

Distributed Algorithms for Dynamic Networks

February 10, 2014

Welles Robinson

Senior Thesis

Department of Computer Science

Georgetown University

Abstract

Abstract

1 High-Level Description of Goal

The goal of this paper is to simplify the creation of distributed algorithms for dynamic networks by demonstrating that any algorithm that works for the broadcast variant of the synchronous model with a star topology can be made to work for the broadcast variant of the synchronous model with any topology. We will do so by describing a simulation algorithm that, if run on every node in the broadcast model, will match perfectly the output of the nodes of the centrally controlled model.

2 Models

Distributed Model

1. We are considering the broadcast variant of the synchronous model, defined with respect to a connected network topology $G=(V,E)$.
2. The broadcast variant is defined such that a given node has no knowledge of its neighbors but may send a single message per round that all of its neighbors will receive.
3. Nodes in the network have unique identification numbers (UIDs) and they have knowledge of their own UID.

3 Problem Definition

1. We are considering the main execution to be the distributed model as defined in the Models section with n processes.
 2. We are also considering a reference execution of the distributed model with a star topology and $n+1$ processes. The processes in this reference execution will run a given algorithm A and each process will keep an individual communication log, in which it will write its input and output for every round.
 3. The main execution will run some algorithm S , which will take as input A , the algorithm run on the reference execution. The processes in the main execution will keep individual communication logs, which do not have the restriction of those in the reference execution that they must contain the given process' input and output for every round.
 4. The problem is defined as solved if for the communication log of every process in the main execution, there is at least one identical communication log in the leaf processes of the reference execution and vice-versa.
1. The Turning-On Only Problem Definition
 2. A node may turn on at the beginning of any round of the main execution. The node should wait to perform any actions until the start of the next round of the reference execution.
 3. As a result, every node in the network must have some conception of when its reference execution round begins and ends so that it may transmit that information to nodes that turn on during the simulation.
1. The Turning-Off Only Problem Definition
 2. A node may turn off at the beginning of any round of the main execution (for now). This means that a node can essentially turn off during the middle of a single round of the reference execution.

4 Algorithm

1. Run unmodified Leader Election without Network Information
 - (a) Assume nodes have UIDs (which is also assumed for the distributed model)
 - (b) Each process runs Terminating Synchronous Breadth First Search and the node that manages to terminate elects itself leader and tells the other nodes to terminate
2. The Simulation

- (a) Assume that each and every round of the given algorithm has a finite repetitions of a step, which is comprised of two parts
 - i. The star node sends a broadcast message to all the leaf nodes (Broadcast part)
 - ii. All leaf nodes send a receive message to the star node in response to the broadcast message (Receive part)
- (b) Each step is simulated by the leader node running a modified instance of terminating synchronous BFS
 - i. The search message sent by a parent node to its children nodes is modified to be “search” plus the broadcast message, which is dictated by the output of the input algorithm A1
 - ii. Upon receiving a search message, a node runs the input algorithm A2 on the broadcast message portion and then creates the receive message portion using the output of A2
 - iii. The done message sent by a child node to its parent node is modified to be “done” plus both the UID and the receive message of the child node as well as any done messages received by the child node
 - iv. Upon receiving a search message, a node writes the broadcast message portion to its communication log
 - v. Upon sending a done message, a node writes the receive message portion to its communication log
 - vi. The leader simulates sending a search message to itself and then simulates sending a done message back to itself
 - vii. This algorithm terminates when the leader has received done messages from all of its children and a simulated done message from itself

Static Model - All the nodes turn on at the same time

Simulation Algorithm takes one input, algorithm A, the algorithm to be simulated. Algorithm A can be broken into two distinct algorithms, A1, the algorithm run by the star process, and A2, the algorithm run by the leaf processes.

Member Variables - maxID (UID); parent (UID); totalChildren (int); childCount (int); wait (int); Message has a root (a round, the UID); a id of the sender (UID); a type search, choose, done; a receiver (UID), defaults to NULL;

Dynamic Addition Model - Nodes turn on at various times but they don't turn off The leader elected will be the node with the highest UID out of all of the nodes that turned on at round 1

Variables - maxRoot - (a round, the UID); parent (UID); totalChildren (int); childCount (int); wait (int); Message has a root (a round, the UID); a id of the sender (UID); a type search, choose, done; a receiver (UID);

5 Proofs

Proof. of Leader Election in the Static Model

This algorithm will eventually elect a leader (see Liveness Proof)
This algorithm will never elect multiple leaders (see Safety Proof)

□

Proof. of Safety in the Static Model

Let U_{max} = the process with the maximum UID
For the breadth-first search instance with source = any node n not equal to U_{max} , it will send a search message to a node, n_1 , that is running a BFS instance, bfs_1 , with source $>n$
 n_1 will, by definition of the algorithm, respond with a search message.
Once n receives a search message from n_1 , it will cease running its current BFS instance and will start running the BFS instance of n_1 .
As a result, bfs_1 will never terminate because n , which was a member of the tree of bfs_1 , will never respond done to its parent and n will not elect itself leader.

□

Proof. of Liveness in the Static Model

Let U_{max} = the process with the maximum UID
After one round, every neighbor (defined as every process within broadcast range) of U_{max} will have received the ID of U_{max} and will, by definition, set its variable $maxID$ equal to ID of U_{max} and be running an instance of BFS with source = U_{max}
After round r , every process within r hops of U_{max} will be running an instance of BFS with source = U_{max} where r is the maximum number of hops of any node away from U_{max}
Now, every node in the network will be running an instance of BFS with source = U_{max}
Eventually, BFS will terminate and U_{max} will set its leader variable=true

□

Step 1: relate rounds of real execution to a single round in the reference execution
Simulation of a round of the reference execution for leader starts when the leader broadcasts its message according to the algorithm and ends when the leader has received messages from all of its children
Simulation of a round of the reference execution for a child node is between when the node broadcasts its message and when the node receives the message from the leader

Step 2: Define what it means for the real execution to correctly implement a round of the reference
Successful simulation of the reference execution means that the leader receives the messages from all of the nodes and the nodes receive the message from the leader node

Step 3: Assume up until Round R in the reference execution, everything in the main execution has matched everything in the reference execution; prove Round R in the reference execution is correctly simulated. So what happens in round R-1? The CNs all receive their message from the leader and the leader receives a message from every node in the network. In the first round of the rounds of the main execution that simulate round R, the leader broadcasts his message to his children and all of the leaf nodes broadcast their message, which is received by their parent. In the second round, the children of the leader broadcast his message, which is received by their children; the nodes who received messages from all of their children sends their messages to their parents.

6 Annotated Bibliography

@inproceedingsBrown et al:???, author = Brown, Matthew and Gilbert, Seth and Lynch, Nancy and Newport, Calvin and Nolte, Tina and Spindel, Michael, title = The Virtual Node Layer: A Programming Abstraction for Wireless Sensor Networks, booktitle = Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, year = ????, location = Cambridge, MA, USA, pages = 1-6, numpages = 6, publisher = MIT Computer Science and Artificial Intelligence Lab, address = Cambridge, MA, USA, keywords = Wireless ad hoc networks, network architecture, virtual infrastructure, annote = This paper addresses the problems around reliable coordination in dynamic, wireless networks. The paper proposes creating a static and reliable abstract layer composed of virtual nodes on top of the unpredictable and unreliable client or physical nodes. The architecture of the emulator is well thought-out and perhaps should be emulated. One key weakness of this approach is the requirement that the locations of the virtual nodes be decided before the start of program execution. This weakness results necessarily from the attempt to model an underlying dynamic network with a perfectly static network. ,

@inproceedingsFriedman:???, author = Friedman, Roy Vaysburd, Alexey, title = Fast Replicated State Machines Over Partitionable Networks, location = Ithaca, NY, USA, pages = 1-8, numpages = 8, publisher = Department of Computer Science Cornell University, keywords = replicated state machine approach, partitionable networks, annote = This paper lays out a faster and more reliable approach to implementing replicated state machines in partitionable networks. The approach is much faster than the pessimistic approach while sacrificing only a little fault-tolerance and much safer than the optimistic approach. The real relevance for my work is how this approach deals with the partitioning of the network, which will be an important in the fully dynamic case. The approach won't be overly relatable because it wants to only allow one partition to truly function and update state while my algorithm wants each partition to be able to function on its own. It requires that a majority of the nodes be in a single network to function, which guarantees the safety of the RSM. ,

@inproceedingsSchneider:1990, author = Schneider, Fred B., title = Implementing Fault-Tolerant Services Using the State Machine Approach: A Tutorial, booktitle = ACM Computing Surveys, series = Vol. 22, No. 4, year = 1990, location = Ithaca, NY, USA, pages = 299–319, numpages = 20, publisher = ACM, keywords = Algorithms, Design, Reliability, Client-server, distributed services, state machine approach, annote = Schneider lays out the replicated state machine approach to achieving fault-tolerance in a distributed system in this tutorial paper. This paper examines two types of faults: Byzantine failures and Fail-stop failures. The paper lays out the core requirements of reliable replica coordination, agreement and order, and techniques for meeting them including a "logical clock" and a synchronized real-time clock. This paper is a recap of much of my distributed systems class. ,

```

initVariables() ;
for round 1 ... r do
    for each message m in Inbox do
        if m.maxID > maxID then
            | updateMaxRoot() ;
        end
        if m.maxID == maxID then
            if m.type == choose AND receiver == myUID then
                | childCount++ ;
                | totalChildren++ ;
            end
            if m.type == done AND receiver == myUID then
                | childCount- ;
                if childCount == 0 then
                    | sendDoneMsg( ) ;
                end
            end
        end
        if m.maxID < maxID then
            | msg = (type=search, sender=myUID, maxID=maxID) ) ;
            | Outbox.enqueue(msg) ;
        end
    end
    if wait != 0 AND childCount == 0 then
        | wait- ;
        if wait == 0 then
            | sendDoneMsg( ) ;
        end
    end
    for each message m in Outbox do
        | broadcast(m);
    end
    myRound++ ;
end

```

Algorithm 1: Simulation Algorithm for Static Model

```

myRound == 0 ;
maxID = myUID ;
message m = (type=search, sender=myUID, maxID=maxID) ;

```

Algorithm 2: initVariables method

```

maxID = m.maxID;
parent = m.sender;
childCount = totalChildren = 0; msg1 = (type=choose, sender=myUID,
maxID=maxID, receiver=m.sender) ;
Outbox.enqueue(msg1) ;
msg2 = (type=search, sender=myUID, maxID=maxID) ;
Outbox.enqueue(msg2) ;
wait = 3 ;

```

Algorithm 3: updateMaxRoot method

```

msg = ( type=done, sender=myUID, maxID=maxID ) ;
Outbox.enqueue( msg ) ;

```

Algorithm 4: sendDoneMsg method

```

for each action in A1 do
  if leader == true then
    message = (type=r.action, sender=myUID) ;
    broadcast ( message ) ;
    msg = ( response to message ) ;
    commLog.write (response to message) ;
  end
  for each round r do
    for each message m in Inbox do
      if m.sender == parent then
        forwardMsgToChildren( ) ;
      end
      if m.receiver == myUID then
        add m to msg ;
        childCount- ;
        if childCount == 0 then
          Outbox.enqueue(msg) ;
        end
      end
    end
    for each message m in Outbox do
      broadcast(m) ;
    end
  end
end
end

```

Algorithm 5: Static Simulation Algorithm


```

msg = (response to m, receiver = parent) ;
commLog.write(response to m) ;
if totalChildren == 0 then
    | Outbox.enqueue(msg) ;
else
    | forwardMsg = (type = m.type, sender=myUID) ;
    | Outbox.enqueue(forwardMsg) ;
end

```

Algorithm 6: forwardMsgToChildren

```

initVariables() ;
for round 1...r do
    for each message m in Inbox do
        if m.root > maxRoot then
            | updateMaxRoot() ;
        end
        if m.root == maxRoot then
            if m.type == choose AND receiver == myUID then
                | childCount++ ;
                | totalChildren++ ;
            end
            if m.type == done AND receiver == myUID then
                | childCount- ;
                | if childCount == 0 then
                    | | sendDoneMsg() ;
                end
            end
        end
        if m.root < maxRoot then
            | msg = (type=search, sender=myUID, root=(r=maxRoot.r+1,
            | id=maxRoot.id) ) ;
            | Outbox.enqueue(msg) ;
        end
    end
    if wait != 0 AND childCount == 0 then
        | wait- ;
        | if wait == 0 then
            | | sendDoneMsg() ;
        end
    end
    for each message m in Outbox do
        | broadcast(m);
    end
    myRound++ ;
    maxRoot = (r=maxRoot.r+1, id=maxRoot.id) ;
end

```

Algorithm 7: Simulation Algorithm for the Dynamic Addition Model

```

myRound == 0 ;
maxRoot = (r=myRound, sender=myUID) ;
message m = (type=search, id=myUID, root=maxRoot) ;

```

Algorithm 8: initVariables method for Dynamic Addition Model

```

maxRoot = m.root;
parent = m.sender;
childCount = totalChildren = 0; msg1 = (type=choose, sender=myUID,
root=(r=maxRoot.r+1, id=maxRoot.id), receiver=m.id) ;
Outbox.enqueue(msg1) ;
msg2 = (type=search, sender=myUID, root=(r=maxRoot.r+1,
id=maxRoot.id) ) ;
Outbox.enqueue(msg2) ;
wait = 3 ;

```

Algorithm 9: updateMaxRoot method for Dynamic Addition Model

```

msg = ( type=done, sender=myUID, root=(r=maxRoot.r+1,
id=maxRoot.id), receiver=parent ) ;
Outbox.enqueue( msg ) ;

```

Algorithm 10: sendDoneMsg method for Dynamic Addition Model