# DS Lecture 7.2
# Leader Election

October 24, 2022

Michele Albano
`mialb@cs.aau.dk`

DEIS
Aalborg University
Denmark

**AALBORG UNIVERSITY**
DENMARK

Based on slides by Peter G. Jensen, AAU.

# What is it?

What is leader election?

**Leader election** algorithms ensure that one, and **only one** process is elected as the **leader** in the event of **lack of a leader**

# What is it?

1

## Examples

- ▶ Algorithms with coordinator
  - ▶ Mutex?
- ▶ Distributed replication
- ▶ DNS-servers in case of network partition
- ▶ Failover mechanism for crashes
- ▶ Parliaments?

# What is it?

## Big Question:
We detected a crash of the leader, now what?

# Agenda

Requirements

Assumptions

Chang-Roberts

Bully Algorithm

# Requirements
Leader Election

3

Let $P = \{p_1, \ldots, p_n\}$ be a set of processes and let
$L(p_i) \in P \cup \{\perp\}$ be the leader as seen from a process $p_i$

1. Safety
   - either $L(p_i) = \perp$ or $L(p_i) = p_j$
     - $p_j$ is the non-crashed process with largest identifier
2. Liveness
   - All process participate, and eventually
     - $p_i$ crashes, or
     - $L(p_i) \neq \perp$

## Note
Processes can crash during the election

# Assumptions

4

- ▶ Processes stay dead
- ▶ Crashes are reliably detected
- ▶ Identifiers are unique

# Chang Roberts

### Idea

- ▶ Pass token of largest ID in a ring
- ▶ Basic algorithm in election
    - ▶ Forward ID to "next" if higher than own,
    - ▶ Forward own ID to "next" otherwise.
    - ▶ Only one active election
- ▶ Two message types:
    - ▶ *election* to vote. If $p_i$ receives it and is not "participating", it starts participating
    - ▶ *elected* to declare who won. It $p_i$ is participating and receives its own ID, it means it won, becomes "non-participating" and send the *elected* message

# Chang-Roberts - Code

```
1    def run(self):
2        while True:
3            nxt = (self.index() + 1) % self.
            number_of_devices()
4            if not self._participated:
5                self.medium().send(
6                    Vote(self.index(), nxt, self.index(),
                        False))
7                self._participated = True
8            ingoing = self.medium().receive()
9            if ingoing is not None:
10               if ingoing.vote() == self.index():
11                   if not ingoing.decided():
12                       self.medium().send(
13                           Vote(self.index(), nxt, self.
                            index(), True))
14                   else:
15                       self._leader = self.index()
16                       return  # this device is the new
                            leader
```

```
17               elif ingoing.vote() < self.index():
18                   continue
19               elif ingoing.vote() > self.index():
20                   # forward the message
21                   self.medium().send(
22                       Vote(self.index(), nxt, ingoing.vote(),
                            ingoing.decided()))
23               if ingoing.decided():
24                   self._leader = ingoing.vote()
25                   return
26           self.medium().wait_for_next_round()
```

## Properties
### Chang Roberts

- ▶ Safe
- ▶ Live
- ▶ $3N - 1$ messages per election

### Crashes?
Can be overcome if reliably detected!

# Bully Algorithm

8

## Idea

- ▶ Bully election requests into silence
- ▶ Priority by ID
- ▶ Basic algorithm in election
    - ▶ Send "shut up" to lower IDs
    - ▶ Higher IDs will "answer"
    - ▶ Finally, highest alive ID broadcasts itself

## NOTE!
**Depends on Synchronous Behavior** (timeout to detect crashes)

# Bully Algorithm - Code

```
1   class Bully(Device):
2       def __init__(...):
3           super().__init__(index, number_of_devices, medium)
4           self._leader = None
5           self._shut_up = False
6           self._election = False
7
8       def run(self):
9           first_round = True
10          while True:
11              got_input = False
12              if not self._shut_up and not self._election:
                    self.start_selection()
13              new_election = False
14              while True:
15                  ingoing = self.medium().receive()
16                  if ingoing is not None:
17                      got_input = True
18                      if ingoing.vote() < self.index():
19                          self.medium().send(
20                              Vote(self.index(), ingoing.source, self
                                .index(), self.largest()))
21                          new_election = True
22                      else:
23                          self._shut_up = True
24                          if ingoing.decided():
25                              self._leader = ingoing.vote()
26                              return
27                  else: break

                    if not self._shut_up and not self._election and              28
                        new_election:
                        self.start_selection()                                   29
                    if not got_input and not first_round:                        30
                        if self._election:                                       31
                            if self._shut_up:                                    32
                                self._shut_up = False                            33
                                self.start_election()                            34
                            else:                                                35
                                for id in self.medium().ids():                   36
                                    if id != self.index():                       37
                                        self.medium().send(Vote(self.index       38
    (), id, self.index(), True))
                                self._leader = self.index()                      39
                                return                                           40
                    self.medium().wait_for_next_round()                          41
                    first_round = False                                          42
                                                                                 43
        def largest(self):                                                       44
            return self.index() == max(self.medium().ids())                      45
                                                                                 46
        def start_election(self):                                               47
            if not self._election:                                               48
                self._election = True                                            49
                for id in self.medium().ids():                                   50
                    if id > self.index():                                        51
                        self.medium().send(Vote(self.index(), id, self           52
    .index(), self.largest()))
```

# Properties
Bully

- ▶ Safe & Live, assuming
    - ▶ Unique IDs
    - ▶ Failure detection is reliable
- ▶ Best-case: $N - 2$ messages pr election
- ▶ Worst-case: $O(N^2)$ messages (the process with the lowest ID starts the election)
- ▶ Election-time: 2 rounds (assuming HW multicast)

## Properties
Bully

11

### Beware
Safety is broken if

- ▶ too tight deadline,
- ▶ process IDs reappear, or
- ▶ system is not synchronous.

Requires

- ▶ **requires synchronous system** , and
- ▶ **ordering of messages**