

Hierarchical Leader Election Algorithm With Remoteness Constraint

Mohamed Tbarka

October 22, 2019

Outline

1 Introduction

- Definition Of A Distributed System
- Election Algorithms

2 Preliminaries

- System Model
- Modeling Asynchronous Dynamic Links
- Configurations and Executions
- Problem Definition

3 H. Leader Election Algorithm

- Informal Description
- Nodes, Neighbors and Heights
- Initial State
- Description Of The Algorithm

Outline

1 Introduction

- Definition Of A Distributed System
- Election Algorithms

2 Preliminaries

- System Model
- Modeling Asynchronous Dynamic Links
- Configurations and Executions
- Problem Definition

3 H. Leader Election Algorithm

- Informal Description
- Nodes, Neighbors and Heights
- Initial State
- Description Of The Algorithm

What's a Distributed System ?

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.

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:

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.

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.

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

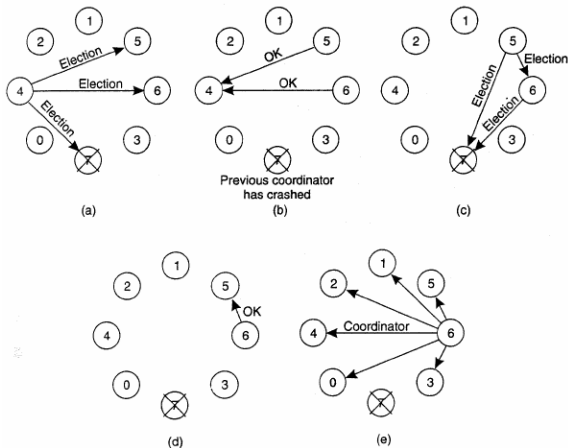


Figure: Bully Algorithm

Outline

1 Introduction

- Definition Of A Distributed System
- Election Algorithms

2 Preliminaries

- System Model
- Modeling Asynchronous Dynamic Links
- Configurations and Executions
- Problem Definition

3 H. Leader Election Algorithm

- Informal Description
- Nodes, Neighbors and Heights
- Initial State
- Description Of The Algorithm

System Model

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

System Model

- We assume a system consisting of a set P of computing nodes and a set χ of directed communication channels from one node to another node. χ 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).

Asynchronous Dynamic Links' Model

The state of $\text{Channel}(u, v)$, which models the communication channel from node u to node v , consists of:

- a status_{uv} variable;

Asynchronous Dynamic Links' Model

The state of $\text{Channel}(u, v)$, which models the communication channel from node u to node v , consists of:

- a status_{uv} variable;
- and a queue mqueue_{uv} of messages.

Configurations & Executions

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

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 χ .

Problem Definition

Each node u in the system has :

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

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

1 Introduction

- Definition Of A Distributed System
- Election Algorithms

2 Preliminaries

- System Model
- Modeling Asynchronous Dynamic Links
- Configurations and Executions
- Problem Definition

3 H. Leader Election Algorithm

- Informal Description
- Nodes, Neighbors and Heights
- Initial State
- Description Of The Algorithm

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;

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.

Nodes, Neighbors and 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 components are referred to as the leader pair (LP). In more detail, the components are defined as follows:

Nodes, Neighbors and 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 components are referred to as the leader pair (LP). In more detail, the components are defined as follows:

- τ , 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.

Nodes, Neighbors and 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 components are referred to as the leader pair (LP). In more detail, the components are defined as follows:

- τ , 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.

- 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.
- δ an integer that is set to ensure that links are directed appropriately to neighbors with the same first three components.

Initial State

- $forming_u$ is empty,
 N_u equals the set of neighbors of u in G_{chan}^{init}

Initial State

- $forming_u$ is empty,
 N_u equals the set of neighbors of u in G_{chan}^{init}
- $height_u[u] = (0, 0, 0, \delta_u, 0, l, u)$ where l is the *id* of a fixed node in u 's connected component in G_{chan}^{init} (the current leader), and δ_u equals the distance from u to l in G_{chan}^{init} ,

Initial State

- $forming_u$ is empty,
 N_u equals the set of neighbors of u in G_{chan}^{init}
- $height_u[u] = (0, 0, 0, \delta_u, 0, l, u)$ where l is the *id* of a fixed node in u 's connected component in G_{chan}^{init} (the current leader), and δ_u equals the distance from u to l in G_{chan}^{init} ,
- for each v in N_u , $height_u[v] = height_v[v]$ (i.e., u has accurate information about v 's height), and

Initial State

- $forming_u$ is empty,
 N_u equals the set of neighbors of u in G_{chan}^{init}
- $height_u[u] = (0, 0, 0, \delta_u, 0, l, u)$ where l is the *id* of a fixed node in u 's connected component in G_{chan}^{init} (the current leader), and δ_u equals the distance from u to l in G_{chan}^{init} ,
- for each v in N_u , $height_u[v] = height_v[v]$ (i.e., u has accurate information about v 's height), and
- \mathcal{T}_u is initialized properly with respect to the definition of causal clocks.

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.

- τ , 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.

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.

- τ , 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).

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.

- τ , 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.

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.

- τ , 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.
- δ

The Code triggered by Update Message

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

```

1. // if  $v$  is in neither forming nor  $N$ , message is ignored
   height[v] :=  $h$ 
2.  $forming := forming \setminus \{v\}$ 
3.  $N := N \cup \{v\}$ 
4.  $myOldHeight := height[u]$ 
5. if  $((nls^u, lid^u) = (nls^v, lid^v))$  // leader pairs are the same
6.   if (SINK)
7.     if  $(\exists (\tau, oid, r) \mid (\tau^w, oid^w, r^w) = (\tau, oid, r) \ \forall \ w \in N)$ 
8.       if  $((\tau > 0) \text{ and } (r = 0))$ 
9.         REFLECTREFLEVEL
10.      else if  $((\tau > 0) \text{ and } (r = 1) \text{ and } (oid = u))$ 
11.        ELECTSELF
12.      else //  $(\tau = 0)$  or  $(\tau > 0 \text{ and } r = 1 \text{ and } oid \neq u)$ 
13.        STARTNEWREFLEVEL
14.      end if
15.    else // neighbors have different ref levels
16.      PROPAGATELARGESTREFLEVEL
17.    end if
18.  // else not sink, do nothing
19.  end if
20. else // leader pairs are different
21.   ADOPTLIFPRIORITY( $v$ )
22. end if
23. if  $(myOldHeight \neq height[u])$ 
24.   send  $Update(height[u])$  to all  $w \in (N \cup forming)$ 
25. end if

```

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

STARTNEWREFLEVEL

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

ADOPTLPIFPRIORITY(v)

1. if $((nlts^v < nlts^u) \text{ or } ((nlts^v = nlts^u) \text{ and } (lid^v < lid^u)))$
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

1 Introduction

- Definition Of A Distributed System
- Election Algorithms

2 Preliminaries

- System Model
- Modeling Asynchronous Dynamic Links
- Configurations and Executions
- Problem Definition

3 H. Leader Election Algorithm

- Informal Description
- Nodes, Neighbors and Heights
- Initial State
- Description Of The Algorithm

Outline

1 Introduction

- Definition Of A Distributed System
- Election Algorithms

2 Preliminaries

- System Model
- Modeling Asynchronous Dynamic Links
- Configurations and Executions
- Problem Definition

3 H. Leader Election Algorithm

- Informal Description
- Nodes, Neighbors and Heights
- Initial State
- Description Of The Algorithm

What's JBotSim ?

JBoTSim is a java library that offers basic primitives for proto-typing, running, and visualizing distributed algorithms in dynamic networks.

Introduction



Preliminaries



H. Leader Election Algorithm



Correctness



Implementation



Conclusion



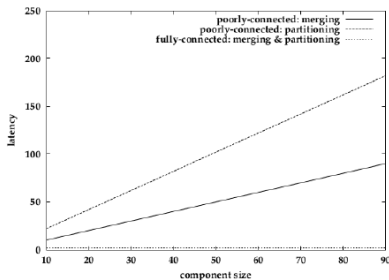
Simulation



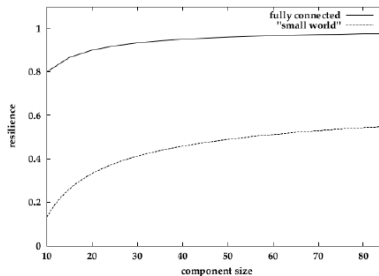
Mohamed Tbarka

Hierarchical Leader Election Algorithm With Remoteness Constraint

How is our algorithm's performance ?



(a) Latency



(b) Resilience

Figure: Simulation Results

Outline

1 Introduction

- Definition Of A Distributed System
- Election Algorithms

2 Preliminaries

- System Model
- Modeling Asynchronous Dynamic Links
- Configurations and Executions
- Problem Definition

3 H. Leader Election Algorithm

- Informal Description
- Nodes, Neighbors and Heights
- Initial State
- Description Of The Algorithm

Is The Algorithm Perfect ?

An open question is how to extend our algorithm and its analysis to handle a wider range of clocks, such as approximately synchronized clocks and vector clocks.

Is The Algorithm Perfect ?

Question ?