

MODELS OF LEADER ELECTIONS AND THEIR APPLICATIONS

by

STEPHEN CURTIS JACKSON

A DISSERTATION

Presented to the Faculty of the Graduate School of the

MISSOURI UNIVERSITY OF SCIENCE AND TECHNOLOGY

In Partial Fulfillment of the Requirements for the Degree

DOCTOR OF PHILOSOPHY

in

COMPUTER SCIENCE

2016

Approved

Dr. Bruce McMillin, Advisor

Dr. One

Dr. Two

Dr. Three

Dr. Four

Copyright 2016

STEPHEN CURTIS JACKSON

All Rights Reserved

ABSTRACT

ACKNOWLEDGMENTS

TABLE OF CONTENTS

| | Page |
|--|------|
| ABSTRACT | iii |
| ACKNOWLEDGMENTS | iv |
| LIST OF ILLUSTRATIONS | ix |
| LIST OF TABLES | xiii |
| SECTION | |
| 1 INTRODUCTION | 1 |
| 2 BACKGROUND | 5 |
| 2.1 Distributed Systems | 5 |
| 2.2 Execution And Communication Models | 5 |
| 2.2.1 Communication Channels | 6 |
| 2.2.2 Clocks | 6 |
| 2.2.3 Execution | 7 |
| 2.3 Faults, Failures, and Errors | 7 |
| 2.3.1 Crash or Fail-Stop Failure | 7 |
| 2.3.2 Omission Failure | 8 |
| 2.3.3 Failure Detectors | 8 |
| 2.3.4 Byzantine Fault | 9 |
| 2.4 Probability | 9 |
| 2.5 Markov Models | 10 |
| 2.5.1 Discrete Time Markov Chain | 11 |

| | | |
|-------|---|----|
| 2.5.2 | Continuous Time Markov Chain | 13 |
| 2.5.3 | Hidden Markov Model | 14 |
| 2.6 | Information Flow | 14 |
| 2.6.1 | Modal Logic | 14 |
| 2.6.2 | Non-Deducible (MSDND) Security | 16 |
| 2.6.3 | BIT Logic | 19 |
| 2.7 | Explicit Congestion Notification | 19 |
| 2.7.1 | Random Early Detection | 20 |
| 2.8 | Distributed Grid Intelligence | 21 |
| 2.8.1 | Real Time | 21 |
| 2.8.2 | Group Management Algorithm | 23 |
| 2.8.3 | Power Management | 25 |
| 3 | PROBLEM STATEMENT AND MOTIVATION..... | 28 |
| 3.1 | Initial Experiments | 29 |
| 3.1.1 | Sequenced Reliable Connection | 29 |
| 3.1.2 | Sequenced Unreliable Connection | 30 |
| 3.2 | Initial Results | 31 |
| 3.2.1 | Sequenced Reliable Connection | 31 |
| 3.2.2 | Sequenced Unreliable Connection | 33 |
| 3.3 | Markov Models | 34 |
| 3.3.1 | Initial Model Calibration | 35 |
| 3.4 | Remarks | 37 |
| 4 | RELATED WORK..... | 39 |
| 4.1 | Analysis of Distributed Systems | 39 |
| 4.2 | Physical Faults Caused by Cyber Entities in CPS | 40 |

| | | |
|-------|--|----|
| 5 | INFORMATION FLOW ANALYSIS OF DISTRIBUTED COMPUTING ... | 42 |
| 5.1 | Methods | 42 |
| 5.2 | Two Armies Problem | 43 |
| 5.3 | Byzantine Generals | 46 |
| 5.3.1 | Example With Two Armies | 48 |
| 5.4 | Election in an Anonymous Complete Network | 51 |
| 6 | ALGORITHM AND MODEL CREATION | 53 |
| 6.1 | Execution Environment | 53 |
| 6.2 | Election Algorithm Overview | 53 |
| 6.3 | Algorithm | 55 |
| 6.4 | Model Construction | 60 |
| 6.5 | Highest Priority Process | 62 |
| 6.5.1 | State Determination | 62 |
| 6.5.2 | Memorylessness | 64 |
| 6.5.3 | Model Construction | 65 |
| 6.5.4 | Model Validation | 66 |
| 6.6 | Profile Chain Analysis | 68 |
| 6.7 | Lower Priority Process | 74 |
| 6.8 | Independent Algorithm | 75 |
| 7 | APPLICATION | 76 |
| 7.1 | Application: ECN Hardening | 76 |
| 7.1.1 | Usage Theory | 76 |
| 7.1.2 | Group Management | 77 |
| 7.1.3 | Cyber-Physical System | 80 |
| 7.1.4 | Relation To Omission Model | 82 |

| | | |
|-------|------------------------------|----|
| 7.1.5 | Calibration | 84 |
| 7.2 | Proof Of Concept | 85 |
| 7.2.1 | Experimental Setup | 85 |
| 7.2.2 | Results | 87 |
| 8 | CONCLUSION..... | 92 |
| | BIBLIOGRAPHY | 94 |
| | VITA..... | 99 |

LIST OF ILLUSTRATIONS

| Figure | Page |
|--|------|
| 2.1 Real Time Scheduler | 22 |
| 2.2 Example of a failed migration. (*) and (**) mark moments when power devices change state to complete the physical component of the migration. In this scenario, the message confirming the demand side made the physical is lost, leaving the supply node uncertain. | 26 |
| 2.3 Example of a failed migration. (*) marks a moment when power devices change state to complete the physical component of the migration. In this scenario, the supply process changes its device state, but the demand process does not. | 26 |
| 2.4 Each migration consumes excess generation capability and removes excess demand. | 27 |
| 3.1 Time in-group over a 10 minute run for a two process system with a 100ms resend time | 31 |
| 3.2 Time in-group over a 10 minute run for a two process system with a 200ms resend time | 31 |
| 3.3 Average size of formed groups for the transient partition case with a 100ms resend time | 32 |
| 3.4 Time in-group over a 10 minute run for the transient partition case with a 100ms resend time | 32 |
| 3.5 Average size of formed groups for the transient partition case with a 200ms resend time | 33 |
| 3.6 Time in-group over a 10 minute run for the transient partition case with a 200ms resend time | 33 |
| 3.7 Time in group over a 10 minute run for two process system with 100ms resend time | 34 |
| 3.8 Time in group over a 10 minute run for two process system with 200ms resend time | 34 |
| 3.9 Comparison of in-group time as collected from the experimental platform and the simulator (1 tick offset between processes). | 35 |

| | | |
|------|--|----|
| 3.10 | Comparison of in-group time as collected from the experimental platform and the simulator (2 tick offset between processes). | 35 |
| 3.11 | Comparison of in-group time as collected from the experimental platform and the time in group from the equivalent Markov chain (128ms between resends). | 36 |
| 3.12 | Comparison of in-group time as collected from the experimental platform and the time in group from the equivalent Markov chain (64ms between resends). | 36 |
| 5.1 | Example of a Markov chain constructed for an election algorithm from information flows. Each state represents the processes that have gone passive. | 51 |
| 6.1 | State machine of a leader election. Processes start as coordinators in the “Normal” state and search for other coordinators to join with. Processes immediately respond to Are You Coordinator (AYC) messages they receive. The algorithm was modified by adding a “Ready Acknowledgment” message as the final step of completing the election. Additionally, processes only accept invites if they have received an “AYC Response” message from the inviting process. | 54 |
| 6.2 | State machine of maintaining a group. The AYC messages are the same as those in Figure 6.1. AYC and Are You There (AYT) are periodically sent by processes, and responses to those messages are immediately sent by the receiving process. In the modified algorithm, the member does not enter the recovery state if they do not receive an AYT response before the timeout expires. | 55 |
| 6.3 | Steady state distribution for 3 processes as well as the average group size (AGS) as a fraction of total processes. | 70 |
| 6.4 | Steady state distribution for 4 processes as well as the AGS as a fraction of total processes. | 70 |
| 6.5 | Steady state distribution for 5 processes as well as the AGS as a fraction of total processes. | 71 |
| 6.6 | Steady state distribution for 6 processes as well as the AGS as a fraction of total processes. | 71 |

| | | |
|-----|--|----|
| 6.7 | Example DGI schedule. Normal operation accounts for a fixed number of migrations each time the load balancing module runs. Message delays reduce the number of migrations that can be completed each round. However, reducing the group size allows more migrations to be completed (because fewer messages are being exchanged) at the cost of flexibility for how those migrations are completed. | 72 |
| 6.8 | Average group size as a percentage of all processes in the system for larger systems. | 73 |
| 7.1 | Example of network queuing during Distributed Grid Intelligence (DGI) operation. DGI modules are semi-synchronous, and create bursty traffic on the network. When there is no other traffic on the network (solid line), the bursty traffic causes a large number of packets to queue quickly, but the queue empties at a similar rate. With background traffic (dashed line), the bursty traffic causes a large number of packets to be queued suddenly. More packets arrive continuously, causing the queue to drain off more slowly. When the background traffic reaches a certain threshold (dotted line), the queue does not empty before the next burst occurs. When this happens, messages will not be delivered in time, and the queue will completely fill. | 78 |
| 7.2 | Example of process organization used in this paper. Two groups of processes are connected by a router. | 79 |
| 7.3 | Plot of the maximum observed average queue size as a function of the overall background traffic. The polynomial estimate is $y = 22.70x^2 - 44.74x + 85.72$ | 87 |
| 7.4 | Plot of the queue size for a queue from switch A to the router when only the DGI generates traffic. | 88 |
| 7.5 | Detailed view of Figure 7.4. The left most peak is from Group Management, and the 3 smaller peaks are from power migrations. | 88 |
| 7.6 | Detailed view of the effect on queue size as other network traffic is introduced. Compared to Figure 7.5, the peaks are taller and wider. Background traffic causes the average queue size to be updated more frequently. | 89 |
| 7.7 | Detailed view of the effect on queue size as other network traffic is introduced. With no Explicit Congestion Notification (ECN) notifications, the peak from Group Management is much larger. The congestion is sufficient that the Group Management and Load Balancing modules are affected. | 89 |

- 7.8 Detailed view of the effect on queue size as other network traffic is introduced. In this scenario, the ECN notifications put Group Management into a maintenance mode that reduces its message complexity and switches Load Balancing to slower migration schedule, preventing undesirable behavior. 89
- 7.9 Detailed view of the effect on queue size as a large amount network traffic is introduced. Groups are unstable and processes occasionally leave the main group. Some migrations are lost due to queueing delays. 89
- 7.10 Effect on queue size as a large amount of network traffic is introduced. Hard notifications cause the groups to divide. As a result of the smaller groups, the group management and load balancing peaks are smaller than those in 7.9. No migrations are lost. 90
- 7.11 Count of lost migrations from all processes over time. Migrations are counted as lost until the second process confirms it has been completed. Without congestion management, a large number of migrations are lost. 90

LIST OF TABLES

| Table | Page |
|---|------|
| 2.1 Logical Statement Formulation Rules | 16 |
| 2.2 The Axiomatic System | 17 |
| 3.1 Error and correlation of experimental data and Markov chain predictions | 36 |
| 6.1 Summary of χ^2 tests performed. | 68 |
| 7.1 Summary of Random Early Detection (RED) parameters. Unspecified values default to the Network Simulator 3 (NS-3) implementation default value | 87 |

1. INTRODUCTION

The design of stochastic models of distributed systems has a long history as a challenging area of interest. Models of distributed systems have to deal with a number of factors. These factors include the various types of failure the system could experience, a lack of tightly synchronized execution, and a large complex state space when there are a high number of agents[24][5]. However, the concept of distributed systems plays a central role in many of the future visions for how critical infrastructure will operate. These critical infrastructures are physical networks whose operation are so vital that if those networks failed to operate correctly it would be highly detrimental to the population that rely on those systems. Cyber-physical systems (CPS) are the integration of computational systems with physical networks. Computational systems already play a critical role in most critical infrastructures, and as demands for security features, such as accessibility, increase distributed systems are becoming an increasingly favorable choice for the computational needs for these systems[53].

The Future Renewable Electric Energy Delivery and Management (FREEDM) center[39], an NSF funded ERC envisions a future power-grid where widely distributed renewable power generation and storage is closely coupled with a distributed system that facilitates the dispatch of power across those areas. Other systems like Vehicular Ad Hoc Networks (VANET)[33][43][25] and air traffic control systems[47][4] also propose similar control systems where many computers must cooperate to ensure both smooth operation and the safety of the people using those systems. As a consequence, ensuring that the computer systems that control those infrastructures behave correctly during fault conditions is critical, especially when those computer systems rely on their interaction with other computers to operate.

A robust CPS should be able to survive and adapt to communication network outages in both the physical and cyber domains. When one of these outages occurs, the physical or cyber components must take corrective action to allow the rest of the system to continue operating normally. Additionally, processes may need to react to the state change of some other process. Managing and detecting when other processes have failed is commonly handled by a leader election algorithm and failure detector.

In a smart-grid system, misbehavior during fault conditions could lead to critical failures such as a blackout or voltage collapse. In a VANET or air traffic control system, vehicles could collide, injuring passengers or destroying property. Additionally, since these systems are a part of critical infrastructure, protecting them from malicious entities is an important consideration.

This work was motivated by observations on the effects of lost messages on the group management module of the DGI used by the FREEDM smart-grid project. These original observations confirmed the need to explore more well defined models for the behaviors of CPS in order for them to better serve the people that use them.

We present a framework for reasoning about inferable state in the context of a distributed system. To do this, we exploit existing work in the field of information flow security. Information flow security has been used to reason about how attacks like STUXNET can manipulate operators beliefs while disrupting a system[21]. In particular, these approaches reason about how the operator in a STUXNET attack has no avenue to verify the reports from a compromised computing device. Using existing modal logic frameworks and using information flow security models[27][21][28], one can formally reason about where information that is not normally known to a domain can be inferred.

We will show in this work that in a system with the correct information flows, an agent in a distributed system can infer the state of other agents in the system. With this information, that agent can then construct a reasonable model of the system

to determine if the current behavior could lead to an undesirable situation with either the cyber or physical network.

Using this framework we present a leader election algorithm that can be modeled with a Markov chain for a known omission fault[16] rate. The presented algorithm maintains the Markov property for the observations of the leader despite omission faults. This approach to considering how a distributed system interacts during a fault condition allows for the creation of new techniques for managing a fault scenario in cyber-physical systems. In the context of FREEDM, these models produce expectations of how much time the DGI will be able to spend coordinating and doing useful work. Using these measures, the behavior of the control system for the physical devices can be adjusted to prevent faults, like blackouts and voltage collapse, in the physical network.

We also propose using existing schemes to detect communication network congestion and inform processes in a CPS of impending congestion. Processes act on this information to change their behavior in anticipation of message delays or loss. This behavior allows them to harden themselves against the congestion, and allows them to continue operating as normally as possible during the congestion. This technique involves changing the behavior of both the leader election[23] and physical device management algorithm during congestion.

To accomplish this, we extend existing networking concepts of RED, ECN[42], and ICMP source quench[6]. When a network device detects congestion, it notifies processes that the network is experiencing congestion and they should react appropriately. We demonstrate an implementation of the FREEDM DGI in a NS-3 simulation environment[15] with our congestion detection feature. The DGI operates normally until the simulation introduces a traffic flow that congests the network devices in the simulation. After congestion has been identified by the RED queuing algorithm, the DGI are informed. When the congestion notifications are introduced, the DGI

maintains configurations which they would normally be unable to maintain during congestion. Additionally, we show a greater amount of work can be done without the work causing unstable power settings to be applied.

2. BACKGROUND

2.1. DISTRIBUTED SYSTEMS

Distributed systems are a computing paradigm characterized by independence of computational units and no universal clock. Components in a distributed system may not directly share computational resources or memory. Instead, computers in a distributed system typically interact through a message passing interface.

As a result, distributed computing is a challenging area of research. In a distributed system, since processes do not share a universal clock, the ordering of messages and events needed to be carefully considered to ensure correct operation of the system. Additionally, since individual components fail, determining which components fail and how they affect the system as a whole is also difficult in a distributed system. Different types of failures can cause different kinds of information to be withheld or changed, disrupting those processes.

2.2. EXECUTION AND COMMUNICATION MODELS

In this work we consider a distributed system where no processes in the system share an address space. All processes must use a message passing interface in order to communicate with other processes to exchange information. As a result of the complexities of distributed systems, various execution models have been developed to define how the execution of a distributed system proceeds. Different execution models affect how easy it is to reason about the system's execution, the types of algorithms that can be executed, and how complicated they are to implement.

2.2.1. Communication Channels. We describe the avenues of communication between processes as channels. A channel is classified as reliable if it meets 3 axioms:

Axiom 1. *Every message sent by a sender is received by a receiver and every received message was sent by a sender in the system. [24]*

Axiom 2. *Every message has an arbitrary but not infinite propagation delay.[24]*

Axiom 3. *Every channel is a First In First Out (FIFO) channel. If process P sends a messages x and y to Q (in that order) then Q receives the messages in order (x then y).[24]*

These are often referred to as synchronous channels. In this work, however, channels are not assumed to be perfectly reliable. Instead, we respect Axiom 3, partially fulfill Axiom 1, and disregard Axiom 2. Therefore, we assume the following about communication channels in the systems modeled (replacing Axiom 1)

Axiom 4. *Every message received by a receiver was sent by a sender in the system.*

As well as Axiom 3. Without the constrained propagation delay from 2, this type of communication channel is typically referred to as an asynchronous channel.

The communication model can be synchronus or asynchronous. In the synchronous communication model, processes can only send a message when the receiving process is ready to receive it. Algorithmically, sending in a synchronous model is usually considered a “blocking” operation, meaning that once a process tries to send a message, it cannot proceed until the message is received. In this work, communication is asynchronous, meaning that a process does not wait for the successful delivery of a message. This is known as a non-blocking send.

2.2.2. Clocks. Clocks are considered synchronized if every clock in the system reads the same time. Since it is impossible for independent clocks to tick at the

same rate, weak synchronization is used to describe when clocks in the system have an upper bound on drift rate from each other.

2.2.3. Execution. In a system with synchronous processes, processes execute in lockstep. At each step a process executes its next available action. Synchronous execution requires tight organization of the processes executing the algorithm. In this work we generally rely on partially synchronous execution by processes. In the partially synchronous model, execution proceeds in rounds or phases. The start of each round or phase is synchronized between processes, using a synchronized clock.

2.3. FAULTS, FAILURES, AND ERRORS

Processes can encounter incorrect behavior or issues during execution. Errors, faults, and failures describe the severity and consequence of the issue.

Definition 1. *An error is a difference between what is considered “correct” output for a given component, and the actual output: an incorrect result.*

Definition 2. *A fault is the manifestation of an error in software, or an incident where a incorrect, step, or data definition is performed in a computer program.*

Definition 3. *Failure is the inability for a component or system to perform its required function or within its specified limits.*

2.3.1. Crash or Fail-Stop Failure. A crash failure or its more generalized form, a fail-stop failure, describes a failure in which a process stops executing. In general, this is considered to be an irreversible failure, since a process typically does not resume from a crashed state. There is a special category of crash failures, called napping failures, where a process will appear to have crashed for a finite amount of time before resuming normal operation. Crash failures are impossible to detect with

absolute certainty in an asynchronous system. Processes can be suspected by other processes through the use of challenge/response messages or heartbeat messages that allow a process to prove that it has not crashed yet. A system that can handle a crash failure is implied to be able to handle a fail-stop failure.[24]

The fail-stop failure has three properties [24]:

1. When a failure occurs the program execution is stopped.
2. A process can detect when another fail-stop process has failed.
3. Volatile storage is lost when the fail-stop process is stopped.

2.3.2. Omission Failure. An omission fault[16] occurs when a message is never received by the intended recipient. Omission faults can occur when the communication medium is unavailable or when the latency of message exceeds a timeout for its expected delivery. Protocols like TCP do not tolerate omission failures: a packet is resent until it is acknowledged by the receiver. If the acknowledgment never comes, the connection is closed. As a contrast, UDP assumes that any datagram could be lost and it is the responsibility of the programmer to handle missing datagrams appropriately. An omission failure can have the same observable affects as a napping failure in some situations.[24]

2.3.3. Failure Detectors. Failure detectors[11] (Sometimes referred to as unreliable failure detectors) are special class of processes in a distributed system that detect other failed processes. Distributed systems use failure detection to identify failed processes for leader election routines. Because it isn't possible to directly detect a failed process in an asynchronous system, there has been a wide breadth of work related to different classifications of failure detectors, with different properties. Some of the properties include[11]:

- Strong Completeness - Every faulty process is eventually suspected by every other working process.

- Weak Completeness - Every faulty process is eventually suspected by some other working process.
- Strong Accuracy - No process is suspected before it actually fails.
- Weak Accuracy - There exists some process is never suspected of failure.
- Eventual Strong Accuracy - There is an initial period where strong accuracy is not kept. Eventually, working processes are identified as such, and are not suspected unless they actually fail.
- Eventual Weak Accuracy - There is an initial period where weak accuracy is not kept. Eventually, working processes are identified as such, and there is some process that is never suspected of failing again.

One class of failure detectors, Omega class Failure detectors, are particularly interesting because of [26]. An eventual weak failure (weak completeness and eventual weak accuracy) detector is the weakest detector which can still solve consensus. It is denoted several ways in various works including $\diamond\mathcal{W}$ [11], \mathcal{W} [10] [12] and Ω (Omega) [26].

2.3.4. Byzantine Fault. Byzantine fault, causes processes in the distributed system send information that is incorrect or misleading to other processes. Constraints for detecting processing that exhibit Byzantine behavior is a famous result in distributed systems.[35]

2.4. PROBABILITY

Several concepts are useful for reasoning about the stochastic properties of a distributed system. The expected value represents the long-term average output of a probability distribution.

$$E[X] = x_1p_1 + x_2p_2 + \dots + x_kp_k \quad (2.1)$$

Conceptually, the expected value is the weighted average of the outcomes of some stochastic system.

Definition 4. *Availability is the probability a system or component is operational and accessible when required for use, denoted A .*

$$A = \frac{E[\text{uptime}]}{E[\text{uptime}] + E[\text{downtime}]} \quad (2.2)$$

Definition 5. *Reliability is the ability of a system or component, in specified conditions, to be able to perform its required functions for a specified period of time.*

$$R(t) = \Pr(T > t) = \int_t^\infty f(x)dx \quad (2.3)$$

Where $R(t)$ is the probability that the system or component functions up until at least time t and $f(x)$ is the probability density function for the component's survival. Unfortunately, there is a clash of the term for reliability in distributed systems and in reliability analysis. We try to constrain the discussion about reliability in the distributed systems sense to communication channels.

2.5. MARKOV MODELS

A Markov chain is a finite set of states $W = \{w_1, w_2, \dots, w_n\}$ and probabilistic transitions between those states. States in a Markov chain are mutually exclusive. In a Markov chain, when a system is some state w_i it has some probability of transitioning to some other state w_j at the next time-step. A Markov chain is a first order chain if the probability of transitioning from i to j does not depend on the

history of transitions that lead to state i . First order chains are described as having a memoryless or Markov property. This formalizes the independence of the next state from the history of previous states. The Markov property describes a Markov chain as a sequence of random variables X_1, X_2, X_3, \dots and states the value of X_{n+1} only depends X_n : [7]

$$\begin{aligned} \Pr(X_{n+1} = x \mid X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) \\ = \Pr(X_{n+1} = x \mid X_n = x_n). \end{aligned} \quad (2.4)$$

An ergodic Markov chain is a chain where it is possible, in some finite number of steps, to go from any state to any other state. A stationary, or time homogeneous, Markov chain is one where the transition probabilities do not change over time. In a stationary Markov chain, the n th visit to a state is indistinguishable from the $n+1$ th visit to a state.

2.5.1. Discrete Time Markov Chain. A Discrete Time Markov Chain (DTMC) is one where the transitions between states happen at discrete time steps. A DTMC with m states can be represented by a $m \times m$ matrix. For simplicity when creating the model, matrices in this work are 1-indexed. In a matrix P , the value of P_{ij} represents the probability of the transition from w_i to w_j . The matrix is row stochastic, meaning the sum of each row in the matrix is equal to one:

$$\sum_{i=1}^m P_{ij} = 1. \quad (2.5)$$

A useful companion to the transition matrix is a state distribution vector. While the transition matrix describes how system will transition between states, the state distribution vector describes the probability of observing a given state.

Definition 6. *A state distribution vector is an m -dimensional vector composed of probability of observing each state in the system at a given instant:*

$$[P_1 \quad P_2 \quad \dots \quad P_m]$$

Where P_i corresponds to the probability of observing state w_i .

A DTMC is a suitable model for a memoryless random process with a finite number states which is observed at fixed time intervals. By utilizing a Markov chain, a variety of statistical analyses can be performed on the modeled system. For example, a Markov chain with the stationary and ergodic properties can be analyzed for its steady state probabilities. The steady state is a state distribution vector that describes the probability a random observation of a long-running process will observe some state w_i . The steady state probability distribution vector can be found via a system of equations: [7]

$$0 \leq \pi_j \leq 1.0 \tag{2.6}$$

$$\sum_{j=1}^m \pi_j = 1.0 \tag{2.7}$$

$$\pi_j = \sum_{i=1}^m \pi_i p_{ij} \tag{2.8}$$

The computation of the steady state will be noted as *Steady()*. A Markov chain can also be used to predict what state a process will be in at some point in the future. Given an initial state and a number of time-steps a matrix operation will yield the likelihood of the process being in each state after the time interval has passed. The mean passage time, a measure of how many time-steps will pass before a process returns or arrives to some state, can also be calculated.

We model a leader election algorithm with a closed form representation of the behavior of the algorithm. This closed form representation is a profile Markov chain (noted as P). The profile Markov chain is validated against a chain generated from execution of the algorithm. The chain constructed from sampled data is known as a test chain (noted as T).

2.5.2. Continuous Time Markov Chain. Transitions in a Continuous Time Markov Chain (CTMC) depend on the amount of time spent in a given state. Let $X(s) = i$ indicate the model is in state i at time s . If the model is time homogeneous, then the probability of transitioning to state j only depends on the time spent in that state (t).

$$P\{X(s+t) = j | X(s) = i\} = P\{X(t) = j | X(0) = i\} \quad (2.9)$$

Each transition has some expected value or holding time which describes the amount of time before a transition occurs. CTMC do not have transitions that return to the same state since the expected value of the transition time describes how long the system remains in the same state. The probability distribution function (pdf) of the exponential distribution can be written as: [41]

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases} \quad (2.10)$$

As a result, the expected or mean value of an exponential distribution, is a function of the parameter λ : [41]

$$E[X] = \frac{1}{\lambda}. \quad (2.11)$$

When there are multiple possible transitions from a state, each with their own expected transition time, the expected amount of time in the state is [40]

$$\sum \lambda(x, y) = \sum \lambda p_{x,y} = \lambda(x) \quad (2.12)$$

where $\lambda(x, y)$ is the expected amount of time before state x transitions to state y . The expected time in a state ($\lambda(x)$) is related to the expected time for an individual transition ($\lambda(x, y)$) by a probability $p_{x,y}$. Each transition lends to an expected amount of time in the state.

2.5.3. Hidden Markov Model. A Hidden Markov Model (HMM) is a Markov chain where the state is not directly visible to an observer. Instead, the HMM outputs observations related to the underlying, hidden, state. The hidden chain is assumed to still meet the Markov property. A HMM is described by the notation $\lambda = (\Pi, P, B)$, where Π is the initial state distribution vector, P is the matrix of state transition probabilities, and B is the observation probability distribution. The B set maps each state in W to a probability distribution for an observation. Let y be some observation, where $y \in Y$, a finite set of discrete, possible values. B then is a $|W| \times |Y|$ matrix where $B_{ij} = \Pr(y_j|W_i)$, the probability of observing y_j given the system is in some hidden state s_i .

2.6. INFORMATION FLOW

2.6.1. Modal Logic. Kripke frames[34][9] play a critical role in this work. The essential concept of this work is that each global state of the distributed system at any given instant can be captured as a countably infinite set of propositional variables. A Kripke frame is a pair $\langle W, R \rangle$ [22] such that W is a set of possible worlds, such each world corresponds to a unique global state of the system. Each

element of R describes a binary relationship for how the described system can move from world to world as events occur in the described system.

In the case of a distributed system, a world could be described as one the possible combinations of values of all boolean state variables $S = \{s_0, s_1, \dots, s_n\}$ in that system. As execution occurs, messages, time, or events cause these variables to change. Each change in boolean variables corresponds to a relationship in R [36]. Therefore, a world w is one possible valuation of all the variables in S and a transition from w to another w' (with its own valuation) can be noted as wRw' . Without loss of generality, each relationship in R must result in the change of at least one variable in S . Additionally, the set of worlds is complete: every possible combination of state variable values is represented in the set of worlds. No relationship in R can lead to a world that does not exist.

Additionally, we can define a set of valuation functions, $\{V\}$. Each function $V_{s_x}^i(w)$ in V describes the value observed by an agent i of a boolean state variable s_x in some world w . If a valuation function for a particular state variable is not defined for an agent, that agent cannot determine the value of that state variable, and cannot determine the value of any logical statement based on that variable. In the case of a distributed system, this concept is analogous to the isolation of memory for each agent. For example, an agent i , cannot simply determine the value of a variable for agent j .

The combination of a Kripke Frame $\langle W, R \rangle$ and a set of valuation functions V is a Kripke model $K = \{W, R, V\}$ sometimes known as a modal model. The complete model describes all the possible worlds, the relation between those worlds and the information available in the domains of the system. A Kripke Model can be made into a Markov chain if a transition probability matrix $P : |W|x|W| \rightarrow [0, 1]$ is mapped to each transition in R . [31][48]

Table 2.1: Logical Statement Formulation Rules

1. if φ is a wff, so are $\neg\varphi$, $\Box\varphi$, and $\Diamond\varphi$.
2. if φ is a wff, so are $B_i\varphi$ and $\neg B_i\varphi$
3. if φ is a wff, so are $T_{i,j}\varphi$ and $\neg T_{i,j}\varphi$
4. if φ is a wff, so are $I_{i,j}\varphi$ and $\neg I_{i,j}\varphi$
5. if φ and ψ are both wff, so are $\varphi \wedge \psi$
6. if φ and ψ are both wff, so are $\varphi \vee \psi$

Let $\varphi \in \Phi_0$ be an atomic proposition in a set of countably many propositions. The set of well-formed formulas (wffs) as defined by the formulation rules in 2.2 is the least set containing Φ_0 . Additionally, we use the modal operator \Box as an abbreviation for $\neg\Diamond\neg\varphi$. The complete axiomatic system is outlined in 2.1. For the uninitiated, the modal box operator (\Box), “it is necessary that” states (in the case of $\Box\varphi$) that in every world w , φ is true. As its dual, the diamond operator (\Diamond) states, that it is not the case that in every world, φ is true.

2.6.2. Non-Deducible (MSDND) Security. In the domain of security there are a wide variety of aspects worth protecting in every system. These are grouped into the core security concepts of integrity, accessibility and privacy. Many traditional security approaches rely heavily on cryptography to provide privacy. However accidental information leakage can still occur, which compromises the privacy of the system. For cyber-physical systems, the leakage is difficult to prevent. Unlike their cyber counterparts, the actions taken by the physical components cannot be easily hidden from an observer. For example, a plane changing altitude or a car turning or changing speed cannot be hidden from an observer. Other, more complicated systems, like the power grid, have actions that are more difficult to observe, but a well motivated attacker can potentially collect critical information about the behavior of the cyber components with observations of the physical network[45].

Information Flow security models are invaluable for assessing what information, if any, is leaked by either the cyber or physical components of the CPS. There

Table 2.2: The Axiomatic System

Definition of logical and modal operators (abbreviations)

- D1. $\varphi \wedge \psi \equiv \neg(\neg\varphi \vee \neg\psi)$
- D2. $\varphi \oplus \psi \equiv (\varphi \vee \psi) \wedge \neg(\varphi \wedge \psi)$ (exclusive or)
- D3. $\varphi \rightarrow \psi \equiv \neg\varphi \vee \psi$
- D4. $\varphi \leftrightarrow \psi \equiv (\varphi \rightarrow \psi) \wedge (\psi \rightarrow \varphi)$
- D5. $\Diamond\psi \equiv \exists w \in W : w \vdash \psi$
- D6. $\Box\varphi \equiv \neg\Diamond\neg\varphi$
- D7. $B_i\varphi$ agent i believes the truth of φ
- D8. $I_{i,j}\varphi$ agent j informs i that $\varphi \equiv \top$
- D9. $T_{i,j}\varphi$ agent i trusts the report from j about φ

Axioms

- P. All the tautologies from the propositional calculus.
- K. $\Box(\varphi \rightarrow \psi) \rightarrow (\Box\varphi \rightarrow \Box\psi)$
- M. $\Box\varphi \rightarrow \varphi$
- A1. $\neg\Box\varphi \rightarrow \Box\neg\Box\varphi$
- A2. $\Diamond(\varphi \vee \psi) \rightarrow \Diamond\varphi \vee \Diamond\psi$
- A3. $\Box\varphi \wedge \Box\psi \rightarrow \Box(\varphi \wedge \psi)$
- B1. $(B_i\varphi \wedge B_i(\varphi \rightarrow \psi)) \rightarrow B_i\psi$
- B2. $\neg B_i\perp$
- B3. $B_i\varphi \rightarrow B_iB_i\varphi$
- B4. $\neg B_i\varphi \rightarrow B_i\neg B_i\varphi$
- I1. $(I_{i,j}\varphi \wedge I_{i,j}(\varphi \rightarrow \psi)) \rightarrow I_{i,j}\psi$
- I2. $\neg I_{i,j}\perp$
- C1. $(B_iI_{i,j}\varphi \wedge T_{i,j}\varphi) \rightarrow B_i\varphi$
- C2. $T_{i,j}\varphi \equiv B_iT_{i,j}\varphi$

Rules of Inference

- R1. From $\vdash \varphi$ and $\vdash \varphi \rightarrow \psi$ infer ψ (Modus Ponens)
- R2. $\neg(\varphi \wedge \psi) \equiv (\neg\varphi \vee \neg\psi)$ (DeMorgan's)
- R3. From $\vdash \varphi$ infer $\vdash \Box\varphi$ (Generalization)
- R4. From $\vdash \varphi \equiv \psi$ infer $\vdash \Box\varphi \equiv \Box\psi$
- R5. From $\vdash \varphi \equiv \psi$ infer $\vdash T_{i,j}\varphi \equiv T_{i,j}\psi$

are a number of information flow security models, all based off similar concepts. Typically, these models partition the system into two domains: the high security domain and the low security domain. However, the MSDND security model allows the system to be partitioned into any number of domains. The MSDND model has been used to describe how the STUXNET attack was able to hide its malicious behavior from the plant operators. The MSDND security model is expressed using modal logic to

determine what information in a domain is deducible to an observer in another domain. This model exploits the possible worlds of modal logic to determine if there are worlds where the value of a logical atom is deducible by someone outside the domain.

This information flow security model can be used to determine what an agent in a distributed system can determine about another agent. The exact specification of timing the distributed system becomes unnecessary as the modal model can express any combination of logical atoms in one of its worlds.[27][21][28]

The MSDND security model can be expressed as follows[21]. Consider a pair of state variables s_x and s_y which may or may not be in the same security domain. The value of s_x and s_y have a logical xor relationship: if s_x is true, s_y must be false. Given an agent i that does not have a valuation function for either of those two variables, the system is MSDND secure for that agent and pair of variables. Written formally:

$$\begin{aligned} MSDND = \exists w \in W : w \vdash \Box[(s_x \vee s_y) \wedge \neg(s_x \wedge s_y)] \\ \wedge [w \models (\neg V_x^i(w) \wedge \neg V_y^i(w))] \end{aligned} \quad (2.13)$$

Of particular interest is the special case where s_x and s_y are relation on the same wff: ($s_x = \varphi$ and $s_y = \neg\varphi$):

$$\begin{aligned} MSDND = \exists w \in W : w \vdash \Box[\varphi \oplus \neg\varphi] \\ \wedge [w \models (\neg V_\varphi^i(w))] \end{aligned} \quad (2.14)$$

In a system where the above logical relationship holds, the agent i cannot determine the value of s_x or s_y . However, if the relationship does not hold, there is some world where the agent can determine the value of s_x and s_y .

2.6.3. BIT Logic. Belief, Information transfer, Trust (BIT) was developed to formalize logic about belief and information transfer. BIT logic has typically been applied to distributed systems, but has also played roles in CPS security. The operations of the BIT logic allow formal definition of how entities pass information, and how they will act on the information passed to them. BIT logic utilizes several modal operators:

- $I_{i,j}\varphi$ defines the transfer of information directly from agent j to an agent i .
- $T_{i,j}\varphi$ defines trust an agent i has in a report from j that φ is true.
- $B_i\varphi$ defines the belief that an agent i has about φ . The actual value of φ is irrelevant: the agent i believes it to be true.

These operators allow reasoning about information transfer between entities. In the context of a distributed system, these operators allow the division of the actual state held by some agent i to what some other agent j believes agent i 's state is.

2.7. EXPLICIT CONGESTION NOTIFICATION

ECN is a technique for managing congestion in IP networks. When an ECN capable network device detects congestion, it can drop the packets or it can signal senders using flags in the packet headers that the network is congested. For a TCP application, the result of the dropped packets causes the slow-start congestion control strategy to reduce the rate packets are sent. A more advanced implementation, using ECN, sets specific bits in the TCP header to indicate congestion. By using ECN, TCP connections can reduce their transmission rate without re-transmitting packets.

UDP applications have not typically utilized ECN. Although the ECN standard has flags in the IPv4 header, access to the IPv4 header is not possible on most systems. Furthermore, there is not a “one size fits all” solution to congestion in UDP algorithms.

2.7.1. Random Early Detection. The RED queuing algorithm is a popular queuing algorithm for switches and routers. It uses a probabilistic model and an Exponentially Weighted Moving Average (EWMA) to determine if the average queue size exceeds predefined values. These values are used to identify potential congestion and manage it. This is accomplished by determining the average size of the queue, and then probabilistically dropping packets to maintain the size of the queue. In RED, when the average queue size avg exceeds a minimum threshold (min_{th}), but is less than a maximum threshold (max_{th}), new packets arriving at the queue may be “marked”. The probability a packet is marked is based on the following relation between p_b and p_a where p_a is the final probability a packet will be marked.

$$p_b = max_p(avg - min_{th})/(max_{th} - min_{th}) \quad (2.15)$$

$$p_a = p_b/(1 - count * p_b) \quad (2.16)$$

Where max_p is the maximum probability a packet will be marked when the queue size is between min_{th} and max_{th} and $count$ is the number of packets since the last marked packet. With RED, the probability a packet is marked varies linearly with the average queue size, and as a function and the time since the last packet was marked. If avg is greater than max_{th} , the probability of marking trends toward 1 as the average queue size approaches $2 * max_{th}$. In the event the queue fills completely, the RED queue operates as a drop-tail queue. In a simple implementation of the RED algorithm, marked packets are dropped.

2.8. DISTRIBUTED GRID INTELLIGENCE

The DGI is a smart grid operating system that organizes and coordinates power electronics. It also negotiates contracts to deliver power to devices and regions that cannot effectively facilitate their own needs. DGI leverages common distributed algorithms to control the power grid, making it an attractive target for modeling a distributed system. Algorithms employed by the DGI and grouped into modules, work in concert to move power from areas of excess supply to excess demand.

DGI utilizes several modules to manage a distributed smart-grid system. Group management, the focus of this work, implements a leader election algorithm to discover which processes are reachable within the cyber domain. Other modules provide additional functionality, such as collecting global snapshots, negotiating the migrations, and giving commands to physical components.

DGI is a real-time system; certain actions (and reactions) involving power system components need to be completed within a pre-specified time-frame to keep the system stable. It uses a round robin scheduler in which each module is given a predetermined window of execution which it may use to perform its duties. When a module's time period expires, the next module in the line is allowed to execute.

The DGI uses the leader election algorithm, "Invitation Election Algorithm," written by Garcia-Molina[23]. This algorithm provides a robust election procedure which allows for transient partitions. Transient partitions are formed when a faulty link inside a group of processes causes the group to divide temporarily. These transient partitions merge when the link becomes more reliable.

2.8.1. Real Time. Real-time requirements were designed to enforce a tight upper bound on the amount of time used creating groups, discovering peers, collecting the global state, and performing migrations.

To enforce these bounds, the real-time DGI has distinct phases which modules were allowed to use for all processing. Each module was given a phase which grants it a specific amount of processor time. Modules used this time to complete any tasks they had prepared. When the allotted time was up the scheduler changed context to the next module. This interaction is illustrated in Figure 2.1

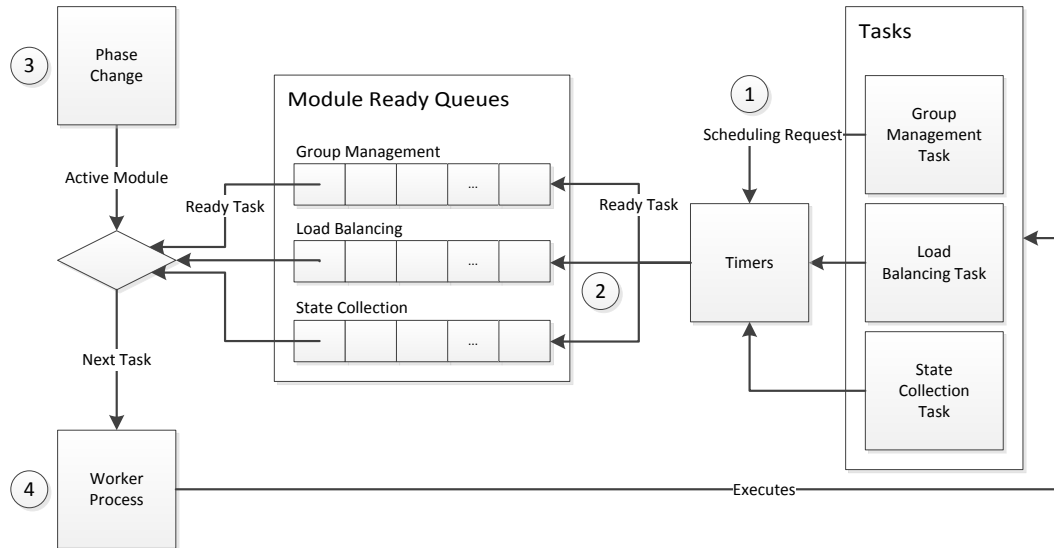


Figure 2.1: The real time scheduler used a round robin approach to allot execution time to modules.

1. Modules requested that a task be executed by specifying a time in the future to execute a task. A timer was set to count down to the specified moment. Modules could place tasks immediately into the ready queue if the task could be executed immediately.
2. When the timer expires the task is placed into the ready queue.
3. Modules were assigned periods of execution (called phases) which were a predetermined length. After the specified amount of time had passed, the module's phase ends and the next module in the schedule began to execute.
4. The worker selected the next ready task for the active module from the ready queue and executed it. These tasks could also schedule other tasks to be run in the future.

Modules informed the scheduler of tasks they wish to perform. The tasks could be scheduled for some point in the future, or scheduled to be executed immediately. When a task became ready it was inserted into a ready queue for the module it was scheduled for.

When that module's phase was active, tasks were pulled from the ready queue and executed. When the phase was complete, the scheduler stopped pulling tasks from the previous module's queue and began pulling from the next module's queue.

This allowed enforcement of upper bound message delay. Modules had a specific amount of processing time allotted. Modules with messages that invoked responses typically required the responses to be received within the same phase. Round numbers enforced that the message was sent within the same phase.

Modules were designed and allotted time to allow for parameters such as maximum query-response time (based on the latency between communicating processes). This implied that a module that engaged in these activities has an upper-bound in latency before messages are considered lost.

2.8.2. Group Management Algorithm. The DGI uses the leader election algorithm, "Invitation Election Algorithm," written by Garcia-Molina[23]. Originally published in 1982, this algorithm provides a robust election procedure that allows for transient partitions. Transient partitions are formed when a faulty link between two or more clusters of DGI causes the groups to divide temporarily. These transient partitions merge when the link becomes more reliable. The election algorithm allows for failures that disconnect two distinct sub-networks. These sub-networks are fully connected, but connectivity between the two sub-networks is limited by an unreliable link.

Since Garcia-Molina's original publication [23], a large number of election algorithms have been created. Each algorithm is designed to be well-suited to the circumstances it will be deployed in. Specialized algorithms exist for wireless sensor

networks[49][19], detecting failures in certain circumstances[52][37], and of course, transient partitions. Work on leader elections has been incorporated into a variety of distributed frameworks: Isis[8], Horus[44], Totem[38], Transis[3], and Spread[2] all have methods for creating groups. Despite this wide array of work, the fundamentals of leader election are consistent across all work. Processes arrive at a consensus of a single peer that coordinates the group. Processes that fail are detected and removed from the group.

The elected leader is responsible for making work assignments, and identifying and merging with other coordinators when they are found, as well as maintaining an up-to-date list of peers for the members of his group. Group members monitor the group leader by periodically checking if the group leader is still alive by sending a message. If the leader fails to respond, the querying nodes will enter a recovery state and operate alone until they can identify another coordinator. Therefore, a leader and each of the members maintain a set of processes which are currently reachable, a subset of all known processes in the system.

Leader election can also be classified as a failure detector[26]. Failure detectors are algorithms which detect the failure of processes within a system; they maintain a list of processes that they suspect have crashed. This informal description gives the failure detector strong ties to the leader election process. The group management module maintains a list of suspected processes which can be determined from the set of all processes and the current membership.

The leader and the members have separate roles to play in the failure detection process. Leaders use a periodic search to locate other leaders in order to merge groups. This serves as a ping / response query for detecting failures within the system. The member sends a query to its leader. The member will only suspect the leader, and not the other processes in their group.

Using a leader election algorithm allows the FREEDM system to autonomously reconfigure rapidly in the event of a failure. Cyber components are tightly coupled with the physical components, and reaction to faults is not limited to faults originating in the cyber domain. Processes automatically react to crash-stop failures, network issues, and power system faults. The automatic reconfiguration allows processes to react immediately to issues, faster than a human operator, without relying on a central configuration point. However, it is important the configuration a leader election supplies is one where the system can do viable work without causing physical faults like voltage collapse or blackouts[13].

A state machine for the election portion of the election algorithm is shown in Figure 6.1. In the normal state, the election algorithm regularly searches for other coordinators to join with. When another coordinator is identified, all other processes will yield to their future coordinator. The method of selecting which process becomes the coordinator of the new group differentiates the modified algorithm from other approaches.

2.8.3. Power Management. In this work we utilize the load balancing algorithm from [1]. The load balancing algorithm performs work by managing power devices with a sequence of migrations[51]. In each migration, a sequence of message exchanges identify processes whose power devices are not sufficient to meet their local demand and other processes supply them with power by utilizing a shared bus. To do this, first processes that cannot meet their demand announce their need to all other processes. Processes with devices that exceed their demand offer their power to processes that announced their need. These processes perform a three-way handshake. At the end of the handshake, the two processes have issued commands to their attached devices to supply power from the shared bus and to draw power from the shared bus. An example of how the power system is affected by migrations is depicted in Figure 2.4. In the chart, processes with net generation (generation >

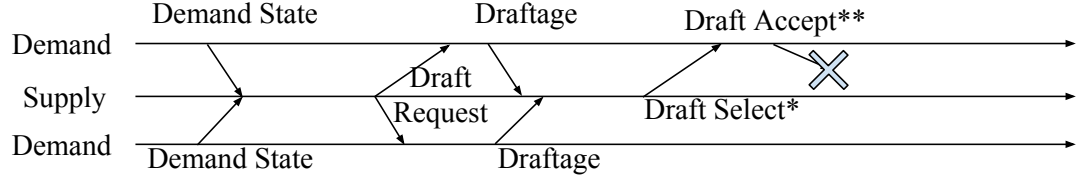


Figure 2.2: Example of a failed migration. (*) and (**) mark moments when power devices change state to complete the physical component of the migration. In this scenario, the message confirming the demand side made the physical is lost, leaving the supply node uncertain.

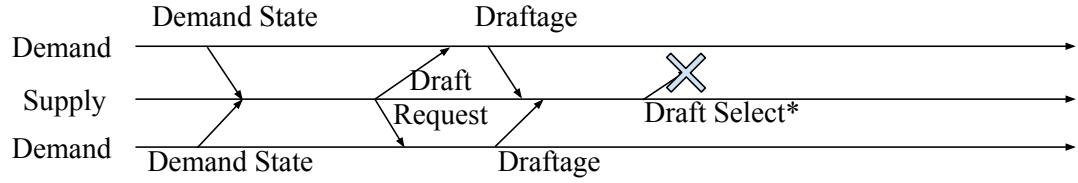


Figure 2.3: Example of a failed migration. (*) marks a moment when power devices change state to complete the physical component of the migration. In this scenario, the supply process changes its device state, but the demand process does not.

0) share power with processes with excess loads. As processes with excess loads are satisfied, both supply processes and demand processes trend toward 0 net generation.

The DGI algorithms can tolerate packet loss and is implemented using UDP to pass messages between DGI processes. Effects of packet loss on the DGI's group management module have been explored in [29] and [30]. The load balancing algorithm can tolerate some message loss, but lost messages can cause migrations to only partially complete, which can cause instability in the physical network. A failed migration is diagrammed in Figures 2.2 and 2.3. With this power migration algorithm, uncompensated actions may occur in the power system. These actions can eventually lead to power instability through issues such as voltage collapse. Additionally, the supply process may not always be certain if the second half of the action was completed or not. If the "Draft Accept" message does not arrive from the demand process, the supply process cannot be certain of whether or not its "Draft Select" message was received. If the supply process takes action to compensate by reversing

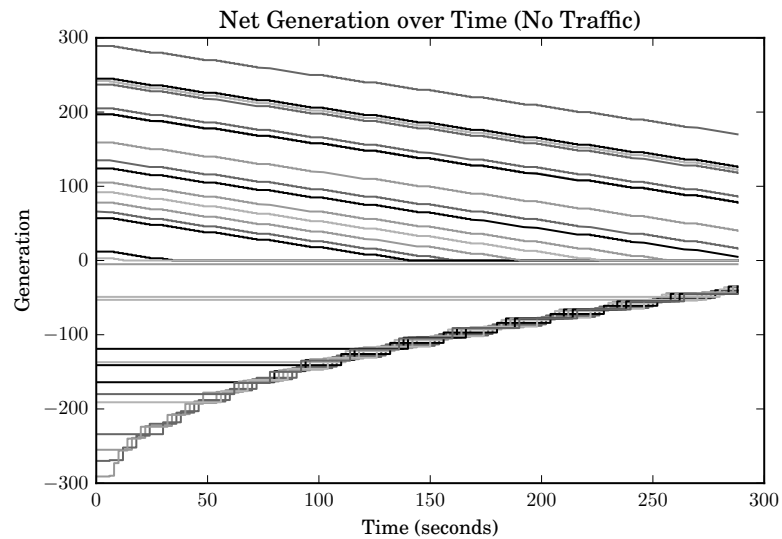


Figure 2.4: Each migration consumes excess generation capability and removes excess demand.

the migration and the confirmation arrives later the system will also be driven towards instability because another process completed an uncompensated action. Processes could confirm the number of failed migrations with a state collection technique. It is therefore desirable to manage the processes to minimize the number failed migrations.

3. PROBLEM STATEMENT AND MOTIVATION

The entry point to research in this area was asking the question “How does the software managing critical infrastructure like FREEDM behave when a critical component, the communication infrastructure, is not operating correctly?” Historically, leader elections have had limited applications in critical systems. However, in the smart grid domain, there is a great opportunity to apply leader election algorithms in a directly beneficial way. [1] presented a simple scheme for performing power distribution and stabilization that relies on formed groups. Algorithms like Zhang, et. al’s Incremental Consensus Algorithm[55], begin with the assumption that there is a group of processes who coordinate to distribute power. In a system where 100% up time is not guaranteed, leader elections are a promising method of establishing these groups. A strong cyber-physical system should be able to survive and adapt to network outages in both the physical and cyber domains. When one of these outages occurs, the physical or cyber components must take corrective action to allow the rest of the network to continue operating normally.

This work observes the effects of message omission on the group management module of the Distributed Grid Intelligence (DGI) used by the FREEDM smart-grid project. This system uses a broker system architecture to coordinate several software modules that form a control system for a smart power grid. These modules include: group management, which handles coordinating processes via leader election; state collection, a module which captures a global system state; and load balancing which uses the captured global state to bring the system to a stable state.

It is important for the designer of a cyber-physical system to consider what effects the cyber components will have on the overall system. Failures in the cyber domain can lead to critical instabilities which bring down the entire system if not

handled properly. In fact, there is a major shortage of work within the realm of the effects cyber outages have on CPS[50][53]. The analysis focuses on quantifiable changes in the amount of time a DGI could spend participating in energy management with other processes.

Using the DGI as a starting point, the analysis of the leader election algorithm in the DGI began with analysis of its behavior when messages were lost. To do this, the DGI software was subjected to omission faults while the state of the algorithm was captured over a period of examination. In the goal of these experiments was to examine what behavior the DGI would exhibit during the fault conditions. Additionally, we hoped to determine the advantages and disadvantages of using a more complicated, reliable message retransmission protocol over a more simple one without retransmission.

3.1. INITIAL EXPERIMENTS

The initial experiments were collected from a non-real time version of the DGI code. Experiments measured the time-in-group for various sets of DGI running the leader election algorithm during omission fault conditions. Additionally, the experiments examined two communication modes, the Sequenced Reliable Connection, and the Sequenced Unreliable Connection. Both methods of communication were valid approaches for the message passing requirements of the DGI.

3.1.1. Sequenced Reliable Connection. The sequenced reliable connection is a modified send and wait protocol with the ability to stop resending messages and move on to the next one in the queue if the message delivery time is too long. When designing this scheme we wanted to achieve several criteria:

- Messages must be accepted in order - Some distributed algorithm rely on the assumption that the underlying message channel is FIFO.

- Messages can become irrelevant - Some messages may only have a short period in which they are worth sending. Outside of that time period, they should be considered inconsequential and should be skipped. To achieve this, we have added message expiration times. After a certain amount of time has passed, the sender will no longer attempt to write that message to the channel. Instead, he will proceed to the next unexpired message and attach a “kill” value to the message being sent, with the number of the last message the sender knows the receiver accepted.
- As much effort as possible should be applied to deliver a message while it is still relevant.

There one adjustable parameter, the resend time, which controls how often the system would attempt to deliver a message it hadn't yet received an acknowledgment for.

3.1.2. Sequenced Unreliable Connection. The SUC protocol is simply a best effort protocol: it employs a sliding window to try to deliver messages as quickly as possible. A window size is decided, and then at any given time, the sender can have up to that many messages in the channel, awaiting acknowledgment. The receiver will look for increasing sequence numbers, and disregard any message that is of a lower sequence number than is expected. The purpose of this protocol is to implement a bare minimum: messages are accepted in the order they are sent.

Like the SRC protocol, the SUC protocol's resend time can be adjusted. Additionally, the window size is also configurable, but was left unchanged for the tests presented in this work.

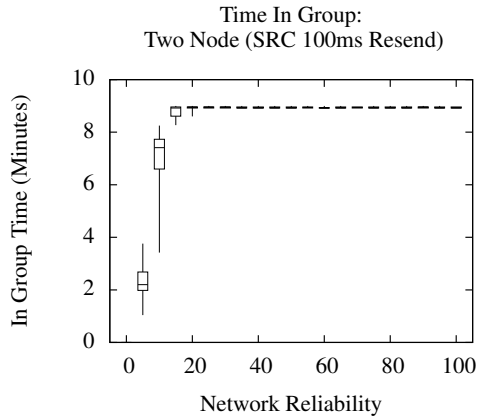


Figure 3.1: Time in-group over a 10 minute run for a two process system with a 100ms resend time

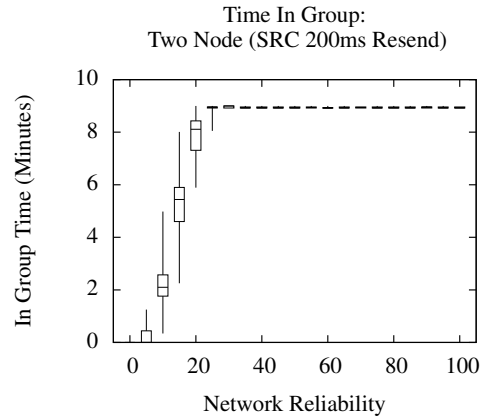


Figure 3.2: Time in-group over a 10 minute run for a two process system with a 200ms resend time

3.2. INITIAL RESULTS

The collected results from the tests were divided into several target scenarios as well as the protocol used.

The first minute of each test in the experimental test was discarded so that any transients in the test could be removed. The tests were run for ten minutes, however the maximum result was 9 minutes of in group time. These graphs first appeared in [29].

3.2.1. Sequenced Reliable Connection.

Two Node Case. The 100ms resend SRC test with two processes was considered a type of control in this study. These tests, pictured in Figure 3.1, highlighted the availability of the DGI with the SRC protocol. The maximum in group time of 9 minutes was achieved with only 15% of datagrams arriving at the receiver.

Figure 3.2 demonstrates that as the rate at which lost datagrams were re-sent was decreased to 200ms, the in-group time decreased. This behavior was expected. Each exchange had a time limit for each message to arrive and the number of attempts was reduced by increasing the resend time.

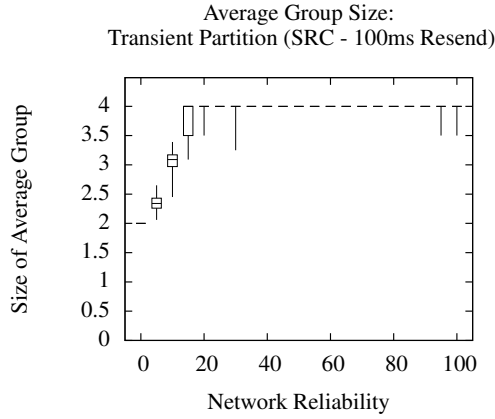


Figure 3.3: Average size of formed groups for the transient partition case with a 100ms resend time

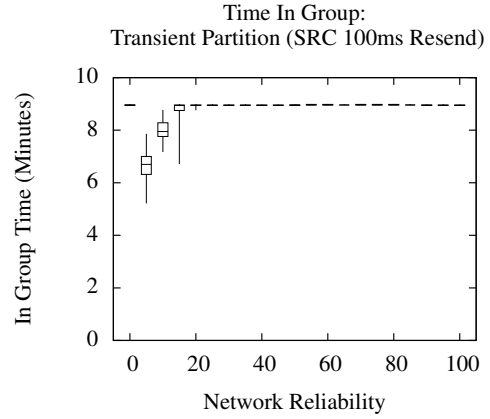


Figure 3.4: Time in-group over a 10 minute run for the transient partition case with a 100ms resend time

Transient Partition Case. The transient partition case was a simple example in which a network partition separates two groups of DGI processes. In the simplest case where the opposite side of the partition was unreachable, processes formed a group with the other processes on the same side of the partition. Two processes were present on each side of the partition. The 100ms case is shown in Figures 3.3 and 3.4.

While messages could not cross the partition, the DGIs stay in a group with the processes on the same side of the partition, leading to an in-group time of 9 minutes (the maximum value possible). As packets began to cross the partition (as the omission rate decreases), DGI instances on either side attempted to complete elections with the processes on the opposite partition and the in group time began to decrease. During this time, however, the mean group size continued to increase. Thus, while the elections were decreasing the amount of time that the module spent in a state where it can actively do work, it typically did not fall into a state where it was in a group by itself. As result, most of the lost in group time came from elections.

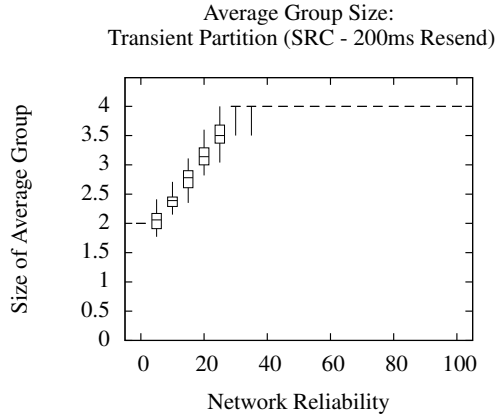


Figure 3.5: Average size of formed groups for the transient partition case with a 200ms resend time

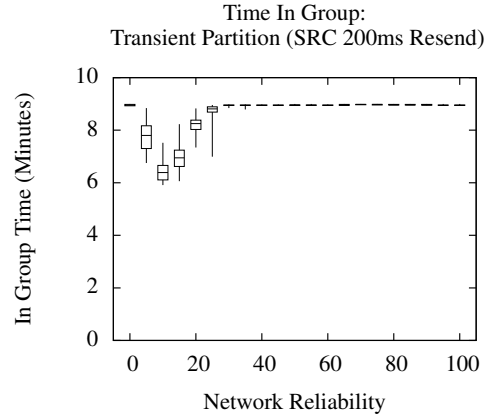


Figure 3.6: Time in-group over a 10 minute run for the transient partition case with a 200ms resend time

The 200ms case (Illustrated in Figures 3.5 and 3.6) suggests a similar behavior to Figures 3.3 and 3.4, with a wider valley produced by the reduced number of datagrams. The mean group size dipped below two in Figure 3.5, possibly because longer resend times allowed for a greater number race conditions between potential leaders. Discussion of these race conditions is shown and discussed during the SUC section since it was more prevalent in those experiments.

3.2.2. Sequenced Unreliable Connection.

Two Node Case. The SUC protocol's experimental tests revealed an immediate problem. There was a general increasing trend for the amount of time in-group shown in Figure 3.7. However, there was a high amount of variance for any particular trial.

In the 200ms resend case (illustrated in Figure 3.8), a greater growth rate occurred in the in group time as the omission rate decreased. When an average was taken across all of the collected data points from the experiment, the average in group time is higher for the 200ms case than it was for the 100ms case (6.86 vs 6.09). There large amount of variance in the collected in group time, however. As a result, it is

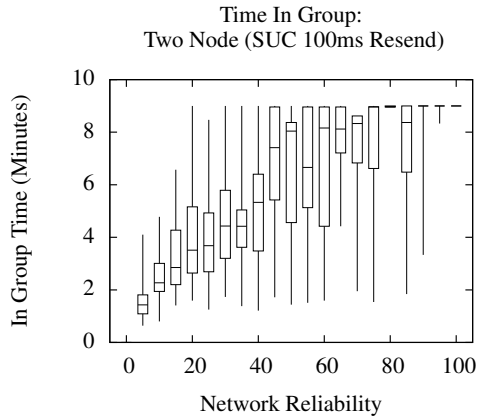


Figure 3.7: Time in group over a 10 minute run for two process system with 100ms resend time

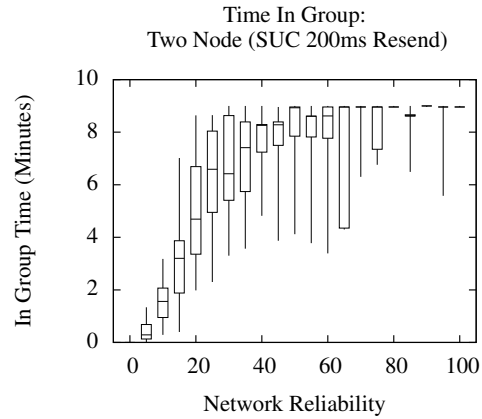


Figure 3.8: Time in group over a 10 minute run for two process system with 200ms resend time

not possible to state with confidence that there is a significant difference between the two cases.

3.3. MARKOV MODELS

After collecting the results from the initial experiments, we sought to describe the observed behavior through the use of continuous-time Markov chains. The behavior of the DGI transition between various states of grouping was calibrated with initial results and applied to other scenarios to validate the results. This approach had several shortcomings. First, for reasons we will demonstrate in subsequent chapters, the leader election algorithm modeled in these chains was not memoryless: the state used in the Markov chain was not sufficient to capture the interaction between processes that was occurring. Secondly, the continuous time model could not accurately capture the execution model of the DGI. In the DGI processes synchronize with each-other to execute in a semi-synchronous manner. As a result, the execution and transition between states is more accurately described with a discrete-time Markov chain.

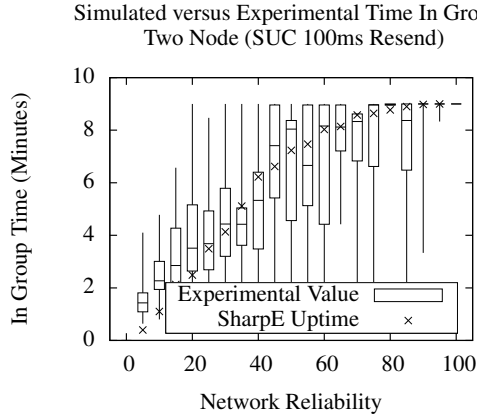


Figure 3.9: Comparison of in-group time as collected from the experimental platform and the simulator (1 tick offset between processes).

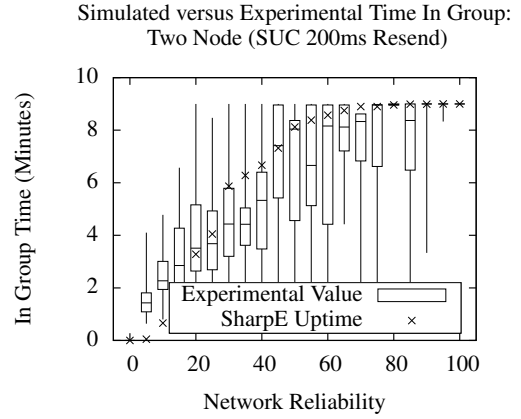


Figure 3.10: Comparison of in-group time as collected from the experimental platform and the simulator (2 tick offset between processes).

3.3.1. Initial Model Calibration. The presented methodology of constructing the model was initially calibrated against the original two-process case. This calibration used a non-real-time version of the DGI code. The resulting Markov chain was processed using SharpE[32][46], a popular tool for reliability analysis. SharpE measured the reward collected in 600 seconds, minus the reward that was collected in the first 60 seconds. Discarding the reward from the first 60 seconds emulated the 60 seconds were discarded in the experimental runs. The SharpE results are plotted along with the experimental results in Figures 3.9 and 3.10.

The race condition between processes during an election is a consideration in the original leader election algorithm, and is an additional factor here. The simulator provided a parameter to allow the operator to select how closely synchronized the peers were. This synchronization was the time difference between when each of them would search for leaders. The exchange of messages, particularly during an election, had a tendency to synchronize processes during elections. Nodes could synchronize even if they did not initially begin in a synchronized state. The simulation results

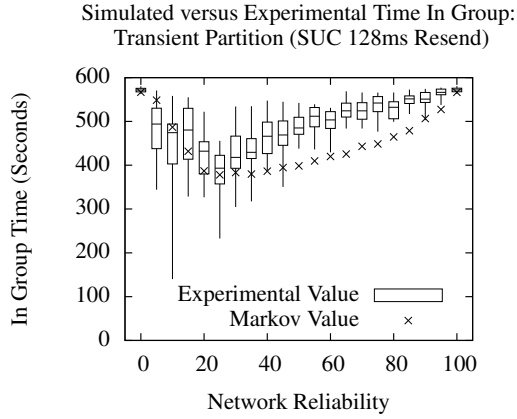


Figure 3.11: Comparison of in-group time as collected from the experimental platform and the time in group from the equivalent Markov chain (128ms between resends).

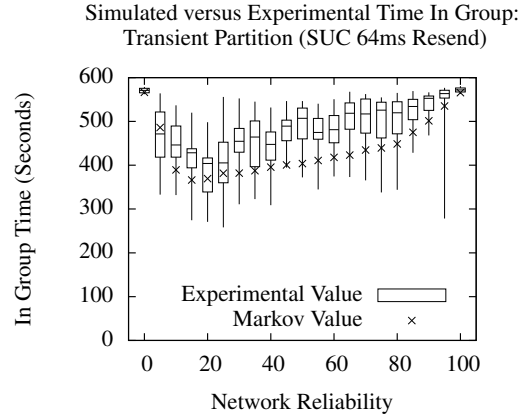


Figure 3.12: Comparison of in-group time as collected from the experimental platform and the time in group from the equivalent Markov chain (64ms between resends).

Table 3.1: Error and correlation of experimental data and Markov chain predictions

| Re-send | Correlation | Error |
|---------|-------------|--------|
| 128 | 0.7656 | 11.61% |
| 64 | 0.8604 | 11.70% |

aligned best for the 100ms resend case with 1 tick (Approximately 100ms difference in synchronization between processes) and 2 ticks (Approximately 400ms) in the 200ms resend case.

The structure of the Markov Chain assumed that processes enter the election state simultaneously. This was an appropriate assumption for the real-time system, since the round-robin scheduler synchronized when processes ran their group management modules. The simulator was set to assume that the synchronization between processes was very tight. New experimental data was collected for the 4 process, transient partition case. The collected data is overlaid with the results from the random walker in Figures 3.11 and 3.12.

As a measure of the strength of the model, the correlation between the predicted value was compared. The average error was also computed for each of the samples taken. This information is presented in Table 3.1. These results were not sufficient to accurately describe the behavior of the system during fault conditions. As a result, we sought to refine the analysis of the model in order to get an accurate portrayal of the behavior of the system during faults.

3.4. REMARKS

In order to ensure that critical infrastructure safely and reliably provides services to those that need the services, it is necessary to understand how the infrastructure behaves during faults. Conceptually, using unvetted critical infrastructure, especially when the infrastructure relies heavily on communication, is the same as stepping into an elevator without a safety brake. Strong analysis of a system's behavior during a failure scenario is paramount to ensuring the safety of those using the infrastructure. We propose knowledge of the correctness of the operation may be more important than the efficiency of the operation of that critical infrastructure. Through this work we demonstrate how distributed algorithms can be improved by reasoning about them using information flow security models. By using these models one can analyze where the certainty of the system is placed, as well as ensuring that the algorithms can be modeled without complete and perfect information of the system.

Additionally, we wish to show that these models can be used to allow the models to go beyond use by those administering the infrastructure. The models created can be analyzed and used while the system is working to provide services ("online"). Using the models created in this way a distributed critical infrastructure system can adjust its behavior, hardening itself against failures that could arise. With

this hardening technique the algorithms can prevent critical failures which could decrease quality of life or endanger that life. Subsequent chapters show how information flow methods can be used to determine the memorylessness of aspects of a distributed system, how those aspects can be used to construct a model, and applications of those models.

4. RELATED WORK

4.1. ANALYSIS OF DISTRIBUTED SYSTEMS

[17] focused on examining the emergent behavior as a collection of processes in a grid computing scenario work to process a large dataset. In their work, the authors focus on analysis derived from a discrete time Markov chain describing a single process completing a task. This chain describes how the process in the grid based system goes through the steps of the process acquiring a task to do, working on that task and subsequently either completing the task or failing. The created chains were “absorbing” chains, meaning that it had one or more states that the process could not leave once it had arrived in those states (task completed and task failed, for obvious reasons). Their analysis used minimal s-t cuts to determine critical paths for the ideal operation of the system. Using this analysis, they considered the Markov chain as a max-flow min-cut problem using the task complete or task failed absorbing states as sinks. By identifying the “critical transitions” of the graph, the areas that would most greatly affect the performance of the system could be identified.

[26] studies an Omega class failure detector using OmNet++[14], a network simulation software package. Instead of omission failures, however, it considers crash failures. Each configuration goes through a predefined sequence of crash failures. OmNet++ is used to count the number of messages sent by each of three different leader election algorithms. Additionally, [26] only considers the system to be in a complete and active state when all participants have consensus on a single leader.

4.2. PHYSICAL FAULTS CAUSED BY CYBER ENTITIES IN CPS

Faults in CPS can originate from many locations. First and most obviously, the traditional physical system being augmented by the CPS is subject to its own failures, either from component failure or the actions of an attacker. Secondly, the CPS must employ sensors to detect the state of the physical components in the system. Again, like the physical components, these sensors are subject to component failure or the actions of an attacker. Lastly, if the CPS communicates between entities using a communication network, the network can be disrupted by any number of issues, including DDoS, attackers, or congestion caused by other users in a shared network.

Several works have shown[45][13][50], that for a computer controlled smart grid, that failures originating at sensors or the communication network have the potential to cause the cyber entities controlling the physical networks to take incorrect actions. Work has been undertaken to identify faulty sensing components in a network, but the identification of bad sensing equipment may not always be possible. Additionally, it is not always possible to identify if the origin of an issue is a faulty sensor or an outside attacker.

If the cyber entity itself has been compromised, it could potentially exhibit Byzantine behavior, causing it to try and trick other components into bad actions[45], or it may try to disrupt the physical network directly. Work has been done to try and identify when an entity in a cyber network is actively working to compromise the physical network by using the underlying physical invariants. However, even if a cyber entity is trying to behave correctly, disruptions to the communication network or the sensors it uses can cause it to take actions similar to a process that is actively attempting to destabilize the system. As before, these actions can be identified by using the underlying invariants of the physical network after they have occurred.

However, it would be ideal to avoid situations where a “good” entity is forced to act badly.

5. INFORMATION FLOW ANALYSIS OF DISTRIBUTED COMPUTING

5.1. METHODS

A model created for a distributed system must have sufficient information to create an accurate model. This information is difficult to obtain because of the circumstances many distributed systems run on. Without exact synchronization, an accurate global snapshot of the system cannot be taken. Instead of attempting to capture exact global snapshots, our approach relies on allowing an agent to reason about the state of the other agents in the system in order to allow that agent to construct a model which can then be distributed to other agents.

To do this we propose the following structure for the execution environment of the distributed system: Each agent has some set logical atoms which it manipulates as its algorithms execute. Each agent belongs to domain unique to that agent (agent i is the only member of logical domain D_i) No agent can directly access a logical atom outside of its domain. The authenticity of any information transfer (using modal operator $I_{i,j}$) is never not trusted. However, the Trust ($T_{i,j}$) operator is still used to describe a message that is lost in transit: in all logical formulas presented the Trust operator describe this circumstance. Agents do not exhibit Byzantine failure, nor do they crash, only messages may be omitted.

If no information is passed between agents, they are MSDND secure (ignoring any sort of leakage from interactions in the physical world). As information is passed, aspects of the agents state are leaked. However, depending on when messages are sent, the agent can be left in doubt as to the state of the other agent. This has a common analogy to the two armies problem.

In order to create algorithms that can be modeled with a Markov chain, first the algorithm must allow at least one agent participating in the algorithm to determine an “image” of the current state of the algorithm. This image is descriptive enough to allow the agent to determine the probabilistic likelihood at arriving at the next image of the state of the algorithm. Secondly, to create modelable algorithms, we restrict it to a class of algorithms where the sequences of worlds that lead from one image to the next is equally likely.

5.2. TWO ARMIES PROBLEM

First, we will show that information flow analysis can be used to determine what state information is deducible to a particular agent in the system. To begin, we use the common two-armies problem to begin an analysis into what states can be determined in a distributed system.

In the two armies problem, two agents, which are generals of their respective armies must cooperate to attack an enemy city. However, the two armies are physically separated by the enemy city and must send messengers to each other to coordinate their plan. If the generals do not make an agreement on the attack, the attack will fail. However, the generals must come to an agreement when their channel for communication (a messenger) is unreliable (they can be captured by the enemy).

After one message has been sent, to one of the two generals, state of the other is MSDND secure. Let φ_0 be a logical atom that indicates “General A will attack at dawn.” For simplicity, we assume that after the time to attack is decided by General A, the agents will not reconsider the plan.

Theorem 1. *If no messages are exchanged, the state of the two generals is mutually MSDND secure.*

Proof. Proof: If no messages are exchanged and no information is leaked from the physical world, the two generals have no way of determining the other's state. \square

Theorem 2. *Once at least one messenger delivers a message to one of the Generals, one of the generals is not MSDND secure.*

Proof. Proof: Let $\{\varphi_i : i \in 1, 2, \dots, n\}$ describe the state that a general has received φ_{i-1} .

Case 2.1. *One messenger is sent by General A and arrives at General B.*

If no confirmations are sent, then General A clearly cannot deduce if General B has received the message and to General A, then to A, B is MSDND secure because it has no way to valuate $B_B\varphi_0$. However, to B, if B believes A's message then A is not MSDND secure to B, because B believes that φ_0 is true.

- | | | |
|----|--|---|
| 1. | φ_0 | General A decides to attack at dawn. |
| 2. | $I_{B,A}\varphi_0$ | General A sends a messenger to B informing them of their army's intent. |
| 3. | $B_B I_{B,A}\varphi_0 \wedge T_{B,A}\varphi_0$ | General B believes the message from general A. |
| 4. | $B_B\varphi_0$ | By C1. |
| 5. | $B_B\varphi_0 \rightarrow \varphi_1$ | General B knows the plan. |
| 6. | $w \models V_{\varphi_0}^B(w)$ | $V_{\varphi_0}^B(w)$ always returns true. |

Therefore, A is not MSDND secure to B. However, $V_{\varphi_1}^A(w) \notin V$, so B is secure to A.

However, to agent A, $V_{\varphi_0}^A(w)$ is undefined, so MSDND holds in that security domain.

Case 2.2. *Any number of messengers are sent and deliver their message, alternating from General A or General B to the other general.*

As each messenger arrives, the receiving general will trust the message and believe, resulting in that general assigning value to φ_i .

In the case $n = 1$, B is now unsure that A has received φ_1 and cannot deduce if $B_A\varphi_1$. B is unsure of A's state and as a consequence A is MSDND secure to B.

- | | | |
|-----|---|--|
| 1. | $B_B\varphi_0$ | Continuing from Case 2.1 |
| 2. | $B_B\varphi_0 \rightarrow \varphi_1$ | General B decides to follow A's plan. |
| 3. | $I_{A,B}\varphi_1$ | General B sends a messenger to A informing them of their army's intent. |
| 4. | $B_AI_{A,B}\varphi_1 \wedge T_{A,B}\varphi_1$ | General A believes the message from general B. |
| 5. | $B_A\varphi_1$ | By C1. |
| 6. | $B_A\varphi_1 \rightarrow \varphi_2$ | General A agrees. |
| ... | | The same logical chain repeats. |
| 7. | $w \models V_{\varphi_n}^x(w)$ | $V_{\varphi_n}^x(w)$ is always true. x is A or B depending on the value of n . |

Therefore, either A or B is not MSDND secure to the other for φ_n .

However, B is not MSDND secure to B because φ_1 is known to A. By extension For $i = 2, 4 \dots n$ B is secure to A, but not A to B. For $i = 3, 5 \dots n$, A is secure to B, but not B to A. □

Corollary 3. *Every message exchange where some atom φ_0 is sent, followed by any number n successful exchanges results in one agent being insecure to the other.*

Theorem 4. *If a messenger is captured, if the message is not resent, both agents will be secure on the last successfully delivered message atom φ_n or φ_0 if the first messenger is captured.*

Proof.

Case 4.1. *One messenger is sent and captured by the enemy.*

It should be obvious and direct that if the messenger does not arrive, it is equivalent to the messenger never being sent. (Theorem 1)

Case 4.2. *If $n - 1$ messengers successfully deliver their message, but messenger n is captured, both are secure on φ_n .*

Suppose General A sends φ_{n-1} to B. It should be obvious that on the delivery of the message φ_{n-1} to B, the value of φ_n is secure in B to A, as A has no way of knowing if φ_{n-1} was delivered, unless B sends φ_n with a messenger. When B does

send φ_n , the messenger never arrives. As a consequence, General A has no way of assigning value to $B_A\varphi_n$ ($V_{\varphi_n}^A \notin V$). However, as before, with φ_{n-1} at A is not secure to B. \square

5.3. BYZANTINE GENERALS

If General A is attempting to coordinate with multiple armies, the problem becomes more complex. If we extend the messenger analogy to cover faulty generals (ones that send incorrect information or omit messengers), the generals can reach consensus if for every faulty general, there are three generals that work correctly. This is a well known result known as the Byzantine Generals problem.

Theorem 5. *In any message exchange that conforms to the constraints of the Byzantine Generals problem, all agents are MSDND insecure on the plan φ .*

Proof. Proof: Suppose agent i decides to use plan φ to attack. Suppose that there is some set of Byzantine Generals T and some set of loyal generals G ($i \in G$). If $|G| > 3|T|$, the algorithm executes successfully, and $B_x\varphi : \forall x \in G$. Therefore, every general in G can valueate φ and the variable is insecure. \square

It is worth noting, however, that if all generals intend to behave well (they are non-byzantine) and messages are lost in transit, consensus can only be reached if the initial distribution of φ is delivered to all parties, and each general still receives enough messages to determine the majority consensus. This would be impractical to ensure in an actual application.

State Determination. When an agent uses the information transfer operator ($I_{i,j}$) to pass information to another agent in the system, it intends for that agent to believe the passed wff. When an agent distributes a wff to many agents with the information transfer operator, it leads to a set of beliefs about the beliefs of the receiving agent. This set N_i is the set of beliefs agent i can have if all the wffs it

passed to the other agents are believed. For example, if an agent distributes a wff φ to a set of agents Ag ($i \notin Ag$), then $N_i = \{B_i B_j \varphi : \forall j \in Ag\}$. Since the belief of each $B_j \varphi$ is outside of the domain of i , the agent i can only value a wff in N_i which has been leaked to i .

Let the set L_i be the subset of N_i for which a valuation function exists in a domain i . L_i can be populated either by direct information transfer or information leakage from interactions with agents. We can similarly define a set M_i which is the subset of N_i and superset of L_i if the agent i had perfect knowledge of the beliefs of information it passed to other agents ($L_i \subseteq M_i \subseteq N_i$).

Theorem 6. *Each member of L_i is MSDND insecure to i .*

Proof. Proof: Each wff in L_i has a valuation function in the domain i .

- | | |
|---|--|
| 1. $I_{j,i}\varphi$ | i informs some agent j of φ |
| 2. $B_j I_{j,i}\varphi \wedge T_{j,i}\varphi$ | j receives φ and believes its authenticity |
| 3. $B_j \varphi$ | By C1. |
| 4. $I_{i,j}\varphi_{ack}$ | j acknowledges $B_j \varphi$ |
| 5. $B_i I_{j,i}\varphi_{ack} \wedge T_{i,j}\varphi_{ack}$ | i receives φ_{ack} |
| 6. $B_i \varphi_{ack}$ | By C1. |
| 7. $B_i \varphi_{ack} \rightarrow B_i B_j \varphi$ | Because j acknowledged φ , i believes j believes φ . |
| 8. $w \models V_{B_i B_j \varphi}^i(w)$ | i does believe φ |
| 9. $w \models V_{B_j \varphi}^j(w)$ | Is always true. |

Therefore, $j \in L_i$, and $B_j \varphi$ is MSDND insecure to i .

□

Theorem 7. *Each member of M_i and N_i but not L_i are MSDND secure to i .*

Proof. Proof: Each wff in M_i and N_i have no valuation in the domain i .

Case 7.1. *The case where j receives some wff φ and is in M_i but not L_i .*

Case 7.2. *The case where j does not receive some wff φ and is in N_i but not M_i .*

□

1. $I_{j,i}\varphi$ i informs some agent j of φ
2. $B_j I_{j,i}\varphi \wedge T_{j,i}\varphi$ j receives φ and believes its authenticity
3. $B_j\varphi$ By C1.
4. $I_{i,j}\varphi_{ack}$ j acknowledges $B_j\varphi$
5. $\neg(B_i I_{j,i}\varphi_{ack} \wedge T_{i,j}\varphi_{ack})$ i does not receive φ_{ack}
6. $w \vdash V_{\varphi_{ack}}^i(w)$ i is uncertain if $B_j\varphi$
7. $w \vdash V_{\varphi}^j(w)$ Is always true

Therefore, $j \in M_i$, and $B_j\varphi$ is MSDND secure to i .

1. $I_{j,i}\varphi$ i informs some agent j of φ
2. $\neg(B_j I_{j,i}\varphi \wedge T_{j,i}\varphi)$ j does not receive φ
3. $w \not\vdash V_{B_j\varphi}^i(w)$ i is uncertain if $B_j\varphi$
4. $w \not\vdash V_{\varphi}^j(w)$ j is uncertain of φ

Therefore, $j \in N_i$, and $B_j\varphi$ is MSDND secure to i .

We assume that in our system, any beliefs an agent hold stem from information transfer from another agent. Therefore, we can assert that the beliefs in L_i for any process i must have a traceable history that stems from a process having a valuation for the referenced atom that aligns with the belief process i holds about that wff. Furthermore, wffs in N_i but not L_i do not have valuation in domain D_i and cannot be used in the construction of the image.

5.3.1. Example With Two Armies. We enforce a requirement that a General cannot reconsider a plan (in this case, φ_0). Then, while the receiver of the last message, φ_n , can construct the same model of the probability to attack as the sender, the recipient can construct a more accurate model given that the message delivery event has occurred. In fact, if the sender is committed to a proposed plan, for this problem the recipient the recipient can be certain the attack will succeed.

Instead, consider a simple algorithm to reach consensus for a system with no byzantine generals, but with omission failure (message loss). Again, each message has the probability p of being delivered. General A distributes the plan φ_0 to k other generals. If General A expects no confirmations, the probability that the message was delivered to all k generals is p^k . If the first message is delivered, the receiving

agents can construct the same model as A if they know how many generals there are and the probability that the messages are delivered.

However, for a longer algorithm with multiple exchanges and states, it is not sufficient to use the consensus probability to determine the state, for an overall view of the system because the state of the system is a probability distribution and not a fixed state.

Instead we obligate the design of the algorithm to rely on its underestimate of the system state L and not on the overestimate for determining the actions taking back the algorithm in the next step. For a leader election this is a good construct: with the correct message passing, the other agent can immediately infer that it was not a part of L if it receives φ_0 again.

Consider the simple consensus algorithm from before. While all the agents that receive φ can construct the same model, that model is not sufficient for any agent to determine the state of the system. However, if a confirmation message $\varphi_{1,i}$ is sent by each agent that received φ_0 then the original sender knows for certain the state of each agent that sent $\varphi_{1,i}$. Therefore, the original agent (a) is certain of a portion of the agents ($i = 1, 2 \dots n$) that have leaked $\varphi_{1,i}$ to it, assuming that $B_a \varphi_{1,i} \rightarrow w \models V_{\varphi_{1,i}}^i(w)$ for the current world, w .

However, this approach is limited because of the set of $\varphi_{1,i}$ that a believes is potentially a subset of all the agent that have a valuation function for $\varphi_{1,i}$ if any of the messages are lost in transit. Based on the information flows presented here and previously, if there is no response mechanism in the algorithm, after the transmission of φ_0 , a has a set of agents which believe might have received φ_0 , but can make no determination of the actual state of the system. If there is an acknowledgment mechanism then a can determine a subset of the total system state.

Constructing a model based on a subset of the complete state may seem like a mistake, but recall that A does not believe that those other entities believe φ_0 .

Assume a simple system where A is attempting to ensure every agent believes φ_0 . If the actions of the agents are divided into discrete steps, then each step A distributes φ_0 to each of the agents that A has not received $\varphi_{1,i}$ from. A Markov chain can easily be constructed for the probability of eventually arriving at state where every agent believes φ_0 . Indeed, this algorithm is a special case of the coin flipping leader election algorithm, an algorithm for leader elections in anonymous networks.

Algorithm 1 Anonymous Coin Flipping Leader Election

```

1:  $b(i) \leftarrow \text{random}(0, 1)$ 
2: Send  $b$  to every active neighbor
3: Receive  $b$  from every active neighbor
4: if  $b(i) = 1 \wedge \forall j \neq i : b(j) = 0$  then
5:    $i$  is the leader.
6: else if  $(b(i) = 1 \wedge \exists j \neq i : b(j) = 1 \vee (\forall k : b(k) = 0))$  then
7:    $b(i) = \text{random}(0, 1)$ 
8:   Go to next round
9: else if  $b(i) = 0 \wedge \exists j : b(j) = 1 : b(j) = 1$  then
10:  Become passive.
11: end if

```

Each round, the transfer of φ_0 to active agents is line 3 of the algorithm. The receive step is equivalent to the expected response. The other agent's $b(i)$ value is decided by the receipt of φ_0 from A. Agents go idle on the receipt of φ_0 , as if they had randomly selected a 0. If the other agents do not receive φ_0 it is logically equivalent to them holding a 1. If their confirmation is not received by A, it is logically equivalent to that agent successfully sending a 1 to A. However, this approach has the advantage of being able to evaluate, (given the omission rate) approximately how many rounds it will take for every agent to acknowledge φ_0 . This construction is presented below in Figure 5.1, which shows the structure of the chain A could generate for the election. In the figure, each state records the set of “passive” processes[54].

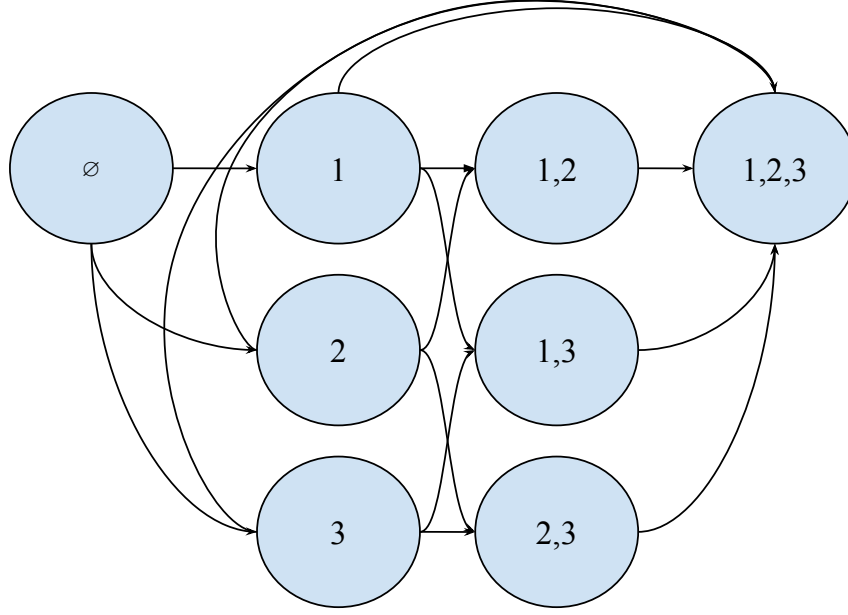


Figure 5.1: Example of a Markov chain constructed for an election algorithm from information flows. Each state represents the processes that have gone passive.

5.4. ELECTION IN AN ANONYMOUS COMPLETE NETWORK

Consider a version of the “coin-flipping” leader election, expressed in terms of BIT logic. This algorithm is functionally identical, but contains additional guards on the conditionals as an additional expression of the knowledge of the agents executing the algorithm. Additionally, note that the construct of labeling each φ that an agent receives is simply assigning labels to make the algorithm easier to reason about. The algorithm only considers the complete collection of φ s for execution and not the source of any of the φ used.

Additionally, let ψ_i be the state where an agent completes the algorithm as the leader, and γ_i is the state where an agent has terminated the algorithm. Correct operation is where $\{\gamma_i = T : \forall i\}$ is true in the same round of execution and exactly one ψ is true ($\sum_i \psi_i = 1$).

Algorithm 2 Anonymous Coin Flipping Leader Election Expressed in BIT logic

```

1:  $\varphi_i \leftarrow \text{random}(T, F)$ 
2:  $\psi_i \leftarrow F$ 
3: Send  $\varphi_i$  to every active neighbor ( $\forall j \neq i : I_j, i\varphi_i$ )
4: Receive  $\varphi_j$  from every active neighbor
5: if  $\varphi_i \wedge (\forall j \neq i, w \models V_{\varphi_j}^i(w) : \neg\varphi_j)$  then
6:    $\psi_i \leftarrow T$ 
7:    $\gamma_i \leftarrow T$ 
8: else if  $(\varphi_i \wedge \exists j \neq i, w \models V_{\varphi_j}^i(\varphi_j)(w) : \varphi_j) \vee (\forall k \text{ s.t. } w \models V_{\varphi_k}^i(w) : \neg\varphi_k)$  then
9:    $\varphi_i \leftarrow \text{random}(T, F)$ 
10:  Go to next round
11: else if  $\exists j \neq i, w \models V_{\varphi_j}^i(w) : \varphi_j$  then
12:    $\gamma_i \leftarrow T$ 
13: end if

```

Theorem 8. *Algorithm 2 may not terminate correctly unless there is detectable, perfect information transfer to all parties in the algorithm.*

Proof: Let i be the agent that would correctly terminate as the leader. Let j be a process that has selected $\varphi_j = F$.

| | |
|---|---|
| 1. $\varphi_i = T, \varphi_j = F$ | Initial conditions |
| 2a. $I_{j,i}\varphi_i$ | i sends φ_i to j |
| 2b. $I_{i,j}\varphi_j$ | j sends φ_j to i |
| 3a. $\neg(B_j I_{j,i}\varphi_i \wedge T_{j,i}\varphi_j)$ | j does not receive φ_i . |
| 3b. $B_i I_{i,j}\varphi_j \wedge T_{i,j}\varphi_i$ | i receives φ_j . |
| 4a. $w \not\models V_{\varphi_j}^i(w)$ | j cannot value φ_i . |
| 4b. $B_i \varphi_j$ | By C1. |
| 5a. $\neg V_{\varphi_j}^i(w) \wedge \neg\varphi_j \rightarrow (\varphi \leftarrow T)$ | j Cannot determine that i can terminate and incorrectly changes φ . |
| 5b. $\varphi_i \wedge w \models V_{\varphi_j}^i(w) \wedge \neg\varphi_j \rightarrow \psi_i$ | i incorrectly terminates as the leader. |

Therefore, i will terminate before j decides to go passive.

Which agrees with results from similar analysis [18]. As with the Byzantine Generals problem, algorithms similar to Byzantine agreement can be constructed for the anonymous networks.

6. ALGORITHM AND MODEL CREATION

6.1. EXECUTION ENVIRONMENT

Execution occurred in a real-time partially synchronous environment. The execution environment was subject to omission failures. Processes synchronized their clocks and executed steps of the election algorithm at predefined intervals. Processes with clocks that were not sufficiently synchronized could not form groups. For this work, process execution occurred in an environment where the clocks were sufficiently synchronized to consistently form groups. In the real-time environment, messages that were delayed and missed their real-time deadlines had the same appearance as an omitted message. The execution environment for the election algorithm had an omission fault occurrence modeled as a Bernoulli trial. In this model, each message had some probability p of being delivered within the timing constraints imposed by the real-time schedule. For the purpose of analyzing the effects of omission failures, processes were not subject to other faults.

6.2. ELECTION ALGORITHM OVERVIEW

A state machine for the election portion of the invitation-election algorithm is shown in Figure 6.1. In the normal state, the election algorithm regularly searches for other coordinators to join with. When another coordinator is identified the identifying coordinator will attempt to invite the other coordinator to their group. In the invitation election algorithm, processes are assigned a priority based on their process ID. The coordinator with the highest priority is the first to send invites. After a brief delay, if it appears that coordinator did not send their invites, the next highest process will send their invites. Coordinators that receive invites will forward the invite

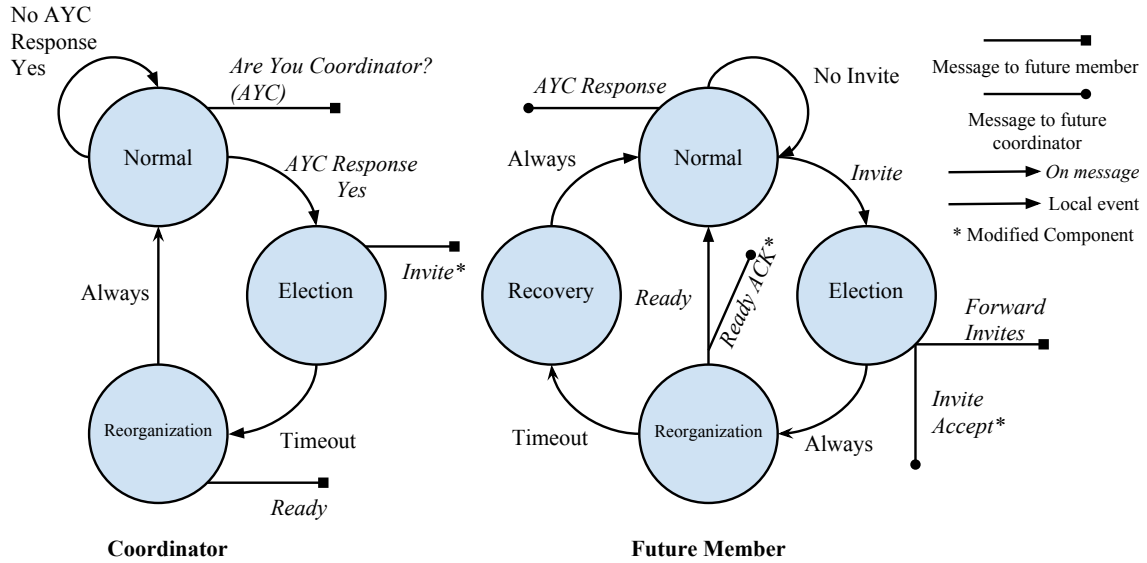


Figure 6.1: State machine of a leader election. Processes start as coordinators in the “Normal” state and search for other coordinators to join with. Processes immediately respond to AYC messages they receive. The algorithm was modified by adding a “Ready Acknowledgment” message as the final step of completing the election. Additionally, processes only accept invites if they have received an “AYC Response” message from the inviting process.

to its group members. Processes that receive an invite that are not already in an election will accept the invite. Once a timeout expires, the coordinator will send a “Ready” message with a list of peers to all processes that accepted the invite. The invited processes have timeouts for when they expect the ready message to arrive. If the message does not arrive in time, the process will enter the recovery state where it resets to a group by itself.

Once a group is formed it must be maintained. To do this, processes occasionally exchange messages to verify the other is still reachable. This interaction is shown in Figure 6.2. Coordinators send “Are You Coordinator” messages to members of its group to check if the process has left the group. Group members send “Are You There” messages to the coordinator to verify they haven’t been removed from the group, and to ensure the coordinator is still alive. If processes fail to reply to

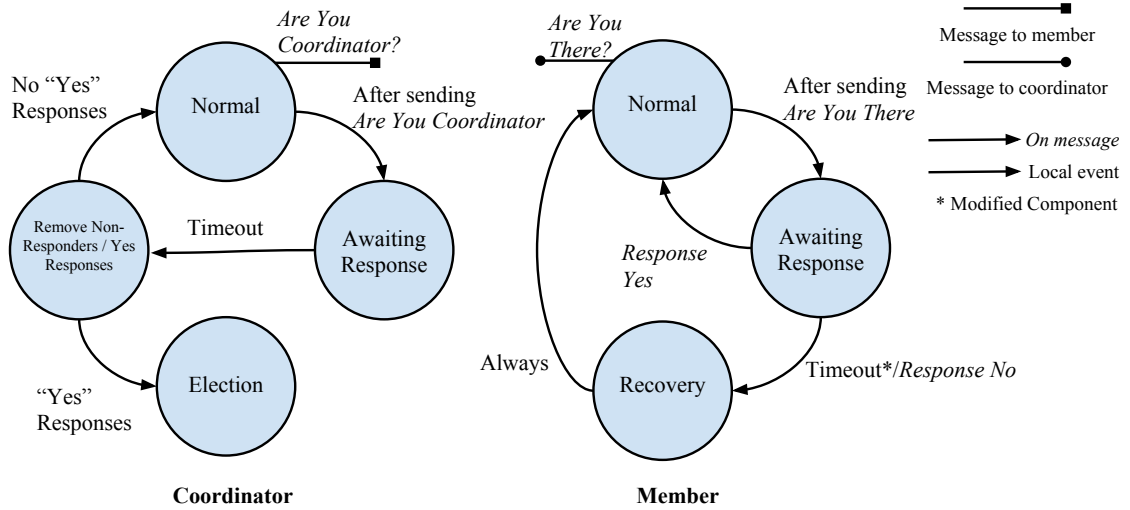


Figure 6.2: State machine of maintaining a group. The AYC messages are the same as those in Figure 6.1. AYC and AYT are periodically sent by processes, and responses to those messages are immediately sent by the receiving process. In the modified algorithm, the member does not enter the recovery state if they do not receive an AYT response before the timeout expires.

received message before a timeout, they will leave the group. Leaving the group can either be caused by the coordinator removing the process, or the process can enter a recovery state and leave the group, forming a new group by itself.

6.3. ALGORITHM

- 1: $AllPeers \leftarrow \{1, 2, \dots, N\}$
- 2: $Coordinators \leftarrow \emptyset$
- 3: $UpPeers \leftarrow Me$
- 4: $State \leftarrow Normal$
- 5: $Coordinator \leftarrow Me$
- 6: $Expected \leftarrow \emptyset$
- 7: $Counter \leftarrow$ A random initial identifier
- 8: $GroupID \leftarrow (Me, Counter)$

```

9:  $PendingID \leftarrow (Me, -1)$ 
10:  $Pending \leftarrow \emptyset$ 
11:
12: function CHECK
13:   This function is called at the start of a round by a leader
14:   if  $State \neq Normal$  or  $Coordinator \neq Me$  then return
15:   end if
16:    $Expected \leftarrow \emptyset$ 
17:   for  $j \in (AllPeers - \{Me\})$  do
18:      $AreYouCoordinator(j)$ 
19:      $Expected \leftarrow Expected \cup j$ 
20:   end for
21:   Peers which respond “Yes” to  $AreYouCoordinator$  are put into the  $Coordinators$ 
      set.
22:   Processes that respond are removed from  $Expected$ .
23:   When an  $AreYouThere$  response is “No” and this process is a coordinator, the
      querying process is put in the  $Coordinators$  set.
24:   Wait for responses, Peers that do not respond are removed from  $UpPeers$ .
25:    $UpPeers \leftarrow (UpPeers - Expected) \cup Me$ 
26:    $UpPeers \leftarrow (UpPeers - Coordinators) \cup Me$ 
27:   if  $Responses = \emptyset$  then return
28:   end if
29:   if  $Me = SelectedProcess$  then
30:     MERGE( $Responses$ )
31:   end if
32: end function
33:
34: function TIMEOUT
35:   This function is called at the start of a round by the group members

```

```

36:   if Coordinator = Me then return
37:   else AREYOU THERE(Coordinator, GroupID, Me)
38:       if Response is No then RECOVERY
39:           Coordinators  $\leftarrow$  Coordinator
40:       end if
41:   end if
42: end function
43:
44: function MERGE(Coordinators)
45:   This function invites all coordinators in Coordinators to join a group led by Me
46:   State  $\leftarrow$  Election
47:   Stop work
48:   Counter  $\leftarrow$  Counter + 1
49:   PendingID  $\leftarrow$  (Me, Counter)
50:   Coordinator  $\leftarrow$  Me
51:   Pending  $\leftarrow$  UpPeers – Me
52:   for  $j \in$  Coordinators do INVITE(j, Coordinator, PendingID)
53:   end for
54:   Wait for responses, Peers that accept the invite are added to Pending.
55:   State  $\leftarrow$  Reorganization
56:   GroupID  $\leftarrow$  PendingID
57:   UpPeers  $\leftarrow$  Pending
58:   for  $j \in$  UpPeers do READY(j, Coordinator, GroupID, UpPeers)
59:   end for
60:   Expected  $\leftarrow$  UpPeers
61:   Wait for responses, Peers that acknowledge are removed from Expected
62:   UpPeers  $\leftarrow$  UpPeers – Expected
63:   State  $\leftarrow$  Normal
64: end function

```

```

65:
66: function RECEIVEREADY(Sender, Leader, Identifier, Peers)
67:   if  $State = Reorganization$  and  $PendingID = Identifier$  then
68:      $UpPeers \leftarrow Peers$ 
69:      $State \leftarrow Normal$ 
70:      $Coordinator \leftarrow Leader$ 
71:      $GroupID \leftarrow Identifier$  READYACKNOWLEDGE(Leader,Identifier)
72:   end if
73: end function
74:
75: function RECEIVEAREYOUCOORDINATOR(Sender)
76:   if  $State = Normal$  and  $Coordinator = Me$  then
77:     Respond Yes
78:   else
79:     Respond No
80:   end if
81: end function
82:
83: function RECEIVEAREYOUTHERE(Sender, Identifier)
84:   if  $GroupID = Identifier$  and  $Coordinator = Me$  and  $Sender \in UpPeers$  then
85:     Respond Yes
86:   else
87:     Respond No
88:      $Coordinators \leftarrow Sender$ 
89:   end if
90: end function
91:
92: function RECEIVEINVITE(Sender,Leader,Identifier)
93:   if  $State \neq Normal$  then return

```

```

94:   end if
95:   if Sender  $\neq$  SelectedProcess then return
96:   end if
97:   Stop Work
98:   PendingID  $\leftarrow$  Identifier
99:   State  $\leftarrow$  Election
100:  Accept(Coordinator, Identifier)
101:  State  $\leftarrow$  Reorganization
102:  if Ready is not received then
103:    Recovery()
104:  end if
105: end function
106:
107: function RECEIVEACCEPT(Sender, Leader, Identifier)
108:   if State  $\leftarrow$  Election and PendingID = Identifier then
109:     Pending  $\leftarrow$  Pending  $\cup$  Sender
110:   end if
111: end function
112: function RECEIVEREADYACKNOWLEDGE(Sender)
113:   Pending  $\leftarrow$  Pending  $\cup$  Sender
114: end function
115:
116: function RECOVERY
117:   State  $\leftarrow$  Election
118:   Stop Work
119:   Counter  $\leftarrow$  Counter + 1
120:   GroupID  $\leftarrow$  (Me, Counter)
121:   Coordinator  $\leftarrow$  Me
122:   UpPeers  $\leftarrow$  Me

```


123: $State \leftarrow Reorganization$

124: $State \leftarrow Normal$

125: **end function**

6.4. MODEL CONSTRUCTION

As part of the execution model described, execution proceeds in a partially synchronous fashion. The state of the system is captured by a process using an imaging function $Im : D_i \times w \rightarrow o$. The imaging function maps a domain and world into an observation. The domain D_i is used to restrict the set of valid valuation functions used to produce the observation of the world. Each observation is taken using the imaging function after each round of execution for the algorithm completes.

We describe the system with two Kripke Models modeling the algorithm at different levels of detail. Both models share a set of valuation functions. In the first model, $K_1 = \langle W_1, R_1, V \rangle$, each action and variable change taken by a process corresponds to a relation in R_1 . With this model, the knowledge of each process can be analyzed to determine which processes have access to which knowledge in their domain.

In the second model, $K_2 = \langle W_2, R_2, V \rangle$, we describe the system with a Kripke Model such that each world is the state of the system after each execution of the algorithm. The set of worlds W_2 in K_2 is a subset of the worlds in W_1 . Each relationship in R_2 is the transition the state makes from the end of the last round of execution to the end of the current round of execution. Therefore, R_2 is a sequential composition of one or more elements of R .

Both K_1 and K_2 can be mapped to Markov chains by assigning each transition a probability based on the likelihood for each particular combination of successes and failures for message delivery that occurred between the two worlds in the model. K_2

is less powerful than K_1 (less things can be proven about execution) but it is more convenient for reasoning about Markov Models.

The election algorithm produces and distributes a set of processes *UpPeers* and a coordinator of that set. *UpPeers* is distributed via message passing and maintained by the coordinator. Different processes may have different versions of *UpPeers* for a given coordinator as processes enter and leave the group. However, a process will eventually receive an up-to-date version of *UpPeers* from the coordinator, or it will leave the group.

The imaging function used in this system captured the cardinality of the *UpPeers* for a process i if that process was a coordinator, or a observation otherwise. One can safely assume that this observation from an imaging function corresponds to an observation for an HMM: the process taking the image does not have complete access to the state. However, for the highest priority process for this algorithm, we will show that the observation of that process corresponds directly to the underlying Markov chain. Consequently, that process does not need to work with an HMM but instead uses a Markov chain. We allow if the observations and the associated probabilities are themselves a Markov chain.

Theorem 9. *For each world w such that $Im(D_i, w) = I_x$ and for each world such that $Im(D_i, w') = I_y$ such that I_y is an immediate successor of I_x , if the transition probability of going from I_x to I_y is the same as the probability of going from w to w' , then observation function and the associated transition probabilities are a Markov chain.*

Proof. Lalalala. □

6.5. HIGHEST PRIORITY PROCESS

6.5.1. State Determination. In the Garcia-Molina version of the algorithm, processes distribute a list of processes that have accepted their invite. Let i be the coordinator of a group distributing ready messages. The process i has a list of processes that have accepted its invite. Let φ_x be a wff that indicates that some process x is part of i 's group. Let $R = \{\varphi_x | x \in \text{AcceptedInvite}\} \cup \{\varphi_i\}$ be the set of processes that have accepted i 's invite. To finalize the group, i will distribute R to each process described in R . If a process does not receive R it will not be in the group.

Theorem 10. *In the original algorithm, each process is MSDND secure on φ_x for each x that is not i or itself (j).*

Proof.

Case 10.1. *The case where $x = i$ or $x = j$.*

If we assume that the authenticity of messages is not questioned in this environment, a ready message from i must mean i is a part of the group. Therefore, i 's inclusion in the group is MSDND insecure to j

If j accepts the ready message, that process will be a part of the group in R and since j knows its own state, j knows it is a part of R .

Case 10.2. *The case where $x \neq i$ and $x \neq j$*

For process i , i distributes R to each other process in R . Let Q be some set of processes that do not receive the R set.

□

Corollary 11. *The set of processes in the ready message is an N_i set, where i is the group's coordinator.*

- | | |
|---|--|
| 1. $I_{j,i}R \forall j \in R$ | i distributes the list of processes that have accepted the invite. |
| 2. $\neg(B_x I_{x,i}R \wedge T_{x,i}R) \forall x \in Q$ | Processes in Q do not receive R . |
| 3. $B_x I_{x,i}R \wedge T_{x,i}R \forall x \notin Q$ | Processes not in Q receive R . |
| 4. $w \not\models V_{B_x R}^j(w) \forall \{x x \notin \{i, j\}\}$ | j does not know which processes received R |

Therefore, the receipt of R by other members of the group is MSDND secure to j .

Since every wff in R is MSDND secure, R must be an N_i set, and consequently, $L_i = \emptyset$. As a consequence, the group leader only has an estimate of the system, which may be an overestimate ($M_i \subseteq N_i$)

With an ready acknowledgment message from members of the group, the Leader can construct an L_i set, which underestimates the state of the system. From previous chapter, we have shown that with message passing and omission, at least one process is left insecure to the other.

Theorem 12. *If a process j sends a ready acknowledgment message to i , j 's inclusion in the i 's view of the system state is MSDND secure to j*

Proof. Proof: When j sends the ready acknowledge message to i , it is uncertain if the acknowledgment message is actually delivered. From Theorem 2.1, we know that j is now MSDND insecure to i , but i is insecure to j . Using the ready acknowledgment messages, i can then valuate $B_i B_j \varphi_j$ for every j that sends an acknowledgment message. □

Once the ready acknowledgment messages have been sent to the leader, the group can begin to interact. As it is impossible for the members of the group to be distributed the L_i set, they must operate using the N_i set they received from the coordinator. Messages from other processes which could only have been sent if they are a part of the M_i sent can leak to the receiving processes their receipt of the R set from i . For the purpose of model construction, the determination of members of

M_i that are not in L_i may be undesirable for a DTMC since it may involve changing the state captured for the chain at the incorrect moment. Instead, for simplicity, while the leader can update their L_i set with leaked information, we avoid doing so to preserve the memorylessness property in the following sections. We defined the state of the system as the cardinality of the processes mentioned in the L_i set.

6.5.2. Memorylessness. The a process can update L_i with leaked information. To ensure the memorylessness property we are interested in capturing when a process receives the ready message, but the acknowledgment is lost, meaning the process is not included in the L_i set. We wish to ensure that the process can determine its own exclusion in L_i so it can behave as though it is not in M_i . To do this, we attach an additional field to the AYC message to indicate if the coordinator i considers the receiving process j a part of the L_i set. This value obviously leaks to the receiving process that it is not part of L_i and it should behave accordingly.

Likewise, for the coordinator, an AYT message leaks to the coordinator that the sending process is a part of the M_i set, if it is not in L_i . Since the behavior of the process in M_i is defined to behave as though it was not in L_i , the coordinator should do the same. This should be the case even if the process sending AYT has already responded negatively to an AYC from i . For the purpose of memorylessness, i should consider j to have responded in the affirmative.

Processes determine which invite to accept based on the exchange of the AYC messages. Processes always seek to reach the “lowest energy state”: they wish to be in a group led by the coordinator with the highest priority. Since the process has determined which invite it will accept, the probability of receiving that invite is the probability that process did not interact with a higher priority process that caused it to expect an invite. Processes will submit AYC messages to higher priority processes that are not its coordinator, even if they are in a group. As a result, the processes will always seek highest priority, regardless of their current group. More importantly,

this releases the process from obligations on its next state from the state of other processes in the system.

A process will only accept an invite if it receives it and it has determined it is the invite it wishes to accept. This result allows a process to determine the likelihood of an invite being accepted. This probability is based on the probability the higher priority process selecting that process to send invites to, and the lower priority process selecting that invite.

6.5.3. Model Construction. Based on the state determination and memorylessness arguments presented, the highest priority process operates independent of the state of the other processes in the system. The highest priority processes invites are always accepted, and the structure of the algorithm prevents any process from being locked out of participating in a round of execution based on the outcome of the previous round. Therefore, for every observation of the system state the highest priority process makes, the next round of execution for all the other processes favors the high priority process, ensuring that its observation corresponds directly with the system state.

In each round, the behavior is described by two components: maintaining a group and inviting other processes into the group. The coordinator will exchange an “Are You Coordinator” message and the peer will respond to verify is still available. To maintain a group of m other processes, the probability is defined as a random variable with the following pdf:

$$\Pr_M(X = k; m) = \begin{cases} \binom{m}{k} p^{2k} (1 - p^2)^{m-k}, & \text{if } 0 \leq k \leq m \\ 0, & \text{otherwise} \end{cases} \quad (6.1)$$

Where k is the number of processes remaining in a group selected from m processes. A process will leave a group if, from the considered process’s perspective, they do not respond to an “Are You Coordinator” message.

To invite other processes to the group, the two processes ultimately exchange up to 8 messages. In a round, a single process can invite many other processes to its group. From a selection of n other coordinators, the probability distribution for joining a new group with k of the n processes is:

$$\Pr_I(Y = k; n) = \begin{cases} \binom{n}{k} p^{8k} (1 - p^8)^{n-k}, & \text{if } 0 \leq k \leq n \\ 0, & \text{otherwise} \end{cases} \quad (6.2)$$

In the profile chain, in a state s that describes the number of processes in a group, the probability of transitioning from s to s' with n total processes (including the considered process) is:

$$\Pr_T(Z = s'; n; s) = \sum_{i=0}^{s-1} \Pr_M(X = i; s - 1) \cdot \Pr_I(X = s' - i; n - s - 1) \quad (6.3)$$

From this distribution, a set of transition probabilities can be calculated for a given omission rate p and number of processes n . This set of transition probabilities forms a profile Markov chain P , which can be evaluated to for any number of processes n and omission rate p . The generated profile chain is ergodic when $0 < p < 1.0$. The profile chain is a stationary Markov chain.

6.5.4. Model Validation. To assert the closed form profile chain accurately represents the implementation of the algorithm, it must be validated. Since T and P are ergodic, they can be checked for equivalence using a goodness-of-fit test. If the goodness-of-fit test indicates the chains are equivalent, they will generate similar sequences and have similar properties when analyzed. Therefore, generated P chains

can be used to analyze the behavior of the algorithm during live execution with changing conditions.

To verify the test chain T is equivalent to the profile chain P , a χ^2 goodness-of-fit test is employed. The null-hypothesis of this test (H_0) asserts the profile chain P is equivalent to the test chain T :

$$H_0 : T = P \quad (6.4)$$

With an alternative hypothesis that the two chains are not equivalent:

$$H_1 : T \neq P \quad (6.5)$$

The χ^2 test measures the goodness-of-fit for a complete chain by combining the measurements of goodness of fit for the transitions away from each state. Therefore, the goodness of fit test for the chain is a summation of tests for each state:[7]

$$\chi^2 = \sum_i^m \sum_j^m = \frac{n_i(P_{ij} - T_{ij})^2}{P_{ij}} \quad (6.6)$$

Where n_i is the number of times the state i was observed in the input sequence used to construct the test chain T . The summation is distributed as χ^2 with $m(m-1)$ degrees of freedom (DF) if all entries in P_{ij} are non-zero. In this work, all transitions in the profile Markov chain P are non-zero when $0 < p < 1.0$. However, the probability of some transitions may be extremely small. The χ^2 value was compared to a critical value (CV) giving a measure of how likely it was H_0 could not be rejected. This work selected an $\alpha = 0.05$ significance level to reject the hypothesis $T = P$.

If the hypothesis H_0 were to be rejected, it would indicate the test chain and profile chain differ significantly. As a consequence of rejecting the hypothesis, the implementation would have behavior from the generated closed form solution.

| Processes | DF | CV | Worst Error | $\Pr(WorstError)$ | p-value |
|-----------|----|------|-------------|-------------------|---------|
| 3 | 6 | 12.6 | 8.90 | 0.80 | 0.18 |
| 4 | 12 | 21.0 | 14.55 | 0.75 | 0.27 |
| 5 | 20 | 31.4 | 23.47 | 0.65 | 0.27 |
| 6 | 30 | 43.8 | 32.69 | 0.85 | 0.34 |

Table 6.1: Summary of χ^2 tests performed.

To verify the model, it was compared to runs of an implementation of the algorithm. Test data was collected for systems with 3, 4, 5 and 6 processes. Information was collected from sufficiently long runs of the system with an omission rates between 0.05 and 0.95 tested at intervals of 0.05. Table 6.1 shows the measured error and p-value for the worst observed error for each number of processes. Since the measured error is less than the critical value and the p-value is greater than 0.05, we cannot reject H_0 . As a consequence, the profile chains (P) are representative of the behavior of the algorithm's implementation.

6.6. PROFILE CHAIN ANALYSIS

Resources can only be managed effectively when multiple DGI coordinate together to manage those resources. Without another DGI to coordinate with, the DGI has a limited range of options to manage power generation, storage and loads. Therefore, the amount of time DGI will spend coordinating with another process is of particular interest. [29] defines an in group time (IGT) metric to measure the amount of time a DGI process spends coordinating with at least one other process. In this work, we define IGT based on the steady state of the profile chain. Let $\pi = Steady(P)$ for some profile chain. The IGT is the sum of all states in π , save the first state where the process is alone:

$$IGT = \sum_{i=2}^m \pi_i \quad (6.7)$$

The IGT is a number between 0 and 1. It represents the probability a random observation sees a group of at least two members. The steady state distributions are presented as stacked bar graphs in Figures 6.3, 6.4, 6.5 and 6.6. Each complete bar in the graph indicates the IGT. Additionally, Figures 6.3, 6.4, 6.5 and 6.6 also contain the average group size (AGS), when the system has reached the steady-state, plotted as a fraction of the total number of processes. Let P be the steady-state distribution vector for some number of processes, n , and a given omission rate:

$$y = \frac{\sum_{i=1}^n P_i * i}{n} \quad (6.8)$$

where y is the plotted AGS as a fraction.

The components of each bar represent the probability the system is in a specific state for a random observation of the system. The height of the component represents the relative probability of observing that state when in a group.

The profile chain can be used to ensure the FREEDM smart grid is able to continue operating despite network issues. The profile chain can be combined with different message sending strategies to maintain service. For example, the DGI can change to a slower mode of operation to ensure operation continues normally despite connectivity issues. By selecting different strategies depending on the message delivery probability the DGI can offer high performance in good network conditions and an acceptable level of service during faults. The profile chain can be extended to an arbitrary number of processes as shown in Figure 6.8. In Figure 6.8, the steady-state of the system is used to compute a weighted average of the group size. To compare the produced steady states, the weighted average was plotted as a percentage of all processes in the system. Values in Figure 6.8, were plotted using Equation 6.8.

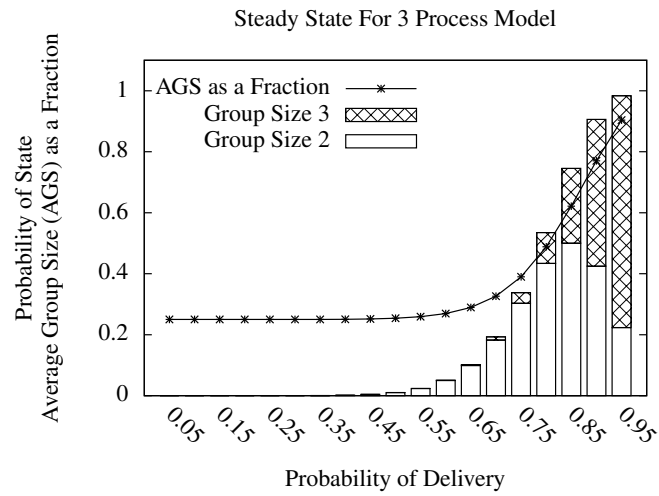


Figure 6.3: Steady state distribution for 3 processes as well as the AGS as a fraction of total processes.

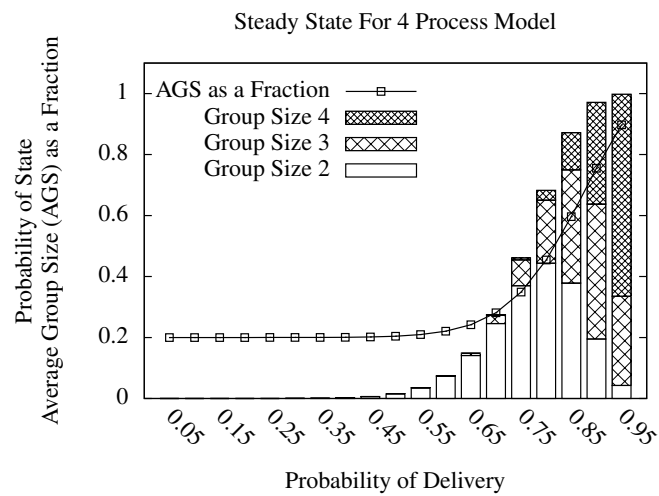


Figure 6.4: Steady state distribution for 4 processes as well as the AGS as a fraction of total processes.

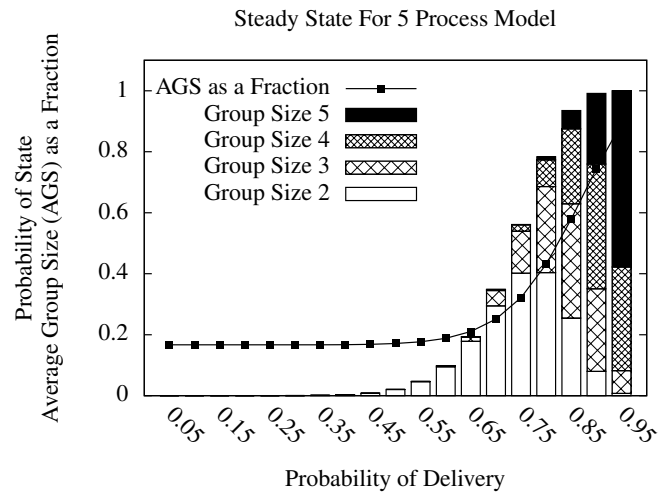


Figure 6.5: Steady state distribution for 5 processes as well as the AGS as a fraction of total processes.

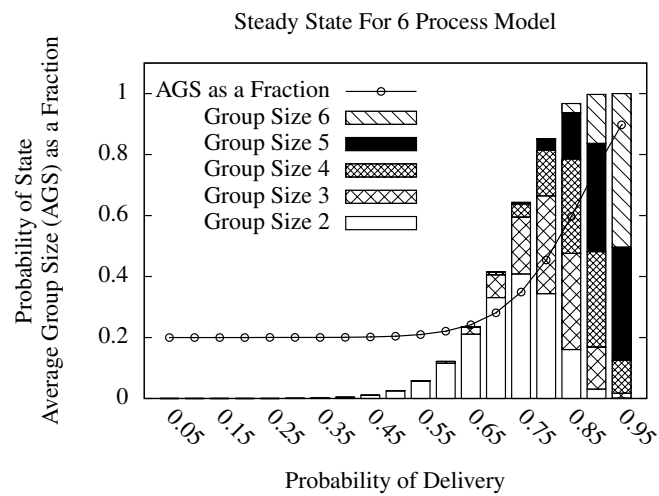


Figure 6.6: Steady state distribution for 6 processes as well as the AGS as a fraction of total processes.

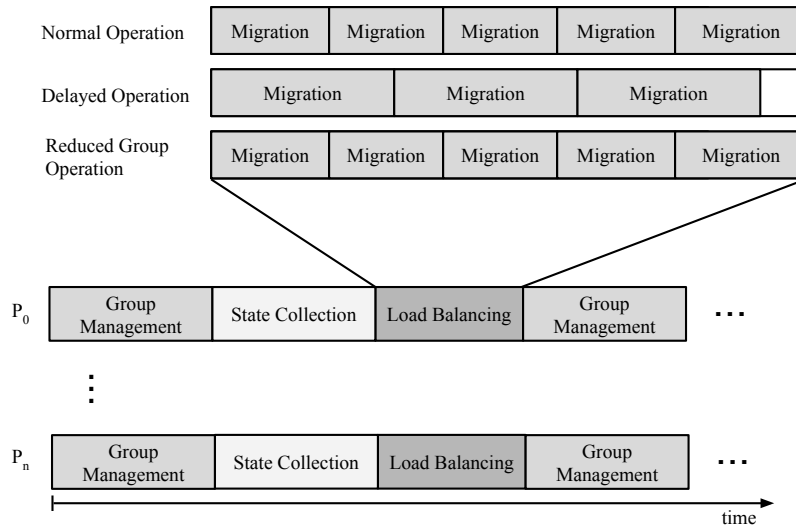


Figure 6.7: Example DGI schedule. Normal operation accounts for a fixed number of migrations each time the load balancing module runs. Message delays reduce the number of migrations that can be completed each round. However, reducing the group size allows more migrations to be completed (because fewer messages are being exchanged) at the cost of flexibility for how those migrations are completed.

Therefore, Figure 6.8, shows the average percentage of total processes that will be in the group in a steady-state system.

The DGI employed a round-robin schedule where each module is given a pre-determined amount of time to execute. This schedule was determined by the number of DGI that could be grouped together and the expected cyber network conditions. Additionally, the Load Balance module, which managed the power resources was scheduled to run multiple times in a long block before the group was evaluated for failure. This schedule is depicted in 6.7. The system began by using group management to organize groups. Next load balancing ran to manage power resources in the created group. Finally, state collection collected a casually consistent state to be used for reporting and offline analysis. If message delays occurred, the number of migrations load balancing could perform was reduced. However, by reducing the

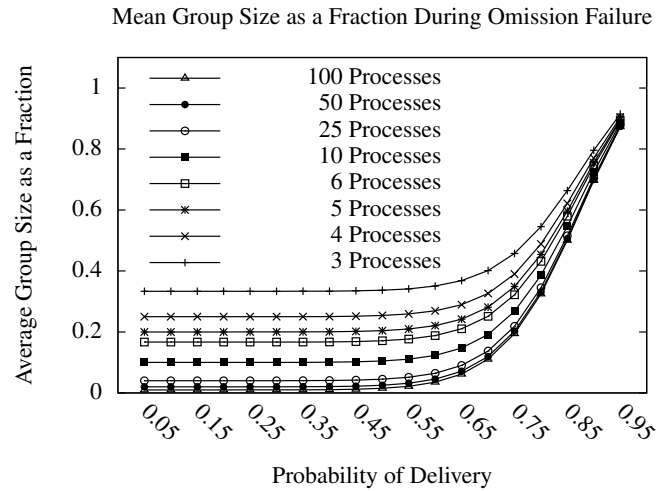


Figure 6.8: Average group size as a percentage of all processes in the system for larger systems.

group size, the number of messages sent by load balancing could be reduced, allowing for a greater amount of work to be done. The outcome of an election had a direct effect on how effectively the DGI can manage resources. Therefore the steady-state analysis gives a good indication of how well the DGI could manage those resources in the future.

This information can also be used to anticipate faults and mitigate them before they occur. Modern routers can supply an expected congestion notification as part of the IP header[20]. When congestion is anticipated, the coordinator can preemptively split the group to reduce congestion. This is possible because the algorithm's message complexity is $O(n^2)$. Future work will combine this analysis with ECN techniques to preemptively change the behavior of DGI to ensure a good level of service. If dropped messages account for an omission rate of even 0.15 from future congestion, performing a coordinated group division can potentially save several rounds of transient states. Breaking the large groups into smaller sets can drastically reduce the number of messages transmitted and help relieve congestion. Additionally, the coordinator can

design the split to ensure work can continue when the groups are split, by mixing supply and demand processes. The division can also account for the placement of congested routers for targeted congestion relief. Furthermore, since the required execution time is a function of group size, the DGI can use the additional execution time within the same real time schedule to use more reliable techniques to deliver messages.

In Figure 6.7, Group Management's execution is broken into four steps: check, merge, ready, and cleanup. Between each step, the DGI waits while messages and their replies are delivered. The leader election algorithm expects that all replies arrive before the next step of the algorithm is executed. If a message does not arrive during the wait period, before the next step of execution the message is considered lost. If it arrives later, it is ignored by the algorithm. While dynamically adjusting the synchronized schedule is not feasible during failure, adjusting the number of messages sent (by sending fewer "Are You Coordinator" messages, for example) is. By sending fewer messages, the number of packets in the communication network is reduced, and the savings in processing can allow multiple delivery attempts in the scheduled wait time.

6.7. LOWER PRIORITY PROCESS

A model cannot be constructed for a process that is not the highest priority for this algorithm. The lower priority process cannot know which processes are already in a group with the higher priority process. In the execution model we have presented, the inclusion of a process in a higher priority process's group is MSDND secure. There is no way for a process that is not in a group to know what processes are in any group except its own.

6.8. INDEPENDENT ALGORITHM

If we restrict the algorithm to be more like the Bully algorithm, running the partially synchronous environment, any process can construct their model. This is accomplished by removing the maintenance portion, releasing each process from the previous state of execution entirely.

7. APPLICATION

In this work, we consider the effects of network congestion on a cyber communication network used by the FREEDM smart-grid. We create a model version of a large number of DGI processes in a simple partitioned setup. In the FREEDM smart-grid, the consequences of this congestion could result in several problematic scenarios. First, if the congestion prevents the DGI from autonomously configuring using its group management system, processes cannot work together to manage power devices. Secondly, if messages arrive too late, or are lost, the DGI could apply settings to the attached power devices that drive the physical network to instability. These unstable settings could lead to problems in the power-grid like frequency instability, blackouts, and voltage collapse.

These techniques allow the DGI to anticipate behavior during a fault and allow it to preemptively harden itself against the fault.

7.1. APPLICATION: ECN HARDENING

7.1.1. Usage Theory. Since the ECN fields in IPv4 are not available to applications running on the system, the notifications are multicast onto the source interface. This application is responsible for generating the multicast ECN message. It also keeps a register of hosts running applications that support reacting to the ECN notification.

There are several reasons for this approach. First, related work has shown an ECN strategy without some other queue management scheme is not sufficient to prevent congestion. By allowing real-time applications that decrease the number of messages for congestion special priority in the RED algorithm, we allow those

applications to continue operating during congestion. Additionally, in later sections, we demonstrate this strategy is effective for managing congestion.

When the RED algorithm identifies congestion it must notify senders of congestion. Since this approach is non-standard and most UDP applications would not understand the notification, we have opted to create an application that runs on switches and routers. Congestion is detected, the application sends a multicast beacon to a group of interfaces informing the attached devices of the level of congestion. For similarity with the RED algorithm and the NS-3 implementation, this notification is classified as either “soft” or “hard.” A soft notification is an indication the congestion in the network is approaching a level where real-time processes can expect message delays that may affect their normal operation. A hard notification indicates the congestion has reached a level where messages are subject to both delay and loss.

7.1.2. Group Management. The group management module’s execution schedule is broken into several periods of message generation and response windows. Because the schedule of the DGI triggers the execution of group management modules approximately simultaneously, the traffic generated by modules is bursty. The number of messages sent is $O(n^2)$ (where n is the number of processes in the system), in a brief window, which is dependent on how well the clocks are synchronized in the system. The duration of the response window is dependent on the amount of time it takes for messages to propagate to the hardest-to-reach process the DGI hopes to group with. Additionally, to contend with congestion, an additional slack must be added to allow the RED algorithm to detect congestion before it reaches a critical level. Figure 7.1 depicts typical queuing behavior for a network device serving DGI processes under different circumstances.

Because the traffic generated by DGI modules is very bursty, the queue experiences a phenomena where the bursty traffic mixed with a steady background traffic causes the queue to fill. With no background traffic, the impulse queues a large

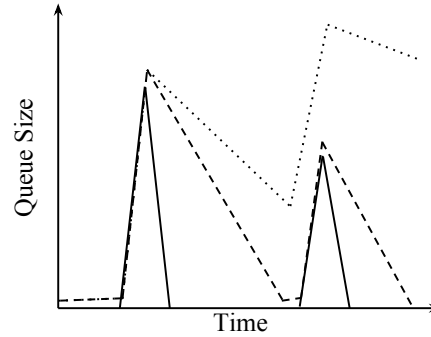


Figure 7.1: Example of network queuing during DGI operation. DGI modules are semi-synchronous, and create bursty traffic on the network. When there is no other traffic on the network (solid line), the bursty traffic causes a large number of packets to queue quickly, but the queue empties at a similar rate. With background traffic (dashed line), the bursty traffic causes a large number of packets to be queued suddenly. More packets arrive continuously, causing the queue to drain off more slowly. When the background traffic reaches a certain threshold (dotted line), the queue does not empty before the next burst occurs. When this happens, messages will not be delivered in time, and the queue will completely fill.

number of messages, but those messages are distributed in a timely manner. When the background traffic is introduced, the queue takes longer to empty. At a critical threshold, the queue does not empty completely before the next burst is generated by the DGI. In this scenario, the queue completely fills and no messages can be distributed. The RED algorithm and ECN are used to delay or prevent the queue from reaching this critical threshold.

The leader election algorithm specified previously was used. This algorithm has a higher message complexity when in a group than the Garcia-Molina algorithm it is based on. However, it does possess a desirable memoryless property that makes it easy to analyze. This work uses an improved version of the algorithm which removes the restrictions in [30] where only one process could become the leader.

Soft ECN. A soft ECN message indicates the network has reached a level of congestion where the router suspects processes will not be able to meet their real time requirements. The soft ECN message encourages the DGI processes to reduce the

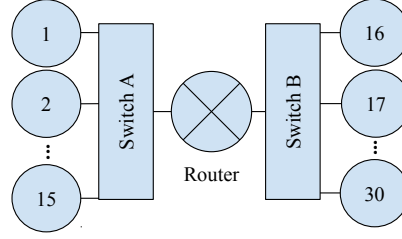


Figure 7.2: Example of process organization used in this paper. Two groups of processes are connected by a router.

number of messages they send to reduce the amount of congestion they contribute to the network, and to allow for reliable distribution techniques to have additional time to deliver messages (since fewer messages are being sent). In the case of potential congestion, the group management module can reduce its traffic bursts by disabling elections during the congestion. When the elections are disabled, messages for group management are only sent to members of the group. Processes do not seek out better or other leaders to merge with. As a consequence, the message complexity for processes responding to the congestion notification reduces from $O(n^2)$ to $O(n)$.

Hard ECN. In a hard ECN scenario, the router will have determined congestion has reached a threshold where the real-time processes will soon not be able to meet their deadlines. In this scenario, the real-time process will likely split its group. In an uncontrolled situation, the split will be random. It is therefore desirable when this level of traffic is reached to split the group. Splitting the group reduces the number of messages sent across the router for modules with $O(n^2)$ (where n is the number of processes in the original group) message complexity. For larger groups, splitting them provides a significant savings in the number of messages that must be queued by the router, especially since the traffic is very bursty.

Suppose a network like one depicted in Figure 7.2, where processes are divided by a router. In Figure 7.2, there are n processes on one side of the network and m on

the other. In normal operation the omission-modelable algorithm has an $O(n^2)$ message complexity. In Soft ECN maintenance mode, the reduced number of messages reduces the complexity to $O(n)$ by disabling elections.

During elections (and with each group update) the leader distributes a fall-back configuration that will coordinate the division of the groups during intense congestion. When the ECN notification is received the processes will halt all current group management operations and enter a splitting mode where they switch to the fall-back configuration. The leader of the group distributes a fall-back notification to ensure all processes in the group apply their new configuration. The complexity of distributing the notification is linear $O(n)$ and processes that already received the notification will have halted their communication. This approach will ideally avoid the burst/drain phenomena from figure 7.1.

The design of the fall-back configuration can be created to optimize various factors. These factors include cyber considerations, such as the likely network path the processes in the group will use to communicate. By selecting the group around the network resources, the group can be selected to minimize the amount of traffic that crosses the congested links in the future. Additionally, considerations from the physical network can be considered. Fall-back groups can be created to ensure they can continue to facilitate the needs of the members. This can take into the consideration the distribution of supply and demand processes in the current group. By having a good mix of process types in the fall-back group the potential for work can remain high.

7.1.3. Cyber-Physical System. For a real-time CPS, message delays could affect coordinated actions. As result, these actions may not happen at the correct moments or at all. Since the two-army problem prevents any process from being entirely certain a coordinated action will happen in concert, problems arising from delay or omission of messages is of particular interest. In particular, we are interested

in the scenario from [13], where only half of a power migration is performed. Other power management algorithms could have similar effects on the power system based on this idea of a process performing an action that is not compensated for by other processes.

Soft ECN. In a soft congestion notification mode, the process being informed of the congestion can reduce its affect on the congestion by changing how often it generates bursty traffic. Processes running the load balancing algorithm make several traffic bursts when they exchange state information and prepare migrations. As shown before, if the interval between these bursts is not sufficient for the queue to drain before the next burst occurs, then critical, overwhelming congestion occurs. Since the schedule of the DGI is fixed at run-time processes cannot simply extend the duration of the load balancing execution phase. However, on notification from the leader, the process can, instead, reduce the number of migrations to increase the message delivery interval. This notification to reduce the schedule originates from the coordinator as part of the message exchange necessary for the process to remain in the group. Every process in the group must receive this message to participate in load balancing, ensuring all processes remain on the same real-time schedule. Using this approach, the amount of traffic generated is unchanged but the time period a process waits for the messages to be distributed is increased.

Hard ECN. When the DGI process receives a hard congestion notification, the processes switch to a predetermined fall-back configuration. This configuration creates a cyber partition. By partitioning the network, the number of messages sent by applications with $O(n^2)$ message complexity can be reduced significantly. Each migration of load balancing algorithm begins with an $O(n^2)$ message burst and so benefits from the reduced group size created by the partition.

Suppose there is a network like the in Figure 7.2 with n processes on one half and m on the other. The number of messages sent across the router for the

undivided group is of the order $2mn$ as the n processes on side A send a message to the m on side B and vice-versa. Let i_1 and j_1 be the number of processes from side A and side B (respectively) in the first group created by the partition. Let i_2 and j_2 be the number of processes in the second group created by the partition under the same circumstances of i_1 and j_1 . The number of messages sent that pass through the router, is then

$$2i_1j_1 + 2i_2j_2 \quad (7.1)$$

For an arbitrary group division, the following can be observed. Suppose i_1 and j_2 are the cardinality of two arbitrarily chosen sets of processes from side A and side B respectively. Following the same cut requirements as before:

$$i_2 = n - i_1 \text{ and } j_2 = m - j_1 \quad (7.2)$$

The number of messages that must pass through the router for this cut is:

$$2i_1j_1 + 2(n - i_1)(m - j_1) \quad (7.3)$$

The benefits of the cut are minimized when i_1 and j_1 are $\frac{n}{2}$ and $\frac{m}{2}$:

$$2(2\frac{mn}{4} + mn - \frac{mn}{2} - \frac{mn}{2}) = mn \quad (7.4)$$

Which is a reduction of half as many messages. For systems with a large number of participating processes this represents a significant reduction in the number of messages sent across the router. As a consequence, this further extends the delivery window for processes sending messages.

7.1.4. Relation To Omission Model. The synchronization of clocks in the environment is assumed to be normally distributed around a true time value provided

by the simulation. The shape of the curve created by plotting the queue resembles that of the Cumulative Distribution Function (CDF) of the normal distribution, noted $F(x)$. A simple description of the traffic behavior can then be described in terms of that curve. First, observe that when the queue hits a specific threshold, even if the queue is drained at an optimal rate, the n th queued packet will not be delivered in time:

$$Qsize - \min(Qsize, (DequeueRate * \Delta t)) \geq 0 \quad (7.5)$$

Where Δt is the deadline for the message to be delivered. If the size of the queue exceeds the number of messages that can be delivered before Δt passes, some messages will not be delivered. The size of the queue during the message bursts created by the DGI depends on the message complexity of the algorithm, the number of messages already in the queue, the other traffic on the network, and any replies that also have to be delivered in that interval. Therefore, let c represent the rate that traffic is generated by other processes. Let $init_q$ represent the number of messages in the queue at the start of a burst. Let $init_m$ represent the number of messages sent in the beginning of the burst. Let $resp$ represent the number of messages sent in response to the burst that must still be delivered before Δt passes. We can then express $Qsize$ as two parts:

$$Qsize = Burst + Obligations \quad (7.6)$$

Where $Burst$ takes the form of the CDF for the normal distribution:

$$Burst = init_m * F(x) \quad (7.7)$$

$$Obligations = c * \Delta t + init_q + resp \quad (7.8)$$

From this we can derive the equation:

$$F(x) \geq \frac{DequeueRate * \Delta t - c * \Delta t - init_q - resp}{init_m} \quad (7.9)$$

Where, from Equation 7.5, $DequeueRate * \Delta t$ is less than or equal to the number of messages in the queue. Solving for $F(x)$ gives a worst case estimate of the omission rate for a specific algorithmic or network circumstance. $DequeueRate$ is affected by the amount of traffic in the system. It should be obvious a greater amount of background traffic corresponds to a greater average queue size. From an relationship between the background traffic and the average queue size and the results presented in [30], Equation 7.9 can be used to select the ECN parameters.

7.1.5. Calibration. Since the distribution of clock synchronization was selected to be a normal distribution, the shape of the curve created by plotting the queue resembles that of the CDF of the distribution, noted $F(x)$. A simple description of the traffic behavior can then be described in terms of that curve. First, observe that when the queue hits a specific threshold, even if the queue is drained at an optimal rate, the n th queued packet will not be delivered in time:

$$Qsize - (DequeueRate * \delta t) \geq 0 \quad (7.10)$$

Where δt is the deadline for the message to be delivered. If the size of the queue exceeds the number of messages that can be delivered before δt passes, some messages will not be delivered. The size of the queue during the message bursts created by the DGI depends on the message complexity of the algorithm, the number of messages already in the queue, the other traffic on the network, and any replies that also have to be delivered in that interval. Therefore, let c represent the rate that traffic is generated by other processes. Let $init_q$ represent the number of messages in the queue at the start of a burst. Let $init_m$ represent the number of messages

sent in the beginning of the burst. Let *resp* represent the number of messages sent in response to the burst that must still be delivered before δt passes. We can then express *Qsize* as two parts:

$$Qsize = Burst + Obligations \quad (7.11)$$

Where *Burst* takes the form of the CDF for the normal distribution:

$$Burst = init_m * F(x) \quad (7.12)$$

$$Obligations = c * \delta t + init_q + resp \quad (7.13)$$

From this we can derive the equation:

$$F(x) \geq \frac{DequeueRate * \delta t - c * \delta t - init_q - resp}{init_m} \quad (7.14)$$

Solving for $F(x)$ gives a worst case estimate of the omission rate for a specific algorithmic or network circumstance.

7.2. PROOF OF CONCEPT

7.2.1. Experimental Setup. Experiments were run in a Network Simulator 3.23[15] test environment. The simulation time replaced the wall clock time in the DGI for the purpose of triggering real-time events. As a result, the computation time on the DGI for processing and preparing messages was neglected. However, to compensate for the lack of processing time, the synchronization of the DGI was instead randomly distributed as a normal distribution. This was done to introduce realism to ensure events did not occur simultaneously. Additionally, the real-time

schedules used by the DGI were adjusted to remove the processing time that was neglected in the simulation.

The DGI were placed into a partitioned environment. The test included 30 nodes. Each of the nodes ran one DGI process. Two sets of 15 DGI were each connect to a switch and each switch was in turn connected to the router. This network is pictured in Figure 7.2. Node identifiers were randomly assigned to nodes in the simulation and used as the process identifier for the DGI.

The links between the router and the switches had a RED enabled queue placed on both network interfaces. The RED parameters for all queues were set identically. A summary of RED parameters are listed in Table 7.1. All links in the simulation were 100Mbps links with a 0.5ms delay. RED was used in packet count mode to determine congestion. ARP tables were populated before the simulation began. RED parameters were selected using results from [30].

The relationship between the background traffic and the average queue size was estimated through runs of the NS-3 simulation. Figure 7.3 demonstrates the observed relationship between the total background traffic and the maximum average queue size for that level of traffic. Additionally, the *DequeueRate* was collected from a run of the simulation without traffic, and was found to be 713.08 packets/second. Therefore, from Equation 7.9, assuming $init_q = 0, resp = 225, init_m = 225$ and $\Delta t = 1$, the maximum traffic rate with no omissions is 263.0 packets/second. The number of packets for the *resp* and $init_m$ were selected from the worst case of the algorithm in [30]. Based on the traffic parameters in Table 7.1, 263.0 packets/second corresponds to 1.077 Mbps of traffic generated at one switch and 2.1545 Mbps traffic overall. From the polynomial estimate in Figure 7.3, the maximum average queue size for that level of traffic is 94.715, estimated as 90 for the RED Min Threshold in Table 7.1. RED Max Threshold is computed using a similar technique, but using

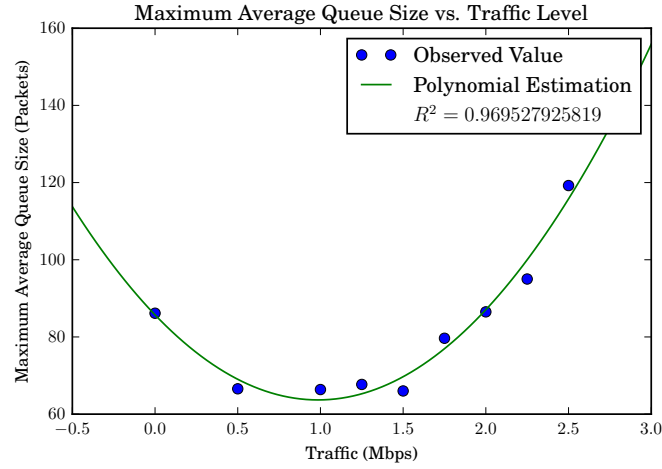


Figure 7.3: Plot of the maximum observed average queue size as a function of the overall background traffic. The polynomial estimate is $y = 22.70x^2 - 44.74x + 85.72$

| Parameter | Value | Parameter | Value |
|-------------------|----------|-------------------|--------|
| RED Queueing Mode | Packet | RED Gentle Mode | True |
| RED Q_w | 0.002 | RED Wait Mode | True |
| RED Min Threshold | 90 | RED Max Threshold | 130 |
| RED Link Speed | 100 Mbps | RED Link Delay | 0.5 ms |

Table 7.1: Summary of RED parameters. Unspecified values default to the NS-3 implementation default value

the message complexity for the Load Balancing algorithm, since it maintains its complexity during Soft ECN mode.

To introduce traffic, processes attached to each of the switches attempted to send a high volume of messages to each other across the router. The number of packets sent per second was a function of the data rate and the size of the packets sent. In each simulation, half of the traffic originated from each switch. Due to the bottleneck due to the properties of the network links, the greatest queueing effect occurred at the switches.

7.2.2. Results. Figures 7.4 and 7.5 show the normal operation of the system. In this configuration, there is no congestion on the network. The DGI start, group

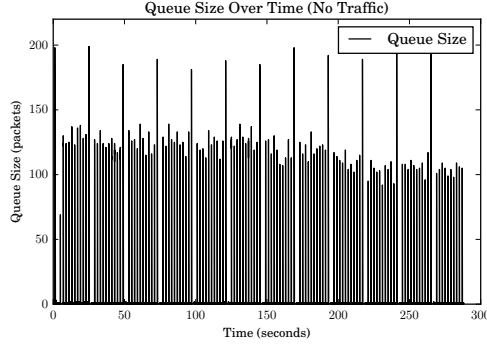


Figure 7.4: Plot of the queue size for a queue from switch A to the router when only the DGI generates traffic.

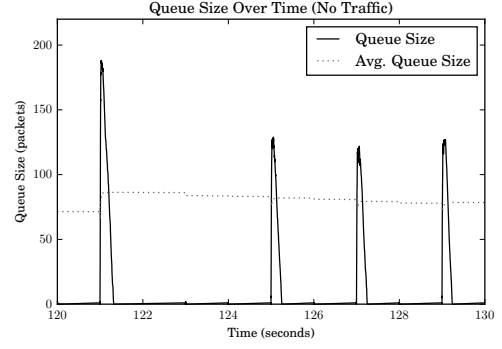


Figure 7.5: Detailed view of Figure 7.4. The left most peak is from Group Management, and the 3 smaller peaks are from power migrations.

together and then begin migrating power between processes. Figure 7.4 plots the queue size over time for a queue used to send packets from a switch to the router. Figure 7.5 is a detailed view of a portion of Figure 7.4. Figure 7.5 shows the queue size during the normal operation of group management as well as the first migration of the load balancing module. The dotted line plots the EWMA of the size of the queue.

From this experiment we establish the min_{th} value used as a RED queue parameter. The traffic generated by each step of the group management algorithm is very bursty. It should be obvious that the tightness of the clock synchronization in the group affect how large this peak is. Like [51], the level of the power at a process is the net sum of its power generation capability and load. As power is shared on the network, processes with excess generation, converge toward zero net power. Demand processes also converge toward zero net power.

Figure 7.6 shows the queue size as the network traffic begins to increase. The DGI in these experiments use a schedule that allows for some congestion to occur before processes are disrupted. This slack gives the network devices the opportunity to identify when the network congestion will go beyond the acceptable threshold.

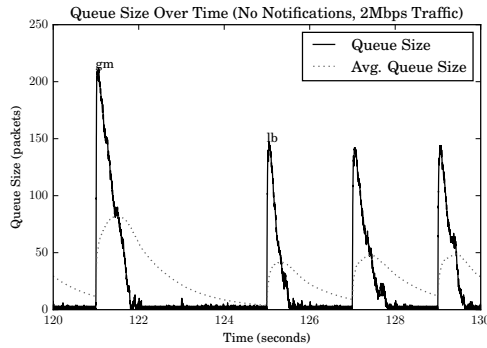


Figure 7.6: Detailed view of the effect on queue size as other network traffic is introduced. Compared to Figure 7.5, the peaks are taller and wider. Background traffic causes the average queue size to be updated more frequently.

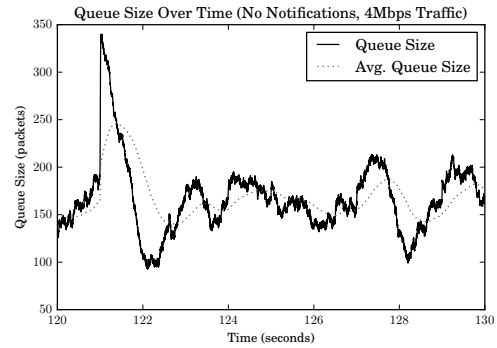


Figure 7.7: Detailed view of the effect on queue size as other network traffic is introduced. With no ECN notifications, the peak from Group Management is much larger. The congestion is sufficient that the Group Management and Load Balancing modules are affected.

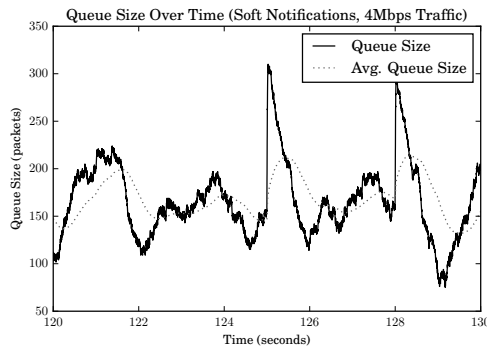


Figure 7.8: Detailed view of the effect on queue size as other network traffic is introduced. In this scenario, the ECN notifications put Group Management into a maintenance mode that reduces its message complexity and switches Load Balancing to slower migration schedule, preventing undesirable behavior.

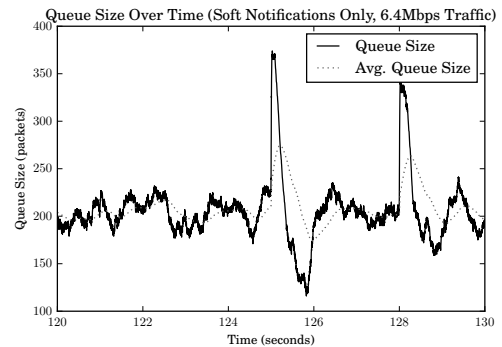


Figure 7.9: Detailed view of the effect on queue size as a large amount network traffic is introduced. Groups are unstable and processes occasionally leave the main group. Some migrations are lost due to queueing delays.

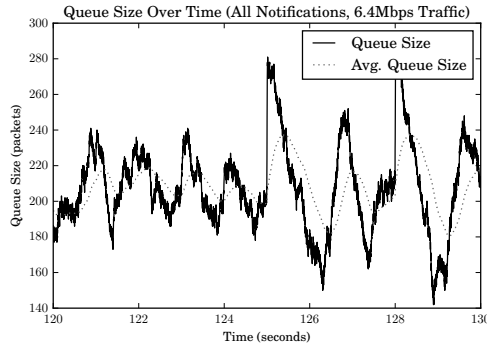


Figure 7.10: Effect on queue size as a large amount of network traffic is introduced. Hard notifications cause the groups to divide. As a result of the smaller groups, the group management and load balancing peaks are smaller than those in 7.9. No migrations are lost.

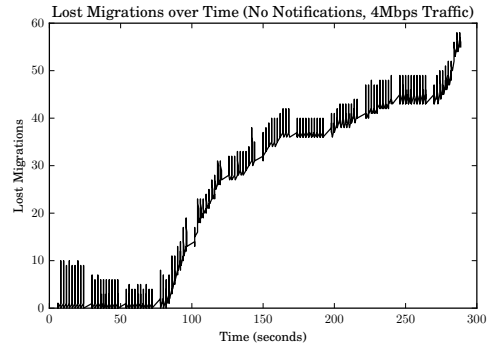


Figure 7.11: Count of lost migrations from all processes over time. Migrations are counted as lost until the second process confirms it has been completed. Without congestion management, a large number of migrations are lost.

Figure 7.7 shows an example of congestion affecting the physical network without ECN. As a result of the congestion in Figure 7.7, processes leave the main group. Additionally power migrations are affected: migrations are lost, or the supply process is left uncertain of migrations completions. Figure 7.11 plots the count of failed migrations over time.

Figure 7.8 shows an example of the ECN algorithm notifying processes of the congestion. Compared to the scenario in Figure 7.7, the ECN algorithm successfully prevents the group from dividing, and increases the number of migrations by reducing the number of attempted migrations each round.

Figures 7.10 and 7.9 show an example of a more extreme congestion scenario. In Figure 7.10, the RED algorithm shares a Hard ECN notification. This notification causes the DGI to switch to a smaller fall-back configuration. This fall-back configuration decreases the queue usage from Figure 7.9 to Figure 7.10. Without this

fall-back configuration behavior, the system is greatly affected by the traffic. However, with the fall-back configuration the system remains stable and no migrations are lost.

8. CONCLUSION

We presented a useful framework for reasoning about creating models of distributed systems that tolerate omission failure. The models and structures presented allow algorithm designers to design algorithms whose behavior can be modeled with a Markov chain which is invaluable for designing critical infrastructure systems that behave reliably during failure conditions.

To do this, applied existing information flow analysis techniques and applied them to a common distributed systems problem (the two-armies problem) and show the information flow analysis was consistent with similar analyses of the problem. We then extended the analysis to show how agents in a distributed system could use that analysis to determine what knowledge each process had. As part of this we described belief sets created by distributing information to several other agents in the system.

Additionally, we defined how information being transferred between agents and the actions they take based on that information could be constrained to hold to the memorylessness property common to Markov chains. Using this concept we demonstrated how a common leader election algorithm could be modified to use this memorylessness property, allowing it to be modeled online during changing conditions.

This work is particularly valuable for the analysis of critical infrastructure systems, where knowledge of their behavior during fault conditions is particularly important. By allowing the ability for algorithms to determine what issues are likely to arise while they are operating, actions can be taken to protect the infrastructure from failure. There are a wide range of possible applications, including actions either undertaken by human operators on site, or autonomous actions taken by the algorithms to harden themselves against failure.

We also presented a technique for hardening a real-time distributed cyber-physical system against network congestion. The RED queueing algorithm and an out-of-band version of explicit congestion notification (ECN) were used to signal an application of congestion. Using this technique the application changed several of its characteristics to ready itself for the increased message delays caused by the congestion.

These techniques were demonstrated on the DGI, a distributed control system for the FREEDM smart-grid project. In particular, this paper demonstrated the hardening techniques were effective in keeping the DGI processes grouped together. Additionally, it helped ensure the changes applied to the DGI through cyber-coordinated actions did not destabilize the physical power network.

This technique will be important to create a robust, reliable CPS for managing future smart-grids. However, this technique could potentially be applied to any CPS that could experience congestion on its network, as long as it has the flexibility to change its operating mode. Potential applications can apply to both the cyber control network and the physically controlled process. For example, in a VANET system, the vehicles could react to congestion by increasing their following distance.

BIBLIOGRAPHY

- [1] R. Akella, Fanjun Meng, D. Ditch, B. McMillin, and M. Crow. Distributed power balancing for the FREEDM system. In *Smart Grid Communications (SmartGridComm), 2010 First IEEE International Conference on*, pages 7–12, October 2010.
- [2] Y. Amir, C. Danilov, M. Miskin-Amir, J. Schultz, and J. Stanton. The spread toolkit: Architecture and performance. Technical report, Johns Hopkins University, 2004.
- [3] Y. Amir, D. Dolev, S. Kramer, and D. Malki. Transis: a communication subsystem for high availability. In *Fault-Tolerant Computing, 1992. FTCS-22. Digest of Papers., Twenty-Second International Symposium on*, pages 76–84, 1992.
- [4] Chongyang Bai and Xuejun Zhang. Aircraft landing scheduling in the small aircraft transportation system. In *Computational and Information Sciences (IC-CIS), 2011 International Conference on*, pages 1019–1022, Oct 2011.
- [5] J. Baillieul and P. J. Antsaklis. Control and communication challenges in networked real-time systems. *Proceedings of the IEEE*, 95(1):9–28, Jan 2007.
- [6] F. Baker. Requirements for IP version 4 routers, 6 1995. RFC 1812.
- [7] U. Narayan Bhat. *Elements of applied stochastic processes, 2nd Edition*. John Wiley & Sons, 1984.
- [8] K. Birman and Renesse R. *Reliable Distributed Computing with the Isis Toolkit*. IEEE Computer Society Press, Los Alamitos, CA 90720-1264, 1994.
- [9] P. Blackburn, M. De Rijke, and Y. Venema. *Modal logic*, volume 53. Cambridge University Press, 2002.
- [10] T. D. Chandra, V. Hadzilacos, and S. Toueg. The weakest failure detector for solving consensus. In *Proceedings of the Eleventh Annual ACM Symposium on Principles of Distributed Computing*, PODC '92, pages 147–158, New York, NY, USA, 1992. ACM.
- [11] T. D. Chandra and S. Toueg. Unreliable failure detectors for reliable distributed systems. *J. ACM*, 43(2):225–267, March 1996.
- [12] Tushar Deepak Chandra, Vassos Hadzilacos, and Sam Toueg. The weakest failure detector for solving consensus. *J. ACM*, 43(4):685–722, July 1996.

- [13] A. Choudhari, H. Ramaprasad, T. Paul, J.W. Kimball, M. Zawodniok, B. McMillin, and S. Chellappan. Stability of a cyber-physical smart grid system using cooperating invariants. In *Computer Software and Applications Conference (COMPSAC), 2013 IEEE 37th Annual*, pages 760–769, July 2013.
- [14] OMNeT++ Community. Omnet++, May 2012. <http://http://www.omnetpp.org/>.
- [15] NS-3 Consortium. Network simulator 3.23. <http://www.nsnam.org/>.
- [16] Flaviu Cristian. Understanding fault-tolerant distributed systems. *COMMUNICATIONS OF THE ACM*, 34:56–78, 1993.
- [17] Christopher Dabrowski and Fern Hunt. Using markov chain analysis to study dynamic behaviour in large-scale grid systems. In *Proceedings of the Seventh Australasian Symposium on Grid Computing and e-Research - Volume 99*, AusGrid '09, pages 29–40, Darlinghurst, Australia, Australia, 2009. Australian Computer Society, Inc.
- [18] Carole Delporte-Gallet, Hugues Fauconnier, and Hung Tran-The. *Distributed Computing and Networking: 14th International Conference, ICDCN 2013, Mumbai, India, January 3-6, 2013. Proceedings*, chapter Uniform Consensus with Homonyms and Omission Failures, pages 161–175. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013.
- [19] Qi Dong and Donggang Liu. Resilient cluster leader election for wireless sensor networks. In *Sensor, Mesh and Ad Hoc Communications and Networks, 2009. SECON '09. 6th Annual IEEE Communications Society Conference on*, pages 1–9, June 2009.
- [20] A. Dracinski and S. Fdida. Congestion avoidance for unicast and multicast traffic. In *Universal Multiservice Networks, 2000. ECUMN 2000. 1st European Conference on*, pages 360–368, 2000.
- [21] N. Falliere, L. O. Murchu, and E. Chien. W32.Stuxnet Dossier. <http://goo.gl/dC8VT>, 2010. [Online; accessed December 2011].
- [22] Timothy French. *Bisimulation Quantifiers for Modal Logics*. PhD thesis, University of Western Australia, 2006.
- [23] H. Garcia-Molina. Elections in a distributed computing system. *Computers, IEEE Transactions on*, C-31(1):48–59, January 1982.
- [24] S. Ghosh. *Distributed Systems: An Algorithmic Approach*. Chapman & Hall, 2007.
- [25] Aniruddha Gokhale, Mark P McDonald, Steven Drager, and William McKeever. A cyber physical systems perspective on the real-time and reliable dissemination

- of information in intelligent transportation systems. Technical report, DTIC Document, 2010.
- [26] C. Gomez-Calzado, M. Larrea, I. Soraluze, A. Lafuente, and R. Cortinas. An evaluation of efficient leader election algorithms for crash-recovery systems. In *Parallel, Distributed and Network-Based Processing (PDP), 2013 21st Euromicro International Conference on*, pages 180–188, 2013.
 - [27] Gerry Howser and Bruce McMillin. Modeling and reasoning about the security of drive-by-wire automobile systems. *International Journal of Critical Infrastructure Protection*, pages 127 – 134, 2012.
 - [28] Gerry Howser and Bruce McMillin. A Multiple Security Domain Model of a Drive-by-Wire System. In *Computer Software and Applications Conference (COMPSAC), 2013 IEEE 37th Annual*, pages 369–374. Computer Software and Applications Conference, 2013.
 - [29] S. Jackson and B. M. McMillin. The effects of network link unreliability for leader election algorithm in a smart grid system. In *Critical Information Infrastructures Security*, pages 59–70. Springer, Berlin, Heidelberg, 2013.
 - [30] Stephen Jackson and Bruce McMillin. Markov models of leader elections in a smart grid system (under review). *Journal of Parallel and Distributed Computing*, 2015.
 - [31] Joost-Pieter Katoen. Model checking meets probability: a gentle introduction. In *Engineering dependable software systems*, volume 34 of *NATO Science for Peace and Security Series - D: Information and Communication Security*, pages 177–205. IOS Press, Amsterdam, 2013.
 - [32] K. S. Kishor. Sharpe, March 2014. <http://sharpe.pratt.duke.edu/>.
 - [33] K. Koscher, A. Czeskis, F. Roesner, S. Patel, T. Kohno, S. Checkoway, D. McCoy, B. Kantor, D. Anderson, H. Shacham, and S. Savage. Experimental security analysis of a modern automobile. In *Security and Privacy (SP), 2010 IEEE Symposium on*, pages 447–462, May 2010.
 - [34] Saul A. Kripke. A Completeness Theorem in Modal Logic. *The Journal of Symbolic Logic*, 24(1):pp. 1–14, 1959.
 - [35] Leslie Lamport, Robert Shostak, and Marshall Pease. The byzantine generals problem. *ACM Trans. Program. Lang. Syst.*, 4(3):382–401, July 1982.
 - [36] Mai Gehrke and Hideo Nagahashi and Yde Venema. A Sahlqvist theorem for distributive modal logic. *Annals of Pure and Applied Logic*, 131(13):65 – 102, 2005.

- [37] N. Mohammed, H. Otrouk, Lingyu Wang, M. Debbabi, and P. Bhattacharya. Mechanism design-based secure leader election model for intrusion detection in MANET. *Dependable and Secure Computing, IEEE Transactions on*, 8(1):89–103, Jan 2011.
- [38] L.E. Moser, P.M. Melliar-Smith, D. A. Agarwal, R. K. Budhia, and C. A. Lingley-papadopoulos. Totem: A fault-tolerant multicast group communication system. *Communications of the ACM*, 39:54–63, 1996.
- [39] NSF FREEDM Systems Center. FREEDM, The Future Renewable Electric Energy Delivery and Management Systems Center. <http://www.freedm.ncsu.edu/>.
- [40] P. Olofsson. *Probability, Statistics, and Stochastic Processes, 2nd Edition*. John Wiley & Sons, 2012.
- [41] N. Privault. *Understanding Markov Chains*. Springer Singapore, 2013.
- [42] K. Ramakrishnan, S. Floyd, and D. Black. The addition of explicit congestion notification (ecn) to ip, 9 2001. RFC 3168.
- [43] D.B. Rawat, C. Bajracharya, and Gongjun Yan. Towards intelligent transportation cyber-physical systems: Real-time computing and communications perspectives. In *SoutheastCon 2015*, pages 1–6, April 2015.
- [44] Robbert Van Renesse, Takako M. Hickey, and Kenneth P. Birman. Design and performance of Horus: A lightweight group communications system. Technical report, Cornell, 1994.
- [45] Thomas Roth and Bruce McMillin. Breaking Nondeducible Attacks on the Smart Grid. In *Seventh CRITIS Conference on Critical Information Infrastructures Security*. Seventh CRITIS Conference on Critical Information Infrastructures Security, 2012. (to appear).
- [46] R.A. Sahner and K.S. Trivedi. Sharpe: a modeler’s toolkit. In *Computer Performance and Dependability Symposium, 1996., Proceedings of IEEE International*, page 58, Sep 1996.
- [47] K. Sampigethaya and R. Poovendran. Cyber-physical system framework for future aircraft and air traffic control. In *Aerospace Conference, 2012 IEEE*, pages 1–9, March 2012.
- [48] Matthias Schmalz, Daniele Varacca, and Hagen Völzer. *CONCUR 2009 - Concurrency Theory: 20th International Conference, CONCUR 2009, Bologna, Italy, September 1-4, 2009. Proceedings*, chapter Counterexamples in Probabilistic LTL Model Checking for Markov Chains, pages 587–602. Springer Berlin Heidelberg, Berlin, Heidelberg, 2009.

- [49] M.M. Shirmohammadi, K. Faez, and M. Chhardoli. Lele: Leader election with load balancing energy in wireless sensor network. In *Communications and Mobile Computing, 2009. CMC '09. WRI International Conference on*, volume 2, pages 106–110, Jan 2009.
- [50] C. Singh and A. Sprintson. Reliability assurance of cyber-physical power systems. In *Power and Energy Society General Meeting, 2010 IEEE*, pages 1–6, July 2010.
- [51] M.J. Stanovich, I. Leonard, K. Sanjeev, M. Steurer, T.P. Roth, S. Jackson, and M. Bruce. Development of a smart-grid cyber-physical systems testbed. In *Innovative Smart Grid Technologies (ISGT), 2013 IEEE PES*, pages 1–6, Feb 2013.
- [52] S. Vasudevan, B. DeCleene, N. Immerman, J. Kurose, and D. Towsley. Leader election algorithms for wireless ad hoc networks. In *DARPA Information Survivability Conference and Exposition, 2003. Proceedings*, volume 1, pages 261–272 vol.1, April 2003.
- [53] Y. Yan, Y. Qian, H. Sharif, and D. Tipper. A survey on smart grid communication infrastructures: Motivations, requirements and challenges. *Communications Surveys Tutorials, IEEE*, PP(99):1–16, 2012.
- [54] Haidi Yue and Joost-Pieter Katoen. *Analytical and Stochastic Modeling Techniques and Applications: 17th International Conference, ASMTA 2010, Cardiff, UK, June 14-16, 2010. Proceedings*, chapter Leader Election in Anonymous Radio Networks: Model Checking Energy Consumption, pages 247–261. Springer Berlin Heidelberg, Berlin, Heidelberg, 2010.
- [55] Ziang Zhang and Mo-Yuen Chow. Incremental cost consensus algorithm in a smart grid environment. In *Power and Energy Society General Meeting, 2011 IEEE*, pages 1–6, July 2011.

VITA