What? How? Where?





Laurent Vanbever nsg.ee.ethz.ch

Dagstuhl Seminar Wed Apr 3 2019

#4

What? How? Where?

Main tasks

#1 provisioning policies specifying user intents

- #2 computing forwarding state deriving compliant paths
- #3 maintaining topology information failure detection, topology discovery, etc.

collecting statistics

flow-level, router-level, network-level

#1 provisioning policies specifying user intents

Main tasks

- #2 computing forwarding state deriving compliant paths
- #3 maintaining topology information failure detection, topology discovery, etc.

#4 collecting statistics

flow-level, router-level, network-level





















How can we centrally provision the forwarding state produced by distributed protocols?

How can we centrally provision the forwarding state produced by distributed protocols?

Fibbing [SIGCOMM'15] NetComplete [NSDI'18] a vanbever_fibbing_sigcomm_2015 (1).pdf (page 1 of 14) a vanbever_netcomplete_nsdi_2018 (1).pdf (page 1 of 16) 🗆 • Q 🕀 🚹 📕 🗸 📩 🛞 Q Search 🗾 🖌 📩 🛞 Q Search Central Control Over Distributed Routing http://fibbing.net NetComplete: Practical Network-Wide Configuration Synthesis with Autocompletion Stefano Vissicchio*, Olivier Tilmans*, Laurent Vanbever[†], Jennifer Rexford[‡] * Université catholique de Louvain. † ETH Zurich. ‡ Princeton University Ahmed El-Hassany, Petar Tsankov, Laurent Vanbever, Martin Vechev * name.surname@uclouvain.be, †lvanbever@ethz.ch, †jrex@cs.princeton.edu ETH Zürich netcomplete ethz.ch ABSTRACT part of the network. As a network operator, you suspect a denial-of-service attack (DoS), but cannot know for Centralizing routing decisions offers tremendous flexisure without inspecting the traffic as it could also be a bility, but sacrifices the robustness of distributed protoflash crowd. Your goal is therefore to: (i) isolate the cols. In this paper, we present Fibbing, an architecture flows destined to these IP addresses, (ii) direct them Abstract tributed protocols. A single misconfiguration can bring that achieves both flexibility and robustness through to a scrubber connected between B and C, in order to down the network infrastructure, or worse, a piece of central control over distributed routing. Fibbing intro-Network operators often need to adapt the configuration "clean" them if needed, and (iii) reduce congestion by the Internet in case of BGP-related misconfigurations. duces fake nodes and links into an underlying link-state of a network in order to comply with changing routload-balancing the traffic on unused links, like (B, E). Every few months downtimes involving major players routing protocol, so that routers compute their own foring policies. Evolving existing configurations, however, such as NYSE [1]. Google [2]. Facebook [3], or United warding tables based on the augmented topology. Fib is a complex task as local changes can have unforeseen Airlines [4] make the news. Actually, studies show that bing is expressive, and readily supports flexible load balglobal effects. Not surprisingly, this often leads to mishuman-induced misconfigurations, not physical failures, ancing, traffic engineering, and backup routes. Based takes that result in network downtimes. on high-level forwarding requirements, the Fibbing conexplain the majority of downtimes [5]. We present NetComplete, a system that assists opertroller computes a compact augmented topology and To address these challenges, recently there has been ators in modifying existing network-wide configurations injects the fake components through standard routingan increased interest in configuration verification [6, 7, to comply with new routing policies. NetComplete takes protocol messages. Fibbing works with any unmodified 8, 9, 10, 11, 12, 13 and synthesis [14, 15, 16, 17, 18, as input configurations with "holes" that identify the commercial routers speaking OSPF. Our experiments (a) Initial topology (b) Augmented topology 19, 20]. Configuration synthesis in particular promises also show that it can scale to large networks with many parameters to be completed and "autocompletes" these to alleviate most of the operator's burdens by deriving forwarding requirements, introduces minimal overhead, with concrete values. The use of a partial configuration Figure 1: Fibbing can steer the initial forwardcorrect configurations out of high-level objectives. addresses two important challenges inherent to existing and quickly reacts to network and controller failures ing paths (see (a)) for D_1 through a scrubber by synthesis solutions; (i) it allows the operators to precisely Challenges in network synthesis While promising, netadding fake nodes and links (see (b)). **CCS Concepts** control how configurations should be changed; and (ii) it work operators can still be reluctant to use existing synallows the synthesizer to leverage the existing configurathesis systems for at least three reasons: (i) interpretabil $\bullet \mathbf{Networks} \to \mathbf{Routing \ protocols}; \ \mathit{Network} \ \mathit{architec-}$ Performing this routine task is very difficult in trations to gain performance. To scale, NetComplete relies ity: the synthesizer can produce configurations that differ tures; Programmable networks; Network management; ditional networks. First, since the middlebox and the on powerful techniques such as counter-example guided wildly from manually provided ones, making it hard to destinations are not adjacent to each other, the coninductive synthesis (for link-state protocols) and partial understand what the resulting configuration does. More-Keywords figuration of multiple devices needs to change. Also, evaluation (for path-vector protocols). over, small policy changes can cause the synthesized since intra-domain routing is typically based on short-Fibbing; SDN; link-state routing est path algorithms, modifying the routing configura-We implemented NetComplete and showed that it can configuration (or configuration templates in the case of tion is likely to impact many other flows not involved autocomplete configurations using static routes, OSPF, PropaneAT [16]) to change radically; (ii) protocol cov-1. INTRODUCTION in the attack. In Fig. 1a, any attempt to reroute flows and BGP. Our implementation also scales to realistic neterage: existing systems [15, 16] are restricted to produc-Consider a large IP network with hundreds of devices, to D_1 would also reroute flows to D_2 since they home works and complex routing policies. Among others, it is ing BGP-only configurations, while most networks rely including the components shown in Fig. 1a. A set of to the same router. Advertising D_1 from the middlebox able to synthesize configurations for networks with up to on multiple routing protocols (e.g., to leverage OSPF's IP addresses (D_1) see a sudden surge of traffic, from would attract the right traffic, but would not necessar-200 routers within few minutes. fast-convergence capabilities); and (iii) scalability: remultiple entry points (A, D, and E), that congests a cent synthesizers such as SyNET [20] handle multiple ily alleviate the congestion, because all D_1 traffic would *S. Vissicchio is a postdoctoral researcher of F.R.S.-FNRS. protocols but do not scale to realistic networks 1 Introduction

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies hear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from

SIGCOMM '15, August 17 - 21, 2015, London, United Kingdom © 2015 ACM. ISBN 978-1-4503-3542-3/15/08...\$15.00 http://dr doi org/10 11/E/070E0E6

traverse (and congest) path (A, D, E, B), leaving (A, B) unused. Well-known Traffic-Engineering (TE) protocols (e.g., MPLS RSVP-TE [1]) could help. Unfortunately, since D_1 traffic enters the network from multiple points, many tunnels (three, on A, D, and E, in our tiny example) would need to be configured and signaled. This increases both control-plane and data-plane overhead.

Software Defined Networking (SDN) could easily solve the problem as it enables centralized and direct control of the forwarding behavior. However, moving away from distributed routing protocols comes at a cost. In

In a world where more and more critical services converge on IP, even slight network downtimes can cause large financial or reputational losses. This strategic importance contrasts with the fact that managing a network is surprisingly hard and brittle. Out of high-level requirements, network operators have to come up (often manually) with low-level configurations specifying the drade of day

NetComplete We present a system, NetComplete, which addresses the above challenges with partial synthesis. Rather than synthesizing a new configuration from scratch, NetComplete allows network operators to express their intent by sketching parts of the existing configuration that should remain intact (capturing a high-level insight) and "holes" represented with symbolic values which the synthesizer should instantia (a g OSPE weights BCP import)



What parts of the CP should we offload (if any) and how?

Blink [NSDI'19]

a vanbever blink nsdi 2019.pdf (page 1 of 16)

Blink: Fast Connectivity Recovery Entirely in the Data Plane

Thomas Holterbach, Edgar Costa Molero, Maria Apostolaki Alberto Dainotti, Stefano Vissicchio, Laurent Vanbever

*ETH Zurich, [†]CAIDA / UC San Diego, [‡]University College London

Abstract

□ ~ Q € 1

We present Blink, a data-driven system that leverages TCPinduced signals to detect failures directly in the data plane. The key intuition behind Blink is that a TCP flow exhibits a predictable behavior upon disruption: retransmitting the same packet over and over, at epochs exponentially spaced in time. When compounded over multiple flows, this behavior creates a strong and characteristic failure signal. Blink efficiently analyzes TCP flows to: (*i*) select which ones to track; (*ii*) reliably and quickly detect major traffic disruptions; and (*iii*) recover connectivity—all this, completely in the data plane.

We present an implementation of Blink in P4 together with an extensive evaluation on real and synthetic traffic traces. Our results indicate that Blink: (*i*) achieves sub-second rerouting for large fractions of Internet traffic; and (*ii*) prevents unnecessary traffic shifts even in the presence of noise. We further show the feasibility of Blink by running it on an actual Tofino switch.

1 Introduction

Thanks to widely deployed fast-convergence frameworks such as IPFFR [35], Loop-Free Alternate [7] or MPLS Fast Reroute [29], sub-second and ISP-wide convergence upon link or node failure is now the norm [6, 15]. At a high-level, these fast-convergence frameworks share two common ingredients: (i) fast detection by leveraging hardware-generated signals (e.g., Loss-of-Light or unanswered hardware keepalive [23]); and (ii) quick activation by promptly activating pre-computed backup state upon failure instead of recomputing the paths on-the-fly.

Problem: Convergence upon remote failures is still slow. These frameworks help ISPs to retrieve connectivity upon internal (or peering) failures but are of no use when it comes to restoring connectivity upon remote failures. Unfortunately, remote failures are both frequent and slow to repair, with average convergence times above 30 s [19, 24, 28]. These failures indeed trigger a *control-plane-driven* convergence through the propagation of BGP updates on a per-router and per-prefix



🖌 🖌 🗂 🛞 Q Search

Figure 1: It can take minutes to receive the *first* BGP update following data-plane failures during which traffic is lost.

basis. To reduce convergence time, SWIFT [19] predicts the entire extent of a remote failure from a few received BGP updates, leveraging the fact that such updates are correlated (*e.g.*, they share the same AS-PATH). The fundamental problem with SWIFT though, is that it can take O(minutes) for the *first* BGP update to propagate after the corresponding data-plane failure.

We illustrate this problem through a case study, by measuring the time the first BGP updates took to propagate after the Time Warner Cable (TWC) networks were affected by an outage on August 27 2014 [1]. We consider as outage time t_0 the time at which traffic originated by TWC ASes observed at a large darknet [10] suddenly dropped to zero. We then collect, for each of the routers peering with RouteViews [27] and RIPE RIS [2], the timestamp t_1 of the first BGP withdrawal they received from the same TWC ASes. Figure 1 depicts the CDFs of $(t_1 - t_0)$ over all the BGP peers (100+ routers, in most cases) that received withdrawals for 7 TWC ASes: more than half of the peers took more than a minute to receive the first update (continuous lines). In addition, the CDFs of the time difference between the outage and the last prefix withdrawal for each AS, show that BGP convergence can be as slow as several minutes (dashed lines).

HW-accelerated CPs [HotNets'18]

a vanbever hw accelerated cps hotnets 2018.pdf (page 1 of 7)

🗾 🔹 🔿 🖉 Q Search

Hardware-Accelerated Network Control Planes

Edgar Costa Molero ETH Zürich U

cedgar@ethz.ch

o Stefano Vissicchio University College London s.vissicchio@cs.ucl.ac.uk Laurent Vanbever ETH Zürich lvanbever@ethz.ch

ABSTRACT

One design principle of modern network architecture seems to be set in stone: a software-based control plane drives a hardware- or software-based data plane. We argue that it is time to revisit this principle after the advent of programmable switch ASICs which can run complex logic at line rate.

We explore the possibility and benefits of accelerating the control plane by offloading some of its tasks directly to the network hardware. We show that programmable data planes are indeed powerful enough to run key control plane tasks including: failure detection and notification, connectivity retrieval, and even policy-based routing protocols. We implement in P4 a prototype of such a "hardware-accelerated" control plane, and illustrate its benefits in a case study.

Despite such benefits, we acknowledge that offloading tasks to hardware is not a silver bullet. We discuss its tradeoffs and limitations, and outline future research directions towards hardware-software codesign of network control planes.

1 INTRODUCTION

As the "brain" of the network, the control plane is one of its most important assets. Among other things, the control plane is responsible for *sensing* the status of the network (e.g., which links are up or which links are overloaded), *computing* the best paths along which to guide traffic, and *updating* the underlying data plane accordingly. To do so, the control plane is composed of many dynamic and interacting processes (e.g., routing, management and accounting protocols) whose operation must scale to large networks. In contrast, the data plane is "only" responsible for forwarding traffic according to the control plane decisions, albeit as fast as possible.

These fundamental differences lead to very different design philosophies. Given the relative simplicity of the data plane and the "need for speed", it is typically entirely implemented in hardware. That said, software-based implementations of data planes are also commonly found (e.g., Open-VSwitch [30]) together with hybrid software-hardware ones (e.g., CacheFlow [20]). In short, data plane implementations

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies 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. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

HotNets-XVII, November 15–16, 2018, Redmond, WA, USA © 2018 Association for Computing Machinery. ACM ISBN 978-1-4503-6120-0/18/11...\$15.00 https://doi.org/10.1145/3286082.3286080 cover the entire implementation spectrum, from pure software to pure hardware. In contrast, there is *much* less diversity in control plane implementations. The sheer complexity of the control plane tasks (e.g., performing routing computations) together with the need to update them relatively frequently (e.g., to support new protocols and features) indeed calls for software-based implementations, with only a few key tasks (e.g., detecting physical failures, activating backup forwarding state) being (sometimes) offloaded to hardware [13, 22]. Yet, we argue that a number of recent developments are

creating both the *need* and *opportunity* for rethinking basic design and implementation choices of network control planes.

Need There is a growing need for faster, more scalable, and yet more powerful control planes. Nowadays, even beefedup and highly-optimized software control planes can only process thousands of (BGP) control plane messages per second [23], and can take minutes to converge upon large failures [17, 36]. Parallelizing only marginally helps: for instance, the BGP specification [31] mandates to lock all Adj-RIBs-In before proceeding with the best-path calculation, essentially preventing the parallel execution of best path computations. A concrete risk is that convergence time will keep increasing with the network size and the number of Internet destinations. At the same time, recent research has repeatedly shown the performance benefits of controlling networks with extremely tight control loops, among others to handle congestion (e.g., [7, 21, 29]).

Opportunity Modern reprogrammable switches (e.g., [1]) can perform complex stateful computations on billions of packets per second [19]. Running (pieces of) the control plane at such speeds would lead to almost "instantaneous" convergence, leaving the propagation time of the messages as the primary bottleneck. Besides speed, offloading control plane tasks to hardware would also help by making them traffic-aware. For instance, it enables to update forwarding entries consistently with real-time traffic volumes rather than in a random order.

Research questions Given the opportunity and the need, we argue that it is time to revisit the control plane's design and implementation by considering the problem of offloading parts of it to hardware. This redesign opens the door to multiple research questions including: Which pieces of the control plane should be offloaded? What are the benefits? and How can we overcome the fundamental hardware limitations ? These fundamental limitations come mainly from the very limited instruction set (e.g., no floating point) and the memory available (i.e., around tens of megabytes [19]) of programmable network hardware. We start to answer these questions in this paper and make two contributions.

How can we centrally provision the forwarding state produced by distributed protocols?

Fibbing [SIGCOMM'15] NetComplete [NSDI'18] a vanbever_fibbing_sigcomm_2015 (1).pdf (page 1 of 14) a vanbever_netcomplete_nsdi_2018 (1).pdf (page 1 of 16) 🗆 • Q 🕀 🚹 📕 🗸 📩 🛞 Q Search 🗾 🖌 📩 🛞 Q Search Central Control Over Distributed Routing http://fibbing.net NetComplete: Practical Network-Wide Configuration Synthesis with Autocompletion Stefano Vissicchio*, Olivier Tilmans*, Laurent Vanbever[†], Jennifer Rexford[‡] * Université catholique de Louvain. † ETH Zurich. ‡ Princeton University Ahmed El-Hassany, Petar Tsankov, Laurent Vanbever, Martin Vechev * name.surname@uclouvain.be, †lvanbever@ethz.ch, †jrex@cs.princeton.edu ETH Zürich netcomplete ethz.ch ABSTRACT part of the network. As a network operator, you suspect a denial-of-service attack (DoS), but cannot know for Centralizing routing decisions offers tremendous flexisure without inspecting the traffic as it could also be a bility, but sacrifices the robustness of distributed protoflash crowd. Your goal is therefore to: (i) isolate the cols. In this paper, we present Fibbing, an architecture flows destined to these IP addresses, (ii) direct them Abstract tributed protocols. A single misconfiguration can bring that achieves both flexibility and robustness through to a scrubber connected between B and C, in order to down the network infrastructure, or worse, a piece of central control over distributed routing. Fibbing intro-Network operators often need to adapt the configuration "clean" them if needed, and (iii) reduce congestion by the Internet in case of BGP-related misconfigurations. duces fake nodes and links into an underlying link-state of a network in order to comply with changing routload-balancing the traffic on unused links, like (B, E). Every few months downtimes involving major players routing protocol, so that routers compute their own foring policies. Evolving existing configurations, however, such as NYSE [1], Google [2], Facebook [3], or United warding tables based on the augmented topology. Fib is a complex task as local changes can have unforeseen Airlines [4] make the news. Actually, studies show that bing is expressive, and readily supports flexible load balglobal effects. Not surprisingly, this often leads to mishuman-induced misconfigurations, not physical failures, ancing, traffic engineering, and backup routes. Based takes that result in network downtimes. on high-level forwarding requirements, the Fibbing conexplain the majority of downtimes [5]. We present NetComplete, a system that assists opertroller computes a compact augmented topology and To address these challenges, recently there has been ators in modifying existing network-wide configurations injects the fake components through standard routingan increased interest in configuration verification [6, 7, to comply with new routing policies. NetComplete takes protocol messages. Fibbing works with any unmodified 8, 9, 10, 11, 12, 13 and synthesis [14, 15, 16, 17, 18, as input configurations with "holes" that identify the commercial routers speaking OSPF. Our experiments (a) Initial topology (b) Augmented topology 19, 20]. Configuration synthesis in particular promises also show that it can scale to large networks with many parameters to be completed and "autocompletes" these to alleviate most of the operator's burdens by deriving forwarding requirements, introduces minimal overhead, with concrete values. The use of a partial configuration Figure 1: Fibbing can steer the initial forwardcorrect configurations out of high-level objectives. ing paths (see (a)) for D_1 through a scrubber by addresses two important challenges inherent to existing and quickly reacts to network and controller failures. synthesis solutions; (i) it allows the operators to precisely Challenges in network synthesis While promising, netadding fake nodes and links (see (b)). **CCS Concepts** control how configurations should be changed; and (ii) it work operators can still be reluctant to use existing synallows the synthesizer to leverage the existing configurathesis systems for at least three reasons: (i) interpretabil $\bullet \mathbf{Networks} \to \mathbf{Routing \ protocols}; \ \mathit{Network} \ \mathit{architec-}$ Performing this routine task is very difficult in trations to gain performance. To scale, NetComplete relies ity: the synthesizer can produce configurations that differ ditional networks. First, since the middlebox and the tures; Programmable networks; Network management; on powerful techniques such as counter-example guided wildly from manually provided ones, making it hard to destinations are not adjacent to each other, the coninductive synthesis (for link-state protocols) and partial understand what the resulting configuration does. More-Keywords figuration of multiple devices needs to change. Also, evaluation (for path-vector protocols). over, small policy changes can cause the synthesized since intra-domain routing is typically based on short-Fibbing; SDN; link-state routing est path algorithms, modifying the routing configura-We implemented NetComplete and showed that it can configuration (or configuration templates in the case of tion is likely to impact many other flows not involved autocomplete configurations using static routes, OSPF, PropaneAT [16]) to change radically; (ii) protocol cov-1. INTRODUCTION in the attack. In Fig. 1a, any attempt to reroute flows and BGP. Our implementation also scales to realistic neterage: existing systems [15, 16] are restricted to produc-Consider a large IP network with hundreds of devices, to D_1 would also reroute flows to D_2 since they home works and complex routing policies. Among others, it is ing BGP-only configurations, while most networks rely including the components shown in Fig. 1a. A set of to the same router. Advertising D_1 from the middlebox able to synthesize configurations for networks with up to on multiple routing protocols (e.g., to leverage OSPF's IP addresses (D_1) see a sudden surge of traffic, from would attract the right traffic, but would not necessar 200 routers within few minutes. fast-convergence capabilities); and (iii) scalability: remultiple entry points (A, D, and E), that congests a cent synthesizers such as SyNET [20] handle multiple ily alleviate the congestion, because all D_1 traffic would traverse (and congest) path (A, D, E, B), leaving (A, B) unused. Well-known Traffic-Engineering (TE) protocols *S. Vissicchio is a postdoctoral researcher of F.R.S.-FNRS. protocols but do not scale to realistic networks 1 Introduction

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies 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. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from

SIGCOMM '15, August 17 - 21, 2015, London, United Kingdom © 2015 ACM. ISBN 978-1-4503-3542-3/15/08...\$15.00 N. http://dr doi org/10 114E/070E0E6

(e.g., MPLS RSVP-TE [1]) could help. Unfortunately, since D_1 traffic enters the network from multiple points, many tunnels (three, on A, D, and E, in our tiny example) would need to be configured and signaled. This increases both control-plane and data-plane overhead.

Software Defined Networking (SDN) could easily solve the problem as it enables centralized and direct control of the forwarding behavior. However, moving away from distributed routing protocols comes at a cost. In

In a world where more and more critical services converge on IP, even slight network downtimes can cause large financial or reputational losses. This strategic importance contrasts with the fact that managing a network is surprisingly hard and brittle. Out of high-level requirements, network operators have to come up (often manually) with low-level configurations specifying the drade of davi

NetComplete We present a system, NetComplete, which addresses the above challenges with partial synthesis. Rather than synthesizing a new configuration from scratch, NetComplete allows network operators to express their intent by sketching parts of the existing configuration that should remain intact (capturing a high-level insight) and "holes" represented with symbolic values which the synthesizer should instantia (a g OSPE weights BCP import)

Goal Centrally control distributed routing protocols

where the computation of the forwarding state is distributed

Goal Centrally control distributed routing protocols where the computation of the forwarding state is distributed

Why? Designing central, scalable and robust control is hard must ensure always-on connectivity to the controller
Distributed protocols are still ruling over networks the vast majority of the networks rely on OSPF, BGP, MPLS, ...

How can we control the network-wide forwarding state produced by distributed protocols?

How can we **control** the network-wide forwarding state produced by distributed protocols?

What are our knobs?

The network-wide forwarding state depends on three parameters

Network-wide Forwarding state





Out of these three parameters, two can be controlled



Given a forwarding state we want to program, we therefore have two ways to provision it Given a network-wide forwarding state to provision, one can synthesize

- way 1 the routing messages shown to the routers
- way 2 the configurations run by the routers

Given a network-wide forwarding state

output to provision, one can synthesize

- inputs the routing messages shown to the routers
- **functions** the configurations run by the routers

Controlling distributed computation through synthesis



Controlling distributed computation through synthesis


Consider this network where a source sends traffic to 2 destinations



As congestion appears, the operator wants to shift away one flow from (C,D)



Moving only one flow is impossible though as both destinations are connected to D



impossible to achieve by reweighing the links



Let's lie to the routers



Let's lie to the routers, by injecting fake nodes, links and destinations





Lies are propagated network-wide by the routing protocol



All routers compute their shortest-paths on the augmented topology



C prefers the virtual node (cost 2) to reach the blue destination...



As the virtual node does not really exist, actual traffic is physically sent to A



Synthesizing routing messages is powerful

TheoremFibbing can programany set of non-contradictory paths

Theorem

Fibbing can program

any set of non-contradictory paths



Synthesizing routing messages is fast and works in practice

We developed efficient algorithms

polynomial in the # of requirements

Compute and minimize topologies in ms independently of the size of the network

We tested them against real routers

works on both Cisco and Juniper

Good news

Lots of lies are not required, some of them are redundant









original shortest-path "down and to the right"



desired shortest-path "up and to the right" Our naive algorithm would create 5 lies—one per router



A single lie is sufficient (and necessary)





















Fibbing computes routing messages to inject in ~1ms



Fibbing minimizes the # of routing messages to inject in ~100ms



Fibbing is fully implemented and works with real routers
Existing routers can easily sustain Fibbing-induced load, even with huge topologies

# fake nodes	router memory (MB)	
1000	0.7	
5 000	6.8	
10 000	14.5	
50 000	76.0	
100 000	153	DRAM is cheap

Because it is entirely distributed, programming forwarding entries is fast

# fake nodes	installation time (s)	
1000	0.9	
5 000	4.5	
10 000	8.9	
50 000	44.7	
100 000	89.50	894.50 µs/entry

Fibbing is limited though, among others by the configurations running on the routers

Works with a single protocol family

Dijkstra-based shortest-path routing

Can lead to loads of messages if the configuration is not adapted

Suffers from reliability issues

need to remove the lies upon failures

Controlling distributed computation through synthesis



NetComplete: Practical Network-Wide Configuration Synthesis with Autocompletion



Ahmed El-Hassany



Petar Tsankov







Laurent Vanbever



Martin Vechev

Networked Systems ETH Zürich — seit 2015







Fewer heart attack patients die when top cardiologists are away at conferences, study finds

Heart attack patients are more likely to survive when top cardiologists are not in the hospital, a new study suggests.Researchers at Harvard Medical School...

FLIP.IT

_____ Like









Curious if the Internet is also better during IETF/NANOG/RIPE...



Fewer heart attack patients die when top cardiologists are away at conferences, study finds

Heart attack patients are more likely to survive when top cardiologists are not in the hospital, a new study suggests.Researchers at Harvard Medical School...

FLIP.IT

Like





....



Curious if the Internet is also better during IETF/NANOG/RIPE...



are away at conferences, study finds

Heart attack patients are more likely to survive when top cardiologists are not in the hospital, a new study suggests.Researchers at Harvard Medical School...

Like



Fewer heart attack patients die when top cardiologists





...





The Internet seems to be better off during week-ends...



0 5 10 15 20

% of route leaks

source: Job Snijders (NTT)

This is a far too common story...

News

'Configuration Error' for AWS Outage

By David Ramel 08/12/2015

A "configuration error" caused this week's

Widespread internet outages affected Comcast, Spectrum, Verizon and AT&T customers

BY: CNN POSTED: 1:45 PM, Nov 6, 2017 UPDATED: 7:35 PM, Nov 6, 2017



Amazon's massive AWS outage was cai Over G+

One incorrect command and the whole internet suffers. BY JASON DEL REY | @DELREY | MAR 2, 2017, 2:20PM EST

f SHARE in LINKEDIN



Sean Gallup / Getty

Amazon today blamed human error for the the big AWS outage that took down a b internet sites for several hours on Tuesday afternoon.

The summer of network misconfigurations

by Joanne Godfrey on August 11, 2016 in Application Connectivity Management, Firewall Change Management, Information Security **Management and Vulnerabilities, Security Policy Management**

EMAIL in Share < 32

😏 Tweet

People around the U.S. experienced some internet downtime on Monday. The outage was brief and service has been restored.

CloudFlare apologizes for Telia screwing you

6

 \mathbf{O}

Ð

(in)



Unhappy about massive outage

5 Oct 2016 at 11:33, Shaun Nichols



Why do we have so many misconfigurations?







a desired network behavior

and



Adapt C so that the network follows the new behavior



a desired network behavior



Adapt C so that the network follows the new behavior



a desired network behavior

Nowadays these adaptations are still mostly done manually, which is error-prone and time-consuming

Cisco IOS

```
ip multicast-routing
interface Loopback0
 ip address 120.1.7.7 255.255.255.255
 ip ospf 1 area 0
interface Ethernet0/0
no ip address
interface Ethernet0/0.17
 encapsulation dot1Q 17
 ip address 125.1.17.7 255.255.255.0
 ip pim bsr-border
 ip pim sparse-mode
router ospf 1
 router-id 120.1.7.7
 redistribute bgp 700 subnets
```

```
router bgp 700
neighbor 125.1.17.1 remote-as 100
!
address-family ipv4
redistribute ospf 1 match internal
external 1 external 2
neighbor 125.1.17.1 activate
!
address-family ipv4 multicast
network 125.1.79.0 mask 255.255.255.0
redistribute ospf 1 match internal
external 1 external 2
neighbor 125.1.17.1 activate
!
```

Nowadays these adaptations are still mostly done manually, which is error-prone and time-consuming

Cisco IOS

```
ip multicast-routing
interface Loopback0
 ip address 120.1.7.7 255.255.255.255
 ip ospf 1 area 0
interface Ethernet0/0
no ip address
interface Ethernet0/0.17
 encapsulation dot1Q 17
 ip address 125.1.17.7 255.255.255.0
 ip pim bsr-border
 ip pim sparse-mode
router ospf 1
 router-id 120.1.7.7
```

```
router bgp 700
neighbor 125.1.17.1 remote-as 100
!
address-family ipv4
redistribute ospf 1 match internal
external 1 external 2
neighbor 125.1.17.1 activate
!
address-family ipv4 multicast
network 125.1.79.0 mask 255.255.255.0
redistribute ospf 1 match internal
external 1 external 2
neighbor 125.1.17.1 activate
!
```

redistribute bgp 700 subnets — Anything else than 700 creates blackholes

Configuration synthesis addresses this problem by deriving low-level configurations from high-level requirements

Configuration synthesis addresses this problem by deriving low-level configurations from high-level requirements

Inputs





Configuration synthesis: a booming research area!

Out of high-level requirements, automatically derive...

Genesis [POPL'17] forwarding rules

Propane [SIGCOMM'16] BGP configurations PropaneAT [PLDI'17]

SyNET [CAV'17]OSPF + BGP configurationsZeppelin [SIGMETRICS'18]

Synthesizing configuration is great, but comes with challenges preventing a wide adoption

Problem #1 interpretability can produce configurations that widely differ from humanly-generated ones

interpretability

can produce configurations that widely differ from humanly-generated ones

Problem #2 continuity

can produce widely different configurations given slightly different requirements

interpretability

can produce configurations that widely differ from humanly-generated ones

continuity

can produce widely different configurations given slightly different requirements

Problem #3 deployability cannot flexibly adapt to operational requirements, requiring configuration heterogeneity

A key issue is that synthesizers do not provide operators with a fine-grained control over the synthesized configurations

Introducing...

NetComplete

NetComplete allows network operators to flexibly express their intents through configuration sketches

A configuration with "holes"

interface TenGigabitEthernet1/1/1 ip address ? ? ip ospf cost 10 < ? < 100</pre>

router ospf 100



router bgp 6500

• • • neighbor AS200 import route-map imp-p1 neighbor AS200 export route-map exp-p1 • • • ip community-list C1 permit ? ip community-list C2 permit ?

route-map imp-p1 permit 10 ? route-map exp-p1 ? 10 match community C2 route-map exp-p2 ? 20 match community C1 • • •

interface TenGigabitEthernet1/1/1 ip address ? ? ip ospf cost 10 < ? < 100</pre>

router ospf 100

?

• • •

• • •

router bgp 6500

• • • neighbor AS200 import route-map imp-p1 neighbor AS200 export route-map exp-p1

ip community-list C1 permit ?

ip community-list C2 permit ?

Holes can identify specific attributes such as:

IP addresses

link costs

BGP local preferences

interface TenGigabitEthernet1/1/1 ip address ? ? ip ospf cost 10 < ? < 100

router ospf 100



router bgp 6500

neighbor AS200 import route-map imp-p1 neighbor AS200 export route-map exp-p1

ip community-list C1 permit ?

ip community-list C2 permit ?

route-map imp-p1 permit 10 ? route-map exp-p1 ? 10 match community C2 route-map exp-p2 ? 20 match community C1

Holes can also identify entire pieces of the configuration



NetComplete "autocompletes" the holes such that the output configuration complies with the requirements

interface TenGigabitEthernet1/1/1 ip address ? ? ip ospf cost 10 < ? < 100</pre>

router ospf 100



router bgp 6500

• • • neighbor AS200 import route-map imp-p1 neighbor AS200 export route-map exp-p1 • • • ip community-list C1 permit ? ip community-list C2 permit ?

route-map imp-p1 permit 10 ? route-map exp-p1 ? 10 match community C2 route-map exp-p2 ? 20 match community C1 • • •

interface TenGigabitEthernet1/1/1 ip address 10.0.0.1 255.255.255.254 ip ospf cost 15

router ospf 100 network 10.0.0.1 0.0.0.1 area 0.0.0.0

router bgp 6500

• • • neighbor AS200 import route-map imp-p1 neighbor AS200 export route-map exp-p1 • • •

ip community-list C1 permit 6500:1

ip community-list C2 permit 6500:2

route-map imp-p1 permit 10 set community 6500:1 set local-pref 50 route-map exp-p1 permit 10 match community C2 route-map exp-p2 deny 20 match community C1

• • •

NetComplete reduces the autocompletion problem to a constraint satisfaction problem
First

- protocol semantics
- Encode the
 in high-level requirements as a logical formula (in SMT) partial configurations

Encode the First

Then

- protocol semantics
- high-level requirements as a logical formula (in SMT) partial configurations

Use a solver (Z3) to find an assignment for the undefined configuration variables s.t. the formula evaluates to True



Outputs







NetComplete

Outputs

Main challenge: Scalability

Insight #1

network-specific heuristics

search space navigation

Insight #2

partial evaluation

search space reduction

NetComplete: Practical Network-Wide Configuration Synthesis with Autocompletion



- 1 BGP synthesis optimized encoding
- 2 OSPF synthesis counter-examples-based
- 3 Evaluation flexible, *yet* scalable

But first... "How to configure routing protocols" 101

inter-domain

routing

BGP

intra-domain

routing

OSPF

But first... "How to configure routing protocols" 101

inter-domain

routing

BGP

intra-domain

routing

Internet

Internet

Internet

A network of *networks*

Border Gateway Protocol (BGP)

Internet

The Internet is a network of networks, referred to as Autonomous Systems (AS)





BGP is the routing protocol "glueing" the Internet together



Using BGP, ASes exchange information about the IP prefixes they can reach, directly or indirectly



129.132.0.0/16 ETH/UNIZH Camp Net

BGP routes carry complete path information instead of distance



ETH/UNIZH Camp Net

Each AS appends itself to the path when it propagates announcements





129.132.0.0/16 ETH/UNIZH Camp Net

Selection

out of all paths a router receives:

along which one should it direct traffic?

Selection

out of all paths a router receives:

along which one should it direct traffic?

control where traffic is going

Selection

out of all paths a router receives:

along which one should it direct traffic?

control where traffic is going

Export

for each selected path:

to which neighbors propagate it?

Selection

out of all paths a router receives:

along which one should it direct traffic?

control where traffic is going

Export

for each selected path:

to which neighbors propagate it?

control where traffic is coming from





. . . .

Prefer routes...

- with higher preference
- with shorter path length

- learned externally rather than internally
- whose egress point is the closest
- with smaller egress IP address (tie-break)



Network operators adapt how a router selects and exports BGP advertisements by configuring inbound/outbound filters

commonly known as BGP policies

Network operators adapt how a router selects and exports BGP advertisements by configuring inbound/outbound filters

BGP filter

 $f: Adv \to (Adv \cup \bot)$

Network operators adapt how a router selects and exports BGP advertisements by configuring inbound/outbound filters

BGP filter

predicate

 $f: Adv \to (Adv \cup \bot)$

action

prefix from Google path matches a regular expression label contains X path received from AS X

set preference X attach/strip label Y drop

. . .







Edge **#B** configuration

router bgp 10

• • •

neighbor AS50 in_filter in_dt neighbor AS50 out_filter out_dt • • •

route-map in_dt

set preference 100





Edge #A configuration

router bgp 10

• • •

neighbor AS30 in_filter in_swiss
neighbor AS30 out_filter out_swiss
...

route-map in_swiss

set preference 50




Edge **#B** configuration

```
router bgp 10
  • • •
 neighbor AS50 in_filter in_dt
 neighbor AS50 out_filter out_dt
  • • •
route-map in_dt
 set preference 100
  set label PROVIDER
  • • •
route-map out_dt
 if(label PROVIDER): drop;
 else allow;
```





Edge #A configuration

router bgp 10 • • • neighbor AS30 in_filter in_swiss neighbor AS30 out_filter out_swiss • • • route-map in_swiss set preference 50 set label **PROVIDER** • • • route-map out_swiss if(label PROVIDER): drop; else allow;



But first... "How to configure routing protocols" 101

inter-domain

routing

intra-domain

routing

OSPF

In OSPF, routers build a precise map of the network by flooding its local view to everyone

- Each router keeps track of its incident links and cost as well as whether they are up or down
- Each router broadcasts its own link state to give every router a complete view of the graph
- Routers run Dijkstra on the corresponding graph to compute their shortest-paths and forwarding tables

OSPF configuration mainly consists in figuring out link weights inducing an intended network-wide forwarding state

intended forwarding state



intended forwarding state



OSPF configuration



NetComplete: Practical Network-Wide Configuration Synthesis with Autocompletion



- 1 BGP synthesis optimized encoding
- 2 OSPF synthesis counter-examples-based
- 3 Evaluation flexible, *yet* scalable

NetComplete: Practical Network-Wide Configuration Synthesis with Autocompletion



1 BGP synthesis optimized encoding

> OSPF synthesis counter-examples-based

Evaluation flexible, *yet* scalable

NetComplete autocompletes router-level BGP policies by encoding the desired BGP behavior as a logical formula

 $M \models Reqs \land BGP_{protocol} \land Policies$

$M \models Reqs \land BGP_{protocol} \land Policies$ how should the network forward traffic concrete, part of the input

 $M \models Reqs \land BGP_{protocol} \land Policies$ R1.BGP_{select}(A1,A2) ^ R1.BGP_{select}(A2,A3) \land ...

concrete, protocol semantic

how do BGP routers select routes

$M \models Reqs \land BGP_{protocol} \land Policies$



$M \models Reqs \land BGP_{protocol} \land Policies$

how routes should be modified symbolic, to be found

M ⊨ Reqs ∧ BGPprotocol ∧ Policies

R1.SetLocalPref(A2) = 200

Solving this logical formula consists in assigning each symbolic variable with a concrete value

R1.BGP_{select}(A1,A2) \land R1.BGP_{select}(A2,A3) \land ...

```
BGP_{select}(X,Y) \Leftrightarrow (X.LocalPref > Y.LocalPref) \lor ...
M \models Reqs \land BGP_{protocol} \land Policies
                                         R1.SetLocalPref(A1) = VarX
                                         R1.SetLocalPref(A2) = 200
```

R1.BGP_{select}(A1,A2) \land R1.BGP_{select}(A2,A3) \land ...

 $BGP_{select}(X,Y) \Leftrightarrow (X.LocalPref > Y.LocalPref) \lor ...$ $M \models Reqs \land BGP_{protocol} \land Policies$ R1.SetLocalPref(A1) = VarX R1.SetLocalPref(A2) = 200

 $BGP_{select}(X,Y) \Leftrightarrow (X.LocalPref > Y.LocalPref) \lor ...$ VarX := 250 — $M \models Reqs \land BGP_{protocol} \land Policies$ R1.BGP_{select}(A1,A2) \land R1.SetLocalPref(A1) = VarXR1.SetLocalPref(A2) = 200R1.BGP_{select}(A2,A3) \land ...

challenges





challenges BGP x OSPF

solutions iterative synthesis



NetComplete encodes reduced policies by relying on the requirements and the sketches

NetComplete encodes reduced policies by relying on the requirements and the sketches

Step 1Capture how announcements should propagateusing the requirements

Output

BGP propagation graph

NetComplete encodes reduced policies by relying on the requirements and the sketches

Step 1

Output **BGP** propagation graph

Step 2 via symbolic execution

partially evaluated formulas

Output

- Capture how announcements should propagate using the requirements

- Combine the graph with constraints imposed by sketches

challenges solutions iterative synthesis



$M \models Reqs \land BGP_{protocol} \land Policies$



BGP Decision Process

- 1 Higher local preference
- 2 Shorter AS Path
- 3 Lowest Origin
- 4 Lowest MED
- 5 eBGP over iBGP
- 6 Lower OSPF weight

BGP Decision Process

- 1 Higher local preference
- 2 Shorter AS Path
- 3 Lowest Origin
- 4 Lowest MED
- 5 eBGP over iBGP
- 6 Lower OSPF weight If we hit this step,

If we hit this step, it means that the BGP decision depends on OSPF

NetComplete first tries to find a BGP-only assignment, one in which the BGP behavior does not depend on OSPF

NetComplete first searches for a solution using solely Step 1 to 5

Decision Process

- 1 Higher local preference
- 2 Shorter AS Path
- 3 Lowest Origin
- 4 Lowest MED
- 5 eBGP over iBGP
- 6 Lower OSPF weight

Constraints

PrefNoOSPF(X,Y)

 $PrefOSPF(X,Y) \Leftrightarrow \neg PrefNoOSPF(X,Y)$

NetComplete first searches for a solution using solely Step 1 to 5

BGP_{select}(X,Y)⇔PrefNoOSPF(X,Y)

$M \models Reqs \land BGP_{protocol} \land Policies$

UNSAT! BGP_{select}(X,Y)⇔PrefNoOSPF(X,Y)

$M \models Reqs \land BGP_{protocol} \land Policies$

If NetComplete cannot find an assignment, it then allows the BGP decisions to depend on OSPF


$M \models Reqs \land BGP_{protocol} \land Policies$

BGPselect(X,Y)⇔PrefNoOSPF(X,Y)

 $BGP_{select}(X,Y) \Leftrightarrow PrefOSPF(X,Y) - generate OSPF-based constraints$ BGP_{select}(X,Y)⇔PrefNoOSPF(X,Y)

$M \models Reqs \land BGP_{protocol} \land Policies$

NetComplete: Practical Network-Wide Configuration Synthesis with Autocompletion



BGP synthesis optimized encoding

2 OSPF synthesis counter-examples-based

> Evaluation flexible, *yet* scalable

As for BGP, Netcomplete phrases the problem of finding weights as a constraint satisfaction problem

Consider this initial configuration in which (A,C) traffic is forwarded along the direct link



For performance reasons, the operators want to enable load-balancing



What should be the weights for this to happen?







synthesis procedure



synthesis procedure

$\forall X \in Paths(A,C) \setminus Reqs$

$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$



synthesis procedure



$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$



synthesis procedure



$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$



Synthesized weights

synthesis procedure



$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$

This was easy, but... it does not scale

$\forall X \in Paths(A,C) \setminus Reqs$

$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$

There can be an exponential number of paths between A and C...

 $\forall X \in Paths(A,C) \setminus Reqs$

$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$

To scale, NetComplete leverages **Counter-Example Guided Inductive Synthesis (CEGIS)**

An contemporary approach to synthesis where a solution is iteratively learned from counter-examples

While enumerating all paths is hard, computing shortest paths given weights is easy!

Instead of considering all paths between X and Y

Instead of considering all paths between X and Y

CEGIS Part 1 Consider a random subset *S* of them and synthesize the weights considering *S* only

CEGIS Part 1

intuition

Instead of considering all paths between X and Y

Consider a random subset 5 of them and synthesize the weights considering *S* only

Fast as *S* is small compared to all paths

Consider a random subset 5 of them and CEGIS Part 1 synthesize the weights considering *S* only

intuition

Fast as *S* is small compared to all paths **but** can be wrong

Instead of considering all paths between X and Y

Consider a random subset 5 of them and CEGIS synthesize the weights considering *S* only Part 1

CEGIS Part 2

If so return

Repeat.

Instead of considering all paths between X and Y

Check whether the weights found comply with the requirements over all paths

Else take a counter-example (a path) that violates the Req and add it to S

Consider a random subset 5 of them and CEGIS synthesize the weights considering *S* only Part 1

CEGIS Part 2

Check whether the weights found comply with the requirements **over all paths**

intuition

Fast too simple shortest-path computation

Instead of considering all paths between X and Y





synthesis procedure



synthesis procedure

$\forall X \in SamplePaths(A,C) \setminus Reqs$



synthesis procedure

∀X ∈ SamplePaths(A,C)\Reqs Sample: { [A,B,D,C] }



synthesis procedure

$\forall X \in SamplePaths(A,C) \setminus Reqs$

$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$



synthesis procedure



$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$



synthesis procedure



$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$



Synthesized weights

synthesis procedure



$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$

The synthesized weights are incorrect: $cost(A \rightarrow B \rightarrow C]) = 250 < cost(A \rightarrow C) = 300$



$\forall X \in SamplePaths(A,C) \setminus Reqs$

$Cost(A \rightarrow C) = Cost(A \rightarrow D \rightarrow C) < Cost(X)$

We simply add the counter example to SamplePaths and repeat the procedure



$\forall X \in SamplePaths(A,C) \setminus Reqs$ \downarrow Sample: { [A,B,D,C] } U { [A,B,C] }

The entire procedure usually converges in few iterations making it very fast in practice

NetComplete: Practical Network-Wide Configuration Synthesis with Autocompletion



BGP synthesis optimized encoding

OSPF synthesis counter-examples-based

3 Evaluation flexible, *yet* scalable

Question #1

Can NetComplete synthesize large-scale configurations?

Question #2

How does the concreteness of the sketch influence the running time?
We fully implemented NetComplete and showed its practicality



~10K lines of Python SMT-LIB v2 and Z3

OSPF, BGP, static routes as partial and concrete configs

Cisco-compatible configurations validated with actual Cisco routers

Methodology



Requirement	Sir
-------------	-----

Sketch

- 15 topologies from Topology Zoo small, medium, and large
- mple, Any, ECMP, and ordered (random) using OSPF/BGP
- Built from a fully concrete configuration from which we made a % of the variables symbolic

NetComplete synthesizes configurations for large networks in few minutes

NetComplete synthesizes configurations for large networks in few minutes

Network size

OSPF synthesis Large time (sec) ~150 nodes

settings 16 reqs, 50% symbolic, 5 repet. CEGIS enabled

Reqs.	Synthesis
type	time
Simplo	1 4 -
Simple	145
ECMP	14s 13s

Without CEGIS, OSPF synthesis is >100x slower and often timeouts

NetComplete synthesis time increases as the sketch becomes more symbolic



0



NetComplete synthesis time increases as the sketch becomes more symbolic



NetComplete: Practical Network-Wide Configuration Synthesis with Autocompletion



BGP synthesis optimized encoding

OSPF synthesis counter-examples-based

Evaluation flexible, *yet* scalable

NetComplete: Practical Network-Wide **Configuration Synthesis with Autocompletion**

Scales to realistic network size

- Autocompletes configurations with "holes"
- leaving the concrete parts intact
- Phrases the problem as constraints satisfaction scales using network-specific heuristics & partial evaluation
- synthesizes configurations for large network in minutes

NetComplete: Practical Network-Wide Configuration Synthesis with Autocompletion



Ahmed El-Hassany



Petar Tsankov







Laurent Vanbever



Martin Vechev

Networked Systems ETH Zürich — seit 2015



Controlling distributed computation through synthesis



Network Control Planes

What? How? Where?



What parts of the CP should we offload (if any) and how?

Blink [NSDI'19]

a vanbever blink nsdi 2019.pdf (page 1 of 16)

Blink: Fast Connectivity Recovery Entirely in the Data Plane

Thomas Holterbach, Edgar Costa Molero, Maria Apostolaki Alberto Dainotti, Stefano Vissicchio, Laurent Vanbever

*ETH Zurich, [†]CAIDA / UC San Diego, [‡]University College London

Abstract

□ ~ Q € 1

We present Blink, a data-driven system that leverages TCPinduced signals to detect failures directly in the data plane. The key intuition behind Blink is that a TCP flow exhibits a predictable behavior upon disruption: retransmitting the same packet over and over, at epochs exponentially spaced in time. When compounded over multiple flows, this behavior creates a strong and characteristic failure signal. Blink efficiently analyzes TCP flows to: (*i*) select which ones to track; (*ii*) reliably and quickly detect major traffic disruptions; and (*iii*) recover connectivity—all this, completely in the data plane.

We present an implementation of Blink in P4 together with an extensive evaluation on real and synthetic traffic traces. Our results indicate that Blink: (*i*) achieves sub-second rerouting for large fractions of Internet traffic; and (*ii*) prevents unnecessary traffic shifts even in the presence of noise. We further show the feasibility of Blink by running it on an actual Tofino switch.

1 Introduction

Thanks to widely deployed fast-convergence frameworks such as IPFFR [35], Loop-Free Alternate [7] or MPLS Fast Reroute [29], sub-second and ISP-wide convergence upon link or node failure is now the norm [6, 15]. At a high-level, these fast-convergence frameworks share two common ingredients: (i) fast detection by leveraging hardware-generated signals (e.g., Loss-of-Light or unanswered hardware keepalive [23]); and (ii) quick activation by promptly activating pre-computed backup state upon failure instead of recomputing the paths on-the-fly.

Problem: Convergence upon remote failures is still slow. These frameworks help ISPs to retrieve connectivity upon internal (or peering) failures but are of no use when it comes to restoring connectivity upon remote failures. Unfortunately, remote failures are both frequent and slow to repair, with average convergence times above 30 s [19, 24, 28]. These failures indeed trigger a *control-plane-driven* convergence through the propagation of BGP updates on a per-router and per-prefix



🖌 🖌 🗂 🛞 Q Search

Figure 1: It can take minutes to receive the *first* BGP update following data-plane failures during which traffic is lost.

basis. To reduce convergence time, SWIFT [19] predicts the entire extent of a remote failure from a few received BGP updates, leveraging the fact that such updates are correlated (*e.g.*, they share the same AS-PATH). The fundamental problem with SWIFT though, is that it can take O(minutes) for the *first* BGP update to propagate after the corresponding data-plane failure.

We illustrate this problem through a case study, by measuring the time the first BGP updates took to propagate after the Time Warner Cable (TWC) networks were affected by an outage on August 27 2014 [1]. We consider as outage time t_0 the time at which traffic originated by TWC ASes observed at a large darknet [10] suddenly dropped to zero. We then collect, for each of the routers peering with RouteViews [27] and RIPE RIS [2], the timestamp t_1 of the first BGP withdrawal they received from the same TWC ASes. Figure 1 depicts the CDFs of $(t_1 - t_0)$ over all the BGP peers (100+ routers, in most cases) that received withdrawals for 7 TWC ASes: more than half of the peers took more than a minute to receive the first update (continuous lines). In addition, the CDFs of the time difference between the outage and the last prefix withdrawal for each AS, show that BGP convergence can be as slow as several minutes (dashed lines).

HW-accelerated CPs [HotNets'18]

a vanbever hw accelerated cps hotnets 2018.pdf (page 1 of 7)

🗾 🔹 🔿 🖉 Q Search

Hardware-Accelerated Network Control Planes

Edgar Costa Molero ETH Zürich U

cedgar@ethz.ch

o Stefano Vissicchio University College London s.vissicchio@cs.ucl.ac.uk Laurent Vanbever ETH Zürich lvanbever@ethz.ch

ABSTRACT

One design principle of modern network architecture seems to be set in stone: a software-based control plane drives a hardware- or software-based data plane. We argue that it is time to revisit this principle after the advent of programmable switch ASICs which can run complex logic at line rate.

We explore the possibility and benefits of accelerating the control plane by offloading some of its tasks directly to the network hardware. We show that programmable data planes are indeed powerful enough to run key control plane tasks including: failure detection and notification, connectivity retrieval, and even policy-based routing protocols. We implement in P4 a prototype of such a "hardware-accelerated" control plane, and illustrate its benefits in a case study.

Despite such benefits, we acknowledge that offloading tasks to hardware is not a silver bullet. We discuss its tradeoffs and limitations, and outline future research directions towards hardware-software codesign of network control planes.

1 INTRODUCTION

As the "brain" of the network, the control plane is one of its most important assets. Among other things, the control plane is responsible for *sensing* the status of the network (e.g., which links are up or which links are overloaded), *computing* the best paths along which to guide traffic, and *updating* the underlying data plane accordingly. To do so, the control plane is composed of many dynamic and interacting processes (e.g., routing, management and accounting protocols) whose operation must scale to large networks. In contrast, the data plane is "only" responsible for forwarding traffic according to the control plane decisions, albeit as fast as possible.

These fundamental differences lead to very different design philosophies. Given the relative simplicity of the data plane and the "need for speed", it is typically entirely implemented in hardware. That said, software-based implementations of data planes are also commonly found (e.g., Open-VSwitch [30]) together with hybrid software-hardware ones (e.g., CacheFlow [20]). In short, data plane implementations

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies 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. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions @acm.org.

HotNets-XVII, November 15–16, 2018, Redmond, WA, USA © 2018 Association for Computing Machinery. ACM ISBN 978-1-4503-6120-0/18/11...\$5.00 https://doi.org/10.1145/3286082.3286080 cover the entire implementation spectrum, from pure software to pure hardware. In contrast, there is *much* less diversity in control plane implementations. The sheer complexity of the control plane tasks (e.g., performing routing computations) together with the need to update them relatively frequently (e.g., to support new protocols and features) indeed calls for software-based implementations, with only a few key tasks (e.g., detecting physical failures, activating backup forwarding state) being (sometimes) offloaded to hardware [13, 22]. Yet, we argue that a number of recent developments are

creating both the *need* and *opportunity* for rethinking basic design and implementation choices of network control planes.

Need There is a growing need for faster, more scalable, and yet more powerful control planes. Nowadays, even beefedup and highly-optimized software control planes can only process thousands of (BGP) control plane messages per second [23], and can take minutes to converge upon large failures [17, 36]. Parallelizing only marginally helps: for instance, the BGP specification [31] mandates to lock all Adj-RIBs-In before proceeding with the best-path calculation, essentially preventing the parallel execution of best path computations. A concrete risk is that convergence time will keep increasing with the network size and the number of Internet destinations. At the same time, recent research has repeatedly shown the performance benefits of controlling networks with extremely tight control loops, among others to handle congestion (e.g., [7, 21, 29]).

Opportunity Modern reprogrammable switches (e.g., [1]) can perform complex stateful computations on billions of packets per second [19]. Running (pieces of) the control plane at such speeds would lead to almost "instantaneous" convergence, leaving the propagation time of the messages as the primary bottleneck. Besides speed, offloading control plane tasks to hardware would also help by making them traffic-aware. For instance, it enables to update forwarding entries consistently with real-time traffic volumes rather than in a random order.

Research questions Given the opportunity and the need, we argue that it is time to revisit the control plane's design and implementation by considering the problem of offloading parts of it to hardware. This redesign opens the door to multiple research questions including: Which pieces of the control plane should be offloaded? What are the benefits? and How can we overcome the fundamental hardware limitations ? These fundamental limitations come mainly from the very limited instruction set (e.g., no floating point) and the memory available (i.e., around tens of megabytes [19]) of programmable network hardware. We start to answer these questions in this paper and make two contributions.

What parts of the CP should we offload (if any) and how?

Blink [NSDI'19]

a vanbever blink nsdi 2019.pdf (page 1 of 16)

Blink: Fast Connectivity Recovery Entirely in the Data Plane

Thomas Holterbach, Edgar Costa Molero, Maria Apostolaki Alberto Dainotti, Stefano Vissicchio, Laurent Vanbever

*ETH Zurich, [†]CAIDA / UC San Diego, [‡]University College London

Abstract

We present Blink, a data-driven system that leverages TCPinduced signals to detect failures directly in the data plane. The key intuition behind Blink is that a TCP flow exhibits a predictable behavior upon disruption: retransmitting the same packet over and over, at epochs exponentially spaced in time. When compounded over multiple flows, this behavior creates a strong and characteristic failure signal. Blink efficiently analyzes TCP flows to: (*i*) select which ones to track; (*ii*) reliably and quickly detect major traffic disruptions; and (*iii*) recover connectivity—all this, completely in the data plane.

We present an implementation of Blink in P4 together with an extensive evaluation on real and synthetic traffic traces. Our results indicate that Blink: (*i*) achieves sub-second rerouting for large fractions of Internet traffic; and (*ii*) prevents unnecessary traffic shifts even in the presence of noise. We further show the feasibility of Blink by running it on an actual Tofino switch.

1 Introduction

Thanks to widely deployed fast-convergence frameworks such as IPFFR [35], Loop-Free Alternate [7] or MPLS Fast Reroute [29], sub-second and ISP-wide convergence upon link or node failure is now the norm [6, 15]. At a high-level, these fast-convergence frameworks share two common ingredients: (i) fast detection by leveraging hardware-generated signals (e.g., Loss-of-Light or unanswered hardware keepalive [23]); and (ii) quick activation by promptly activating pre-computed backup state upon failure instead of recomputing the paths on-the-fly.

Problem: Convergence upon remote failures is still slow. These frameworks help ISPs to retrieve connectivity upon internal (or peering) failures but are of no use when it comes to restoring connectivity upon remote failures. Unfortunately, remote failures are both frequent and slow to repair, with average convergence times above 30 s [19,24,28]. These failures indeed trigger a *control-plane-driven* convergence through the propagation of BGP updates on a per-router and per-prefix



🖊 🖌 📩 🛞 Q Search

Figure 1: It can take minutes to receive the *first* BGP update following data-plane failures during which traffic is lost.

basis. To reduce convergence time, SWIFT [19] predicts the entire extent of a remote failure from a few received BGP updates, leveraging the fact that such updates are correlated (*e.g.*, they share the same AS-PATH). The fundamental problem with SWIFT though, is that it can take O(minutes) for the *first* BGP update to propagate after the corresponding data-plane failure.

We illustrate this problem through a case study, by measuring the time the first BGP updates took to propagate after the Time Warner Cable (TWC) networks were affected by an outage on August 27 2014 [1]. We consider as outage time to, the time at which traffic originated by TWC ASes observed at a large darknet [10] suddenly dropped to zero. We then collect, for each of the routers peering with RouteViews [27] and RIPE RIS [2], the timestamp t_1 of the first BGP withdrawal they received from the same TWC ASes. Figure 1 depicts the CDFs of $(t_1 - t_0)$ over all the BGP peers (100+ routers, in most cases) that received withdrawals for 7 TWC ASes: more than half of the peers took more than a minute to receive the first update (continuous lines). In addition, the CDFs of the time difference between the outage and the last prefix withdrawal for each AS, show that BGP convergence can be as slow as several minutes (dashed lines).

HW-accelerated CPs [HotNets'18]

🗾 🖌 🛅 🕢 Q Sea

Hardware-Accelerated Network Control Planes

dgar Costa Molero ETH Zürich

Stefano Vissicchio University College Londor s.vissicchio@cs.ucl.ac.uk Laurent Vanbever ETH Zürich lvanbever@ethz.ch

STRACT

One design principle of modern network architecture seems to be set in stone: a software-based control plane drives a hardware- or software-based data plane. We argue that it is time to revisit this principle after the advent of programmable switch ASICs which can run complex logic at line rate.

We explore the possibility and benefits of accelerating the control plane by offloading some of its tasks directly to the network hardware. We show that programmable data planes are indeed powerful enough to run key control plane tasks including: failure detection and notification, connectivity retrieval, and even policy-based routing protocols. We implement in P4 a prototype of such a "hardware-accelerated" control plane, and illustrate its benefits in a case study.

Despite such benefits, we acknowledge that offloadin tasks to hardware is not a silver bullet. We discuss its tradeof and limitations, and outline future research directions towarc hardware-software codesign of network control planes.

1 INTRODUCTION

As the "brain" of the network, the control plane is one of its most important assets. Among other things, the control plane is responsible for *sensing* the status of the network (e.g., which links are up or which links are overloaded), *computing* the best paths along which to guide traffic, and *updating* the underlying data plane accordingly. To do so, the control plane is composed of many dynamic and interacting processes (e.g., routing, management and accounting protocols) whose operation must scale to large networks. In contrast, the data plane is "only" responsible for forwarding traffic according to the control plane decisions, albeit as fast as possible.

These fundamental differences lead to very different de sign philosophies. Given the relative simplicity of the data plane and the "need for speed", it is typically entirely imple mented in hardware. That said, software-based implementa tions of data planes are also commonly found (e.g., Open VSwitch [30]) together with hybrid software-hardware ones (e.g., CacheFlow [20]). In short, data plane implementations

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies 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. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

INVESSAVIT, NOVEMBER 15-10, 2018, Reamona, WA, OSA
2018 Association for Computing Machinery.
'M ISBN 978-1-4503-6120-0/18/11...\$15.00
oc://doi.org/10.1145/2286062.3286080

cover the entire implementation spectrum, from pure software to pure hardware. In contrast, there is *much* less diversity in control plane implementations. The sheer complexity of the control plane tasks (e.g., performing routing computations) together with the need to update them relatively frequently (e.g., to support new protocols and features) indeed calls for software-based implementations, with only a few key tasks (e.g., detecting physical failures, activating backup forwarding state) being (sometimes) offloaded to hardware [13, 22].

Yet, we argue that a number of recent developments ar creating both the *need* and *opportunity* for rethinking basi design and implementation choices of network control plane

Need There is a growing need for faster, more scalable, and yet more powerful control planes. Nowadays, even beefedup and highly-optimized software control planes can only process thousands of (BGP) control plane messages per second [23], and can take minutes to converge upon large failures [17, 36]. Parallelizing only marginally helps: for instance, the BGP specification [31] mandates to lock all Adj-RIBs-In before proceeding with the best-path calculation, essentially preventing the parallel execution of best path computations. A concrete risk is that convergence time will keep increasing with the network size and the number of Internet destinations. At the same time, recent research has repeatedly shown the performance benefits of controlling networks with extremely tight control loops, among others to handle congestion (e.g., [7, 21, 29]).

Opportunity Modern reprogrammable switches (e.g., [1]) can perform complex stateful computations on billions of packets per second [19]. Running (pieces of) the control plane at such speeds would lead to almost "instantaneous" convergence. leaving the propagation time of the messages as the primary bottleneck. Besides speed, offloading control plane tasks to hardware would also help by making them traffic-aware. For instance, it enables to update forwarding entries consistently with real-time traffic volumes rather than in a random order.

Research questions Given the opportunity and the need, we argue that it is time to revisit the control plane's design and implementation by considering the problem of offloading parts of it to hardware. This redesign opens the door to multiple research questions including: *Which pieces of the control plane should be offloaded? What are the benefits?* and *How can we overcome the fundamental hardware limitations*? These fundamental limitations come mainly from the very limited instruction set (e.g., no floating point) and the memory available (i.e., around tens of megabytes [19]) of programmable network hardware. We start to answer these questions in this paper and make two contributions.

Blink: Fast Connectivity Recovery Entirely in the Data Plane



Joint work with

Edgar Costa Molero Maria Apostolaki Stefano Vissicchio Alberto Dainotti Laurent Vanbever

ETH Zürich ETH Zürich University College London CAIDA, UC San Diego ETH Zürich

Thomas Holterbach ETH Zürich

NSDI 26th February 2019

https://blink.ethz.ch

Fire at AT&T facility causes widespread outage in North Texas

Time Warner Cable comes back from nationwide Internet outage by Brian Stelter @brianstelter



November 6, 2017 | Emerging Threats





Your network



Your network





Remote

Remote

Remote

Your network



Remote

Remote

Local



Upon local failures, connectivity can be quickly restored

Upon local failures, connectivity can be quickly restored

Fast failure detection using *e.g.*, hardware-generated signals

Fast traffic rerouting using *e.g.*, Prefix Independent Convergence or MPLS Fast Reroute

Upon remote failures, the only way to restore connectivity is to wait for the Internet to converge

Upon remote failures, the only way to restore connectivity is to wait for the Internet to converge

... and the Internet converges very slowly*



*Holterbach et al. SWIFT: Predictive Fast Reroute ACM SIGCOMM, 2017

Fire at AT&T facility causes widespread outage in North Texas

Time Warner Cable comes back from nationwide Internet outage by Brian Stelter @brianstelter

August 27, 2014: 11:07 PM ET

November 6, 2017 | Emerging Threats





BGP took minutes to converge upon the Time Warner Cable outage in 2014 1.0 8.0 0.6 CDF over the BGP peers 0.4 0.2 0.0 100 200 300 400 500 600 0 Time difference between the outage and the BGP withdrawals (s)



Control-plane (*e.g.*, BGP) based techniques typically converge slowly upon remote outages

Control-plane (e.g., BGP) based techniques typically converge slowly upon remote outages

What about using data-plane signals for fast rerouting?

Blink: Fast Connectivity Recovery Entirely in the Data Plane



Thomas Holterbach ETH Zürich

NSDI 26th February 2019

https://blink.ethz.ch

Outline

- 1. Why and how to use data-plane signals for fast rerouting
- 2. *Blink* infers more than 80% of the failures, often within 1s
- 3. Blink quickly reroutes traffic to working backup paths
- 4. *Blink* works in practice, on existing devices

Outline

1. Why and how to use data-plane signals for fast rerouting

2. *Blink* infers more than 80% of the failures, often within 1s

3. *Blink* quickly reroutes traffic to working backup paths

4. *Blink* works in practice, on existing devices

TCP flows exhibit the same behavior upon failures

TCP flows exhibit the same behavior upon failures





destination

TCP flows exhibit the same behavior upon failures destination source S:500 A:1000





failure

TCP flows exhibit the same behavior upon failures










We simulated a failure affecting 100k flows with NS3

Same RTT distribution than in a real trace*

We simulated a failure affecting 100k flows with NS3

Same RTT distribution than in a real trace*

Number of retransmissions



We simulated a failure affecting 100k flows with NS3

Same RTT distribution than in a real trace*

Number of retransmissions



We simulated a failure affecting 100k flows with NS3

Same RTT distribution than in a real trace*

70K 60K Failure 50K 40K Number of retransmissions 30K 20K 10K 0 2 3 5 7 0 4 6 Time (s)

We simulated a failure affecting 100k flows with NS3

Same RTT distribution than in a real trace*

70K 60K Failure 50K 40K Number of retransmissions 30K 20K 10K 0 3 5 2 7 0 4 6 Time (s)

We simulated a failure affecting 100k flows with NS3

Same RTT distribution than in a real trace*

70K 60K Failure 50K 40K Number of retransmissions 30K 20K 10K 0 5 2 3 7 0 4 6 Time (s)

We simulated a failure affecting 100k flows with NS3

Same RTT distribution than in a real trace*

70K 60K Failure 50K 40K Number of retransmissions 30K 20K 10K 0 3 5 7 2 0 4 6 Time (s)

Outline

1. Why and how to use data-plane signals for fast rerouting

2. *Blink* infers more than 80% of the failures, often within 1s

3. *Blink* quickly reroutes traffic to working backup paths

4. *Blink* works in practice, on existing devices

To detect failures, *Blink* looks at TCP retransmissions

number of retransmissions



number of retransmissions



number of retransmissions



number of retransmissions



Solution #1: *Blink* looks at consecutive packets with the same sequence number

Solution #1: Blink looks at consecutive packets with the same sequence number destination source S:500 A:1000 $= SRTT + 4 * RTT_VAR$ failure RTO: 200ms cwnd:4 pkts(=congestion window) S:1000 S:2100-5:3100 5:4100 *t* + 200ms cwnd:1 S:1000 exponential backoff cwnd:1 *t* + 600ms S:1000 cwnd:1 *t* + 1400ms S:1000



Solution #2: *Blink* monitors the number of flows experiencing retransmissions over time using a sliding window

number of retransmissions



number of flows experiencing retransmissions

number of retransmissions





number of flows experiencing retransmissions

number of retransmissions





number of flows experiencing retransmissions

number of retransmissions



number of flows experiencing retransmissions



number of retransmissions



number of flows experiencing retransmissions



number of retransmissions



number of flows experiencing retransmissions



number of retransmissions



number of flows experiencing retransmissions



number of retransmissions



number of flows experiencing retransmissions



Blink is intended to run in programmable switches

Blink is intended to run in programmable switches Problem: those switches have very limited resources

Solution #1: Blink focuses on the popular prefixes, *i.e.*, the ones that attract data traffic

Solution #1: Blink focuses on the popular prefixes, *i.e.*, the ones that attract data traffic

As Internet traffic follows a Zipf-like distribution* (1k pref. account for >50%), **Blink** covers the vast majority of the Internet traffic

*Sarra et al. Leveraging Zipf's Law for Traffic offloading ACM CCR, 2012



Solution #2: Blink monitors a sample of the flows for each monitored prefix

TCP flows



Traffic to a destination prefix

Solution #2: Blink monitors a sample of the flows for each monitored prefix



TCP flows

default 64 flows monitored

Traffic to a destination prefix
To monitor active flows, *Blink* evicts a flow from the sample if it does not send a packet for a given time (default 2s)

To monitor active flows, **Blink** evicts a flow from the sample if it does not send a packet for a given time (default 2s)

and selects a new one in a *first-seen, first-selected* manner

Blink infers a failure for a prefix when the majority of the monitored flows experience retransmissions

Blink infers a failure for a prefix when the majority of the monitored flows experience retransmissions

number of flows experiencing retransmissions



Time

Blink infers a failure for a prefix when the majority of the monitored flows experience retransmissions

number of flows experiencing retransmissions



Time

We evaluated *Blink* failure inference using 15 real traces, 13 from CAIDA, 2 from MAWI, covering a total of 15.8 hours

We evaluated *Blink* failure inference using 15 real traces, 13 from CAIDA, 2 from MAWI, covering a total of 15.8 hours

We are interested in:





Speed: How long does Blink take to infer failures

Accuracy: True Positive Rate vs False Positive Rate

As we do not have ground truth, we generated synthetic traces following the traffic characteristics extracted from the real traces

As we do not have ground truth, we generated synthetic traces following the traffic characteristics extracted from the real traces

Step #1 - We extracted the RTT, Packet rate, Flow duration from the real traces

Step #2 - We used NS3 to replay these flows and simulate a failure

Step #3 - We ran a Python-based version of **Blink** on the resulting traces

Blink failure inference accuracy is above 80% for 13 real traces out of 15



Blink failure inference accuracy is above 80% for 13 real traces out of 15

8.0 0.6 **True Positive Rate** 0.4 0.2 0.0 6 4 2 3 5



Real traces ID

Blink avoids incorrectly inferring failures when packet loss is below 4%





Blink avoids incorrectly inferring failures when packet loss is below 4%



2	3	4	5	 8	9	
	0	0.67	0.67	 1.3	2.7	

Blink infers a failure within 1s for the majority of the cases



10 12 14 9 11 13 15 15 Real traces ID

Blink infers a failure within 1s for the majority of the cases



Outline

- 1. Why and how to use data-plane signals for fast rerouting
- 2. *Blink* infers more than 80% of the failures, often within 1s
- 3. *Blink* quickly reroutes traffic to working backup paths
- 4. *Blink* works in practice, on existing devices

Upon detection of a failure, *Blink* immediately activates backup paths pre-populated by the control-plane









Solution: As for failures, *Blink* uses data-plane signals to pick a working backup path

Solution: As for failures, *Blink* uses data-plane signals to pick a working backup path







Solution: As for failures, *Blink* uses data-plane signals to pick a working backup path





As for failures, **Blink** compares the sequence number of consecutive packets to detect blackholes or loops*



*See the paper for an evaluation of the rerouting

Outline

- 1. Why and how to use data-plane signals for fast rerouting
- 2. *Blink* infers more than 80% of the failures, often within 1s
- 3. *Blink* quickly reroutes traffic to working backup paths
- 4. *Blink* works in practice, on existing devices

We ran *Blink* on the 15 real traces (15.8 hours)

We ran *Blink* on the 15 real traces (15.8 hours) and it detected 6 outages, each affecting at least 42% of all the flows









Blink works on a real Barefoot Tofino switch



Blink works on a real Barefoot Tofino switch



Blink works on a real Barefoot Tofino switch


Blink works on a real Barefoot Tofino switch



RTTs in [10ms; 300ms]

Blink works on a real Barefoot Tofino switch



RTTs in [10ms; 300ms]

Blink works on a real Barefoot Tofino switch



RTTs in [10ms; 300ms]

Blink: Fast Connectivity Recovery Entirely in the Data Plane

Infers failures from data-plane signals with more than 80% accuracy, and often within 1s

Fast reroutes traffic at line rate to working backup paths

Works on real traffic traces and on existing devices









https://blink.ethz.ch

Blink: Fast Connectivity Recovery Entirely in the Data Plane



Joint work with

Edgar Costa Molero Maria Apostolaki Stefano Vissicchio Alberto Dainotti Laurent Vanbever

ETH Zürich ETH Zürich University College London CAIDA, UC San Diego ETH Zürich

Thomas Holterbach ETH Zürich

NSDI 26th February 2019

What parts of the CP should we offload (if any) and how?

Blink [NSDI'19] a vanbever hw accelerated cps hotnets 2018.pdf (page 1 of 7) 🗾 🖌 🚯 Q Search Hardware-Accelerated Network Control Planes Edgar Costa Molero Stefano Vissicchio Laurent Vanbever ETH Zürich University College London ETH Zürich cedgar@ethz.ch s.vissicchio@cs.ucl.ac.uk lvanbever@ethz.ch Blink: Fast Connectivity Recovery Entirely in the Data Plane ABSTRACT cover the entire implementation spectrum, from pure software to pure hardware. In contrast, there is much less diversity in One design principle of modern network architecture seems control plane implementations. The sheer complexity of the to be set in stone: a software-based control plane drives a control plane tasks (e.g., performing routing computations) hardware- or software-based data plane. We argue that it is together with the need to update them relatively frequently time to revisit this principle after the advent of programmable (e.g., to support new protocols and features) indeed calls for switch ASICs which can run complex logic at line rate. software-based implementations, with only a few key tasks We explore the possibility and benefits of accelerating the (e.g., detecting physical failures, activating backup forwardcontrol plane by offloading some of its tasks directly to the neting state) being (sometimes) offloaded to hardware [13, 22]. work hardware. We show that programmable data planes are Yet, we argue that a number of recent developments are indeed powerful enough to run key control plane tasks includcreating both the need and opportunity for rethinking basic ing: failure detection and notification, connectivity retrieval, design and implementation choices of network control planes. and even policy-based routing protocols. We implement in P4 a prototype of such a "hardware-accelerated" control plane, Need There is a growing need for faster, more scalable, and and illustrate its benefits in a case study. yet more powerful control planes. Nowadays, even beefed-Despite such benefits, we acknowledge that offloading up and highly-optimized software control planes can only tasks to hardware is not a silver bullet. We discuss its tradeoffs process thousands of (BGP) control plane messages per secand limitations, and outline future research directions towards ond [23], and can take minutes to converge upon large failhardware-software codesign of network control planes. ures [17, 36]. Parallelizing only marginally helps: for instance, the BGP specification [31] mandates to lock all Adj-RIBs-In **1 INTRODUCTION** before proceeding with the best-path calculation, essentially preventing the parallel execution of best path computations. As the "brain" of the network, the control plane is one of A concrete risk is that convergence time will keep increasing its most important assets. Among other things, the control with the network size and the number of Internet destinations plane is responsible for sensing the status of the network (e.g.,

1 Introduction

Problem: Convergence upon *remote* failures is still slow.

HW-accelerated CPs [HotNets'18]

which links are up or which links are overloaded), computing the best paths along which to guide traffic, and updating the underlying data plane accordingly. To do so, the control plane is composed of many dynamic and interacting processes (e.g., routing, management and accounting protocols) whose operation must scale to large networks. In contrast, the data plane is "only" responsible for forwarding traffic according to the control plane decisions, albeit as fast as possible.

These fundamental differences lead to very different design philosophies. Given the relative simplicity of the data plane and the "need for speed", it is typically entirely implemented in hardware. That said, software-based implementations of data planes are also commonly found (e.g., Open-VSwitch [30]) together with hybrid software-hardware ones (e.g., CacheFlow [20]). In short, data plane implementations

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies 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. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org. HotNets-XVII, November 15-16, 2018, Redmond, WA, USA

© 2018 Association for Computing Machinery ACM ISBN 978-1-4503-6120-0/18/11...\$15.00 https://doi.org/10.1145/3286062.328608

At the same time, recent research has repeatedly shown the performance benefits of controlling networks with extremely tight control loops, among others to handle congestion (e.g., [7, 21, 29]).

Opportunity Modern reprogrammable switches (e.g., [1]) can perform complex stateful computations on billions of packets per second [19]. Running (pieces of) the control plane at such speeds would lead to almost "instantaneous" convergence, leaving the propagation time of the messages as the primary bottleneck. Besides speed, offloading control plane tasks to hardware would also help by making them traffic-aware. For instance, it enables to update forwarding entries consistently with real-time traffic volumes rather than in a random order.

Research questions Given the opportunity and the need, we argue that it is time to revisit the control plane's design and implementation by considering the problem of offloading parts of it to hardware. This redesign opens the door to multiple research questions including: Which pieces of the control plane should be offloaded? What are the benefits? and How can we overcome the fundamental hardware limitations? These fundamental limitations come mainly from the very limited instruction set (e.g., no floating point) and the memory available (i.e., around tens of megabytes [19]) of programmable network hardware. We start to answer these questions in this paper and make two contributions.

Hardware-Accelerated Network Control Planes

Edgar Costa Molero⁽¹⁾, Stefano Vissicchio⁽²⁾, Laurent Vanbever⁽¹⁾

> (1)ETHzürich

(2)



Software-based control planes have room for improvements



Software-based control planes have room for improvements

1 Reaction time



It can take seconds to minutes to detect failures



Software-based control planes have room for improvements

Reaction time

2 Compute

It can take minutes to recompute an entire forwarding table



Software-based control planes have room for improvements

Reaction time

2 Compute

3 Update

It takes ~100us to update a single forwarding entry



Modern programmable devices can perform computations on billions of packets per second



Modern programmable devices can perform computations on billions of packets per second

- Read & modify packet headers e.g. to update network state
- Perform (simple) operations e.g. min & max
- Add or remove custom headers e.g. to carry routing information
- Maintain state
 - e.g. to save best paths



Yes... *but...*

sensing, notification, computation

sensing, notification, computation

Switches can precisely "sense" the network by synchronously exchanging packet counts







Switches can precisely "sense" the network by synchronously exchanging packet counts





14

Upstream switch starts probing campaigns



15

Traffic for some prefixes gets dropped





Downstream switch sends counters to upstream



17

Upstream switch detects the failure by comparing counters





sensing, notification, computation

Upon detecting a failure, switches can flood notifications network-wide

Avoid broadcast storms

Simple reliable communication

Send notification duplicates

Use maximum priority queues

Use per switch broadcast sequence numbers



sensing, notification, computation

Switches can run distributed routing protocols in hardware

22

Switches can run distributed routing protocols in hardware







Statically configured tables map prefixes to registers in memory







Registers store best paths and its attributes







Switches periodically advertise vectors to neighbors







Switches periodically advertise vectors to neighbors







Switches periodically advertise vectors to neighbors







Computing new forwarding state after a a link failure







Computing new forwarding state after a a link failure






Computing new forwarding state after a a link failure







Does it actually work?

32

Does it actually work? Yes!

33

We built a P4₁₆ prototype (we're working on a Tofino implementation)

Implementation

Implemented in P4₁₆ Compiled it to bmv2 2k LoC

Capabilities

Intra-domain destinations path-vector routing

Inter-domain destinations **BGP-like route selection**

34

We tested our implementation in a simple case study



35

Only the internal switches run the hardware-based control plane







Each switch is connected to an external peer or customer



37

We generate two TCP flows from AS1 and AS2



38

Switches monitor the traffic







39

Internal link fails, triggering the path-vector algorithm



Traffic S1- AS3







Traffic S1-AS3





External link failure triggers a prefix withdrawal









Network computes new egress and applies new policies







Could we offload control-plane tasks to the data plane?

Yes... *but...*

Programmable hardware is not limitless



Programmable hardware is not limitless

Some tasks *cannot* be offloaded

while offloading others is *not desirable*

Reliable protocols e.g. TCP requires too much state

Poor scalability of control plane tasks

hardware memory is scarce and expensive



Can we have the best of both worlds?



Hardware-software codesign

Specification



Optimization

Synthesis



Hardware-software codesign



Synthesis



Hardware-software codesign





Find out more about our "quest" https://nsg.ee.ethz.ch



Network Control Planes

How? Where? Why?!





Laurent Vanbever nsg.ee.ethz.ch

Dagstuhl Seminar Wed Apr 3 2019