CS 230 - Distributed Systems

Lecture 1 - Introduction to Distributed Systems Mondays, Wednesdays 12:30-1:50p.m.

Prof. Nalini Venkatasubramanian nalini@ics.uci.edu

Course logistics and details

- **#** Course Web page
 - △ http://www.ics.uci.edu/~cs230
- # Lectures MW 12:30-1:50p.m
- # Must Read: Course Reading List
 - **区**Collection of Technical papers and reports by topic
- **#** Reference Books
 - **☑ Distributed Systems: Concepts & Design**, 4th ed. by Coulouris et al. ISBN: 0-321-26354-5.
 - **☑ Distributed Systems: Principles and Paradigms**, 2nd ed. by Tanenbaum & van Steen. ISBN: 0-132-39227-5.
 - **Systems**, 1st ed. by Kshemkalyani & Singhal. ISBN: 0-521-87634-6

Prerequisite Knowledge

- ** Necessary Operating Systems Concepts and Principles, basic computer system architecture
- # Highly Desirable Understanding of Computer Networks, Network Protocols
- ★ Necessary Basic programming skills in Java, C++,...

Course logistics and details

- **# Homeworks**
- **# Midterm Examination**
- **#Course Project**
 - - **☑**Project proposal due end of Week 2
 - **区**Survey of related research due end of Week 6
 - **区Final Project presentations/demos/reports − Finals** week

CompSci 230 Grading Policy

- #Homeworks 30% of final grade
 - **≥ 1** paper summary due every week from Week 2
- #Midterm 30% of final grade
 - **区**Tentatively in Week 7
- #Class Project 40% of the final grade
- #Final assignment of grades will be based on a curve.

Lecture Schedule

- - **☑Introduction Needs/Paradigms**
 - Basic Concepts and Terminology, Concurrency
 - **区Time and State in Distributed Systems**
 - **☑ Messaging/Communication in Distributed Systems**
 - Naming, Directory Services
- - **区**Communication
 - Remote Procedure Calls, Remote Method Invocation
 - **☑ Distributed Process and Resource Management**
 - Task Migration, Load Balancing, SOA
 - **☑ Distributed Process Synchronization**
 - Distributed Mutual Exclusion, Distributed Deadlocks
 - **☑** Distributed I/O and Storage Subsystems
 - Distributed FileSystems

Lecture Schedule

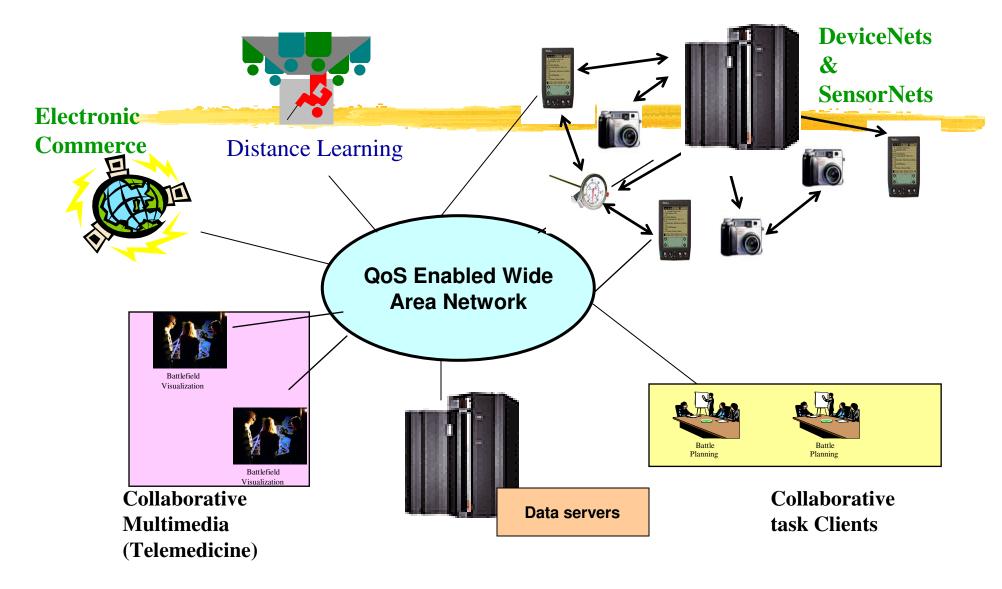
- ✓ Weeks 7,8,9: Non-functional "ilities" in distributed systems
 - **X** Reliability and Fault Tolerance
 - **区 Example 2 Example 3 Example 3 Example 4 Example 4 Example 4 Example 5 Example 5 Example 6 Example 7 Example 7**
 - **Scalability**
 - **区** Security and Privacy
- - **P2P, Grid and Cloud Computing**
 - **⋈** Mobile and Pervasive Systems/Applications

Introduction

Distributed Systems

- △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."

Next Generation Information Infrastructure



Requirements - Availability, Reliability, Quality-of-Service, Cost-effectiveness, Security

Characterizing Distributed Systems

Multiple Autonomous Computers

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

Interconnections

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

Shared State

- No shared memory
- systems cooperate to maintain shared state
- maintaining global invariants requires correct and coordinated operation of multiple computers.

Examples of Distributed Systems

- # Transactional applications Banking systems
- **# Manufacturing and process control**
- **#** Inventory systems
- # General purpose (university, office automation)
- # Communication email, IM, VoIP, social networks
- **#** Distributed information systems

 - Cloud Computing Infrastructures

Why Distributed Computing?

- **#Inherent distribution**
 - ☑Bridge customers, suppliers, and companies at different sites.
- **Speedup** improved performance
- #Fault tolerance
- ****** Resource Sharing
- **#**Scalability
- # Flexibility

Why are Distributed Systems Hard?

#Scale

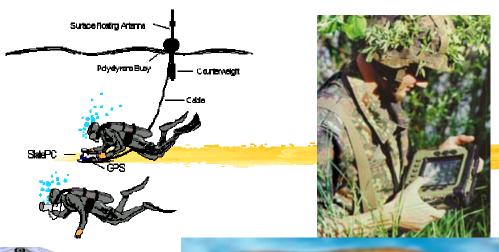
- numeric, geographic, administrative
- **#Loss** of control over parts of the system
- #Unreliability of message passing
 - Ounreliable communication, insecure communication, costly communication

Failure

- □ Parts of the system are down or inaccessible

Design goals of a distributed system

- - avoids centralization
- #Fault tolerance/availability
- **#**Transparency
 - △location, migration, replication, failure, concurrency



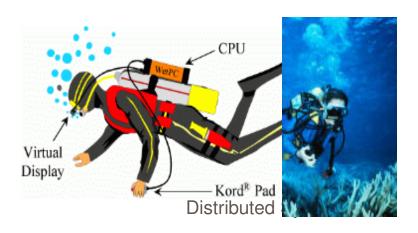


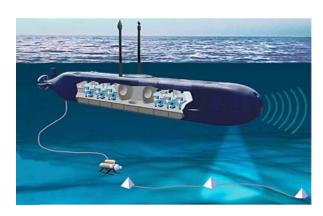






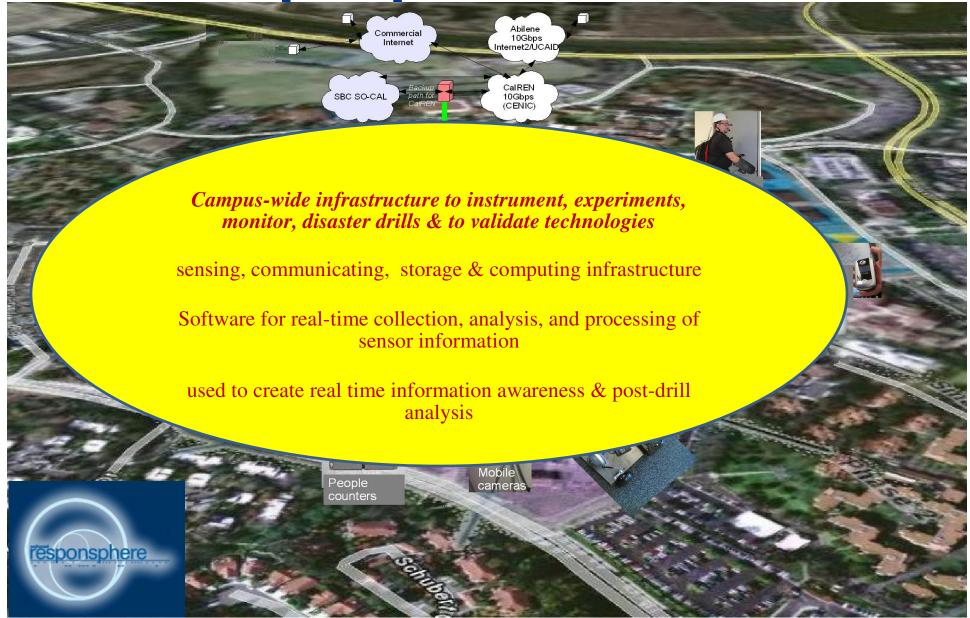






Pervasive Sensing and Computing

– UCI Responsible Pervasive Sensing and Computing



Classifying Distributed Systems

- **#**Based on degree of synchrony
 - **△**Synchronous
 - △ Asynchronous
- **#Based on communication medium**
- #Fault model

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

- □ 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

 - □ Isolated as if there were one serial database
 - □ Durable effects are not lost
- #General correctness of distributed computation
 - **△**Safety

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

- - **Processes communicate by sharing messages**
- □ Distributed Shared Memory (DSM)
 - **区**Communication through a virtual shared memory.

Message Passing

Basic communication primitives

- □ Receive message

Modes of communication

- **⊠**atomic action requiring the participation of the sender and receiver.
- 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 memory CPU₂ **CPU** Memory CPU1 Memory • •• CPU3 • •• CPU4

Distributed Shared 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 **CPU CPU CPU** Memory **CPU CPU CPU CPU** Memory Cache Cache Cache Cache

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
- - **™**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.
 - **区** 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

- - ☑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.
- - **☑** processor fails by transmitting only a subset of the messages it actually attempts to send.
- □ General Omission
 - **区**Receive and/or send omission

Other Distributed System issues

- **#Concurrency and Synchronization**
- **#Distributed Deadlocks**
- **X** Time in distributed systems
- **# Naming**
- ****** Replication
 - improve availability and performance
- **#**Migration
- **#**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:

 - △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

Management and Support

Network Management

Enterprise Systems: Perform enterprise activities

Application Systems: support enterprise systems

Distributed Computing Platform

- Application Support Services (OS, DB support, Directories, RPC)
- Communication Network Services (Network protocols, Physical devices)
- Hardware

Interoperability

Portability Integration

Managemen and Support

Management Network

Enterprise Systems:

- •Engineering systems Manufacturing
- Business systems
- Office systems

Application Systems:

User

<u>Interfaces</u>

Processing programs

Data files & <u>Databases</u>

Distributed Computing Platform

Application Support Services

C/S Support

Dist. Data Trans. Mgmt. Distributed OS

Common Network Services

• Network protocols & interconnectivity

OSI

protocols

TCP/IP

SNA

Interoperability Portability Integration

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 | | | ement |
|-------------------------------------|------------------------------|-----------|-------------------------|------------|
| | DCE Distributed File Service | | | |
| | DCE | DCE | Other Basic Services | Management |
| | Distributed | Directory | | |
| | Time Service | Service | | |
| | DCE Remote Procedure Calls | | | |
| DCE Threads Services | | | | |
| Operating System Transport Services | | | | |

Distributed Systems Middleware

- #Middleware is the software between the application programs and the operating System and base networking
- #Integration Fabric that knits together applications, devices, systems software, data
- # Middleware provides a comprehensive set of higher-level distributed computing capabilities and a set of interfaces to access the capabilities of the system.

The Evergrowing Middleware

Distributed Computing **Environment (DCE)**





LDAP

EAI

XQuery

XPath



WSDL

BPEL

JNDI

IOP

IIOP GIOP

Object Request Broker (ORB)



ZEN





Remote Procedure Call (RPC)





Extensible Markup Language (XML) Distributed Sys WebSphere. software



JMS

BEA Tuxedo®

RTCORBA

SOAP

IDL

ORBlite

Message Queuing (MSMQ)

opalORB

Remote Method Invocation

(RMI)

Encina/9000

BEA WebLogic®







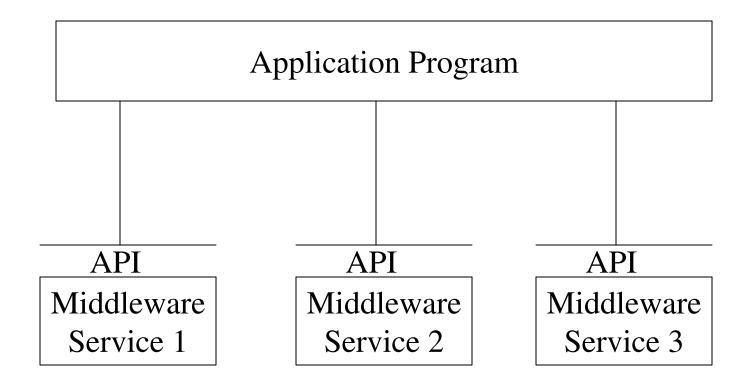
The Evergrowing Alphabet Soup

- ★ CORBA, OMG, CanCORBA, ORBIX, JavaORB, ORBLite, TAO, Zen, RTCORBA, FTCORBA, DCOM, POA, IDL, IOP, IIOP,
- ₩ ObjectBroker, Visibroker, Orbix, ObjectBus, ESBs
- # MOM TIBCO TIB/Rendezvous, BEA MessageQ, Microsoft MSMQ, ActiveWorks
- ₩ JVM, JINI, RMI, J2EE, EJB, J2ME, JDBC, JTA, JTS, JMS, JNDI,
- # Enterprise Middleware Technologies -- BEA WebLogic, IBM WebSphere, TivoliBeans
- # ENCINA, Tuxedo, CICS
- **XML**, XQuery,
- # SOAP, Web Services, WSDL, BPEL
- ₩

Distributed Systems Middleware

- □ Enables the modular interconnection of distributed software
 - **⊠**abstract over low level mechanisms used to implement resource management services.
- Computational Model
 - **区**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.
 - **区**Concurrent execution of multiple distributed systems policies.

Modularity in Middleware Services



Useful Middleware Services

- **△**State Capture Service
- **△**Event Service
- **△**Transaction Service
- **△** Fault Detection Service
- **△**Replication Service

Integration Frameworks Middleware

- #Integration frameworks are integration environments that are tailored to the needs of a specific application domain.
- **#** Examples

Distributed Object Computing

- ****Combining distributed computing with an object model.**
 - △Allows software reusability and a more abstract level of programming
 - □ The use of a broker like entity or bus that keeps track of processes, provides messaging between processes and other higher level services
 - **△**Examples
 - **XCORBA**
 - **⊠JINI, EJB, J2EE**
 - **E-SPEAK**
 - Note: DCE uses a procedure-oriented distributed systems model, not an object model.

Issues with Distributed Objects

- △ Abstraction
- **△**Latency
- **△**Synchronization
- **△**Complexity

Techniques for object distribution

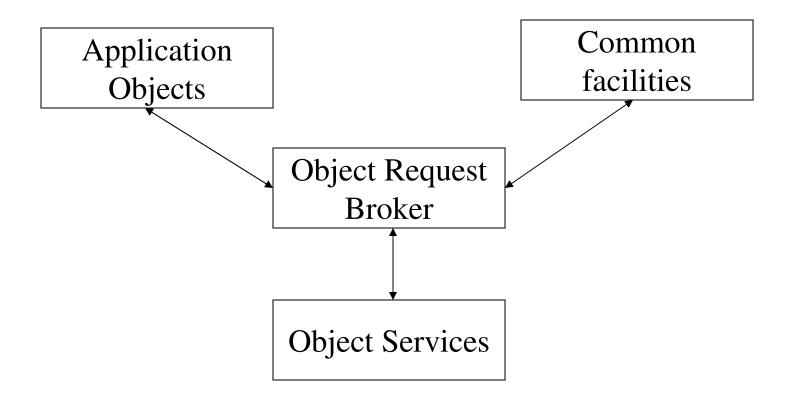
- - **⊠Object knows about network; Network data is** minimum
- △Argument/Return Passing
 - **区Like RPC. Network data** = args + return result + names
- **△**Shared Memory
 - **⊠based on DSM implementation**
 - Network Data = Data touched + synchronization info
 Distributed Systems

 45

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.

Virtual Time & Global States of Distributed Systems

- Asynchronous distributed systems consist of several *processes* without common memory which communicate (solely) via *messages* with unpredictable transmission delays
- # Global time & global state are hard to realize in distributed systems
 - - **Simulate synchronous** distributed system on a given asynchronous systems
 - **区** Simulate a global time
 - **区** Simulate a global state

Simulate Synchronous Distributed Systems

Synchronizers [Awerbuch 85]

- △ Simulate clock pulses in such a way that a message is only generated at a clock pulse and will be received before the next pulse
- □ Drawback
 - **☑Very high message overhead**

Clock Synchronization in Distributed Systems

#Clocks in a distributed system drift:

- □ Relative to each other
 - **区** Logical Clocks are clocks which are synchronized relative to each other.
- □ Relative to a real world clock
 - **区** 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
- **X** Notion of time is critical
- **#**"Happens Before" notion.
- ** "Happens Before" notion is not straightforward in distributed systems

Event Ordering

- ****Lamport defined the "happens before" (=>)**relation
 - \triangle If a and b are events in the same process, and a occurs before b, then a => b.
 - \triangle 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.

 \triangle 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

Logical Clocks

- ****** Monotonically increasing counter
- **%** No relation with real clock
- #Each process keeps its own logical clock Cp used to timestamp events

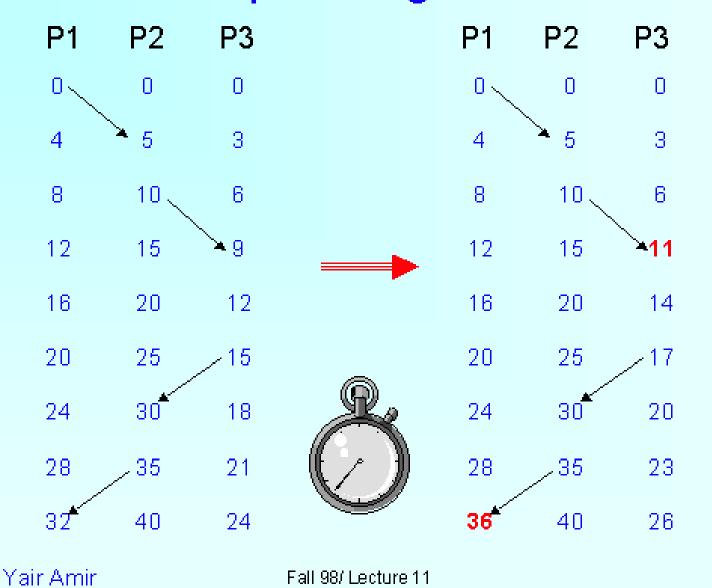
Causal Ordering and Logical Clocks

#Cp is incremented before each event.

$$\triangle$$
Cp = Cp + 1

- # When p sends a message m, it piggybacks a logical timestamp t = Cp.
- #When q receives (m,t) it computes:
 - \triangle Cq = max(Cq,t) before timestamping the message receipt event.
- #Results in a partial ordering of events.

Lamport Logical Clock



DISHIDULEU OYSLEHIS

6

Total Ordering

#Extending partial order to total order



#Global timestamps:

```
\triangle(Ta,Pa) < (Tb,Pb) iff \triangle(Ta < Tb) or ((Ta = Tb) and (Pa < Pb))
```

Problems with Total Ordering

- **#** A linearly ordered structure of time is not always adequate for distributed systems
 - captures dependence of events
 - □ loses independence of events artificially enforces an ordering for events that need not be ordered.
 - Mapping partial ordered events onto a linearly ordered set of integers it is losing information
 - Events which may happen simultaneously may get different timestamps as if they happen in some definite order.
- # A partially ordered system of *vectors* forming a *lattice* structure is a natural representation of time in a distributed system
- # It resembles Minkowski's relativistic space-time (see Special Theory of Relativity)

Physical Clocks

#How do we measure real time?

- - **Solar Day Transit of the sun**
 - **Solar Seconds Solar Day/(3600*24) Solar Seconds Solar Day/(3600*24)**
- □ Problem (1940) Rotation of the earth varies (gets slower)

Atomic Clocks

#1948

- - **≥9192631779** transitions = 1 mean solar second in 1948
- - **⋉**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

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

- □ Comparing 2 times may yield 3 answers

⊻t1 < t2

⋉t2 < t1

⊠not determined

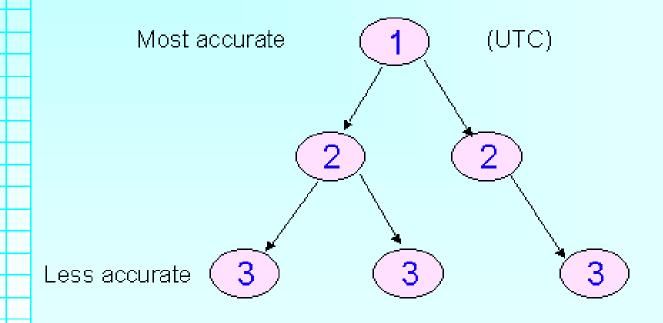
- Periodically a clerk contacts all the time servers on its LAN
- □ Based on their answers, it computes a new time and gradually converges to it.

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

 - □ Uses a hierarchy of servers located across the Internet (Primary servers connected to a UTC time source).

Hierarchy in NTP



Yair Amir

Fall 98/ Lecture 11

18

Time Manager Operations

```
#Logical Clocks
```

```
\triangleC.adjust(L,T)
```

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

Improvement Service In the current value of clock C

X Timers

Messages

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