

## CS 455: INTRODUCTION TO DISTRIBUTED SYSTEMS [ELECTION ALGORITHMS]

Shrideep Pallickara  
Computer Science  
Colorado State University

April 12, 2018

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.1

### Frequently asked questions from the previous class survey

- Does a process always say "yes" to itself?
- What is the initial process state before any requests to enter the critical section have been made? Released?
- Approximate voting is  $\sim 2\sqrt{N}$  in the matrix approximation?
- What happens if multiple processes initiate an election for the same event, but at different times? What ensures that they are part of the same round?
- Other interesting CS chairs?

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.2

### Topics covered in this lecture

- Election algorithms
  - ▣ Ring based algorithm
- Failure detectors

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.3

## ELECTION ALGORITHMS

April 12, 2018

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.4

### Managing the identity of the elected process

- Each process  $p_i$  ( $i=1, 2, \dots, N$ ) has a variable  $elect_i$ 
  - ▣ Contains identifier of the elected process
- When a process first becomes a participant in an election
  - ▣ Set this variable to  $\perp$  indicating that it is undefined

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.5

### Requirements for the election algorithm

- E1 (safety)
  - ▣ Participant process has  $elect_i = \perp$  or  $elect_i = P$ 
    - $P$  is a non-crashed process at the end of run with the largest identifier
- E2 (liveness)
  - ▣ All processes  $p_i$  participate and eventually either set  $elect_i \neq \perp$  or crash

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.6

### Measuring performance of election algorithms

- Network **bandwidth utilization**
  - How many messages are sent?
- **Turnaround time** for the algorithm
  - Number of message transmissions between the initiation and termination of a run

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.7

### RING-BASED ELECTION ALGORITHM

April 12, 2018

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.8

### Ring-based elections

- Each process  $p_i$  has a communication channel to the next process ( $p_{(i+1) \bmod N}$ ) in the next ring
- All messages are sent **clockwise** around the ring

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.9

### Conducting elections

- Any process can begin an election
- Process marks itself as a participant and:
  - ① Places its identifier in the election message
  - ② Sends it to its clockwise neighbor
- **On forwarding** an election message
  - Process marks itself as a **participant**

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.10

When a process receives an election message it compares identifier with its own

- If the arrived identifier is **greater**
  - **Forward** message to its neighbor
- If the arrived identifier is **smaller**
  - If the process is not a participant
    - **Substitute** with own identifier and forward the message
  - If the process is already a participant
    - **Do not forward** the message

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.11

### Selecting the coordinator

- If the received identifier is that of the receiver itself?
  - That process' identifier must be the greatest
  - Becomes the **coordinator**
- Coordinator marks itself as a non-participant
  - Sends an **elected** message to its neighbor

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.12

### When a process $p_i$ receives an elected message

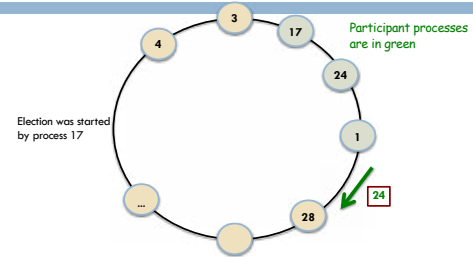
- Marks itself as a **non-participant**
- Sets  $elect_i$  to the identifier in the message
- Unless it is the coordinator, it forwards message to its neighbor

April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.13

### Ring based elections



April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.14

### Satisfying the requirements: Safety

- All identifiers are compared
  - A process must receive its own identifier back before sending an elected message
- For any two processes, the one with the **larger** identifier *will not pass* the other's identifier
  - Impossible that 2 processes will receive their own identifier back

April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.15

### Satisfying the requirements: Liveness

- Message traversals are guaranteed around the ring
  - The basic algorithm assumes no failures
- Duplicate messages arising when 2 processes start election at the same time?
  - Participant and non-participant states *extinguish* this ASAP and ...
  - *Always before* winning election result has been announced

April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.16

### Performance analysis

- If only a single process starts an election
  - Worst case is when **anti-clockwise neighbor** has the highest identifier
    - Total of  $N-1$  messages to reach this neighbor
    - Will not announce election till its identifier completes another circuit ...  $N$  messages
    - Elected message is sent  $N$  times
    - Total of  $3N-1$  messages
  - Turnaround time also is  $3N-1$  in the worst case
    - Messages are sent *sequentially*

April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.17

### FAILURE DETECTORS

April 12, 2018

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.18

## Failure detection

- In order to properly **mask** failures ...
  - ▣ We generally need to detect them well
- Failure detection is one of the cornerstones of fault tolerance in distributed systems

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.19

## What it all boils down to ...

- For a group of processes, non-faulty members must be able to decide:
  - ▣ Who is **still** a member and who is **not**
- We need to be able to detect when a member has failed

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.20

## The two mechanisms for detecting failures

- Processes actively send "**are you alive?**" messages to each other
  - ▣ For which they obviously expect an answer
- Passively wait until heartbeats ("**I am alive!**") come in from different processes
  - ▣ In practice, active pinging occurs often as well

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.21

## Huge body of theoretical work on failure detectors

- **Timeout** mechanisms are used to check whether a process has failed
- In real settings:
  - ▣ Due to unreliable networks, just because a process does not respond to a ping that does not mean it failed
  - ▣ So ... false positives can occur quite easily
    - A healthy process may be removed from the membership

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.22

## Disambiguating network failures from node failures

- Multiple nodes participate in failure detection
- When a node notices a timeout on a ping message
  - ▣ The node contacts other nodes to see if they can reach the **presumed** failing node

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.23

## Communication styles

- Asynchronous communications
    - ▣ No timing assumptions
  - Synchronous communications have bounds on
    - ▣ Maximum message transmission delay
    - ▣ Time to execute each step of a process
    - ▣ Clock drift rates
- } Allow us to use timeouts to detect process crashes

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.24

### Failure assumptions & failure detectors

- Processes use reliable channels to communicate
  - Underlying protocol handles corruptions and retransmissions
- Failures of processes are independent

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.25

### Failure assumptions & failure detectors

- Network partitions** are possible between the set of communicating processes
- Over the internet with complex topologies and independent routing choices
  - Connectivity may be **asymmetric**
    - Communication from  $p$  to  $q$  is possible, but not vice versa
  - Connectivity may be **intransitive**
    - Communication is possible from  $p$  to  $q$  to  $r$ , but  $p$  cannot communicate directly with  $r$

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.26

### Failure detectors

- Service** that processes queries about whether a process failed
- Often implemented as an object **local to each process**
  - Runs a failure detection algorithm in conjunction with counterparts at other processes
- Not necessarily accurate
  - Only as good as the information available at that process

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.27

### Unreliable failure detectors

- Produces one of two values when given the identity of a process
  - Suspected or Unsuspected
  - These values are just **hints** and may not accurately reflect if a process has failed
- Unsuspected**
  - Detector recently received evidence suggesting process has not failed
- Suspected**
  - Detector has some indication that process probably failed
    - Message not received for more than the nominal silence interval

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.28

### Unreliable failure detectors

- Suspicious may be **misplaced**
  - Process may be functioning, but on the other side of a network partition
  - Process runs slower than expected

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.29

### Reliable failure detectors

- Answers liveness queries with
  - Unsuspected
  - Failed

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.30

## BUILDING FAILURE DETECTORS

April 12, 2018

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.31

## Implementing an unreliable failure detector

- Each process  $p$ , sends a **heartbeat** every  $T$  seconds
- Failure detector uses estimate of **maximum message transmission delay** of  $D$  seconds
- If failure detector at  $q$  does not receive heartbeat from  $p$  within  $T + D$  seconds of the last one?
  - ▢ Detector reports to  $q$  that  $p$  is Suspected
- If heartbeat is received within  $T + D$  then detector at  $q$  deems  $p$  to be OK

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.32

## Choosing values for $T$ and $D$

- If we choose **small values** for  $T$  and  $D$ ?
  - ▢ Failure detector is likely to **suspect non-crashed** processes many times
  - ▢ Bandwidth will be consumed by heartbeat messages
- If we choose a large  $D$ ?
  - ▢ **Crashed processes** will often be reported as **Unsuspected**

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.33

## Practical solution to the problem

- Use timeout values that reflect **observed network delay** conditions
- If failure detector at  $q$  receives heartbeats from  $p$  every 20 seconds instead of 10 seconds?
  - ▢ Reset timeout for  $p$  to 20 seconds
- Failure detector would still be unreliable
  - ▢ But **probability of accuracy** increases

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.34

## Building reliable failure detectors

- Possible **only in synchronous systems**
- $D$  in this case is not an estimate, but an **absolute bound**
  - ▢ Absence of heartbeat from  $p$  within  $T+D$  seconds entitles detector at  $q$  to **conclude** that  $p$  has failed

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.35

## Contrasting reliable and unreliable failure detectors

- Unreliable failure detectors can be
  - ▢ **Inaccurate**
    - Suspects process that has not failed
  - ▢ **Incomplete**
    - May not suspect a process that has failed
- Reliable failure detectors require a system that is synchronous
  - ▢ Few practical systems are

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.36

So why did we look at failure detectors?

- They help us think about failures in distributed systems
- Any practical system designed to cope with failures, must detect them
  - However imperfectly!
- Unreliable failure detectors with well-defined properties help us to provide practical solutions for coordinating processes

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.37

## THE BULLY ALGORITHM

April 12, 2018

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.38

### Bully algorithm (Garcia-Molina): Key features

- Allows processes to crash during an election
- Assumptions:
  - Message delivery between processes is reliable
  - Synchronous system
    - Uses **timeouts** to detect a failure
  - Each process **knows processes that have higher identifiers**
    - Can communicate with them

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.39

### Message types

- Election
  - Sent to announce an election
- Answer
  - Sent in response to an election message
- Coordinator
  - Sent to announce the identity of the elected process

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.40

### Initiating elections

- A process begins this when it **notices** that the **coordinator has failed**
- Several processes may discover this concurrently

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.41

### Reliable failure detectors are possible because the system is synchronous

- $T_{trans}$ : Maximum transmission delay
- $T_{process}$ : Maximum delay for processing a message
- Upper bound on elapsed time between sending a message to a process & receiving a response
  - $T = 2T_{trans} + T_{process}$
  - If no response arrives within **T**, local failure detector tags intended recipient as having failed

April 12, 2018  
Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
Dept. Of Computer Science, Colorado State University

L24.42

### In the case of a failure

- Process that knows it has the highest identifier can elect itself as the coordinator
  - Simply send a coordinator message to processes with lower identifiers

April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.43

When a process with a lower identifier detects coordinator failure it initiates an election

- Send an election message to processes with higher identifiers
  - Await answer messages in response
- If no response within time  $T$ , process considers itself the coordinator
- If an answer does arrive, wait for additional time  $T'$  for coordinator message to arrive
  - If this does not arrive ... start another election

April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.44

### How a process responds to messages that it receives

- If a process  $p_i$  receives a coordinator message, it sets its variable *elected<sub>i</sub>* to the coordinator ID
- If a process receives an election message
  - ① Sends back an answer message and ...
  - ② **Begins another election**
    - Unless it has started one already

April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.45

### But why is this called the bully algorithm?

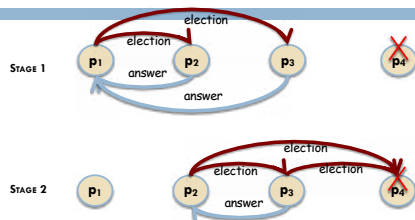
- When a process is started to replace a crashed process ... it starts an election
- If this new process has the highest identifier?
  - It decides that it is the new coordinator and announces this
- The new process becomes the coordinator **even though the current coordinator is functioning**

April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.46

### Election of a coordinator after the failure of p4

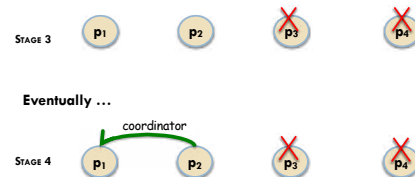


April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.47

### Election of a coordinator after the failure of p4 and then p3



April 12, 2018  
 Instructor: SHRIDEEP PALICKARA

CS455: Introduction to Distributed Systems [Spring 2018]  
 Dept. Of Computer Science, Colorado State University

L24.48



The contents of this slide set are based on the following references

- *Distributed Systems: Concepts and Design*. George Coulouris, Jean Dollimore, Tim Kindberg, Gordon Blair. 5th Edition. Addison Wesley. ISBN: 978-0132143011 [Chapter 15]
- *Distributed Systems: Principles and Paradigms*. Andrew S. Tanenbaum and Maarten Van Steen. 2nd Edition. Prentice Hall. ISBN: 0132392275/978-0132392273 [Chapter 6, 8]