

Ultimate Forwarding Resilience in OpenFlow Networks

Christopher Hannon
Illinois Institute of Technology
channon@hawk.iit.edu

Jianhui Wang
Argonne National Laboratory
jianhui.wang@anl.gov

Dong Jin
Illinois Institute of Technology
dong.jin@iit.edu

Cheol Won Lee
National Security Research
Institute, South Korea
cheolee@nsr.re.kr

Chen Chen
Argonne National Laboratory
morningchen@anl.gov

Jong Cheol Moon
National Security Research
Institute, South Korea
jcmoon@nsr.re.kr

ABSTRACT

Software defined networking is a rapidly expanding networking paradigm that aims to separate the control logic from the forwarding devices. Through centralized control, network operators are able to deploy and manage more efficient forwarding strategies. Traditionally, when the network undergoes a change through maintenance, failure, or cyber attack, the centralized controller processes these events and deploys new forwarding rules reactively. This work provides a strategy that does not require a controller in order to maintain connectivity while only using features within the existing OpenFlow protocol version 1.3 or greater. In this paper we illustrate why forwarding resiliency is desired in OpenFlow networks and provide an algorithm that computes the flow entries required to achieve maximal forwarding resiliency in presence of both multiple link and controller failures on any arbitrary network.

1. INTRODUCTION

OpenFlow is a standard that defines the communication and interaction between the forwarding and control layers of the software-defined networking (SDN) architecture. The protocol defines the technical specifications required by current SDN switches to support the interaction from OpenFlow controllers to the forwarding switches. SDN provides many features to improve network operations, including enhanced global visibility, greater customizability, and easier deployment. While historically motivated through data center management to optimize performance, minimize cost, and maximize utilization [12], [2], [9], SDN is being proposed for various other networks such as industrial control systems for increased efficiency, enhanced cyber security, and industry specific applications [6], [18].

Industrial control systems, such as the modern power grid, require networks that are resilient to cyber attacks as well as natural disasters that may affect both the communication network as well as the power network. Fault tolerance is a

well studied problem in computer networks. When a network is said to be fault tolerant, a communication channel from source to destination can remain established even if an intermediate link fails. In OpenFlow, a network relies on the central controller to make decisions for the data plane and then install the forwarding logic onto the switches with commands issued from the controller.

One widespread and growing networked industrial control system, the modern power system, is composed of heterogeneous networks, including power line communication (PLC), Ethernet, cellular, and radio communication media. Redundant links are required to enable backup paths for traffic to ensure resiliency. In such heterogeneous networks, SDN is being proposed for resilient communications [3] and as a design architecture [13] to ensure successful and secure operations of the modern power system in the event of cyber attacks and cyber errors. Additionally, other applications have been evaluated utilizing SDN in smart grid communications such as multi-rate multicast for Phasor Measurement Unit communications [8] further motivating the need for resilient communication in such networks that rely on accurate and reliable data streams for control purposes.

Communication resilience in the smart grid is important to the power applications that require high availability and minimal packet loss in the event of communication disruption. For example, power protection elements such as relays require low latency and are highly affected by long recovery times in the event of link failure during management events. As communications in the grid evolve to connectionless UDP traffic as the California ISO reports in [18], high availability becomes a significant design goal to minimize packet loss.

In this paper, we present a forwarding framework with the following objectives: maximal resiliency, high availability, fast recovery time, and ease of deployment. Although fault tolerance is native to SDN, the centralized controller opens up a new interface for cyber attack or failure. Thus, in presence of controller failure, SDN is at a disadvantage compared to traditional decentralized network recovery algorithms such as spanning tree protocols.

Our solution generalizes maximal resiliency for communication between hosts for mission critical traffic within OpenFlow networks. We guarantee that the same number of links may fail as with a controller-based recovery strategy, specifically $|\text{min-cut-st}| - 1$, which is the minimum cut of the network, i.e., fewest number of edges required to fail to create two disjoint subgraphs with s and t on separate subgraphs. Therefore, as long as there exists a path from source to destination, our modified depth-first search (DFS) algorithm,

ACM acknowledges that this contribution was authored or co-authored by an employee, or contractor of the national government. As such, the Government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for Government purposes only. Permission to make digital or hard copies for personal or classroom use is granted. Copies must bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. To copy otherwise, distribute, republish, or post, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SDN-NFV Sec'17, March 22-24, 2017, Scottsdale, AZ, USA

© 2017 ACM. ISBN 978-1-4503-4908-6/17/03...\$15.00

DOI: <http://dx.doi.org/10.1145/3040992.3041002>

Strategy	# Edges needed to remove in order to disconnect s and t	Recovery time	Description
This paper	$ \text{min-cut-st} $	$O(ms)$	requires size of min-cut with s and t on disjoint subgraphs
k-disjoint-backup	$1 \leq k \leq \text{min-cut-st} $	$O(ms)$	can be between 1 and size of $ \text{min-cut-st} $
Backup at every intermediate switch	$2 \leq k \leq \text{min-cut-st} $	$O(ms)$	if a backup path has a link failure, then at minimum 2 edges can be removed; if there are no backup paths then $k = 1 = \text{min-cut-st} $
Controller Based	$ \text{min-cut-st} $	$O(ms)$	requires use of controller
RSTP / STP	$ \text{min-cut-st} $	$O(seconds)$	slow and network cannot perform multi-path forwarding

Table 1: Comparison of the existing failover schemes in event of link failure. Our approach guarantees the same number of links as a controller based recovery scheme as well as traditional decentralized mechanisms. The advantage of our approach is that we do not require a controller to guarantee this performance.

which transforms the network graph to a tree structure, will enable the correct forwarding independent of network topology. This work shows that the same resiliency properties can be established during controller failure as without or compared to decentralized spanning tree solutions. A controller based approach for recovery is a run-time solution for network link failure recovery. In this work, we provide a compile-time recovery strategy that predetermines all backup paths and deploys them proactively. The related work only guarantees resilience in presence of single link failure. Therefore, for mission critical traffic in control systems, our strategy guarantees communication can continue. In this work, we design the algorithm and data structure to convert the graph representation of the network into flow tables and group tables on the OpenFlow switches that enforce the maximally resilient network property.

Section 2 discusses the related work, Section 3 shows how our strategy for providing maximal forwarding resiliency which extends from a single link failure and the challenges associated with multiple link failure. Section 4 presents our all-paths algorithm that changes the graph representation into an N-ary tree representation, which allows us to easily determine the OpenFlow rules necessary for the network. Section 5 compares our strategy with related work, while Section 6 concludes our work and proposes future research directions for resilient forwarding in OpenFlow networks.

2. RELATED WORK

Many of the current strategies for fault tolerance and link failure rely on the central controller to calculate the optimal path from source to destination using a global view of the network topology updated with the failed edge. In [14] and [10] the authors show how to utilize the controller to update the network in event of link failure. When a link fails, the event is sent to the controller through the OpenFlow protocol, the controller then computes new paths for each flow and modifies the necessary flow rules on the switches to create new flows. This is the most straightforward approach using a centralized control architecture. Calculating the new optimal path however requires time to process the updated graph and time to install the new forwarding logic onto the corresponding switches. During this time packets may be lost.

In [17] the authors utilize a hybrid approach through backup routes and centralized recomputation by using two steps. In the first step, the forwarding devices locally detect a fail-

ure and resort to a failure mode where precomputed backup paths are taken. The second step is for the controller to compute the optimal paths based on that link failure. This ensures that a functional while non-optimal path can take over faster than the controller to recompute a new path. The authors extend the Open vSwitch to include the Bidirectional Forwarding Detection to report in the fast failover liveness check. The design for the fast failover rules is to provide an alternative path, i.e., backup port to send the traffic to in the event of a failure. If the switch does not have a backup path, then a backtracking method is used to return the packet to the previous switch, where the process is repeated recursively. The authors compute backup paths so that each intermediate hop computes a backup path locally, implemented using the failover group table entry. If no backup path exists, the traffic gets sent back one hop until a backup path is located. They however do not consider the case of multiple link failures as our work does. The minimum number of link failures to disconnect the source and destination is two.

Similarly, [1] also implements a hybrid approach that uses the same mechanisms of [17] and show they can achieve recovery in 50 ms. These works share the same limitations of single link failure.

For the reliability of SDN controllers, ONOS [4], and [5], illustrate the benefits of distributed controllers. One such benefit is reliability during controller failure. In our work, we consider failure from a single SDN controller.

This paper presents a strategy that computes every possible path from source to destination and creates a flow table and group tables that provide ultimate resiliency in the face of multiple link failures. The contributions of the algorithm in this paper are as follows:

- rules are proactively installed, ensuring that the recovery process does not require a controller
- the approach enables multiple link failure resiliency in the network
- the approach minimizes packet loss in the event of multiple link failures

In this paper, we describe an algorithm that enables maximal resilience to link failures and evaluate the effect it has on the switches' forwarding tables, as well as connectivity and latency in the network. This algorithm requires no changes to the OpenFlow protocol and is computed offline and proactively installs rules onto the switches.

3. MOTIVATION

This section considers a single path from source to destination rather than every flow on the network simultaneously. The s node represents the source while the t node is the destination node. G represents the network with switches as nodes and links as edges. This notation will be used throughout the paper. The following examples provide the motivation and explain the mechanisms of the maximally resilient forwarding algorithm described in the next section which generalizes the process for arbitrary networks.

3.1 Single Link Failure Recovery

During a link failure, the switches with ports connected to the failed link, detect the failure. Upon detection, fast failover group tables will allow for a different rule set to be placed on the packets destined to failed link using the liveness monitor feature. Group table actions consist of sending the packet to a new port, a new group table, or even back to the flow tables [16]. In this work, the group table for fast failover is used to monitor the outgoing port and forwards to a backup in case of failure on the main path.

In the network depicted in Figure 1, the main path is $s \rightarrow s1 \rightarrow s2 \rightarrow s3 \rightarrow s4 \rightarrow s8 \rightarrow t$. When the edge from $s1$ to $s2$ fails, and switch $s1$ detects this through its liveness monitor for port y in the failover table and resorts to port x . Therefore the traffic can now be routed through $s \rightarrow s1 \rightarrow s5 \rightarrow s6 \rightarrow s7 \rightarrow s8 \rightarrow t$ and can do this without any reactive influence from a centralized controller. The link failure can be detected using a mechanism such as bidirectional forwarding detection.

Only relying on failover tables to chose an alternate route locally may not satisfy every case. For instance, if a different edge experiences a failure ($s2 \rightarrow s3$ on the primary path), $s2$ does not have any backup routes locally as it only has a degree of 2. Therefore a more complex scheme must be used involving backtracking. This backtracking is inspired from the crankback algorithm used in MPLS. Crankback [7] allows for sending a message backward when setting up a connection in routing between networks. By sending the message backwards, a new path can be established that meets the requirements of the flow at the cost of latency and bandwidth consumption.

In our case, if a path becomes infeasible (i.e., no failover paths locally) the packet gets forwarded back to the previous switch and is forwarded from there. The rule for a returned packet is the same as it would be as if the edge it just returned on was failed. This process of crankback can happen recursively if there are no backup links to send to. If the edge ($s2 \rightarrow s3$) fails, $s2$ requires a group table to monitor the liveness of port x and send through y as a backup. $s1$ would need a rule in the flow table to match on incoming port y and send to x .

In this case, the traffic flows from $s \rightarrow s1 \rightarrow s2 \rightarrow s1 \rightarrow s5 \rightarrow s6 \rightarrow s7 \rightarrow s8 \rightarrow t$. Although $s1$ must process the packet twice, the information stream from s to t remains intact and does this without reactive influence from a controller. Additionally, if the link from $s2 \rightarrow s3$ comes back online, this optimal path would be taken automatically without any configuration changes due to the liveness status of group tables. In all the related work, the network resiliency without controller influence stops here, however the examples presented only experience a single link failure. Our work shows how this approach can be extended for link failure in

$|min - cut(G, s, t)| - 1$ edges regardless of network size or topology where $|min - cut(G, s, t)|$ is the size of the min-cut where s and t are on separate disjoint subgraphs of G .

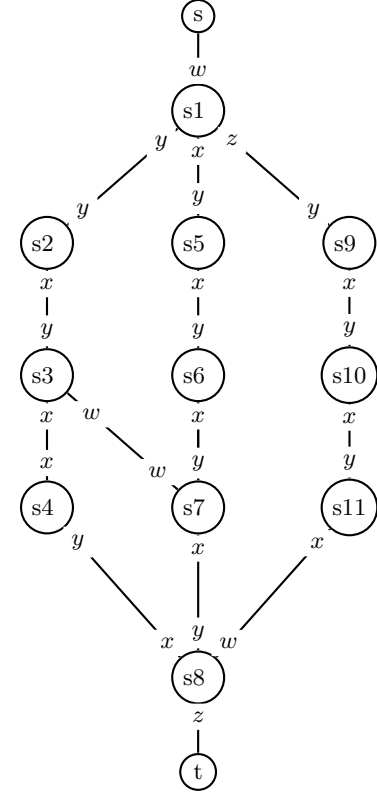


Figure 1: The main path is $s \rightarrow s1 \rightarrow s2 \rightarrow s3 \rightarrow s4 \rightarrow s8 \rightarrow t$. Case 1: $s1 \rightarrow s2$ failed. In this case, the group table installed on $s1$ of type failover will be sent to the next live port, x , to complete the flow $s \rightarrow s1 \rightarrow s5 \rightarrow s6 \rightarrow s7 \rightarrow s8 \rightarrow t$. Case 2: $s2 \rightarrow s3$ failed. In this case, the group table installed on $s2$ returns the packet to $s1$ in crankback style routing. Case 3: Multiple link failure. $s4 \rightarrow s8$ and $s7 \rightarrow s8$ failed. In this case, the packet has the ultimate path $s \rightarrow s1 \rightarrow s2 \rightarrow s3 \rightarrow s4 \rightarrow s3 \rightarrow s7 \rightarrow s3 \rightarrow s2 \rightarrow s1 \rightarrow s9 \rightarrow s10 \rightarrow s11 \rightarrow s8 \rightarrow t$. Table 2 shows rules to perform this forwarding. Case 4: Multiple link failure. If edges $s2 \rightarrow s3$, $s7 \rightarrow s8$ and $s11 \rightarrow s8$ failed, the rules in Table 2 do not find the successful path. A more general approach is necessary.

3.2 Multiple Link Failure Recovery

The introduction of loops and more complicated networks requires additional logic when faced with multiple link failures. For example, consider if the network in Figure 1 contains multiple link failures, $s4 \rightarrow s8$ and $s7 \rightarrow s8$ with corresponding flow and group tables defined in Table 2. Switches $s2, s4, s9, s10, s11$, all follow the same format as the table $s1$:FT1 and are excluded for brevity in Table 2. The failover rule returns the packet to the previous switch to support crankback routing.

In this network the main path from s to t is $s \rightarrow s1 \rightarrow s2 \rightarrow s3 \rightarrow s4 \rightarrow s8 \rightarrow t$. The path from $s4 \rightarrow s8$ is failed, therefore $s4$ routes the traffic back according to its failover group. $s3$ has a rule to route $s3 \rightarrow s7 \rightarrow s8 \rightarrow t$. However, upon reaching $s7$, port x reports that the edge from $s7 \rightarrow s8$ is also down. The fast failover rule routes the frame back to

Table 2: Switch Configurations for the Example in Figure 1

Flow Table s_1 :FT1			Flow Table s_2 :FT1			Flow Table s_7 :FT1		
FlowID	Match Field	Action	FlowID	Match Field	Action	FlowID	Match Field	Action
Flow 1	InPort: w	Group GT1	Flow 1	InPort: y	Group GT1	Flow 1	InPort: w	Group GT1
Flow 2	InPort: y	OutPort: z	Flow 2	InPort: x	OutPort: y	Flow 2	InPort: x	OutPort: w
Group Table s_1 :GT1 Type: FF			Group Table s_2 :GT1 Type: FF			Group Table s_7 :GT1 Type: FF		
Bucket ID	WatchPort	Action	Bucket ID	WatchPort	Action	Bucket ID	WatchPort	Action
Bucket 1	y	OutPort: y	Bucket 1	x	OutPort: x	Bucket 1	x	OutPort: x
Bucket 2	z	OutPort: z	Bucket 2	y	OutPort: y	Bucket 2	w	OutPort: w
Flow Table s_3 :FT1			Group Table s_3 :GT1Type: FF			Flow Table s_8 :FT1		
Flow ID	Match Field	Action	Bucket ID	WatchPort	Action	Bucket ID	WatchPort	Action
Flow 1	InPort: y	Group GT1	Bucket 1	x	OutPort: x	Flow 1	InPort: *	OutPort z
Flow 2	InPort: x	OutPort w	Bucket 2	w	OutPort: w			
Flow 3	InPort: w	OutPort y	Bucket 3	y	OutPort: y			

s_3 according to its failover rule. s_3 has nowhere to send the packet so using Flow 3 returns it to s_1 via s_2 . At this point since the path between s_5 and t (without going through s_1) is blocked, the switch s_1 should forward the frame to $s_9 \rightarrow s_{10} \rightarrow s_{11} \rightarrow 8 \rightarrow t$. Therefore, s_1 has rule FT2 to perform this action. In this case, the packet is able to arrive at t with an ultimate route of $s \rightarrow s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_7 \rightarrow s_3 \rightarrow s_2 \rightarrow s_1 \rightarrow s_9 \rightarrow s_{10} \rightarrow s_{11} \rightarrow s_8 \rightarrow t$. While this is not the most direct route, this strategy successfully enables forwarding in the presence of two failed links on the primary and backup paths.

This strategy is specifically tailored to work with the failed edges. The previous rule set fails if the following edges ($s_2 \rightarrow s_3$, $s_7 \rightarrow s_8$, $s_{11} \rightarrow s_8$) are down, even though there exists a path $s \rightarrow s_1 \rightarrow s_5 \rightarrow s_6 \rightarrow s_7 \rightarrow s_3 \rightarrow s_4 \rightarrow s_8 \rightarrow t$. s_7 can check to see if the route through s_3 is a candidate, and a few additional rules can be added to the current strategy allowing for this path to be tried.

However a problem quickly emerges. With the additional rules, specifically the $s_7 \rightarrow s_3$, an infinite loop is easily created. A loop can occur if the link from $s_4 \rightarrow s_8$ is down. The usual behavior is to send this packet back to s_2 using the crankback technique. Eventually this packet will return to s through $s_3 \rightarrow s_2$ or be dropped if that link is down. This brings motivation for adding a mechanism to remember the history of the switches that have already processed the packet to treat the switch differently depending on the source to ensure a proper crankback. The challenge is that OpenFlow does not support switch-to-switch communication for this problem. We utilize the packet header space, specifically the VLAN id field as a ledger, which we show in the following section.

4. ALGORITHM

We present a modified depth-first-search algorithm that finds all paths from source to destination recursively. In this modified DFS algorithm, the data structure created is an N-ary tree that represents all paths that can lead to the destination node. In this N-ary tree, the root node is the source and each leaf is the destination t . Each intermediate node has a set of children ≥ 1 that can reach the destination t . Not every edge from an intermediate node is represented

in the data structure nor necessarily all nodes as some may not have a path to the destination. Loops in the network may introduce multiple paths p_1, p_2 from source to destination, where $p_1 \cap p_2 \neq \emptyset$. This implies that p_1 and p_2 are not disjoint. Such occurrences create duplicate nodes in the tree that are marked accordingly.

The first path on node n is from the parent to child 1, if child 1 port is down, then it goes to child 2 and this process continues until there are no children remaining of n . If this is the case, there is no path from n to t so the packet is sent back to the parent in the style of crankback routing. Similarly, if there is a downstream link down from $n.child_k$ to t , then the packet returns to n and should be sent to $n.child_{k+1}$. However, a different fast failover table must be used to ensure that if port $n.child_{k+1}$ is down that the fast failover sends the packet to $n.child_{k+2}$ if possible. Therefore, there exists a flow entry for each port sending it to a fast failover table. The flow table for intermediate node i follows the format in Table 3.

As seen in Table 3, each InPort sends the packet to a different group Table. Because failover rules can only forward in case of local link failures, there needs to be a different group table for fast failover for each port on a switch to prevent loops. The group table 1 on a node m is the table that is arrived from the port connected to the parent of m . Subsequent group tables are for each child of m . Each child's group table is a subset of group table 1. For example, if a node m has three children, the first failover rule in group table 1 will be child 1 then child 2 then child 3. If downstream of child 1 there is a link failure and the packet returns to the node m , the packet should be forwarded to child 2. Therefore, the group table for m that matches child 1's incoming port is group table 1 without the entry for child 1. If this link is broken, the fast failover should send it to child 3. If the original fast failover group is used, then the packet would be sent back down to child 1, where it is known that there is no path thus causing a loop. Hence, the group table for child 2 is group table 2 without the entry for child 2. Tables 4-5 illustrate the general case for forwarding with multiple children of a node with a total degree of $n + 1$.

Many networks that enable resiliency contain loops. The major difficulty for multiple paths routing through loops is

Algorithm 1 Modified DFS Algorithm

```

1: function ALGO( $G = (V, E), s, t, M$ )
2:   for all  $v \in N(s)$  do
3:     if  $v \in T$  then ▷ Node visited before
4:       if  $v \in M$  then ▷ Case 1: backedge
5:         Return
6:       else if  $v \in R$  then ▷ Case 2: crossed edge
7:          $Tag(v, s, M)$ 
8:       else ▷ Case 3: no path to  $t$ 
9:         Return
10:      end if
11:    else
12:       $M = M \cup \{v\}$  ▷ Add  $v$  to  $M$ 
13:       $M[s].child = v$ 
14:      if  $v = t$  then ▷ Case 4: reached  $t$ 
15:         $R = R \cup M$  ▷ Add  $M$  to  $T$ 
16:      else ▷ Case 5: unexplored node
17:         $T = T \cup \{v\}$ 
18:         $Algo(G, v, t, M)$ 
19:      end if
20:    end if
21:  end for
22: end function
23:
24:  $R = \{ \}$  ▷ N-ary tree, contains all paths
25:  $M = \{ \}$  ▷ tmp working tree
26:  $T = \{ \}$  ▷ list of marked nodes
27:  $s = \text{source}$ 
28:  $t = \text{destination}$ 
29:  $Algo(G, s, t, M)$ 
30:
31: function TAG( $v, s, M$ )
32:    $A = \text{copy}(v)$  ▷ creates copy of  $v$ 
33:    $M = M \cup \{A\}$  ▷ add  $A$  to  $M$ 
34:    $M[s].child = A$ 
35:    $Copy\_down(v, A, M)$  ▷ child of  $v$  will reach  $t$ 
36:    $Copy\_up(v, A, M)$  ▷ parent of  $v$  may reach  $t$ 
37: end function
38:
39: function COPY_DOWN( $v, A, M$ )
40:   for all  $v' \in \text{children}(v) \notin M$  do
41:      $A' = \text{copy}(v')$ 
42:      $M[A].child = A'$ 
43:     if  $v = t$  then
44:        $R = R \cup M$ 
45:     else
46:        $Copy\_down(v, A', M)$ 
47:     end if
48:   end for
49: end function
50:
51: function COPY_UP( $v, A, M$ )
52:   if  $v.parent \in M$  then
53:     Return
54:   else
55:      $A' = \text{copy}(v.parent)$ 
56:      $M[A].child = A'$ 
57:      $Tag(A', v.parent, M)$  ▷ treat as crossed edge
58:   end if
59: end function

```

that there are already rules on the switches making it difficult to keep track of which routes have been traversed before. In the Appendix, we present the full algorithm and technique to utilize the VLAN id field as a ledger between switches.

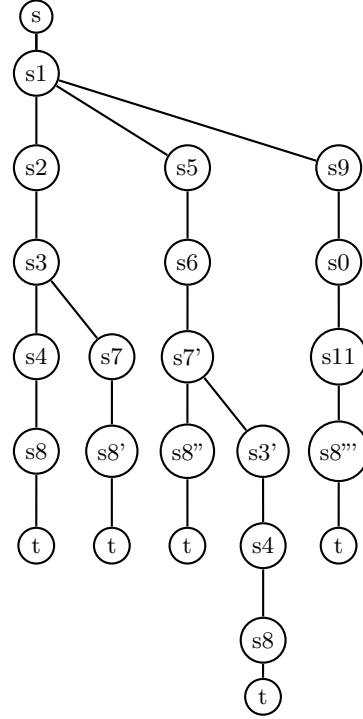


Figure 2: The resulting N-ary tree from Algorithm 1 on the network in Figure 1. The nodes that are depicted with ' (prime) are backedges. The rule from a node with child and with prime notation checks the VLAN id field to see if the switch has processed the packet previously; if not, uses the VLAN tag to return the packet to the correct parent, thus preventing forwarding loops.

Table 3: Flow Table for intermediate node si

FlowID	Match Field(s)	Action
FT1	InPort: $si.port.parent$	Group $GT1$
FT2	InPort: $si.port.child_1$	Group $GT2$
FT $k+1$	InPort: $si.port.child_k$	Group $GTk + 1$
FT $n+1$	InPort: $si.port.child_n$	OutPort $si.port.parent$

Table 4: Group Table 1 for switch sn

Group Table ID: $GT1$		Type: FF
Bucket ID	WatchPort	Actions
Bucket 1	$sn.port.child_1$	OutPort: $sn.port.child_1$
Bucket i	$sn.port.child_i$	OutPort: $sn.port.child_i$
Bucket n	$sn.port.child_n$	OutPort: $sn.port.child_n$
Bucket n+1	$sn.port.parent$	OutPort: $sn.port.parent$

Table 5: Group Table k for switch sn

Group Table ID: $GT2$		Type: FF
Bucket ID	WatchPort	Actions
Bucket 1	$sn.port.child_k$	OutPort: $sn.port.child_k$
Bucket i-k	$sn.port.child_i$	OutPort: $sn.port.child_i$
Bucket n-k	$sn.port.child_n$	OutPort: $sn.port.child_n$
Bucket n-k+1	$sn.port.parent$	OutPort: $sn.port.parent$

5. EVALUATION

We implement three forwarding algorithms using the Ryu controller [15] in Mininet [11] for comparison. Algorithm 1 is the shortest path algorithm that offers no resilience to link failure in the primary path. Algorithm 2 is the proactive technique from related work [17] generalized for arbitrary network sizes. Algorithm 3 is the solution provided in this paper. To understand the performance under general link failure, we evaluate the connectivity under *every* combination of failed links from 1 to 10.

In this work, we do not consider the transient state of links failing, we consider the steady state where each link may be in one of two possible states: working and failed. If links are reported as failed and later reported as alive again, then the subsequent packets will be processed as normal because no rule is ever deleted, i.e., timeout on all rules set to infinity. The controller proactively installs all flow entries and group tables onto the switches and then shuts down. Links are brought offline using the Mininet API and the connectivity and latency test is measured using the Ping tool. The experimental results are shown in Table 6 and Figure 3, and it is clear that our algorithm is most resilient to link failures. Under two link failures, our algorithm has an extended ping time but does not lose connectivity. Before the number of failed links reaches 10 (out of 13 links in total), our algorithm always results in more working scenarios than the other two algorithms do.

Given the graph in Figure 1, the size of the minimum cut of the graph with s and t on disjoint subgraphs is 3, if we exclude the link $s \rightarrow s1$ and $s8 \rightarrow t$. This means that the minimum number of links that can be removed to disconnect the source and destination is 3 links, i.e., $s1 \rightarrow s2$, $s5 \rightarrow s6$ and $s11 \rightarrow s8$. The maximum number of edges that can be removed is 9. As long as there exists a path from source to destination, our forwarding strategy will always find the path. Compared to a controller based recovery strategy, our strategy has equal performance in terms of link failures, and we can guarantee this at compile-time rather than at run-time since our backup strategy is installed proactively. In addition to SDN based strategies, the traditional mechanisms for network recovery, namely spanning tree protocols, can guarantee the same performance, but take much longer

# Fail Links	# Combinations	Algo 1		Algo 2		Algo 3	
		#	%	#	%	#	%
1	13	9	69.2	13	100	13	100
2	78	36	46.2	58	74.4	78	100
3	286	84	29.4	139	48.6	253	88.5
4	715	126	17.6	210	29.4	449	62.8
5	1287	126	9.8	213	16.6	486	37.8
6	1716	84	4.9	146	8.4	334	19.5
7	1716	36	2.1	65	3.8	144	8.4
8	1287	9	0.7	17	1.3	36	2.8
9	715	1	0.1	2	0.3	4	1
10	286	0	0	0	0	0	0

Table 6: Given # Failed Links, # Combination is the total number of possible network scenarios. Below each Algo, # indicates the number of scenarios with a working forwarding path, and % indicates the percentage of working scenarios. Comparison of the algorithms show that our technique has the greatest resilience to link failures, i.e., outperforms the other two algorithms when the number of failed links $\in [2, 9]$.

on the order of seconds to re-converge after a link failure. For mission critical traffic this is too long, our recovery strategy using OpenFlow meets the objective of high availability and fast recovery time. Comparing our strategy in this example with the related work, the backup paths approach [17] $s1 \rightarrow s8$ only finds three paths, while our scheme provides all 5 backup paths.

On the other hand, because no rules are ever added or deleted (unless administratively), there will be greater overhead in the number of rules needed to be installed on a given switch. For each switch v in the network, there will be $deg v$ flow table entries and group tables required for a single flow $s \rightarrow t$. Crossedges will require additional $deg v$ group tables and flow entries to enable correct forwarding.

The total number of group tables and flow entries is at most $2E$, which is at maximum $O(n^2)$ flow entries and group tables for highly connected networks. Additionally, due to the crankback style of routing, there will be additional traffic within the network, because a packet never returns to the same switch except during crankback. If $flow(s, t)$ represents the quantity of traffic sent from s to t , then at maximum each link can experience $2 * flow(s, t)$ traffic. Finally, due to the crankback routing and sub-optimal path selection, if there remains a path between s and t , then the latency is at most

$$O\left(\sum_{e \in E} 2 * prop(e) + \sum_{v \in V} deg(v) * process(v)\right)$$

$prop$ is the propagation delay in the links, $process$ is the processing delay at the switch including the queue delay. The maximum number of rules in our approach is on $s3$ with 6 group tables and 6 flow table entries for one host-to-host flow, algorithm 2 requiring 4 group tables and 4 flow table entries, and algorithm 1 requiring only 2 flow entries.

6. CONCLUSION

We have presented a general scheme to implement resilient forwarding in OpenFlow networks with minimal packet loss. Our scheme works without reactive influence from a controller, enabling greater resilience than traditional SDN approaches. Our algorithm provides the OpenFlow rules for multiple link failures equivalent to a centralized controller approach enabling the hosts to remain connected. We show with an example that our approach has greater resilience for all link failures. Due to the overhead incurred by this method of forwarding in presence of multiple link failures, it is designed for mission critical communication applications as to prevent congestion within the network. Our future work includes researching techniques to reduce the overhead in terms of number of tables on the switches for multiple source-destination pairs through rule compression, increasing the scalability of the VLAN ledger mechanisms, and evaluating the trade-offs between overhead and resiliency by formulating this problem as an optimization problem.

7. ACKNOWLEDGMENTS

This work is partly sponsored by the Maryland Procurement Office under Contract No. H98230-14-C-0141, the Air Force Office of Scientific Research (AFOSR) under grant FA9550-15-1-0190, the U.S. Department of Energy (DOE)'s Office of Electricity Delivery and Energy Reliability, and a cooperative agreement between IIT and National Security Research Institute (NSRI) of Korea. Any opinions, findings and conclusions or recommendations expressed in this

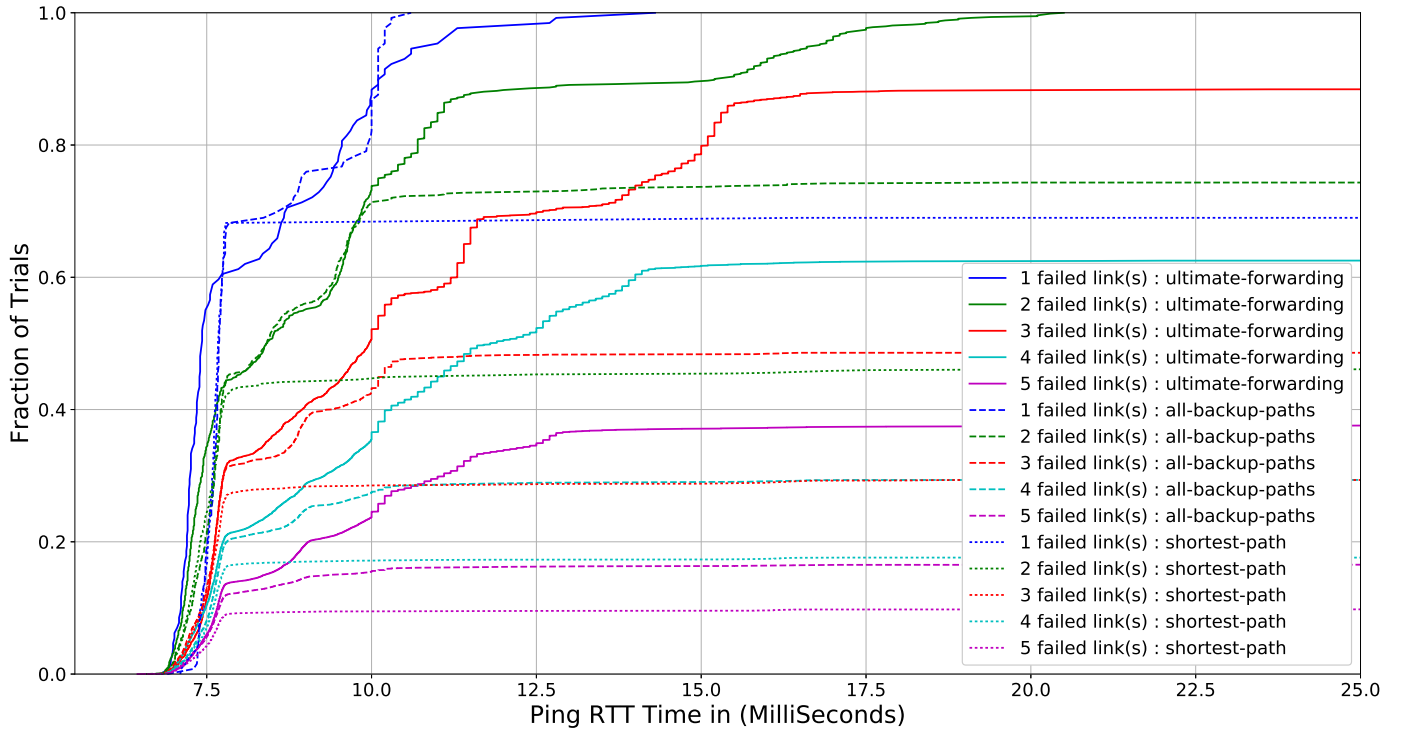


Figure 3: Comparison of latency under link failure of the three implemented algorithms

material are those of the author(s) and do not necessarily reflect the views of the Maryland Procurement Office, AFOSR, DOE and NSRI. The authors also thank Xu Yang for his technical assistance.

8. REFERENCES

- [1] R. Ahmed, E. Alfaki, and M. Nawari. Fast failure detection and recovery mechanism for dynamic networks using software-defined networking. In *2016 Conference of Basic Sciences and Engineering Studies (SGCAC)*, pages 167–170, Feb 2016.
- [2] M. Al-Fares, S. Radhakrishnan, B. Raghavan, N. Huang, and A. Vahdat. Hedera: Dynamic flow scheduling for data center networks. In *Proceedings of the 7th USENIX Conference on Networked Systems Design and Implementation, NSDI’10*, pages 19–19, Berkeley, CA, USA, 2010. USENIX Association.
- [3] A. Aydeger, K. Akkaya, M. H. Cintuglu, A. S. Uluagac, and O. Mohammed. Software defined networking for resilient communications in smart grid active distribution networks. In *2016 IEEE International Conference on Communications (ICC)*, pages 1–6, May 2016.
- [4] P. Berde, M. Gerola, J. Hart, Y. Higuchi, M. Kobayashi, T. Koide, B. Lantz, B. O’Connor, P. Radoslavov, W. Snow, et al. Onos: towards an open, distributed sdn os. In *Proceedings of the third workshop on Hot topics in software defined networking*, pages 1–6. ACM, 2014.
- [5] A. Dixit, F. Hao, S. Mukherjee, T. Lakshman, and R. Kompella. Towards an elastic distributed sdn controller. In *ACM SIGCOMM Computer Communication Review*, volume 43, pages 7–12. ACM, 2013.
- [6] X. Dong, H. Lin, R. Tan, R. K. Iyer, and Z. Kalbarczyk. Software-defined networking for smart grid resilience: Opportunities and challenges. In *Proceedings of the 1st ACM Workshop on Cyber-Physical System Security, CPSS ’15*, pages 61–68, New York, NY, USA, 2015. ACM.
- [7] A. Farrel, A. Satyanarayana, A. Iwata, N. Fujita, and G. Ash. Crankback signaling extensions for mpls and gmpls rsvp-te. *Network Working Group*, RFC 4920, 2007.
- [8] A. Goodney, S. Kumar, A. Ravi, and Y. H. Cho. Efficient PMU networking with software defined networks. In *2013 IEEE International Conference on Smart Grid Communications (SmartGridComm)*, pages 378–383, Oct 2013.
- [9] B. Heller, S. Seetharaman, P. Mahadevan, Y. Yakoumis, P. Sharma, S. Banerjee, and N. McKeown. Elastictree: Saving energy in data center networks. In *Proceedings of the 7th USENIX Conference on Networked Systems Design and Implementation, NSDI’10*, pages 17–17, Berkeley, CA, USA, 2010. USENIX Association.
- [10] H. Kim, M. Schlansker, J. R. Santos, J. Tourrilhes, Y. Turner, and N. Feamster. Coronet: Fault tolerance for software defined networks. In *2012 20th IEEE International Conference on Network Protocols (ICNP)*, pages 1–2, Oct 2012.
- [11] B. Lantz, B. Heller, and N. McKeown. A Network in a Laptop: Rapid Prototyping for Software-defined Networks. In *Proceedings of the 9th ACM SIGCOMM*

- Workshop on Hot Topics in Networks*, Hotnets-IX, pages 19:1–19:6, New York, NY, USA, 2010. ACM.
- [12] H. H. Liu, X. Wu, M. Zhang, L. Yuan, R. Wattenhofer, and D. Maltz. zupdate: Updating data center networks with zero loss. In *ACM SIGCOMM Computer Communication Review*, volume 43, pages 411–422. ACM, 2013.
 - [13] S. Rinaldi, P. Ferrari, D. Brandão, and S. Sulis. Software defined networking applied to the heterogeneous infrastructure of smart grid. In *2015 IEEE World Conference on Factory Communication Systems (WFCS)*, pages 1–4, May 2015.
 - [14] S. Sharma, D. Staessens, D. Colle, M. Pickavet, and P. Demeester. Enabling fast failure recovery in openflow networks. In *the 8th International Workshop on the Design of Reliable Communication Networks (DRCN)*, pages 164–171, Oct 2011.
 - [15] R. P. Team. Component-Based Software Defined Networking Framework. <https://osrg.github.io/ryu/>, 2017. Online; accessed 23-Jan-2017.
 - [16] The Open Networking Foundation. OpenFlow Switch Specification, Jun. 2016.
 - [17] N. L. Van Adrichem, B. J. Van Asten, and F. A. Kuipers. Fast recovery in software-defined networks. In *the Third European Workshop on Software Defined Networks*, pages 61–66, Sept 2014.
 - [18] T. Williams. The role of software-defined networking and wan virtualization in securing scada systems: An alternative to pki and vpns for securing telemetry data. <https://www.youtube.com/watch?v=5w9eL3asTrQ>. Online; accessed 23-Jan-2017.