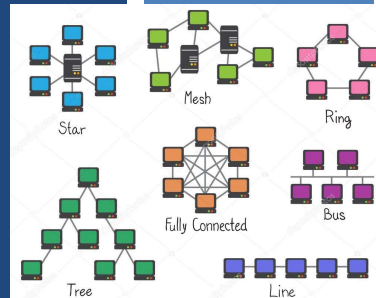


TCSS 558: APPLIED DISTRIBUTED COMPUTING

Chapter 4 – Communication Chapter 6 – Coordination

Wes J. Lloyd
School of Engineering
and Technology
University of Washington - Tacoma



OBJECTIVES

- Homework 2
- Active Reading Quiz
- Chapter 4 Communication
 - 4.4 Multicast communication
- Ch. 6 – Coordination
 - 6.1 Clock synchronization
 - 6.2 Logical clocks, Lamport clocks, Vector clocks
 - 6.3 Distributed mutual exclusion

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.2

FEEDBACK – 3/4

- How do you give MPI failure transparency?
- Failure transparency involves hiding from the user, the fact that the system (or some aspect of it) has failed
- Providing failure transparency requires a system to implement fault tolerance
- Here is an FAQ on fault tolerance in OpenMPI:
- <https://www.open-mpi.org/faq/?category=ft>
- A number of techniques for fault tolerance have been employed previously in OpenMPI, but are not deprecated
- OpenMPI is said to mimic the fault tolerance provided by the FT-MPI framework

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.3

FEEDBACK - 2

- What types of messages are usually sent with gossip spreading?
- Gossip, in the context of Ch. 4.4, refers to multicast communication (one to many) across unstructured peer-to-peer network
- These are ad hoc connections where the structure of the network is unknown
- Multicast messages could be anything
- Multicast often concerns data dissemination – spreading data to many peer nodes as quickly and efficiently as possible

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.4

CH. 4.4: MULTICAST COMMUNICATION

Multicast

one to many

X = subscriber

Apache ActiveMQ

L14.5

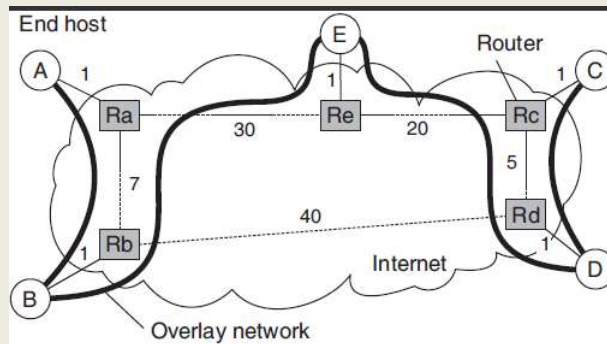
MULTICAST COMMUNICATION

- Sending data to multiple receivers
- Many **failed** proposals for network-level / transport-level protocols to support multicast communication
- **Problem:** How to set up communication paths for information dissemination?
- **Solutions:** require huge management effort, human intervention
- Focus shifted more recently to **peer-to-peer** networks
 - Structured overlay networks can be setup easily and provide efficient communication paths
 - Application-level multicasting techniques more successful
 - Gossip-based dissemination: unstructured p2p networks

March 6, 2019	TCSS558: Applied Distributed Computing [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma	L14.6
---------------	---	-------

NETWORK STRUCTURE

- **Overlay network**
 - Virtual network implemented on top of an actual physical network
- **Underlying network**
 - The actual physical network that implements the overlay



March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.7

FLOOD-BASED MULTICASTING

- **Broadcasting:** every node in overlay receives message
- **Key design issue:** minimize the use of intermediate nodes for which the message is not intended
- **Tree:** if only the leaf nodes are to receive the multicast message, many intermediate nodes are involved
- **Solution:** construct an overlay network for each multicast group
- **Flooding:** each node simply forwards a message to each of its neighbors, except to the message originator

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

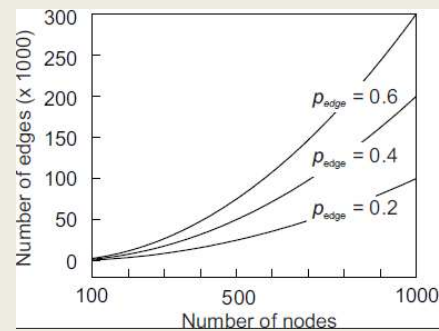
L14.8

RANDOM GRAPHS

- When no information on the structure of the overlay network
- Assume network can be represented as a **Random graph**
- Probability P_{edge} that two nodes are joined
- Overlay network will have: $\frac{1}{2} * P_{\text{edge}} * N * (N-1)$ edges

Random graphs allow us to assume some structure (# of nodes, # of edges) regarding the network by scaling the P_{edge} probability

Assumptions may help then to reason or rationalize about the network...



March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.9

PROBABILISTIC FLOODING



-*Washington state in winter?*
- When a node is flooding a message, concept is to enforce a probability of message spread (p_{flood})
- Throttles message flooding based on a probability
- Implementation needs to consider # of neighbors to achieve various p_{flood} scores
- With lower p_{flood} messages may not reach all nodes
- **USEFULNESS:** For random network with 10,000 nodes
- With $p_{\text{edge}} = 0.1$ and $p_{\text{flood}} = .01$
- Achieves 50-fold reduction in messages vs. full flooding

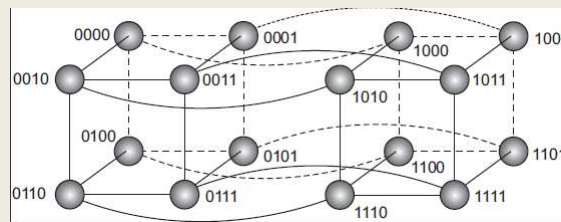
March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.10

MESSAGE FLOODING

- For deterministic topologies (such as hypercube), design of efficient flooding scheme is much simpler
- If the overlay network is structured, this gives us a deterministic topology
- Schlosser et al [2002] – offer simple and efficient **broadcasting scheme** that relies on keeping track of neighbors per dimension

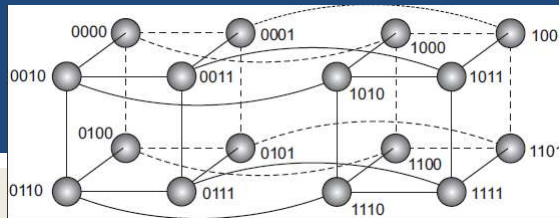


March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.11

MESSAGE FLOODING - 2



- **Hypercube Broadcast**
- N(1001) starts the network broadcast
- N(1001) neighbors {0001, 1000, 1011, 1101}
- N(1001) Sends message to all neighbors
- **Edge Labels (which bit is changed, 1st, 2nd, 3rd, 4th...)**
- Edge to 0001 – labeled 1 – change the 1st bit
- Edge to 1000 – labeled 4 – change the 4th bit
- Edge to 1011 – labeled 3 – change the 3rd bit
- Edge to 1101 – labeled 2 – change the 2nd bit
- **RULE: nodes only forward along edges with a higher dimension**
- Node 1101 receives message on edge labeled 2
- Broadcast msg is only forwarded on higher dimension edges

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.12

MESSAGE FLOODING - 3

- **Hypercube:** forward msg along edges with higher dimension
- Node(1101)-neighbors {0101,1100,1001,1111}
- Node (1101) - incoming broadcast edge = 2
- Label Edges:
 - Edge to 0101 - labeled 1 - change the 1st bit
 - Edge to 1100 - labeled 4 - change the 4th bit *<FORWARD>*
 - Edge to 1001 - labeled 2 - change the 2nd bit
 - Edge to 1111 - labeled 3 - change the 3rd bit *<FORWARD>*
- N(1101) broadcast - forward only to N(1100) and N(1111)
- (1100) and (1111) are the higher dimension edges
- Broadcast requires just: N-1 messages, where nodes $N=2^n$,
n=dimensions of hypercube

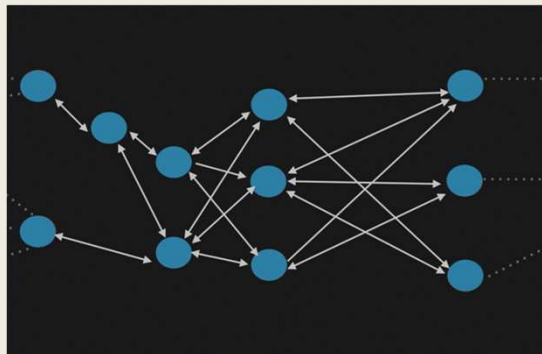
March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.13

GOSSIP BASED DATA DISSEMINATION

- When structured peer-to-peer topologies are not available
- Gossip based approaches support multicast communication over unstructured peer-to-peer networks
- General approach is to leverage how gossip spreads across a group
- This is also called "epidemic behavior"...
- Data updates for a specific item begin at a specific node



March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.14

INFORMATION DISSEMINATION

- **Epidemic algorithms**: algorithms for large-scale distributed systems that spread information
- Goal: “infect” all nodes with new information as fast as possible
- **Infected**: node with data that can spread to other nodes
- **Susceptible**: node without data
- **Removed**: node with data that is unable to spread data

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.15

ANTI ENTROPY DISSEMINATION MODEL

- **Anti-entropy**: Propagation model where node P picks node Q at random and exchanges message updates
- Akin to random walk
- **PUSH**: P only **pushes** its own updates to Q
- **PULL**: P only **pulls** in new updates from Q
- **TWO-WAY**: P and Q send updates to each other (i.e. a push-pull approach)
- Push only: hard to propagate updates to last few hidden susceptible nodes
- Pull: better because susceptible nodes can pull updates from infected nodes
- Push-pull is better still

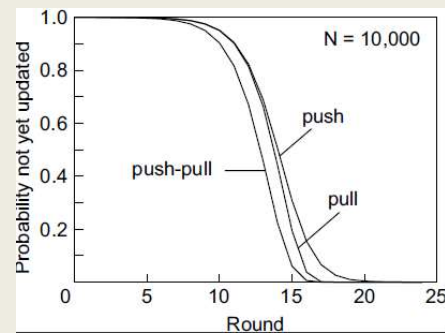
March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.16

ANTI ENTROPY EFFECTIVENESS

- **Round:** span of time during which every node takes initiative to exchange updates with a randomly chosen node
- The number of rounds to propagate a single update to all nodes requires $O(\log(N))$, where N =number of nodes
- Let p_i denote probability that node P has not received msg m after the i^{th} round.
- For pull, push, and push-pull based approaches:



March 6, 2019

TCS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.17

RUMOR SPREADING

- Variant of epidemic protocols
- Provides an approach to “stop” message spreading
- Mimics “gossiping” in real life
- **Rumor spreading:**
- **Node P** receives new data item **X**
- Contacts an arbitrary **node Q** to push update
- **Node Q** reports already receiving **item X** from another node
- **Node P** may loose interest in spreading the rumor with probability = p_{stop} , let's say 20% . . . (or 0.20)

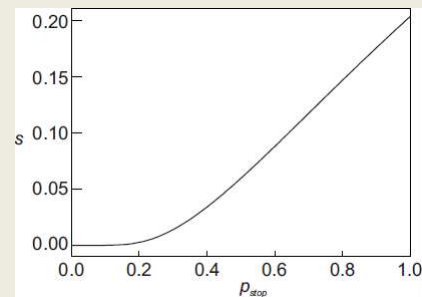
March 6, 2019

TCS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.18

RUMOR SPREADING - 2

- Does not guarantee all nodes will be updated
- The fraction of nodes s , that remain susceptible grows relative to the probability that node P stops propagating when finding a node already having the message
- Fraction of nodes not updated remains < 0.20 with high p_{stop}
- Susceptible nodes (s) vs. probability of stopping \rightarrow



March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.19

DIRECTIONAL GOSSIPING

- Taking network topology into account can help
- When gossiping, nodes connected to only a few other nodes are more likely to be contacted
- **Epidemic protocols assume:**
 - For gossiping, nodes are randomly selected
 - One node, can randomly select any other node in the network
 - Complete set of nodes is known to each member

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.20

REMOVING DATA

- Gossiping is good for spreading data
- But how can data be removed from the system?
- Idea is to issue ***“death certificates”***
- Act like data records, which are spread like data
- When death certificate is received, data is deleted
- Certificate is held to prevent data element from reinitializing from gossip from other nodes
- Death certificates time-out after expected time required for data element to clear out of entire system
- A few nodes maintain death certificates forever

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.21

DEATH CERTIFICATE EXAMPLE

- For example:
- **Node P** keeps death certificates forever
- **Item X** is removed from the system
- **Node P** receives an update request for **Item X**, but also holds the death certificate for **Item X**
- **Node P** will recirculate the death certificate across the network for **Item X**

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.22

CHAPTER 6 - COORDINATION

- 6.1 Clock Synchronization
 - Physical clocks
 - Clock synchronization algorithms
- 6.2 Logical clocks
 - Lamport clocks
 - Vector clocks
- 6.3 Mutual exclusion
- 6.4 Election algorithms
- 6.6 Distributed event matching (*light*)
- 6.7 Gossip-based coordination (*light*)

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.23

CHAPTER 6 - COORDINATION

- How can processes synchronize and coordinate data?
- Process synchronization
 - Coordinate cooperation to grant individual processes temporary access to shared resources (e.g. a file)
- Data synchronization
 - Ensure two sets of data are the same (data replication)
- Coordination
 - Goal is to manage interactions and dependencies between activities in the distributed system
 - Encapsulates synchronization

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.24

COORDINATION - 2

- Synchronization challenges begin with time:
 - How can we synchronize computers, so they all agree on the time?
 - How do we measure and coordinate when things happen?
- Fortunately, for synchronization in distributed systems, it is often sufficient to only agree on a relative ordering of events
 - E.g. not actual time

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.25

COORDINATION - 3

- Groups of processes often appoint a coordinator
- Election algorithms can help elect a leader
- Synchronizing access to a shared resource is achieved with distributed mutual exclusion algorithms
- Also in chapter 6:
 - Matching subscriptions to publications in publish-subscribe systems
 - Gossip-based coordinate problems:
 - Aggregation
 - Peer sampling
 - Overlay construction

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.26



CH. 6.1: CLOCK SYNCHRONIZATION

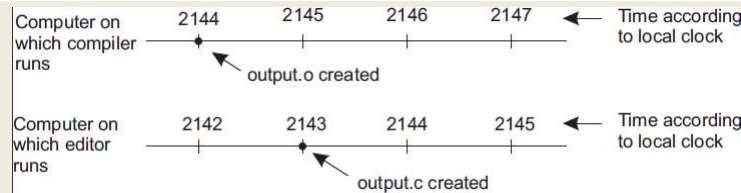
L14.27

CLOCK SYNCHORNIZATION

- **Example:**
- “make” is used to compile source files into binary object and executable files
- As an optimization, make only compiles files when the “last modified time” of source files is more recent that object and executables
- Consider if files are on a shared disk of a distributed system where there is no agreement on time
- Consider if the program has 1,000 source files

March 6, 2019	TCSS558: Applied Distributed Computing [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma	L14.28
---------------	---	--------

TIME SYNCHRONIZATION PROBLEM FOR DISTRIBUTED SYSTEMS



- Updates from different machines, may have clocks set to different times
- Make becomes confused with which files to recompile

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.29

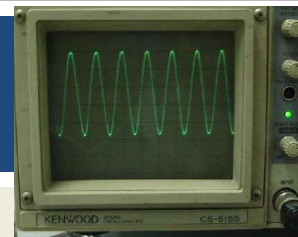


PHYSICAL CLOCKS

- **Computer timers:** precisely machined quartz crystals
- When under tension, they oscillate at a well defined frequency
- In analog electronics/communications crystals once used to set the frequency of two-way radio transceivers for
- Today, crystals are associated with a counter and holding register on a digital computer.

1960s ERA radio crystal →

- Each oscillation decrements a counter by one
- When counter gets to zero, an interrupt fires
- Can program timer to generate interrupt, let's say 60 times a second, or another frequency to track time

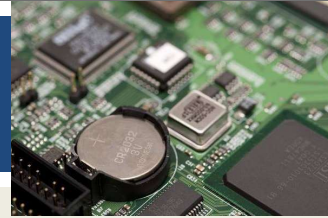


March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.30

COMPUTER CLOCKS



- Digital clock on computer sets base time
- Crystal clock tracks forward progress of time
 - Translation of wave “ticks” to clock pulses
- CMOS battery on motherboard maintains clock on power loss
- **Clock skew**: physical clock crystals are not exactly the same
- Some run at slightly different rates
- Time differences accumulate as clocks drift forward or backward slightly
- In an automobile, where there is no clock synchronization, clock skew may become noticeable over months, years

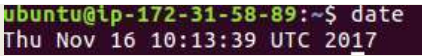


March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.31

UNIVERSAL COORDINATED TIME

- **Universal Coordinated Time (UTC)** 
- Worldwide standard for time keeping
- Equivalent to Greenwich Mean Time (United Kingdom)
- 40 shortwave radio stations around the world broadcast a short pulse at the start of each second (WWV)
- World wide “atomic” clocks powered by constant transitions of the non-radioactive caesium-133 atom
 - 9,162,631,770 transitions per second
- Computers track time using UTC as a base
 - Avoid thinking in local time, which can lead to coordination issues
 - Operating systems may translate to show local time

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.32

COMPUTING: CLOCK CHALLENGES

- How do we synchronize computer clocks with real-world clocks?
- How do we synchronize computer clocks with each other?

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.33

CLOCK SYNCHRONIZATION

- **UTC services:** use radio and satellite signals to provide time accuracy to 50ns
- **Time servers:** Server computers with UTC receivers that provide accurate time
- **Precision (π):** how close together a set of clocks may be
- **Accuracy:** how correct to actual time clocks may be
- **Internal synchronization:** Sync local computer clocks
- **External synchronization:** Sync to UTC clocks
- **Clock drift:** clocks on different machines gradually become out of sync due to crystal imperfections, temperature differences, etc.
- **Clock drift rate:** typical is 31.5s per year
- **Maximum clock drift rate (ρ):** clock specifications include one

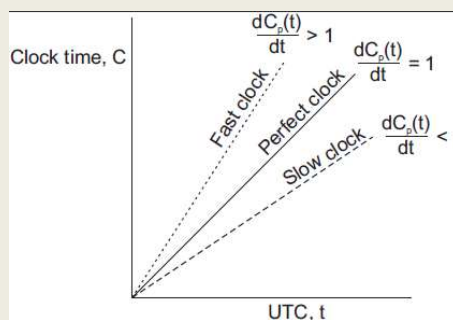
March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.34

CLOCK SYNCHRONIZATION - 2

- If two clocks drift from UTC in opposite directions, after time Δt after synchronization, they may be 2ρ apart.
- Clocks must be resynchronized every $\pi/2\rho$ seconds
- Network time protocol
- Provide coordination of time for servers
- Leverage distributed network of time servers



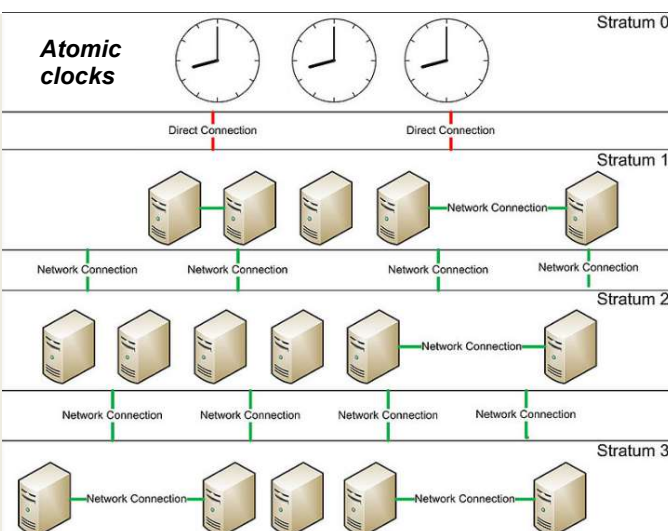
March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.35

NETWORK TIME PROTOCOL

- Servers organized into stratum
- Stratum-1 servers have UTC receivers and are sync'd with atomic clocks
- Servers connect with closest NTP server for time synchronization
- Servers assume role as NTP server at stratum+1



March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

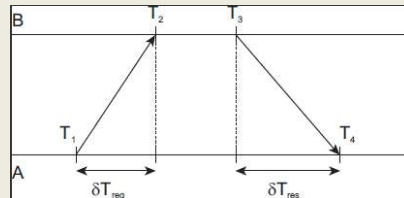
L14.36

NTP - 2

- Must estimate network delays when synchronizing with remote UTC receiver clocks / time servers

Time server B

Client A



- A sends message to B, with timestamp T1
- B records time of receipt T2 (from local clock)
- B returns response with send time T3, and receipt time T2
- A records arrival of T4

- Assuming propagation delay of A→B→A is the same

- Estimate propagation delay:

$$\theta = T_3 + \frac{(T_2 - T_1) + (T_4 - T_3)}{2} - T_4 = \frac{(T_2 - T_1) + (T_3 - T_4)}{2}$$

- Add delay to time

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.37

NTP - 3

- Cannot set clocks backwards (recall “make” file example)
- Instead, temporarily slow the progress of time to allow fast clock to align with actual time
- Change rate of clock interrupt routine
- Slow progress of time until synchronized
- NTP accuracy is within 1-50ms
- In Ubuntu Linux, to quickly synchronize time:
`$apt install ntp ntpdate`
- Specify local timeservers in /etc/ntp.conf
`server time.u.washington.edu iburst`
`server bigben.cac.washington.edu iburst`
- Shutdown service (sudo service ntp stop)
- Run ntpdate: (sudo ntpdate time.u.washington.edu)
- Startup service (sudo service ntp start)

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.38

BERKELEY ALGORITHM

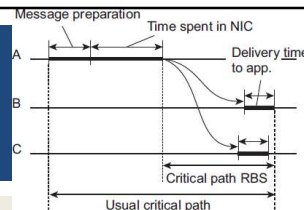
- Berkeley time daemon server actively polls network to determine average time across servers
- Suitable when no machine has a UTC receiver
- Time daemon instructs servers how much to adjust clocks to achieve precision
- Accuracy can not be guaranteed
- Berkeley is an internal clock synchronization algorithm

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.39

CLOCK SYNCHRONIZATION IN WIRELESS NETWORKS



- Sensor networks bring unique challenges for clock synchronization
 - **Address resource constraints:** limited power, multihop routing slow
- **Reference broadcast synchronization (RBS)**
- Provides precision of time, not accuracy as in Berkeley
- No UTC clock available
- RBS sender broadcasts a reference message to allow receivers to adjust clocks
- No multi-hop routing
- Time to propagate a signal to nodes is roughly constant
- Message propagation time does not consider time spent waiting in NIC for message to send
 - Wireless network resource contention may force wait before message even can be sent

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.40

REFERENCE BROADCAST SYNCHRONIZATION (RBS)

- Node broadcasts reference message m
- Each node p records time $T_{p,m}$ when m is received
- $T_{p,m}$ is read from node p 's clock
- Two nodes p and q can exchange delivery times to estimate mutual relative offset
- Then calculate relative average offset for the network:

$$\text{Offset}[p, q] = \frac{\sum_{k=1}^M (T_{p,k} - T_{q,k})}{M}$$

- Where M is the total number of reference messages sent
- Nodes can simply store offsets instead of frequently synchronizing clocks to save energy

March 6, 2019

TCCS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.41

REFERENCE BROADCAST SYNCHRONIZATION (RBS) - 2

- Cloud skew: over time clocks drift apart
- Averages become less precise
- Elson et al. propose using standard linear regression to predict offsets, rather than calculating them
- IDEA: Use node's history of message times in a simple linear regression to continuously refine a formula with coefficients to predict time offsets:

$$\text{Offset}[p, q](t) = \alpha t + \beta$$

March 6, 2019

TCCS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.42

CH. 6.2: LOGICAL CLOCKS

L14.43

LOGICAL CLOCKS

- In distributed systems, synchronizing to actual time may not be required...
- It may be sufficient for every node to simply agree on a current time (e.g. logical)
- **Logical clocks** provide a mechanism for capturing chronological and causal relationships in a distributed system
- Think **counters** . . .
- Leslie Lamport [1978] seminal paper showed that absolute clock synchronization often is not required
- Processes simply need to agree on the order in which events occur

March 6, 2019	TCSS558: Applied Distributed Computing [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma	L14.44
---------------	---	--------

LOGICAL CLOCKS - 2

- Happens-before relation
- $A \rightarrow B$: **Event A**, happens before **event B**...
- All processes must agree that **event A** occurs first
- Then afterward, **event B**
- Actual time not important. . .

- If **event A** is the event of proc P1 sending a msg to a proc P2, and **event B** is the event of proc P2 receiving the msg, then $A \rightarrow B$ is also true. . .
- The assumption here is that message delivery takes time
- Happens before is a transitive relation:
- $A \rightarrow B, B \rightarrow C$, therefore $A \rightarrow C$

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.45

LOGICAL CLOCKS - 3

- If two events, say event X and event Y do not exchange messages, not even via third parties, then the sequence of $X \rightarrow Y$ vs. $Y \rightarrow X$ can not be determined!!
- Within the system, these events appear concurrent
- **Concurrent**: nothing can be said about when the events happened, or which event occurred first
- Clock time, C, must always go forward (increasing), never backward (decreasing)
- Corrections to time can be made by adding a positive value, but never by subtracting one

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.46

LOGICAL CLOCKS - 4

- Three processes each with local clocks
- Lamport's algorithm corrects their values

P ₁	P ₂	P ₃
0	0	0
6	8	10
12	16	20
18	24	30
24	32	40
30	40	50
36	48	60
42	56	70
48	64	80
54	72	90
60	80	100

P ₁	P ₂	P ₃
0	0	0
6	8	10
12	16	20
18	24	30
24	32	40
30	40	50
36	48	60
42	61	70
48	69	80
70	77	90
76	85	100

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.47

LOGICAL CLOCKS

- Events:
 - 6: P1 send m1 to P2
 - 16: P2 receives m1
 - 24: P2 sends m2 to P3
 - 40: P3 receives m2
 - 60: P3 sends m3 to P2
 - 56: P2 receives m3
 - 56: P2 clock reset=61
 - 64: P2 sends m4 to P1
 - 54: P1 receives m4
 - 70: P1 clock reset=70

P ₁	P ₂	P ₃
0	0	0
6	8	10
12	16	20
18	24	30
24	32	40
30	40	50
36	48	60
42	56	70
48	64	80
54	72	90
60	80	100

P ₁	P ₂	P ₃
0	0	0
6	8	10
12	16	20
18	24	30
24	32	40
30	40	50
36	48	60
42	61	70
48	69	80
70	77	90
76	85	100

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.48

LAMPORT LOGICAL CLOCKS - IMPLEMENTATION

- Negative values not possible
 - When a message is received, and the local clock is before the timestamp when then message was sent, the local clock is updated to $\text{message_sent_time} + 1$
1. Clock is incremented before an event: sending a message, receiving a message, some other internal event
 P_i increments C_i : $C_i \leftarrow C_i + 1$
 2. When P_i send msg m to P_j , m 's timestamp is set to C_i
 3. When P_j receives msg m , P_j adjusts its local clock
 $C_j \leftarrow \max\{C_j, \text{ts}(m)\}$
 4. Ties broken by considering Proc ID: $i < j$; $\langle 40, i \rangle < \langle 40, j \rangle$

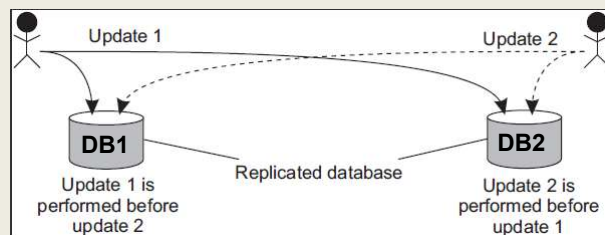
March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.49

TOTAL-ORDERED MULTICASTING

- Consider concurrent updates to a replicated database
- Communication latency between DB1 and DB2 is 250ms



- **Initial Account balance: \$1,000**
- **Update #1: Deposit \$100**
- **Update #2: Add 1% Interest**
- **Total Ordered Multicasting needed**

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

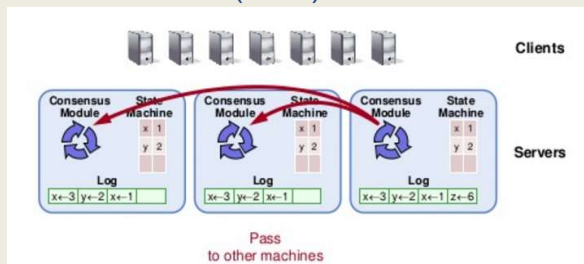
L14.50



- Each message timestamped with local logical clock of sender
- Multicast message is conceptually sent to the sender
- Assumptions:
 - Messages from same sender received in order they were sent
 - No messages are lost
- When messages arrive they are placed in local queue ordered by timestamp
- Receiver multicasts acknowledgement of message receipt to other processes
 - Time stamp of message receipt is lower the acknowledgement
- This process replicates queues across sites
- Process delivers messages to application only when message at the head of the queue has been acknowledged by every process in the system

TOTAL-ORDERED MULTICASTING - 3

- Can be used to provide replicated state machines (RSMs)
- Concept is to replicate event queues at each node
- (1) **Using logical clocks** and (2) **exchanging acknowledgement messages**, allows for events to be “**totally**” ordered in replicated event queues
- Events can be applied “**in order**” to each (distributed) replicated state machine (RSM)



March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.53

VECTOR CLOCKS

- Lamport clocks don't help to determine causal ordering of messages
- Vector clocks capture causal histories and can be used as an alternative
- What is causality?

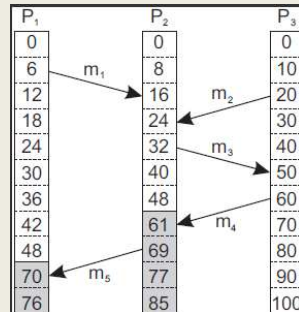
March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.54

WHAT IS CAUSALITY?

- Consider the messages:



- P₂ receives m₁, and subsequently sends m₃
- Causality:** Sending m₃ may depend on what's contained in m₁
- P₂ receives m₂, receiving m₂ is **not** related to receiving m₁
- Is sending m₃ causally dependent on receiving m₂?**

March 6, 2019

TCCS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.55

VECTOR CLOCKS

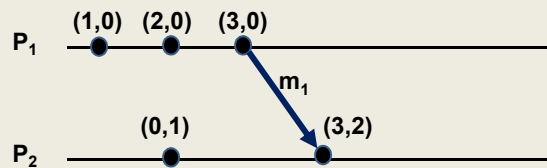
- Vector clocks keep track of **causal history**
- If two local events happened at process P, then the causal history H(p₂) of event p₂ is {p₁, p₂}
- P sends messages to Q (event p₃)
- Q previously performed event q₁
- Q records arrival of message as q₂
- Causal histories merged at Q H(q₂) = {p₁, p₂, p₃, q₁, q₂}
- Fortunately, can simply store history of last event, as a vector clock → H(q₂) = (3, 2)
- Each entry corresponds to the last event at the process

March 6, 2019

TCCS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.56

VECTOR CLOCKS - 2



- Each process maintains a vector clock which
 - Captures number of events at the local process (e.g. logical clock)
 - Captures number of events at all other processes
- Causality is captured by:
 - For each event at P_i , the vector clock (VC_i) is incremented
 - The msg is timestamped with VC_i ; and sending the msg is recorded as a new event at P_i
 - P_j adjusts its VC_j choosing the max of: the message timestamp –or– the local vector clock (VC_j)

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.57

VECTOR CLOCKS - 3

- P_j knows the # of events at P_i based on the timestamps of the received message
- P_j learns how many events have occurred at other processes based on timestamps in the vector
- These events *“may be causally dependent”*
- In other words: they may have been necessary for the message(s) to be sent...

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.58

VECTOR CLOCKS EXAMPLE

Local clock is underlined

CAUSALITY

Diagram illustrating a causality scenario in a distributed system with three processes (P_1 , P_2 , P_3). The diagram shows the execution of messages m_1 , m_2 , m_3 , and m_4 with their respective vector timestamps. The local clock is underlined in the timestamp components.

$ts(m_2)$	$ts(m_4)$	$ts(m_2) < ts(m_4)$	$ts(m_2) > ts(m_4)$	Conclusion
$(2, \underline{1}, 0)$	$(4, \underline{3}, 0)$	Yes	No	m_2 may causally precede m_4

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.59

VECTOR CLOCKS EXAMPLE - 2

Diagram illustrating a potential conflict scenario in a distributed system with three processes (P_1 , P_2 , P_3). The diagram shows the execution of messages m_1 , m_2 , m_3 , and m_4 with their respective vector timestamps. The local clock is underlined in the timestamp components.

$ts(m_2)$	$ts(m_4)$	$ts(m_2) < ts(m_4)$	$ts(m_2) > ts(m_4)$	Conclusion
$(4, \underline{1}, 0)$	$(2, \underline{3}, 0)$	No	No	m_2 and m_4 may conflict

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.60

- P_3 can't determine if m_4 may be causally dependent on m_2
- Is m_4 causally dependent on m_3 ?

VECTOR CLOCKS - 4

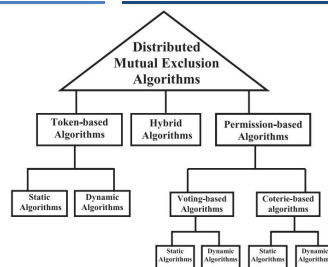
- **Disclaimer:**
- Without knowing actual information contained in messages, it is not possible to state with certainty that there is a causal relationship or perhaps a conflict
- Vector clocks can help us suggest possible causality
- We never know for sure...

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.61

CH. 6.3: DISTRIBUTED MUTUAL EXCLUSION



L14.62

DISTRIBUTED MUTUAL EXCLUSION

- Coordinating access among distributed processes to a shared resource requires **Distributed Mutual Exclusion**
- Token-based algorithms:
- Mutual exclusion by passing a “token” between nodes
- Nodes often organized in ring
- Only one token, holder has access to shared resource
- Avoids starvation: *everyone gets a chance to obtain lock*
- Avoids deadlock: easy to avoid

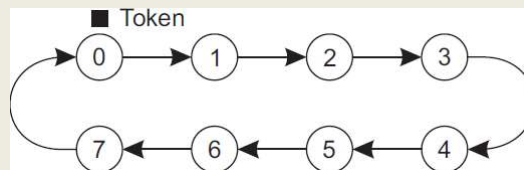
March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.63

TOKEN-RING ALGORITHM

- Construct overlay network
- Establish logical ring among nodes



- Single token circulated around the nodes of the network
- Node having token can access shared resource
- If no node accesses resource, token is constantly circulated around ring

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.64

TOKEN-RING CHALLENGES

1. If token is lost, token must be regenerated
 - **Problem:** may accidentally circulate multiple tokens
2. Hard to determine if token is lost
 - What is the difference between token being lost and a node holding the token for a long time?
3. When node crashes, circular network route is broken
 - Dead nodes can be detected by adding a receipt message for when the token passes from node-to-node
 - When no receipt is received, node assumed dead
 - Dead process can be “jumped” in the ring

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.65

DISTRIBUTED MUTUAL EXCLUSION - 2

- **Permission-based algorithms**
- Processes must require permission from other processes before first acquiring access to the resource
- **Centralized algorithm**
- Elect a single leader node to coordinate access to shared resource(s)
- Manage mutual exclusion on a distributed system similar to how it mutual exclusion is managed for a single system
- Nodes must all interact with leader to obtain “*the lock*”

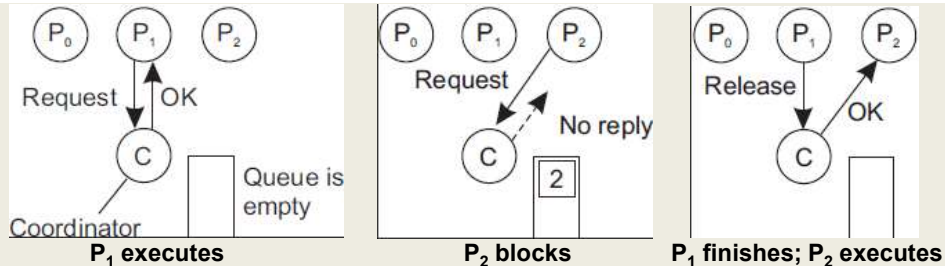
March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.66

CENTRALIZED MUTUAL EXCLUSION

Permission granted from coordinator \vee No response from coordinator



- When resource not available, coordinator can block the requesting process, or respond with a reject message
- P_2 must poll the coordinator if it responds with reject otherwise can wait if simply blocked
- Requests granted permission fairly using FIFO queue
- Just three messages: (request, grant, release)

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.67

CENTRALIZED MUTUAL EXCLUSION - 2

- Issues
 - Coordinator is a single point of failure
 - Processes can't distinguish dead coordinator from "permission denied"
 - No difference between CRASH and Block (*for a long time*)
 - Large systems, coordinator becomes performance bottleneck
 - Scalability: Performance does not scale
- Benefits
 - Simplicity:
Easy to implement compared to distributed alternatives

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
 School of Engineering and Technology, University of Washington - Tacoma

L14.68

DISTRIBUTED ALGORITHM

- Ricart and Agrawala [1981], use total ordering of all events
 - Leverages Lamport logical clocks
- Package up resource request message (AKA Lock Request)
- Send to all nodes
- Include:
 - Name of resource
 - Process number
 - Current (logical) time
- Assume messages are sent reliably
 - No messages are lost

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.69

DISTRIBUTED ALGORITHM - 2

- When each node receives a request message they will:
 1. Say OK (*if the node doesn't need the resource*)
 2. Make no reply, queue request (*node is using the resource*)
 3. Perform a timestamp comparison (*if node is waiting to access the resource*), then:
 1. Send OK if requester has lower logical clock value
 2. Make no reply if requester has higher logical clock value
- Nodes sit back and wait for all nodes to grant permission
- Requirement: every node must know the entire membership list of the distributed system

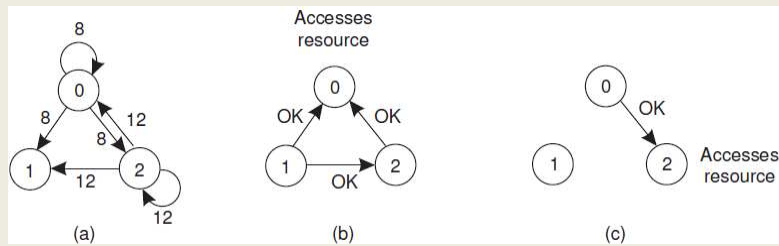
March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.70

DISTRIBUTED ALGORITHM - 3

- If Node 0 and Node 2 simultaneously request access
- Node 0's time stamp is lower (8) than Node 2 (12)
- Node 1 and Node 2 grant Node 0 access
- Notice that Node 1 also grants Node 2 permission



- In case of conflict, lowest timestamp wins!

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.71

CHALLENGES WITH DISTRIBUTED ALGORITHM

- **Problem:** Algorithm has N points of failure !
- Where N = Number of Nodes in the system
- **Problem:** When node is accessing the resource, it does not respond
 - Lack of response can be confused with **failure**
 - **Solution:** When node receives request for resource it is accessing, always send a reply either granting or denying permission (ACK)
 - Enables requester to determine when nodes have died

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.72

CHALLENGES WITH DISTRIBUTED ALGORITHM - 2

- **Problem:** Multicast communication required –or- each node must maintain full group membership
 - Track nodes entering, leaving, crashing...
- **Problem:** Every process is involved in reaching an agreement to grant access to a shared resource
 - This approach may not scale on resource-constrained systems
- **Solution:** Can relax total agreement requirement and proceed when a simple majority of nodes grant permission
 - Presumably any one node locking the resource prevents agreement
- Distributed algorithm for mutual exclusion works best for:
 - Small groups of processes
 - When memberships rarely change

March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.73

QUESTIONS




March 6, 2019

TCSS558: Applied Distributed Computing [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma

L14.74

EXTRA SLIDES



75