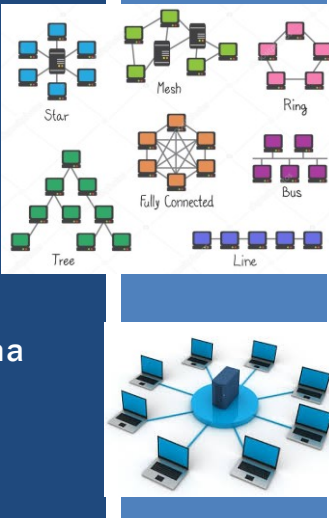


# TCSS 558: APPLIED DISTRIBUTED COMPUTING

## Chapter 6 – Coordination - II

Wes J. Lloyd  
School of Engineering  
& Technology (SET)  
University of Washington - Tacoma



The slide features a dark blue background with white text. On the right side, there is a collage of network topology diagrams. The diagrams are labeled: Star (a central node connected to multiple peripheral nodes), Mesh (a grid of interconnected nodes), Ring (nodes connected in a closed loop), Tree (a hierarchical structure of nodes), Fully Connected (every node connected to every other node), Bus (all nodes connected to a single central line), and Line (nodes connected in a straight line). Below these diagrams is a 3D illustration of a central server tower connected to several laptops arranged in a circle.

1

## OBJECTIVES – 2/28

- **Questions from 2/23**
- Assignment 2: Replicated Key Value Store
- Chapter 4.4 - Review / Finish
- Chapter 6: Coordination
  - Chapter 6.1: Clock Synchronization
  - Chapter 6.2: Logical Clocks  
Vector Clocks
- Class Activity – Total Ordered Multicasting (Thursday)
  - Chapter 6.3: Distributed Mutual Exclusion

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.2
-------------------	---	-------

2

## ONLINE DAILY FEEDBACK SURVEY

- Daily Feedback Quiz in Canvas – Available After Each Class
- Extra credit available for completing surveys **ON TIME**
- Tuesday surveys: due by ~ Wed @ 10p
- Thursday surveys: due ~ Mon @ 10p

TCSS 558 A > Assignments

Winter 2021

Home

Announcements

**Assignments**

Zoom

Chat

Search for Assignment

Upcoming Assignments

**TCSS 558 - Online Daily Feedback Survey - 1/5**  
Not available until Jan 5 at 1:30pm | Due Jan 6 at 10pm | -/1 pts

February 28, 2023    TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma    L16.3

3

### TCSS 558 - Online Daily Feedback Survey - 1/5

Due Jan 6 at 10pm    Points 1    Questions 4  
Available Jan 5 at 1:30pm - Jan 6 at 11:59pm 1 day    Time Limit None

Question 1    0.5 pts

On a scale of 1 to 10, please classify your perspective on material covered in today's class:

1	2	3	4	5	6	7	8	9	10
Mostly Review To Me				Equal New and Review					Mostly New to Me

Question 2    0.5 pts

Please rate the pace of today's class:

1	2	3	4	5	6	7	8	9	10
Slow				Just Right					Fast

February 28, 2023    TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma    L16.4

4

## MATERIAL / PACE

- Please classify your perspective on material covered in today's class (31 respondents):
  - 1-mostly review, 5-equal new/review, 10-mostly new
  - **Average - 6.57** (↑ - *previous 6.42*)
- Please rate the pace of today's class:
  - 1-slow, 5-just right, 10-fast
  - **Average - 5.94** (↑ - *previous 5.81*)

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.5

5

## TENURE TRACK FACULTY CANDIDATE RESEARCH SEMINARS – EXTRA CREDIT

- Wednesday March 1 – 12:30pm – Cherry Parkes room 106
  - Dr. Jiaxuan You CS PhD Stanford University  
Talk title: Learning from the Interconnected World with Graphs
- Friday March 3 – 12:30pm – Cherry Parkes room 106
  - Dr. Dongfang Zhao CS PhD Illinois Institute of Technology  
Talk title: High-Performance Data-Intensive Computing Systems
- Earn up to 2.5% extra credit added to the overall course grade
- Scored out of 5 total points
- First seminar – earn 2 points
- 2<sup>nd</sup>, 3<sup>rd</sup>, 4<sup>th</sup> seminar – earn 1 point each

Add to final course grade:  
 $(\text{Total\_seminar\_points} / 5 \text{ points}) * 2.5\%$

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

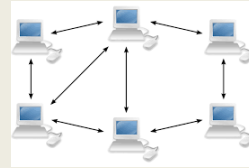
L16.6

6

## FEEDBACK FROM 2/23

- **Where are the message queues located in peer-to-peer systems?**

- “**Peer-to-peer**” refers to a network topology, which is how nodes are connected



- A **message queue** is a component of a component-based or n-tier application

- A message queue is an architectural component that facilitates implementation of distributed applications

- *Where could a message queue be located for a distributed system that exists across a group of nodes with a peer-to-peer network topology?*

- *Is the message queue replicated?*
- *If numbering nodes from left-to-right, which is most accessible?*

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.7

7

## FEEDBACK - 2

- **Is rumor spreading a subset of gossiping spreading, just with the addition of the stopping mechanism?**

- YES

- **Does rumor spreading also apply the push/pull/push-pull model?**

- Yes, it could

- Push/pull/push-pull refers to how data is spread in the system between nodes

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.8

8

## FEEDBACK - 3

- **About message flooding on page 59 of slides. What's the meaning of a higher dimension? Why just choose node 1101...?**
- If the message arrives at node 1101 on an edge labeled as level 2
- We then label all outbound edges from 1101, and only forward the message along edges that are a higher order than that of the edge that received the message
- This way we can minimize the number of messages for broadcast of the n-dimensional hypercube
- We will go over this again...

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.9

9

## OBJECTIVES - 2/28

- Questions from 2/23
- **Assignment 2: Replicated Key Value Store**
- Chapter 4.4 - Review / Finish
- Chapter 6: Coordination
  - Chapter 6.1: Clock Synchronization
  - Chapter 6.2: Logical Clocks  
Vector Clocks
- Class Activity - Total Ordered Multicasting (Thursday)
  - Chapter 6.3: Distributed Mutual Exclusion

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.10

10

## SHORT-HAND-CODES FOR MEMBERSHIP TRACKING APPROACHES

- Include readme.txt or doc file with instructions in submission
- Must document membership tracking method

>> *please indicate which types to test* <<

ID	Description
F	Static file membership tracking - file is not reread
FD	Static file membership tracking DYNAMIC - file is periodically reread to refresh membership list
T	TCP membership tracking - servers are configured to refer to central membership server
U	UDP membership tracking - automatically discovers nodes with no configuration

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.11

11

## ASSIGNMENT 2

- **Sunday March 12<sup>th</sup>**
- **Goal: Replicated Key Value Store**
- **Team signup posted on Canvas under 'People'**
- **Build off of Assignment 1 GenericNode**
- **Focus on TCP client/server w/ replication**
- **How to track membership for data replication?**
  - Can implement multiple types of membership tracking for extra credit

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.12

12

## OBJECTIVES - 2/28

- Questions from 2/23
- Assignment 2: Replicated Key Value Store
- **Chapter 4.4 - Review / Finish**
- Chapter 6: Coordination
  - Chapter 6.1: Clock Synchronization
  - Chapter 6.2: Logical Clocks  
Vector Clocks
- Class Activity - Total Ordered Multicasting (Thursday)
  - Chapter 6.3: Distributed Mutual Exclusion

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

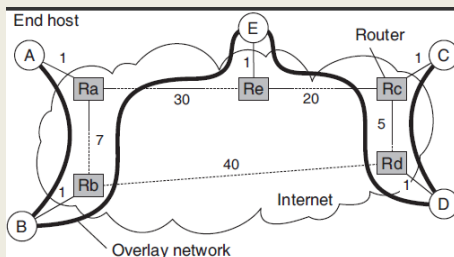
L16.13

13

## REVIEW

- Consider sending a message from node B to D:
- What is the Link Stress at physical node Re ?
- What is the delay from B to D via the underlying network ?
- What is the delay from B to D via the overlay network ?
- What is the stretch from B to D ?

Numbers represent  
network delay  
between nodes



February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L15.14

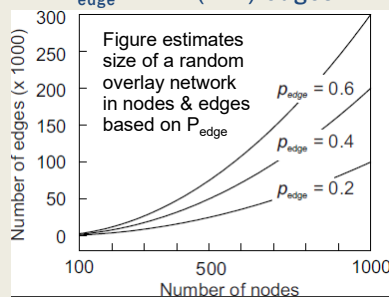
14

## RANDOM GRAPHS

- Used when no information on the structure of the overlay network
- Assume network can be represented as a **Random graph**
- Random graphs are described by a probability distribution:
  1. Given a probability  $P_{\text{edge}}$  that two nodes are joined
  2. Size of a random overlay network is:  $\frac{1}{2} * P_{\text{edge}} * N * (N-1)$  edges

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

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



February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L15.15

15

## PROBABILISTIC FLOODING



*Washington state in winter?*

- When a node is flooding a message  $m$ :  
( $p_{\text{flood}}$ ) is the probability that the message is spread to a specific neighbor  $= (p_{\text{flood}})$
- Throttle message flooding based on a probability
- Implementation needs to consider # of neighbors to achieve various  $p_{\text{flood}}$  scores
- With lower  $p_{\text{flood}}$  messages may not reach all nodes
- **Efficiency of probabilistic broadcasting:** For random network with 10,000 nodes
  - With  $p_{\text{edge}} = 0.1$  and  $p_{\text{flood}} = 0.01$
  - Achieves 50-fold reduction in messages vs. full flooding

February 28, 2023


TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L15.16

16



## PROBABILISTIC FLOODING




*Washington state in winter?*

- When a node is flooding a message  $m$ :  
( $p_{\text{flood}}$ ) is the probability that the message is spread to a specific neighbor  $= (p_{\text{flood}})$
- Throttling
- Implementation
- With  $p_{\text{edge}} = 0.1$
- Efficiency of probabilistic broadcasting:** For random network with 10,000 nodes
  - With  $p_{\text{edge}} = 0.1$  and  $p_{\text{flood}} = 0.01$
  - Achieves 50-fold reduction in messages vs. full flooding

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L15.17
-------------------	---	--------

17

## PROBABILISTIC FLOODING



*Washington state in winter?*


- When a node is flooding a message  $m$ :  
( $p_{\text{flood}}$ ) is the probability that the message is spread to a specific neighbor  $= (p_{\text{flood}})$
- Throttling
- Implementation
- With  $p_{\text{edge}} = 0.1$
- Efficiency of probabilistic broadcasting:** For random network with 10,000 nodes
  - With  $p_{\text{edge}} = 0.1$  and  $p_{\text{flood}} = 0.01$
  - Achieves 50-fold reduction in messages vs. full flooding

Edges =  $\frac{1}{2} * P_{\text{edge}} * N * (N-1)$

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L15.18
-------------------	---	--------

18

# PROBABILISTIC FLOODING



*Washington state in winter?*

- When a node is flooding a message  $m$ :  
( $p_{\text{flood}}$ ) is the probability that the message is spread to a specific node
- Throughput
- Implementation achieves
- With
- Efficient 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


**How many edges does network with 10,000 nodes have with  $p_{\text{edge}}=0.1$ ?**

$$\text{Edges} = \frac{1}{2} * P_{\text{edge}} * N * (N-1)$$
$$\frac{1}{2} * (.1) * (10000) * (9999)$$

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L15.19
-------------------	---	--------

19

# PROBABILISTIC FLOODING



*Washington state in winter?*

- When a node is flooding a message  $m$ :  
( $p_{\text{flood}}$ ) is the probability that the message is spread to a specific node
- Throughput
- Implementation achieves
- With
- Efficient 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

**How many edges does network with 10,000 nodes have with  $p_{\text{edge}}=0.1$ ?**


$$\text{Edges} = \frac{1}{2} * P_{\text{edge}} * N * (N-1)$$
$$\frac{1}{2} * (.1) * (10000) * (9999)$$

**4,999,500 edges**

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L15.20
-------------------	---	--------

20

## PROBABILISTIC FLOODING



Washington state in winter?

- When a node is flooding a message  $m$ :  
( $p_{\text{flood}}$ ) is the probability that the message is spread to a specific neighbor = ( $p_{\text{flood}}$ )
- Throt
- Imple
- achie
- With
- Effici**
- netwo
- With  $p_{\text{edge}} = 0.1$  and  $p_{\text{flood}} = .01$
- Achieves 50-fold reduction in messages vs. full flooding

**What does it mean to have  $p_{\text{flood}} = .01$ ?**

*With lower  $p_{\text{flood}}$  messages may not reach all nodes*


*If there are 10,000 nodes and ~5 million edges, how many messages is full flooding?*

*> every node sends to all neighbors across all edges*

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L15.21
-------------------	---	--------

21

## PROBABILISTIC FLOODING



Washington state in winter?

- When a node is flooding a message  $m$ :  
( $p_{\text{flood}}$ ) is the probability that the message is spread to a specific neighbor = ( $p_{\text{flood}}$ )

**What does it mean to have  $p_{\text{flood}} = .01$ ?**

node  $Q$  has  $n$  neighbors  
Probability that all neighbors don't forward message to  $Q$  is  $p = (1 - p_{\text{flood}})^n$


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

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L15.22
-------------------	---	--------

22

## PROBABILISTIC FLOODING



*Washington state in winter?*

What does it mean to have  $p_{\text{flood}} = .01$ ?

node Q has n neighbors  
Probability that no neighbors of Q forward the message to Q:  
 $p = (1 - p_{\text{flood}})^n$  ← probability of Q not getting the message

10 nodes: if  $n=10$ ,  $p = (1 - .01)^{10} = .904$  (small network, few edges, likely Q doesn't get msg)  
100 nodes: if  $n=100$ ,  $p = (1 - .01)^{100} = .366$  (less likely Q doesn't get msg)  
298 nodes: if  $n=298$ ,  $p = (1 - .01)^{298} = .05$  (Q probably get msg)

- Achieves 50-fold reduction in messages vs. full flooding

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L15.23
-------------------	---	--------

23

## MESSAGE FLOODING W/ HYPERCUBE

- **Hypercube:** for broadcast send minimum # of messages  
Only forward msg along edges with **higher dimension**
- Node(1101)-neighbors {0101,1100,1001,1111}
- Node (1101) - receives msg on incoming broadcast edge = 2
- **Broadcast from 1101 - Label Edges:**
  - Edge to 0101 - labeled 1 - change the 1<sup>st</sup> bit
  - Edge to 1100 - labeled 4 - change the 4<sup>th</sup> bit **\*<FORWARD>\***
  - Edge to 1001 - labeled 2 - change the 2<sup>nd</sup> bit
  - Edge to 1111 - labeled 3 - change the 3<sup>rd</sup> 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

February 23, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L15.24
-------------------	---	--------

24

## RUMOR SPREADING

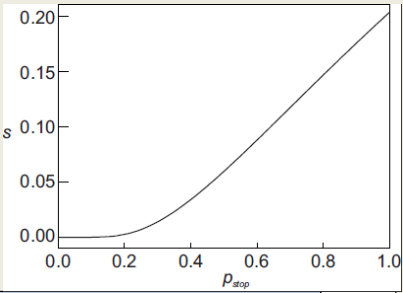
- Variant of epidemic protocols
- Provides an approach to “**stop**” message spreading
- Mimics “gossiping” in real life
  
- **Rumor spreading:**
- **Node P** receives new data **Item X**
- Contacts an arbitrary **node Q** to push update
- **Node Q** reports already receiving **Item X** from another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

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L15.25
-------------------	---	--------

25

## RUMOR SPREADING - 2

- $p_{\text{stop}}$  is the probability node will stop spreading once contacting a node that already has the message
- Rumor spreading does not guarantee all nodes will be updated
- 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  $p_{\text{stop}}$
- Susceptible nodes ( $s$ ) vs. probability of stopping  $\rightarrow$



February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L15.26
-------------------	---	--------

26

## 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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L15.27

27

## 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**

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L15.28

28

## OBJECTIVES - 2/28

- Questions from 2/23
- Assignment 2: Replicated Key Value Store
- Chapter 4.4 - Review / Finish
- Chapter 6: Coordination
  - **Chapter 6.1: Clock Synchronization**
  - Chapter 6.2: Logical Clocks  
Vector Clocks
- Class Activity - Total Ordered Multicasting (Thursday)
  - Chapter 6.3: Distributed Mutual Exclusion

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.29
-------------------	---	--------

29

## 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*)

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.30
-------------------	---	--------

30



## CH. 6.1: CLOCK SYNCHRONIZATION

L16.31

31

## CLOCK SYNCHRONIZATION

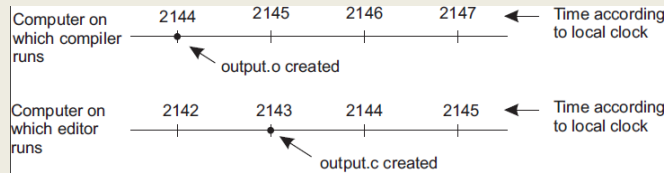
- **Example:**
- “make” is used to compile source files into binary object and executable files
- As an optimization, make only compiles files when the “last modified time” of source files is more recent than object and executables
- Consider if files are on a shared disk of a distributed system where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

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.32
-------------------	---	--------

32



## 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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.33

33

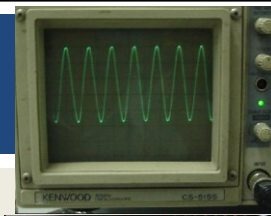


## PHYSICAL CLOCKS

- **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



February 28, 2023



TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.34

34

## 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



February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.35
-------------------	---	--------

35

## UNIVERSAL COORDINATED TIME

- **Universal Coordinated Time (UTC)** `ubuntu@ip-172-31-58-89:~$ date`  
`Thu Nov 16 10:13:39 UTC 2017`
  - 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

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.36
-------------------	---	--------

36

## COMPUTING: CLOCK CHALLENGES

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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.37

37

## 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 ( $\pi$ ):** 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 ( $\rho$ ):** clock specifications include one

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.38

38

## CLOCK SYNCHRONIZATION - 2

- If two clocks drift from UTC in opposite directions, after time  $\Delta t$  after synchronization, they may be  $2\rho$  apart.
  - $\rho$  - clock drift rate,  $\pi$  - clock precision (max 50ns)
- Clocks must be resynchronized every  $\pi/2\rho$  seconds
- Network time protocol
- Provide coordination of time for servers
- Leverage distributed network of time servers

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.39
-------------------	---	--------

39

## 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

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.40
-------------------	---	--------

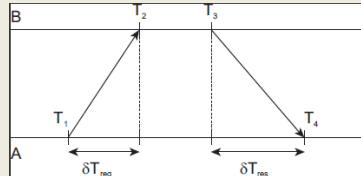
40

## NTP - 2

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

Time server B

Client A



1. A sends message to B, with timestamp  $T_1$
2. B records time of receipt  $T_2$  (from local clock)
3. B returns response with send time  $T_3$ , and receipt time  $T_2$
4. A records arrival of  $T_4$

- Assuming propagation delay of  $A \rightarrow B \rightarrow A$  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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.41

41

## 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`)

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.42

42

## 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

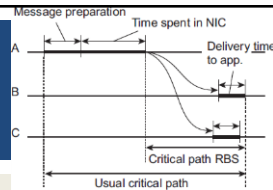
February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.43

43

## 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 time precision, 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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.44

44

## REFERENCE BROADCAST SYNCHRONIZATION (RBS)

- Node broadcasts reference message  $k$
- Each node  $p$  records time  $T_{p,k}$  when  $k$  is received
- $T_{p,k}$  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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.45

45

## 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:

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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.46

46

WE WILL RETURN AT  
3:01 PM



47

## OBJECTIVES - 2/28

- Questions from 2/23
- Assignment 2: Replicated Key Value Store
- Chapter 4.4 - Review / Finish
- Chapter 6: Coordination
  - Chapter 6.1: Clock Synchronization
  - **Chapter 6.2: Logical Clocks**  
Vector Clocks
- Class Activity - Total Ordered Multicasting (Thursday)
  - Chapter 6.3: Distributed Mutual Exclusion

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.48
-------------------	---	--------

48



## CHAPTER 6 - COORDINATION

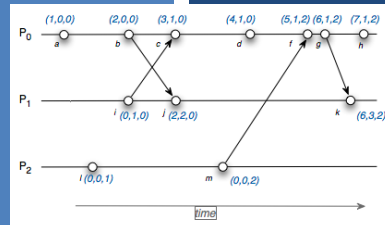
- 6.1 Clock Synchronization
  - Physical clocks
  - Clock synchronization algorithms
- 6.2 Logical clocks
  - Lamport clocks
  - Vector clocks
- 6.3 Mutual exclusion
- 6.4 Election algorithms
- 6.6 Distributed event matching (*light*)
- 6.7 Gossip-based coordination (*light*)

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.49

49



## CH. 6.2: LOGICAL CLOCKS

L16.50

50

## 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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.51

51

## 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  $A \rightarrow C$

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.52

52

## 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

February 28, 2023

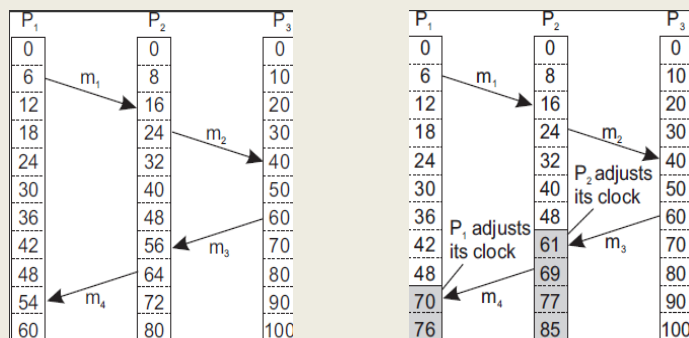
TCSS558: Applied Distributed Computing [Winter 2023]  
 School of Engineering and Technology, University of Washington - Tacoma

L16.53

53

## LOGICAL CLOCKS - 4

- Three processes each with local clocks
- **Lamport's algorithm** corrects process clock values
- Always propagate the most recent known value of logical time



February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
 School of Engineering and Technology, University of Washington - Tacoma

L16.54

54

## 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**

69: P2 sends m4 to P1

54: P1 receives m4

**70: P1 clock reset=70**

P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>	P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>
0	0	0	0	0	0
6	8	10	6	8	10
12	16	20	12	16	20
18	24	30	18	24	30
24	32	40	24	32	40
30	40	50	30	40	50
36	48	60	36	48	60
42	56	70	42	61	70
48	64	80	48	69	80
54	72	90	70	77	90
60	80	100	76	85	100

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.55

55

## 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 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  
 $C_j \leftarrow \max\{C_j, \text{timestamp}(m)\}$
4. Ties broken by considering Proc ID:  $i < j$ ;  $\langle 40, i \rangle < \langle 40, j \rangle$   
 Both Lamport clocks are = 40  
 The winner has a higher alphanumeric Process ID  
 $J$  (winner) is greater than  $i$ , alphabetically

February 28, 2023

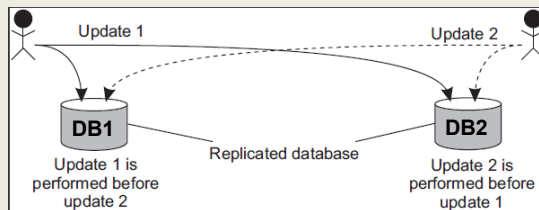
TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.56

56

## 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**

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.57

57

## TOTAL-ORDERED MULTICASTING EXAMPLE

- Two messages ( $m_1, m_2$ ) must be distributed, to two processes ( $p_1, p_2$ )
- We assume messages have correct lamport clock timestamps
- $m_1(10, p_1, \text{add } \$100)$
- $m_2(12, p_2, \text{add } 1\% \text{ interest})$
- Each process maintains a queue of messages
- Arriving messages are placed into queues ordered by the Lamport clock timestamp
- In each queue, each message must be acknowledged by every process in the system before operations can be applied to the local database

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.58

58

## TOTAL-ORDERED MULTICASTING EXAMPLE

- Two messages ( $m_1, m_2$ ) must be distributed, to two processes ( $p_1, p_2$ )
- We assume messages have correct lamport clock timestamps
- $m_1(10, p_1, \text{add } \$100)$

### Key point:

Multicast messages are also received by the sender (*itself*)

Arriving messages are placed into queues ordered by the Lamport clock timestamp

- In each queue, each message must be acknowledged by every process in the system before operations can be applied to the local database

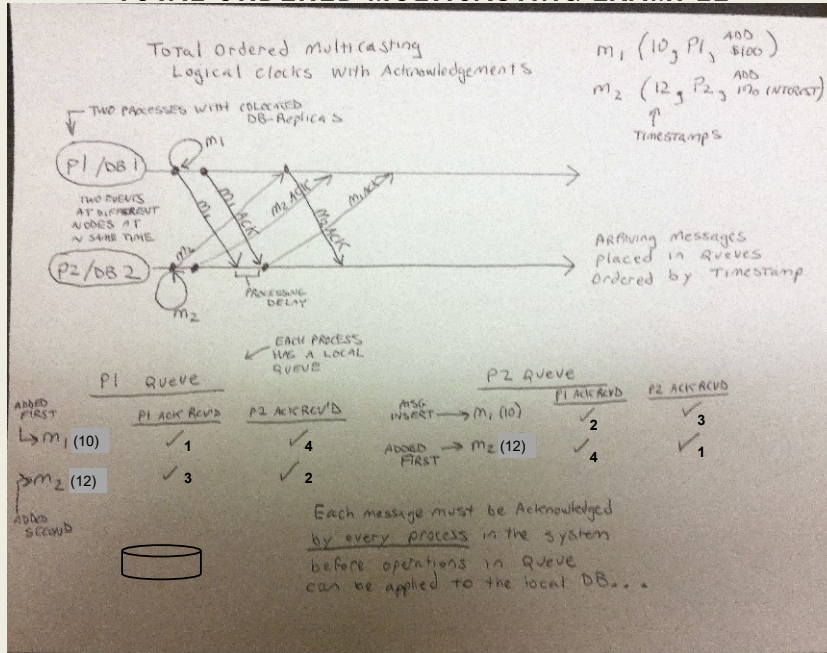
February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
 School of Engineering and Technology, University of Washington - Tacoma

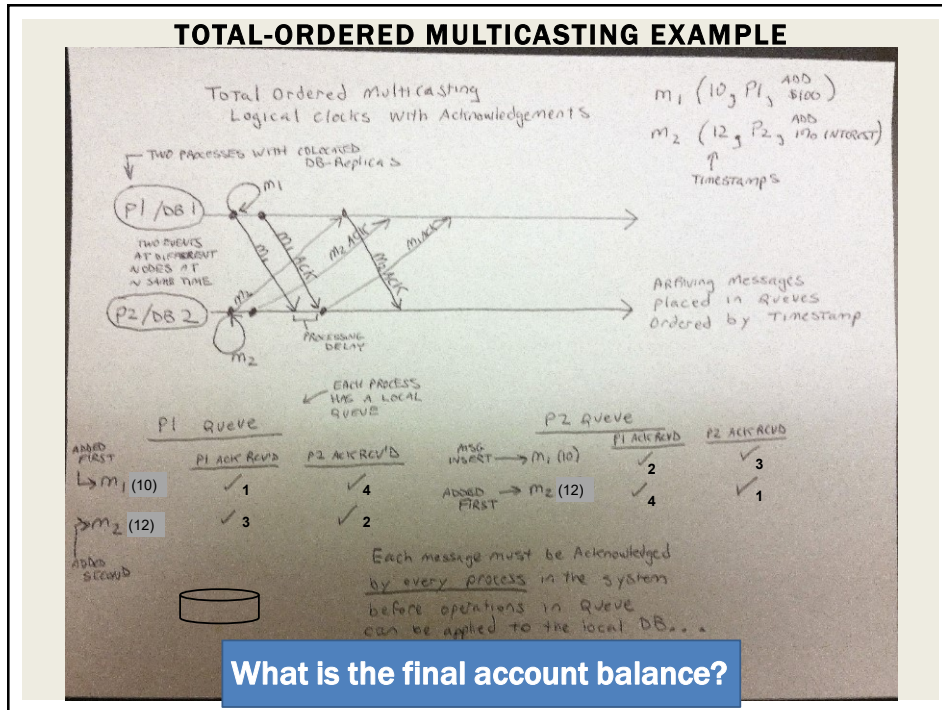
L16.59

59

## TOTAL-ORDERED MULTICASTING EXAMPLE



60



61

TOTAL-ORDERED MULTICASTING - 2

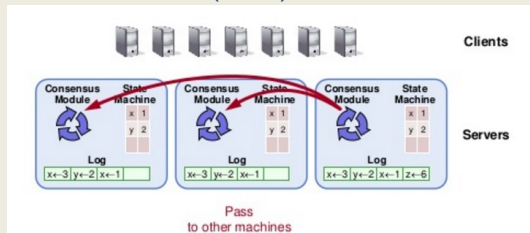
- Each message timestamped with local logical clock of sender
- **Multicast messages are also received by the sender (itself)**
- Assumptions:
  - Messages from same sender received in order they were sent
  - No messages are lost
- When messages arrive they are placed in local queue ordered by timestamp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

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.62
-------------------	---	--------

62

## TOTAL-ORDERED MULTICASTING - 3

- Can be used to implement 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)



February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.63

63

## OBJECTIVES - 2/28

- Questions from 2/23
- Assignment 2: Replicated Key Value Store
- Chapter 4.4 - Review / Finish
- Chapter 6: Coordination
  - Chapter 6.1: Clock Synchronization
  - Chapter 6.2: Logical Clocks
  - Vector Clocks**
  - Class Activity - Total Ordered Multicasting (Thursday)
    - Chapter 6.3: Distributed Mutual Exclusion

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.64

64



## 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
- But what is causality? ...

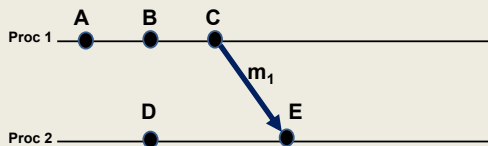
February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.65

65

## WHAT IS CAUSALITY?



- Having a causal relationship between two events (A and E) indicates that event E results from the occurrence of event A.
- When one event results from another, there is a causal relationship between the two events.
- This is also referred to as ***cause and effect***.

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.66

66

## CAUSALITY - 2

- **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
  
- Lamport/Vector clocks can help us suggest possible causality
- But we never know for sure...

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.67

67

## CAUSALITY - 3

- 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?**

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.68

68

## VECTOR CLOCKS

- Vector clocks help keep track of **causal history**
- If two local events happened at process P, then the causal history  $H(p_2)$  of event  $p_2$  is  $\{p_1, p_2\}$
- P sends messages to Q (event  $p_3$ )
- Q previously performed event  $q_1$
- Q records arrival of message as  $q_2$
- Causal histories merged at Q  $H(q_2) = \{p_1, p_2, p_3, q_1, q_2\}$
- Fortunately, can simply store history of last event, as a vector clock  $\rightarrow H(q_2) = (3, 2)$
- Each entry corresponds to the last event at the process

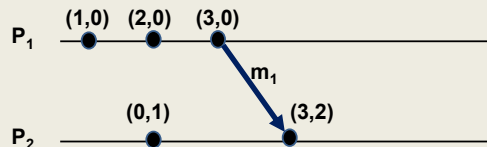
February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.69

69

## 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 ( $VC_j$ )

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.70

70

## VECTOR CLOCKS - 3

- Pj knows the # of events at Pi based on the timestamps of the received message
- Pj 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...

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.71

71

## VECTOR CLOCKS EXAMPLE

- Local clock is underlined

CAUSALITY

$m_2$	$m_4$	$m_2 < m_4$	$m_2 > m_4$	Conclusion
(2,1,0)	(4,3,0)	Yes	No	m2 may causally precede m4

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.72

72

## VECTOR CLOCKS EXAMPLE - 2

$m_2$	$m_4$	$m_2 < m_4$	$m_2 > m_4$	Conclusion
(4,1,0)	(2,3,0)	No	No	m2 and m4 may conflict

- P3 can't determine if m4 may be causally dependent on m2
- **Is m4 causally dependent on m3 ?**

February 28, 2023
TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma
L16.73

73

## VECTOR CLOCKS EXAMPLE - 3

- Provide a vector clock label for unlabeled events

February 28, 2023
TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma
L16.74

74

## VECTOR CLOCKS EXAMPLE - 4

- TRUE/FALSE:
- The sending of message  $m_3$  is causally dependent on the sending of message  $m_1$ .
- The sending of message  $m_2$  is causally dependent on the sending of message  $m_1$ .

February 28, 2023	TCCS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.75
-------------------	---	--------

75

## VECTOR CLOCKS EXAMPLE - 5

- TRUE/FALSE:
- $P_1 (1,0,0)$  and  $P_3 (0,0,1)$  may be concurrent events.
- $P_2 (0,1,1)$  and  $P_3 (0,0,1)$  may be concurrent events.
- $P_1 (1,0,0)$  and  $P_2 (0,1,1)$  may be concurrent events.

February 28, 2023	TCCS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.76
-------------------	---	--------

76

## OBJECTIVES - 2/28

- Questions from 2/23
- Assignment 2: Replicated Key Value Store
- Chapter 4.4 - Review / Finish
- Chapter 6: Coordination
  - Chapter 6.1: Clock Synchronization
  - Chapter 6.2: Logical Clocks
  - **Vector Clocks**
- Class Activity - Total Ordered Multicasting (Thursday)
  - Chapter 6.3: Distributed Mutual Exclusion

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.77
-------------------	---	--------

77

## OBJECTIVES - 2/28

- Questions from 2/23
- Assignment 2: Replicated Key Value Store
- Chapter 4.4 - Review / Finish
- Chapter 6: Coordination
  - Chapter 6.1: Clock Synchronization
  - Chapter 6.2: Logical Clocks
  - Vector Clocks
- **Class Activity - Total Ordered Multicasting (Thursday)**
  - Chapter 6.3: Distributed Mutual Exclusion

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.78
-------------------	---	--------

78

## OBJECTIVES - 2/28

- Questions from 2/23
- Assignment 2: Replicated Key Value Store
- Chapter 4.4 - Review / Finish
- Chapter 6: Coordination
  - Chapter 6.1: Clock Synchronization
  - Chapter 6.2: Logical Clocks  
Vector Clocks
- Class Activity - Total Ordered Multicasting (Thursday)
  - Chapter 6.3: Distributed Mutual Exclusion

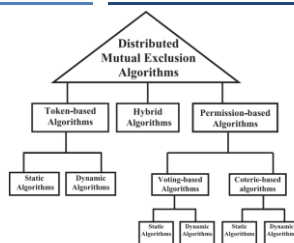
February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.79

79

## CH. 6.3: DISTRIBUTED MUTUAL EXCLUSION



L16.80

80



## DISTRIBUTED MUTUAL EXCLUSION ALGORITHMS

- Coordinating access among distributed processes to a shared resource requires **Distributed Mutual Exclusion**
- **Algorithms in 6.3**
- Token-ring algorithm
- **Permission-based algorithms:**
- Centralized algorithm
- Distributed algorithm (Ricart and Agrawala)
- Decentralized voting algorithm (Lin et al.)

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.81

81

## 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**

February 28, 2023

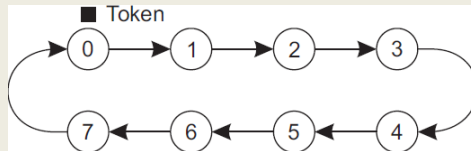
TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.82

82

## 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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.83

83

## 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 (**lock**) 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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.84

84

## DISTRIBUTED MUTUAL EXCLUSION ALGORITHMS - 3

- **Permission-based algorithms**
  - Processes must require permission from other processes before first acquiring access to the resource
    - CONTRAST: Token-ring did not ask nodes for permission
- **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”*

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.85
-------------------	---	--------

85

## CENTRALIZED MUTUAL EXCLUSION

Permission granted from coordinator ∨ No response from coordinator

**P<sub>1</sub> executes**

**P<sub>2</sub> blocks**

**P<sub>1</sub> finishes; P<sub>2</sub> executes**

- When resource not available, coordinator can block the requesting process, or respond with a reject message
- P2 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 (OK), release)

February 28, 2023	TCSS558: Applied Distributed Computing [Winter 2023] School of Engineering and Technology, University of Washington - Tacoma	L16.86
-------------------	---	--------

86

## CENTRALIZED MUTUAL EXCLUSION - 2

- **Issues**
- Coordinator is a single point of failure
- Processes can't distinguish dead coordinator from **"blocking"** when resource is unavailable
  - 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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.87

87

## 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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.88

88

## 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. *If node is also waiting to access the resource*: perform a timestamp comparison -
    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

February 28, 2023

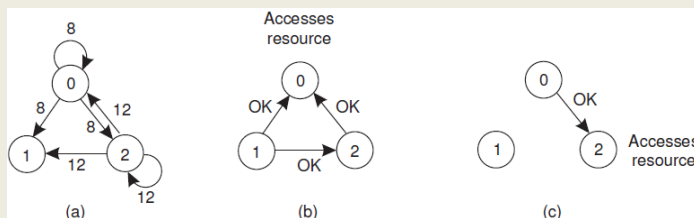
TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.89

89

## DISTRIBUTED ALGORITHM - 3

- Node 0 and Node 2 simultaneously request access to resource
- Node 0's time stamp is lower (8) than Node 2 (12)
- Node 1 and Node 2 grant Node 0 access
- Node 1 is not interested in the resource, it OKs both requests



- **In case of conflict, lowest timestamp wins!**
  - Node 2 rejects its own request (1@) in favor of node 0 (8)

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.90

90

## CHALLENGES WITH DISTRIBUTED ALGORITHM

- **Problem:** Algorithm has  $N$  points of failure !
- Where  $N$  = Number of Nodes in the system
  
- **No Reply Problem:** When node is accessing the resource, it does not respond
  - Lack of response can be confused with **failure**
  - **Possible 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

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.91

91

## 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*
  - *If one node gets majority of acknowledges no other can*
  - *Requires every node to know size of system (# of nodes)*
- Distributed algorithm for mutual exclusion works best for:
  - Small groups of processes
  - When memberships rarely change

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.92

92

## DECENTRALIZED ALGORITHM

- Lin et al. [2004], decentralized voting algorithm
- Resource is replicated N times
- Each replica has its own coordinator ... (N coordinators)
- Accessing resource requires majority vote:  
total votes (m) > N/2 coordinators
- **Assumption #1:** When coordinator does not give permission to access a resource (because it is busy) it will inform the requester

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.93

93

## DECENTRALIZED ALGORITHM - 2

- **Assumption #2:** When a coordinator crashes, it recovers quickly, but will have forgotten votes before the crash.
- Approach assumes coordinators reset arbitrarily at any time
- **Risk:** on crash, coordinator forgets it previously granted permission to the shared resource, and on recovery it errantly grants permission again
- **The Hope:** if coordinator crashes, *upon recovery, the node granted access to the resource has already finished before the restored coordinator grants access again . . .*

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.94

94

## DECENTRALIZED ALGORITHM - 3

- With 99.167% coordinator availability (30 sec downtime/hour) chance of violating correctness **is so low** it can be neglected in comparison to other types of failure
- Leverages fact that a new node must obtain a majority vote to access resource, **which requires time**

N	m	p	Violation	N	m	p	Violation
8	5	3 sec/hour	$< 10^{-15}$	8	5	30 sec/hour	$< 10^{-10}$
8	6	3 sec/hour	$< 10^{-18}$	8	6	30 sec/hour	$< 10^{-11}$
16	9	3 sec/hour	$< 10^{-27}$	16	9	30 sec/hour	$< 10^{-18}$
16	12	3 sec/hour	$< 10^{-36}$	16	12	30 sec/hour	$< 10^{-24}$
32	17	3 sec/hour	$< 10^{-52}$	32	17	30 sec/hour	$< 10^{-35}$
32	24	3 sec/hour	$< 10^{-73}$	32	24	30 sec/hour	$< 10^{-49}$

N = number of resource replicas, m = required "majority" vote  
 p=seconds per hour coordinator is offline

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
 School of Engineering and Technology, University of Washington - Tacoma

L16.95

95

## DECENTRALIZED ALGORITHM - 4

- Back-off Polling Approach for *permission-denied*:**
- If permission to access a resource is denied via majority vote, process can poll to gain access again with a **random** delay (**known as back-off**)
- Node waits for a random amount, retries...
- If too many nodes compete to gain access to a resource, majority vote can lead to low resource utilization
  - No one can achieve majority vote to obtain access to the shared resource**
  - Mimics elections where with too many candidates, where no one candidate can get >50% of the total vote**
- Problem Solution detailed in [Lin et al. 2014]

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
 School of Engineering and Technology, University of Washington - Tacoma

L16.96

96



## DISTRIBUTED MUTUAL EXCLUSION ALGORITHMS REVIEW

- Which algorithm offers the best scalability to support distributed mutual exclusion in a large distributed system?
- (A) Token-ring algorithm
- (B) Centralized algorithm
- (C) Distributed algorithm
- (D) Decentralized voting algorithm

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.97

97

## DISTRIBUTED MUTUAL EXCLUSION ALGORITHMS REVIEW - 2

- Which algorithm(s) involve blocking when a resource is not available?
- (A) Token-ring algorithm
- (B) Centralized algorithm
- (C) Distributed algorithm
- (D) Decentralized voting algorithm

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.98

98

## DISTRIBUTED MUTUAL EXCLUSION ALGORITHMS REVIEW - 3

- Which algorithm(s) involve arriving at a consensus to determine whether a node should be granted access to a resource?
- (A) Token-ring algorithm
- (B) Centralized algorithm
- (C) Distributed algorithm
- (D) Decentralized voting algorithm

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.99

99

## DISTRIBUTED MUTUAL EXCLUSION ALGORITHMS REVIEW - 4

- Which algorithm(s) have  $N$  points of failure, where  $N$  = Number of Nodes in the system?
- (A) Token-ring algorithm
- (B) Centralized algorithm
- (C) Distributed algorithm
- (D) Decentralized voting algorithm

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.100

100

**QUESTIONS**

February 28, 2023

TCSS558: Applied Distributed Computing [Winter 2023]  
School of Engineering and Technology, University of Washington - Tacoma

L16.10  
1

101