Self-Adjusting Networks: Vision, Solutions, Challenges Stefan Schmid

"We cannot direct the wind, but we can adjust the sails." (Folklore)



Acknowledgements:





It`s a Great Time to Be a Networking Researcher!



Credits: George Varghese

It`s a Great Time to Be a Networking Researcher!



Enables and motivates self-adjusting networks!



Credits: George Varghese

It's High Time! Explosive Traffic

~

and for

î P

NETFLIX

Datacenters ("hyper-scale")



Interconnecting networks:
a critical infrastructure
of our digital society.



It's High Time! Explosive Traffic

M

a se l'

î P

NETFLIX

Datacenters ("hyper-scale")



Interconnecting networks:
a critical infrastructure
of our digital society.



Credits: Marco Chiesa

It's High Time!

Reality vs Requirements

Today, dependability requirements stand in contrast with reality:

Countries disconnected

Data Centre > Networks

Google routing blunder sent Japan's Internet dark on Friday

Another big BGP blunder

By Richard Chirgwin 27 Aug 2017 at 22:35	40 🖵	SHARE V

Last Friday, someone in Google fat-thumbed a border gateway protocol (BGP) advertisement and sent Japanese Internet traffic into a black hole.

The trouble began when The Chocolate Factory "leaked" a big route table to Verizon, the result of which was traffic from Japanese giants like NTT and KDDI was sent to Google on the expectation it would be treated as transit.

Passengers stranded

British Airways' latest Total Inability To Support Upwardness of Planes* caused by Amadeus system outage

Stuck on the ground awaiting a load sheet? Here's why

By Gareth Corfield 19 Jul 2018 at 11:16

109 SHARE V



Even 911 affected

Officials: Human error to blame in Minn. 911 outage

According to a press release, CenturyLink told department of public safety that human error by an employee of a third party vendor was to blame for the outage

Aug 16, 2018

Duluth News Tribune

SAINT PAUL, Minn. — The Minnesota Department of Public Safety Emergency Communication Networks division was told by its 911 provider that an Aug. 1 outage was caused by human error.

Even tech-savvy companies struggle:



It's High Time!

Reality vs Requirements

Today, dependability requirements stand in contrast with reality:

Countries disconnected

Data Centre > Networks

Google routing blunder sent Japan's Internet dark on Friday

Another big BGP blunder

|--|

Last Friday, someone in Google fat-thumbed a border gateway protocol (BGP) advertisement and sent Japanese Internet traffic into a black hole.

The trouble began when The Chocolate Factory "leaked" a big route table to Verizon, the result of which was traffic from Japanese giants like NTT and KDDI was sent to Google on the expectation it would be treated as transit.

Passengers stranded

British Airways' latest Total Inability To Support Upwardness of Planes* caused by Amadeus system outage

Stuck on the ground awaiting a load sheet? Here's why

By Gareth Corfield 19 Jul 2018 at 11:16

109 🖵 SHARE 🔻



Even 911 affected

Officials: Human error to blame in Minn. 911 outage

According to a press release, CenturyLink told department of public safety that human error by an employee of a third party vendor was to blame for the outage

Aug 16, 2018

Duluth News Tribune

SAINT PAUL, Minn. — The Minnesota Department of Public Safety Emergency Communication Networks division was told by its 911 provider that an Aug. 1 outage was caused by human error.



Even tech-savvy companies struggle:



Agenda Three Use Cases



Passau, Germany

Agenda Three Use Cases



Especially Under Failures (Policy Compliance)



Especially Under Failures (Policy Compliance)



Cluster with globally reachable services

Cluster with internally accessible services

Especially Under Failures (Policy Compliance)



reachable services

accessible services

Especially Under Failures (Policy Compliance)



reachable services

accessible services

Especially Under Failures (Policy Compliance)



Especially Under Failures (Policy Compliance)



If link (G,X) fails and traffic from G is rerouted via Y and C to X: X announces (does not block) G and H as it comes from C. (Note: BGP.)





→ Reachability: Can traffic from ingress port A reach B?



→ Reachability: Can traffic from ingress port A reach B?

…> Loop-freedom: Do forwarding
rules imply loop-free routes?



→ Reachability: Can traffic from ingress port A reach B?

--> Loop-freedom: Do forwarding
rules imply loop-free routes?

→ Policy: Does traffic from A to B never go via C?



→ Reachability: Can traffic from ingress port A reach B?

---> Loop-freedom: Do forwarding
rules imply loop-free routes?

→ Policy: Does traffic from A to B never go via C?

→ Waypoint enforcement: Is traffic from A to B always routed via a node C (e.g., an IDS)?



→ Reachability: Can traffic from ingress port A reach B?

---> Loop-freedom: Do forwarding rules imply loop-free routes?

→ Policy: Does traffic from A to B never go via C?

→ Waypoint enforcement: Is traffic from A to B always routed via a node C (e.g., an IDS)?

... and everything even under failures?!



Router configurations (Cisco, Juniper, etc.)

Formal language which supports automated analysis



Router configurations (Cisco, Juniper, etc.)

Formal language which supports automated analysis





Challenge: Hard Even for Computers?

→ NORDUnet: provider for Nordic countries

→ 24 MPLS routers, running Juniper OS, >30,000 labels!



Challenge: Hard Even for Computers?

NORDUnet: provider for Nordic countries

→ 24 MPLS routers, running Juniper OS, >30,000 labels!



For specific networks such as MPLS: feasible and fast! Tools such as P-Rex or AalWiNes do it in secs for MPLS: reduction to automata theory, polynomial-time.



Hard Even for Computers?

NORDUnet: provider for Nordic countries

→ 24 MPLS routers, running Juniper OS, >30,000 labels!



For specific networks such as MPLS: feasible and fast! Tools such as P-Rex or AalWiNes do it in secs for MPLS: reduction to automata theory, polynomial-time. But general networks more challenging.



- → Approaches: Petri games, Stackelberg games, UPPAAL Stratego...
- ---> But synthesis slower than verification

- → Approaches: Petri games, Stackelberg games, UPPAAL Stratego...
- ---> But *synthesis slower* than verification
- \dashrightarrow An opportunity for using AI!
- ideally AI+FM: guarantees from formal
 methods, performance from AI
- For example: synthesize with AI then verify with formal methods
- ---> Examples: DeepMPLS, DeepBGP, ...



- → Approaches: Petri games, Stackelberg games, UPPAAL Stratego...
- ---> But *synthesis slower* than verification
- \dashrightarrow An opportunity for using AI!
- ideally AI+FM: guarantees from formal
 methods, performance from AI
- For example: synthesize with AI then verify with formal methods
- ---> Examples: DeepMPLS, DeepBGP, ...





- → Approaches: Petri games, Stackelberg games, UPPAAL Stratego...
- ---> But *synthesis slower* than verification
- \dashrightarrow An opportunity for using AI!
- ideally AI+FM: guarantees from formal
 methods, performance from AI
- For example: synthesize with AI then verify with formal methods
- → Examples: DeepMPLS, DeepBGP, …
- ... and what about quantitative properties?





A Possible Starting Point: The AalWiNes Tool



Online demo: <u>https://demo.aalwines.cs.aau.dk/</u> Source code: <u>https://github.com/DEIS-Tools/AalWiNes</u> Paper: <u>https://schmiste.github.io/conext20.pdf</u>

Agenda Three Use Cases



Passau, Germany

Agenda Three Use Cases



Let's go back to datacenter use case: Moore's Law of Datacenters

---> Recall: explosive growth of demand

→ Problem: network equipment reaching capacity limits → Transistor density rates stalling

- → "End of Moore's Law in networking"
- Hence: more equipment, larger networks
- Resource intensive and:
 inefficient



Annoying for companies, opportunity for researchers
Root Cause

Fixed and Demand-Oblivious Topology

How to interconnect?



Root Cause

Fixed and Demand-Oblivious Topology



Root Cause

Fixed and Demand-Oblivious Topology



© 	© 	© 	© ■ 	© ■ 	© ■ 	© 	©













The Motivation

Much Structure in the Demand

Empirical studies:

traffic matrices sparse and skewed



Microsoft

destinations

traffic bursty over time



Hypothesis: this can be exploited.

Sounds Crazy? Emerging Enabling Technology.



H2020:

"Photonics one of only five key enabling technologies for future prosperity."

US National Research Council: "Photons are the new Electrons."

Enabler

Novel Reconfigurable Optical Switches

---> **Spectrum** of prototypes

- \rightarrow Different sizes, different reconfiguration times
- → From our last years' ACM **SIGCOMM** workshop OptSys



Example

Optical Circuit Switch

---> Optical Circuit Switch rapid adaption of physical layer



\rightarrow Based on rotating mirrors

Optical Circuit Switch

By Nathan Farrington, SIGCOMM 2010

The Big Picture



Now is the time!

Unique Position

Demand-Aware, Self-Adjusting Systems



Question 1:

How to Quantify such "Structure" in the Demand?

Complexity Map



temporal complexity



25

Complexity Map



25

Complexity Map





Our approach: iterative randomization and compression of trace to identify dimensions of structure.

Complexity Map





Our approach: iterative randomization and compression of trace to identify dimensions of structure.

Griner et al., Sigmetrics 2020 Question 2:

Given This Structure, What Can Be Achieved? Metrics and Algorithms?

A first insight: entropy of the demand.









More than an analogy!





Agenda Three Use Cases



Passau, Germany

- Automation and programmability: enables more adaptable networks
- ---> Attractive for:
 - ---> Fine-grained traffic engineering (e.g., at Google)
 - ---> Accounting for changes in the demand (spatio-temporal structure)
 - ---> Security policy changes
 - ---> Service relocation
 - ---> Maintenance work
 - ---> Link/node failures





Enabled by SDN, it has become "easy" to quickly change route to blue route.



But still need clever algorithms! Updates are asynchronous, may lead to temporal inconsistencies.



But still need clever algorithms! Updates are asynchronous, may lead to temporal inconsistencies.

Again: Formal Methods for Self-Adjusting Updates



Vision: self-adjusting networks could synthesize even their algorithms! "Ex machina": e.g., parametrized.

Examples: NetSynth, Latte

- Already "in the making"!
- NetSynth (PLDI'15): supports any LTL property and hence operator preferences. Then: standard framework to synthesize schedule.
- ---> Latte (PER'20): fast Petri net model and synthesis



Examples: NetSynth, Latte

- Already "in the making"!
- NetSynth (PLDI'15): supports any LTL property and hence operator preferences. Then: standard framework to synthesize schedule.
- ---> Latte (PER'20): fast Petri net model and synthesis


Examples: NetSynth, Latte

- Already "in the making"!
- NetSynth (PLDI'15): supports any LTL property and hence operator preferences. Then: standard framework to synthesize schedule.
- ---> Latte (PER'20): fast Petri net model and synthesis



Challenges of Self-Adjusting Networks

Challenges

- ---> **Performance** of formal methods? Opportunity: algorithm engineering!
- ---> Use of **AI**: to speed up synthesis and deal with complexity?
- ---> Limitations of automation: can networks detect themselves, when the need "help from the operator"?
- ---> **Data:** How to learn about and/or predict demand? Telemetry?
- ---> Programmability vs *security*?

Example: Security

---> Enabling technology like SDN often deployed "in software" ---> E.g., virtual switches in datacenters



Virtual Switches



Virtual switches reside in the server's virtualization layer (e.g., Xen's Dom0). Goal: provide connectivity and isolation.

Complexity: Parsing



Parsing must be fast!



Unified packet parser is fast, but complex.

Credits:

Marco Chiesa









Conclusion

- ---> A vision: self-adjusting networks
- ---> Example 1: policy-compliant networks
 - \rightarrow self-verifying
 - \rightarrow self-repairing
- ---> Example 2: demand-aware topologies
- ---> On both fronts: tip of the iceberg!
- → E.g., self-adjusting networks further supported by telemetry (data) and AI (e.g., prediction)



Thank you!

References

AalWiNes: A Fast and Quantitative What-If Analysis Tool for MPLS Networks

Peter Gjøl Jensen, Morten Konggaard, Dan Kristiansen, Stefan Schmid, Bernhard Clemens Schrenk, and Jiri Srba.

16th ACM International Conference on emerging Networking EXperiments and Technologies (**CoNEXT**), Barcelona, Spain, December 2020.

On the Complexity of Traffic Traces and Implications

Chen Avin, Manya Ghobadi, Chen Griner, and Stefan Schmid. ACM **SIGMETRICS** and ACM Performance Evaluation Review (**PER**), Boston, Massachusetts, USA, June 2020.

<u>Toward Demand-Aware Networking: A Theory for Self-Adjusting Networks</u> (Editorial) Chen Avin and Stefan Schmid.

ACM SIGCOMM Computer Communication Review (CCR), October 2018.

<u>Cerberus: The Power of Choices in Datacenter Topology Design (A Throughput Perspective)</u> Chen Griner, Johannes Zerwas, Andreas Blenk, Manya Ghobadi, Stefan Schmid, and Chen Avin. ACM **SIGMETRICS** and ACM Performance Evaluation Review (**PER**), Mumbai, India, June 2022.

Latte: Improving the Latency of Transiently Consistent Network Update Schedules Mark Glavind, Niels Christensen, Jiri Srba, and Stefan Schmid. 38th International Symposium on Computer Performance, Modeling, Measurements and Evaluation (PERFORMANCE) and ACM Performance Evaluation Review (PER), Milan, Italy, November 2020.

Taking Control of SDN-based Cloud Systems via the Data Plane (Best Paper Award) Kashyap Thimmaraju, Bhargava Shastry, Tobias Fiebig, Felicitas Hetzelt, Jean-Pierre Seifert, Anja Feldmann, and Stefan Schmid. ACM Symposium on SDN Research (**SOSR**), Los Angeles, California, USA, March 2018.

Backup Slides







Original Routing

One failure: push 30: route around (v_2, v_3)



Original Routing

One failure: push 30: route around (v_2, v_3)



Original Routing

One failure: push 30: route around (v_2, v_3)

Two failures: first push 30: route around (v₂,v₃) Push recursively 40: route around (v₂,v₆)

Which demand has more structure?

Which demand has more structure?

More uniform

More structure

Spatial vs temporal structure

- ---> Two different ways to generate same traffic matrix:
 - \rightarrow Same non-temporal structure
- ---> Which one has more structure?



Spatial vs temporal structure

- ---> Two different ways to generate same traffic matrix:
 - \rightarrow Same non-temporal structure
- ---> Which one has more structure?



Systematically?



Information-Theoretic Approach
"Shuffle&Compress"



Increasing complexity (systematically randomized)

More structure (compresses better)







Bonus Material



Hogwarts Stair

Bonus Material



Golden Gate Zipper

Bonus Material



In HPC