RELIABILITY IN SDN CONTROL PLANE

Sedef Savas

June 30, 2017 Netlab Friday Group Meeting



Networks from topology-zoo.org









Affect of Disaster Radius on # of Disasters





Affect of Disaster Radius on # of Disasters



r = 500, walk = $100 \rightarrow 1080$ disaster circles, max 3 nodes down



Eliminate same-effect disasters and disasters affect < 2 links

46 out of 1080 affects > 1 element



														Observed Values			
											Monthly N	etwork Averages	Target Values	Мау	April	March	
	RTTs don't vary substantially								stant	ally	U.S. Network Averages						
it is don't vary substantiany.								11 y 500	Jocarre	any.	Roundtrip	Latency	< 37 ms	32.6	32.9	33.1	
											Roundtrip	Loss*	< 0.05%	0.00%	0.00%	0.01%	
											Network R	Reliability	> 99.95%	99.9991%	99.9999%	99.9997%	
	_	_									Modem Co	nnect Success Rate	> 95%	99.99%	99.99%	99.99%	
	t o	+									<u>Network J</u>	itter	< 1 ms	0.57	0.58	0.58	
									SF	to NYC is 4500km	$n \rightarrow 4$	14 ms only	r transmissi	on RTT			
U.S. Network Latency International Path Latency																	
rigures are in ms. Thresholds are distance sensitive.																	
						_								Wash	91	Frankfurt	
Cambridge	28 49	Cam				C	urre	ent		W/	e car	n assume overall	RTTe	Tuon			
Chicago	24 27	23 (Chi			C)vera	all		v v	c cai		11115	Trans-F	Pacific Path		
Cleveland	18 35	18	7 Cl	<u> </u>		A	vera	ge:		or	- 1 - 5	v transmission	John	SF	141	Hong Kong	
Dallas	20 6	49	28 2	9 Dal			32 m	າຣ		al	5 I.J		Jelays		Legend Increasin Latency		
Denver	35 24	44	20 2	7 19	Den												
Detroit	22 38	22	7	4 34	26 D)et									15 minutes to en	esh this page every	
Houston	19 7	50	32 3	4 6	25	39 Ho	u								current data	a is displayed.	
Indianapolis			5				Ind										
Kansas City	23 15	32	12 1	8 10	15	19 1	5	Kan									
	51 34	70	16 5	2 22	27	52 3	6					A single contro	ller ca	n not keer	flow cotun	time	
Madicon	51 04	10	5	2 02	21	02 C		10 M	od			A single contro	nor ca		stup		
Nachsilla	0 04		40 4	4 40	20	45 0			au			appristant or w	ithin a	agantahla	limita which	h ia	
Nasriville	0 24		10 1	1 10	32	15 2	.4	17 49	Nas			consistent of w	iuiiii a	cceptable	minus, wind		
New Orleans	12 13		36 3	3 13	31	36	1	22 43	16 NG)		man ant ad to lea		in famme	al mantamatic		
New York	23 45	6	19 1	3 40	37	22 4	1	27 64	25 3	5 NY		reported to be	LUUMS	in for me	sn restoratio)n.	
Orlando	10 28	37	33 2	7 27	45	32 2	22	32 58	1	4 34 Orl		[Poutoba at al "Du	namic C	ontrollar Dra	wicioning in SD	N" 121	
Philadelphia	21 46	8	18 1	1 39		14 4	1	25 63	21 3	2 3 32	Pa	[Doutaba et. al, Dy		Untroller Pro		, ISJ	
Phoenix 4	41 21	65	44 4	7 19	37	53 2	24	29 11	39 3	1 62 45	Phx						
San Antonio	22 3	54	35 3	77	26	40	5	17 30	26 1	2 46 26	47 20 S	A					
San Diego	43 29	73	49 5	6 27	30	55 3	31	37 4	45 3	8 <u>675</u> 3	67 8 2	26 SD					
San Francisco	58 45	70	47 5	1 41	30	54 4	4	44 9	61 5	1 67 (5	61 20 3	39 13 SF					
St. Louis	19 25	27	7 1	4 20	21	14 2	25 13	6 49	11 12 3	2 22 27	20 39 2	26 44 50 StL					
Seattle	71 58	70	47 5	3 52	34	53 5	18	49 27	64 6	4 65 78	63 36 5	i9 29 16 53 Sea					
Washington	19 41	11	22 1	4 36	42	23 4	4	22 67	25 3	0 5 20	4 55 4	1 65 66 17 68 Was htt	p://ipnetv	vork.bgtmo.ip.	att.net/pws/netv	vork delay.html	
a a a a miguri	10 11			- 00	74	20 4		22 01	20 0	0 2 2 3	-1 00 4		,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,				

Flow Setup Latency

[*] Setting up routes in **cellular networks** (when a device becomes active, or during handoff) must complete within 40ms to ensure users can interact with Web services timely.

Flow setup time = speed of control programs, and latency to/from controller + switch modifying forwarding state as dictated

Inbound latency is involved in switch generating events (e.g., when a flow is seen for the first time) can be high (8 ms per rule).

Outbound latency is involved in the switch installing/modifying/deleting forwarding rules provided by control applications, is high as well (3ms - 30ms per rule for insertion and modification).

7 [*] Mazu: Taming Latency in Software Defined Networks, Bell Labs, 2011. UCDAVIS



Create Reachability Circles for Selecting Controllers

Node 8: 1, 6, 7, 8, 9, 10

Node 9: 0, 1, 2, 7, 8, 9, 10

At least 2 closeby controllers will be up for each swich, and at least 1 of them will survive after any possible disaster

Create disaster-joint sets for reachability circles:

Node 8: {1, 6, 7}, {8, 9, 10} **Node 9**: {0, 1}, {2, 7}, {8, 9}, 10

Less than 2 disaster-joint sets or less than 2 controllers makes solution infeasible. So, low S-C latency cause infeasible solutions.



Feasibility Score of DC Nodes

Prioritize more reachable (in terms of latency), less disaster-prone, and more connected DC nodes to place controllers.

Feasibility score =
$$(\alpha * R + \beta * D + \gamma * C)$$

R = Rank of nodes based on reachability (more reachable nodes are favorable.) D = Rank of nodes based on # of disasters that affects them. C = Edge-connectivity of nodes.

Based on this score, initial # and location of controllers are determined.



Switch-Controller Assignment

Least-reachable-selected controllers will be assigned first.

Set a max value for number of switch per controller.



How to Guarantee Enough Capacity After a Disaster?

Load-balancing is important but increase # of controllers.



CONTROL PATH MANAGEMENT FRAMEWORK FOR ENHANCING SDN RELIABILITY

IEEE TRANSACTIONS ON NETWORK AND SERVICE MANAGEMENT, VOL. 14, NO. 2, JUNE 2017

Sejun Song¹, Hyungbae Park¹, Baek-Young Choi¹, Taesang Choi, and Henry Zhu², University of Missouri–Kansas City¹, Cisco²



Introduction

In SDN due to new and multi-lateral network domains, critical challenges to achieve the same reliability services as existing networks.

Control path network lies between a control and a data plane to connect them through in-band SDN or an out-of-band traditional network.

They propose a **control path management framework** to enhance SDN reliability.



New Reliability Challenges

- Traditional networks use distributed reliability protocols (heartbeat mechanisms, no heartbeat = failure).
- SDN creates a new network plane between control/data planes.



Fig. 1. Multi-lateral SDN reliability domains.



Proposed Solution

Control plane reliability is more crucial.

To address reliability challenges, they propose a control path management framework. Strategies:

- 1) ensure a redundant control connection between the data plane and control plane networks
- 2) virtualize the control plane and control path networks to enable a logically centralized cluster (pool) of controllers.
- 3) a fast and accurate failure detection and isolation mechanism in SDN.
- 4) build a control message orchestration mechanism



Traditional Network Reliablity

TRADITIONAL NETWORK RELIABILITY MECHANISMS

Mechanisms	Protocols
Link bundling	Link Aggregation Control Protocol (LACP) [6], EtherChannel [7]
Multipath routing	Equal-Cost Multi-Path routing (ECMP) [8]
System redundancy	Virtual Router Redundancy Protocol (VRRP) [9],
	Host Standby Router Protocol (HSRP) [10], Resilient Packet Ring (RPR) [11]
State synchronization	Non-Stop Routing (NSR) [12], Non-Stop Forwarding (NSF) [13], Stateful Switch-Over (SSO) [14]
Failure detection and handling	Ethernet Automatic Protection Switching (EAPS) [15],
	Ethernet Ring Protection Switching (ERPS) [16], Fast Re-Routing (FRR) [17]

Reliability protocols are embedded in dedicated network devices and treat both data and control failures as interrelated problems according to the physical network topologies.





Fig. 2. Classifications of SDN reliability solutions.

[*] uses switch's link signal to check for fast failure detection, faster than controller that identifies failed link through heartbeat messages and sent out an update. No recovery.

[**] Detection and recovery. Extends OpenFlow protocol to support a monitoring function on switches (similar to fault management of MPLS).

- Recover the data plane network from multiple link failures using back-up routes
- Offload control functionality to a switch and achieved fast recovery and better data plane reliability

[*] M. Desai and T. Nandagopal, "Coping with link failures in centralized control plane architecture," in *Proc. IEEE COMSNET*.
 17 Bengaluru, India, 2010, pp. 1–10.
 [**] J. Kempf *et al.*, "Scalable fault management for OpenFlow," in *Proc. IEEE ICC*, Ottawa, ON, Canada, 2012, pp. 6606–6610.

Control Plane Relialibility

Distributed control plane designs (HyperFlow, ONIX, CPRecovery, B4, and ONOS). Main concerns are scalability and synchronization of network status among multiple physical controllers.

[*] A task manager distributes incoming computations to each controller instance.

These reliability solutions do not consider the correlation between failures of the control plane network and *control path* that are newly introduced in SDN.

Although there are several studies handling the reliability between a controller cluster and a data plane, the path between them has been largely assumed as a **logical connection**. In this work, *control path* is over a 'network' that can be established as in-band or out-of-band of the existing data plane and controller networks.

[*] Z. Cai, A. L. Cox, and T. S. E. Ng, "Maestro: A system for scalable OpenFlow control," Dept. Comput. Sci., Rice Univ., Houston, TX, USA, Tech. Rep. TR10-11, Dec. 2010.



Availability

Availability is formally defined as the fraction of time that a system is operational (Mean Time To Failure (MTTF)),

To improve availability:

- 1) increase the uptime of a system (Mean Time Between Failure (MTBF)) or
- 2) reduce the downtime/outage of a system (Mean Time To Repair (MTTR)).

Little can be done to increase MTTF, since the commodity systems do fail in real operations. Focus reducing MTTR by improving the failure detection and isolation process.

Traditional networks use heartbeat: difficult to identify the exact root cause: absence of heartbeats could have possibly originated from various scenarios of a failure, , thus its recovery mechanism may not be effective.





SENitional Network

Node failandozevallertortsefailuer ho.port 1, also means a failure Bhardolle's seets 20 8 edit fails for failur failur failur

and node 2 and sends an LLDP (Link Layer Discovery node 2 also detects failure on port 2, also means a failure on Protocol) message to regenerate the network topology. node 4's port 1, node 4, or link. Each node blocks the failed ports. According to a new network topology, the SDN controller regesterwork each flag trables be eageneithed. Intervolational to the failed Intervolation of the comparison of the comparison of the failed intervolation of the comparison of the comparison of the failed intervolation of the comparison of t

Illustration of heartbeat-based SDN reliability models: (a) SPoF on control path (b) no SPoF on control path.

UCDAVIS



SDNaditional Network

Node Alfails droode if detestarfaile then port 1, also means a failure Controlle a's east potentialitable failukes on them node 1

and node 2 and sends an LLDP (Link Layer Discovery node 2 also detects failure on port 2, also means a failure on Protocol) message to regenerate the network topology. node 4's port 1, node 4, or link. Each node blocks the failed ports.

AGCO2 ding, tak, and the the terms of terms of the terms of te

Illustration of heartbeat-based SDN reliability models: (a) SPoF on control path (b) no SPoF on control path.



SDN Reliability Challanges: Observations A. SPoF With Multiple Logical Control Path Connections



Recovers from controller failure

Fig. 4. Illustration of unnecessary single points of failure (see the yellow numbers): Multiple logical connections are overlapped such as (1) the legacy switch between the controller and the OpenFlow switch, (2) the link between the legacy switch and the OpenFlow switch, and (3) the interface of the OpenFlow switch.

- network may experience not only a long recovery time but ultimately a service disruption as well.
- effectively disperse logical connections to fully exploit available physical redundancy, so that a failure detection and a switch-over would take place seamlessly without requiring a reconnection process.



B. Configuration of Multiple Controllers

have multiple concurrent logical connections from switches to multiple controllers to minimize switch-over time.

Current OpenFlow switch has to know its controllers. No dynamicity. Manual configuration per switch for the changes in control cluster.



Management Cost = $\sum_{Time} \sum_{c}^{\#of controllers} (1 - P_c)$

where *Pc* denotes the probability controller topology can last without any changes (add/remove controller)

Fig. 6. Relative management costs for given network sizes: the management cost increases as the number of OpenFlow switches and the probability of a cluster configuration change increases.



C. Unrecoverable Control Path Failure Case

Slave only receives port-status messages but not packet-in/flow-removed. Slave is able to detect the network topology/status changes.

if a slave does not receive heartbeat messages consecutively, it initiates a process to become master controller, sends a role change request message to switches.

However, if slave keeps receiving heartbeat messages from master while the *control path* towards the OpenFlow switches is in a failure mode.



C. Unrecoverable Control Path Failure Case (cont.)



Fig. 7. Scenario when an OpenFlow switch loses its master controller: the connection between the slave controller and the OpenFlow switch transfers only port-status messages.

Master and slave controllers are connected through a legacy switch for synchronization

Current specification does not allow an OpenFlow switch to initiate its controller's role change.

This is because current reliability feature does not consider the correlation between failures of the *control plane network* and *inter-connection network*.



D. Control Message Scalability Issues

SDN imposes excessive control traffic overhead, controller platforms allow a variety of heterogeneous application interfaces and protocols to the data plane.

It can cause various scalability and reliability problems: slow message processing, potential message drops, delayed root cause analysis, and late responses against urgent problems.

In traditional routers, internal packet prioritization is used.

Current OpenFlow specification, the SDN controllers drop packets randomly regardless of the importance and urgency of the packets.



PROPOSED SDN RELIABILITY MANAGEMENT FRAMEWORK A. Aligning Logical and Physical Control Path Redundancy

Route logical connections through physical disjoint paths to alleviate/remove the SPoF problem.

Modified the OpenFlow reference implementation by adding an interface selection feature.



B. Controller Cluster Structure Agnostic Virtualization

Each switch does not have to know the distinct IP addresses or port numbers of controllers. To automate adaptation to control plane changes:

- Virtualize physically distributed multiple controllers into one logically centralized controller with 1 virtual IP address.
- Associate Virtual IP with **cluster information broadcaster** (*CIBroadcaster*). *CIBroadcaster* will send the up-to-date cluster information to new switches.
- Other controllers are backup broadcasters and listen to the heartbeat messages from CIB.



B. Controller Cluster Structure Agnostic Virtualization (cont.) How a controller cluster maintains consistency of the cluster information: Hello and update messages.



C. Fast and Accurate Failure Detection and Recovery



Fig. 12. Fast and accurate failure detection and recovery using topology awareness and link signals: (1) the master controller initiates the recovery (Algorithm 3) (2) the OpenFlow switch initiates the recovery (Algorithm 4).



C. Fast and Accurate Failure Detection and Recovery (cont.)



Fig. 13. Failure recovery initiated by the master controller (Algorithm 3).



Fig. 14. Failure recovery initiated by an OpenFlow switch (Algorithm 4).



Fig. 15. Comparison of recovery schemes initiated by an OpenFlow switch and a controller.



D. Control Message Orchestration Module

classification/prioritization system for creating, handling network control messages.



Set 2 bits of the type of service (ToS) field in IPv4 header according to importance of classified control message.

Hence, controllers and switches can differentiate processing sequence and selectively drop received control messages.

Impact of prioritization on CPU utilization.

