
Advanced Operating Systems Lecture notes

Dr. Clifford Neuman

**University of Southern California
Information Sciences Institute**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

CSci555: Advanced Operating Systems
Lecture 10 – October 28 2011
Case Studies: Locus, Athena,
Andrew, HCS, others

**Dr. Clifford Neuman
University of Southern California
Information Sciences Institute**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

The LOCUS System

- ❖ **Developed at UCLA in early 80's**
 - ❑ **Essentially a distributed Unix**
- ❖ **Major contribution was transparency**
 - ❑ **Transparency took many forms**
- ❖ **Environment:**
 - ❑ **VAX 750's and/or IBM PCs connected by an Ethernet**
- ❖ **UNIX compatible.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

LOCUS

- ❖ **Network/location transparency:**
 - ❑ **Network of machines appear as single machine to user.**
 - ❑ **Hide machine boundaries.**
 - ❑ **Local and remote resources look the same to user.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Transparency in Locus

- ❖ **Network Transparency**
 - ❑ **Ability to hide boundaries**
- ❖ **Syntactic Transparency**
 - ❑ **Local and remote calls take same form**
- ❖ **Semantic Transparency**
 - ❑ **Independence from Operand Location**
- ❖ **Name Transparency**
 - ❑ **A name always refers to the same object**
 - ❑ **No need for closure, only one namespace**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Transparency in Locus (cont)

- ❖ **Location Transparency**
 - ❑ **Location can't be inferred from name**
 - ❑ **Makes it easier to move objects**
- ❖ **Syntactic Transparency**
 - ❑ **Local and remote calls take same form**
- ❖ **Performance Transparency**
 - ❑ **Programs with timing assumptions work**
- ❖ **Failure Transparency**
 - ❑ **Remote errors indistinguishable from local**
- ❖ **Execution Transparency**
 - ❑ **Results don't change with location**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

LOCUS Distributed File System

- ❖ **Tree-structured file name space.**
 - ❑ **File name tree covers all file system objects in all machines.**
 - ❑ **Location transparency.**
 - ❑ **File groups (UNIX file systems) “glued” via mount.**
- ❖ **File replication.**
 - ❑ **Varying degrees of replication.**
 - ❑ **Locus responsible for consistency: propagate updates, serve from most up-to-date copy, and handle partitions.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Replication in LOCUS

- ❖ **File group replicated at multiple servers.**
- ❖ **Replicas of a file group may contain different subsets of files belonging to that file group.**
- ❖ **All copies of file assigned same descriptor (i-node #).**
 - ❑ **File unique name: <file group#, i-node #).**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Replica Consistency

❖ Version vectors.

- ❑ Version vector associated with each copy of a file.
- ❑ Maintain update history information.
- ❑ Used to ensure latest copies will be used and to help updating outdated copies.
- ❑ Optimistic consistency.
 - Potential inconsistencies.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

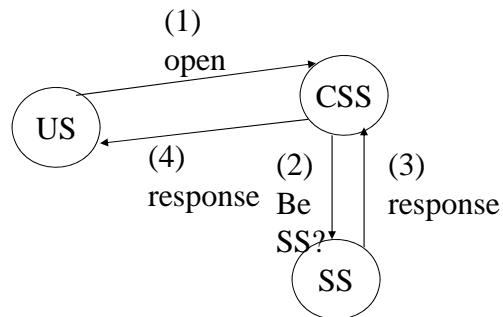
File System Operations 1

- ❖ Using site (US): client.
- ❖ Storage site (SS): server.
- ❖ Current synchronization site (CSS): synchronization site; chooses the SS for a file request.
 - ❑ Knowledge of which files replicated where.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

File System Operations 2

❖ Open:



Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

File Modification

❖ At US:

- ❑ After each change, page sent to SS.
- ❑ At file close, all modified pages flushed to SS.

❖ At SS: atomic commit.

- ❑ Changes to a file handled atomically.
- ❑ No changes are permanent until committed.
- ❑ *Commit* and *abort* system calls.
- ❑ At file close time, changes are committed.
- ❑ Logging and shadow pages.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

CSS

- ❖ **Can implement variety of synchronization policies.**
 - ❑ **Enforce them upon file access.**
 - ❑ **E.g., if sharing policy allows only read-only sharing, CSS disallows concurrent accesses.**

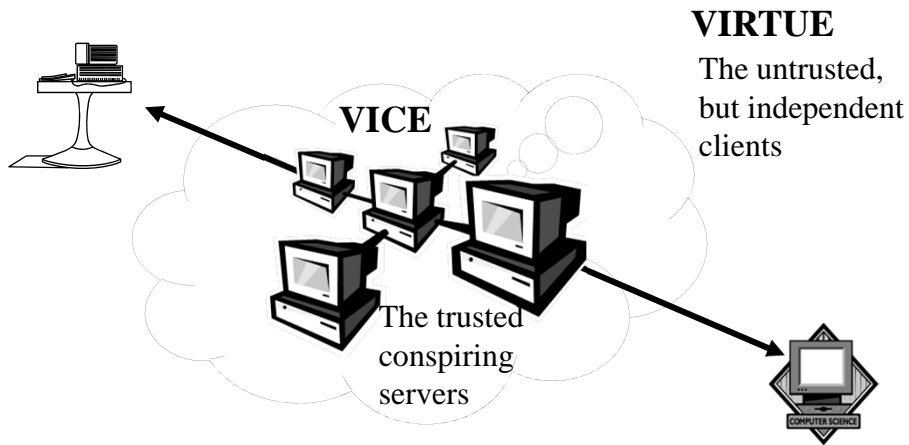
Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Andrew System

- ❖ **Developed at CMU starting in 1982**
 - ❑ **With support from IBM**
 - ❑ **To get computers used as a tool in basic curriculum**
- ❖ **The 3M workstation**
 - ❑ **1 MIP**
 - ❑ **1 MegaPixel**
 - ❑ **1 MegaByte**
 - ❑ **Approx \$10K and 10 Mbps network, local disks**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Vice and Virtue



Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Andrew System (key contributions)

- ❖ **Network Communication**
 - ❑ **Vice (trusted)**
 - ❑ **Virtue (untrusted)**
 - ❑ **High level communication using RPC w/ authentication**
 - ❑ **Security has since switched to Kerberos**
- ❖ **The File System**
 - ❑ **AFS (led to DFS, Coda)**
- ❖ **Applications and user interface**
 - ❑ **Mail and FTP subsumed by file system (w/ gateways)**
- ❖ **Window manager**
 - ❑ **similar to X, but tiled**
 - ❑ **toolkits were priority**
 - ❑ **Since moved to X (and contributed to X)**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Project Athena

- ❖ **Developed at MIT about same time**
 - ❑ **With support from DEC and IBM (and others)**
 - **MIT retained all rights**
 - ❑ **To get computers used as a tool in basic curriculum**
- ❖ **Heterogeneity**
 - ❑ **Equipment from multiple vendors**
- ❖ **Coherence**
 - ❑ **None**
 - **Protocol**
 - **Execution abstraction (e.g. programming environment)**
 - ❑ **Instruction set/binary**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Mainframe/WS vs Unified Model (athena)

- ❖ **Unified model**
 - ❑ **Services provided by system as a whole**
- ❖ **Mainframe / Workstation Model**
 - ❑ **Independent hosts connected by e-mail/FTP**
- ❖ **Athena**
 - ❑ **Unified model**
 - ❑ **Centralized management**
 - ❑ **Pooled resources**
 - ❑ **Servers are not trusted (as much as in Andrew)**
 - ❑ **Clients and network not trusted (like Andrew)**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Project Athena - File system evolution

- ❖ **Remote Virtual Disk (RVD)**
 - ❑ Remotely read and write blocks of disk device
 - ❑ Manage file system locally
 - ❑ Sharing not possible for mutable data
 - ❑ Very efficient for read only data
- ❖ **Remote File System (RFS)**
 - ❑ Remote execution of file system calls
 - ❑ Target host is part of argument (no syntactic transparency).
- ❖ **SUN's Network File System (NFS) - covered**
- ❖ **The Andrew File System (AFS) - covered**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Project Athena - Other Services

- ❖ **Security**
 - ❑ Kerberos
- ❖ **Notification/location**
 - ❑ Zephyr
- ❖ **Mail**
 - ❑ POP
- ❖ **Printing/configuration**
 - ❑ Hesiod-Printcap / Palladium
- ❖ **Naming**
 - ❑ Hesiod
- ❖ **Management**
 - ❑ Moira/RDIST

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Heterogeneous Computer Systems Project

❖ **Developed**

- ❑ **University of Washington, late 1980s**

❖ **Why Heterogeneity**

- ❑ **Organizational diversity**
- ❑ **Need for capabilities from different systems**

❖ **Problems caused by heterogeneity**

- ❑ **Need to support duplicate infrastructure**
- ❑ **Isolation**
- ❑ **Lack of transparency**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

HCS Approach

❖ **Common service to support heterogeneity**

- ❑ **Common API for HCS systems**
- ❑ **Accommodate multiple protocols**

❖ **Transparency**

- ❑ **For new systems accessing existing systems**
- ❑ **Not for existing systems**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

HCS Subsystems

- ❖ **HRPC**
 - ❑ **Common API, modular organization**
 - ❑ **Bind time connection of modules**
- ❖ **HNS (heterogeneous name service)**
 - ❑ **Accesses data in existing name service**
 - ❑ **Maps global name to local lower level names**
- ❖ **THERE**
 - ❑ **Remote execution (by wrapping data)**
- ❖ **HFS (filing)**
 - ❑ **Storage repository**
 - ❑ **Description of data similar to RPC marshalling**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

CORBA (Common Object Request Broker Architecture)

- ❖ **Distributed Object Abstraction**
 - ❑ **Similar level of abstraction as RPC**
- ❖ **Correspondence**
 - ❑ **IDL vs. procedure prototype**
 - ❑ **ORB supports binding**
 - ❑ **IR allows one to discover prototypes**
 - ❑ **Distributed Document Component Facility vs. file system**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

- ❖ **A case study in binding**
 - ❑ **The virtual service is a key abstraction**
- ❖ **Nodes claim ownership of resources**
 - ❑ **Including IP addresses**
- ❖ **On failure**
 - ❑ **Server is restarted, new node claims ownership of the IP resource associated with failed instance.**
 - ❑ **But clients must still retry request and recover.**

CSci555: Advanced Operating Systems

Lecture 11 – November 4 2011

Kernels

Dr. Clifford Neuman
University of Southern California
Information Sciences Institute

Kernels

- ❖ **Executes in supervisory mode.**
 - ❑ **Privilege to access machine's physical resources.**
- ❖ **User-level process: executes in "user" mode.**
 - ❑ **Restricted access to resources.**
 - ❑ **Address space boundary restrictions.**

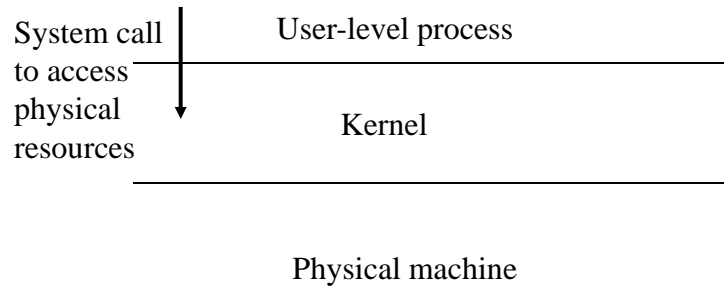
Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Kernel Functions

- ❖ **Memory management.**
 - ❑ **Address space allocation.**
 - ❑ **Memory protection.**
- ❖ **Process management.**
 - ❑ **Process creation, deletion.**
 - ❑ **Scheduling.**
- ❖ **Resource management.**
 - ❑ **Device drivers/handlers.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

System Calls



System call: implemented by hardware interrupt (trap) which puts processor in supervisory mode and kernel address space; executes kernel-supplied handler routine (device driver) executing with interrupts disabled.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Kernel and Distributed Systems

- ❖ **Inter-process communication: RPC, MP, DSM.**
- ❖ **File systems.**
- ❖ **Some parts may run as user-level and some as kernel processes.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Be or not to be in the kernel?

❖ **Monolithic kernels versus microkernels.**

Monolithic kernels

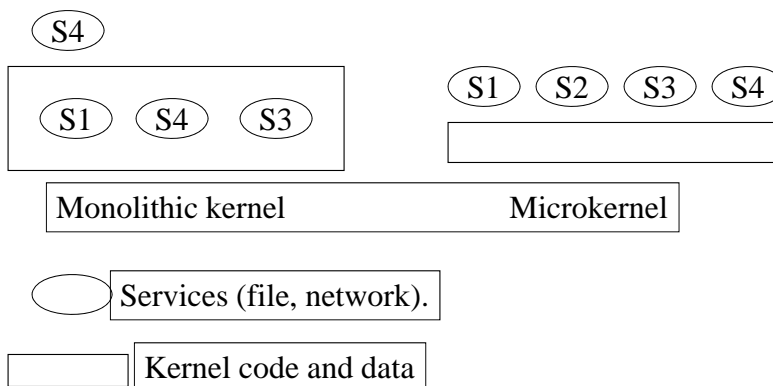
- **Examples: Unix, Sprite.**
- **“Kernel does it all” approach.**
- **Based on argument that inside kernel, processes execute more efficiently and securely.**
- **Problems: massive, non-modular, hard to maintain and extend.**

Microkernels

- ❖ Take as much out of the kernel as possible.
- ❖ Minimalist approach.
- ❖ Modular and small.
 - 10KBytes -> several hundred Kbytes.
 - Easier to port, maintain and extend.
 - No fixed definition of what should be in the kernel.
 - Typically process management, memory management, IPC.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Micro- versus Monolithic Kernels



Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Microkernel

Application

OS Services

Microkernel

Hardware

- Services dynamically loaded at appropriate servers.

- Some microkernels run service processes only @ user space; others allow them to be loaded into either kernel or user space.

The V Distributed System

- ❖ **Stanford (early 80's) by Cheriton et al.**
- ❖ **Distributed OS designed to manage cluster of workstations connected by LAN.**
- ❖ **System structure:**
 - **Relatively small kernel common to all machines.**
 - **Service modules: e.g., file service.**
 - **Run-time libraries: language support (Pascal I/O, C stdio)**
 - **Commands and applications.**

V's Design Goals

- ❖ **High performance communication.**
 - ❑ **Considered the most critical service.**
 - **Efficient file transfer.**
 - ❑ **“Uniform” protocol approach for open system interconnection.**
 - **Interconnect heterogeneous nodes.**
 - ❑ **“Protocols, not software, define the system”.**

The V Kernel

- ❖ **Small kernel with basic protocols and services.**
- ❖ **Precursor to microkernel approach.**
- ❖ **Kernel as a “software backplane”.**
 - ❑ **Provides “slots” into which higher-level OS services can be “plugged”.**

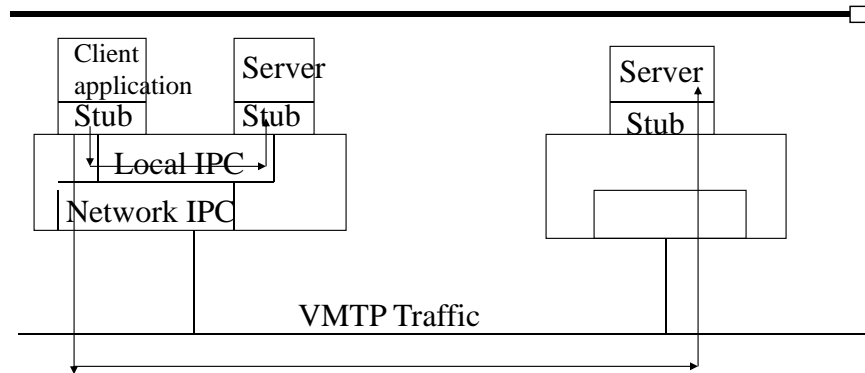
Distributed Kernel

- ❖ **Separate copies of kernel executes on each node.**
- ❖ **They cooperate to provide “single system” abstraction.**
- ❖ **Services: address spaces, LWP, and IPC.**

V's IPC Support

- ❖ **Fast and efficient transport-level service.**
 - **Support for RPC and file transfer.**
- ❖ **V's IPC is RPC-like.**
 - **Send primitive: send + receive.**
 - **Client sends request and blocks waiting for reply.**
 - **Server: processes request serially or concurrently.**
 - **Server response is both ACK and flow control.**
 - **It authorizes new request.**
 - **Simplifies transport protocol.**

V's IPC



Support for short, fixed size messages of 32 bytes with optional data segment of up to 16 Kbytes; simplifies buffering, transmission, and processing.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

VMTP (1)

- ❖ Transport protocol implemented in V.
- ❖ Optimized for request-response interactions.
 - ❑ No connection setup/teardown.
 - ❑ Response ACKs request.
 - ❑ Server maintains state about clients.
 - Duplicate suppression, caching of client information (e.g., authentication information).

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

VMTP (2)

-
- ❖ **Support for group communication.**
 - **Multicast.**
 - **Process groups (e.g., group of file servers).**
 - **Identified by group id.**
 - **Operations: send to group, receive multiple responses to a request.**

VMTP Optimizations

-
- ❖ **Template of VMTP header + some fields initialized in process descriptor.**
 - **Less overhead when sending message.**
 - ❖ **Short, fixed-size messages carried in the VMTP header: efficiency.**

V Kernel: Other Functions

- ❖ **Time, process, memory, and device management.**
- ❖ **Each implemented by separate kernel module (or server) replicated in each node.**
 - **Communicate via IPC.**
 - **Examples: kernel process server creates processes, kernel disk server reads disk blocks.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Time

- ❖ **Kernel keeps current time of day (GMT).**
- ❖ **Processes can get(time), set(time), delay(time), wake up.**
- ❖ **Time synchronization among nodes: outside V kernel using IPC.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Process Management

- ❖ **Create, destroy, schedule, migrate processes.**
- ❖ **Process management optimization.**
 - ❑ **Process initiation separated from address space allocation.**
 - **Process initiation = allocating/initializing new process descriptor.**
 - ❑ **Simplifies process termination (fewer kernel-level resources to reclaim).**
 - ❑ **Simplifies process scheduling: simple priority based scheduler; 2nd. level outside kernel.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Memory Management 1

- ❖ **Protect kernel and other processes from corruption and unauthorized access.**
- ❖ **Address space: ranges of addresses (regions).**
 - ❑ **Bound to an open file (UIO like file descriptor).**
 - ❑ **Page fault references a portion of a region that is not in memory.**
 - ❑ **Kernel performs binding, caching, and consistency services.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Memory Management 2

- ❖ **Virtual memory management: demand paging.**
 - ❑ **Pages are brought in from disk as needed.**
 - ❑ **Update kernel page tables.**
- ❖ **Consistency:**
 - ❑ **Same block may be stored in multiple caches simultaneously.**
 - ❑ **Make sure they are kept consistent.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Device Management

- ❖ **Supports access to devices: disk, network interface, mouse, keyboard, serial line.**
- ❖ **Uniform I/O interface (UIO).**
 - ❑ **Devices are UIO objects (like file descriptors).**
 - ❑ **Example: mouse appears as an open file containing x & y coordinates & button positions.**
 - ❑ **Kernel mouse driver performs polling and interrupt handling.**
 - ❑ **But events associated with mouse changes (moving cursor) performed outside kernel.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

More on V...

- ❖ **Paper talks about other V functions implemented using kernel services.**
 - ❑ **File server.**
 - ❑ **Printer, window, pipe.**
- ❖ **Paper also talks about classes of applications that V targets with examples.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

The X-Kernel

- ❖ **UofArizona, 1990.**
- ❖ **Like V, communication services are critical.**
- ❖ **Machines communicating through internet.**
 - ❑ **Heterogeneity!**
 - ❑ **The more protocols on user's machine, the more resources are accessible.**
- ❖ **The x-kernel philosophy: provide infrastructure to facilitate protocol implementation.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Virtual Protocols

- ❖ **The x-kernel provide library of protocols.**
 - ❑ **Combined differently to access different resources.**
 - ❑ **Example:**
 - **If communication between processes on the same machine, no need for any networking code.**
 - **If on the same LAN, IP layer skipped.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

The X-Kernel : Process and Memory

- ❖ **ability to pass control and data efficiently between the kernel and user programs**
 - ❑ **user data is accessible because kernel process executes in same address space**
- ❖ **kernel process -> user process**
 - ❑ **sets up user stack**
 - ❑ **pushes arguments**
 - ❑ **use user-stack**
 - ❑ **access only user data**
- ❖ **kernel -> user (245 usec), user -> kernel 20 usec on SUN 3/75**

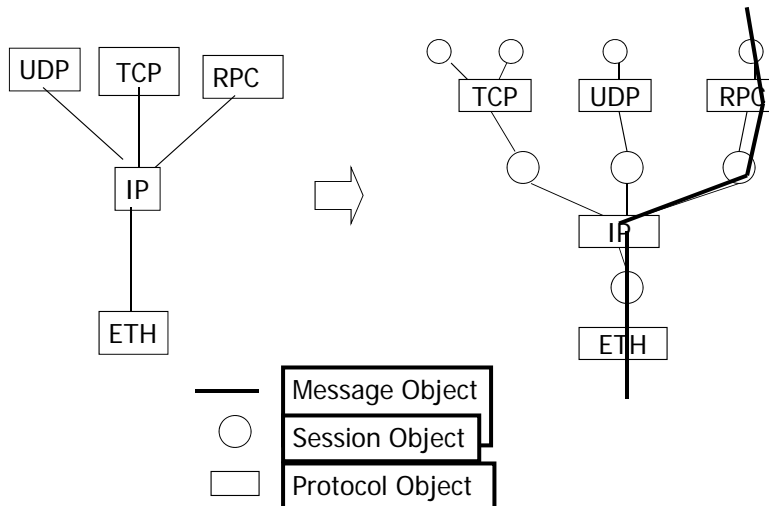
Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Communication Manager

- ❖ Object-oriented infrastructure for implementing and composing protocols.
- ❖ Common protocol interface.
- ❖ 2 abstract communication objects:
 - Protocols and sessions.
 - Example: TCP protocol object.
 - TCP open operation: creates a TCP session.
 - TCP protocol object: switches each incoming message to one of the TCP session objects.
 - Operations: demux, push, pop.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

X-kernel Configuration



Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Message Manager

- ❖ **Defines single abstract data type: message.**
 - ❑ **Manipulation of headers, data, and trailers that compose network transmission units.**
 - ❑ **Well-defined set of operations:**
 - **Add headers and trailers, strip headers and trailers, fragment/reassemble.**
 - ❑ **Efficient implementation using directed acyclic graphs of buffers to represent messages + stack data structure to avoid data copying.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Mach

- ❖ **CMU (mid 80's).**
- ❖ **Mach is a microkernel, not a complete OS.**
- ❖ **Design goals:**
 - ❑ **As little as possible in the kernel.**
 - ❑ **Portability: most kernl code is machine independent.**
 - ❑ **Extensibility: new features can be implemented/tested alongside existing versions.**
 - ❑ **Security: minimal kernel specified and implemented in more secure way.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Mach Features

- ❖ **OSs as Mach applications.**
- ❖ **Mach functionality:**
 - ❑ **Task and thread management.**
 - ❑ **IPC.**
 - ❑ **Memory management.**
 - ❑ **Device management.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Mach IPC

- ❖ **Threads communicate using ports.**
- ❖ **Resources are identified with ports.**
- ❖ **To access resource, message is sent to corresponding port.**
 - ❑ **Ports not directly accessible to programmer.**
 - ❑ **Need handles to “port rights”, or capabilities (right to send/receive message to/from ports).**
- ❖ **Servers: manage several resources, or ports.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Mach: ports

- ❖ *process port* is used to communicate with the kernel.
- ❖ *bootstrap port* is used for initialization when a process starts up.
- ❖ *exception port* is used to report exceptions caused by the process.
- ❖ *registered ports* used to provide a way for the process to communicate with standard system servers.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Protection

- ❖ **Protecting resources against illegal access:**
 - ❑ **Protecting port against illegal sends.**
- ❖ **Protection through *capabilities*.**
 - ❑ **Kernel controls port capability acquisition.**
 - ❑ **Different from Amoeba.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Capabilities 1

- ❖ **Capability to a port has field specifying port access rights for the task that holds the capability.**
 - ❑ **Send rights: threads belonging to task possessing capability can send message to port.**
 - ❑ **Send-once rights: allows at most 1 message to be sent; after that, right is revoked by kernel.**
 - ❑ **Receive rights: allows task to receive message from port's queue.**
 - **At most 1 task, may have receive rights at any time.**
 - **More than 1 task may have send/send-once rights.**

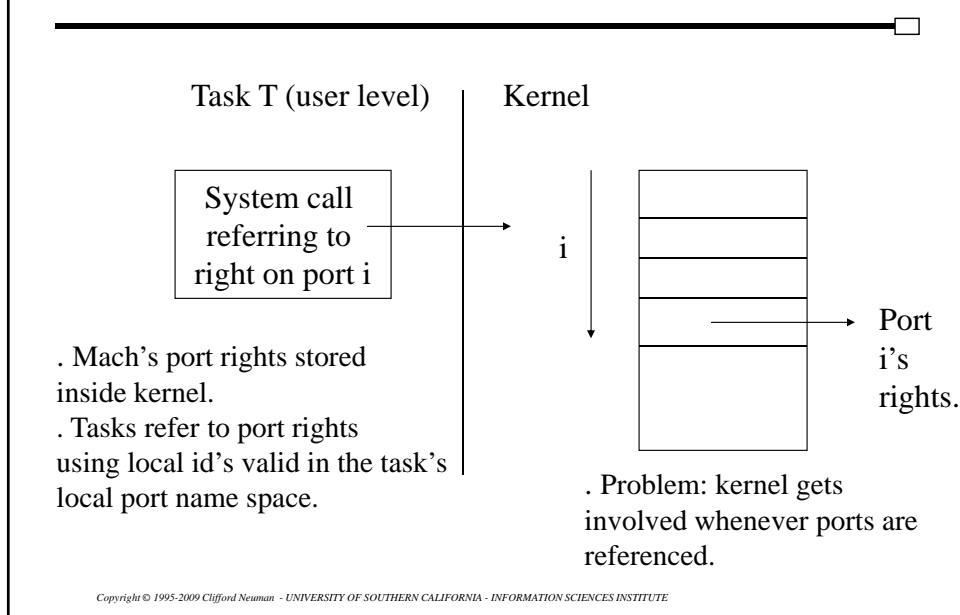
Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Capabilities 2

- ❖ **At task creation:**
 - ❑ **Task given bootstrap port right: send right to obtain services of other tasks.**
 - ❑ **Task threads acquire further port rights either by creating ports or receiving port rights.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Port Name Space



Communication Model

❖ Message passing.

❖ Messages: fixed-size headers + variable-length list of data items.

Header	T	Port rights	T	In-line data	T	Pointer to out-of-line data
--------	---	-------------	---	--------------	---	-----------------------------

Header: destination port, reply port, type of operation.

T: type of information.

Port rights: send rights: receiver acquires send rights to port.

Receive rights: automatically revoked in sending task.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Ports

- ❖ **Mach port has message queue.**
 - ❑ **Task with receive rights can set port's queue size dynamically: flow control.**
 - ❑ **If port's queue is full, sending thread is blocked; send-once sender never blocks.**
- ❖ **System calls:**
 - ❑ **Send message to kernel port.**
 - ❑ **Assigned at task creation time.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Task and Thread Management

- ❖ **Task: execution environment (address space).**
- ❖ **Threads within task perform action.**
- ❖ **Task resources: address space, threads, port rights.**
- ❖ **PAPER:**
 - ❑ **How Mach microkernel can be used to implement other OSs.**
 - ❑ **Performance numbers comparing 4.3 BSD on top of Mach and Unix kernels.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

CSci555: Advanced Operating Systems
Lecture 12 - November 11 2011
Scheduling, Fault Tolerance
Real Time, Database Support

Dr. Clifford Neuman
University of Southern California
Information Sciences Institute

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Scheduling and Real-Time systems

❖ **Scheduling**

- ❑ **Allocation of resources at a particular point in time to jobs needing those resources, usually according to a defined policy.**

❖ **Focus**

- ❑ **We will focus primarily on the scheduling of processing resources, though similar concepts apply the the scheduling of other resources including network bandwidth, memory, and special devices.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

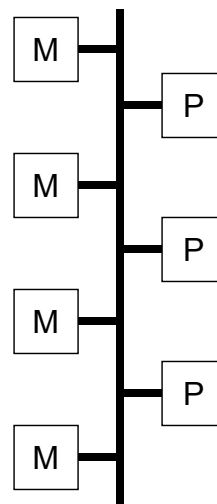
Parallel Computing - General Issues

- ❖ **Speedup - the final measure of success**
 - ❑ **Parallelism vs Concurrency**
 - **Actual vs possible by application**
 - ❑ **Granularity**
 - **Size of the concurrent tasks**
 - **Reconfigurability**
 - ❑ **Number of processors**
 - ❑ **Communication cost**
 - ❑ **Preemption v. non-preemption**
 - ❑ **Co-scheduling**
 - **Some things better scheduled together**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Shared Memory Multi-Processing

- ❖ **Distributed shared memory, and shared memory multi-processors**
- ❖ **Processors usually tightly coupled to memory, often on a shared bus. Programs communicated through shared memory locations.**
- ❖ **For SMPs cache consistency is the important issue. In DSM it is memory coherence.**
 - ❑ **One level higher in the storage hierarchy**
- ❖ **Examples**
 - **Sequent, Encore Multimax, DEC Firefly, Stanford DASH**



Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Where is the best place for scheduling

- ❖ **Application is in best position to know its own specific scheduling requirements**
 - ❑ Which threads run best simultaneously
 - ❑ Which are on Critical path
 - ❑ But Kernel must make sure all play fairly
- ❖ **MACH Scheduling**
 - ❑ Lets process provide hints to discourage running
 - ❑ Possible to hand off processor to another thread
 - Makes easier for Kernel to select next thread
 - Allow interleaving of concurrent threads
 - ❑ Leaves low level scheduling in Kernel
 - ❑ Based on higher level info from application space

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

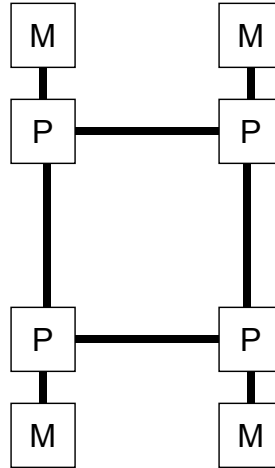
Scheduler activations

- ❖ **User level scheduling of threads**
 - ❑ Application maintains scheduling queue
- ❖ **Kernel allocates threads to tasks**
 - ❑ Makes upcall to scheduling code in application when thread is blocked for I/O or preempted
 - ❑ Only user level involved if blocked for critical section
- ❖ **User level will block on kernel calls**
 - ❑ Kernel returns control to application scheduler

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Distributed-Memory Multi-Processing

- ❖ **Processors coupled to only part of the memory**
 - ❑ Direct access only to their own memory
- ❖ **Processors interconnected in mesh or network**
 - ❑ Multiple hops may be necessary
- ❖ **May support multiple threads per task**
- ❖ **Typical characteristics**
 - ❑ Higher communication costs
 - ❑ Large number of processors
 - ❑ Coarser granularity of tasks
- ❖ **Message passing for communication**



Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Condor

- ❖ **Identifies idle workstations and schedules background jobs on them**
- ❖ **Guarantees job will eventually complete**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Condor

- ❖ **Analysis of workstation usage patterns**
 - ❑ **Only 30%**
- ❖ **Remote capacity allocation algorithms**
 - ❑ **Up-Down algorithm**
 - **Allow fair access to remote capacity**
- ❖ **Remote execution facilities**
 - ❑ **Remote Unix (RU)**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Condor

- ❖ **Leverage: performance measure**
 - ❑ **Ratio of the capacity consumed by a job remotely to the capacity consumed on the home station to support remote execution**
- ❖ **Checkpointing: save the state of a job so that its execution can be resumed**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Condor - Issues

- ❖ **Transparent placement of background jobs**
- ❖ **Automatically restart if a background job fails**
- ❖ **Users expect to receive fair access**
- ❖ **Small overhead**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Condor - scheduling

- ❖ **Hybrid of centralized static and distributed approach**
- ❖ **Each workstation keeps own state information and schedule**
- ❖ **Central coordinator assigns capacity to workstations**
 - ❑ **Workstations use capacity to schedule**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

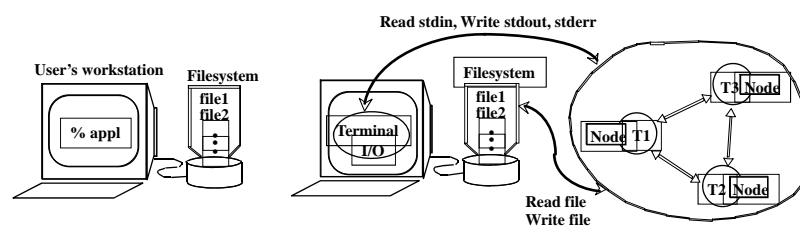
Prospero Resource Manager

Prospero Resource Manager - 3 entities

- ❖ One or more system managers
 - *Each manages subset of resources*
 - *Allocates resources to jobs as needed*
- ❖ A job manager associated with each job
 - *Identifies resource requirements of the job*
 - *Acquires resources from one or more system managers*
 - *Allocates resources to the job's tasks*
- ❖ A Node manager on each node
 - *Mediates access to the nodes resources*

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

The Prospero Resource Manager



A) User invokes an application program on his workstation.

b) The program begins executing on a set of nodes. Tasks perform terminal and file I/O on the user's workstation.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Advantages of the PRM

❖ Scalability

- ❑ *System manager does not require detailed job information*
- ❑ *Multiple system managers*

❖ Job manager selected for application

- ❑ *Knows more about job's needs than the system manager*
- ❑ *Alternate job managers useful for debugging, performance tuning*

❖ Abstraction

- ❑ *Job manager provides a single resource allocator for the job's tasks*
- ❑ *Single system model*

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Real time Systems

❖ Issues are scheduling and interrupts

- ❑ **Must complete task by a particular deadline**
- ❑ **Examples:**
 - **Accepting input from real time sensors**
 - **Process control applications**
 - **Responding to environmental events**

❖ How does one support real time systems

- ❑ **If short deadline, often use a dedicated system**
- ❑ **Give real time tasks absolute priority**
- ❑ **Do not support virtual memory**
 - **Use early binding**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Real time Scheduling

- ❖ **To initiate, must specify**
 - ❑ **Deadline**
 - ❑ **Estimate/upper-bound on resources**
- ❖ **System accepts or rejects**
 - ❑ **If accepted, agrees that it can meet the deadline**
 - ❑ **Places job in calendar, blocking out the resources it will need and planning when the resources will be allocated**
- ❖ **Some systems support priorities**
 - ❑ **But this can violate the RT assumption for already accepted jobs**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

CSci555: Advanced Operating Systems

Lecture 12 - November 11, 2011

Fault Tolerant Computing

Dr. Clifford Neuman
University of Southern California
Information Sciences Institute

NOTE: This is a very short lecture, with much of the discussion integrated with the material on scheduling from the previous lecture.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Fault-Tolerant systems

❖ Failure probabilities

- ❑ Hierarchical, based on lower level probabilities
- ❑ Failure Trees
- ❑ Add probabilities where any failure affects you
 - Really $(1 - ((1 - \lambda)(1 - \lambda)(1 - \lambda)))$
- ❑ Multiply probabilities if all must break
 - Since numbers are small, this reduces failure rate
- ❑ Both failure and repair rate are important

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Making systems fault tolerant

❖ Involves masking failure at higher layers

- ❑ Redundancy
- ❑ Error correcting codes
- ❑ Error detection

❖ Techniques

- ❑ In hardware
- ❑ Groups of servers or processors execute in parallel and provide hot backups

❖ Space Shuttle Computer Systems examples

❖ RAID example

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Types of failures

- ❖ **Fail stop**
 - ❑ **Signals exception, or detectably does not work**
- ❖ **Returns wrong results**
 - ❑ **Must decide which component failed**
- ❖ **Byzantine**
 - ❑ **Reports difficult results to different participants**
 - ❑ **Intentional attacks may take this form**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Recovery

- ❖ **Repair of modules must be considered**
 - ❑ **Repair time estimates**
- ❖ **Reconfiguration**
 - ❑ **Allows one to run with diminished capacity**
 - ❑ **Improves fault tolerance (from catastrophic failure)**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

OS Support for Databases

- ❖ **Example of OS used for particular applications**
- ❖ **End-to-end argument for applications**
 - ❑ **Much of the common services in OS's are optimized for general applications.**
 - ❑ **For DBMS applications, the DBMS might be in a better position to provide the services**
 - **Caching, Consistency, failure protection**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

CSci555: Advanced Operating Systems

Lecture 13 - November 18, 2011

Grid and Cloud Computing

Dr. Clifford Neuman
University of Southern California
Information Sciences Institute

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Announcements

- ❖ **No lecture November 25, Thanksgiving recess.**
- ❖ **Next lecture 2 December – final lecture**
 - ❑ **Need volunteer to do distribute class evaluations**
 - ❑ **Send email or suggest topics on discussion forum**
 - ❑ **Will include review for final exam**
 - ❑ **Research paper due 2 December**
- ❖ **Final exam on Friday December 9 – 2:00PM - 4:00PM**
 - ❑ **Location TBD**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Grids

- ❖ **Computational grids apply many distributed system techniques to meta computing (parallel applications running on large numbers of nodes across significant distances).**
 - ❑ **Libraries provide a common base for managing such systems.**
 - ❑ **Some consider grids different, but in my view the differences are not major, just the applications are.**
- ❖ **Data grids extend the grid “term” into other classes of computing.**
 - ❑ **Issues for data grids are massive storage, indexing, and retrieval.**
 - ❑ **It is a file system, indexing, and ontological problem.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

The Cloud

- ❖ **The cloud is many things to many people**
 - ❑ **Software as a service and hosted applications**
 - ❑ **Processing as a utility**
 - ❑ **Storage as a utility**
 - ❑ **Remotely hosted servers**
 - ❑ **Anything beyond the network card**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

The Cloud

- ❖ **Clouds are hosted in different ways**
 - ❑ **Private Clouds**
 - ❑ **Public Clouds**
 - ❑ **Hosted Private Clouds**
 - ❑ **Hybrid Clouds**
 - ❑ **Clouds for federated enterprises**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

The Cloud

❖ Clouds are hosted in different ways

- ❑ Private Clouds
- ❑ Public Clouds
- ❑ Hosted Private Clouds
- ❑ Hybrid Clouds
- ❑ Clouds for federated enterprises

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

The Paper

❖ Cloud Computing and Grid Computing 360 Degree compared.

❖ Written by one of the principal “architectures” of grid computing and provides one perspective.

- ❑ Basically the paper is trying to frame cloud computing in terms of grid computing so that cloud computing does not steal the credit for many of the technological advances that was claimed by grid-computing.
- ❑ In reality, many of the advances are from distributed systems research that predated the grid, and the grid did much of the same to distributes system research as cloud computing is doing to the grid.

❖ In both cases the innovation is/will be engineering and standardization in the context of particular classes of applications.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Issues in the Grid and Cloud

- ❖ **Common interfaces and middleware**
 - ❑ **Directory services**
 - ❑ **Security services**
 - ❑ **File services**
 - ❑ **Scheduling services / allocation**
- ❖ **Support for federated environments**
 - ❑ **Security in such environments**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Directory Services

- ❖ **Need for a catalog of cloud or grid resources.**
- ❖ **Directory services also map locations for services once allocated to a computation.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Security Services

❖ Virtualization

- ❑ Separation of “platform”

❖ VPN's

- ❑ Brings remote resources “inside”

❖ Federated Identity

- ❑ Or separate identity for cloud

❖ Policy services

- ❑ Much work is needed

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

File Services

❖ Performance often dictates storage near the computation.

- ❑ But the data must be migrated
- ❑ Alternatively, data accessed through callbacks to originating system.
- ❑ Or in a separate storage cloud.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Scheduling/Migration and Allocation

- ❖ **Characterize Node Capabilities in the Cloud**
 - ❑ **Security Characteristics**
 - Accreditation of the software for managing nodes and data
 - ❑ **Legal and Geographic Characteristics**
 - Includes data on managing organizations and contractors
 - ❑ **Need language to characterize**
 - ❑ **Need endorers to certify**
- ❖ **Define Migration Policies**
 - ❑ **Who is authorized to handle data**
 - ❑ **Any geographic constraints**
 - ❑ **Necessary accreditation for servers and software**
 - Each node that accepts data must be capable for enforcing policy before data can be redistributed.
 - ❑ **Languages needed to describe**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Federation

- ❖ **Resources provided by parties with different interests.**
 - ❑ **No single chain of authority**
 - ❑ **Resources acquired from multiple parties and must be interconnected.**
- ❖ **Policy issues dominate**
 - ❑ **Who can use resources**
 - ❑ **Which resources is one willing to use.**
 - ❑ **Translating ID's and policies at boundaries**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

CSci555: Advanced Operating Systems
Lecture 13 – November 18, 2011
Selected Topics and Scalable Systems

Dr. Clifford Neuman
University of Southern California
Information Sciences Institute

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Announcements

- ❖ **No lecture November 25, Thanksgiving recess.**
- ❖ **Next lecture 2 December – final lecture**
 - ❑ **Send email or suggest topics on discussion forum**
 - ❑ **Will include review for final exam**
 - ❑ **Research paper due 2 December**
- ❖ **Final exam on Friday December 9 – 2:00PM - 4:00PM**
 - ❑ **Location TBD**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Hints for building scalable systems

❖ From Lampson:

- ❑ **Keep it simple**
- ❑ **Do one thing at a time**
- ❑ **If in doubt, leave it out**
- ❑ **But no simpler than possible**
- ❑ **Generality can lead to poor performance**
- ❑ **Make it fast and simple**
- ❑ **Don't hide power**
- ❑ **Leave it to the client**
- ❑ **Keep basic interfaces stable**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Hints for building scalable systems

❖ From Lampson:

- ❑ **Plan to throw one away**
- ❑ **Keep secrets**
- ❑ **Divide and conquer**
- ❑ **Use a good idea again**
- ❑ **Handle normal and worst case separately**
- ❑ **Optimize for the common case**
- ❑ **Split resources in a fixed way**
- ❑ **Cache results of expensive operations**
- ❑ **Use hints**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Hints for building scalable systems

❖ From Lampson:

- ❑ When in doubt use brute force
- ❑ Compute in the background
- ❑ Use batch processing
- ❑ Safety first
- ❑ Shed load
- ❑ End-to-end argument
- ❑ Log updates

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Mid-Term Review

Replication, Distribution, and Caching

❖ Briefly describe the use of replication, distribution and caching in each of the following systems. State whether each of the techniques (replication, distribution, and caching) is used, and if used, briefly explain where and how it is used.

1. Grapevine
2. Domain name system
3. Kerberos
4. ISIS
5. Quorum Consensus

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Mid-Term Review

Naming - For each of systems listed below, indicate the style of naming that is supported from the set of host based naming, global naming, attributed based naming, and user/system/object centered. Note that some systems will employ more than one style of naming, and if that is the case you should identify each style and the part of the system that employs each. Explain in no more than two sentences why you classify the naming in each system as you have. In one or two additional sentences identify the technique(s) used to implement (including techniques for improving performance if relevant) for the naming mechanism you described.

1. **Host names (fully qualified domain names) from the domain name system:**
2. **File names in the sprite file system:**
3. **URLs on the world wide web:**
4. **File name resolution in Amoeba, as well as object names/handles in Amoeba:**
5. **File names on various computer running the Unix file system (including files mounted through NFS, Samba, or other distributed file systems):**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Mid-Term Review

You have been hired to design a “cloud” based collaboration service which is supposed to allow users to store files remotely, and share selected files with other users. Files to be shared can be identified by inclusion in “published” directories, but files are individually identifiable, and shared files might exist in the “directories” of multiple users sharing the file.

1. **Naming** – Discuss the naming approaches you would use for such a system. Is there only a single naming mechanisms, and if so what is it. If there are multiple naming mechanisms to be used, describe each such mechanism and the purpose for each in your design (8 points)
2. **Synchronization** – Discuss the tradeoffs of different approaches to synchronization in your system. Consider both approaches that allow only serial sharing (i.e. files open for write can only be accessed by one user at a time), vs. files supporting concurrent writes, or reads and writes. In answering this question, describe the synchronization methods needed to support both kinds of sharing and the performance and other implications of the synchronization methods to be deployed (12 points).

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Mid-Term Review

You have been hired to design a “cloud” based collaboration service which is supposed to allow users to store files remotely, and share selected files with other users. Files to be shared can be identified by inclusion in “published” directories, but files are individually identifiable, and shared files might exist in the “directories” of multiple users sharing the file.

1. **Replication and caching** – Discuss the use of replication and caching in your design? How will your approaches improve performance? How do the synchronization issues discussed in [b] affect the replica and cache consistency techniques that must be applied – and what impact will these issues have on performance? (10 points)
2. **Security** – Discuss approaches for providing security of the files stored in the cloud storage service. How are policies for access to files managed? Are these policies managed in access control lists, or in capability lists (either is fine – explain how it works in your design and why your approach is sufficient). What kinds of permissions can be provided to authorized users, and how would one grant access to a shared file? (10 points)

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

CSci555: Advanced Operating Systems

Lecture 14 – December 2nd, 2011

Scale in Distributed Systems

Dr. Clifford Neuman
University of Southern California
Information Sciences Institute

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Announcements

- ❖ **Research paper due today**
 - ❑ **Late submissions with small penalty**
- ❖ **Class evaluations at Break**
 - ❑ **DEN students please return online**
- ❖ **Final Exam**
 - ❑ **Friday December 9, 2PM-4PM**
 - ❑ **Location to be determined**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Scale in Distributed Systems - Neuman

- ❖ **A system is said to be scalable if it can handle the addition of users and resources without suffering a noticeable loss of performance or increase in administrative complexity.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Three dimensions of scale

❖ Numerical

- ❑ Number of objects, users

❖ Geographic

- ❑ Where the users and resources are

❖ Administrative

- ❑ How many organizations own or use different parts of the system

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Effects of Scale

❖ Reliability

- ❑ Autonomy, Redundancy

❖ System Load

- ❑ Order of growth

❖ Administration

- ❑ Rate of change
- ❑ Heterogeneity

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Techniques - Replication

- ❖ **Placement of replicas**
 - ❑ **Reliability**
 - ❑ **Performance**
 - ❑ **Partition**
 - ❑ **What if all in one place**
- ❖ **Consistency**
 - ❑ **Read-only**
 - ❑ **Update to all**
 - ❑ **Primary Site**
 - ❑ **Loose Consistency**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Techniques - Distribution

- ❖ **Placement of servers**
 - ❑ **Reliability**
 - ❑ **Performance**
 - ❑ **Partition**
- ❖ **Finding the right server**
 - ❑ **Hierarchy/iteration**
 - ❑ **Broadcast**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

❖ **Placement of Caches**

- ❑ **Multiple places**

❖ **Cache consistency**

- ❑ **Timeouts**
- ❑ **Hints**
- ❑ **Callback**
- ❑ **Snooping**
- ❑ **Leases**

CSci555: Advanced Operating Systems

Lecture 14 – December 2nd, 2011
Selected Topics and Discussions

Dr. Clifford Neuman
University of Southern California
Information Sciences Institute

Is the OS still relevant

- ❖ **What is the role of an OS in the internet**
 - ❑ **Are today's computers appliances for accessing the web?**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Is the OS still relevant

- ❖ **OS Manages local resources**
 - ❑ **Provides protection between applications**
 - ❑ **Though the role seems diminished, it is actually increasing in importance**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Today's File Systems

- ❖ **Network Attached Storage**
- ❖ **Cloud Storage**
- ❖ **Content Distribution Systems**
- ❖ **Peer to Peer File Systems**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Content Delivery

- ❖ **Pre-staging of content**
- ❖ **Techniques needed to redirect to local copy.**
- ❖ **Ideally need ways to avoid central bottleneck.**
- ❖ **Use of URN's can help, but needs underlying changes to browsers.**
 - ❑ **For dedicated apps, easier to deploy**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Naming Today

- ❖ **URL's vs URN's**
- ❖ **System based identifiers**
 - **Facebook**
 - **Twitter**
 - **Tiny URL's**
 - **These make the problem worse in the interest of locking users into their system.**
- ❖ **Internationalized Domain Names**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Multi-Core Systems

- ❖ **Shared Memory Multiprocessor**
 - **But few apps know how to take advantage of it**
 - **But modern OS – many processes**
- ❖ **Still leaves contention for other resources**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Internet Search Techniques

❖ Issues

- ❑ How much of the net to index
 - How much detail
 - How to select
 - ❑ Relevance of results
 - Ranking results – avoiding spam
 - Context for searching
 - Transitive indexing
- ### ❖ Scaling the search engines

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Internet Search Techniques - Google

❖ Data Distribution

- ❑ Racks and racks of servers running Linux – key data is replicated
 - Some for indices
 - Some for storing cached data
 - ❑ Query distributed based on load
 - ❑ Many machines used to for single query
- ### ❖ Page rank
- ❑ When match found, ranking by number and quality of links to the page.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

The Structure of Distributed Systems

- ❖ **Client server**
- ❖ **Object Oriented**
- ❖ **Peer to Peer (additional discussion)**
- ❖ **Cloud Based**
- ❖ **Federated**
- ❖ **Agent Based**
- ❖ **Virtualized**
- ❖ **Embedded**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Peer to Peer

- ❖ **Peer to peer systems are client server systems where the client is also a server.**
- ❖ **The important issues in peer to peer systems are really:**
 - ❑ **Trust – one has less trust in servers**
 - ❑ **Unreliability – Nodes can drop out at will.**
 - ❑ **Management – need to avoid central control (a factor caused by unreliability)**
- ❖ **Ad hoc network related to peer to peer**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Future of Distributed Systems

- ❖ **More embedded systems (becoming less “embedded”).**
 - ❑ **Process control / SCADA**
 - ❑ **Real time requirements**
 - ❑ **Protection from the outside**
 - ❑ **Are they really embedded?**
- ❖ **Stronger management of data flows across applications.**
- ❖ **Better resource management across organizational domains.**
- ❖ **Multiple views of available resources.**
- ❖ **Hardware abstraction**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Hardware Abstraction

- ❖ **Many operating systems are designed today to run on heterogeneous hardware**
- ❖ **Hardware abstraction layer often part of the internal design of the OS.**
 - ❑ **Small set of functions**
 - ❑ **Called by main OS code**
- ❖ **Usually limited to some similarity in hardware, or the abstraction code becomes more complex and affects performance.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Emulation and Simulation

- ❖ **Need techniques to test approaches before system is built.**
 - ❑ **Simulations**
 - ❑ **Need real data sets to model assumptions.**
- ❖ **Need techniques to test scalability before system is deployed.**
 - ❑ **Deployment harder than implementation**
 - ❑ **Emulations and simulations beneficial**
- ❖ **Issues in emulation and simulation**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Windows

- ❖ **XP, Win2K and successors based loosely on Mach Kernel.**
- ❖ **Techniques drawn from many other research systems.**
- ❖ **Backwards compatibility has been an issue affecting some aspects of its architecture.**
- ❖ **Despite common criticism, the current versions make a pretty good system for general computing needs.**

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

Miscellaneous

- ❖ **Security issues with the Domain Name System**
 - ❑ A result of multi-level caching
 - ❑ And security not considered up front
- ❖ **Neutrality in Distributed Systems**
 - ❑ Protocols
 - ❑ Net Neutrality
 - ❑ Application frameworks / middleware
- ❖ **Unix and Linux**
 - ❑ Kernel Structure
 - ❑ Filesystems

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

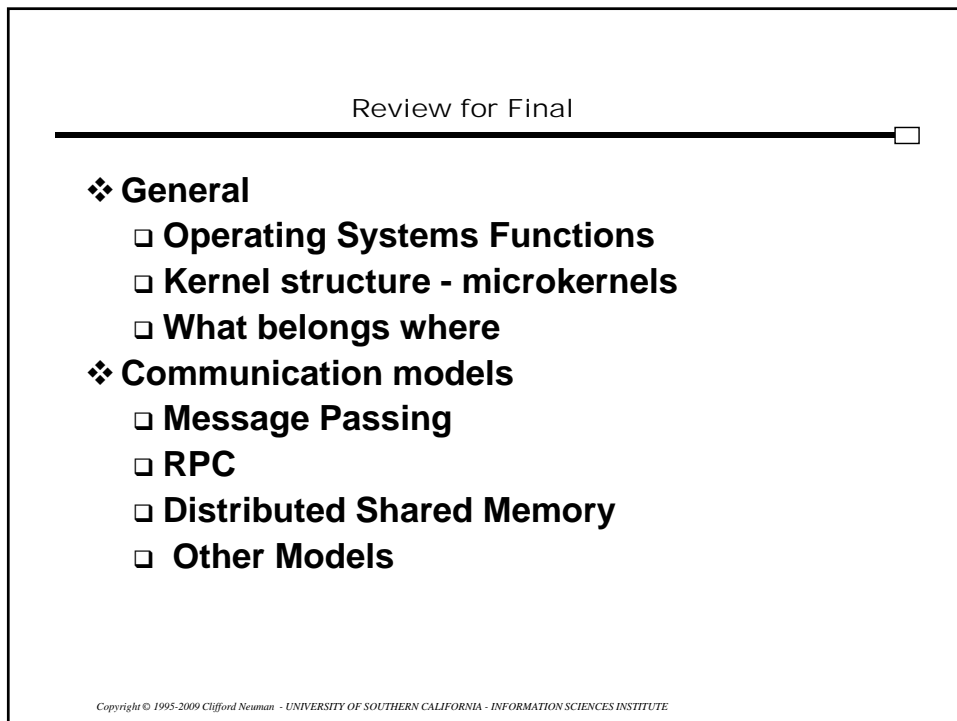
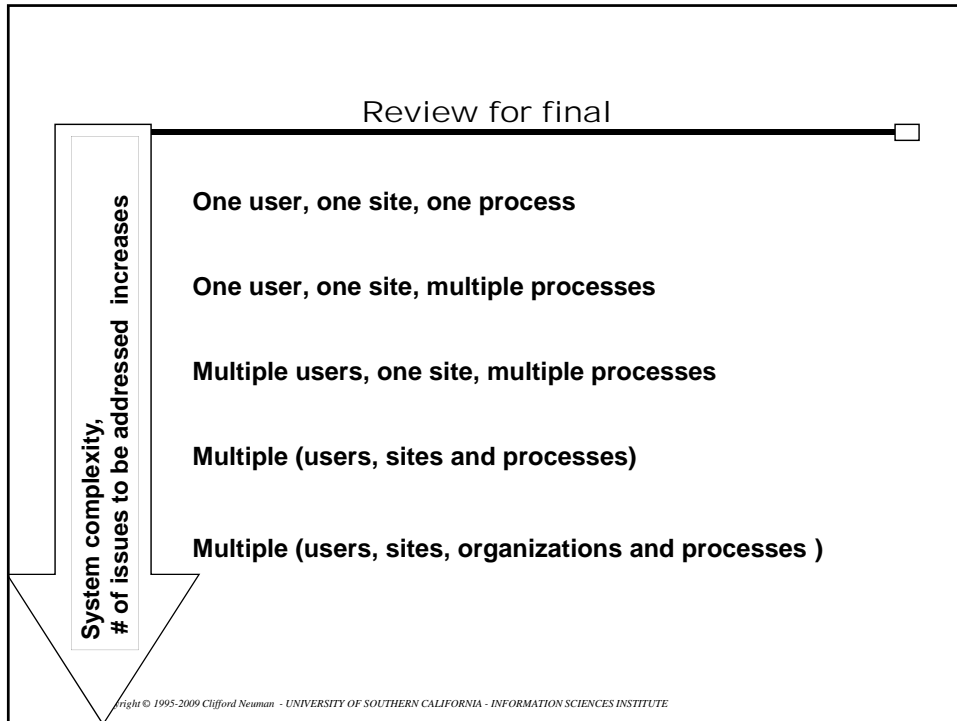
CSci555: Advanced Operating Systems

Lecture 14 - December 2nd, 2011

REVIEW

Dr. Clifford Neuman
University of Southern California
Information Sciences Institute

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE



❖ **Synchronization - Transactions**

- ❑ Time Warp
- ❑ Reliable multicast/broadcast

❖ **Naming**

- ❑ Purpose of naming mechanisms
- ❑ Approaches to naming
- ❑ Resource Discovery
- ❑ Scale

❖ **Security – Requirements**

- ❑ Confidentiality
- ❑ Integrity
- ❑ Availability

❖ **Security mechanisms (prevention/detection)**

- ❑ Protection
- ❑ Authentication
- ❑ Authorization (ACL, Capabilities)
- ❑ Intrusion detection
- ❑ Audit

❖ **Cooperation among the security mechanisms**

❖ **Scale**

❖ **Distributed File Systems - Caching**

- ❑ **Replication**
- ❑ **Synchronization**
 - **voting, master/slave**
- ❑ **Distribution**
- ❑ **Access Mechanism**
- ❑ **Access Patterns**
- ❑ **Availability**

❖ **Other file systems**

- ❑ **Log Structured**
- ❑ **RAID**

❖ **Case Studies**

- ❑ **Locus**
- ❑ **Athena**
- ❑ **Andrew**
- ❑ **V**
- ❑ **HCS**
- ❑ **Amoeba**
- ❑ **Mach**
- ❑ **CORBA**

- ❖ **Resource Allocation**
- ❖ **Real time computing**
- ❖ **Fault tolerant computing**

SCALE

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

2006 Exam – 1a Scalability

1a) System load (10 points) – Suggest some techniques that can be used to reduce the load on individual servers within a distributed system? Provide examples of how these techniques are used from each of the following systems: The Domain Name System, content delivery through the world wide web, remote authentication in the Kerberos system. Note that some of the systems use more than one technique.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

2006 Exam – 1b Scalability

1b) Identifying issues (20 points) for each of the techniques described in part (a) there are issues that must be addressed to make sure that the system functions properly (I am interested in the properly aspect here, not the most efficiently aspect). For each technique identify the primary issues that needs to be addressed and explain how it is addressed in each of the listed systems that uses the technique.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

2006 Exam – 2 Kernel

2) For each of the operating system functions listed below list the benefits and drawbacks to placing the function in the Kernel, leaving the function to be implemented by the application, or providing the function in users space through a server (the server case includes cases where the application selects and communicates with a server, and also the case where the application calls the kernel, but the processing is redirected by the kernel to a server). For each function, suggest the best location(s) to provide this function. If needed you can make an assumption about the scenario for which the system will be used. Justify your choice for placement of this function. There may be multiple correct answers for this last part – so long as your justification is correct.

- File System
- Virtual Memory
- Communications
- Scheduling
- Security

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

2006 Exam – 3 Design Problem – Fault Tolerance

- 3) You are designing a database system that requires significant storage and processing power. Unfortunately, you are stuck using the hardware that was ordered by the person whose job you just filled. This morning, the day after you first arrived at work, a truck arrived with 10 processors (including memory, network cards, etc), 50 disk drives, and two uninterruptible power supplies. The failure rates of the processors (including all except the disk drives and power supplies) is λ_p . The failure rates on the disk drives is λ_d , and the failure rate for the power supplies is λ_e .
- You learned from your supervisor that the reason they let the last person go is that he designed the system so that the failure of any of the components would cause the system to stop functioning. In terms of $\lambda_p, d,$ and e , what is the failure probability for the system as a whole. (5 points)
 - The highest expected load on your system could be handled by about half the processors. The largest expected dataset size that is expected is about 1/3 the capacity of the disks that arrived. Suggest a change to the structure of the system, using the components that have already arrived, that will yield better fault tolerance. In terms of $\lambda_p, d,$ and e , what is the failure probability for the new system? (note, there are easy things and harder things you can do here, I suggest describing the easing things, generating the probability based on that approach, and then just mentioning some of the additional steps that could be taken to further improve the fault tolerance (15 points)
 - List some of the problems that you would need to solve or some of the assumptions you would need to make, in order to construct the system described in part b from the components that arrived this morning (things like number of network interfaces per processor, how the disks are connected to processors or the network). Discuss also any assumptions you need to make regarding detect ability of failures, and describe your approach to failover (how will the failures be masked, what steps are taken when a failure occurs). (15 points)

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

2007 Exam – 1a Leases

For each of the following approaches to consistency, if they were to be implemented as a lease, list the corresponding lease term, and the rules for breaking the lease (i.e. if the normal rules for breaking a lease are not provided by the system, what are the effective rules of the mechanism. (16 points)

- AFS-2/3 Callback
- AFS-1 Check-on-use
- Time to live in the domain name system
- Locks in a transaction system

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

2007 Exam – 1b Log Structured File Systems

- A. Discuss the similarity between a transaction system and the log structure file system.
- B. How does the log structure file system improve the performance of writes to the file system?
- C. Why does it take so much less time to recover from a system crash in a log structured file system than it does in the traditional Unix file system? How is recovery accomplished in the log structure approach?

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

2007 Exam – 2 Kernels

For a general purpose operating system such Linux, discuss the placement of services, listing those functions that should be provided by the kernel, by the end application itself, and by application level servers. Specifically, what OS functions should be provided in each location? Justify your answer and state your assumptions.

- a) In the Kernel itself
- b) In the application itself
- c) In servers outside the kernel

For a system supporting embedded applications, such as process control, what changes would you make in the placement of OS functions (i.e. what would be different than what you described in a-c). Justify your answer.

Copyright © 1995-2009 Clifford Neuman - UNIVERSITY OF SOUTHERN CALIFORNIA - INFORMATION SCIENCES INSTITUTE

2007 Exam – 3 Design Problem

You have been hired to build a system to manage ticket sales for large concerts. This system must be highly scalable supporting near simultaneous request from the “flash crowds” accessing the system the instant a new concert goes on sale. The system must accept requests fairly, so that ticket consolidators are unable to “game the system” to their advantage through automated programs on well placed client machines located close to the servers in terms of network topology. To handle the load will require multiple servers all with access to the ticketing database, yet synchronization is a must as we can’t sell the same seat to more than one person. The system must support several functions, among which are providing venue and concert information to potential attendees, displaying available seats, reserving seats, and completing the sale (collecting payment, recording the sale, and enabling the printing of a barcode ticket).

- a) Describe the architecture of your system in terms of the allocation of functions across processors. Will all processors be identical in terms of their functionality, or different servers provide different functions, and if so which ones and why?
- b) Explain the transactional characteristics of your system. In particular, when does a transaction begin, and when does it commit or abort, and which processors (according to the functions described by you in part a) will be participants in the transaction.
- c) What objects will have associated locks and when will these object be locked and unlocked.
- d) How will you use replication in your system and how will you manage consistency of such replicated data
- e) How will you use distribution in your system