

Distributed Systems Middleware

A thick, horizontal yellow brushstroke with a textured, painterly appearance, spanning the width of the slide below the title.

Prof. Nalini Venkatasubramanian
Dept. of Information & Computer
Science
University of California, Irvine

ICS 280 - Distributed Systems Middleware



Lecture 1 - Introduction to Distributed
Systems Middleware

Mondays, Wednesdays 3:30-5:00p.m.

Prof. Nalini Venkatasubramanian

nalini@ics.uci.edu

Course logistics and details

- Course Web page -
 - <http://www.ics.uci.edu/~ics243f>
- Lectures - MW 3:30-4:50p.m, Course Laboratories - machines on 3rd floor CS labs
- ICS 280 Reading List
 - Technical papers and reports
 - Reference Books

Course logistics and details



- Homeworks
 - Paper summaries
 - Survey paper
- Course Presentation
- Course Project
 - Maybe done individually, in groups of 2 or 3(max)
 - Potential projects on webpage

ICS 280 Grading Policy

- Homeworks - 30%
 - 1 paper summary due every week
 - (3 randomly selected each worth 10% of the final grade). -
- Project Survey Paper - 10%
- Class Presentation - 10%
- Class Project - 50% of the final grade
- Final assignment of grades will be based on a curve.

Lecture Schedule

- **Weeks 1 and 2:**
 - | **Middleware and Distributed Computing Fundamentals**
 - | **Fundamentals of Concurrency**
 - | **General Purpose Middleware - Technical challenges**
 - | **Adaptive Computing**
- **Weeks 3 and 4: Distributed Systems Management**
 - | **Distributed Operating Systems**
 - | **Messaging and Communication in Distributed Systems**
 - | **Naming and Directory Services**
 - | **Distributed I/O and Storage Subsystems**
 - | **Distributed Resource Management**
- **Week 5 and 6: Distributed Object Models**
 - | **Concurrent Objects – Actors, Infospheres**
 - | **Common Object Services**
 - | **Synchronization with Distributed Objects**
 - | **Composing Distributed Objects**

Course Schedule

- Weeks 7 and 8: Middleware Frameworks - Case Studies
 - | **DCE**
 - | **CORBA, RT-CORBA**
 - | **Jini**
 - | **Espeak, XML based middleware**
- Weeks 9 and 10: Middleware for Distributed Application Environments
 - | **QoS-enabled middleware**
 - | **Fault tolerant applications**
 - | **Secure applications**
 - | **Transaction Based applications**
 - | **Ubiquitous and Mobile Environments**

Introduction



■ Distributed Systems

- Multiple independent computers that appear as one
- Lamport's Definition
 - | " You know you have one when the crash of a computer you have never heard of stops you from getting any work done."
- "A number of interconnected autonomous computers that provide services to meet the information processing needs of modern enterprises."

Characterizing Distributed Systems

■ Multiple Computers

- each consisting of CPU's, local memory, stable storage, I/O paths connecting to the environment

■ Interconnections

- some I/O paths interconnect computers that talk to each other

■ Shared State

- systems cooperate to maintain shared state
- maintaining global invariants requires correct and coordinated operation of multiple computers.

Examples of Distributed Systems



- Banking systems
- Communication - email
- Distributed information systems
 - WWW
 - Federated Databases
- Manufacturing and process control
- Inventory systems
- General purpose (university, office automation)

Why Distributed Computing?

- Inherent distribution
 - Bridge customers, suppliers, and companies at different sites.
- Speedup - improved performance
- Fault tolerance
- Resource Sharing
 - Exploitation of special hardware
- Scalability
- Flexibility

Why are Distributed Systems Hard?

- Scale
 - numeric, geographic, administrative
- Loss of control over parts of the system
- Unreliability of message passing
 - unreliable communication, insecure communication, costly communication
- Failure
 - Parts of the system are down or inaccessible
 - Independent failure is desirable

Design goals of a distributed system

- Sharing
 - HW, SW, services, applications
- Openness(extensibility)
 - use of standard interfaces, advertise services, microkernels
- Concurrency
 - compete vs. cooperate
- Scalability
 - avoids centralization
- Fault tolerance/availability
- Transparency
 - location, migration, replication, failure, concurrency

END-USER

Application Developer

- Personalized Environment
- Predictable Response
- Location Independence
- Platform Independence

- Code Reusability
- Interoperability
- Portability
- Reduced Complexity
- Flexibility
- Real-Time Access to information
- Scalability
- Faster Developmt. And deployment of Business Solutions

- Increased Complexity
- Lack of Mgmt. Tools
- Changing Technology

System Administrator

ORGANIZATION

[Khanna94]

Classifying Distributed Systems

- Based on degree of synchrony
 - Synchronous
 - Asynchronous
- Based on communication medium
 - Message Passing
 - Shared Memory
- Fault model
 - Crash failures
 - Byzantine failures

Computation in distributed systems

- Asynchronous system
 - no assumptions about process execution speeds and message delivery delays
- Synchronous system
 - make assumptions about relative speeds of processes and delays associated with communication channels
 - constrains implementation of processes and communication
- Models of concurrency
 - Communicating processes
 - Functions, Logical clauses
 - Passive Objects
 - Active objects, Agents

Concurrency issues

- Consider the requirements of transaction based systems
 - Atomicity - either all effects take place or none
 - Consistency - correctness of data
 - Isolated - as if there were one serial database
 - Durable - effects are not lost
- General correctness of distributed computation
 - Safety
 - Liveness

Communication in Distributed Systems

- Provide support for entities to communicate among themselves
 - Centralized (traditional) OS's - local communication support
 - Distributed systems - communication across machine boundaries (WAN, LAN).
- 2 paradigms
 - Message Passing
 - Processes communicate by sharing messages
 - Distributed Shared Memory (DSM)
 - Communication through a virtual shared memory.

Message Passing

■ Basic communication primitives

- | Send message
- | Receive message

■ Modes of communication

| Synchronous

- atomic action requiring the participation of the sender and receiver.
- Blocking send: blocks until message is transmitted out of the system send queue
- Blocking receive: blocks until message arrives in receive queue

| Asynchronous

- Non-blocking send: sending process continues after message is sent
- Blocking or non-blocking receive: Blocking receive implemented by timeout or threads. Non-blocking receive proceeds while waiting for message. Message is queued (BUFFERED) upon arrival.

Reliability issues

■ Unreliable communication

- Best effort, No ACK's or retransmissions
- Application programmer designs own reliability mechanism

■ Reliable communication

- Different degrees of reliability
- Processes have some guarantee that messages will be delivered.
- Reliability mechanisms - ACKs, NACKs.

Reliability issues

■ Unreliable communication

- Best effort, No ACK's or retransmissions
- Application programmer designs own reliability mechanism

■ Reliable communication

- Different degrees of reliability
- Processes have some guarantee that messages will be delivered.
- Reliability mechanisms - ACKs, NACKs.

Distributed Shared Memory

- Abstraction used for processes on machines that do not share memory
 - Motivated by shared memory multiprocessors that do share memory
- Processes read and write from virtual shared memory.
 - Primitives - read and write
 - OS ensures that all processes see all updates
- Caching on local node for efficiency
 - Issue - cache consistency

Remote Procedure Call

- Builds on message passing
 - extend traditional procedure call to perform transfer of control and data across network
 - Easy to use - fits well with the client/server model.
 - Helps programmer focus on the application instead of the communication protocol.
 - Server is a collection of exported procedures on some shared resource
 - Variety of RPC semantics
 - “maybe call”
 - “at least once call”
 - “at most once call”

Fault Models in Distributed Systems

■ Crash failures

- A processor experiences a crash failure when it ceases to operate at some point without any warning. Failure may not be detectable by other processors.
 - Failstop - processor fails by halting; detectable by other processors.

■ Byzantine failures

- completely unconstrained failures
- conservative, worst-case assumption for behavior of hardware and software
- covers the possibility of intelligent (human) intrusion.

Other Fault Models in Distributed Systems

■ Dealing with message loss

■ Crash + Link

- | Processor fails by halting. Link fails by losing messages but does not delay, duplicate or corrupt messages.

■ Receive Omission

- | processor receives only a subset of messages sent to it.

■ Send Omission

- | processor fails by transmitting only a subset of the messages it actually attempts to send.

■ General Omission

- | Receive and/or send omission

Other distributed system issues

- Concurrency and Synchronization
- Distributed Deadlocks
- Time in distributed systems
- Naming
- Replication
 - improve availability and performance
- Migration
 - of processes and data
- Security
 - eavesdropping, masquerading, message tampering, replaying

Client/Server Computing

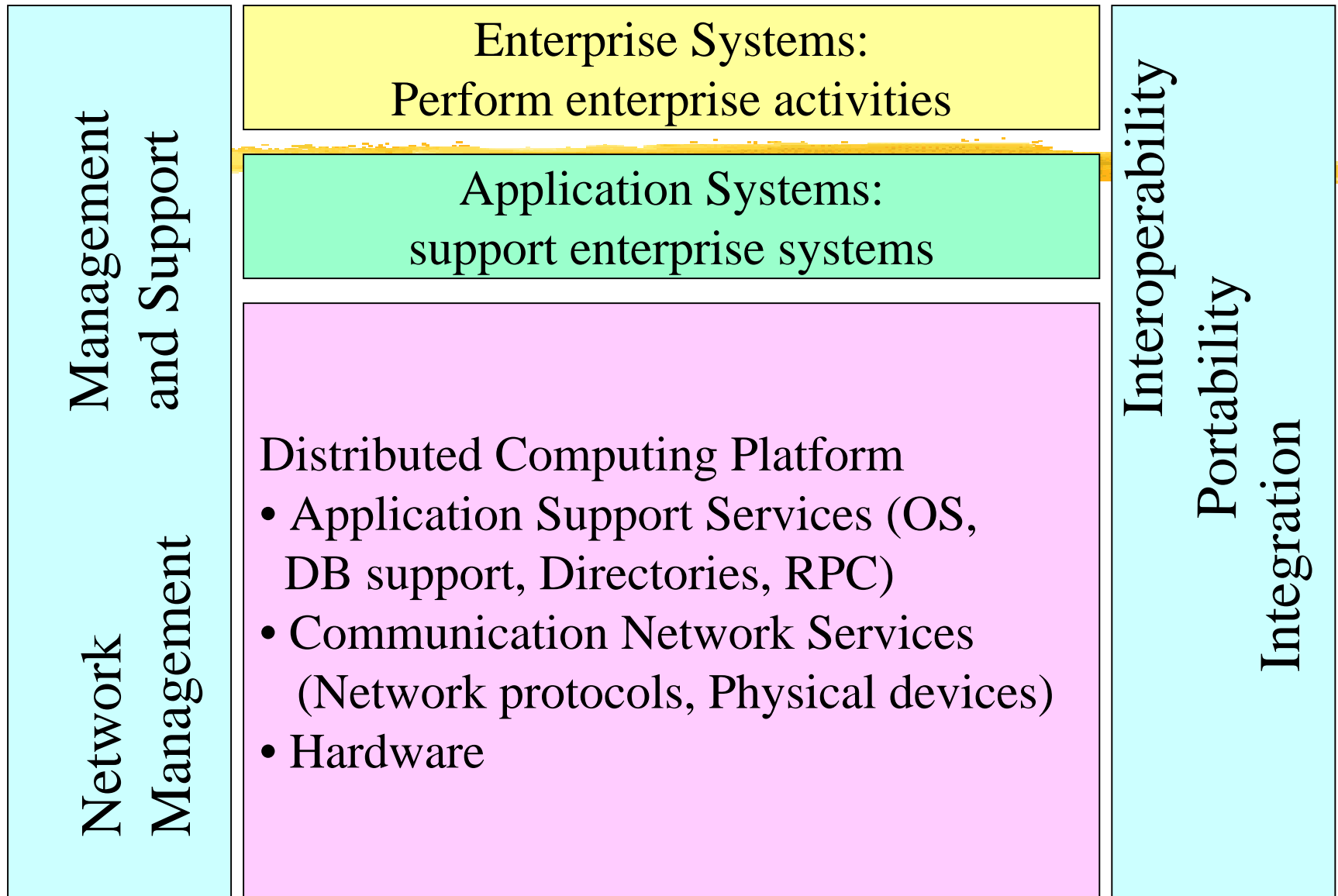


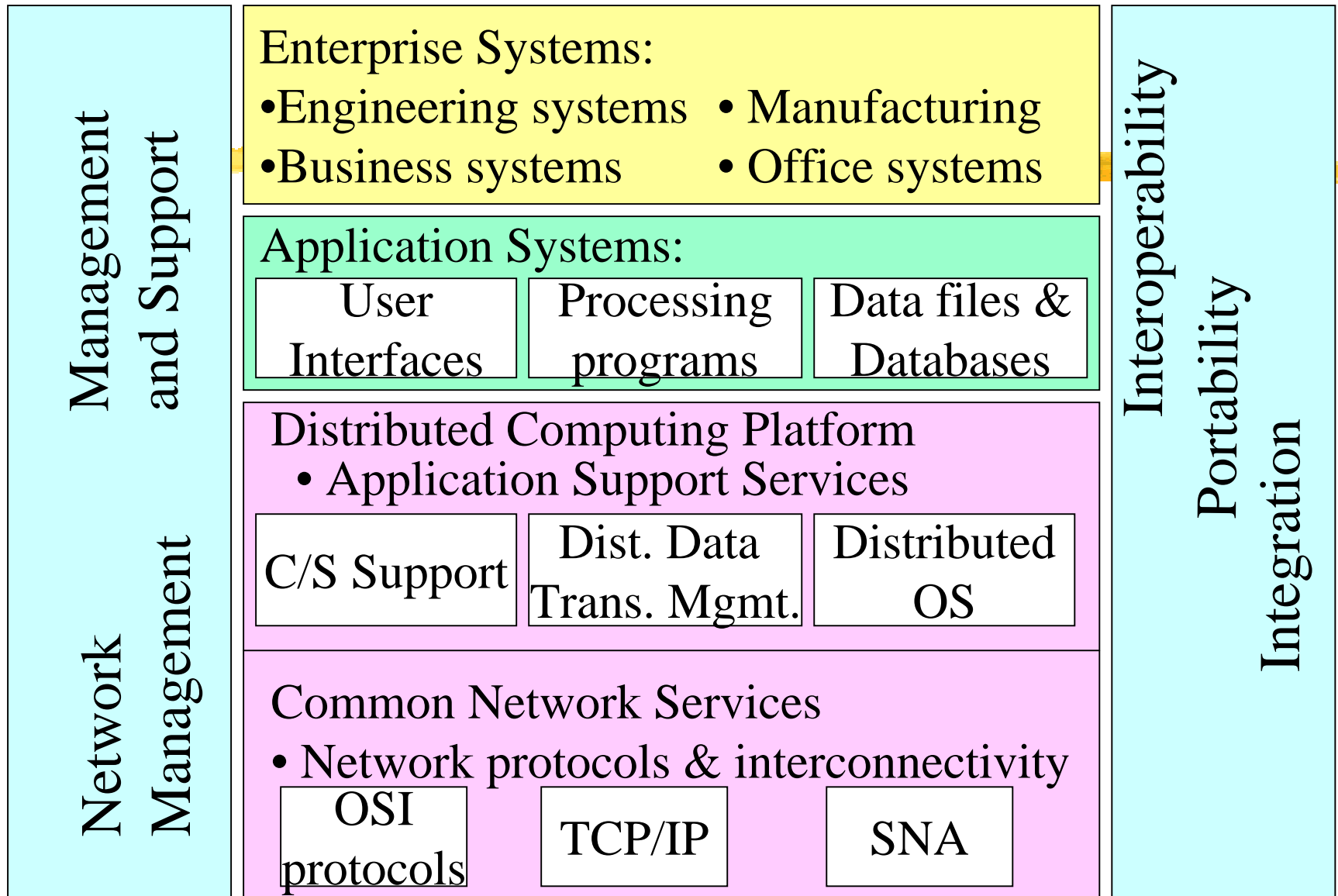
- Client/server computing allocates application processing between the client and server processes.
- A typical application has three basic components:
 - Presentation logic
 - Application logic
 - Data management logic

Client/Server Models



- There are at least three different models for distributing these functions:
 - Presentation logic module running on the client system and the other two modules running on one or more servers.
 - Presentation logic and application logic modules running on the client system and the data management logic module running on one or more servers.
 - Presentation logic and a part of application logic module running on the client system and the other part(s) of the application logic module and data management module running on one or more servers





What is Middleware?

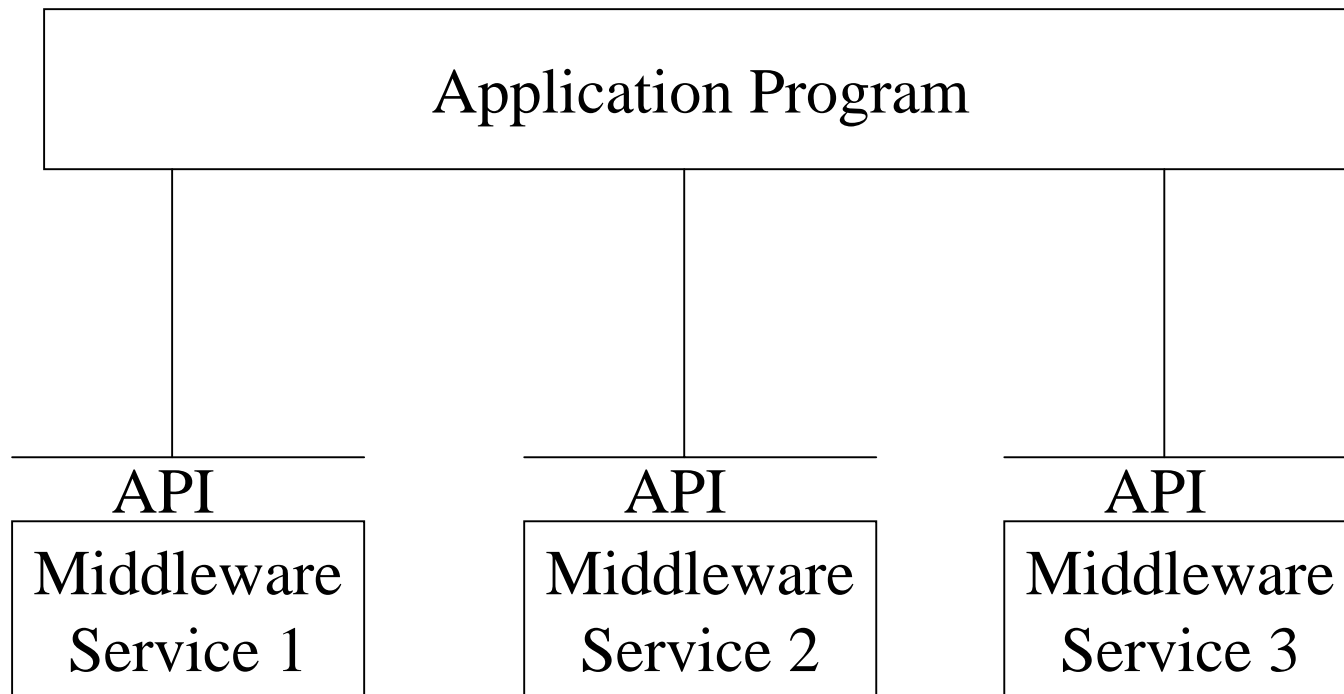


- Middleware is the software between the application programs and the operating System and base networking
- Middleware provides a comprehensive set of higher-level distributed computing capabilities and a set of interfaces to access the capabilities of the system.

Distributed Systems Middleware

- Enables the modular interconnection of distributed software
 - | abstract over low level mechanisms used to implement resource management services.
- Computational Model
 - | Support separation of concerns and reuse of services
- Customizable, Composable Middleware Frameworks
 - | Provide for dynamic network and system customizations, dynamic invocation/revocation/installation of services.
 - | Concurrent execution of multiple distributed systems policies.

Modularity in Middleware Services



Useful Middleware Services



- Naming and Directory Service
- State Capture Service
- Event Service
- Transaction Service
- Fault Detection Service
- Trading Service
- Replication Service
- Migration Service

Types of Middleware Services

- Component services
 - Provide a specific function to the requestor
 - Generally independent of other services
 - Presentation, Communication, Control, Information Services, computation services etc.
- Integrated Sets
- Integration frameworks

Integrated Sets Middleware

- An Integrated set of services consist of a set of services that take significant advantage of each other.
- Example: DCE

Distributed Computing Environment (DCE)

- DCE is from the Open Software Foundation (OSF), and now X/Open, offers an environment that spans multiple architectures, protocols, and operating systems.
 - DCE supported by major software vendors.
- It provides key distributed technologies, including RPC, a distributed naming service, time synchronization service, a distributed file system, a network security service, and a threads package.

DCE

DCE Security Service	Applications			Management
	DCE Distributed File Service			
	DCE Distributed Time Service	DCE Directory Service	Other Basic Services	
	DCE Remote Procedure Calls			
	DCE Threads Services			
Operating System Transport Services				

Integration Frameworks Middleware



- Integration frameworks are integration environments that are tailored to the needs of a specific application domain.
- Examples
 - Workgroup framework - for workgroup computing.
 - Transaction Processing monitor frameworks
 - Network management frameworks

Distributed Object Computing

- Combining distributed computing with an object model.
 - Allows software reusability
 - More abstract level of programming
 - The use of a broker like entity that keeps track of processes, provides messaging between processes and other higher level services
 - Examples
 - CORBA
 - JINI
 - E-SPEAK
 - Note: DCE uses a procedure-oriented distributed systems model, not an object model.

Issues with Distributed Objects



- Abstraction
- Performance
- Latency
- Partial failure
- Synchronization
- Complexity

Techniques for object distribution

■ Message Passing

- | Object knows about network; Network data is minimum

■ Argument/Return Passing

- | Like RPC. Network data = args + return result + names

■ Serializing and Sending Object

- | Actual object code is sent. Might require synchronization. Network data = object code + object state + sync info

■ Shared Memory

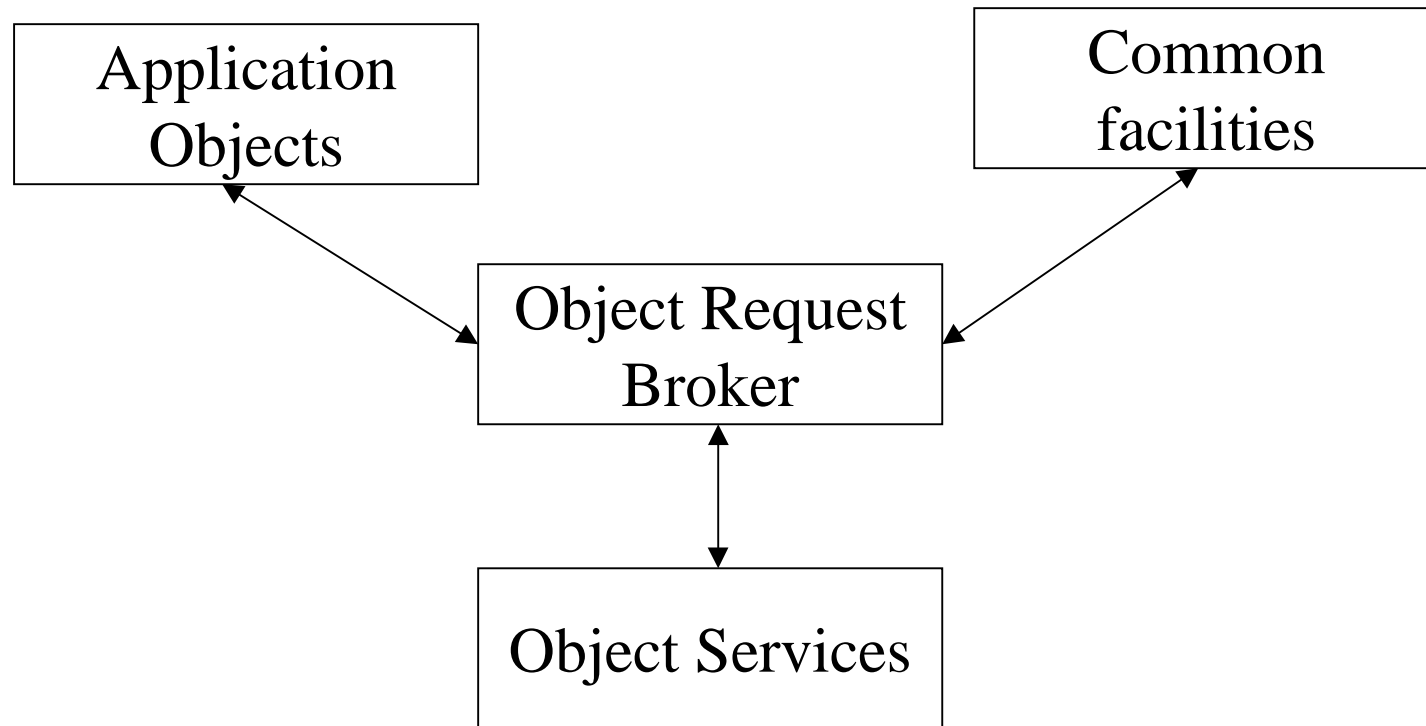
- | based on DSM implementation
- | Network Data = Data touched + synchronization info

CORBA



- CORBA is a standard specification for developing object-oriented applications.
- CORBA was defined by OMG in 1990.
- OMG is dedicated to popularizing Object-Oriented standards for integrating applications based on existing standards.

The Object Management Architecture (OMA)



OMA

- ORB: the communication hub for all objects in the system
- Object Services: object events, persistent objects, etc.
- Common facilities: accessing databases, printing files, etc.
- Application objects: document handling objects.

Clock Synchronization in Distributed Systems

- Clocks in a distributed system drift:
 - Relative to each other
 - | Logical Clocks are clocks which are synchronized relative to each other.
 - Relative to a real world clock
 - | Determination of this real world clock may be an issue
 - | Physical clocks are logical clocks that must not deviate from the real-time by more than a certain amount.

Synchronizing Logical Clocks

- Need to understand the ordering of events
- Notion of time is critical
- “Happens Before” notion.
 - E.g. Concurrency control using timestamps
- “Happens Before” notion is not straightforward in distributed systems
 - No guarantees of synchronized clocks
 - Communication latency

Event Ordering

- Lamport defined the “happens before” ($=>$) relation
 - If a and b are events in the same process, and a occurs before b , then $a => b$.
 - If a is the event of a message being sent by one process and b is the event of the message being received by another process, then $a => b$.
 - If $X => Y$ and $Y => Z$ then $X => Z$.
- If $a => b$ then $time(a) => time(b)$*

Causal Ordering

- “Happened Before” also called causal ordering
- Possible to draw a causality relation between 2 events if
 - They happen in the same process
 - There is a chain of messages between them

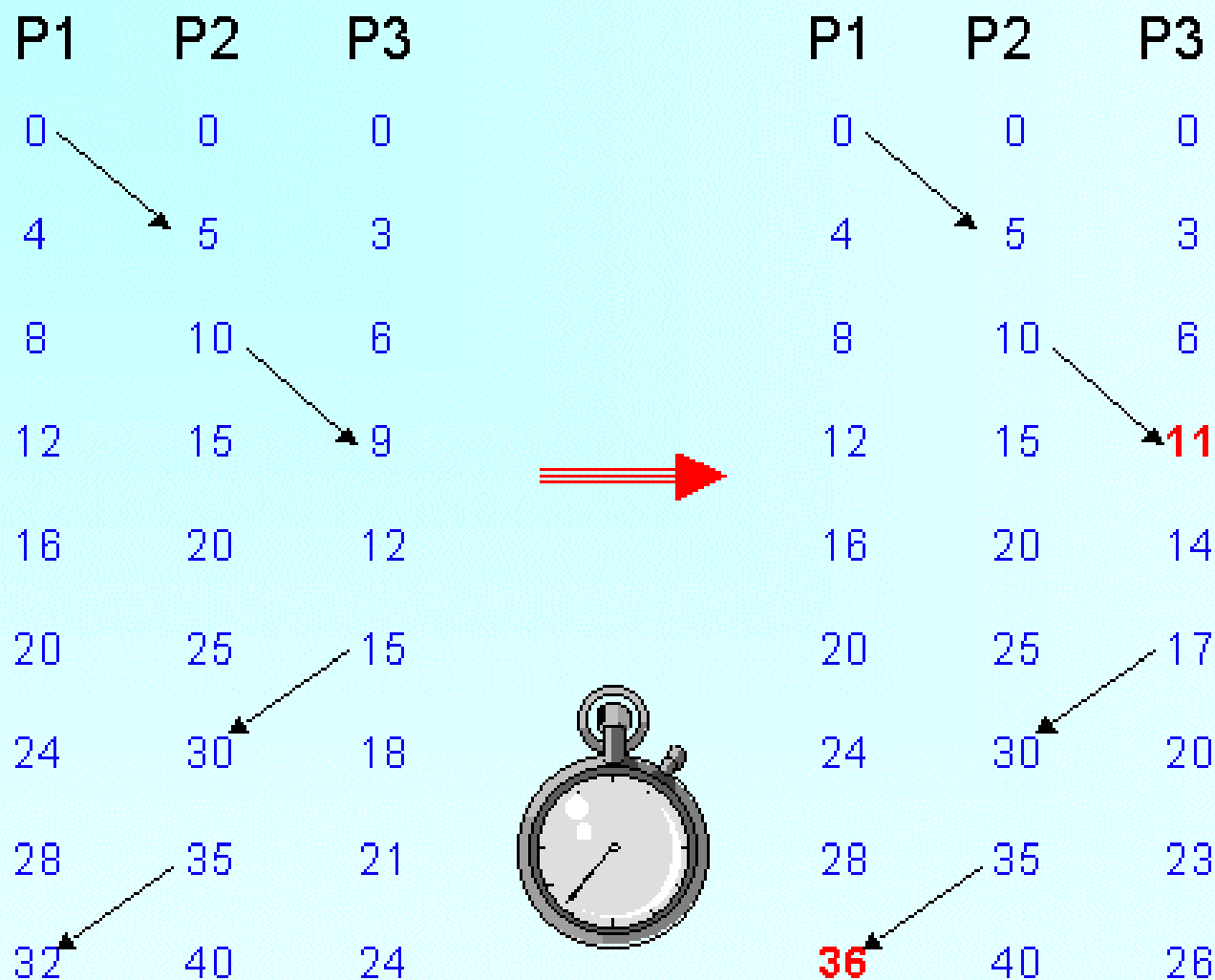
Logical Clocks

- Monotonically increasing counter
- No relation with real clock
- Each process keeps its own logical clock C_p used to timestamp events

Causal Ordering and Logical Clocks

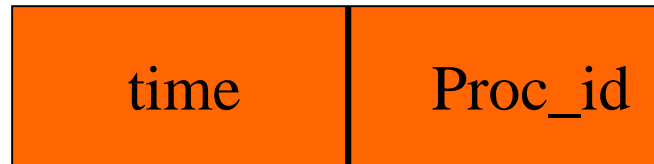
- C_p is incremented before each event.
 - $C_p = C_p + 1$
- When p sends a message m , it piggybacks a logical timestamp $t = C_p$.
- When q receives (m, t) it computes:
 - $C_q = \max(C_q, t)$ before timestamping the message receipt event.
- Results in a partial ordering of events.

Lamport Logical Clock



Total Ordering

- Extending partial order to total order



- Global timestamps:
 - (Ta, Pa) where Ta is the local timestamp and Pa is the process id.
 - $(Ta, Pa) < (Tb, Pb)$ iff
 - $(Ta < Tb)$ or $((Ta = Tb) \text{ and } (Pa < Pb))$
 - Total order is consistent with partial order.

Physical Clocks

- How do we measure real time?
 - 17th century - Mechanical clocks based on astronomical measurements
 - Solar Day - Transit of the sun
 - Solar Seconds - $\text{Solar Day} / (3600 * 24)$
 - Problem (1940) - Rotation of the earth varies (gets slower)
 - Mean solar second - average over many days

Atomic Clocks

■ 1948

- counting transitions of a crystal (Cesium 133) used as atomic clock
- TAI - International Atomic Time
 - 9192631779 transitions = 1 mean solar second in 1948
- UTC (Universal Coordinated Time)
 - From time to time, we skip a solar second to stay in phase with the sun (30+ times since 1958)
 - UTC is broadcast by several sources (satellites...)

Accuracy of Computer Clocks

- Modern timer chips have a relative error of $1/100,000$ - 0.86 seconds a day
- To maintain synchronized clocks
 - Can use UTC source (time server) to obtain current notion of time
 - Use solutions without UTC.

Berkeley UNIX algorithm

- One daemon without UTC
- Periodically, this daemon polls and asks all the machines for their time
- The machines respond.
- The daemon computes an average time and then broadcasts this average time.

Decentralized Averaging Algorithm



- Each machine has a daemon without UTC
- Periodically, at fixed agreed-upon times, each machine broadcasts its local time.
- Each of them calculates the average time by averaging all the received local times.

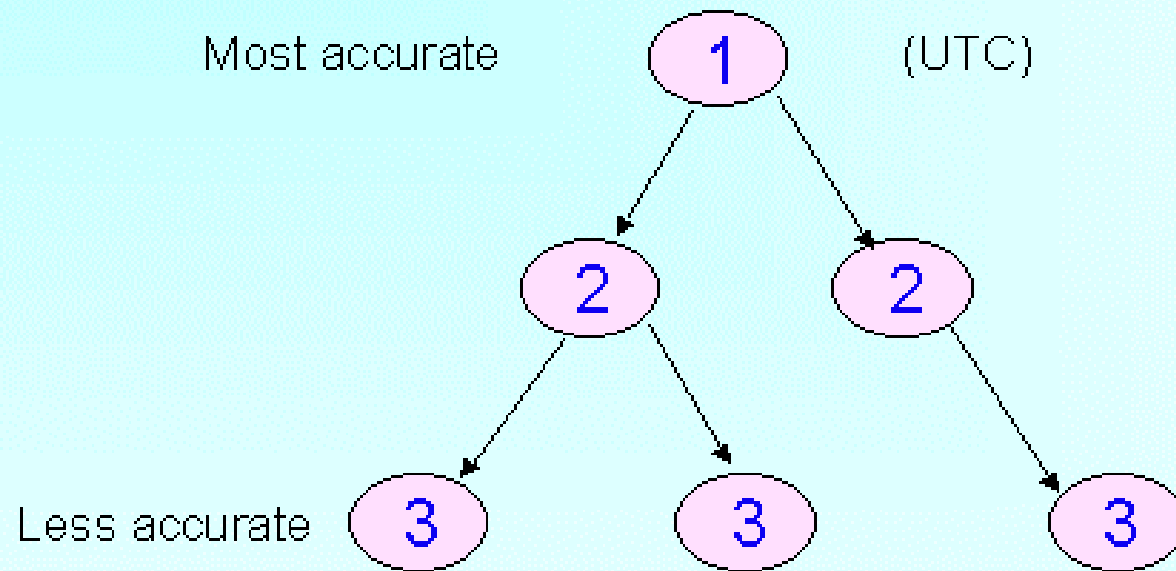
Clock Synchronization in DCE

- DCE's time model is actually in an interval
 - I.e. time in DCE is actually an interval
 - Comparing 2 times may yield 3 answers
 - $t1 < t2$
 - $t2 < t1$
 - not determined
 - Each machine is either a time server or a clerk
 - Periodically a clerk contacts all the time servers on its LAN
 - Based on their answers, it computes a new time and gradually converges to it.

The Network Time Protocol

- Enables clients across the Internet to be synchronized accurately to the UTC
 - Overcomes large and variable message delays
 - Statistical techniques for filtering can be applied
 - based on past behavior of server
 - Can survive lengthy losses of connectivity
 - Enables frequent synchronization
 - Provides protection against interference
 - Uses a hierarchy of servers located across the Internet (Primary servers connected to a UTC time source).

Hierarchy in NTP



Time Manager Operations

■ Logical Clocks

■ C.adjust(L,T)

- | adjust the local time displayed by clock C to T (can be gradually, immediate, per clock sync period)

■ C.read

- | returns the current value of clock C

■ Timers

- TP.set(T) - reset the timer to timeout in T units

■ Messages

- receive(m,l); broadcast(m); forward(m,l)