

Ground Control to Major Faults: Towards a Fault Tolerant and Adaptive SDN Control Network

Liron Schiff[†] Stefan Schmid[‡] Marco Canini^{*}

[†]Tel Aviv University [‡]Aalborg University & TU Berlin ^{*}Université catholique de Louvain

Abstract—To provide high availability and fault-tolerance, SDN control planes should be distributed. However, distributed control planes are challenging to design and bootstrap, especially if to be done in-band, without dedicated control network, and without relying on legacy protocols. This paper promotes a distributed systems approach to build and maintain connectivity between a distributed control plane and the data plane. In particular, we make the case for a *self-stabilizing* distributed control plane, where from any initial configuration, controllers self-organize, and quickly establish a communication channel among themselves. Given the resulting managed control plane, arbitrary network services can be implemented on top.

This paper presents a model for the design of such self-stabilizing control planes, and identifies fundamental challenges. Subsequently, we present techniques which can be used to solve these challenges, and implement a plug & play distributed control plane which supports automatic topology discovery and management, as well as flexible controller membership: controllers can be added and removed dynamically. Interestingly, we argue that our approach can readily be implemented in today’s OpenFlow protocol. Moreover, our approach comes with interesting security features.

I. INTRODUCTION

In traditional networks, the control software is distributed across all devices, which run routing protocols to compute forwarding state. An advantage of this design is that legacy networks can use in-band control: that is, control plane packets are carried over the same data plane network as with the regular traffic. Because legacy routing protocols are designed to exchange data between neighboring routers and do not require global state, they are self-stabilizing in the sense that they can automatically bootstrap the network and converge to a valid operating state from any initial conditions.

In contrast, Software Defined Networking (SDN) advocates for software-based control of a network based on a logically-centralized global network view. A number of research and development efforts have illustrated that this approach provides several benefits in terms of improved flexibility and performance, easy of management and decreased operational complexity, and lower costs [1]–[5]. While the literature has well articulated several benefits of the separation between control and data planes, the question of how connectivity between these planes is maintained – that is, paths for supporting switch-controller or controller to controller communication – has not received much attention. In practice, most SDN deployments [4]–[6] use out-of-band control, where control plane packets are carried by a dedicated management network. The management network runs its own routing system, which

typically is realized using traditional routing protocols such as STP or OSPF.

In-band control is desirable for several reasons including its economical benefits (in certain contexts, such as carrier networks, out-of-band control would be prohibitively expensive). In essence, with in-band control there is no reason to build, operate, and ensure the reliability of a separate network. Also, out-of-band networks are typically underprovisioned and have limited redundancy. In these conditions, congestion losses are possible and link failures can lead to a partitioned SDN control plane even though the data plane network is not disconnected. This raises several concerns regarding the availability of the SDN architecture — for instance, can we guarantee that, if the data plane network is a connected graph, the SDN control plane will always establish a route between every pair of nodes in the network?

Moreover, in-band control in SDN is an appealing prospective for transitioning from legacy networks [7]. However it is also a challenging one, as others have noted, “*by extricating the control plane out of network devices and implementing it using a distributed system, SDN inherits the weaknesses associated with reliable distributed services*” [8]. These challenges are further amplified by the fact that the logically-centralized SDN control plane actually means physically distributed, for a plethora of reasons, including reliability, availability, scalability and latency [9]–[11].

This paper makes a first step towards the design of a fault-tolerant (namely self-stabilizing), distributed SDN control plane. We first present a model that captures the underlying fundamental challenges of self-stabilizing control networks. We then build upon this model to reason about and develop basic self-stabilizing mechanisms. Our main contribution is the design of an in-band control system that coordinates distributed controllers in a self-organizing manner, to bootstrap a connectivity service between controllers and switches. Our approach is based on a plug & play paradigm: it supports a dynamic membership of controllers, in the sense that new controllers can be added to the network, or old controllers removed from the network, at any time and without explicit notice; despite these changes, our approach will ensure that every unmanaged switch is assigned to one controller, and that a communication network between controllers is established. Once this “bootstrap problem” has been solved, SDN control planes can be realized on top of our system. We note that our approach, unlike prior work [12], is self-reliant and represents an “SDN-only” approach: it is based on OpenFlow only, and does not depend on any legacy protocols.

II. MOTIVATION

The SDN control plane, that is, a collection of network-attached servers, must have connectivity to the data plane. The question we explore in this paper can be stated as: is it possible to rely on the same and already existing primitives by which SDN applications govern the network to bootstrap and provide connectivity between control and data planes? We find a positive answer by building upon the concept of self-stabilization. A distributed system that is self-stabilizing will end up in a correct state independently of its initial state [13], [14]. That correct state is reached after a finite number of steps (the convergence time).

We advocate for a software-driven, in-band control mechanism to support the broader goal of building distributed SDN control planes. To make this case, let us consider what services are required by SDN controllers such as Onix [9] or STN [10].

Connectivity: To allow communication between controllers and switches, and between controllers.

Controller discovery: To allow individual controllers discover other existing controllers and detect when some are no longer reachable.

Switch discovery: To detect switches that are not yet associated with a controller and establish a control channel with them. Also, to re-establish communication with switches in case of link failure, network partitions, or controller failures.

Security: To ensure controller and switch authentication as well as encryption of control traffic. In addition to a secure bootstrap, control traffic may have to be given priority over other traffic, *e.g.*, to prevent some types of DoS attacks.

Self-stabilization is a natural approach to meet these goals while coping with dynamic conditions, such as arrival and departure of controllers [15], arbitrary topology changes (switch or link failures), and communication errors (packet losses or delays). We emphasize that, as usual in the framework of self-stabilizing algorithms, anytime before, during, and after a self-stabilization, additional changes may occur: the system will simply reconverge starting from the current state.

III. MODEL AND CHALLENGES

The design of a self-stabilizing mechanism to build and maintain connectivity between a distributed control plane and the data plane raises some fundamental challenges. This section introduces a model that captures some of the particularities of such a system, and which serves as a basis for reasoning about self-stabilizing mechanisms. As we will see, the postulated mechanisms can be implemented in today's OpenFlow protocol, and lie at the heart of our control plane, which instantiates the principles presented in this paper.

We consider a network connecting two kinds of components, *controllers* and *switches*, henceforth also called *nodes*:

- 1) **Controllers:** must be able to communicate and coordinate. To enable these, we assume each controller has a unique address, but controllers also respond to a virtual anycast address. Controllers must be able to solve consensus. In principle, a controller can perform arbitrary local computations (*i.e.*, it is a Turing machine).

- 2) **Switches:** to communicate, they also have unique addresses. Switches are simple and “passive” devices: they cannot perform any computations on their own. In principle, a switch can be in one of two possible states: *managed* or *unmanaged*. When unmanaged, a switch periodically attempts to connect with a controller using a pre-configured remote controller address, *i.e.*, the virtual anycast address. The configuration of a switch is defined by a list of match-action rules with priorities and (possibly infinite) timeout values. The configuration can be changed either (1) by a controller or (2) through a timeout.

The fundamental problem arises from the tension of what a (passive) switch needs and what an (active) controller can do. Concretely, we identify the following challenges for the design of a self-stabilizing distributed control plane:

- C1) **Fault-Tolerance:** The system should tolerate any failures which may lead one or multiple switches to become unmanaged: link failures, controller failures, TCP timeouts, congestion, etc. A central challenge here regards the detection of inactivity: a switch must notice when it is unmanaged. This is difficult in our model where switches are passive.
- C2) **Consensus:** Without consensus, a control plane protocol may never converge: A switch must eventually be managed by exactly one controller. Accordingly, a coordination problem needs to be solved by the controllers, and a unique controller elected for each unmanaged switch.
- C3) **Establishing In-Band Connectivity:** A controller may not only have to manage switches to which it is directly (physically) connected, but also “remote switches”. Accordingly, communication paths need to be established between controllers and switches. As switches are passive, the responsibility to establish communication paths to switches must lie at the controller.
- C4) **Interference:** As a consequence of the need to establish in-band connectivity, a controller must be able to install forwarding rules on “relay switches”. Switches relaying the control traffic to remote switches must be able to classify traffic, *e.g.*, differentiating between regular traffic, control traffic from or to different controllers, etc. More importantly, the rules installed by the controller to establish the in-band management should not interfere with the existing rules used for controlling the regular production traffic in the data plane.
- C5) **No Prior Knowledge:** Controllers and switches can be added and removed at runtime, and switches should not have to be pre-configured with any specific controller addresses. Rather, the control plane should be fully “plug & play”.

IV. SELF-STABILIZING CONTROL PLANE

This section first presents our main ideas to design a distributed control plane which addresses the challenges identified before. In particular, we note that OpenFlow readily provides the features we postulate for our solution.

A. Overview and Building Blocks

Our goal is to place each switch in the network under the control of a single controller and to establish routes that support connectivity to the switches and connectivity between controllers. To do so, each controller runs the same algorithm continuously, reacting to any change in the network (e.g., due to failures or additions of switches, links, or controllers) in a self-stabilizing manner.

First, we introduce a pre-configured *anycast controller address* at the switches: a logical IP address shared by *all controllers*. This solves Challenge C5: in contrast, relying on statically pre-configured *controller IP addresses* would violate the controller plug-and-play requirement.

However, anycast addresses come with a problem in that they may lead to ambiguity and collisions. Accordingly, we require each controller to try to manage a switch exclusively, and configure it to accept connections only from it: A controller takes over control of a switch in an atomic operation (cf. [16]): the consensus challenge (Challenge C2) is solved using an in-band, “first-come-first-served” approach.

To achieve fault tolerance (Challenge C1), when a switch is disconnected from its controller, it needs to detect the controller inactivity and transition to the unmanaged state, allowing other controllers to connect and re-establish the in-band control plane. To detect inactivity, we leverage *rule timeouts* at the switches. In particular, when rules time out, a switch falls back to a set of “safe a priori rules”: these rules can be implemented as low priority rules, which are hidden as long as other, managed rules are installed (and regularly refreshed) by a controller. Timeouts combined with low priority rules are the main mechanism to realize fault tolerance: these low priority rules will eventually appear once all other rules time out, and prevent the switch from being permanently cut off from the network.

Concretely, we can define an unmanaged switch configuration with a set of a priori rules that make sure that neighboring switches learn about the unmanaged switch. While the switch is not able to take actions toward finding and reconnecting to a controller, by informing the neighboring, managed switches, a controller will eventually learn about unmanaged switches.

To manage a remote switch (Challenge C3), the controller installs rules on its neighboring switches, which can then be used to install rules on switches two hops away, etc., iteratively expanding the controller domain. Moreover, the relay mechanism described above and the use of special VLANs for controller traffic solves Challenge C4: our approach ensures that the in-band rules installed to provide controller connectivity never conflict with the rules for regular data plane traffic.

A standard and efficient approach to provide connectivity between nodes in a network are *spanning trees*. And indeed, spanning trees are a useful tool also in our settings, to connect controllers and managed switches. However, in our case, it is not sufficient to use standard Layer-2 STP, since such trees do not give us control to match switches to controllers, and moreover would require more functionality at the switches. In fact, as we will see, for our self-stabilizing algorithm we even use two kinds of trees: the first tree type is used by each controller to communicate with its switches, and the

second type is constructed by each controller to allow all other controllers to reach it.

B. Mechanisms

With these concepts in mind, we can now provide more details on the internals of our control plane. As discussed, we establish connectivity in the control network by creating and maintaining two distinct spanning trees:

- 1) A *per-region spanning tree* (Figure 1a): a bi-directional spanning tree that spans over the *region* owned by the controller. The region owned by a controller is a connected graph containing the controller and switches it controls.
- 2) A *network-wide spanning tree* (Figure 1b): a spanning tree directed and rooted at the controller that spans over the whole network. This second tree supplies each switch and all other controllers with a path to reach the controller. The aim of this second spanning tree is precisely to enable each controller to reach any other controller.

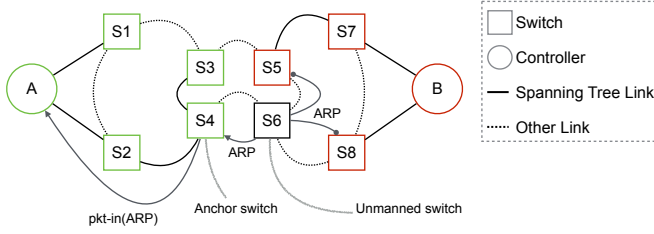
At any moment in time, a switch is either *managed* or *unmanaged*. When a switch is unmanaged, it broadcasts (only) its connection attempts (i.e., ARP requests for the controller anycast address) to all ports. Controller responses are forwarded from all ports to the switch management stack.

When the switch becomes managed, the controller atomically installs a set of match-action rules binding the communication to the controller tree. This configuration has higher priority than the unmanaged configuration but it has timeouts, so in case the switch becomes unmanaged, the managed configuration expires.

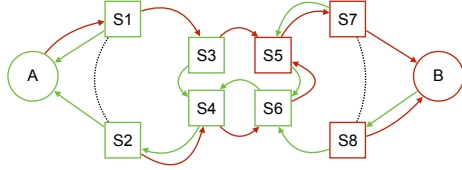
Note that when a switch becomes unmanaged, it no longer forwards any packets of other switches along the region tree, thereby causing them to become unmanaged as well. While this may seem suboptimal, many failures (e.g., link failures) on the path between the controller and a switch will also affect the children and descendants of that switch (with respect to the tree topology). Allowing an unmanaged switch to serve other switches, conflicts with the natural desire to provide contiguous and rooted management trees, and renders it hard to ensure self stabilization.

Algorithm 1 presents a pseudo-code version of the algorithm run by each controller. This algorithm – in combination with the switch behavior defined by the OpenFlow forwarding model – seeks to create two spanning trees to support connectivity for two classes of communications: (1) between switches and controllers, and (2) among controllers.

Example. Figure 1 illustrates at a high level an example. Figure 1 (a) shows the region’s spanning trees of two controllers *A* and *B*. *A*’s region comprises switches *S1 – S4*, and *B*’s region all other switches except *S6*. *S6* has yet to connect to a controller, as denoted by the fact that this switch is broadcasting an ARP packet to its neighbors in order to resolve the virtual controller IP. Figure 1 (b) shows the two fully established network-wide spanning trees as colored arrows that indicate the path towards the two color-coded controllers.



(a) Controllers “conquer” switches adjacent to their regions of control and build a spanning tree for controller-to-switch connectivity.



(b) Per-controller global spanning trees provide controller-to-controller connectivity.

Fig. 1: Controllers iteratively exploring the network, taking control over unmanaged switches (a) and building per-controller spanning trees (b).

name	match	output	prio
flood_mgmt	in_port = MGMT_PORT	ALL	1
all_to_us	mac_dst = MAC _S	MGMT_PORT	1
port_probe	ip_dst = PORT_PROBE	CTRL	1
path_probe	ip_dst = PATH_PROBE	CTRL	1
net_wide_probe	ip_dst = NET_PROBE	CTRL	1

TABLE I: A priori rules. These rules are installed on every switch *MGMT_PORT* is the switch’s management port, out of which the control traffic is sent. *CTRL* is a special OpenFlow port that sends the packet to the controller as packet-in. The other constants designate pre-defined global IP addresses used to tag special packets used in our algorithm.

A Priori Configuration. Each switch is pre-configured with a unique management MAC address, a unique IP address, and a set of OpenFlow rules (called the **a priori rules**). These rules are identical for all switches, except for where they refer to the switch’s unique MAC or IP address. Table I illustrates the a priori rules.

The OpenFlow specifications [17] state that switches are responsible for connecting to their controller, given its IP address; however, it does not specify the mechanism to do so. In practice, switches have a special management port that they use to communicate with the controller. We assume that OpenFlow switches can be configured to ensure that control traffic passes through the flow tables.¹ We refer to the switch port that processes switch-local control traffic according to OpenFlow rules as **mgmt port**.

Building the Region Spanning Tree. Initially, the region of each controller only includes the controller itself. At this point, the switches that are directly connected to the controller can be added to that controller’s region. When this is done, switches that are 1-hop away can be added to the controller’s region, etc.

A switch initiates a connection to a controller by resolving the pre-configured controller IP address via an ARP request

¹A practical way to achieve this despite current vendors’ limitations, is to patch the management port back to a regular switch port.

```

1 process event at ctrl :
  // Region Spanning Tree Events
2 on ARP request from directly-connected switch S:
3   - send ARP reply to S
4   - accept the TCP connection initiated by S
5 on ARP request from switch S (via packet-in):
6   - save S’s anchor switch R and anchor port
7   - configure all packets from I to S to be relayed
  through packet-in / packet-out messages via R
8   - send ARP reply to S
9   - accept the TCP connection initiated by S
10 on OpenFlow session established with switch S:
11   - send a port probe packet to S
12   - install the ctrl → S path by sending downstream
  rules to every switch on the path
13   - periodically send a path probe packet to S – but not
  via packet-out – until a packet-in with the probe is
  received
14 on port probe answer from S:
15   - send the ownership rules to S
16 on path probe answer from S:
17   - configure all packets from I to S to use the normal
  path (don not tunnel using packet-in and packet-out
  messages)
18 on OpenFlow session terminated with switch S:
19   - remove downstream rules from switches on the
  ctrl → S path

  // Network-Wide Spanning Tree Events
20 on C’s network-wide probe as packet-in from S:
21   - install network-wide spanning tree rules for C on S
22 on C’s network-wide probe:
23   - periodically attempt establishing a connection with
  C, until it succeeds or network-wide probes are no
  longer seen for a certain amount of time
24 on timeout of rule net_wide_flood on switch S:
25   - remove rule net_wide_block from S

  // Periodic Events
26 on periodically:
27   - refresh the ownership rules and downstream rules on
  the switches we control
28 on periodically:
29   - send a network-wide probe to all directly-connected
  switches we control

```

Algorithm 1: Pseudo-code executed by controller *ctrl*. *C* and *S* refer to a controller and switch, respectively.

that is broadcast to all ports, except the mgmt port (rule *flood_mgmt*). In the absence of failures, and with common ARP table cache implementations, the switch connects to the first controller whose ARP reply it receives. When the switch connects to a controller, we say that the switch becomes connected. Informally, the switch becomes managed once the controller has taken ownership of it (see details below). A switch is unmanaged if it is neither connected nor managed.

Once the controller establishes the OpenFlow session (a TCP connection) with the switch (see details below), the controller proceeds as a first step to install into the connected

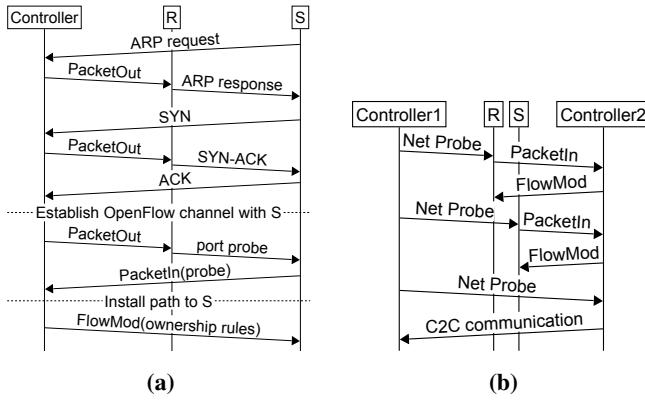


Fig. 2: (a) The management process of newly detected switch S through already managed switch R . (b) The creation of network wide spanning trees and controller to controller channels.

switch the ownership rules described in Table II. These rules override the a priori ones and are associated with an activity timeout; upon timer expiration (e.g., due to failures) the switch returns to the a priori configuration from which the region spanning tree can be reconstructed. We discuss failures below.

Note how this mechanism maps to region growth: initially, controllers can only answer the ARP request of switches directly connected to them. After taking ownership of a switch S , a controller installs rules on the switch. In particular, these rules enable ARP requests of switches connected to S to reach the controller.

To establish its region spanning tree, a controller needs to discover how switches in the tree are connected to one another, down to the port numbers being used. When a controller takes ownership of a switch, the *arp_to_ctrl* rule installed on the switch instructs it to send any ARP request it receives to the controller encapsulated in an OpenFlow packet-in message. The packet-in meta-data contains the input port number on which the ARP packet was received.

When an unmanaged switch attempts to connect to a controller, we call **anchor switch** the first switch of each controller’s region that receives the ARP request of the connecting switch. Hence, there is at most one anchor switch per switch-controller pair. Unless an unmanaged switch is directly connected to the controller (Lines 2-4), it is always the anchor switch that forwards the ARP request to the controller via a packet-in message. When a controller receives an ARP request via an anchor switch, the ARP reply is then delivered to the unmanaged switch via an OpenFlow packet-out message sent to the anchor switch (Lines 5-9).

At this point, the controller still does not have a direct path to send traffic to the unmanaged switch, i.e., it must send all packets intended for the switch through the anchor switch’s OpenFlow session. In the reverse direction (switch to controller), a path exists between the anchor switch and the controller, but the control traffic sent by the switch is still flooded on all its ports.

Note that even though the flooding is still ongoing, all switches other than those in the controller’s region will drop these packets, as they carry the virtual controller IP address and the unique MAC address of the controller, and only

switches in the controller’s region can match on this combination. It is possible, however, that the packet will be relayed by multiple switches in the region, but this will be handled properly by the controller, using the TCP sequence number.

To complete the setup and make that switch part of the controller’s region (hence allowing that switch to serve as anchor switch itself), the controller must perform two steps. First, to stop relaying traffic over the anchor switch, the controller must establish a path from itself to the switch. It does so by installing a rule on every switch on its path to the connecting switch in the region’s spanning tree (Line 12). This rule, called *forward_to_owned_switch*, is shown in Table III.

To confirm that all such rules have been installed properly, the controller sends a *path probe packet* on the path to the switch, and makes the switch to this path once it receives an answer from the switch (Lines 13 and 16-17). The answer from the switch is simply a packet-in encapsulating the original probe packet, resulting from the a priori rule called *path_probe*.

Second, to stop the now-managed switch from flooding all its control packets, the controller must learn the *upstream port* of the switch: the port number through which the connecting switch is connected to its anchor switch. Once this number is known, it can be used to route all packets going from the switch to the controller via the anchor switch (rule *us_to_controller*).

To learn the upstream port, the controller send a *port probe packet* to the switch, via the anchor port of the anchor switch (Line 11). The a priori rule *port_probe* ensures this packet is sent back to the controller via packet-in. The controller can determine the upstream port by the in-port in the packet-in meta-data. Once the upstream port is known, the controller installs the *ownership rules* on the switch (Lines 14-15 of the algorithm). These rules enable the switch to serve as an anchor switch itself (rule *arp_to_virtual_controller*); enable it to forward packets towards the controller (*all_to_controller*), and directs its control traffic towards the controller instead of flooding (*us_to_controller*).

Building the Network-Wide Spanning Tree. With the region’s spanning tree, a controller can communicate with all the switches it controls. In order to be able to communicate with all other controllers, we make use of a per-controller network-wide spanning tree. Each controller constructs its network-wide spanning tree concurrently with the region’s spanning tree.

A controller A constructs its network-wide spanning tree by periodically sending out a special packet: **network-wide probe** (Lines 28-29). When a managed switch receives controller A ’s network-probe for the first time, it will encapsulate it in a packet-in message and send it to its owning controller (rule *net_wide_probe*, say B . Note that A and B can be the same controller. B will then install A ’s network-wide spanning tree rules (Table IV) on the switch.

The network-wide spanning tree rules ensure that further probes received from A on the same port as the first probe will be flooded to all other ports (rule *net_wide_flood*). Furthermore, any probe received from A on any other port is dropped (rule *net_wide_block*) — effectively, this prevents loops in the

name	match	output	prio	overrides
arp_to_ctrl	eth_type = ARP target = CTRL_IP	CTRL	2	/
us_to_ctrl	in_port = MGMT_PORT	<i>S</i> 's upstream port	3	flood_mgmt
all_to_ctrl	in_port = MAC _C	<i>S</i> 's upstream port	3	/

TABLE II: Ownership rules. These rules are installed by a controller C on a switch S after C and S have established an OpenFlow session and C has received the port probe answer to learn S 's upstream port (the port of S through which it is attached to C 's region). $CTRL_IP$ is the global virtual controller IP.

name	match	output	prio
forward_to_owned_switch	ip_dst = IP_{S2}	$S2$'s anchor port on $S1$	2

TABLE III: The downstream rule. This rule is installed by a controller C on each switch $S1$ that lays on the path between C and a switch $S2$ with whom C has established an OpenFlow session.

name	match	output	prio	overrides
net_wide_block	ip_dst = NET_PROBE mac_src = MAC _{C2}	/	2	net_wide_probe (part)
net_wide_flood	in_port = <i>parent port</i> ip_dst = NET_PROBE mac_src = MAC _{C2}	FLOOD	3	net_wide_block (part)
ctrl_to_ctrl	ip_dst = IP_{C2}	<i>parent port</i>	2	/

TABLE IV: Network-wide spanning tree rules. These rules are installed by a controller $C1$ on a switch S it manages whenever the first network-wide probe packet from controller $C2$ is received by S . We call *parent port* the port of S on which the first network-wide probe from $C2$ was received. NET_PROBE is the IP address signaling a network-wide probe packet.

topology from causing broadcast storms.

Once these rules are installed, subsequent probes from A will reach one hop further in the network. Eventually, A 's network-wide spanning tree will span the entire network. Whenever the spanning tree of a controller A extends to another controller, say B , B learns of a path to talk to A . This path results from all the instances of the *ctrl_to_ctrl* rule installed in the network-wide spanning tree. B will then periodically attempt to establish a connection with A (e.g., TCP) (Lines 22-23). These attempts will succeed as soon as the spanning tree of B reaches A , hence supplying A with a return path to communicate with B .

Handling Failures. Informally, our failure model considers all the usual failures that ultimately lead to the termination of an OpenFlow session. Failure of an OpenFlow session is detected by failing to receive an answer to an OpenFlow echo-request message, which are sent periodically. To deal with a failed OpenFlow session, the affected switch reverts the flow tables to the a priori rules (using timeouts and rule priorities) and re-attempts to connect to a controller. By setting the rule timeout to be larger than the timeout of the echo-request message, we can ensure that the rules will never be dropped while the OpenFlow session is live.

Note that the way the switch handles data plane traffic, depends entirely on the controller. We do not impose a pre-defined behavior for this. For instance, a controller application that needs to react to Packet In messages will be unavailable for a brief period of time. But existing data plane rules will still match the ongoing traffic.

To avoid the overhead of refreshing a large number of rules, the non-initial rules could be grouped into a separate secondary flow table.² Then, a single rule in the first flow table could be used to pass all matching packets to this secondary table; and

so, only this catch-all rule needs to be refreshed.

Note that a controller failure is just a special case of an OpenFlow session failure. When that happens, the switch will try to reconnect to a controller; it just will not be able to reconnect to the same as before. The OpenFlow specification does not specify the delay. Assuming we can control it, it should be set higher than the rule refresh timeout; therefore ensuring we only try to reconnect after reverting to the a priori rules. The controller also detects the termination of the session with a switch S and uninstalls the downstream rules for S from switches on the controller- S path (Lines 18-19).

Furthermore, due to failures, the network-wide spanning-tree may need to be altered. Unlike the region's spanning tree, we cannot rely on switches reverting to their a priori rules after a spanning tree experiences a partition. To solve this problem, we set an idle timeout on the *net_wide_flood* rule. As long as there are probe packets, the rule persists. After the probes stop, the rule eventually expires. The controller of the switch will be notified of the rule expiration and will also remove the *net_wide_block* rule from the switch, so that other controllers have a chance to re-include the switch in their network-wide spanning trees.

Finally, we note that since the control traffic is sent in-band, steps should be taken to prevent congestion caused by regular traffic to negatively impact the control traffic. We have not implemented this yet; but we believe a solution based on strict priority queues could be made to work without major difficulties.

V. RELATED WORK

Closest to our work, is the approach of Sharma *et al.* [12] to bootstrap connectivity in an OpenFlow network. However, their work does not consider how to support multiple controllers nor how to establish the control network. Moreover,

²Flow table pipelining is available since OpenFlow 1.1.

their approach relies on switch support for traditional STP and requires modifying DHCP on the switches.

Our approach nicely complements ongoing related research. In particular, our control plane can be used together with and support distributed systems such as ONIX [9], ElastiCon [15], Beehive [18], Kandoo [11], STN [10], to just name a few. Our paper also provides missing links for the interesting work by Akella and Krishnamurthy [8], whose switch-to-controller and controller-to-controller communication mechanisms rely on strong primitives such as consensus protocols, consistent snapshot and reliable flooding, which are not currently available in OpenFlow switches.

We also note that our approach is not limited to a specific technology, but offers flexibilities and can be configured with additional robustness mechanisms, such as warm backups, local fast failover [19], or alternatives spanning trees [20], [21].

Our paper also contributes to the active discussion of which functionality can and should be implemented in OpenFlow. DevoFlow [22] was one of the first works proposing a modification of the OpenFlow model, namely to push responsibility over most flows to switches and adding efficient statistics collection mechanisms. SmartSouth [23] shows that in recent OpenFlow versions, interesting network functions (such as anycast or network traversals) can readily be implemented in-band. More closely related to our paper, [16] shows that it is possible to implement atomic read-modify-write operations on an OpenFlow switch, which can serve as a powerful synchronization and coordination primitive also for distributed control planes; however, such an atomic operation is not required in our system: a controller can claim a switch with a simple write operation. In this paper, we presented a first discussion of how to implement a strong notion of fault-tolerance, namely self-stabilization [13], [14], [24], in SDN.

Bibliographic Note. Our approach was demoed under the name *Medieval* at ACM SOSR 2015 [25].

VI. CONCLUSION

While the benefits of the separation between control and data planes have been studied intensively in the SDN literature, the important question of how to connect these planes has received much less attention.

In this paper, we presented a model and list of challenges for the design of an in-band distributed control plane. Moreover, we presented techniques that allow us to address these challenges. While we yet have to work out the full details and formal proofs, we hope that our work indicates the feasibility of a self-reliant, self-stabilizing distributed control plane.

We also note that our approach comes with interesting security features. First, it can be used in parallel with other security measures (public key infrastructure) to bootstrap authentication and encrypted communication channels. Moreover, it can be used to give control traffic a higher priority than other traffic, thereby preventing some types of DoS attacks. The global spanning trees can allow each switch to sense and communicate with all controllers in its connected component, which can be useful to bootstrap more advanced security protocols, such as SNBI [26], which needs to deliver encryption keys and NTP packets.

Acknowledgments. We are thankful to Nicolas Laurent for his contributions to prior versions of this work.

REFERENCES

- [1] M. Casado, M. J. Freedman, J. Pettit, J. Luo, N. McKeown, and S. Shenker, "Ethane: Taking Control of the Enterprise," in *SIGCOMM*, 2007.
- [2] M. Reitblatt, N. Foster, J. Rexford, C. Schlesinger, and D. Walker, "Abstractions for Network Update," in *SIGCOMM*, 2012.
- [3] N. Vasić, D. Novaković, S. Shekhar, P. Bhurat, M. Canini, and D. Kostić, "Identifying and Using Energy-Critical Paths," in *CoNEXT*, 2011.
- [4] C.-Y. Hong, S. Kandula, R. Mahajan, M. Zhang, V. Gill, M. Nanduri, and R. Wattenhofer, "Achieving High Utilization with Software-Driven WAN," in *SIGCOMM*, 2013.
- [5] S. Jain, A. Kumar, S. Mandal, J. Ong, L. Poutievski, A. Singh, S. Venkata, J. Wanderer, J. Zhou, M. Zhu, J. Zolla, U. Hözlze, S. Stuart, and A. Vahdat, "B4: Experience with a Globally-Deployed Software Defined WAN," in *SIGCOMM*, 2013.
- [6] T. Koponen, K. Amidon, P. Balland, M. Casado, A. Chanda, B. Fulton, I. Ganichev, J. Gross, N. Gude, P. Ingram, E. Jackson, A. Lambeth, R. Lenglet, S.-H. Li, A. Padmanabhan, J. Pettit, B. Pfaff, R. Ramanathan, S. Shenker, A. Shieh, J. Stribling, P. Thakkar, D. Wendlandt, A. Yip, and R. Zhang, "Network Virtualization in Multi-tenant Datacenters," in *NSDI*, 2014.
- [7] D. Levin, M. Canini, S. Schmid, F. Schaffert, and A. Feldmann, "Panopticon: Reaping the Benefits of Incremental SDN Deployment in Enterprise Networks," in *USENIX ATC*, 2014.
- [8] A. Akella and A. Krishnamurthy, "A Highly Available Software Defined Fabric," in *HotNets*, 2014.
- [9] T. Koponen, M. Casado, N. Gude, J. Stribling, L. Poutievski, M. Zhu, R. Ramanathan, Y. Iwata, H. Inoue, T. Hama, and S. Shenker, "Onix: A Distributed Control Platform for Large-scale Production Networks," in *OSDI*, 2010.
- [10] M. Canini, P. Kuznetsov, D. Levin, and S. Schmid, "A Distributed and Robust SDN Control Plane for Transactional Network Updates," in *INFOCOM*, 2015.
- [11] S. Hassas Yeganeh and Y. Ganjali, "Kandoo: A Framework for Efficient and Scalable Offloading of Control Applications," in *HotSDN*, 2012.
- [12] S. Sharma, D. Staessens, D. Colle, M. Pickavet, and P. Demeester, "In-Band Control, Queuing, and Failure Recovery Functionalities for OpenFlow," *IEEE Network*, vol. 30, no. 1, January 2016.
- [13] E. W. Dijkstra, "Self-stabilizing systems in spite of distributed control," *Commun. ACM*, vol. 17, no. 11, pp. 643–644, Nov. 1974.
- [14] M. Schneider, "Self-stabilization," *ACM Comput. Surv.*, vol. 25, no. 1, pp. 45–67, Mar. 1993.
- [15] A. Dixit, F. Hao, S. Mukherjee, T. Lakshman, and R. Kompella, "Towards an Elastic Distributed SDN Controller," in *HotSDN*, 2013.
- [16] L. Schiff, P. Kuznetsov, and S. Schmid, "In-Band Synchronization for Distributed SDN Control Planes," *SIGCOMM Comput. Commun. Rev.*, vol. 46, no. 1, Jan. 2016.
- [17] Open Networking Foundation, "OpenFlow Switch Specification Version 1.3.4." [Online]. Available: <https://www.opennetworking.org/images/stories/downloads/sdn-resources/onf-specifications/openflow/openflow-switch-v1.3.4.pdf>
- [18] S. H. Yeganeh and Y. Ganjali, "Beehive: Simple Distributed Programming in Software-Defined Networks," in *SOSR*, 2016.
- [19] M. Reitblatt, M. Canini, A. Guha, and N. Foster, "FatTire: Declarative Fault Tolerance for Software-defined Networks," in *HotSDN*, 2013.
- [20] M. Borokhovich, L. Schiff, and S. Schmid, "Provable Data Plane Connectivity with Local Fast Failover: Introducing OpenFlow Graph Algorithms," in *HotSDN*, 2014.
- [21] M. Parter, "Dual Failure Resilient BFS Structure," in *PODC*, 2015.
- [22] A. R. Curtis, J. C. Mogul, J. Tourrilhes, P. Yalagandula, P. Sharma, and S. Banerjee, "DevoFlow: Scaling Flow Management for High-performance Networks," in *SIGCOMM*, 2011.
- [23] L. Schiff, M. Borokhovich, and S. Schmid, "Reclaiming the Brain: Useful OpenFlow Functions in the Data Plane," in *HotNets*, 2014.
- [24] B. Awerbuch, B. Patt-Shamir, and G. Varghese, "Self-stabilization by local checking and correction," in *FOCS*, 1991.
- [25] L. Schiff, S. Schmid, and M. Canini, "Medieval: Towards A Self-Stabilizing, Plug & Play, In-Band SDN Control Network," in *SOSR (Demo)*, 2015.
- [26] OpenDaylight, "Secure network bootstrapping infrastructure," 2015. [Online]. Available: https://wiki.opendaylight.org/view/Project_Proposals:Secure_Network_Bootstrapping_Infrastructure