Networks in the Age of Distributed Computation

Stefan Schmid (TU Berlin)

Networks in the Age of Distributed Computation

Stefan Schmid (TU Berlin)



Slides credits: Marco Chiesa, Torsten Höfler

Not in this talk!

Appentation in the

In this talk!



Wired networks?

- "Cosy living room": wellunderstood and just works
- Passed test of time
- Should and cannot be changed

Wired networks!

- Place where *fantastic innovations* are happening ☺ On all layers.
- For *performance* and *dependability*
- Still: specific and interesting *contraints* due to simple but fast hardware
- DISC bonus (compared to wireless): *simple* and discrete models ⁽ⁱ⁾



Why do networks evolve? The Internet 50 years ago...



Now we live in a different era: Age of Computation

Datacenters ("hyperscale")



Data intensive applications requiring significant processing.

Age of Computation: Evidence

Datacenters ("hyperscale")

Data intensive applications requiring significant processing.

Nvidia: fastest growing company ever



Amazon buys nuclear-powered data center from Talen



Susquehanna nuclear plant in Salem Township, Penn., along with the data center in foreground. (Photo: Talen Energy)

Training even across *multiple datacenters* (and *powerplants*)!

Age of Computation: Evidence

Data intensive applications requiring significant processing.

Energy consumption and probably also computation trends will likely stay. *Kardashev Scale* even classifies civilizations by their energy use!



Nvidia: fastest growing company ever



Amazon buys nuclear-powered data center from Talen



Susquehanna nuclear plant in Salem Township, Penn., along with the data center in foreground. (Photo: Talen Energy

Training even across *multiple datacenters* (and *powerplants*)!

Datacenters ("hyperscale")

Age of Computation: More Evidence

Nobel Prizes in Physics and Chemistry...





Age of Computation: More Evidence

... and soon also in Economics and Literature?!





Datacenters ("hyperscale")

Distributed applications...



Image: state state

Datacenters ("hyperscale")

Distributed applications...



... require networks!





Image: state of the state of the

Datacenters ("hyperscale")

Networks are a critical infrastructure of digital society. Especially *to*, *from*, and *inside* datacenter networks!

Distributed applications...



... require networks!





Image: system of the system of th

Networks are a critical infrastructure of digital society. Especially *to*, *from*, and *inside* datacenter networks!

Datacenters ("hyperscale")

Distributed applications...









Challenge and Opportunity: Networks become larger and larger

- Also here: end of *Moore's Law in* networking
 - Transistor density rates stalling
- Hence: need more equipment, larger networks
- Opportunity: network itself forms large distributed system! With specialized but fast hardware.
 - E.g., in-network processing to speed up allreduce?



- A highly *decentralized problem*!
- How much packets dropped in Internet today?

- A highly *decentralized problem*! ۰
- How much packets dropped in Internet today? •
 - Not negligible.
- Because of the way we *control* congestion! •
 - A TCP sender cannot directly "see" traffic load in network...
 - ... so *opportunistically probes*: increases sending rate until loss

is this?

So **TCP** needs packet loss to determine their sending rate



- Well, huge success for decades: additive increase, multiplicative decrease (AIMD)
 - No congestion collapse since 1990s
 - Same mechanism since 30+ years, while *traffic increased by factor 1 billion*!



- Well, huge success for decades: additive increase, multiplicative decrease (AIMD)
 - No congestion collapse since 1990s
 - Same mechanism since 30+ years, while *traffic increased by factor 1 billion*!



- A little bit better: Linux' TCP CUBIC
 - Idea: increase sending rate faster until "near last packet loss"-rate



Can we do better? Significant efforts right now!

- Still: performance could be better
 - Google's BBR, QUIC, Netflix, ECN, etc.: additional signals about congestion (e.g., *latency*)
 - Also: congestion control in *datacenters* (e.g., to handle ML workloads)
- Opportunity for DISC: Many of these protocols have no theoretical underpinnings!
 - And indeed, have issues, e.g., regarding *fairness*
 - Often hard to *catch issues* empirically and or in simulations!



Theory needed! Example BBR.



- BBR: relatively *fast and large* deployment
- But with *fairness* and other *issues*
- **Needed several adjustments** and new versions still under development

Literature: Model-Based Insights on the Performance, Fairness, and Stability of BBR. Scherrer et al., ACM IMC 2022.

Example Innovation on Network Layer: Segment Routing ("Valiant Routing")

Example Innovation on Network Layer: Segment Routing ("Valiant Routing")



Example Innovation on Network Layer: Segment Routing ("Valiant Routing")





Credits: Marco Chiesa



Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only $s_2 \rightarrow d = 1.5$

all link capacities of 1



Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only $s_2 \rightarrow d = 1.5$

all link capacities of 1

S1 Cl

Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only s₂ \rightarrow d = 1.5

How to set link weights to serve this traffic? Without violating capacities and to minimize load.

all link capacities of 1



all link capacities of 1

Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only s₂ \rightarrow d = 1.5

- operator sets link weights > 0
- per-destination routing



all link capacities of 1

-> shortest path DAG

Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only s₂ \rightarrow d = 1.5

- operator sets link weights > 0
- per-destination routing
- shortest-path DAGs



all link capacities of 1

- --> shortest path DAG
- 1/2 splitting ratio

Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only s₂ \rightarrow d = 1.5

- operator sets link weights > 0
- per-destination routing
- shortest paths DAGs
- equal-split



Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only $s_2 \rightarrow d = 1.5$

- operator sets link weights > 0
- per-destination routing
- \leq 100% link utilization
- shortest paths DAGs
 - equal-split
Traffic Engineering



Only two possible demand matrices:

1. only
$$s_1 \rightarrow d = 1.5$$

2. only $s_2 \rightarrow d = 1.5$

Traditional traffic engineering:

- operator sets link weights > 0
- per-destination routing
- shortest paths DAGs
- equal-split

Traffic Engineering



Only two possible demand matrices:

1. only
$$s_1 \rightarrow d = 1.5$$

2. only $s_2 \rightarrow d = 1.5$

Traditional traffic engineering:

No link-weight assignment can attain \leq 100% link utilization!

(for both demand matrices, although in principle enough capacity available!)

What about this?!



Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only $s_2 \rightarrow d = 1.5$

Traditional traffic engineering:

- operator sets link weights > 0
- per-destination routing
- shortest paths DAGs
- equal-split

What about this?!



Careful: first flow now splits *twice*! Two more shortest paths later.

Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only $s_2 \rightarrow d = 1.5$

Traditional traffic engineering:

- operator sets link weights > 0
- per-destination routing
- shortest paths DAGs
- equal-split

Powerful Extension: Segment Routing "Valiant Routing for IP Networks"

- Can define waypoints between source and destination
 - Like Valiant routing: important technique in oblivious routing (but random waypoint)
- Shortest paths on "segments" between waypoints (and source and destination)



Traffic Engineering with Segment Routing

- Good! All links
- \leq 100% utilization



Half of traffic from s₂ via waypoint v!

Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only $s_2 \rightarrow d = 1.5$

Segment Routing:

- Can push a waypoint w between source s2 and destination d
- Then: shortest path from *s* to *w*, and shortest path from *w* to *d*

Traffic Engineering with Segment Routing

- Good! All links
- \leq 100% utilization



Half of traffic from s₂ via waypoint v!

Only two possible demand matrices:

- 1. only $s_1 \rightarrow d = 1.5$
- 2. only $s_2 \rightarrow d = 1.5$

Segment Routing:

- Can push a waypoint w between source s2 and destination d
- Then: shortest path from *s to w*,

and shortest path from *w to d*

Literature: Traffic Engineering with Joint Link Weight and Segment Optimization. Parham et al., ACM CONEXT, 2021.

Example: Many more...

• New Ethernet versions

...

- Automotive Ethernet
- Ethernet for *datacenters*



- Hollow-fiber: faster speed of light!
 - Cost(*latency*)>>>Cost(bandwidth)
- Optical and reconfigurable networks

ENTERPRISE NETWORKING

Hollow Fiber: The New Option for Low Latency

Very low latency hollow fiber services target the financial services industry today and offer another option for serving latency-sensitive applications in the broader business market.



Roadmap: Two Examples

• Resilient routing

• Datacenter networks

















Why is slow bad? Packet drops until restored!



How can a switch/router locally decide how to handle an arriving packet?



Nodes Locally Store A Forwarding *Match -> Action Table*



And what information is locally available to decide how to handle an arriving packet?



Locally Available Information: The Packet Header (e.g., Source, *Destination*)



Locally Available Information: The *Inport* of the Received Packet





Raises an Interesting Question

Can we *pre-install* local fast failover rules which ensure reachability under multiple failures? *In particular: How many failures* can be tolerated by static forwarding tables? So: How many failures can be tolerated by static forwarding tables?



If we partition the network, there is not much to do





Definition: *Connectivity k* of a network *N*: the minimum number of link deletions that partitions *N*



Ideal resilience

Given a *k*-connected graphs, we can tolerate *any k-1 link failures*.

Ideal resilience

Given a *k*-connected graphs, we can tolerate *any k-1 link failures*.

Perfect resilience

Any source *s* can always reach any destination *t* as long as the unterlying network is *physically connected*.

Ideal resilience

Given a *k*-connected graphs, we can tolerate *any k-1 link failures*.

Perfect resilience

Any source *s* can always reach any destination *t* as long as the unterlying network is *physically connected*.

Can this be achieved? Assume undirected link failures.

Ideal resilience

Given a *k*-connected graphs, we can tolerate *any k-1 link failures*.

Perfect resilience

Any source *s* can always reach any destination *t* as long as the unterlying network is *physically connected*.

Can this be achieved? Assume undirected link failures.

Spectrum of Models

Recall our switch model:



Achievable resilience depnds on *what can be matched*:

Per- destination	Per source	Incoming port	Probabilistic forwarding	Packet header rewriting
destination	Per source	port	forwarding	header rewriting

Credits: Marco Chiesa

Spectrum of Models



Example: Which level of resiliency?

Per- destination	Per source	Incoming port	Probabilistic forwarding	Packet header rewriting	Resiliency
Х					

Per-destination routing *cannot cope* with *even one* link failure


Per- destination	Per source	Incoming port	Probabilistic forwarding	Packet header rewriting	Resiliency
Х	Х	Х			?



29

Per- destination	Per source	Incoming port	Probabilistic forwarding	Packet header rewriting	Resiliency
Х	Х	Х			Yes



k disjoint paths: try one after the other, routing *back to source* each time.

Per- destination	Per source	Incoming port	Probabilistic forwarding	Packet header rewriting	Resiliency
Х		Х			?

What about this scenario? Practically important. But open problem since many years...

Per- destination	Per source	Incoming port	Probabilistic forwarding	Packet header rewriting	Resiliency
Х		Х			?

What about this scenario? Practically important. But open problem since many years...

For some special graphs we know: the answer is positive!

Ideal Resilience: Example 2-dim Torus?







Decompose torus into 2edge-disjoint Hamilton Cycles (HC)

1st Hamilton cycle



Decompose torus into 2edge-disjoint Hamilton Cycles (HC)

1st Hamilton cycle2nd Hamilton cycle



- Decompose torus into 2edge-disjoint Hamilton Cycles (HC)
- Can route in both directions: *4-arc-disjoint* HCs



- Decompose torus into 2edge-disjoint Hamilton Cycles (HC)
- Can route in both directions: *4-arc-disjoint* HCs

3-resilient routing to destination d:

- go along 1st directed HC, if hit failure, reverse direction
- if again failure switch to 2nd HC, if again failure reverse direction
- No more failures possible!

Ideal Resilience with Hamilton Cycles

Chiesa et al.: if k-connected graph has k arc disjoint Hamilton Cycles, k-1 resilient routing can be constructed!

> Chiesa et al. **On the Resiliency of Static Forwarding Tables.** IEEE/ACM Transactions on Networking (ToN), 2017.

Ideal Resilience with Hamilton Cycles

Chiesa et al.: if k-connected graph has k arc disjoint Hamilton Cycles, k-1 resilient routing can be constructed!

What about graphs which cannot be decomposed into Hamilton cycles?

Chiesa et al. **On the Resiliency of Static Forwarding Tables.** IEEE/ACM Transactions on Networking (ToN), 2017.

Ideal Resilience in General k-Connected Graphs

- Use directed trees (i.e. *arborescences*) instead of Hamilton cycles
 - Arc-disjoint, spanning, and rooted at destination
- Classic result: k-connectivity guarantees karborescence decomposition

Basic idea:

- Idea: route towards root on one arborescence
- After failure: change arborescence (e.g. in circular fashion)
- Incoming port defines current arborescence
- After k-1 failures: At least one arborescence intact



Ideal Resilience in General k-Connected Graphs

- Use directed trees (i.e. *arborescences*) instead of Hamilton cycles
 - Arc-disjoint, spanning, and rooted at destination
- Classic result: k-connectivity guarantees karborescence decomposition

Basic idea:

- Idea: route towards root on one arborescence
- After failure: change arborescence (e.g. in circular fashion)
- Incoming port defines current arborescence
- After k-1 failures: At least one arborescence intact



J. Edmonds, **Edge-disjoint branchings**. Combinatorial Algorithms, 1972.











Credits: Marco Chiesa



Credits: Marco Chiesa

General technique: routing along the same tree



When a failed link is hit...



35

... how do we choose the next arborescence?



35

But how do we choose the next arborescence?

Circular-arborescence routing:

- compute an order of the arborescences
- switch to the next arborescence when hitting a failed link

Arborescence order





Arborescence order



Go along arborescence 1 to destination...



Arborescence order



Go along arborescence 2 to destination...



Arborescence order

1 2 3 4

Go along arborescence 3 to destination...



Arborescence order



Go along arborescence 4 to destination...



Arborescence order





Intuition: each single failure may affect two arborescences

All k=4 arborescences used (2 failures disconnected affected all four): LOOP!

Resilience Criteria

Ideal resilience

Given a *k*-connected graphs, we can tolerate *any k-1 link failures*.

Perfect resilience

Any source *s* can always reach any destination *t* as long as the unterlying network is *physically connected*.

Can this be achieved? Assume undirected link failures.

Resilience Criteria

Perfect resilience is impossible to achieve in general.

Already on simple planar graphs, proof by case distinction (and indistinguishability).



Related to several DISC problems but with twist!

- Geometric routing
 - E.g., a left-hand rule can be used in planar graphs
- Local algorithms without communication
 - E.g., Balanced Incomplete Block Design (*BIBD*) can be used to minimize congestion!
- Graph exploration and connectivity problems
 - E.g., Omer Reingold's "undirected connectivity in log-space"

Many Open Questions...

- Big open question: ideal resilience conjecture
 - False? DISC experts!
- What if we can *rewrite* some header bits?
 - With log(n) bits it is easy: can remember all failures. What about less?
- What about fast rerouting in **Segment Routing** networks?
- What about *special graph classes*?
- Automated **synthesis** of tables (e.g., BDDs)

Many Open Questions...

- Big open question: ideal resilience conjecture
 - False? DISC experts!
- What if we can *rewrite* some header bits?
 - With log(n) bits it is easy: can remember all failures. What about less?
- What about fast rerouting in **Segment Routing** networks?
- What about *special graph classes*?
- Automated **synthesis** of tables (e.g., BDDs)

Literature:

On the Feasibility of Perfect Resilience with Local Fast Failover. Foerster et al., SIAM APOCS, 2021. Randomized Local Fast Rerouting for Datacenter Networks with Almost Optimal Congestion. Bankhamer et al. DISC, 2021.

Local Reroute with Segment Routing?

• Recall segment routing: shortest path routing on segments

- Fast rerouting currently under standardization at IETF
 - Good time to have impact!


How to handle at least 1 failure?

• When a *node v* on *route from s to t locally* detects failure on *link e,* it can *push a waypoint w*.



How to handle at least 1 failure?

- When a *node v* on *route from s to t locally* detects failure on *link e,* it can *push a waypoint w*.
- <u>Rule:</u> v should push a w such that the shortest path s1 (from v to w) and the shortest path s2 (from w to t) does not include e again! So can route around.



A Local Solution

- We need two definitions:
 - P-Space: the nodes which v can reach on shortest paths without using e
 - Q-Space: the nodes which can reach t on shortest paths without using e



Then: choose any waypoint w at intersection* for rerouting!

*If intersection empty, spaces must be adjacent and there is also a (different) solution.

A Local Solution...

What about

2 failures?

- We need two definitions:
 - P-Space: the nodes which v can reach on shortest paths without using e
 - **Q-Space**: the nodes which can reach *t* on shortest paths without using e



Then: choose *any waypoint w at intersection** for rerouting!

*If intersection empty, spaces must be adjacent and there is also a (different) solution.

Literature: TI-MFA: Keep Calm and Reroute Segments Fast. Foerster et al., IEEE Global Internet Symposium (GI), 2018.

Roadmap: Two Examples

• Resilient routing

• Datacenter networks



Datacenter Networks

How to interconnect racks?

© ■ 	© 	© *	© *	∎	© *	© *	© *

Datacenter Networks



Datacenter Networks





demand:





1 2 3 4 5 6 7 8

Matches demand!









new demand:





new demand:

Matches demand!



1 2 3

4 5 6 7 8



Self-adjusting networks: adapt

in a demand-aware manner!



Empirical Motivation

- Workloads have much spatial and temporal structure
 - That is, low entropy
- Can be exploited for optimization



Enabler

- Optical circuit switch
 - E.g., Google
- Adapt in microsecs!



Self-Adjusting Networks













First Deployments and a Challenge

- Google's *demand-aware* reconfigurable datacenter
- Key challenge according to Amin Vahdat: scalable and distributed control



- Optimal static network for a source
 - Huffman tree or biased binary search tree



- Optimal static network for a source
 - Huffman tree or biased binary search tree
- For entire demand: take union
 - But reduce degree



- Optimal static network for a source
 - Huffman tree or biased binary search tree
- For entire demand: take union
 - But reduce degree
- Dynamic: replace with splay tree



- Optimal static network for a source
 - Huffman tree or biased binary search tree
- For entire demand: take union
 - But reduce degree
- Dynamic: replace with splay tree
- Distributed?
 - Distributed version of splay trees?



Conclusion

- Wired networks: *different* from what you may think! And *evolving*.
- Much control is *distributed*
 - Congestion control, local fast re-routing, demand-aware networks
- A good moment to *contribute*: on publications..
 - DISC expertise where other communities got stuck?
- ... and in practice: have *impact*, e.g., at standardizations at IETF, initiatives like Ultra Ethernet Consortium



Thank you! Questions?

<u>A Survey of Fast-Recovery Mechanisms in Packet-Switched Networks</u> Marco Chiesa, Andrzej Kamisinski, Jacek Rak, Gabor Retvari, and Stefan Schmid. IEEE Communications Surveys and Tutorials (**COMST**), 2021.

On the Price of Locality in Static Fast Rerouting

Klaus-Tycho Foerster, Juho Hirvonen, Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan.

52nd IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), Baltimore, Maryland, USA, June 2022.

The Hazard Value: A Quantitative Network Connectivity Measure Accounting for Failures

Pieter Cuijpers, Stefan Schmid, Nicolas Schnepf, and Jiri Srba.

52nd IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), Baltimore, Maryland, USA, June 2022.

On the Feasibility of Perfect Resilience with Local Fast Failover

Klaus-Tycho Foerster, Juho Hirvonen, Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan.

SIAM Symposium on Algorithmic Principles of Computer Systems (APOCS), Alexandria, Virginia, USA, January 2021.

Brief Announcement: What Can(not) Be Perfectly Rerouted Locally

Klaus-Tycho Foerster, Juho Hirvonen, Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan. International Symposium on Distributed Computing (**DISC**), Freiburg, Germany, October 2020.

Improved Fast Rerouting Using Postprocessing

Klaus-Tycho Foerster, Andrzej Kamisinski, Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan. IEEE Transactions on Dependable and Secure Computing (**TDSC**), 2020.

Resilient Capacity-Aware Routing

Stefan Schmid, Nicolas Schnepf and Jiri Srba.

27th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS), Virtual Conference, March 2021.

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.

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.

Randomized Local Fast Rerouting for Datacenter Networks with Almost Optimal CongestionGregor Bankhamer, Robert Elsässer, and Stefan Schmid..International Symposium on Distributed Computing (DISC), Freiburg, Germany, October 2021.Bonsai: Efficient Fast Failover Routing Using Small ArborescencesKlaus-Tycho Foerster, Andrzej Kamisinski, Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan.49th IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), Portland, Oregon, USA, June 2019.CASA: Congestion and Stretch Aware Static Fast ReroutingKlaus-Tycho Foerster, Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan38th IEEE Conference on Computer Communications (INFOCOM), Paris, France, April 2019.Load-Optimal Local Fast Rerouting for Dense NetworksMichael Borokhovich, Yvonne-Anne Pignolet, Gilles Tredan, and Stefan Schmid.IEEE/ACM Transactions on Networking (TON), 2018.PURR: A Primitive for Reconfigurable Fast RerouteMarco Chiesa, Roshan Sedar, Gianni Antichi, Michael Borokhovich, Andrzej Kamisinski, Georgios Nikolaidis, and Stefan Schmid.15th ACM International Conference on emerging Networking EXperiments and Technologies (CoNEXT), Orlando, Florida, USA, December 2019.

Artefact Evaluation: Available, Functional, Reusable.

On the Resiliency of Static Forwarding Tables

In IEEE/ACM Transactions on Networking (ToN), 2017

M. Chiesa, I. Nikolaevskiy, S. Mitrovic, A. Gurtov, A. Madry, M. Schapira, S. Shenker

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. 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. Demand-Aware Network Design with Minimal Congestion and Route Lengths Chen Avin, Kaushik Mondal, and Stefan Schmid. IEEE/ACM Transactions on Networking (TON), 2022. 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.

Backup Slides

Intelligent Routers: A Use Case



Intelligent Routers: A Use Case


Assume: shared memory *size 3*.



Assume: shared memory *size 3*.

Scenario 1: assign buffer *opportunistically*!



Assume: shared memory *size 3*.

Scenario 1: assign buffer *opportunistically*!



Assume: shared memory *size 3*.

Scenario 1: assign buffer *opportunistically*!



Assume: shared memory *size 3*.

Scenario 1: assign buffer *opportunistically*!



Suboptimal: green packets could be transmitted in parallel, but there is no more space!

Assume: shared memory *size 3*.

Scenario 1: assign buffer *opportunistically*!



Suboptimal: green packets could be transmitted in parallel, but there is no more space!

Assume: shared memory *size 3*.



Assume: shared memory *size 3*.



Assume: shared memory *size 3*.



Assume: shared memory *size 3*.



Assume: shared memory *size 3*.



Assume: shared memory *size 3*.

Scenario 2: assign buffer *conservatively* and *keep space*.



Suboptimal: drops were unnecessary, buffer not needed for green packets!

Credence

- Traffic at switch can be *predicted* fairly well
- AI/ML could significantly *improve buffer management*...
- ... and hence *admission control and throughput*!
- Further reading:

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**), 2024.