
```

mount request

header info



```
1 Network File System
2   Program Version: 4
3   V4 Procedure: COMPOUND (1)
4   Status: NFS4_OK (0)
5   Tag: mount
6     length: 12
7     contents: mount
8   Operations (count: 5)
9     Opcode: PUTROOTFH (24)
10      Status: NFS4_OK (0)
11     Opcode: GETFH (10)
12      Status: NFS4_OK (0)
13      ...
14      ...
15     Opcode: LOOKUP (15)
16      Status: NFS4_OK (0)
17     Opcode: GETFH (10)
18      Status: NFS4_OK (0)
19      ...
20      ...
21     Opcode: GETATTR (9)
22      Status: NFS4_OK (0)
23      obj_attributes
24      attrmask
25      mand_attr: FATTR4_SUPPORTED_ATTRS (0)
26      attrmask
27      mand_attr: FATTR4_SUPPORTED_ATTRS (0)
28      mand_attr: FATTR4_TYPE (1)
29      ...
30      ...
```

# (mount) Reply



# NFS V4 setclientid Request

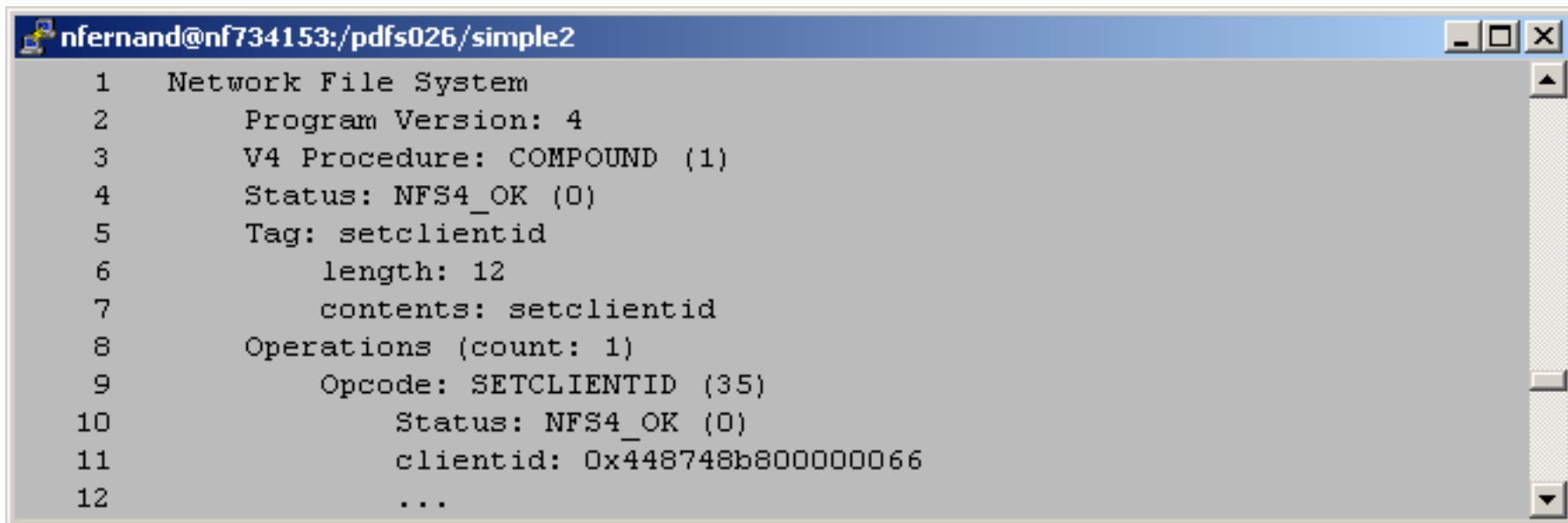
---

```
nfernand@nf734153:/pdfs026/simple2
 2      Program Version: 4
 3      V4 Procedure: COMPOUND (1)
 4      Tag: setclientid
 5          length: 12
 6          contents: setclientid
 7      Operations (count: 1)
 8          Opcode: SETCLIENTID (35)
 9          client
10              ...
11              ...
12          callback
13              cb_program: 0x00000000
14              cb_location
15                  ...
16              callback_ident: 0x00000000
```



# NFS V4 setclientid Reply

---



A terminal window titled "nfernand@nf734153:/pdfs026/simple2" displays the output of an NFS V4 setclientid operation. The output is a list of 12 lines, numbered 1 through 12, showing the details of the reply. The window has a blue title bar and standard window controls (minimize, maximize, close) on the right side.

```
1 Network File System
2   Program Version: 4
3   V4 Procedure: COMPOUND (1)
4   Status: NFS4_OK (0)
5   Tag: setclientid
6     length: 12
7     contents: setclientid
8   Operations (count: 1)
9     Opcode: SETCLIENTID (35)
10    Status: NFS4_OK (0)
11    clientid: 0x448748b800000066
12    ...
```



# NFS V4 Open Request

```
hpdfs026
1 Network File System
2   Program Version: 4
3   V4 Procedure: COMPOUND (1)
4   Tag: open
5     length: 12
6     contents: open
7   minorversion: 0
8   Operations (count: 4)
9     Opcode: PUTFH (22)
10    ...
11    ...
12    Opcode: OPEN (18)
13      seqid: 0x00000001
14      share_access: OPEN4_SHARE_ACCESS_BOTH (3)
15      share_deny: OPEN4_SHARE_DENY_NONE (0)
16      clientid: 0x448748b800000066
17      ...
18      ...
19      Opcode: GETFH (10)
20      Opcode: GETATTR (9)
21      ...
```



# NFS V4 Open Reply

```
hpdfs026
1 Network File System
2   Program Version: 4
3   V4 Procedure: COMPOUND (1)
4   Status: NFS4_OK (0)
5   Tag: open
6     length: 12
7     contents: open
8   Operations (count: 4)
9     Opcode: PUTFH (22)
10      Status: NFS4_OK (0)
11     Opcode: OPEN (18)
12      Status: NFS4_OK (0)
13      stateid
14        seqid: 0x00000001
15        other: 44D52AE40000006500000000
16      ...
17     Opcode: GETFH (10)
18      Status: NFS4_OK (0)
19      ...
20     Opcode: GETATTR (9)
21      Status: NFS4_OK (0)
22     ...
```





**Operation v3 v4**

**Beschreibung**

Create	Ja	Nein	Erstellen einer regulären Datei
Create	Nein	Ja	Erstellen einer irregulären Datei
Link	Ja	Ja	Erstellen einer direkten Verknüpfung zu einer Datei
Symlink	Ja	Nein	Erstellen einer symbolischen Verknüpfung zu einer Datei
Mkdir	Ja	Nein	Erstellen eines Unterverzeichnisses in einem gegebenen Verzeichnis
Mknod	Ja	Nein	Erstellen einer Spezialdatei
Rename	Ja	Ja	Ändern einer Dateibezeichnung
Remove	Ja	Ja	Entfernen einer Datei aus einem Dateisystem
Rmdir	Ja	Nein	Entfernen eines leeren Unterverzeichnisses aus einem Verzeichnis
Open	Nein	Ja	Öffnen einer Datei
Close	Nein	Ja	Schließen einer Datei
Lookup	Ja	Ja	Suchen einer Datei anhand ihrer Bezeichnung
Readdir	Ja	Ja	Lesen der Einträge eines Verzeichnisses
Readlink	Ja	Ja	Auslesen der in einer symbolischen Verknüpfung gespeicherten Pfadangabe
Getattr	Ja	Ja	Auslesen der Attributwerte einer Datei
Setattr	Ja	Ja	Setzen eines oder mehrerer Attributwerte für eine Datei
Read	Ja	Ja	Auslesen der in einer Datei enthaltenen Daten
Write	Ja	Ja	Schreiben von Daten in eine Datei



# AFS Andrew File System

---

**Scalability as primary design goal.**

**As much as possible local accesses to files.**

**Any accessed file transferred to the client.**

**Files stored persistently on local disc cache.**

**Large files are transferred in large chunks (64 kB).**

**Active notification mechanisms to approximate one-copy consistency.**



# AFS Andrew File System

---

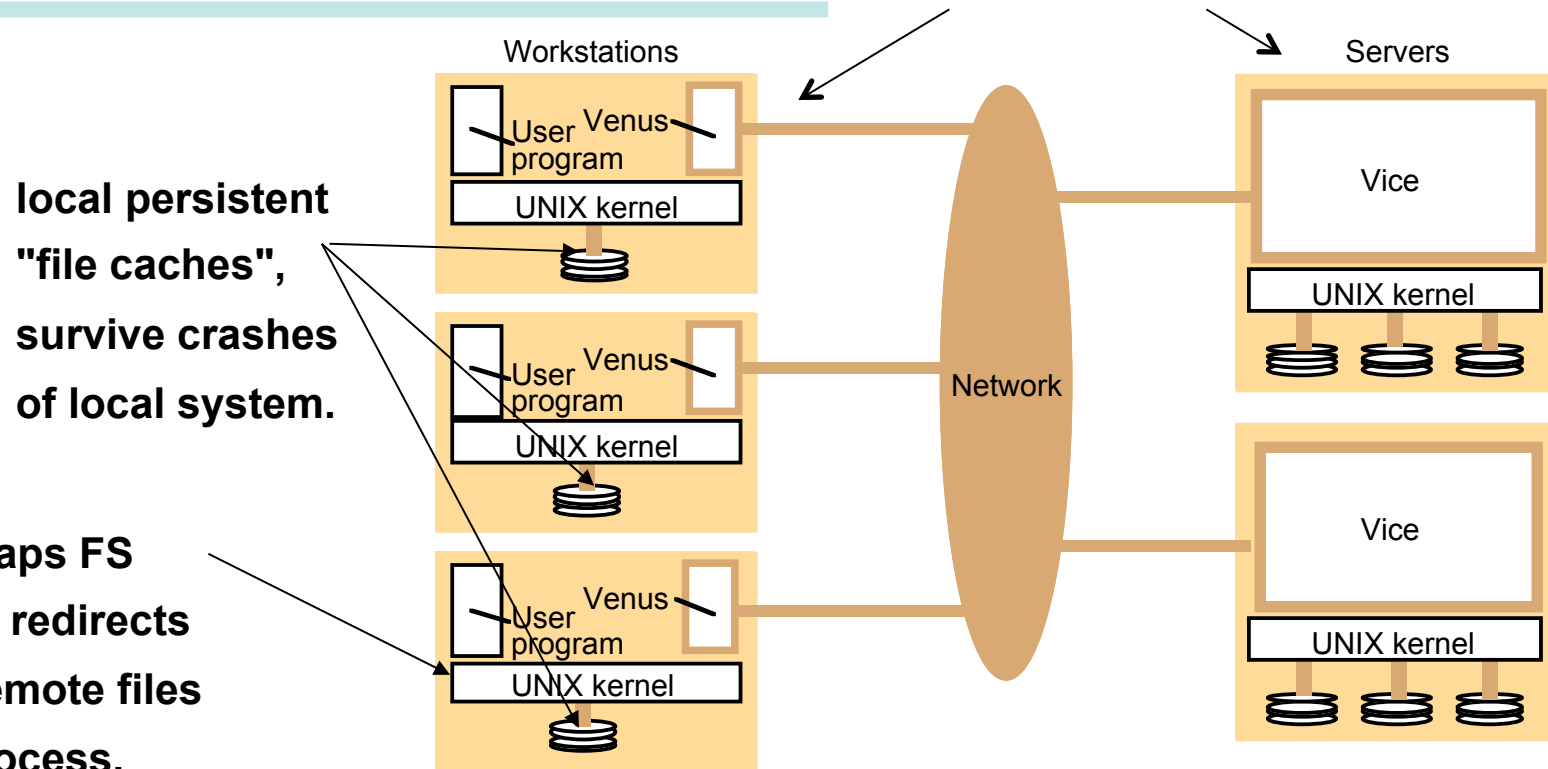
## Questions:

- **How to find the server that holds the copy of the file?**
- **How does AFS obtain control if a client issues an open or close in the shared file space?**
- **Which memory space will be reserved for the cached files?**
- **How to ensure that the cached files constitute the most recent version of the file if more than one client has a copy?**



# AFS Architecture

venus and vice



local persistent "file caches", survive crashes of local system.

Unix kernel traps FS accesses and redirects requests to remote files to a Venus Process.

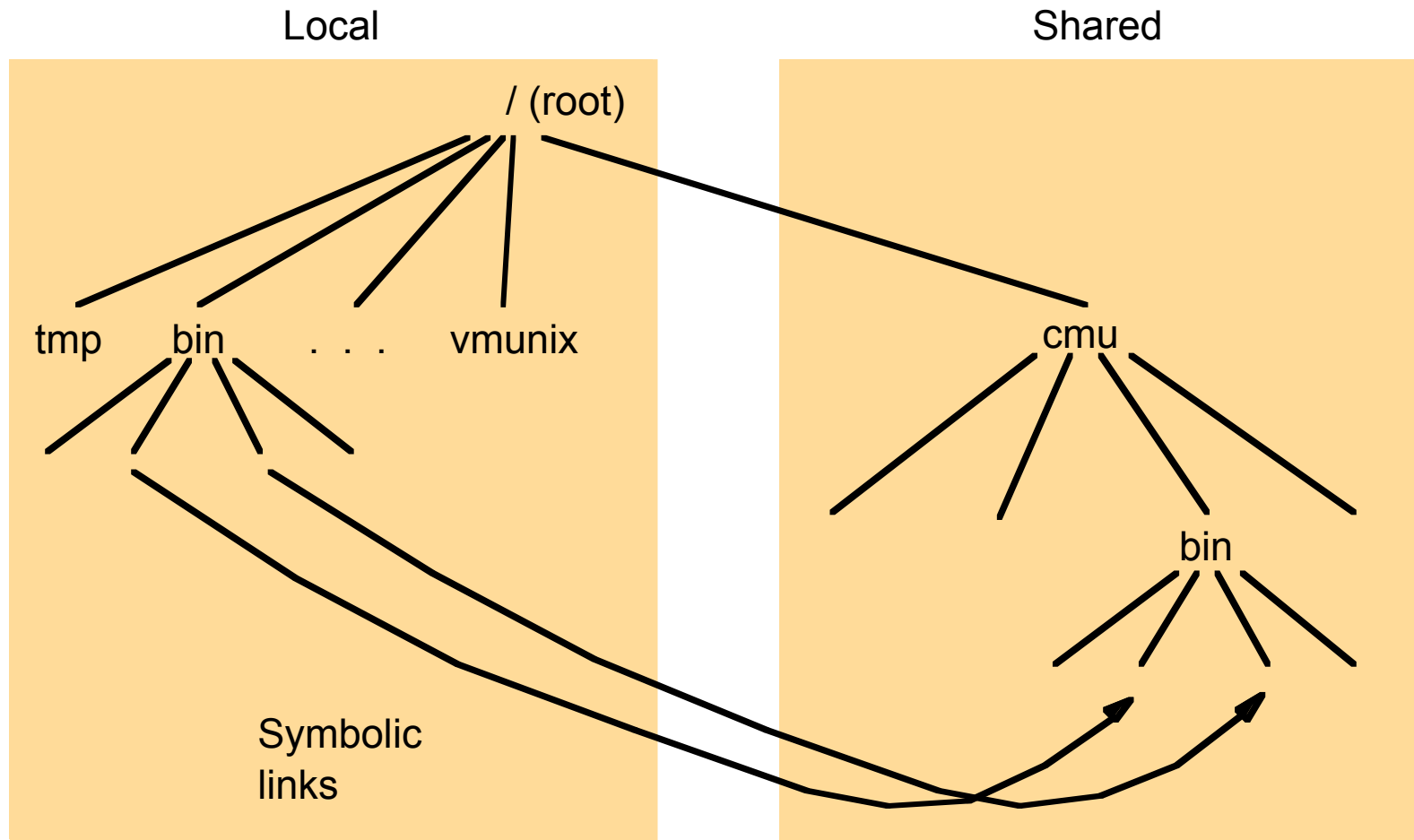
Files are organized in migratable "Volumes" (smaller entities compared to file systems in NFS).

Flat File Service, hierarchical view is established by the Venus Processes

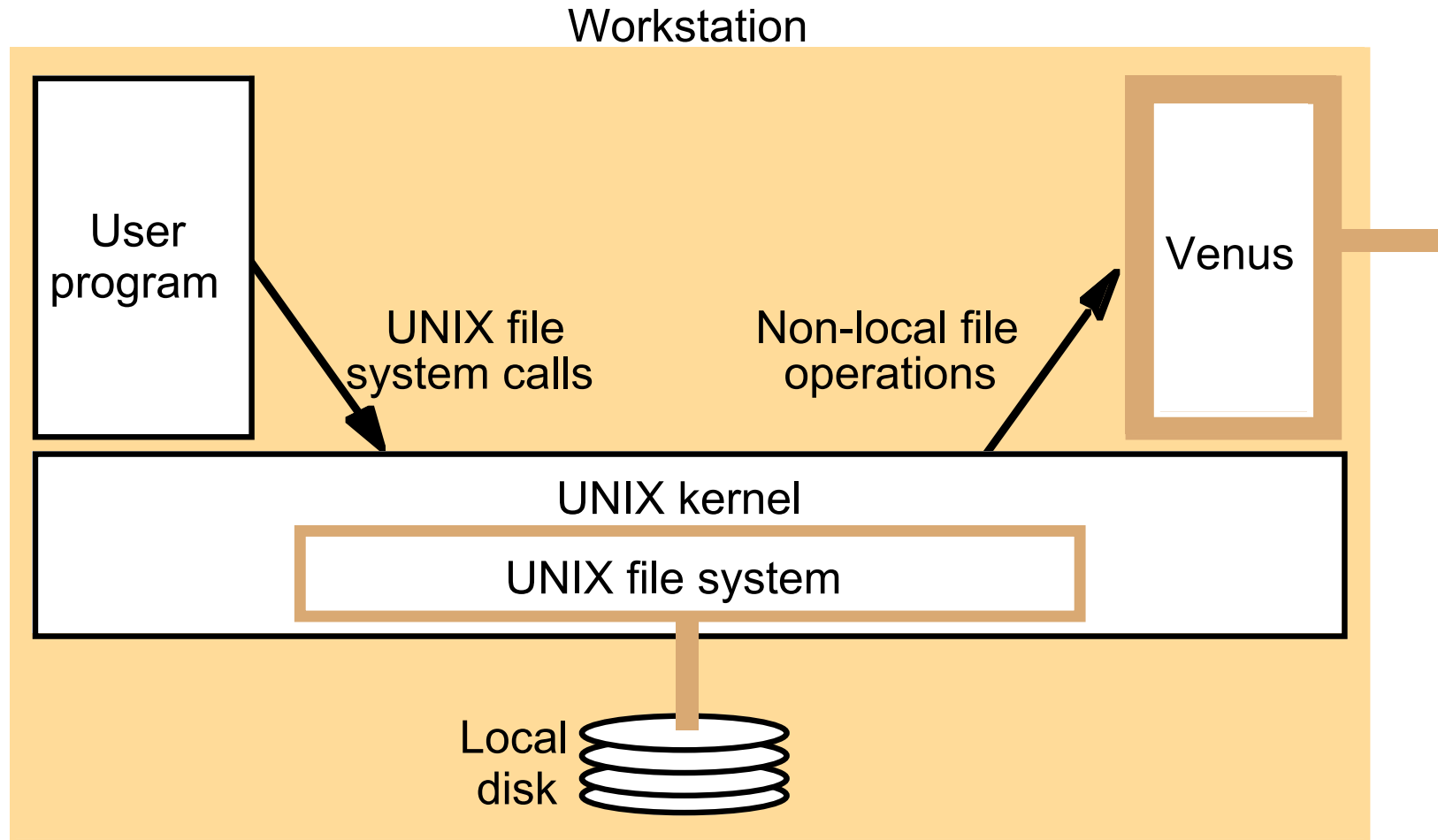
Every File has a unique 96-Bit ID (fid). Path names are translated in FIDs by Venus processes.



# File name space seen by clients of AFS



# System call interception in AFS



# AFS: file system calls

<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>

# The main components of the Vice service interface

---

<i>Fetch(fid) -&gt; attr, data</i>	Returns the attributes (status) and, optionally, the contents of file identified by the <i>fid</i> and records a callback promise on it.
<i>Store(fid, attr, data)</i>	Updates the attributes and (optionally) the contents of a specified file.
<i>Create() -&gt; fid</i>	Creates a new file and records a callback promise on it.
<i>Remove(fid)</i>	Deletes the specified file.
<i>SetLock(fid, mode)</i>	Sets a lock on the specified file or directory. The mode of the lock may be shared or exclusive. Locks that are not removed expire after 30 minutes.
<i>ReleaseLock(fid)</i>	Unlocks the specified file or directory.
<i>RemoveCallback(fid)</i>	Informs server that a Venus process has flushed a file from its cache.
<i>BreakCallback(fid)</i>	This call is made by a Vice server to a Venus process. It cancels the callback promise on the relevant file.

---





# AFS: Basic Consistency Mechanism

---

**AFS consistency relies on a notification concept. The consistency mechanism is based on "Callback Promises" (similar to write invalidate in caches).**

**Callbacks are RPCs issued by the VICE server to the respective remote Venus processes with a Callback Promise Token as parameter. It guarantees that the VENUS process is notified if a client changed a file.**

**A Callback Promise Token may have the values:**

- valid**
- cancelled**

**The Server is responsible to invoke the respective remote Venus process when a file was modified with the value "cancelled".**

**A subsequent local "read" or "open" on the client must request a new file copy.**



# AFS properties

---

## Replicated position database for volumes

- each server holds a full local copy of the position database

## Write-protected replicas

- read-only volumes like /usr/bin or /man are replicated on multiple **servers**

## Transfer of large chunks of data (64kb)

## Caching of parts of the file

- since version 3 only needed parts of files (64k) are transferred, consistency semantics is maintained.

## Performance: Standard Benchmark with 18 server nodes

- Server load AFS: 40%      one reason is notification mechanism in AFS compared to time-out in NFS for consistency maintenance
- Server load NFS: 100%

