Theory of Demand-Aware Networks

Stefan Schmid (University of Vienna)

IEEE ComSoc Autumn School on Network Slicing and Data Driven Communication

The Trend: Flexibilities

Flexibilities: Along 3 Dimensions



Somewhere in beautiful Germany...

Flexibilities: Along 3 Dimensions



Somewhere in beautiful Germany...

Flexibilities: Along 3 Dimensions



Flexibilities Enable Demand-Aware Networks

But when are they useful?

A Simple Answer

Demand-Oblivious Networks =



Seriously: We believe, often, in practice!



In theory: traffic matrix uniform and static In practice: skewed and dynamic

Empirical Motivation

Observation 1:

- Many rack pairs exchange little traffic
- Only some *hot rack pairs* are active

Observation 2:

 Some source racks send large amounts of traffic to many other racks



Microsoft data: 200K servers across 4 production clusters, cluster sizes: 100 - 2500 racks. Mix of workloads: MapReduce-type jobs, index builders, database and storage systems.

So: How much structure is there?



How to *measure* it? And which *types of structures*? E.g., temporal structure in addition to non-temporal structure? *Tricky*!

Often only intuitions in the literature...

"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"



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*?

... and it *is* intuitive! 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*?



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

- 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.



- Traditional networks are optimized for the "worst-case" (all-to-all communication traffic)
- Example, fat-tree topologies: provide full bisection bandwidth





Good in the worst case **but**: cannot leverage different **temporal** and **non-temporal** structures of traffic traces!



To exploit **temporal** structure, plexity Map need adaptive demand-aware ("self-adjusting") networks. Uniform Bursty 1.00.9 Non-temporal complexity *Good* in the worst case *but*: 0.8 cannot leverage different 0.7 temporal and non-temporal 0.6 structures of traffic traces! 0.5 Skewed 0.4 Non-temporal structure could 0.3 0.3 0.9 be exploited already with *static* 0.4 0.81.00 Temporal complexity demand-aware networks!



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**



Goal: Design self-adjusting networks which leverage *both* dimensions of structure!





Further Reading

Measuring the Complexity of Packet Traces. Avin, Ghobadi, Griner, Schmid. **ArXiv** 2019.

So: How to design networks which exploit this structure? How good can they be?

Metrics again!

A Simple Example: Network Design

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

Output: Constant-Degree DAN

Destinations



 \mathcal{D}





Sources



Sources

Input: Workload

Output: Constant-Degree DAN





Sources

Objective: Expected Route Length



Remark

• Can represent demand matrix as a **demand graph**

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

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







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

Remark: Hardness Proof

- Example Δ = 2: A Minimum Linear Arrangement (MLA) problem
 - A "Virtual Network Embedding Problem", VNEP
 - *Minimize sum* of lengths of virtual edges



- Example Δ = 2: A Minimum Linear Arrangement (MLA) problem
 - A "Virtual Network Embedding Problem", VNEP
 - *Minimize sum* of lengths of virtual edges



- Example Δ = 2: A Minimum Linear Arrangement (MLA) problem
 - A "Virtual Network Embedding Problem", VNEP
 - *Minimize sum* of lengths of virtual edges



- Arrangement (March design) ir.
 A "Virtual March and so is bedding Problem"
 Minim March and so is bedding Problem • Example Δ = 2: A Mipi inear
 - Jedding Problem", VNEP



- Example Δ = 2: A Mini design inear Arrangement (Management (Management (Management Arrangement (Management Arrangement Arrangement Arrangement Arrangement Arrangement Arrangement Arrangement (Management Arrangement (Management Arrangement Arrangement Arrangement Arrangement Arrangement Arrangement Arrangement (Management Arrangement (Management Arrangement Arran
- But what about > 2? *Embedding* problem still hard, but we have an additional degree of freedom:

Do topological flexibilities make problem easier or harder?!



A new knob for optimization!



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



Is there a better tradeoff in DANs?

Sometimes, DANs can be much better!

Example 1: low-degree demand



If **demand graph** is of degree Δ , it is trivial to design a **DAN** of degree Δ 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 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!

		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	

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 $\geq \Omega(\max\{H_{\Delta}(Y|X), H_{\Delta}(X|Y)\})$



 \mathcal{D}

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*?



From Trees to Networks



Idea: Degree Reduction



Taking union of ego-trees results in **high degree** u and v will appear in many ego-trees

Demand-Aware Network Designs of Bounded Degree. Chen Avin, Kaushik Mondal, and Stefan Schmid. **DISC**, 2017. Node *h* **helps edge (u, v)** by participating in *ego-tree(u)* as a relay node toward *v* and *in ego-tree(v)* as a relay toward *u*
But: How to design DANs which also leverage *temporal structure*?



Inspiration from self-adjusting datastructures again!

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:

 $1, \dots, 1, 3, \dots, 3, 5, \dots, 5, 7, \dots, 7, \dots, \log(n), \dots, \log(n)$ $\longleftrightarrow \qquad \longleftrightarrow \qquad \longleftrightarrow \qquad \longleftrightarrow \qquad \longleftrightarrow \qquad \longleftrightarrow \qquad$ *many many many many many many many*

 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.

Algorithms for Self-Adjusting Networks



From trees to networks!







A Dynamic Ego-Tree: Splay Tree



Uncharted Landscape!

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





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.

A survey!

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*