# Hierarchical Leader Election Algorithm With Remoteness Constraint

## Mohamed Tbarka

## October 21, 2019

# Outline

# Outline

Mohamed Tbarka

Hierarchical Leader Election Algorithm With Remoteness Constraint

What's Distributed Systems ?

# What's Distributed Systems ?

A distributed system is a network that consists of autonomous computers that are connected using a distribution middleware. They help in sharing different resources and capabilities to provide users with a single and integrated coherent network.

Election Algorithms

# The Bully Algorithm

As a first example, consider the bully algorithm devised by Garcia-Molina (1982). When any process notices that the coordinator is no longer responding to requests, it initiates an election. A process, P, holds an election as follows:

Election Algorithms

# The Bully Algorithm

As a first example, consider the bully algorithm devised by Garcia-Molina (1982). When any process notices that the coordinator is no longer responding to requests, it initiates an election. A process, P, holds an election as follows:

- *P* sends an *ELECTION* message to all processes with higher numbers.

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|

Election Algorithms

# The Bully Algorithm

As a first example, consider the bully algorithm devised by Garcia-Molina (1982). When any process notices that the coordinator is no longer responding to requests, it initiates an election. A process, P, holds an election as follows:

- *P* sends an *ELECTION* message to all processes with higher numbers.
- If no one responds, P wins the election and becomes coordinator.

Election Algorithms

# The Bully Algorithm

As a first example, consider the bully algorithm devised by Garcia-Molina (1982). When any process notices that the coordinator is no longer responding to requests, it initiates an election. A process, P, holds an election as follows:

- *P* sends an *ELECTION* message to all processes with higher numbers.

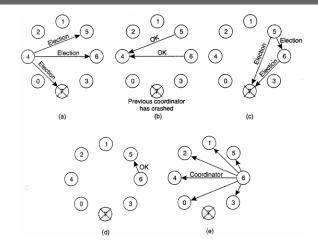- If no one responds, P wins the election and becomes coordinator.

- If one of the higher-ups answers, it takes over. P's job is done.

Election Algorithms



Figure: Bully Algorithm

# Outline

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|
| ○ | ○ ○ ○ | ● ○ ○ ○ | ○ ○ ○ ○ ○ ○ ○ | ○ | ○ ○ | ○ ○ |

System Model

# System Model

- We assume a system consisting of a set $P$ of computing nodes and a set $\chi$ of directed communication channels from one node to another node. $\chi$ consists of one channel.

System Model

# System Model

- We assume a system consisting of a set $P$ of computing nodes and a set $\chi$ of directed communication channels from one node to another node. $\chi$ consists of one channel.

- We model the whole system as a set of (infinite) state machines that interact through shared events (a specialization of the IOA model [17]).

Modeling Asynchronous Dynamic Links

# Asynchronous Dynamic Links' Model

The state of Channel(u, v), which models the communication channel from node u to node v, consists of:

- a $status_{uv}$ variable;

Modeling Asynchronous Dynamic Links

# Asynchronous Dynamic Links' Model

The state of Channel(u, v), which models the communication channel from node u to node v, consists of:

- a $status_{uv}$ variable;

- and a queue $mqueue_{uv}$ of messages.

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|

Configurations and Executions

# Configurations & Executions

- The notion of configuration is used to capture an instantaneous snapshot of the state of the entire system.

Configurations and Executions

# Configurations & Executions

- The notion of configuration is used to capture an instantaneous snapshot of the state of the entire system.
- A configuration is a vector of node states, one for each node in $P$, and a vector of channel states, one for each channel in $\chi$.

Problem Definition

# Problem Definition

Each node $u$ in the system has :

- a local variable $lid_u$ to hold the id of the supreme leader;

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|
| ○ | ○ ○ ○ ○ | ○ ○ ○ ○ ● ○ ○ | ○ ○ ○ ○ ○ ○ ○ | ○ | ○ ○ | ○ ○ |

Problem Definition

# Problem Definition

Each node $u$ in the system has :

- a local variable $lid_u$ to hold the id of the supreme leader;
- another local variable $slid_u$ to hold the identifier of the sub-leader whose remoteness towards $u$ obeys the constraint.

# Outline

Informal Description

# Heights

After a leader is gone, the algorithm consists on three waves:

- First wave : initiated by one of the lost leader's neighbors looking for it;

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|

Informal Description

# Heights

After a leader is gone, the algorithm consists on three waves:

- First wave : initiated by one of the lost leader's neighbors looking for it;

- Second wave : initiated by the node located at the edge of the network if the search has hit a dead-end;

# Heights

After a leader is gone, the algorithm consists on three waves:

- First wave : initiated by one of the lost leader's neighbors looking for it;

- Second wave : initiated by the node located at the edge of the network if the search has hit a dead-end;

- Third wave : initiated by the same node which initiated the first wave updating the other nodes' heights.

Outline          Introduction          Preliminaries          Leader Election Algorithm          Correctness          Implementation          Conclusions
○                ○                     ○                       ○                                 ○                     ○                       ○
                 ○                     ○                       ○                                                       ○                       ○
                 ○○                    ○                       ●
                                       ○                       ○
                                       ○                       ○
                                                               ○○

Nodes, Neighbors and Heights

# Heights

The height for each node is a 7-tuple of integers $((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of the causal clock time when the current search for an alternate path to the leader was initiated.

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|
| ○ | ○ | ○ | ○ | ○ | ○ | ○ |
| | ○ | ○ | ○ | | ○ | ○ |
| | ○○ | ○ | ● | | | |
| | | ○ | ○ | | | |
| | | ○ | ○ | | | |
| | | ○○ | ○○ | | | |

Nodes, Neighbors and Heights

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the
  node that started the current search (we assume node ids are
  positive integers).

Nodes, Neighbors and Heights

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of the causal clock time when the current search for an alternate path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the node that started the current search (we assume node ids are positive integers).
- $r$, a bit that is set to 0 when the current search is initiated and set to 1 when the current search hits a dead end.

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|

Nodes, Neighbors and Heights

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level $(RL)$ and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the
  node that started the current search (we assume node ids are
  positive integers).
- $r$, a bit that is set to 0 when the current search is initiated
  and set to 1 when the current search hits a dead end.
- $\delta$

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|

Initial State

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the
  node that started the current search (we assume node ids are
  positive integers).

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|
| ○       | ○            | ○             | ○                         | ○           | ○              | ○           |
|         | ○            | ○             | ○                         |             | ○              | ○           |
|         | ○○           | ○             | ●                         |             |                |             |
|         |              | ○             | ○                         |             |                |             |
|         |              | ○             | ○                         |             |                |             |
|         |              |               | ○○                        |             |                |             |

Initial State

## Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the
  node that started the current search (we assume node ids are
  positive integers).
- $r$, a bit that is set to 0 when the current search is initiated
  and set to 1 when the current search hits a dead end.

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|

Initial State

# Heights

The height for each node is a 7-tuple of integers $((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of the causal clock time when the current search for an alternate path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the node that started the current search (we assume node ids are positive integers).
- $r$, a bit that is set to 0 when the current search is initiated and set to 1 when the current search hits a dead end.
- $\delta$

Goal Of The Algorithm

# Heights

The height for each node is a 7-tuple of integers $((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of the causal clock time when the current search for an alternate path to the leader was initiated.

Goal Of The Algorithm

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the
  node that started the current search (we assume node ids are
  positive integers).

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|-------------|---------------|---------------------------|-------------|----------------|-------------|

Goal Of The Algorithm

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the
  node that started the current search (we assume node ids are
  positive integers).
- $r$, a bit that is set to 0 when the current search is initiated
  and set to 1 when the current search hits a dead end.

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
| --- | --- | --- | --- | --- | --- | --- |

Goal Of The Algorithm

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the
  node that started the current search (we assume node ids are
  positive integers).
- $r$, a bit that is set to 0 when the current search is initiated
  and set to 1 when the current search hits a dead end.
- $\delta$

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|-------------|---------------|---------------------------|-------------|----------------|-------------|

Description Of The Algorithm

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.

Description Of The Algorithm

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level $(RL)$ and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the
  node that started the current search (we assume node ids are
  positive integers).

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|

Description Of The Algorithm

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the
  node that started the current search (we assume node ids are
  positive integers).
- $r$, a bit that is set to 0 when the current search is initiated
  and set to 1 when the current search hits a dead end.

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|-------------|---------------|--------------------------|-------------|----------------|-------------|
| ○ | ○ | ○ | ○ | ○ | ○ | ○ |
|   | ○ | ○ | ○ |   | ○ | ○ |
|   | ○○ | ○ | ○ |   |   |   |
|   |   | ○ | ○ |   |   |   |
|   |   | ○ | ○ |   |   |   |
|   |   |   | ● |   |   |   |
|   |   |   | ○○ |   |   |   |

Description Of The Algorithm

# Heights

The height for each node is a 7-tuple of integers
$((\tau, oid, r), \delta, (nlts, lid), id)$, where the first three components are
referred to as the reference level ($RL$) and the fifth and sixth.

- $\tau$, a non-negative timestamp which is either 0 or the value of
  the causal clock time when the current search for an alternate
  path to the leader was initiated.
- $oid$, is a non-negative value that is either 0 or the id of the
  node that started the current search (we assume node ids are
  positive integers).
- $r$, a bit that is set to 0 when the current search is initiated
  and set to 1 when the current search hits a dead end.
- $\delta$

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|
| ○ | ○ | ○ | ○ | ○ | ○ | ○ |
|   | ○ | ○ | ○ |   | ○ | ○ |
|   | ○○ | ○ | ○ |   | ○ | ○ |
|   |   | ○ | ○ |   |   |   |
|   |   | ○ | ○ |   |   |   |
|   |   |   | ●○ |   |   |   |

Sample Execution

# The Code triggered by Update Message

**When node $u$ receives** $Update(h)$ **from node** $v \in forming \cup N$:

    // if $v$ is in neither *forming* nor $N$, message is ignored

1.   $height[v] := h$

2.   $forming := forming \setminus \{v\}$

3.   $N := N \cup \{v\}$

4.   $myOldHeight := height[u]$

5.   if $((nlts^u, lid^u) = (nlts^v, lid^v))$ // leader pairs are the same

6.       if (SINK)

7.           if $(\exists\ (\tau, oid, r)\ |\ (\tau^w, oid^w, r^w) = (\tau, oid, r)\ \forall\ w \in N)$

8.               if $((\tau > 0)$ and $(r = 0))$

9.                   REFLECTREFLEVEL

10.             else if $((\tau > 0)$ and $(r = 1)$ and $(oid = u))$

11.                  ELECTSELF

12.             else // $(\tau = 0)$ or $(\tau > 0$ and $r = 1$ and $oid \neq u)$

13.                  STARTNEWREFLEVEL

14.             end if

15.        else // neighbors have different ref levels

16.             PROPAGATELARGESTREFLEVEL

17.        end if

       // else not sink, do nothing

18.       end if

19.   else // leader pairs are different

20.       ADOPTLPIFPRIORITY $(v)$

21.   end if

22.   if $(myOldHeight \neq height[u])$

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|
| ○ | ○ ○ ○○ | ○ ○ ○ ○ ○ ○ | ○ ○ ○ ○ ○ ○ ○● | ○ | ○ ○ | ○ ○ ○ |

Sample Execution

# Subroutines

ELECTSELF
1.   $height[u] := (0,0,0,0, -\mathcal{T}_u, u, u)$

REFLECTREFLEVEL
1.   $height[u] := (\tau, oid, 1, 0, nlts^u, lid^u, u)$

PROPAGATELARGESTREFLEVEL
1.   $(\tau^u, oid^u, r^u) := max\{(\tau^w, oid^w, r^w)|\ w \in N\}$
2.   $\delta^u := min\{\ \delta^w\ |\ w \in N\ \text{and}\ (\tau^u, oid^u, r^u) = (\tau^w, oid^w, r^w)\} - 1$

STARTNEWREFLEVEL
1.   $height[u] := (\mathcal{T}_u, u, 0, 0, nlts^u, lid^u, u)$

ADOPTLPIFPRIORITY$(v)$
1.   if $((nlts^v < nlts^u)$ or $((nlts^v = nlts^u)$ and $(lid^v < lid^u)))$
2.        $height[u] := (\tau^v, oid^v, r^v, \delta^v + 1, nlts^v, lid^v, u)$
3.   else send $Update(height[u])$ to $v$
4.   end if

# Outline

Mohamed Tbarka

Hierarchical Leader Election Algorithm With Remoteness Constraint

# Outline

# What's JBotSim ?

# Outline

Outline    Introduction    Preliminaries    Leader Election Algorithm    Correctness    Implementation    Conclusions
○          ○               ○                ○                           ○             ○                 ○
           ○               ○                ○                                         ○                 ●
           ○○              ○                ○
                           ○                ○
                           ○                ○
                                            ○
                                            ○
                                            ○○

Where can I learn more?

## Questions and Answers

Want to know more?

- Browse `http://web.mit.edu/smoot/history.htm`.

# Questions and Answers

Want to know more?

- Browse http://web.mit.edu/smoot/history.htm.

- Smoot's Legacy http://alum.mit.edu/news/AlumniNews/
  Archive/smoots_legacy.

| Outline | Introduction | Preliminaries | Leader Election Algorithm | Correctness | Implementation | Conclusions |
|---------|--------------|---------------|---------------------------|-------------|----------------|-------------|
| ○ | ○ | ○ | ○ | ○ | ○ | ○ |
| | ○ | ○ | ○ | | ○ | ● |
| | ○○ | ○ | ○ | | | |
| | | ○ | ○ | | | |
| | | ○ | ○ | | | |
| | | | ○○ | | | |

Where can I learn more?

## Questions and Answers

Want to know more?

- Browse http://web.mit.edu/smoot/history.htm.

- Smoot's Legacy http://alum.mit.edu/news/AlumniNews/Archive/smoots_legacy.

- Smoot Salute!
  http://web.mit.edu/spotlight/smoot-salute.