

# The FREEDM Architecture of Fault Tolerant Network Routing through Software Overlays \*

Christopher Zimmer

North Carolina State University  
cjzimme2@ncsu.edu

Frank Mueller

North Carolina State University  
mueller@cs.ncsu.edu

## Abstract

Control decisions of intelligent devices in critical infrastructure can have a significant impact on human life and the environment. Insuring that the appropriate data is available is crucial in making informed decisions. Such considerations are becoming increasingly important in today's cyber-physical systems that combine computational decision making on the cyber side with physical control on the device side. In the FREEDM system, power management of green energy is provided in a highly distributed and scalable manner. The system has to insure that Intelligent Energy Management (IEM) and Intelligent Fault Management devices have the appropriate data to make control decisions for microgrids and with respect of microgrid connectivity to an upstream utility power grid. The job of insuring the timely arrival of the data falls onto the network designed to support these intelligent devices. This network needs to be fault tolerant. When nodes, devices or communication links fail along a default route of a message from A to B, the underlying hardware and software layers should ensure that this message will actually be delivered as long as alternative routes exist. Insuring multi-route pathways and discovery of these pathways is critical in insuring delivery of critical data. In this work, we present methods of developing network topologies of smart devices that will enable multi-route discovery in an intelligent power grid. This will be accomplished through the utilization of software overlays (1) that maintain a digital representation of the physical network and (2) allow new route discovery in the case of fault. Also, in this work we aim to present a visualization of the connection states and pathways through the network aimed at helping external entities to understand the states of the network. Our vision is that the application of this approach in an intelligent power grid will enable IEM and IFM devices to make automated, decentralized decisions and to maintain state of lower-level devices.

**Keywords** Distributed Networking, Distributed Fault Tolerance

## 1. Introduction

Failures of network equipment in intelligent systems can result in

- incorrect decisions regarding device failure,
- faulty decisions made due to lack of data,
- system reconfigurations, or
- degradation of system performance.

In modern network topologies, network failures resulting in these issues may be avoided through smart routing technologies that can take faulty equipment out of the loop. However, such fault tolerance

is only feasible in situations where the faulty equipment does not constitute a single point of failure of communication within the network. Therefore, it is important to maintain redundant pathways through networks.

Routing decisions are an important part of networking. Concrete routes are configured statically in many networks. When a networking device on a static route fails, any messages sent along that route will timeout and result in communication failure with respect to this end point. In these scenarios, many systems will assume the end point to be out of service. This does not have to be the case. Networks of devices can be designed to contain multiple pathways to connect clusters of nodes in a redundant manner. If these pathways exist, a network needs to be able to alternate and utilize them in times of fault.

Another consideration to make when designing a network is its shape / topology of connectivity. The shape of the network can have a significant impact on its performance. For example, networks with a ring topology can only sustain a single link failure. Fully connected mesh networks offer the greatest amount of fault tolerance but this comes at the cost of one connection per pair of nodes, which imposes exponential resource needs.

In this work, we present a method of utilizing software network overlays to provide shape and meta information about connectivity. Utilizing knowledge of shape/topology and meta information, the network is able to react in case of faults and generate new routes through the network in manner that is transparent to the user by providing a software overlay middleware.

## 2. Software Overlay Network

Using software overlays to improve network resilience is an idea first described by Anderson et al. [1]. In their work, they presented the basis for a resilient overlay network (RON) partitioning distributed nodes that may contain a different topological perspective than the external, physical network topology. Their work assumed nodes to potentially be separated geographically across the Internet. Our work utilizes a similar partitioning for the routing of messages but deviates in that it utilizes this approach in a much smaller local area network (LAN) to facilitate fault-tolerant communication. In our approach, devices are organized into software partitions that are calculated locally based on their IP address. Partitions are created as a side effect of subnet masks. Each partition is assumed to be fully connected. These partitions are then grouped together in clusters of a certain static size. The combined group of clusters and partitions are interconnected with horizontal and vertical uplinks. An example is depicted in Figure 1. Utilizing vertical uplinks, we then organize our software overlay network into a tree-based topology. Similar work has been performed in the High Performance Computing domain by Varma et al. [4]. Uplinks will serve as the default routing path for general message communication in the absence of

\*This work was supported in part by NSF grant EEC-0812121 and U.S. Army Research Office (ARO) grant W911NF-08-1-0105 managed by NCSU Secure Open Systems Initiative (SOSI).

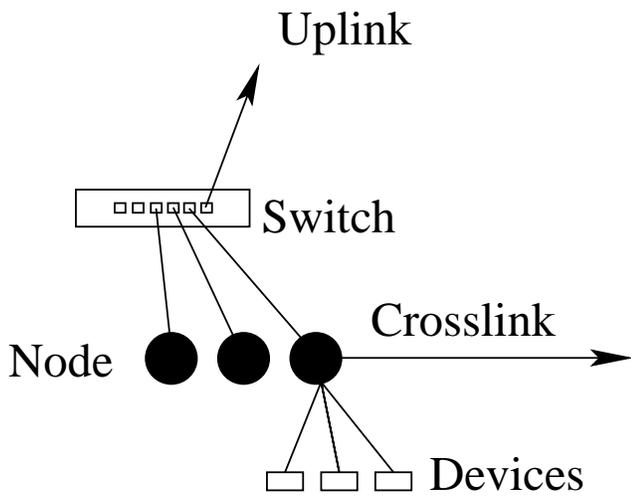


Figure 1. Device Cluster

failures. Figure 2 depicts the vertical uplinks and shows the resulting tree formed by them. Uplinks are necessary to provide inter-cluster communication. They constitute the network backbone of the system. To increase fault tolerance, it is necessary to introduce horizontal crosslinks that will serve as secondary paths through the network, as depicted in Figure 2.

This abstract software overlay can fit onto arbitrary intelligent power grids. Most importantly, it provides redundant communication pathways and the potential to connect the network in alternate ways in case of faults in the system via its software middleware layer. This capability is crucial for allowing intelligent nodes in the system to maintain appropriate state and to coordinate the actions of system control tasks.

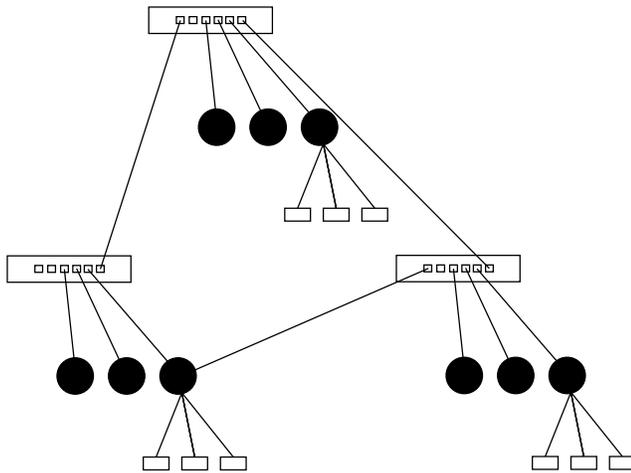
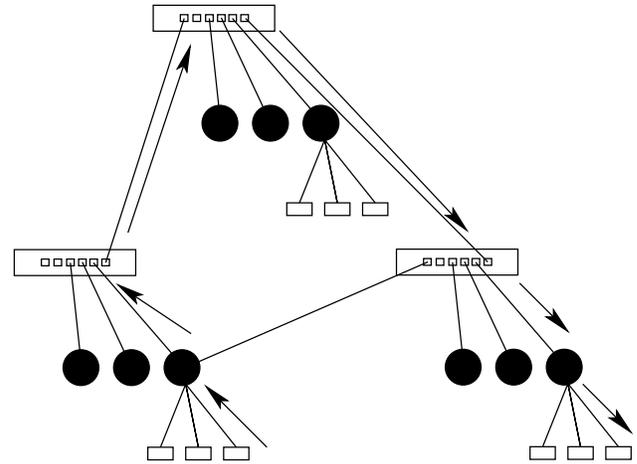


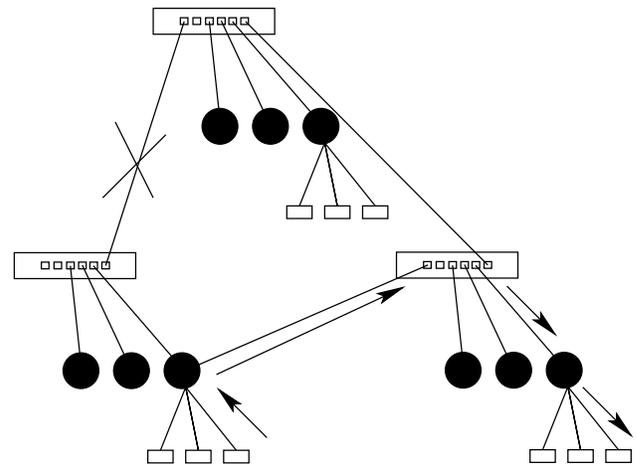
Figure 2. Cluster Tree

In Figure 3 communication pathways are primarily used through the switching interface composed of uplinks as depicted in the first half of the figure. Our system differs from a regular network in the composition of a series of intelligently placed crosslinks that can occur as node to switch or node to node lines. In the event of a loss of an uplink, the abstract network will enter into a reorganization mode. Reorganization mode is typically defined in the system as what to do in the

case of a sent message timeout. In this work, the reorganization mode will explore the possibilities of alternate routing in the network by using meta-information describing the characteristic of the network. In using software overlays part of the network information that is provided to a node is its partition information that the node can use to determine its neighbors on a switch and the partitions above and below it in a tree. From this information a node in reorganization mode can begin communicating with its neighbors to determine important information. The most important of this information being the location of crosslinks, and through using the crosslinks determining if this is a node failure on the receiving end or a link failure along the switching path. The second half of Figure 3 depicts the utilization of a crosslink in an attempt to resend a previously failed message. A proper response from the receiving node would indicate to the reorganized node that the failure was in the switch link. A system like this could be very useful



Network Flow of Fault Free Tree



Network Flow in Presence of Line Fault

Figure 3. Message Pathways

in an intelligent powergrid utilizing a distributed network. This system will aid the distributed grid intelligence of the intelligent power grid to insure more stable reorganizations in the case of

**wide area faults in the power grid. One of the primary issues at hand is distributed leader election of intelligent energy management nodes. A central theme in leader election is dealing with non-determinism in the network. This work by enabling more effective communication will aid in creating node reorganizations that provide a more stable service.**

### 3. API

To aid in the creation and testing of this work, a unified message passing API has been created that facilitates coordination between nodes ranging from the large and sophisticated IEM nodes to small ZigBee devices. The API is a non-blocking, asynchronous approach similar to Active Messages [5] as implemented in Tiny OS [3]. In this message passing API, a device or node registers a message type to receive a message handler. The handler is then used in sending and receiving messages. The current API provides constructs for

- non-blocking send,
- non-blocking receive,
- handle generation,
- conditioned waiting and,
- condition signalling.

Due to frequent faults and large timeouts in a distributed network, non-blocking network abstractions were designed to facilitate this work. This allows devices within the network to send messages without having to wait for acknowledgements before proceeding with other work. The same approach is applied to receives to avoid a need for actively monitoring a queue. In a non-blocking approach, a received message is handled by the network API. When a new message is received, the application is able to use it or ignore it until a later time. It is thus possible to create blocking semantics if desired. This is done using the conditioned wait and condition signalling methods in the API. Through these methods, a running process on a device sends a message and then blocks until begin woken up by the receipt of a response message.

The API is being developed using the Mace distributed prototyping language [2]. Mace is a C++ abstraction that enables the low-level network details to be abstracted from the programmer while leaving significant amounts of flexibility in the message handling abilities and supporting timeout-based fault detection, which is central to our fault tolerance network overlay approach. We have developed a universal FREEDM messaging passing API and a basic prototype of our proposed system on top of Mace.

### 4. Distributed Visualizer

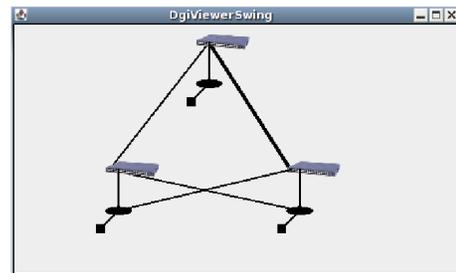
In a distributed network it is important to be able to understand the structure and status of the network. This information may often be difficult to obtain in a distributed network. To aid the maintainers of the system in identifying problems and correcting them we are developing a visualization tool that enables a real-time view of the state of the networked devices in the system and the dynamic routing through our software overlays. This system will provide information regarding a nodes current running status as well as generate a topological layout of the network based on data provided to the running model.

Utilizing a centralized server approach, nodes can communicate with the interface server. The server provides a graphical representation of the status of the nodes and current messages in flight. The projected design of the visualizer will present a fully detailed representation of the underlying LAN. This enables one to monitor system activity, to detect failed components, to observe alternate

routing activity, and to sustain partial functionality in the presence of partitioning / islanding of power microgrids. As such, one can determine which links have failed and, more specifically, which nodes have failed. The information is provided by working nodes that report the status of successful and failed communication attempts within the network.

This work differs from using a commodity tool such as Cacti or Hobbitt in that it will be able to provide the visualization for non-IP devices such as those used in a ZigBee platform that are MAC addressable. This visualizer will interface with IP addressable nodes to detect any MAC-based devices connected to it and display their current status.

In the figure below, the visualizer is being used to display the communication paths of three separate devices. The visualizer can be provided with information detailing the locations of software links between nodes to create a graphical representation of the network. In this figure, the network structure is being coded as a tree network resembling the shape of the network utilized in our Mace prototype.



**Figure 4.** Visualizer Swing Window

The visualization tool is developed using the Java Swing graphics packages over a network socket API to connect with our distributed prototype written in C++ using the Mace distributed program library [2]. In its current status, it supports:

- setting status,
- defining links,
- defining partitions,
- visualizing path, and
- visualizing messages.

The prototype of the message-passing system is instrumented with calls that relay messages during each critical step in the program communication path. At each step, a command message is sent to the visualizer that renders the information on screen. Failed nodes may not be able to report their current status to the system visualizer. In such a case (*i.e.*, when timeouts occur in the system), the node status is updated by the node that receives the timeout.

### 5. Conclusion

The vision of this work is to enable IEM and IFM nodes to communicate, even in the events of multiple link failures. The first step to accomplish this is through introducing increased but intelligent redundancy in the links of the network. We introduced a middleware framework that utilizes software overlays to support fault-tolerant communication in a network with redundancy through alternate communication paths. Our development of low overhead route detection algorithms to assist in the presence of single and multiple link failures constitutes the key contribution to provide such fault tolerance in a transparent manner to other control software. Our middleware layer will provide the means for higher-

level distributed grid intelligence (DGI), such as providing hierarchical control with leader-election schemes within this software overlay architecture. Overall, our software middleware architecture for fault tolerant network overlays realize the vision of sustainable and decentralized energy management on the software side in the FREEDM system and beyond.

## References

- [1] D. G. Andersen, H. Balakrishnan, M. F. Kaashoek, and R. Morris. Abstract the case for resilient overlay networks.
- [2] C. Killian, J. Anderson, R. Braud, R. Jhala, and A. Vahdat. Mace: language support for building distributed systems. In *ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 179–188, 2007.
- [3] P. Levis, S. Madden, J. Polastre, R. Szewczyk, K. Whitehouse, A. Woo, D. Gay, J. Hill, M. Welsh, E. Brewer, and D. Culler. Tinyos: An operating system for sensor networks. pages 115–148. 2005.
- [4] J. Varma, C. Wang, F. Mueller, C. Engelmann, and S. L. Scott. Scalable, fault-tolerant membership for MPI tasks on hpc systems. In *International Conference on Supercomputing*, pages 219–228, June 2006.
- [5] T. von Eicken, D. E. Culler, S. C. Goldstein, and K. E. Schauer. Active messages: a mechanism for integrated communication and computation. In *ISCA '92: Proceedings of the 19th annual international symposium on Computer architecture*, pages 256–266, New York, NY, USA, 1992. ACM.