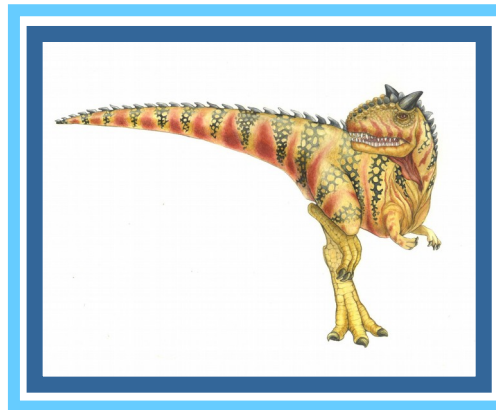# Chapter 15:  Advanced File System

# Chapter 15: Advanced File System

- File-system Structure
- Disk Structure
- File-System Mounting
- Partitions and Mounting
- Virtual File Systems
- Remote File Systems
- Consistency Semantics
- NSF
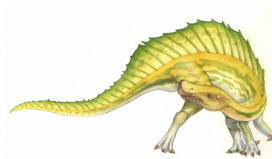- Internal File Structure
- File Sharing

# Objectives

- To describe the details of implementing local file systems and directory structures

- To describe the implementation of remote file systems

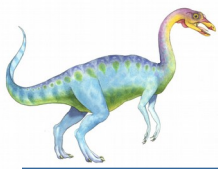- To discuss block allocation and free-block algorithms and trade-offs
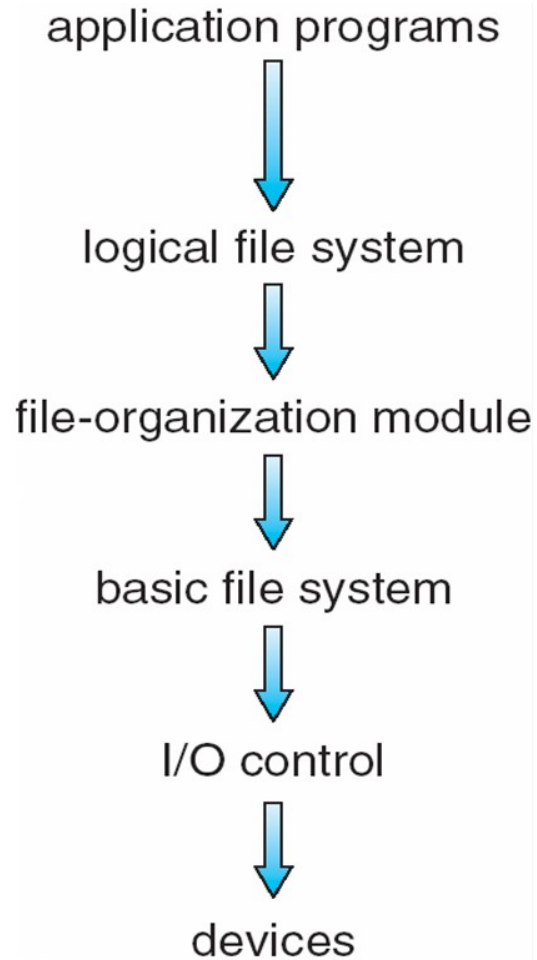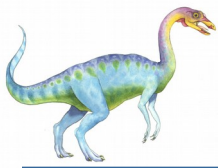
# File-System Structure

- File structure
  - Logical storage unit
  - Collection of related information
- **File system** resides on secondary storage (disks)
  - Provided user interface to storage, mapping logical to physical
  - Provides efficient and convenient access to disk by allowing data to be stored, located retrieved easily
- Disk provides in-place rewrite and random access
  - I/O transfers performed in **blocks** of **sectors** (usually 512 bytes)
- **File control block** – storage structure consisting of information about a file
- **Device driver** controls the physical device
- File system organized into layers

# Layered File System



application programs

↓

logical file system

↓

file-organization module

↓

basic file system

↓

I/O control

↓

devices

# File System Layers

- **Device drivers** manage I/O devices at the I/O control layer
  - Given commands like "read drive1, cylinder 72, track 2, sector 10, into memory location 1060" outputs low-level hardware specific commands to hardware controller

- **Basic file system** given command like "retrieve block 123" translates to device driver
  - Also manages memory buffers and caches (allocation, freeing, replacement)
    - Buffers hold data in transit
    - Caches hold frequently used data

- **File organization module** understands files, logical address, and physical blocks
  - Translates logical block # to physical block #
  - Manages free space, disk allocation

# File System Layers (Cont.)

- **Logical file system** manages metadata information

    - Translates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)

    - Directory management

    - Protection

- Layering useful for reducing complexity and redundancy, but adds overhead and can decrease performance. Translates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)

    - Logical layers can be implemented by any coding method according to OS designer
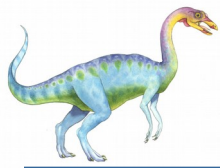
# File System Layers (Cont.)

- Many file systems, sometimes many within an operating system

    - Each with its own format (CD-ROM is ISO 9660; Unix has **UFS**, FFS; Windows has FAT, FAT32, NTFS as well as floppy, CD, DVD Blu-ray, Linux has more than 40 types, with **extended file system** ext2 and ext3 leading; plus distributed file systems, etc.)

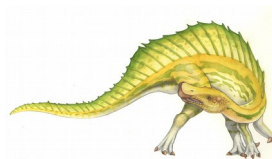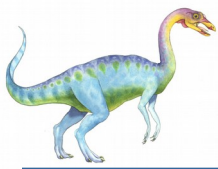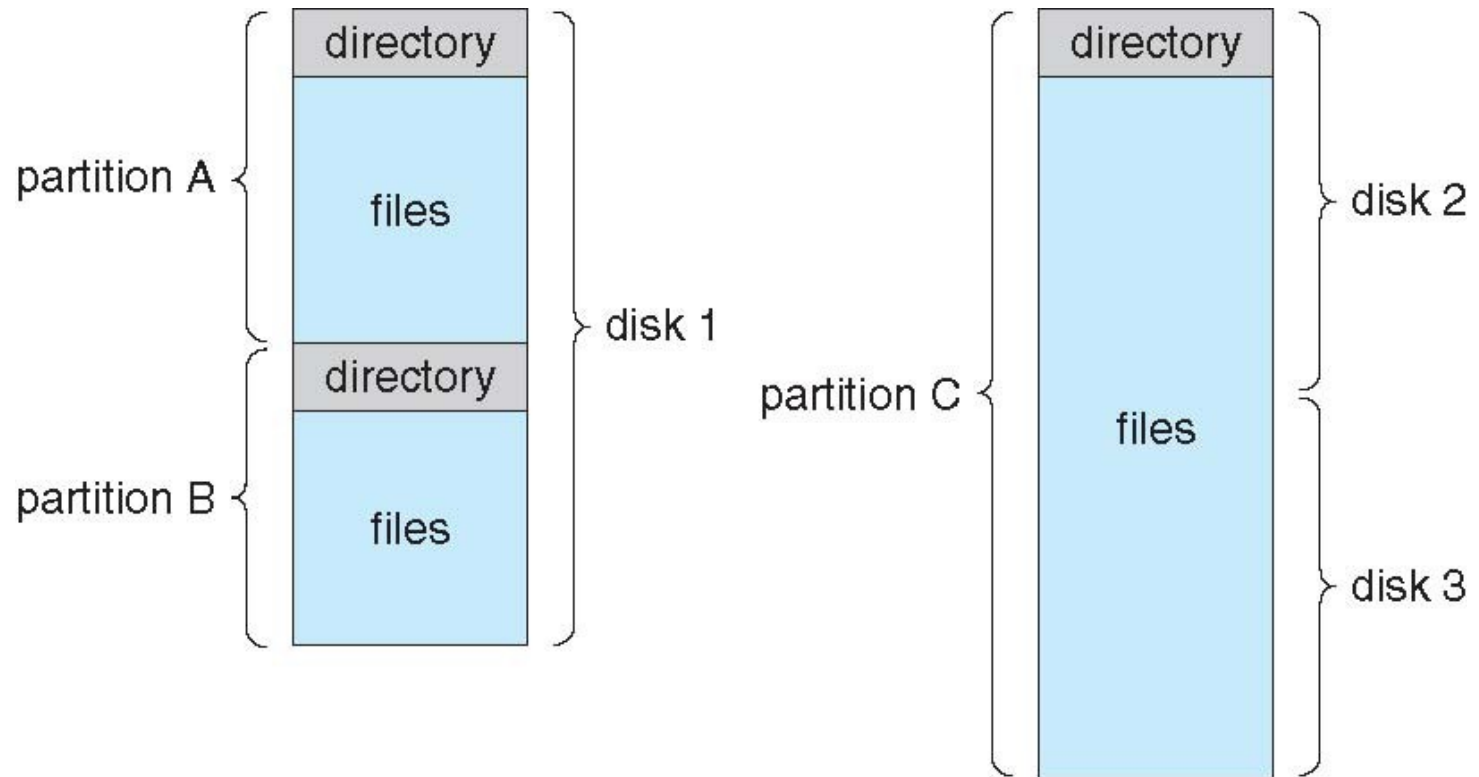    - New ones still arriving – ZFS, GoogleFS, Oracle ASM, FUSE

# Disk Structure

- Disk can be subdivided into **partitions**

- Disk or partition can be used **raw** – without a file system, or **formatted** with a file system

- An entity containing a file system is known as a **volume**

- Each volume containing file system also tracks that file system's info in **device directory** or **volume table of contents**

- A device directory (or simply directory) records information; e.g., name, location, size, etc, about all the files that are kept on the volume.

# A Typical File-system Organization
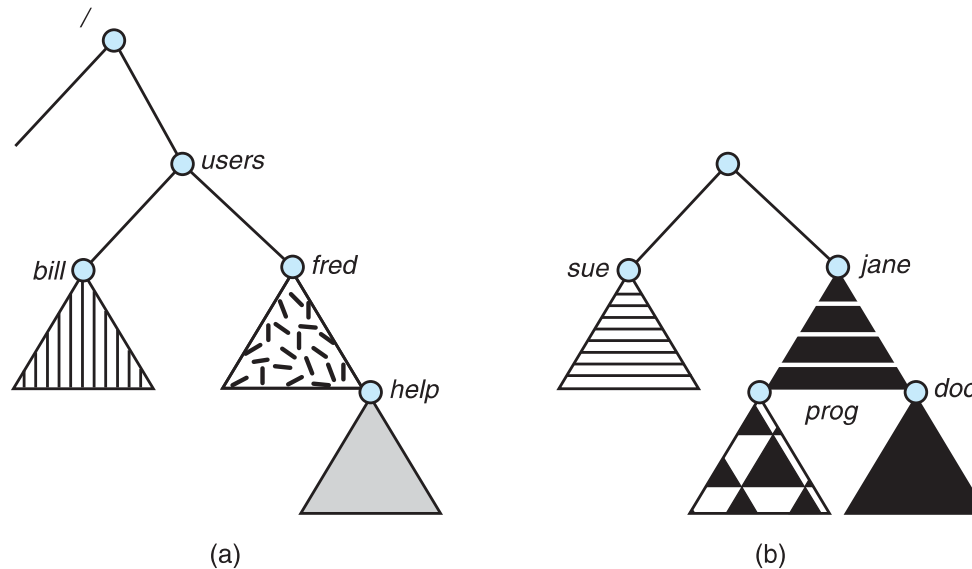
# File System Mounting

- A file system must be mounted before it can be accessed by the various processes executing in the system

- Mount procedure. The operating system is given the name of the device and the mount point -- a location within the file structure where the file system is to be attached.

- Typically, a mount point is an empty directory.

  - For instance, on a UNIX system, a file system containing a user's home directories might be mounted as:

    ‣ /home

  - To access the directory structure within that file system, we precede the directory names with /home, as in
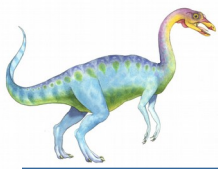
    ‣ /home/jane.

# File System Mounting Example

- Consider the file system depicted below



(a)          (b)

- The triangles represent subtrees of directories.

- Figure (a) shows an existing file system. Figure (b) shows an unmounted volume residing on /device/dsk.

- At this point, only the files on the existing file system can be accessed.

# File System Mounting Example

- The figure blow shows the effect of mounting the volume residing on /device/dsk over /users.



- If the volume is unmounted, the file system is restored to the situation depicted in previous slide.

# File-System Implementation

- We have system calls at the API level, but how do we implement their functions?
  - On-disk and in-memory structures

- **Boot control block** contains info needed by system to boot OS from that volume
  - Needed if volume contains OS, usually first block of volume

- **Volume control block** (**superblock, master file table**) contains volume details
  - Total number of blocks, number of free blocks, block size, free block pointers or array

- Directory structure organizes the files
  - Names and inode numbers, master file table

# File-System Implementation (Cont.)

- **File Control Block** (**FCB**) contains many details about a particular file.
  - inode number, permissions, size, dates
  - NFTS (Window 7 file system) stores into in master file table using relational DB structures

| |
|---|
| file permissions |
| file dates (create, access, write) |
| file owner, group, ACL |
| file size |
| file data blocks or pointers to file data blocks |

# Partitions and Mounting

- Partition can be a volume containing a file system ("cooked") or **raw** – just a sequence of blocks with no file system

- Boot block can point to boot volume or boot loader set of blocks that contain enough code to know how to load the kernel from the file system
  - Or a boot management program for multi-os booting

- **Root partition** contains the OS, other partitions can hold other Oses, other file systems, or be raw
  - Mounted at boot time
  - Other partitions can mount automatically or manually

- At mount time, file system consistency checked
  - Is all metadata correct?
    - If not, fix it, try again
    - If yes, add to mount table, allow access
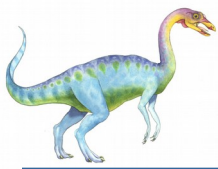
# Virtual File Systems

- **Virtual File Systems** (**VFS**) on Unix provide an object-oriented way of implementing file systems

- VFS allows the same system call interface (the API) to be used for different types of file systems
  - Separates file-system generic operations from implementation details
  - Implementation can be one of many file systems types, or network file system
    - Implements **vnodes** which hold inodes or network file details
  - Then dispatches operation to appropriate file system implementation routines

# Virtual File Systems (Cont.)

- The API is to the VFS interface, rather than any specific type of file system
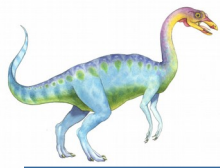
# Virtual File System Implementation

- For example, Linux has four object types:
  - inode, file, superblock, dentry
- VFS defines set of operations on the objects that must be implemented
  - Every object has a pointer to a function table
    - Function table has addresses of routines to implement that function on that object
    - For example:
    - • `int open(. . .)`—Open a file
    - • `int close(. . .)`—Close an already-open file
    - • `ssize t read(. . .)`—Read from a file
    - • `ssize t write(. . .)`—Write to a file
    - • `int mmap(. . .)`—Memory-map a file

# File Sharing

- Sharing of files on multi-user systems is desirable

- Sharing may be done through a **protection** scheme

- On distributed systems, files may be shared across a network

- Network File System (NFS) is a common distributed file-sharing method

- If multi-user system

  - **User IDs** identify users, allowing permissions and protections to be per-user
    **Group IDs** allow users to be in groups, permitting group access rights

  - Owner of a file / directory

  - Group of a file / directory

# File Sharing – Remote File Systems

- Uses networking to allow file system access between systems
  - Manually via programs like FTP
  - Automatically, seamlessly using **distributed file systems**
  - Semi automatically via the **world wide web**

- **Client-server** model allows clients to mount remote file systems from servers
  - Server can serve multiple clients
  - Client and user-on-client identification is insecure or complicated
  - **NFS** is standard UNIX client-server file sharing protocol
  - **CIFS** is standard Windows protocol
  - Standard operating system file calls are translated into remote calls

- Distributed Information Systems (**distributed naming services**) such as LDAP, DNS, NIS, Active Directory implement unified access to information needed for remote computing

# File Sharing – Failure Modes

- All file systems have failure modes
  - For example corruption of directory structures or other non-user data, called **metadata**
- Remote file systems add new failure modes, due to network failure, server failure
- Recovery from failure can involve **state information** about status of each remote request
- **Stateless** protocols such as NFS v3 include all information in each request, allowing easy recovery but less security

# File Sharing – Consistency Semantics

- Specify how multiple users are to access a shared file simultaneously
  - Similar to Ch 5 process synchronization algorithms
    - Tend to be less complex due to disk I/O and network latency (for remote file systems
  - Andrew File System (AFS) implemented complex remote file sharing semantics
  - Unix file system (UFS) implements:
    - Writes to an open file visible immediately to other users of the same open file
    - Sharing file pointer to allow multiple users to read and write concurrently
  - AFS has session semantics
    - Writes only visible to sessions starting after the file is closed

# Network File System (NFS)

- An implementation and a specification of a software system for accessing remote files across LANs (or WANs)

- The implementation is part of the Solaris and SunOS operating systems running on Sun workstations using an unreliable datagram protocol (UDP/IP protocol) and Ethernet

- Interconnected workstations viewed as a set of independent machines with independent file systems, which allows sharing among these file systems in a transparent manner
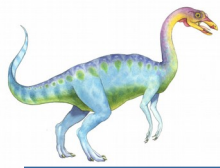
# NFS -- Interconnected workstations

- A remote directory is mounted over a local file system directory

  - The mounted directory looks like an integral sub-tree of the local file system, replacing the sub-tree descending from the local directory

- Specification of the remote directory for the mount operation is nontransparent; the host name of the remote directory has to be provided

  - Files in the remote directory can then be accessed in a transparent manner

- Subject to access-rights accreditation, potentially any file system (or directory within a file system), can be mounted remotely on top of any local directory

# NFS (Cont.)

- NFS is designed to operate in a heterogeneous environment of different machines, operating systems, and network architectures; the NFS specifications independent of these media

- This independence is achieved through the use of RPC primitives built on top of an External Data Representation (XDR) protocol used between two implementation-independent interfaces

- The NFS specification distinguishes between the services provided by a mount mechanism and the actual remote-file-access services
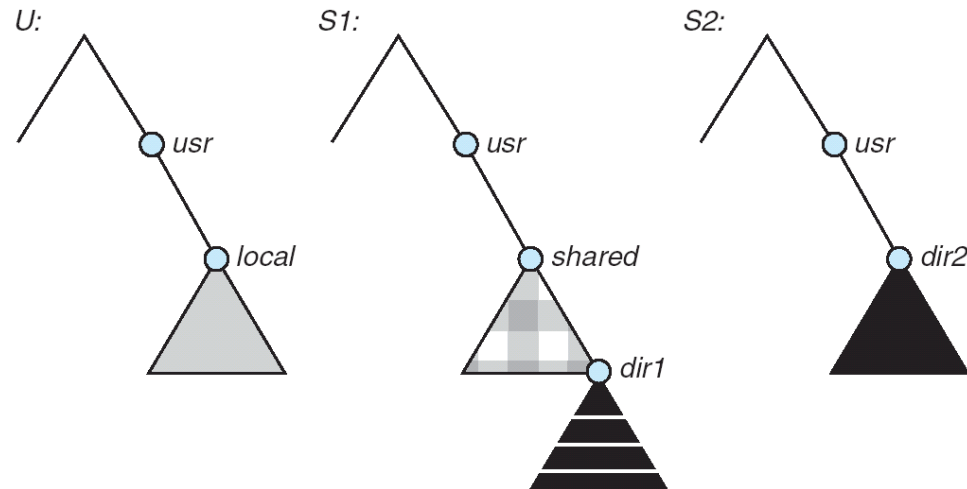
# NFS Mount

- For a remote directory to be accessible in a transparent manner from a particular machine (say M1), a client of M1 must first carry out a mount operation.

- The semantics of the operation involve mounting a remote directory over a directory of a local file system.

- Once the mount operation is completed, the mounted directory looks like an integral sub-tree of the local file system, replacing the sub-tree descending from the local directory.

- The local directory becomes the name of the root of the newly mounted directory.

- Specification of the remote directory as an argument for the mount operation is not done transparently; the location (or host name) of the remote directory has to be provided.

- However, from then on, users on machine M1 can access files in the remote directory in a totally transparent manner.

# NFS Mount Example

- Consider the file system depicted in the figure below, where the triangles represent subtrees of directories that are of interest.
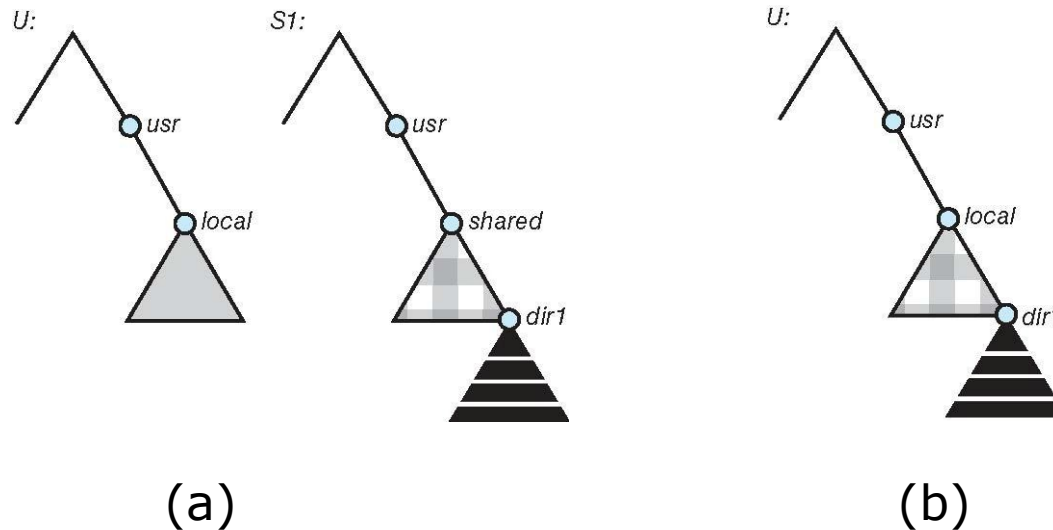


- The figure shows three independent file systems of machines named U, S1, and S2.

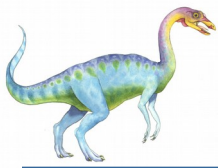- At this point, on each machine, only the local files can be accessed.

# NFS Mount Example

- Figure (b) below shows the effects of mounting S1:/usr/shared over U:/usr/local. This figure (b) depicts the view users U have of their file system



(a)            (b)

- After the mount is complete, the users can access any file within the dir1 directory using the prefix /usr/local/dir1.

- The original directory /usr/local on that machine is no longer visible.

# NFS Mount Protocol

- Establishes initial logical connection between server and client

- Mount operation includes name of remote directory to be mounted and name of server machine storing it

  - Mount request is mapped to corresponding RPC and forwarded to mount server running on server machine

  - Export list – specifies local file systems that server exports for mounting, along with names of machines that are permitted to mount them

- Following a mount request that conforms to its export list, the server returns a file handle—a key for further accesses

- File handle – a file-system identifier, and an inode number to identify the mounted directory within the exported file system

- The mount operation changes only the user's view and does not affect the server side

# NFS Protocol

- Provides a set of remote procedure calls for remote file operations. The procedures support the following operations:
  - searching for a file within a directory
  - reading a set of directory entries
  - manipulating links and directories
  - accessing file attributes
  - reading and writing files
- NFS servers are **stateless**; each request has to provide a full set of arguments (NFS V4 is just coming available – very different, stateful)
- Modified data must be committed to the server's disk before results are returned to the client (lose advantages of caching)
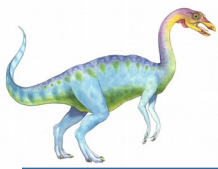- The NFS protocol does not provide concurrency-control mechanisms
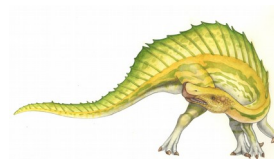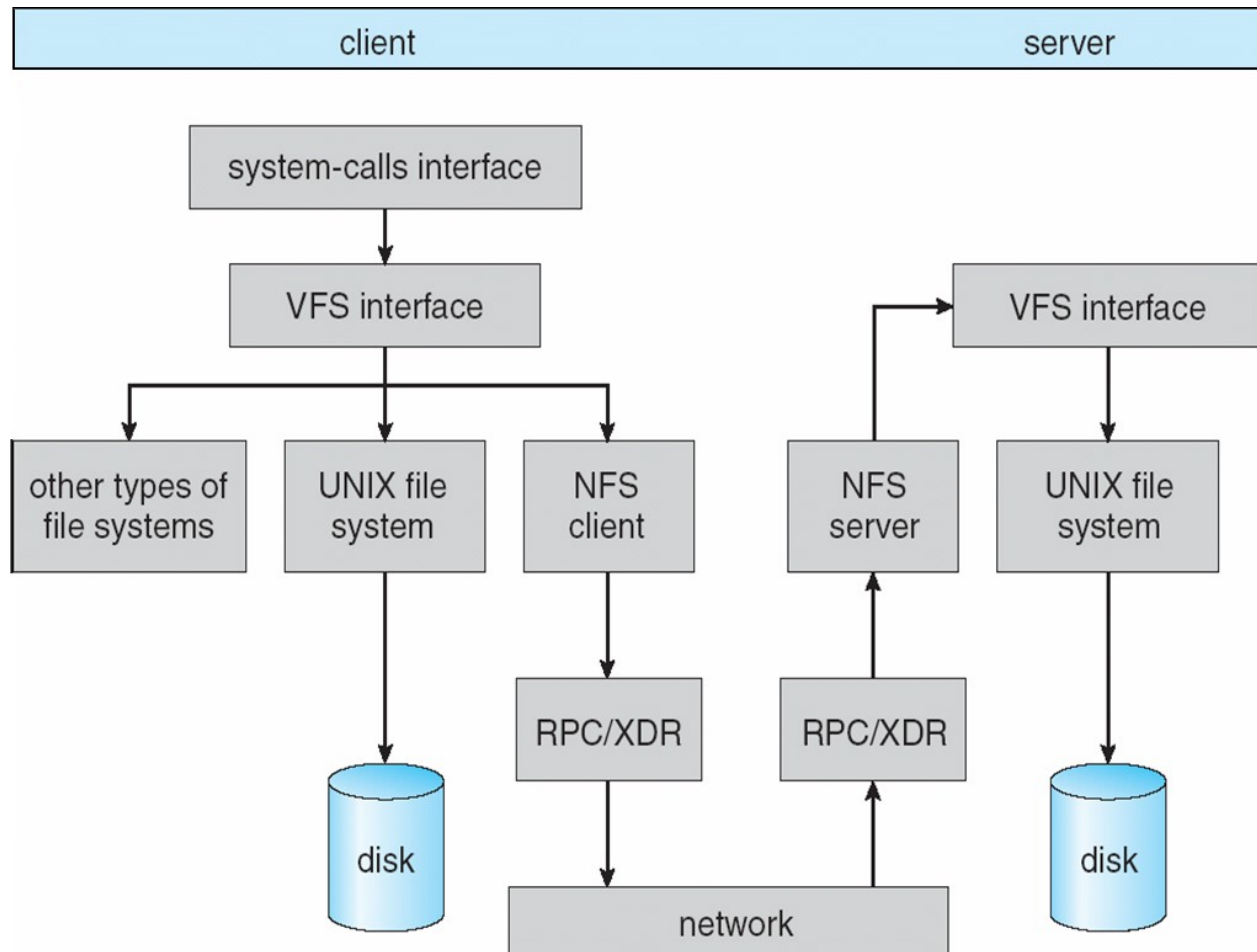
# Three Major Layers of NFS Architecture

- UNIX file-system interface (based on the **open, read, write**, and **close** calls, and **file descriptors**)

- Virtual File System (VFS) layer – distinguishes local files from remote ones, and local files are further distinguished according to their file-system types

  - The VFS activates file-system-specific operations to handle local requests according to their file-system types

  - Calls the NFS protocol procedures for remote requests

- NFS service layer – bottom layer of the architecture

  - Implements the NFS protocol

# Schematic View of NFS Architecture

# NFS Path-Name Translation

- Performed by breaking the path into component names and performing a separate NFS lookup call for every pair of component name and directory vnode

- To make lookup faster, a directory name lookup cache on the client's side holds the vnodes for remote directory names

# NFS Remote Operations

- Nearly one-to-one correspondence between regular UNIX system calls and the NFS protocol RPCs (except opening and closing files)

- NFS adheres to the remote-service paradigm, but employs buffering and caching techniques for the sake of performance

- File-blocks cache – when a file is opened, the kernel checks with the remote server whether to fetch or revalidate the cached attributes

  - Cached file blocks are used only if the corresponding cached attributes are up to date

- File-attribute cache – the attribute cache is updated whenever new attributes arrive from the server

- Clients do not free delayed-write blocks until the server confirms that the data have been written to disk

# End of Chapter 15