





# FEEDBACK - 3/6

- For the static N-dimensional hypercube, broadcast requires N-1 messages
- Does it mean that N-1 is the minimal message it can spread?
- Message spreading algorithms are typically concerned with how to disseminate messages across unstructured adhoc topologies
- As the size of the network is unknown, the goal is to experiment with different approaches to message spreading and spread termination
- N-1 is the minimum messages to fully disseminate a message starting at one node, to the entire hypercube
   March 11, 2019
   TCSSSS: Appled Distributed Computing [Winter 2019]
   TCSSSS: Appled Distributed Computing [Winter 2019]
   Starting at Distributed Computing [























































#### 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:

March 11, 2019

Easy to implement compared to distributed alternatives

TCSS558: Applied Distributed Computing [Winter 2019] School of Engineering and Technology, University of Was L15.34















3

 $\bigtriangledown$ 







- When <u>any</u> process notices the coordinator is no longer responding to requests, it initiates an election
- Process P<sub>k</sub> initiates an election as follows:
- P<sub>k</sub> sends an ELECTION message to all processes with higher process IDs (P<sub>k+1</sub>, P<sub>k+2</sub>, ... P<sub>N-1</sub>)
- 2. If no one responds, P<sub>k</sub> wins the election and becomes coordinator
- 3. If one of the higher-ups answers, it takes over and runs the election.
- When the higher numbered process receives an ELECTION message from a lower-numbered colleague, it responds with "OK", indicating it's alive, and it takes over the election.

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

L15.47

BULLY ALGORITHM - 2
The higher numbered process then holds an election with <u>only</u> higher numbered processes (nodes).
Eventually all processes give up except one, and the remaining process becomes the new coordinator.
The coordinator announces victory by sending all processes a message stating it is starting as the coordinator.
If a higher numbered node that was previously down comes back up, it holds an election, and ultimately takes over the coordinator role.
The process with the "biggest" ID in town always wins.
Hence the name, <u>bully algorithm</u>

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

March 11, 2019

March 11, 2019

L15.48





















L15.59





- Give N tokens to N randomly chosen nodes
- No node can hold more than (1) token
- Tokens are "repelling force". Other tokens move away
- All tokens exert the same repelling force
- This automates token distribution across an overlay network

March 11, 2019 TCSSS58: Applied Distributed Computing [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma







## PROBLEMS WITH PAXOS - 2

- Other attempts to flesh out details are divergent from Lamport's own sketches
- Problem 3: Paxos architecture is poor for building practical systems
- Paxos' notion of consensus is for a single log entry
- Consensus approach can be designed around a sequential log
- Problem 4: Paxos approach uses a symmetric peer-topeer approach vs. a leader-based approach
  - Works when just (1) decision

```
    Having a leader simplifies making multiple decisions
```

March 11, 2019 TCSS558: Applied Distributed Computing [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma

L15.65

# RESULTING PROBLEMS Implementations of Paxos typically diverge as each develops a different architecture for solving the difficult problem(s) of implementing Paxos Paxos formulation is good for proving theorems about correctness, but challenging to use for implementing real systems Though it has been used a fair bit See paper: Consensus In the Cloud: Paxos Systems Demystified

 Observation: significant gaps between the description of the algorithm and the needs of a real-world system, result in final systems based on divergent, unproven protocols March 11, 2019
 TCSSS8: Applied Distributed Computing (Where 2019) Stand Technology, University of Washington-Tacoma
 Lisce







TERMS							
Raft divides f Terms are nu Terms start w If election re- <b>term</b> is starte There is only Terms act as Each server s Terms are exc	time into <u>TERMS</u> of arbitrary length mbered with consecutive integers <i>v</i> ith an election (term # is incremented) sults in a SPLIT VOTE, term ends, and a <b>n</b> ed with an election (1) <u>Leader</u> in any given term a <u>logical clock</u> stores current term number changed in communication						
March 11, 2019	TCSS558: Applied Distributed Computing [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma						



























TCSS558: Applied Distributed Computing [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma December 5, 2017











# REPLICATION TRADEOFF EXAMPLE

- Process P accesses a local replica N times per second
- Replica is updated M times per second
- Updates involve complete refreshes of the data
- If N << M (very low access rate) many updates M are never accessed by P.
- Network communication overhead for most updates is useless.
- TRADEOFFS:
- Either move the replica away from P
- So the total number of accesses from multiple processes is higher
- Or, apply a different strategy for updating the replica
   i.e. less frequent updates, possibly need based
- BALANCE TRADEOFF BETWEEN REPLICA ACCESS FREQUENCY
  AND COSTS OF REPLICATION (communication overhead)

   December 7, 2017
   Stobiol Engineering and Enchnology. Univer 2019
   Stobiol Engineering and Enchnology. University of Washington Tacoma

## **REPLICATION: SCALABILITY ISSUES**

#### TIGHT CONSISTENCY

- Reads must return same result
- Replication must occur after an update, before a read
- Provided by synchronous replication
- Update is performed across all copies as a single atomic operation (or transaction)
- Assignment 2 replication is with tight consistency.
- Keeping multiple copies consistent is subject to scalability problems
- May need global ordering of operations (e.g. Lamport clocks), or the use of a coordinator to assign order
- Global synchronization across a wide area network is time consuming (network latency)

December 7, 2017 TCSS558: Applied Distributed Computing [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma

L19.96





## DATA-CONSISTENCY MODELS

### CONSISTENCY MODEL

- Rules that must be followed to ensure consistency
- Represents a contract between processes and data store
- If processes agree to obey certain rules, store promises to work correctly
- No general rules for loosening consistency
- What can be tolerated is highly application dependent
- Three types of inconsistencies
- Data variation
- Staleness December 7, 2017
- Ordering of update operations TCSS558: Applied Distributed Computing School of Engineering and Technology, Ur





L19.100



	SEQUE	ENTIAL	CONS	SISTENCY		
Result process operati in the c	of any exect ses were exe ons of each order specif	ution is th ecuted in s individua ied by its j	e same a some seq I process program.	s if the operat uential order, appear <u>In thIs</u>	ions of all and the s <b>sequence</b>	
Sequentially Consistent			NOT Sequentially Consistent			
P1: W(x)a			P1: W(	x)a		
P2:	W(x)b		P2:	W(x)b		
P3:	R(x)b	R(x)a	P3:	R(x)b	R(x)a	
P4:	F	R(x)b R(x)a	P4:	R	R(x)a R(x)b	
<ul> <li>Exact o</li> <li>As long</li> <li>Process</li> </ul>	rder seen b as they all ses here mu	y processe agree st see: R(	es <b>DOES M</b> x)b, then	R(x)a		













