## Network programmability

A primer on routing synthesis



Laurent Vanbever ETH Zürich (D-ITET)

Schloss Dagstuhl January 18 2017 Human factors are responsible for 50% to 80% of network outages

Juniper Networks, What's Behind Network Downtime?, 2008

JUL 8, 2015 @ 03:36 PM 11,261 VIEWS

### United Airlines Blames Router for Grounded Flights



FULL BIO  $\sim$ 

After a computer problem caused nearly two hours of grounded flights for United Airlines this morning and ongoing delays throughout the day, the airline announced the culprit: a faulty router.

Spokeswoman Jennifer Dohm said that the router problem caused "degraded network connectivity," which affected various applications.

A computer glitch in the airline's reservations system caused the Federal Aviation Administration to impose a groundstop at 8:26 a.m. E.T. Planes that were in the air continued to operate, but all planes on the ground were held. There were reports of agents writing tickets by hand. The ground stop was lifted around 9:47 a.m. ET.



The outage was due to one faulty Internet device

### Facebook, Tinder, Instagram suffer widespread issues



MAGE: GETTY IMAGES



BY JENNI RYALL AUSTRALIA **UPDATED:** Tuesday, Jan. 27 / 4:32 a.m. EST — A Facebook spokeswoman told Mashable that the outage was due to a change to the site's configuration systems, and not a hacker attack. "Earlier this evening many people had trouble accessing Facebook and Instagram. This was not the result of a third party attack but instead occurred after we introduced a change that affected our configuration systems. We moved quickly to fix the problem, and both services are back to 100% for everyone.", she said.

JAN 27, 2015

**UPDATED**: Tuesday, Jan. 27 / 2:14 a.m. EST — Facebook, Tinder and Twitter appear to be back to normal after a 40 minute outage and mass freak out.

The outage was due to a change to the site's configuration systems



Traders work on the floor of the New York Stock Exchange (NYSE) in July 2015. (Photo by Spencer Platt/Getty Images)

#### DOWNTIME

### UPDATED: "Configuration Issue" Halts Trading on NYSE

The article has been updated with the time trading resumed.

A second update identified the cause of the outage as a "configuration issue."

A third update added information about a software update that created the configuration issue. NYSE network operators identified the culprit of the 3.5 hour outage, blaming the incident on a "network configuration issue"

### The Internet Under Crisis Conditions Learning from September 11

Committee on the Internet Under Crisis Conditions: Learning from September 11

Computer Science and Telecommunications Board Division on Engineering and Physical Sciences

NATIONAL RESEARCH COUNCIL OF THE NATIONAL ACADEMIES

National Research Council. The Internet Under Crisis Conditions: Learning from September 11

### The Internet Under Crisis Conditions Learning from September 11

Internet advertisements rates suggest that The Internet was more stable than normal on Sept 11

Committee on the Internet Under Crisis Conditions: Learning from September 11

Computer Science and Telecommunications Board Division on Engineering and Physical Sciences

NATIONAL RESEARCH COUNCIL OF THE NATIONAL ACADEMIES

### The Internet Under Crisis Conditions Learning from September 11

Internet advertisements rates suggest that The Internet was more stable than normal on Sept 11

Committee on the Internet Under Crisis Conditions: Learning from September 11

Computer Science and Telecommunications Board Division on Engineering and Physical Sciences

NATIONAL RESEARCH COUNCIL OF THE NATIONAL ACADEMIES Information suggests that operators were watching the news instead of making changes to their infrastucture



\$	👤 Follow
----	----------

Fun fact: most BGP route leaks happen on Wednesdays, but in the weekend us humans collectively take a break! :-)



# Think of the network as a distributed system running a distributed algorithm



### This algorithm produces the forwarding state which drives Internet traffic to its destination



Operators adapt their network forwarding behavior by configuring each network device individually

### Configuring each element is often done manually, using arcane low-level, vendor-specific "languages"

#### 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
```

Juniper JunOS

```
interfaces {
   so-0/0/0 {
        unit 0 {
            family inet {
                address 10.12.1.2/24;
            family mpls;
        }
    }
   ge-0/1/0 {
        vlan-tagging;
        unit 0 {
            vlan-id 100;
            family inet {
                address 10.108.1.1/24;
            family mpls;
        }
        unit 1 {
            vlan-id 200;
            family inet {
                address 10.208.1.1/24;
            }
        }
    }
}
protocols {
    mpls {
        interface all;
    bgp {
```

### A single mistyped line is enough to bring down the entire network

#### Cisco IOS

```
redistribute bgp 700 subnets — Anything else than 700 creates blackholes family inet {
```

How can we program network-wide forwarding state in existing networks?

# The forwarding state computed by a router depends on two inputs

Forwarding state

	prefix	next-hop
1	1.0.0/24	0
2	1.0.1.0/16	1
300k	100.0.0/8	0
600k	200.99.0.0/2	4 1



# The router configuration specifies how the router compute its state

```
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
   atuant 125 1 70 0 made 255 255 255 0
```



The routing messages sent by neighboring devices



*"I can reach 1.0.0.0/24"* 

()

1

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

### Network programmability

through synthesis

Fibbing "the inputs" SyNET "the functions"

### Network programmability

through synthesis

Fibbing "the inputs" SyNET "the functions"

#### [SIGCOMM'15]

Joint work with:

Stefano Vissicchio, Olivier Tilmans and Jennifer Rexford

## Fibbing

## Fibbing

= lying

## Fibbing

to **control** router's forwarding table

### **Central Control Over Distributed Routing**



- 1 Fibbing lying made useful
- 2 Expressivity any path, anywhere
- 3 Scalability 1 lie is better than 2

### **Central Control Over Distributed Routing**



Fibbing lying made useful

1

Expressivity any path, anywhere

Scalability 1 lie is better than 2

### A router implements a function from routing messages to forwarding paths



# The forwarding paths are known, provided by the operators or by the controller



# The function is known, from the protocols' specification & the configuration



Given a path and a function, our framework computes corresponding routing messages by inverting the function



### The type of input to be computed depends on the routing protocol

Protocol	Family	Algorithm/ Function	Router Input
IGP	Link-State	Dijkstra	Network graph
BGP	Path-Vector	Decision process	Routing paths
# We focus on routers running link-state protocols that take the network graph as input and run Dijkstra

Protocol	Family	Algorithm/ Function	Router Input
IGP	Link-State	Dijkstra	Network graph
BGP	Path-Vector	Decision process	Routing paths

# 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 router



#### Let's lie to the router



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



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



# Lies are propagated network-wide by the protocol



After the injection, this is the topology seen by all routers, on which they compute Dijkstra



Now, 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



Fibbing workflow Fibbing starts from the operators requirements and a up-to-date representation of the network



path network reqs. graph

# Operators requirements are expressed in a high-level language

Syntax of Fibbing's path requirements language

$$pol$$
 $:::=$  $(s_1; \ldots; s_n)$ Fibbing Policy $s$  $:::=$  $p \mid b$ Requirement $r$  $:::=$  $p_1$  and  $p_2 \mid p_1$  or  $p_2 \mid p$ Path Req. $p$  $:::=$ Path $(n^+)$ Path Expr. $n$  $:::=$  $id \mid * \mid n_1$  and  $n_2 \mid n_1$  or  $n_2$ Node Expr. $n$  $:::=$  $id \mid * \mid n_1$  and  $n_2 \mid n_1$  or  $n_2$ Node Expr. $b$  $:::=$  $r$  as backupof  $((id_1, id_2)^+)$ Backup Req.

# Out of these, the compilation stage produces DAGs

#### Compilation



path	network	forwarding
reqs.	graph	DAGs

The augmentation stage augments the network graph with lies to implement each DAG

#### Augmentation



The optimization stage reduces the amount of lies necessary





The injection stage injects the lies in the production network



# **Central Control Over Distributed Routing**



Fibbing lying made useful

2 Expressivity any path, anywhere

> Scalability 1 lie is better than 2

TheoremFibbing can programany set of non-contradictory paths

Theorem

Fibbing can program any set of non-contradictory paths



# Fibbing can load-balance traffic on multiple paths







Links have a capacity of 1



Links have a capacity of 1

With such demands and forwarding, the lower path is congested (1.25)



Congestion can be alleviated by splitting the orange flow into two equal parts (.25)



This is impossible to achieve using a link-state protocol



#### This is easily achievable with Fibbing



# One lie is introduced, announcing the orange destination



Now E has two equal cost paths (7) to reach only the orange destination and use them both



## **Central Control Over Distributed Routing**



Fibbing lying made useful

Expressivity any path, anywhere

3 Scalability1 lie is better than 2

### Scalability
## Scalability

time to compute lies space # of lies

## Scalability

time to compute lies

space # of lies



initial

desired



For each router *r* whose next-hop for a destination *d* changes to *j*:

For each router *r* whose next-hop

for a destination *d* changes to *j*:

- Let *w* be the current path weight between *r* and *d*
- Create one virtual node v advertising d
  with a weight x < w</li>
- Connects it to *r* and *j*

Create one virtual node v advertising d
 with a weight x < w</li>

#### always possible

by reweighting the initial graph

Create one virtual node v advertising d
 with a weight x < w</li>





The resulting topology can be huge and each router needs to run Dijkstra on it

Dijkstra's algorithm complexity



## Scalability

time to compute lies

space # of lies Good news

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

# Let's us consider a simple example









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)



We can minimize the topology size using an Integer Linear Program

# While efficient, an ILP is inherently slow



Computation time matters in case of network failures





A loop is created as C starts to use A which still forwards according to the lie



### The solution is to remove the lie



### The solution is to remove the lie



Upon failures, the network topology has to be recomputed, **fast** 





Merger iteratively tries to merge lies produced by the Naive algorithm



Merger iteratively tries to merge lies produced by the Naive algorithm



Merger iteratively tries to merge lies produced by the Naive algorithm














## Let's compare the performance of Naive and Merger





Naive computes entire virtual topologies in ms



Merger is relatively slower, but still, sub-second





Naive introduces one lie per changing next-hop



Merger reduces the size of the topology by 25% on average (50% in the best case)



We implemented a fully-fledged Fibbing prototype and tested it against real routers

# We implemented a fully-fledged Fibbing prototype and tested it against real routers

2 measurements

How many lies can a router sustain?

How long does it take to process a lie?

## 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

### **Central Control Over Distributed Routing**



Fibbing lying made useful

Expressivity

any path, anywhere

Scalability 1 lie is better than 2

# Fibbing realizes some of the SDN promises today, on an existing network

Facilitate SDN deployment

SDN controller can program routers and SDN switches

Simplify controller implementation

most of the heavy work is still done by the routers

Maintain operators' mental model

good old protocols running, easier troubleshooting

### Check out our webpage http://fibbing.net







AMS-IX and 1 other follow

Michiel Appelman @michielappelman · Aug 21 Interesting concept and cool webpage: fibbing.net – Central Control Over Distributed Routing

Olivier Bonaventure Retweeted



ACM SIGCOMM @ACMSIGCOMM · Aug 20

SIGCOMM 2015 best paper award: "Central Control Over Distributed Routing" by Vissicchio et. Al., conferences.sigcomm.org/sigcomm/2015/p... #sigcomm2015



Brian Krent @BrianKrent · Nov 13 "Central Control Over Distributed Routing" fibbing.net/files/sig15.pdf

"Fibbing: Small Lies for Better Networks"

CSAIL at MIT follows

John Evdemon @jevdemon · Nov 13



Fibbing is an architecture that enables central control over distributed routing. Interesting idea. fibbing.net

### Network programmability

through synthesis

Fibbing "the inputs"



# Fibbing is limited 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



#### Outputs



Network specification (N)

Physical topology (φ<sub>N</sub>)

High-level requirements (φ<sub>R</sub>)

#### Network specification (N)

Physical topology (φ<sub>N</sub>)

High-level requirements (φ<sub>R</sub>)

#### A set of Datalog *rules* that formalize how routers build their forwarding state

Fwd(Net, Node, Next) :Route(Net, Node, Next, Proto),
SetAD(Protocol, Node, Cost)
minAD(Net, Node, Cost)

Network specification (N)

Physical topology (φ<sub>N</sub>)

High-level requirements (φ<sub>R</sub>)

A set of constraints over the *input* predicates of the Datalog program

Network specification (N)

Physical topology (φ<sub>N</sub>)

High-level requirements (φ<sub>R</sub>)

A set of constraints over the *output* predicates of the Datalog program

"Traffic from R1 to R5 should flow via R2 and R3"

#### Outputs



problem

Given  $N, \varphi_{N, \varphi_{R}}$ 

SyNet generates a Datalog input *I* such that the topology  $\varphi_N$  and routing  $\varphi_R$  constraints are satisfied for the given network specification *N* 

problem	Given $N, \varphi_N, \varphi_R$
	SyNet generates a Datalog input <i>I</i> such that the topology $\varphi_N$ and routing $\varphi_R$ constraints are satisfied for the given network specification <i>N</i>
challenge	this is undecidable (in general)

problem	Given <i>N</i> , φ <sub>N</sub> , φ <sub>R</sub>
	SyNet generates a Datalog input <i>I</i> such that the topology $\varphi_N$ and routing $\varphi_R$ constraints are satisfied for the given network specification <i>N</i>
challenge	this is undecidable (in general)
key ideas	make the problem finite, use divide-and-conquer convert into a satisfiability question (on SMT constraints) "scale" using domain-specific heuristics

# SyNET can generate configurations for (small) networks

# routers



## SyNET can generate configurations for (small) networks

# routers

		4	9	16
	static	1.8s	18.2s	116.1s
# protocols	static, OSPF	4.2s	37.0s	197.0s
	static, OSPF, BGP	13.8s	189.4s	577.4s

### SyNET can generate configurations for (small) networks (new version goes to 81 routers)



		4	9	16
	static	1.8s	18.2s	116.1s
# protocols	static, OSPF	4.2s	37.0s	197.0s
	static, OSPF, BGP	13.8s	189.4s	577.4s

## Check out our webpage synet.ethz.ch



#### Automatic vs. Manual Configuration

Bouting protocols are complex. Moreover, protocols often have complex interdependencies. For example, BGP uses interdomain routing costs as input for selecting the best route. Not surprisingly, the majority of network downtimes are caused by incorrect
### Network programmability

through synthesis

Fibbing "the inputs" SyNET "the functions"

# Raw network programmability is only the beginning, what about configuration updates?

#### [SIGCOMM'11, INFOCOM'12, TON'12, TON'13, INFOCOM'13, TON'17]



# Raw network programmability is only the beginning, what about network visibility?

vanbever\_millefeuille\_hotnets\_2018.pdf (page 1 of 7) ~ EV Q Q A Q Search Mille-Feuille: Putting ISP traffic under the scalpel Olivier Tilmans <sup>1</sup>\*, Tobias Bühler <sup>8</sup>, Stefano Vissicchio <sup>†</sup>, Laurent Vanbever <sup>8</sup> <sup>‡</sup> Université catholique de Louvain, <sup>§</sup> ETH Zürich, <sup>†</sup> University College London <sup>‡</sup>olivier.tilmans@uclouvain.be,<sup>§</sup>{buehlert, lvanbever}@ethz.ch, †s.vissicchio@cs.ucl.ac.uk Mille-Feuille 🛛 🖇 🕺 😽 ABSTRACT Output For Internet Service Provider (ISP) operators, getting an ac-11 ms (>10 ms) for traffix to pf (Google) Rees Tondery Statistics curate picture of how their network behaves is challenging. mirror 🗾 for y ma 52 footional <u>مي</u> Given the traffic volumes that their networks carry and the mirror pi for a ms impossibility to control end-hosts, ISP operators are typically forced to randomly sample traffic, and rely on aggre-Figure 1: From high-level requirements, Mille-Feuille actigated statistics. This provides coarse-grained visibility, at a

In this paper, we present *Mille-Feuille*, a novel monitoring architecture that provides fine-grained visibility over ISP traffic. Mille-Feuille schedules activation and deactivation

time resolution that is far from ideal (seconds or minutes).

vates and deactivates fine-grained mirroring rules networkwide, to capture thin slices of traffic, at scale, in  $\mathcal{O}(ms)$ .

[Hotnets'16]

operators lack control of which aggregated statistics are renorted by each device, and when This non-determinism (in

#### Network controller



### Network programmability

A primer on routing synthesis



Laurent Vanbever www.vanbever.eu

Schloss Dagstuhl January 18 2017