Principles of self-\* communication networks: Data-driven optimization and verification Stefan Schmid (TU Berlin)

> Al in Networking Summer School 2022

#### Self-\* Networks: Why?

#### Self-\* Networks: Why?

... where \*={configuring, repairing, optimizing, "driving", ...}

#### Self-\* Networks: Why?

... where \*={**configuring**, repairing, optimizing, "driving", ...}

#### Networks Are Complex: Even Tech-Savvy Companies Struggle



We discovered a misconfiguration on this pair of switches that caused what's called a *"bridge loop"* in the network.

A network change was [...] executed incorrectly [...] more "stuck" volumes and added more requests to the *re-mirroring storm*.





Service outage was due to a series of internal network events that corrupted router data tables.

Experienced a network connectivity issue [...] interrupted the airline's flight departures, airport processing and reservations systems



#### Networks Are Complex: Even Tech-Savvy Companies Struggle



We discovered a misconfiguration on this pair of switches that caused what's called a *"bridge loop"* in the network.

A network change was [...] executed incorrectly [...] more "stuck" volumes and added more requests to the *re-mirroring storm*.





Service outage was due to a series of internal network events that corrupted router data tables.

Experienced a network connectivity issue [...] interrupted the airline's flight departures, airport processing and reservations systems



Most of them: due to human errors.

#### Particularly Challenging for Humans: Reasoning about Policy-Compliance under Failures



#### Particularly Challenging for Humans: Reasoning about Policy-Compliance under Failures







#### Particularly Challenging for Humans: Reasoning about Policy-Compliance under Failures



Routers and switches store list of forwarding rules, and conditional failover rules.





Sysadmin responsible for:

• **Reachability:** Can traffic from ingress port A reach egress port B?



Sysadmin responsible for:

- **Reachability:** Can traffic from ingress port A reach egress port B?
- **Loop-freedom:** Are the routes implied by the forwarding rules loop-free?



Sysadmin responsible for:

- **Reachability:** Can traffic from ingress port A reach egress port B?
- **Loop-freedom:** Are the routes implied by the forwarding rules loop-free?
- **Policy:** Is it ensured that traffic from A to B never goes via C?



Sysadmin responsible for:

- **Reachability:** Can traffic from ingress port A reach egress port B?
- **Loop-freedom:** Are the routes implied by the forwarding rules loop-free?
- **Policy:** Is it ensured that traffic from A to B never goes via C?
- Waypoint enforcement: Is it ensured that traffic from A to B is always routed via a node C?



k failures = possibilities ٠ ٠ А E.g. IDS, firewall

Sysadmin responsible for:

- **Reachability:** Can traffic from ingress port A reach egress port B?
- **Loop-freedom:** Are the routes implied by the forwarding rules loop-free?
- **Policy:** Is it ensured that traffic from A to B never goes via C?
- Waypoint enforcement: Is it ensured that traffic from A to B is always routed via a node C?

... and everything even under multiple failures?!

# Can we automate such tests or even self-repair?

#### Another Reason:

### Flexibilities and Optimization Opportunities

#### Another Reason:

## Flexibilities and Optimization Opportunities

... where \*={configuring, repairing, **optimizing**, "driving", ...}

#### Flexibilities: Along 3 Dimensions



Somewhere in beautiful Germany...

#### Flexibilities: Along 3 Dimensions



Somewhere in beautiful Germany...

#### Flexibilities: Along 3 Dimensions





 Based on freespace optics



 Based on freespace optics



- Based on freespace optics
- Reconfiguration in ~10 μs:



Digital Micromirror Devices (DMDs)





#### ProjecToR in More Details: DMDs



- Each micromirror can be turned on/off
- Essentially a 0/1-image: e.g., array size 768 x 1024
- Direction of the diffracted light can be finely tuned





#### ProjecToR in More Details: Coupling DMDs with angled mirrors



#### ProjecToR in More Details: Coupling DMDs with angled mirrors



#### **Empirical Motivation**

## **Empirical Motivation: Structure**

*"less than 1% of the rack pairs account for 80% of the total traffic"* 

*"only a few ToRs switches are hot and most of their traffic goes to a few other ToRs"* 

*"over 90% bytes flow in elephant flows"* 

## Spatial and Temporal Locality

#### traffic matrices sparse and skewed

traffic bursty over time



#### Facebook



Time (seconds)
#### The Vision



# Roadmap

- How much structure is there in the data? A systematic approach.
- Exploiting structure in data: an example
- Self-repairing networks



# Roadmap

- How much structure is there in the data? A systematic approach.
- Exploiting structure in data: an example
- Self-repairing networks





Traffic matrix of two different **distributed ML** applications (GPU-to-GPU): Which one has *more structure*?



Traffic matrix of two different **distributed ML** applications (GPU-to-GPU): Which one has *more structure*?

# **Temporal Structure**



Two different ways to generate *same traffic matrix* (same non-temporal structure): Which one has *more structure*?

## **Temporal Structure**



Two different ways to generate *same traffic matrix* (same non-temporal structure): Which one has *more structure*?



Two different ways to generate *same traffic matrix* (same non-temporal structure): Which one has *more structure*?

# A Principled Approach: The Trace Complexity

- An information-theoretic approach: how can we measure the entropy (rate) of a traffic trace?
- Henceforth called the trace complexity
- Simple approximation: *"shuffle&compress"* 
  - Remove structure by iterative *randomization*
  - Difference of compression *before and after* randomization: structure

















# **Complexity Map**: Entropy ("complexity") of traffic traces.



# **Complexity Map**: Entropy ("complexity") of traffic traces.



Size = product of entropy

# **Complexity Map**: Entropy ("complexity") of traffic traces.



**Observation**: different applications feature quite significant (and different!) temporal and nontemporal structures.

- Facebook clusters: DB, WEB, HAD
- HPC workloads: CNS, Multigrid
- Distributed Machine Learning (ML)
- Synthetic traces like pFabric

# Roadmap

- How much structure is there in the data? A systematic approach.
- Exploiting structure in data: an example
- Self-repairing networks



# A Simple Example

Demand-Aware Network Designs of Bounded Degree Chen Avin, Kaushik Mondal, and Stefan Schmid. 31st International Symposium on Distributed Computing (**DISC**), Vienna, Austria, October 2017.









# More Formally

#### Input: Workload **Destinations** $\frac{2}{65}$ $\frac{2}{65}$ $\frac{3}{65}$ $\frac{1}{65}$ $\frac{2}{65}$ $\frac{2}{65}$ $\frac{1}{65}$ $\frac{1}{65}$ $\frac{2}{65}$ Sources $\frac{1}{65}$ $\frac{2}{65}$ $\frac{1}{65}$ 2 ( $\frac{3}{65}$ $\frac{1}{65}$ 65 $\frac{3}{65}$ $\frac{3}{65}$ $\frac{2}{65}$ $\frac{3}{65}$

 $\mathcal{D}$ 

#### Output: Constant-Degree DAN



Ν



Sources

#### Input: Workload

#### Output: Constant-Degree DAN





Sources

# **Objective: Expected Route Length**



# Remark

• Can represent demand matrix as a demand graph

sparse distribution  $\boldsymbol{\mathcal{D}}$ 

Destinations

sparse graph  $G(\mathcal{D})$ 

2





# Some Examples

- DANs of  $\Delta = 3$ :
  - E.g., complete binary tree
  - $d_N(u,v) \le 2 \log n$
  - Can we do **better** than **log n**?



- DANs of  $\Delta$  = 2:
  - E.g., set of lines and cycles





Rewinding the Clock: Degree-Diameter Tradeoff  $\mathcal{E}$  ach network with n nodes and max degree  $\Delta$  >2 must have a diameter of at least  $\log(n)/\log(\Delta-1)-1$ . Example: constant  $\Delta$ ,  $\log(n)$  diameter

K .udos to: Pedro Gasa


#### Is there a better tradeoff in DANs?

#### Sometimes, DANs can be much better!

Example 1: low-degree demand



If **demand graph** is of degree  $\Delta$ , it is trivial to design a **DAN** of degree  $\Delta$  which achieves an *expected route length of 1*.

Just take DAN = demand graph!

#### Sometimes, DANs can be much better!

#### **Example 2: skewed demand**



If **demand** is highly skewed, it is also possible to achieve an *expected route length of O(1)* in a constant-degree DAN.



#### Sometimes, DANs can be much better!

#### Example 2: skewed demand



Toward Demand-Aware Networking: A Theory for Self-Adjusting Networks. Chen Avin and Stefan Schmid. ACM SIGCOMM CCR, October 2018 If **demand** is highly skewed, it is also possible to achieve an *expected route length of O(1)* in a constant-degree DAN.



#### So on what does it depend?

### So on what does it depend?



We argue (but still don't fully know!): on the "entropy" of the demand!





#### Intuition: Entropy Lower Bound



### Lower Bound Idea: Leverage Coding or Datastructure

#### Destinations

		1	2	3	4	5	6	7	
Sources	1	0	$\frac{2}{65}$	$\frac{1}{13}$	$\frac{1}{65}$	$\frac{1}{65}$	$\frac{2}{65}$	$\frac{3}{65}$	
	2	$\frac{2}{65}$	0	$\frac{1}{65}$	0	0	0	$\frac{2}{65}$	
	3	$\frac{1}{13}$	$\frac{1}{65}$	0	$\frac{2}{65}$	0	0	$\frac{1}{13}$	
	4	$\frac{1}{65}$	0	$\frac{2}{65}$	0	$\frac{4}{65}$	0	0	
	5	$\frac{1}{65}$	0	$\frac{3}{65}$	$\frac{4}{65}$	0	0	0	
	6	$\frac{2}{65}$	0	0	0	0	0	$\frac{3}{65}$	
	7	$\frac{3}{65}$	$\frac{2}{65}$	$\frac{1}{13}$	0	0	$\frac{3}{65}$	0	

• DAN just for a *single (source) node 3* 



- How good can this tree be? Cannot do better than Δ-ary Huffman tree for its destinations
- Entropy lower bound on ERL known for binary trees, e.g. *Mehlhorn* 1975

### Lower Bound Idea: Leverage Coding or Datastructure

An optimal "ego-tree" for this source!

 $\frac{2}{65}$  $\overline{13}$  $\overline{65}$ Sources  $\frac{2}{65}$  $\frac{1}{65}$ 

Destinations

- DAN just for a *single (source) node 3*
- How good can this tree be? Cannot do better than Δ-ary Huffman tree for its destinations
- Entropy lower bound on ERL known for binary trees, e.g. *Mehlhorn* 1975

### So: Entropy of the Entire Demand



- Compute ego-tree for each source node
- Take *union* of all ego-trees
- Violates *degree restriction* but valid lower bound



### Entropy of the *Entire* Demand: Sources *and* Destinations

#### Do this in **both dimensions**: EPL $\geq \Omega(\max\{H_{\Delta}(Y|X), H_{\Delta}(X|Y)\})$



### Entropy of the *Entire* Demand: Sources *and* Destinations

#### Do this in **both dimensions**: EPL $\geq \Omega(\max\{H_{\Delta}(Y|X), H_{\Delta}(X|Y)\})$



Demand-Aware Network Designs of Bounded Degree. Chen Avin, Kaushik Mondal, and Stefan Schmid. **DISC**, 2017.

#### Achieving Entropy Limit: Algorithms



### **Ego-Trees Revisited**

 ego-tree: optimal tree for a row (= given source)



## **Ego-Trees Revisited**

 ego-tree: optimal tree for a row (= given source)





Can we merge the trees *without distortion* and *keep degree low*?



### An Analogy

Static vs dynamic demandaware networks!?

DANs vs SANs?

# An Analogy to Coding



if demand *arbitrary* and *unknown* 



















#### Analogous to *Datastructures*: Oblivious...

- Traditional, **fixed** BSTs do not rely on any assumptions on the demand
- Optimize for the worst-case
- Example demand:

 Items stored at O(log n) from the root, uniformly and independently of their

frequency

Corresponds to max possible demand!



#### ... Demand-Aware ...

- Demand-aware fixed BSTs can take advantage of *spatial locality* of the demand
- E.g.: place frequently accessed elements close to the root
- E.g., Knuth/Mehlhorn/Tarjan trees
- Recall example demand: 1,...,1,3,...,3,5,...,5,7,...,7,...,log(n),...,log(n)
  - Amortized cost O(loglog n)

Amortized cost corresponds to *empirical entropy of demand*!



#### ... Self-Adjusting!

- Demand-aware reconfigurable BSTs can additionally take advantage of temporal locality
- By moving accessed element to the root: amortized cost is *constant*, i.e., O(1)
  - Recall example demand:
    1,...,1,3,...,3,5,...,5,7,...,7,...,log(n),...,log(n)



#### Datastructures

Oblivious

#### Demand-Aware

#### Self-Adjusting



Lookup **O(log n)**  Exploit spatial locality: empirical entropy O(loglog n) Exploit temporal locality as well: O(1)

#### Analogously for Networks



#### DAN









Const degree (e.g., expander): route lengths *O(log n)* 

Exploit spatial locality

Exploit temporal locality as well

Avin, S.: Toward Demand-Aware Networking: A Theory for Self-Adjusting Networks. **SIGCOMM CCR** 2018.

### Roadmap

- How much structure is there in the data? A systematic approach.
- Exploiting structure in data: an example
- Self-repairing networks





k failures = possibilities ٠ ٠ А E.g. IDS, firewall

Sysadmin responsible for:

- **Reachability:** Can traffic from ingress port A reach egress port B?
- **Loop-freedom:** Are the routes implied by the forwarding rules loop-free?
- **Policy:** Is it ensured that traffic from A to B never goes via C?
- Waypoint enforcement: Is it ensured that traffic from A to B is always routed via a node C?

... and everything even under multiple failures?!

# Can we automate such tests or even self-repair?

# Can we automate such tests or even self-repair?



Yes! Sometimes even *fast*: with formal methods (enhanced with AI...).

### How (MPLS) Networks Work

• Forwarding based on top label of label stack



Default routing of two flows

### How (MPLS) Networks Work

• Forwarding based on top label of label stack



### How (MPLS) Networks Work



Default routing of two flows

### Fast Reroute Around 1 Failure

• Forwarding based on top label of label stack (in packet header)



Default routing of two flows

• For failover: push and pop label



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

### Fast Reroute Around 1 Failure

• Forwarding based on top label of label stack (in packet header)


# Fast Reroute Around 1 Failure

• Forwarding based on top label of label stack (in packet header)



### 2 Failures: Push *Recursively*



**Original** Routing

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

Two failures: first push 30: route around (v<sub>2</sub>,v<sub>3</sub>) *Push recursively* 40: route around (v<sub>2</sub>,v<sub>6</sub>)

### 2 Failures: Push *Recursively*







## Leveraging Automata-Theoretic Approach



MPLS configurations, Segment Routing etc. Pushdown Automaton and Prefix Rewriting Systems Theory



MPLS configurations, Segment Routing etc. Pushdown Automaton and Prefix Rewriting Systems Theory

# **Network Model**

• Network: a 7-tuple  $N = (V, E, I_v^{in}, I_v^{out}, \lambda_v, L, \delta_v^F)$  Nodes Nodes

# **Network Model**

• Network: a 7-tuple

$$N = (V, E, I_v^{in}, I_v^{out}, \lambda_v, L, \delta_v^F)$$

Interface function: maps outgoing interface to next hop node and incoming interface to previous hop node  $\lambda_v: I_v^{in} \cup I_v^{out} \to V$ 

That is: 
$$(\lambda_v(in), v) \in E$$
 and  $(v, \lambda_v(out)) \in E$ 

# **Network Model**

• Network: a 7-tuple

$$N = (V, E, I_v^{in}, I_v^{out}, \lambda_v, L, \delta_v^F)$$

**Routing function**: for each set of failed links  $F \subseteq E$ , the routing function

$$\delta_v^F: I_v^{in} \times L^* \to 2^{(I^{out} \times L^*)}$$

defines, for all incoming interfaces and packet headers, outgoing interfaces together with modified headers.

# Routing



• Example: routing (in)finite sequence of tuples



# **MPLS Network Model**

• MPLS supports three specific operations on header sequences:

 $Op = \{swap(\ell) \mid \ell \in L\} \cup \{push(\ell) \mid \ell \in L\} \cup \{pop\}$ 

• The local routing table can then be defined as



• Local link protection function defines backup interface



### Example Rules: *Regular Forwarding* on Top-Most Label



### A Complex and Big Formal Language: Why Polynomial Time?



Arbitrary number k of failures: How can I avoid checking all  $\binom{n}{k}$  many options?!

### **Classic Result in Automata Theory**

- Classic result by **Büchi** 1964: the set of all reachable configurations of a pushdown automaton a is regular set
- Hence, we can operate only on Nondeterministic Finite Automata (NFAs) when reasoning about the pushdown automata



Julius Richard Büchi 1924-1984 Swiss logician

- The resulting regular operations are all polynomial time
  - Important result of model checking

### **Tool and Query Language**

Part 1: Parses query and constructs Push-Down System (PDS)

• In Python 3

Part 2: Reachability analysis of constructed PDS

• Using *Moped* tool

# failures header path header

### Regular query language



### query processing flow

## **Example: Traversal Testing With 2 Failures**

Traversal test with k=2: Can traffic starting with [] go through s5, under up to k=2 failures?



## Speeding things up with Deep Learning? And synthesis.

DeepMPLS: Fast Analysis of MPLS Configurations Using Deep Learning. Fabien Geyer and Stefan Schmid. **IFIP Networking**, Warsaw, Poland, May 2019.

# Deep Learning for MPLS: DeepMPLS (s. talk by Fabien Geyer)

- Yes sometimes without losing guarantees
- Extend graph-based neural networks





Label for Swa

Input label

Network topologies and MPLS rules

Predict counter-examples and fixes



### Challenges of Self-\* Networks

- Can a self-\* network realize its limits?
- E.g., when quality of **input data** is not good enough?
- When to hand over to human? Or fall back to "safe/oblivious mode"?
- Can we learn from self-driving cars?





## Thank you! Questions?



#### **Demand-aware networks**

Survey of Reconfigurable Data Center Networks: Enablers, Algorithms, Complexity Klaus-Tycho Foerster and Stefan Schmid. SIGACT News, June 2019. Toward Demand-Aware Networking: A Theory for Self-Adjusting Networks (Editorial) Chen Avin and Stefan Schmid. ACM SIGCOMM Computer Communication Review (CCR), October 2018. Measuring the Complexity of Network Traffic Traces Chen Griner, Chen Avin, Manya Ghobadi, and Stefan Schmid. arXiv, 2019. Demand-Aware Network Design with Minimal Congestion and Route Lengths Chen Avin, Kaushik Mondal, and Stefan Schmid. 38th IEEE Conference on Computer Communications (INFOCOM), Paris, France, April 2019. **Distributed Self-Adjusting Tree Networks** Bruna Peres, Otavio Augusto de Oliveira Souza, Olga Goussevskaia, Chen Avin, and Stefan Schmid. 38th IEEE Conference on Computer Communications (INFOCOM), Paris, France, April 2019. Efficient Non-Segregated Routing for Reconfigurable Demand-Aware Networks Thomas Fenz, Klaus-Tycho Foerster, Stefan Schmid, and Anaïs Villedieu. IFIP Networking, Warsaw, Poland, May 2019. DaRTree: Deadline-Aware Multicast Transfers in Reconfigurable Wide-Area Networks Long Luo, Klaus-Tycho Foerster, Stefan Schmid, and Hongfang Yu. IEEE/ACM International Symposium on Quality of Service (IWQoS), Phoenix, Arizona, USA, June 2019. Demand-Aware Network Designs of Bounded Degree Chen Avin, Kaushik Mondal, and Stefan Schmid. 31st International Symposium on Distributed Computing (DISC), Vienna, Austria, October 2017. SplayNet: Towards Locally Self-Adjusting Networks Stefan Schmid, Chen Avin, Christian Scheideler, Michael Borokhovich, Bernhard Haeupler, and Zvi Lotker. IEEE/ACM Transactions on Networking (TON), Volume 24, Issue 3, 2016. Early version: IEEE IPDPS 2013. Characterizing the Algorithmic Complexity of Reconfigurable Data Center Architectures Klaus-Tycho Foerster, Monia Ghobadi, and Stefan Schmid. ACM/IEEE Symposium on Architectures for Networking and Communications Systems (ANCS), Ithaca, New York, USA, July 2018.

### What-if analysis

#### P-Rex: Fast Verification of MPLS Networks with Multiple Link Failures

Jesper Stenbjerg Jensen, Troels Beck Krogh, Jonas Sand Madsen, Stefan Schmid, Jiri Srba, and Marc Tom Thorgersen. 14th ACM International Conference on emerging Networking EXperiments and Technologies (**CoNEXT**), Heraklion/Crete, Greece, December 2018.

Polynomial-Time What-If Analysis for Prefix-Manipulating MPLS Networks

Stefan Schmid and Jiri Srba.

37th IEEE Conference on Computer Communications (INFOCOM), Honolulu, Hawaii, USA, April 2018.

### Secure sampling and dataplane

Preacher: Network Policy Checker for Adversarial Environments
Kashyap Thimmaraju, Liron Schiff, and Stefan Schmid.
38th International Symposium on Reliable Distributed Systems (SRDS), Lyon, France, October 2019.
MTS: Bringing Multi-Tenancy to Virtual Switches
Kashyap Thimmaraju, Saad Hermak, Gabor Retvari, and Stefan Schmid.
USENIX Annual Technical Conference (ATC), Renton, Washington, USA, July 2019.
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**

# How Predictable is Traffic?

Even if reconfiguration fast, control plane (e.g., data collection) can become a bottleneck. However, many good examples:

- Machine learning applications
- Trend to disaggregation (specialized racks)
- Datacenter communication dominated by elephant flows
- Etc.



ML workload (GPU to GPU): deep convolutional neural network *Predictable from their dataflow graph*