Revolutionizing Datacenter Networks via Self-Adjusting Optical Topologies

Stefan Schmid (INET @ TU Berlin)

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

Acknowledgements:





Revolutionizing Datacenter Networks via Self-Adjusting Optical Topologies



#### INET @ TU Berlin

Acknowledgements:





## The Challenge

Growing Traffic, e.g., due to AI/ML



Fixed and Demand-Oblivious Topology

How to interconnect? Focus on this talk: scale-out network.

© 	© ■ 	© ■ 	© ■ ■ 	© ■ 	© ■ 	© 	© 

Fixed and Demand-Oblivious Topology

- Example: fat-tree topology (bi-regular)
  - → 2 types of switches: top-of-rack (ToR) connect to hosts, additional switches connecting switches to increase throughput



Fixed and Demand-Oblivious Topology

- …> Example: expander topology (uni-regular)
  - ightarrow Only 1 type of switches:

lower installation and management overheads



Fixed and Demand-Oblivious Topology

…> Example: expander topology (uni-regular)

→ Only 1 type of switches: lower installation and management overheads



Highway which ignores actual traffic: frustrating!



Fixed and Demand-Oblivious Topology

Example: expander topology (uni-regular) --->

 $\rightarrow$  Only 1 type of switches: lower installation and management overheads



Highway which ignores actual traffic: frustrating!

Many flavors, but in common: fixed and **oblivious** ("ignorant") to actual demand.

00

]©	∎©°	∎©°	<b>□</b> ◎ <sup>∞</sup>	<b>≣</b> ⊚ %	∎©°	∎©°	
		=		=	=	=	

∎© °	∎© °° 	© °	© *	© 	© 	© 	© 













## Analogy



Golden Gate Zipper

### The Motivation

Much Structure in the Demand: Complexity Map



#### Traffic is also clustered: Small Stable Clusters



#### Opportunity: *exploit* with little reconfigurations!

7

### Sounds crazy?

Optical circuit switches are already deployed



## The Big Picture



Now is the time!

### Potential Gain



### Potential Gain



#### Reality more complicated than that...

#### Challenge: Traffic Diversity

#### Diverse patterns:

- → Shuffling/Hadoop: all-to-all
- → All-reduce/ML: ring or tree traffic patterns → Elephant flows
- → Query traffic: skewed → Mice flows
- → Control traffic: does not evolve but has non-temporal structure

#### Diverse requirements:

→ ML is bandwidth hungry, small flows are latencysensitive



#### Diverse topology components:

→ demand-oblivious and demand-aware

> Demandoblivious Demandaware

















13







### A Solution: Cerberus



We have a first approach:

*Cerberus*\* serves traffic on the "best topology"! (Optimality open)

\* Griner et al., ACM SIGMETRICS 2022

#### Flow Size Matters

On what should topology type depend? We argue: flow size.

#### Flow Size Matters

On what should topology type depend? We argue: flow size.



---- **Observation 1:** Different apps have different flow size distributions.

#### Flow Size Matters Similar tradeoff for 400Gbps or 800Gpbs Flow transmission time (40Gbps) 100ns 100µs 1ms **1**S 100ms 1μS 10µS 10ms 1 CDF of bytes Websearch- 2010 0.75 Datamining- 2011 Hadoop- 2015 0.5 Pareto distribution 0.25 0 $10^{3}$ 10<sup>5</sup> $10^{6}$ 10<sup>8</sup> 10<sup>9</sup> $10^{4}$ $10^{7}$ **10<sup>10</sup>** Flow size (bytes)

---> **Observation 1:** Different apps have different flow size distributions.

---- Observation 2: The transmission time of a flow depends on its size.

### Flow Size Matters



- ---> Observation 1: Different apps have different flow size distributions.
- ---> Observation 2: The transmission time of a flow depends on its size.
- ••• Observation 3: For small flows, flow completion time suffers if network needs to be reconfigured first.
- ---> Observation 4: For large flows, reconfiguration time may amortize.

### Flow Size Matters



- ---> Observation 1: Different apps have different flow size distributions.
- ----> Observation 2: The transmission time of a flow depends on its size.
- ••• Observation 3: For small flows, flow completion time suffers if network needs to be reconfigured first.
- ---> Observation 4: For large flows, reconfiguration time may amortize.











15





Scheduling: Small flows go via static switches...





Scheduling: ... medium flows via rotor switches...





Scheduling: ... and large flows via demand-aware switches (if one available, otherwise via rotor).

#### Cerberus Framework



#### vs Rotor-Net and Expander-Net

### Throughput Analysis

Demand Matrix



Metric: throughput of a demand matrix...

> Abdu et al., SC 2016 Namyar et al., SIGCOMM 2021

### Throughput Analysis

Demand Matrix





Met	:ri	<b></b>	thro	oughput
of	а	dem	and	matrix

... is the maximal scale down factor by which traffic is feasible  $0 \le \theta(T) \le 1$ .

> Abdu et al., SC 2016 Namyar et al., SIGCOMM 2021

### Throughput Analysis

Demand Matrix



 $\times \theta(T) =$ 



Metric: throughput
of a demand matrix...

... is the maximal scale down factor by which traffic is feasible  $0 \le \theta(T) \le 1$ .

Throughput of network  $\theta^*$ : worst case T

Abdu et al., SC 2016 Namyar et al., SIGCOMM 2021

### Throughput: Expander-Net

Demand Matrix



Permutation matrix





Namyar et al., SIGCOMM 2021

### Throughput: Demand-Aware

Demand Matrix





Permutation matrix

Permutation matrix is the best demand matrix for demand-aware net!

## Throughput: Cerberus



### Throughput: Summary

Demand Matrix



	expander-net	rotor-net	Cerberus
BW-Tax	<ul> <li>✓</li> </ul>	1	×
LT-Tax	×	1	✓
$\theta(T)$	Thm 2	Thm 3	Thm 5
$ heta^*$	0.53	0.45	Open
Datamining	0.53	0.6	0.8 (+33%)
Permutation	0.53	0.45	$\approx 1 \ (+88\%)$
Case Study	0.53	0.66	0.9 (+36%)

For the given input parameters:

 $n, k, R_d, R_r$ 

### Throughput Bounds

- → Throughput bounds for many designs not fully understood yet
- → Particularly simple demand-aware network design: Vermillion\*
- → Periodic reconfigurations (like Sirius) which can be adapted
- $\rightarrow$  How close can we approximate self-adjusting netowrks?



Addanki et al., arXiv 2025: https://arxiv.org/pdf/2405.20869

Addanki et al., Vermillion: https://arxiv.org/pdf/2504.09892

# More benefits of optical & reconfigurable switching

- Reconfigurable datacenter networks naturally support heterogeneous network elements
- ---> And therefore also *incremental* hardware upgrades

See interview with Amin Vahdat, Google in CACM: https://www.youtube.com/watch?v=IxcV1gu8ETA



### Research at INET (1)

---> Experimental frameworks

#### ExReC: Experimental Framework for Reconfigurable Networks Based on Off-the-Shelf Hardware

Johannes Zerwas TU Munich, Germany

Stefan Schmid TU Berlin, University of Vienna & Fraunhofer SIT Germany & Austria

#### ABSTRACT

In order to meet the increasingly stringent throughput and latency requirements on datacenter networks, several innovative network architectures based on reconfigurable optical topologies have been proposed. Examples include demandChen Avin Ben-Gurion University of the Negev Beer-Sheva, Israel

Andreas Blenk TU Munich & University of Vienna Germany & Austria

#### ACM Reference Format:

Johannes Zerwas, Chen Avin, Stefan Schmid, and Andreas Blenk. 2021. ExReC: Experimental Framework for Reconfigurable Networks Based on Off-the-Shelf Hardware. In Symposium on Architectures for Networking and Communications Systems (ANCS '21), December 13–16. 2021. Lavfette. IN. USA. ACM. New York. NY. USA.

### Research at INET (2)

---> How to control RDCNs on the network layer (using local routing)

#### Duo: A High-Throughput Reconfigurable Datacenter Network Using Local Routing and Control

JOHANNES ZERWAS, TUM School of Computation, Information and Technology, Technical University of Munich, Germany CSABA GYÖRGYI, University of Vienna and ELTE Eötvös Loránd University, Austria and Hungary ANDREAS BLENK, Siemens AG, Germany STEFAN SCHMID, TU Berlin & Fraunhofer SIT, Germany CHEN AVIN, Ben-Gurion University, Israel

The performance of many cloud-based applications critically depends on the capacity of the underlying datacenter network. A particularly innovative approach to improve the throughput in datacenters is enabled by emerging optical technologies, which allow to dynamically adjust the physical network topology, both in an oblivious or demand-aware manner. However, such topology engineering, i.e., the operation and control of dynamic datacenter networks, is considered complex and currently comes with restrictions and overheads.

We present Duo, a novel demand-aware reconfigurable rack-to-rack datacenter network design realized with a simple and efficient control plane. Duo is based on the well-known de Bruijn topology (implemented using a small number of optical circuit switches) and the key observation that this topology can be enhanced using dynamic ("opportunistic") links between its nodes.

In contract to previous systems. Duo has several desired features: i) It makes effective use of the network

### Research at INET (3)

---> Congestion control for highly dynamic networks

#### **POWERTCP:** Pushing the Performance Limits of Datacenter Networks<sup>\*</sup>

Vamsi Addanki University of Vienna TU Berlin Oliver Michel University of Vienna Princeton University Stefan Schmid University of Vienna TU Berlin

#### Abstract

Increasingly stringent throughput and latency requirements in datacenter networks demand fast and accurate congestion control. We observe that the reaction time and accuracy of existing datacenter congestion control schemes are inherently limited. They either rely only on explicit feedback about the network state (e.g., queue lengths in DCTCP) or only on varistringent performance requirements are introduced by today's trend of resource disaggregation in datacenters where fast access to remote resources (e.g., GPUs or memory) is pivotal for the overall system performance [36]. Building systems with strict performance requirements is especially challenging under bursty traffic patterns as they are commonly observed in datacenter networks [12, 16, 47, 53, 55].

These requirements introduce the need for fast and ecou

## Research at INET (4)

→ Making demand-aware RDCNs more practical (direct routing only)

#### Vermilion: A Traffic-Aware Reconfigurable Optical Interconnect with Formal Throughput Guarantees

Vamsi Addanki	Chen Avin	Goran Dario Knabe
<sup>TU Berlin</sup>	Ben-Gurion University of the Negev	<sup>TU Berlin</sup>
Giannis Patronas	Dimitris Syrivelis	Nikos Terzenidis
NVIDIA	NVIDIA	NVIDIA
Paraskevas Bakopoulos	Ilias Marinos	Stefan Schmid
<sub>NVIDIA</sub>	NVIDIA	<sup>TU Berlin</sup>
STRACT	Throughput*	

#### AB

The increasing gap between datacenter traffic volume and the capacity of electrical switches has driven the development of reconfigurable network designs utilizing optical circuit switching. Recent advancements, particularly those featuring periodic fixed-duration reconfigurations, have achieved practical end-to-end delays of just a few microseconds. However, current designs rely on multi-hop



### Research at INET (5)

#### ---> Buffering aspects

#### **ABM: Active Buffer Management in Datacenters**

Vamsi Addanki" TU Berlin Maria Apostolaki\* Princeton University

Stefan Schmid TU Berlin

#### ABSTRACT

Today's network devices share buffer across queues to avoid drops during transient congestion and absorb bursts. As the buffer-perbandwidth-unit in datacenter decreases, the need for optimal buffer utilization becomes more pressing. Typical devices use a hierarchical packet admission control scheme: First, a Buffer Management (BM) scheme decides the maximum length per queue at the device level and then an Active Oueue Management (AOM) scheme decides which packets will be admitted at the queue level. Unfortunately, the lack of cooperation between the two control schemes leads to (i) harmful interference across queues, due to the lack of isolation; (ii) increased queueing delay, due to the obliviousness to the per-queue drain time; and (iii) thus unpredictable burst tolerance. To overcome these limitations, we propose ABM, Active Buffer Management which incorporates insights from both BM and AQM. Concretely, ABM accounts for both total buffer occupancy (typically used by BM) and queue drain time (typically used by AQM). We analytically prove that ABM provides isolation, bounded buffer drain time and achieves predictable burst tolerance without sacrificing throughput. We empirically find that ABM improves the 99th nercentile FCT for short flows by up to 94% compared to the



Laurent Vanbever

Manya Ghobadi

MIT

Figure 1: BM and AQM are orthogonal in their goals, and the hierarchical scheme fundamentally limits the burst absorption capabilities of the buffer.

#### Reverie: Low Pass Filter-Based Switch Buffer Sharing for Datacenters with RDMA and TCP Traffic

Vamsi Addanki Wei Bai Stefan Schmid Maria Apostolaki TU Berlin Microsoft Research TU Berlin Princeton University

#### Abstract

The switch buffers in datacenters today are dynamically shared by traffic classes with different loss tolerance and reaction to congestion signals. In particular, while legacy applications use loss-tolerant transport, e.g., DCTCP, newer applications require lossless datacenter transport, e.g., RDMA over At a high level, the goal of a buffer-sharing scheme is to provide isolation between traffic classes, while maximizing the benefit of the buffer e.g., by absorbing bursts and achieving high throughput. Existing buffer management schemes (even recent ones) [1, 8, 15, 25] were designed considering exclusively loss-tolerant traffic (e.g., TCP variants). However,

## Research at INET (6)

---> How to efficiently collect and exploit information about flows

#### Credence: Augmenting Datacenter Switch Buffer Sharing with ML Predictions

#### TCP's Third Eye: Leveraging eBPF for Telemetry-Powered Congestion Control

Jörn-Thorben Hinz TU Berlin Vamsi Addanki TU Berlin Csaba Györgyi University of Vienna

Theo Jepsen Intel

#### ABSTRACT

For years, congestion control algorithms have been navigating in the dark, blind to the actual state of the network. They were limited to the course-grained signals that are visible from the OS kernel, which are measured locally (e.g., RTT) or hints of imminent congestion (e.g., packet loss and ECN). As applications and OSs are

#### Stefan Schmid TU Berlin

#### 1 INTRODUCTION

The volume of traffic in datacenters is increasing rapidly over time [6, 31, 33]. The throughput and latency offered by the underlying architecture and the set of protocols plays a critical role in the performance of modern cloud-based applications [26]. To this end, major research efforts over the past decade have been in

#### Vamsi Addanki Maciej Pacut TU Berlin TU Berlin

Stefan Schmid TU Berlin

#### Abstract

Packet buffers in datacenter switches are shared across all the switch ports in order to improve the overall throughput. The trend of shrinking buffer sizes in datacenter switches makes buffer sharing extremely challenging and a critical performance issue. Literature suggests that push-out buffer sharing algorithms have significantly better performance guarantees compared to drop-tail algorithms. Unfortunately, switches are unable to benefit from these algorithms due to lack of support for push-out operations in hardware. Our key observation is that drop-tail buffers can emulate push-out buffers if the future packet arrivals are known ahead of time. This suggests that aug-



### Thank you! Questions?



Slides available here:



#### Online Video Course





### YouTube Interview & CACM

Check out our **YouTube interviews** on Reconfigurable Datacenter Networks:



Revolutionizing Datacenter Networks via Reconfigurable Topologies Chen Avin and Stefan Schmid. Communications of the ACM (CACM), 2025. Watch here: <u>https://www.youtube.com/@self-adjusting-networks-course</u>



#### Websites



http://self-adjusting.net/ Project website





https://trace-collection.net/
Trace collection website



#### June'25 CACM Article

#### **Revolutionizing Datacenter Networks via Reconfigurable Topologies**

CHEN AVIN, is a Professor at Ben-Gurion University of the Negev, Beersheva, Israel STEFAN SCHMID, is a Professor at TU Berlin, Berlin, Germany

With the popularity of cloud computing and data-intensive applications such as machine learning, datacenter networks have become a critical infrastructure for our digital society. Given the explosive growth of datacenter traffic and the slowdown of Moore's law, significant efforts have been made to improve datacenter network performance over the last decade. A particularly innovative solution is reconfigurable datacenter networks (RDCNs): datacenter networks whose topologies dynamically change over time, in either a demand-oblivious or a demand-aware manner. Such dynamic topologies are enabled by recent optical switching technologies and stand in stark contrast to state-of-the-art datacenter network topologies, which are fixed and oblivious to the actual traffic demand. In particular, reconfigurable demand-aware and "self-adjusting" datacenter networks are motivated empirically by the significant spatial and temporal structures observed in datacenter communication traffic. This paper presents an overview of reconfigurable datacenter networks. In particular, we discuss the motivation for such reconfigurable architectures, review the technological enablers, and present a taxonomy that classifies the design space into two dimensions: static vs. dynamic and demand-oblivious vs. demand-aware. We further present a formal model and discuss related research challenges. Our article comes with complementary video interviews in which three leading experts, Manya Ghobadi, Amin Vahdat, and George Papen, share with us their perspectives on reconfigurable datacenter networks.

#### KEY INSIGHTS

- Datacenter networks have become a critical infrastructure for our digital society, serving explosively growing communication traffic.
- Reconfigurable datacenter networks (RDCNs) which can adapt their topology dynamically, based on innovative
  optical switching technologies, bear the potential to improve datacenter network performance, and to simplify
  datacenter planning and operations.
- Demand-aware dynamic topologies are particularly interesting because of the significant spatial and temporal structures observed in real-world traffic, e.g., related to distributed machine learning.
- The study of RDCNs and self-adjusting networks raises many novel technological and research challenges related to their design, control, and performance.

## References (1)

Revolutionizing Datacenter Networks via Reconfigurable Topologies Chen Avin and Stefan Schmid. Communications of the ACM (CACM), 2025.

Cerberus: The Power of Choices in Datacenter Topology Design (A Throughput Perspective) 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.

Mars: Near-Optimal Throughput with Shallow Buffers in Reconfigurable Datacenter Networks Vamsi Addanki, Chen Avin, and Stefan Schmid. ACM SIGMETRICS and ACM Performance Evaluation Review (PER), Orlando, Florida, USA, June 2023.

#### Duo: A High-Throughput Reconfigurable Datacenter Network Using Local Routing and Control

Johannes Zerwas, Csaba Györgyi, Andreas Blenk, Stefan Schmid, and Chen Avin. ACM SIGMETRICS and ACM Performance Evaluation Review (PER), Orlando, Florida, USA, June 2023.

#### 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.

#### Toward Demand-Aware Networking: A Theory for Self-Adjusting Networks (Editorial)

Chen Avin and Stefan Schmid. ACM SIGCOMM Computer Communication Review (CCR), October 2018.

#### Credence: Augmenting Datacenter Switch Buffer Sharing with ML Predictions

Vamsi Addanki, Maciej Pacut, and Stefan Schmid. 21st USENIX Symposium on Networked Systems Design and Implementation (NSDI), Santa Clara, California, USA, April 2024.

#### TCP's Third Eye: Leveraging eBPF for Telemetry-Powered Congestion Control

Jörn-Thorben Hinz, Vamsi Addanki, Csaba Györgyi, Theo Jepsen, and Stefan Schmid. SIGCOMM Workshop on eBPF and Kernel Extensions (eBPF), Columbia University, New York City, New York, USA, September 2023.

ABM: Active Buffer Management in Datacenters

Vamsi Addanki, Maria Apostolaki, Manya Ghobadi, Stefan Schmid, and Laurent Vanbever. ACM **SIGCOMM**, Amsterdam, Netherlands, August 2022.

## References (2)

ExRec: Experimental Framework for Reconfigurable Networks Based on Off-the-Shelf Hardware
Johannes Zerwas, Chen Avin, Stefan Schmid, and Andreas Blenk.
16th ACM/IEEE Symposium on Architectures for Networking and Communications Systems (ANCS), Virtual Conference,
December 2021.

Demand-Aware Network Design with Minimal Congestion and Route Lengths

Chen Avin, Kaushik Mondal, and Stefan Schmid. IEEE/ACM Transactions on Networking (TON), 2022.

A Survey of Reconfigurable Optical Networks

•

Matthew Nance Hall, Klaus-Tycho Foerster, Stefan Schmid, and Ramakrishnan Durairajan. Optical Switching and Networking (**OSN**), Elsevier, 2021.

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.