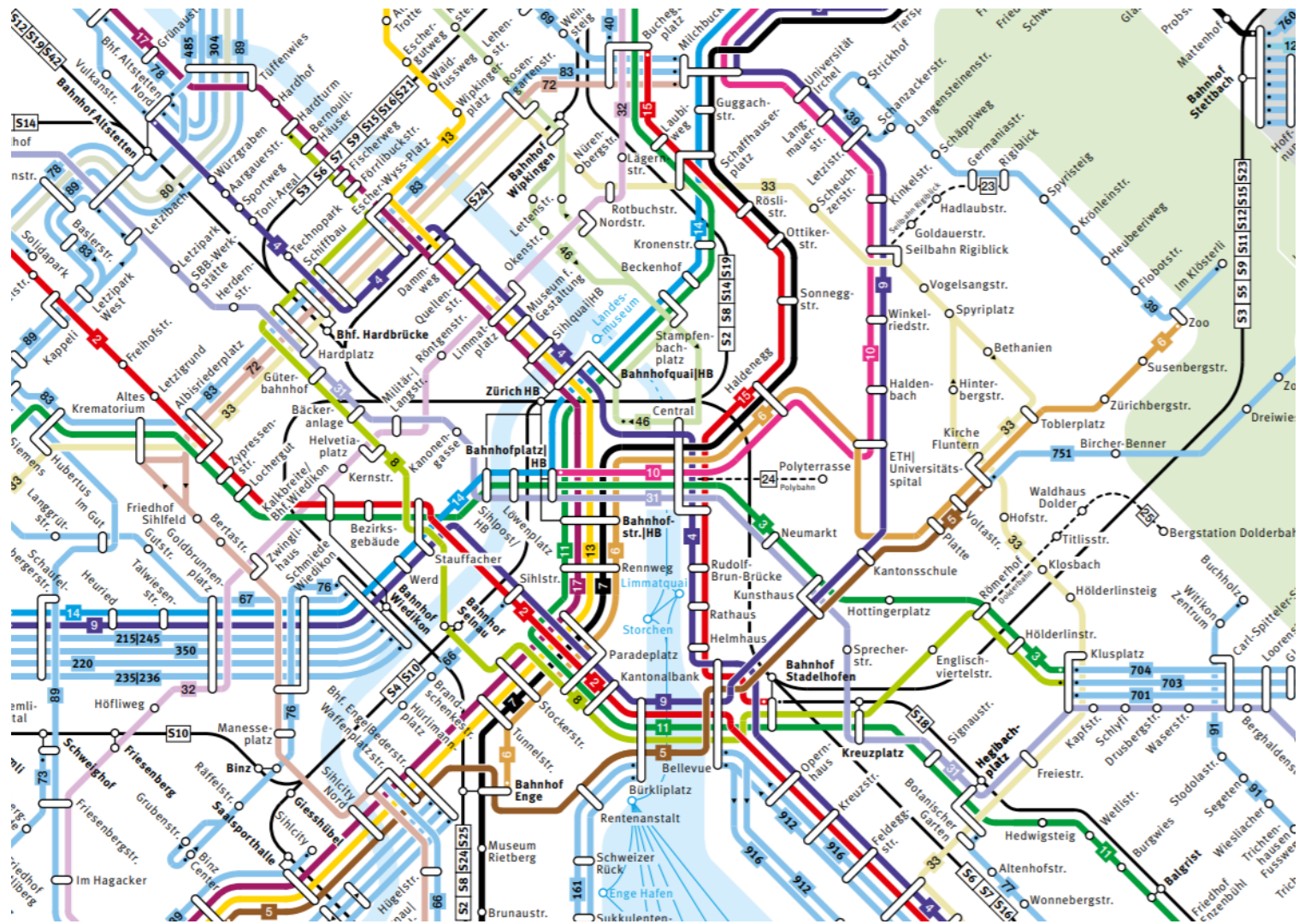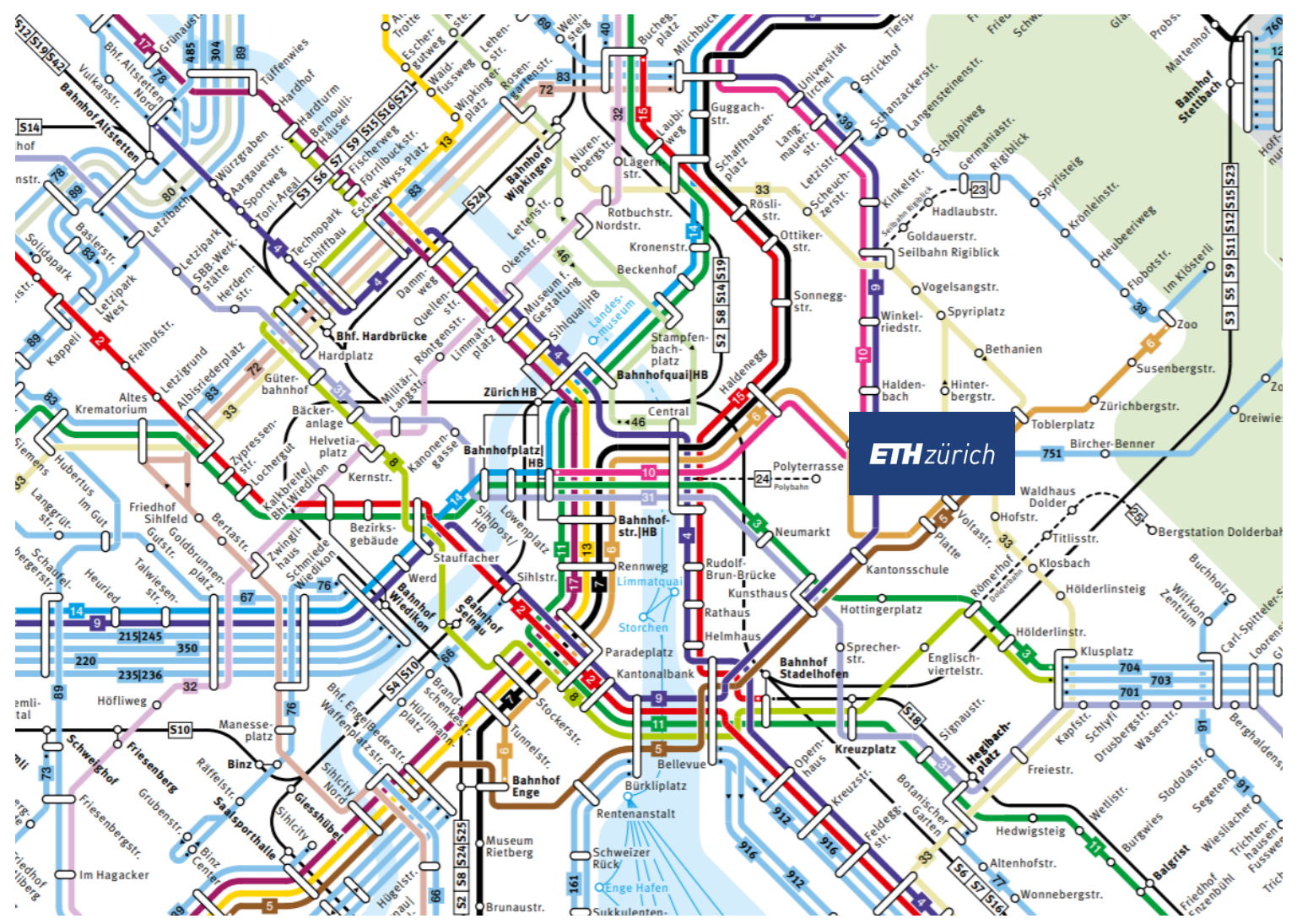# Congestion and Stretch Aware Static Fast Rerouting [appeared @INFOCOM'19]
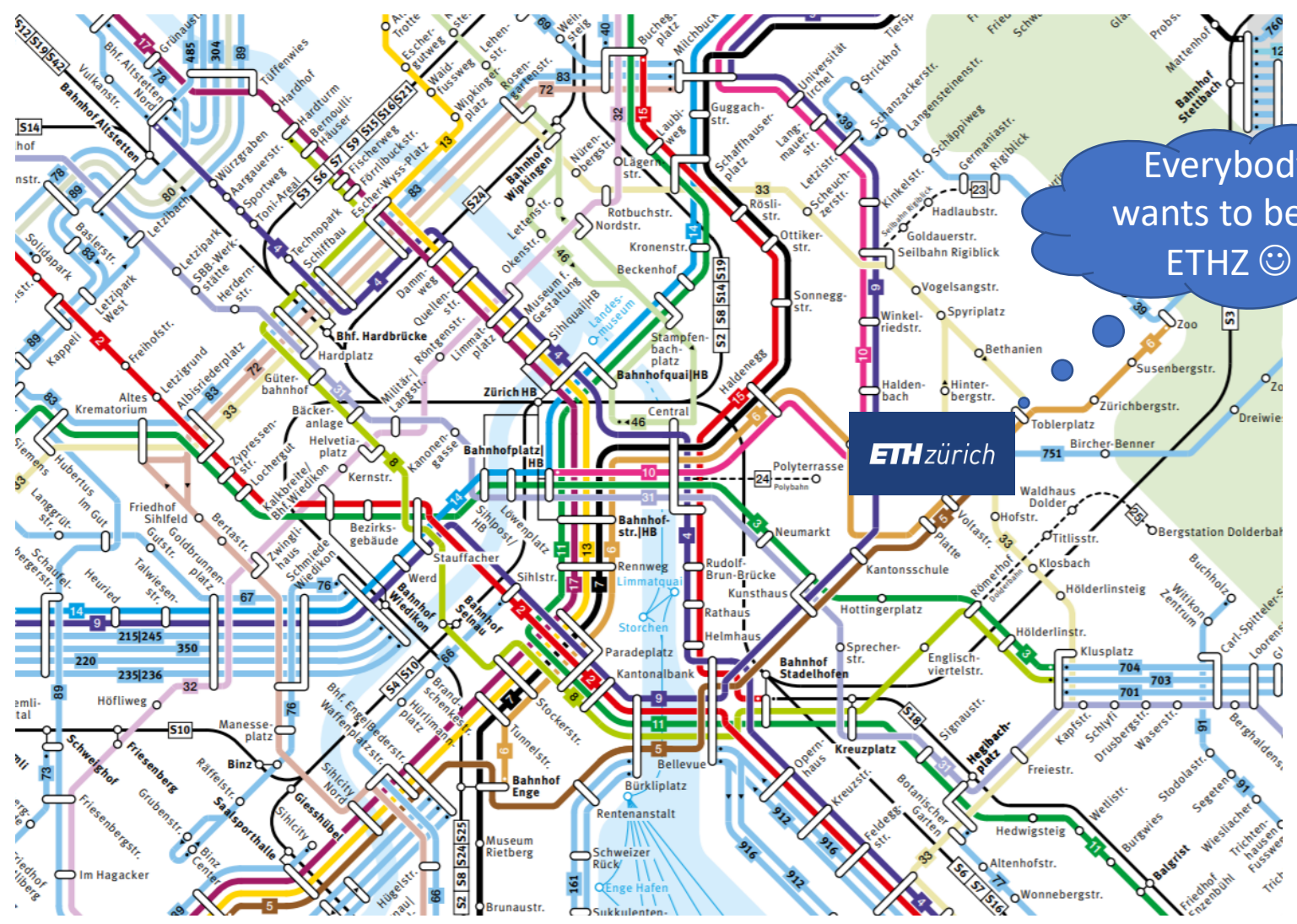
Klaus-Tycho Foerster, Yvonne-Anne Pignolet (DFINITY), Stefan Schmid, and Gilles Tredan (LAAS-CNRS)
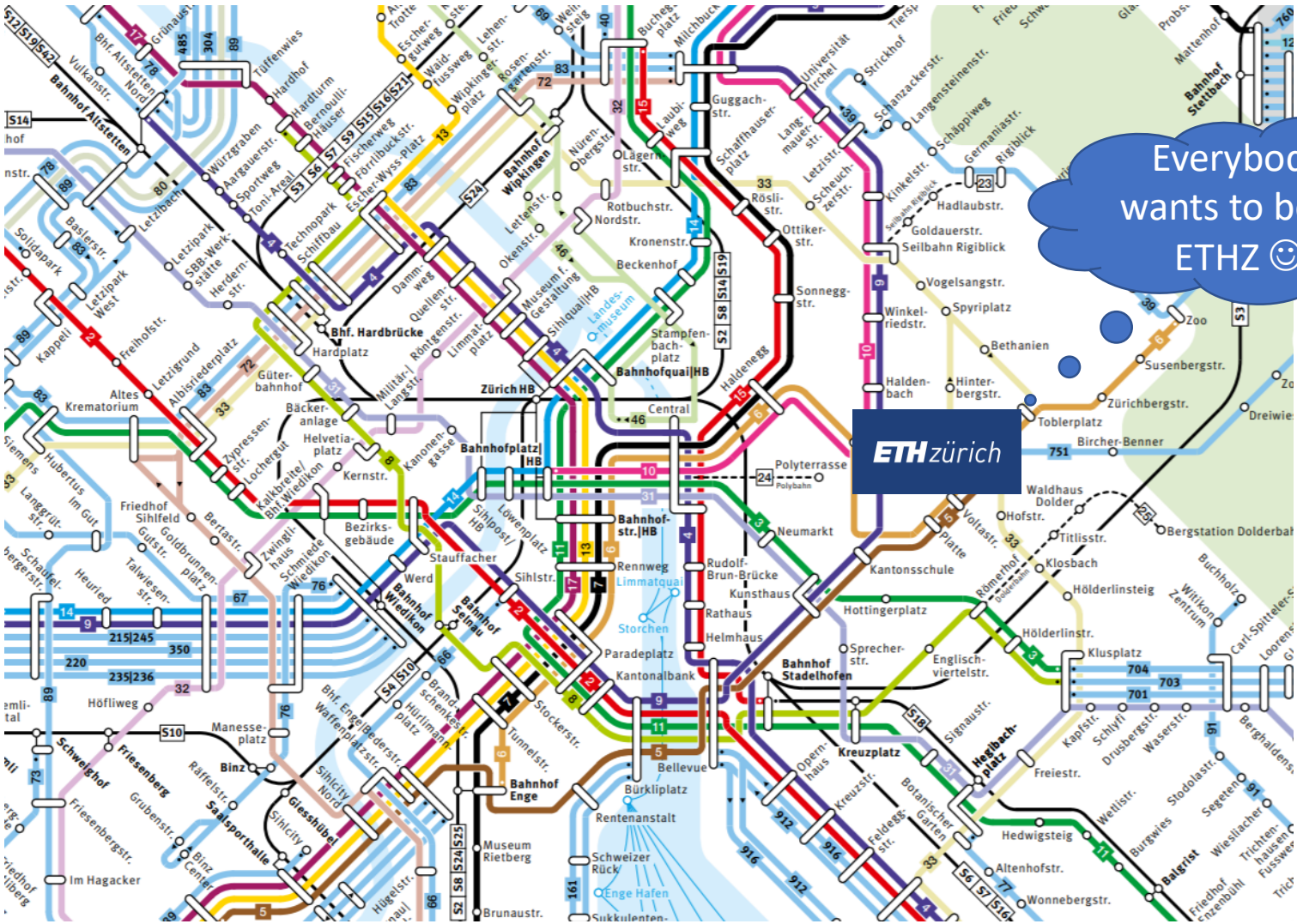
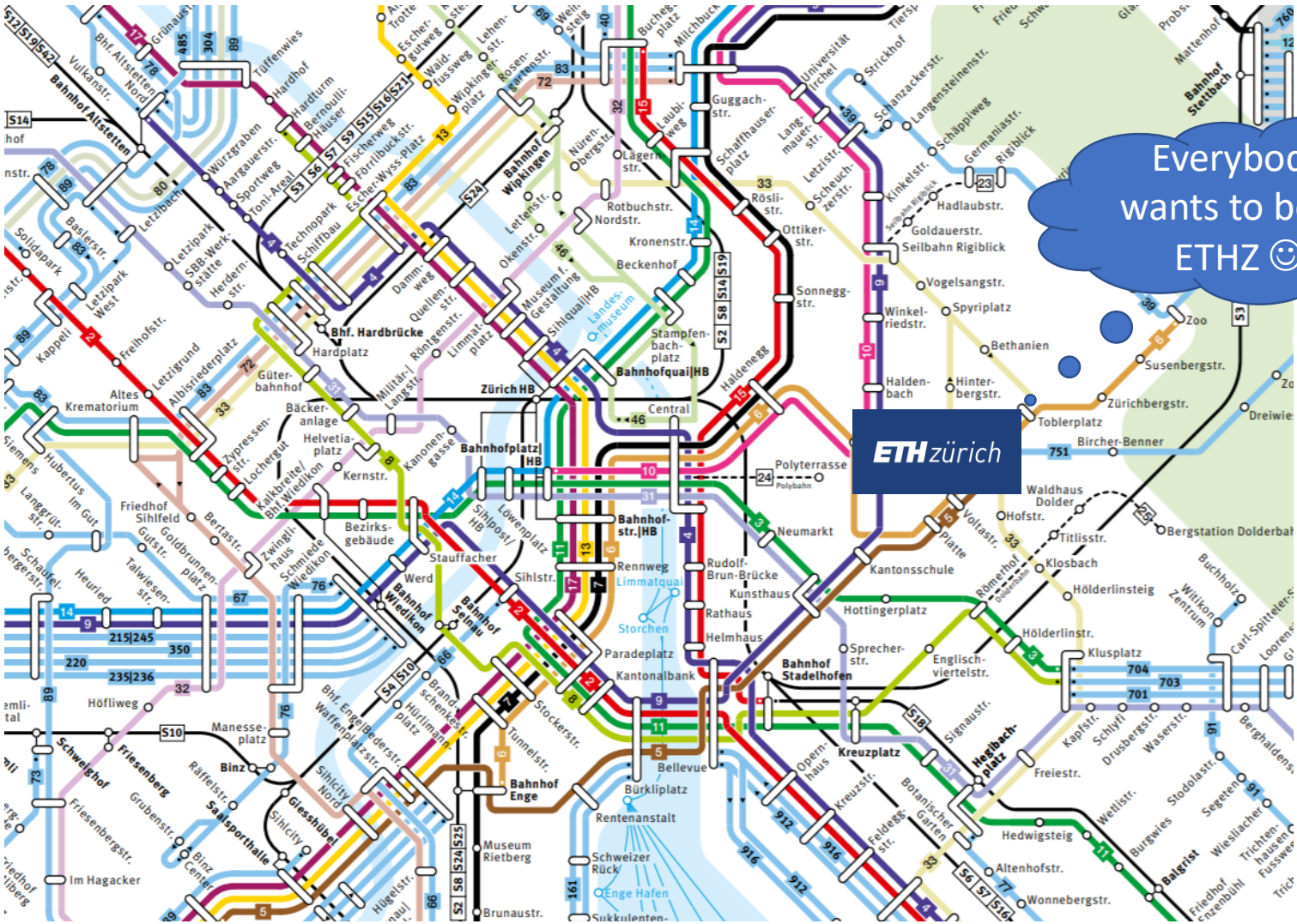[Idea taken from Gilles Tredan]

Everybody wants to be at ETHZ ☺

# What if a link fails?



Everybody wants to be at ETHZ ☺

# What if a link fails? Take a detour ☺

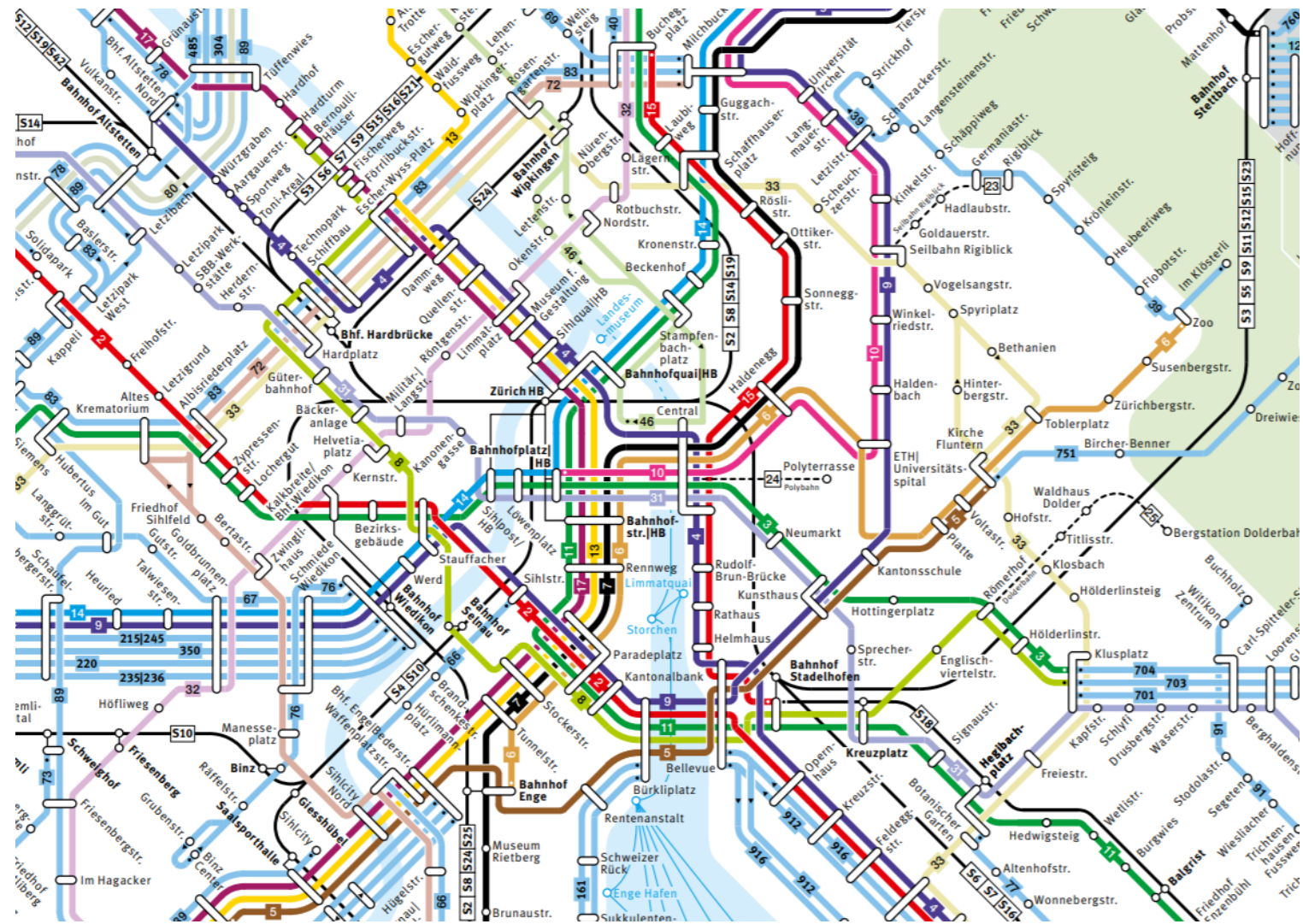

Everybody wants to be at ETHZ ☺

# Everybody takes the same detour? High load!

# Distribute people over all detours? High path stretch!

https://www.elle.com/beauty/health-fitness/news/a35632/why-we-fall-asleep-on-trains/

S12|S19|S42 Grünau 17 485 304 89 Tüffenwies
Bhf. Altstetten Nord 17 78 Hardturm Bernoulli-Häuser
Vulkanstr. Würzgraben Aargauerstr. Sportweg S3 S6 S7 S9 S15 S16 S21 Fischerweg Förrlibuckstr. 13
Bahnhof Altstetten S14 80 Toni-Areal Technopark Schifbau Escher-Wyss-Platz 83
Letzibach Letzipark SBB-Werkstätte Herdern 4 Dammweg Quellenstr. Limmatplatz Museum f. Gestaltung SihlquaiIHB
Baslerstr. 83 89 Letzipark West 4 Müllerstr. Röntgenstr. S2 S8 S14 S19
Solidapark Kappeli 2 Freihofstr. 72 Güterbahnhof Hardplatz Militär-/Langstr. Bahnhofquai|HB Haldenegg
Altes Krematorium 83 33 Bhf. Hardbrücke 31 Bäckeranlage Zürich HB Central
Helvetiaplatz Kanonengasse Bahnhofplatz|HB
Siemens Hubertus Langgrüt Im Gut Friedhof Sihlfeld Kernstr. 8 Bezirksgebäude Polyterrasse 24 ETH|Universitätsspital
Schauflebergerweg Heurled Goldbrunnenplatz Zwinglihaus Löwenplatz 10 31 Bahnhofstr.|HB Neumarkt
14 9 215|245 350 Bahnhof Wiedikon Stauffacher Sihlstr. 17 13 Rennweg Limmatquai Rudolf-Brun-Brücke Kunsthaus
220 235|236 32 76 Bahnhof Selnau Werd 2 Paradeplatz Storchen Rathaus Helmhaus
Höfliweg Manesseplatz S4 S10 66 Brandschenkestr. Kantonalbank Bahnhof Stadelhofen
Schweighof Friesenberg 73 Binz Bhf. Enge Böslederstr. Tunnelstr. 6 Bellevue Opernhaus Kreuzplatz
Giesshübel Sihlcity Nord Bahnhof Enge Bürkliplatz Rentenanstalt 912 916
Saalsporthalle S2 S8 S24 S25 Museum Rietberg Schweizer Rück Enge Hafen 161
Im Hagacker Brunaustr. 5 Sukkulenten-

Universität Irchel Strickhof Langensteinstr. Schäppiweg Rigiblick
Guggachstr. 15 Schafhauserplatz 39 Kinkelstr. Germaniastr.
Rösli-str. 33 Letzistr. Seilbahn Rigiblick Goldauerstr. 23 Hadlaubstr.
Rotbuchstr. Nord Kronenstr. Ottikerstr. Seilbahn Rigiblick Vogelsangstr. Spyristeig
Beckenhof Sonneggstr. Winkelriedstr. 9 Spyriplatz Heuberiweg Im Klösterli
Stampfenbachplatz Winkelriedstr. 10 Bethanien Flobotstr. S3 S5 S9 S11 S12 S15 S23
Haldenbach Hinterbergstr. 6 Susenbergstr. Zoo
Kirche Fluntern Zürichbergstr. 33 Toblerplatz
ETH|Universitätsspital Waldhaus Dolder 751 Bircher-Benner
3 Platte Voltastr. 33 Römerhof Klosbach Bergstation Dolderbahn
Kantonsschule Hofstr. Titlisstr. 25
Hottingerplatz Hölderlinsteig Witikon Zentrum
Sprecherstr. Englischviertelstr. 3 Hölderlinstr. Klusplatz 704 703
Bahnhof Hegibachplatz S18 Signaustr. 91 701
Kreuzplatz Freiestr. Kapfstr. Schlyfi Drusbergstr. Waserstr.
916 Feldeggstr. Botanischer Garten Wetlistr. Stadelastr. Segeten 91
33 Hedwigsteig S6 S7 S16 Wonnebergstr. Altenhofstr. Burgwies Wieslacher Trichtenhausen Balgrist

Bahnhof Stettbach 12
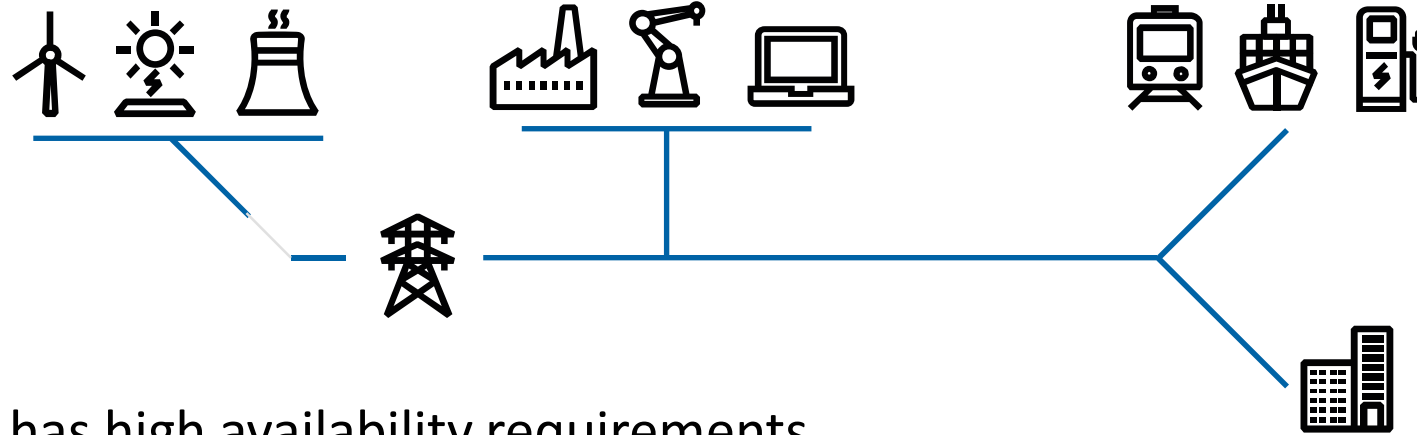Bahnhof Wipkingen S24

Bahnhof Wiedikon Werdstr.

"*The disparity in timescales between packet forwarding (which can be less than a microsecond) and control plane convergence (which can be as high as hundreds of milliseconds) means that failures often lead to unacceptably long outages*"

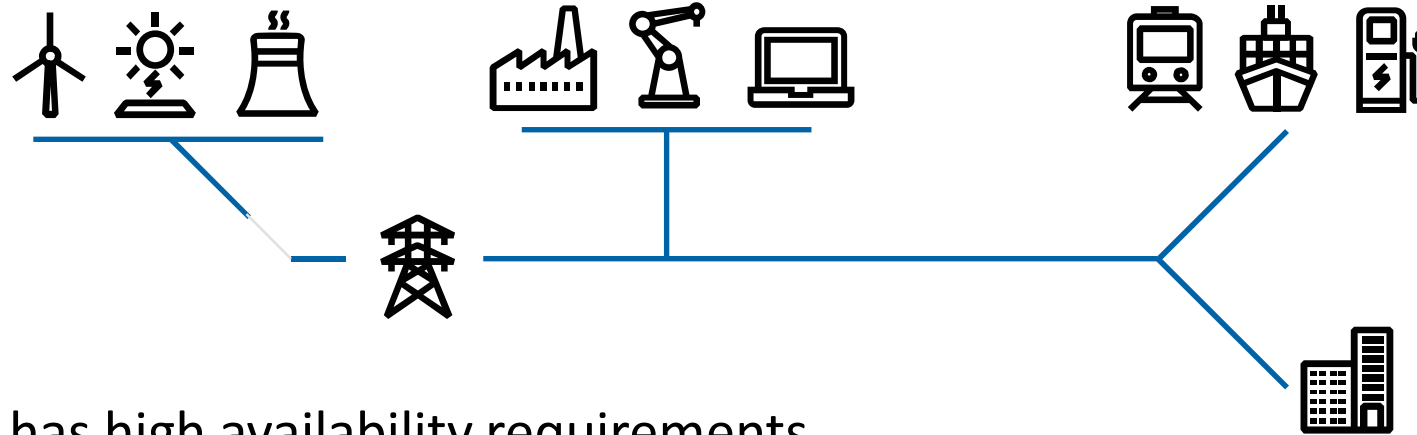Ensuring Connectivity via Data Plane Mechanisms: NSDI'13

# Motivation



- Critical infrastructure has high availability requirements

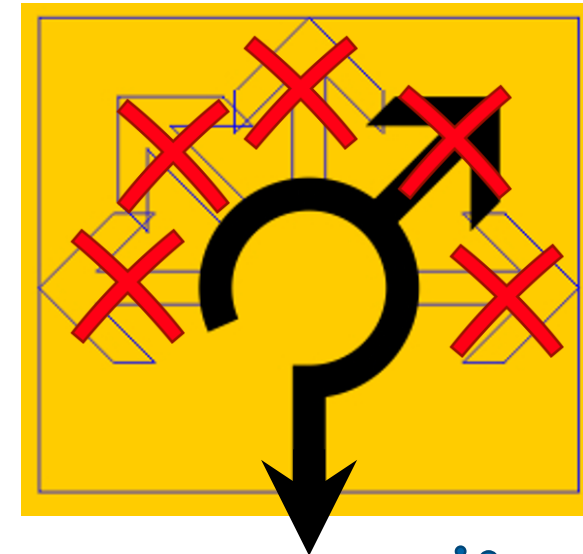- Industrial systems are more and more connected

- Hard real-time requirements

# Motivation



- Critical infrastructure has high availability requirements

- Industrial systems are more and more connected

- Hard real-time requirements

  ⇒ How to provide dependability guarantee despite link failures in networks?
  ⇒ Possible without communication between nodes?
  ⇒ With low load? With low stretch?

# Talk Structure

1. Model and Objectives

2. Background and Lower Bounds

3. Algorithms and Upper Bounds

4. Simulation Results

5. Conclusion and Outlook

# Model I/II: Routing and Network

- Network is a strongly connected directed graph

- Forwarding may only match on:
  1. Source
  2. Destination
  3. Incident failures
  4. Incoming port

- No packet (header) changes allowed, no communication

- Static routing tables, deterministic behaviour

- Single destination routing, uniform flow sizes

Route can be a *walk*

# Model II/II: Quality *from a Worst-Case Perspective*

1. **Resilience**
   - How many link failures can we survive and still guarantee delivery?
   - Upper bound: (*r+1*)-link-connected graph: at most *r*

2. **Load**
   - Maximum additional link utilization due to rerouting

3. **Stretch**
   - Maximum additional hops due to rerouting

# Background: Static Fast Rerouting for Multiple Failures

## Resiliency on General Graphs

- Elhourani *et al.* [ToN'16] / Chiesa *et al.* [INFOCOM'16 etc]:
  - Employ directed link-disjoint arborescences
    - *i.e.* disjoint spanning routing trees
    - after failure: change tree (*e.g.* in circular fashion)
    - incoming port defines current tree



From Chiesa *et al.* 2016

## Resiliency & Load on Complete Graphs

- Borokhovich & Schmid [OPODIS'13]
  - Bounds and handcrafted schemes
- Pignolet *et al.* [DSN'17]
  - Connection to Balanced Incomplete Block Designs (BIBDs)
    - General scheme how to distribute well after failures



From Pignolet *et al.* 2017

**Resiliency & Load on General Graphs**
*this paper*

With improved BIBDs!

# The Price of Locality (for *every* Scheme and Graph)

**Stretch** under *r* failures:

- Adversary can force to visit *r+1* neighbors of destination

  Fail *r* links incident to the destination

**Load** under *r* failures:

- Adversary can force additional load of $\sqrt{r}$

  Previously only weaker bound known, without incoming port

  Let's try to meet this bound for many flows

# CASA: Rerouting on Arborescences

- Takes arborescences as input *e.g.* generated by Chiesa *et al.*
  - Influences the stretch, we get good bounds for *e.g.* so-called *independent spanning trees*

Algorithm

1: *Determine current arborescence T from in-port*

2: *If next hop in T alive, use it, else*

3: *Pick next arborescence T' from* **BIBD-Matrix**

until the next
hop is alive

different flows
use different *T'*

We re-structure BIBD-matrix to be good for many flows

# CASA: Example *without* BIBD

# CASA: Example *without* BIBD



Use same detour 🙁

# CASA: Example *with* BIBD

How much extra load?

- Up to $O(\sqrt{r})$ · · · • Lower bound: $\sqrt{r}$

- For more flows than #arborescences

$$\sqrt{\#failures} < \frac{(\#arborescences)^{\frac{3}{2}}}{\#flows}$$

## *Beyond CASA*

- *r+1* arborescences give *r*-resiliency under directed link failures
  - But unclear how to obtain *r*-resiliency under bi-directed link failures


- Motivation for a simplified heuristic: *SquareOne*
  - Pick *r+1* bi-directed link-disjoint source-destination paths
    - Under failure: bounce back to the source, pick next path



https://Netflix.com

# SquareOne

# SquareOne

How good in practice?

No theoretical guarantees beyond resiliency

Easy to compute via *e.g.* max-flow formulations. Order path priority *e.g.* by length

# Selected Evaluations

Setting from prior work

- 8-connected 8-regular random graphs (**RR**, 100 routers each)

- well-connected cores of real-world ASes (*Rocketfuel*) (204-387 routers, 1667-4736 links)

- Three arborescence methods (using the *same* arborescences)
  - *CASA* **(BIBD)**
  - Deterministic Circular (**DetCirc**) from Chiesa *et al.*
  - Random (**PRNB**) from Chiesa *et al.*

Thanks to Marco Chiesa and Ilya Nikolaevskiy for their support
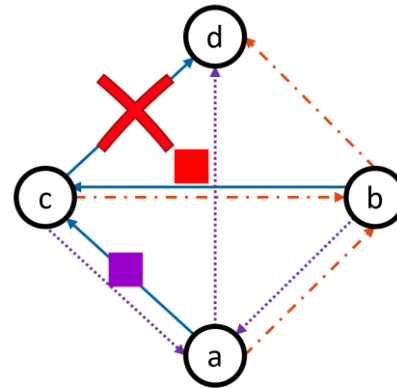
- Also: **SquareOne**

Issues in practice:
Real randomness on routers?
Packet reordering?

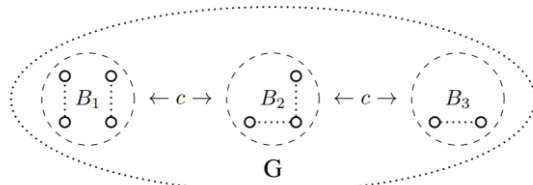# Deterministic Worst-Case Failures

# Conclusion

- We present **efficient static fast failover schemes** on general graphs

  ○ **CASA**: Combines **arborescences** and improved **block-designs** (BIBDs)
    - With theoretical guarantees

  ○ **SquareOne**: Well performing resilient heuristic
    - Based on edge-disjoint paths

- Next slide: Further related problems we work on

# Some More Related Problems

- **Improving arborescence decompositions**
  - #1: Build small stretch arborescences in parallel
    - Current approach: build sequentially in greedy fashion
    - Benefit: Resilient to more failures under nice distributions

  

  - #2: Account for e.g. Shared Risk Link Groups (SRLGs)
    - Leverage post-processing according to objective function
    - Ideally: A SRLG is contained in a single arborescence

  **Appears at #1: DSN 2019, #2: SRDS 2019**

- **Allowing packet header modification (MPLS, SR)**
  - #1: More powerful, but harder to verify correctness?
    - MPLS w. multiple link failures: verification in polynomial time!

|  | P-Rex | NetKAT | HSA | VeriFlow | Anteater |
|---|---|---|---|---|---|
| **Protocol Support** | SR/MPLS | OF | Agn. | OF | Agn. |
| **Approach** | Autom. | Alg. | Geom. | Tries | SAT |
| **Complexity** | Polynom. | PSPACE | Polynom. | NP | NP |
| **Static** | ✓ | ✓ | ✓ | χ | ✓ |
| **Reachability** | ✓ | ✓ | ✓ | ✓ | ✓ |
| **Loop Queries** | ✓ | ✓ | ✓ | ✓ | ✓ |
| **What-if** | ✓ | N/A | ✓ | N/A | χ |
| **Unlim. Header** | ✓ | N/A | χ | χ | N/A |
| **Performance** | ✓ | ✓ [1] | ✓ | ✓ | ✓ |
| **Waypointing** | ✓ | ✓ | ✓ | ✓ | χ |
| **Language** | Py., C | OCaml | Py., C | Py. | C++, Ruby |

  - #2: Leverage Segment Routing (in Linux kernel for IPv6)
    - Allows maximal link protection e.g. in Hypercubes

  **Appears at #1: CoNEXT 2018, #2: OPODIS 2018**

# Papers

- *Improved Fast Rerouting Using Postprocessing*
  Klaus-T. Foerster, Andrzej Kamisinski, Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan. *SRDS 2019*

- *Bonsai: Efficient Fast Failover Routing Using Small Arborescences*
  Klaus-T. Foerster, Andrzej Kamisinski, Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan. *DSN 2019*

- *CASA: Congestion and Stretch Aware Static Fast Rerouting*
  Klaus-T. Foerster, Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan. *INFOCOM 2019*

- *P-Rex: Fast Verification of MPLS Networks with Multiple Link Failures*
  Jesper S. Jensen, Troels B. Krogh, Jonas S. Madsen, S. Schmid, Jiri Srba, and Marc T. Thorgersen. *CoNEXT 2018*

- *Local Fast Segment Rerouting on Hypercubes*
  Klaus-T. Foerster, Mahmoud Parham, Stefan Schmid, and Tao Wen. *OPODIS 2018*

# Congestion and Stretch Aware Static Fast Rerouting [appeared @INFOCOM'19]

Klaus-Tycho Foerster, Yvonne-Anne Pignolet (DFINITY), Stefan Schmid, and Gilles Tredan (LAAS-CNRS)
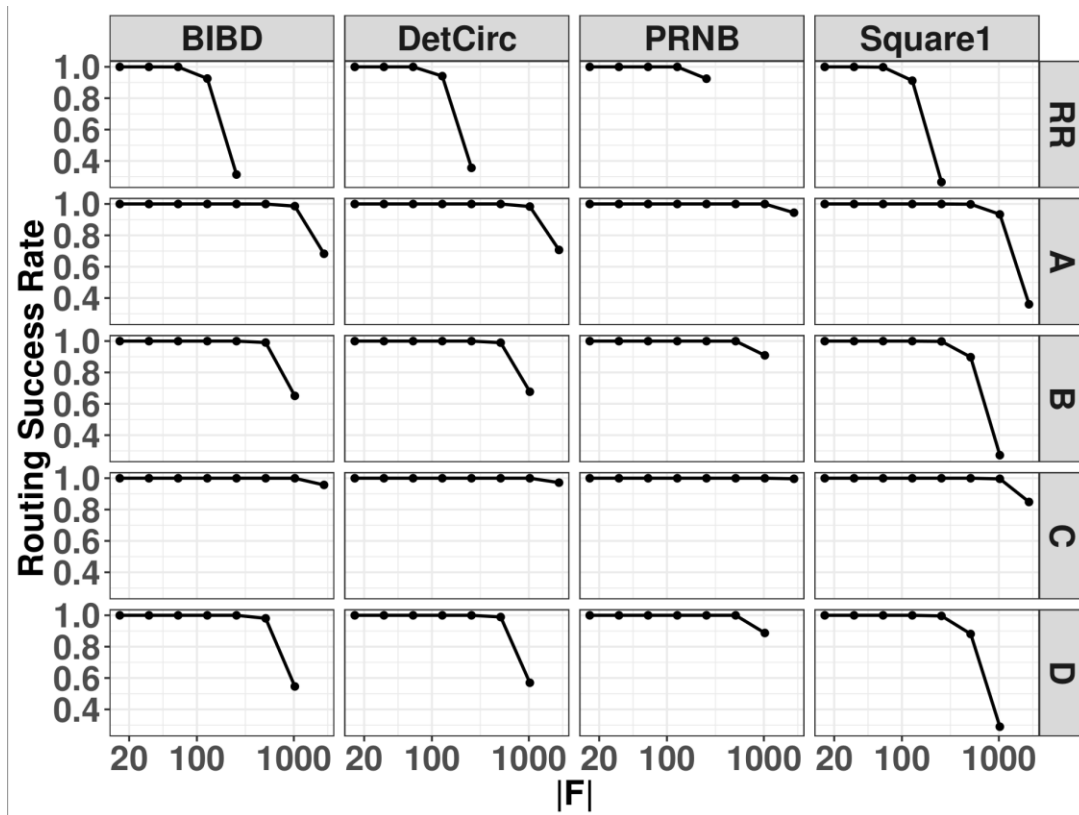
# Papers Referenced

- *How (Not) to Shoot in Your Foot with SDN Local Fast Failover: A Load-Connectivity Tradeoff*
  Michael Borokhovich and Stefan Schmid. *OPODIS 2013*

- *Load-Optimal Local Fast Rerouting for Dependable Networks*
  Yvonne-Anne Pignolet, Stefan Schmid, and Gilles Tredan. *DSN 2013*

- *IP Fast Rerouting for Multi-Link Failures*
  Theodore Elhourani, Abishek Gopalan, Srinivasan Ramasubramanian.
  *IEEE/ACM Trans. Netw.* 24(5): 3014-3025 (2016)

- *The Quest for Resilient (Static) Forwarding Tables*
  Marco Chiesa and  Ilya Nikolaevskiy *et al. INFOCOM 2016*
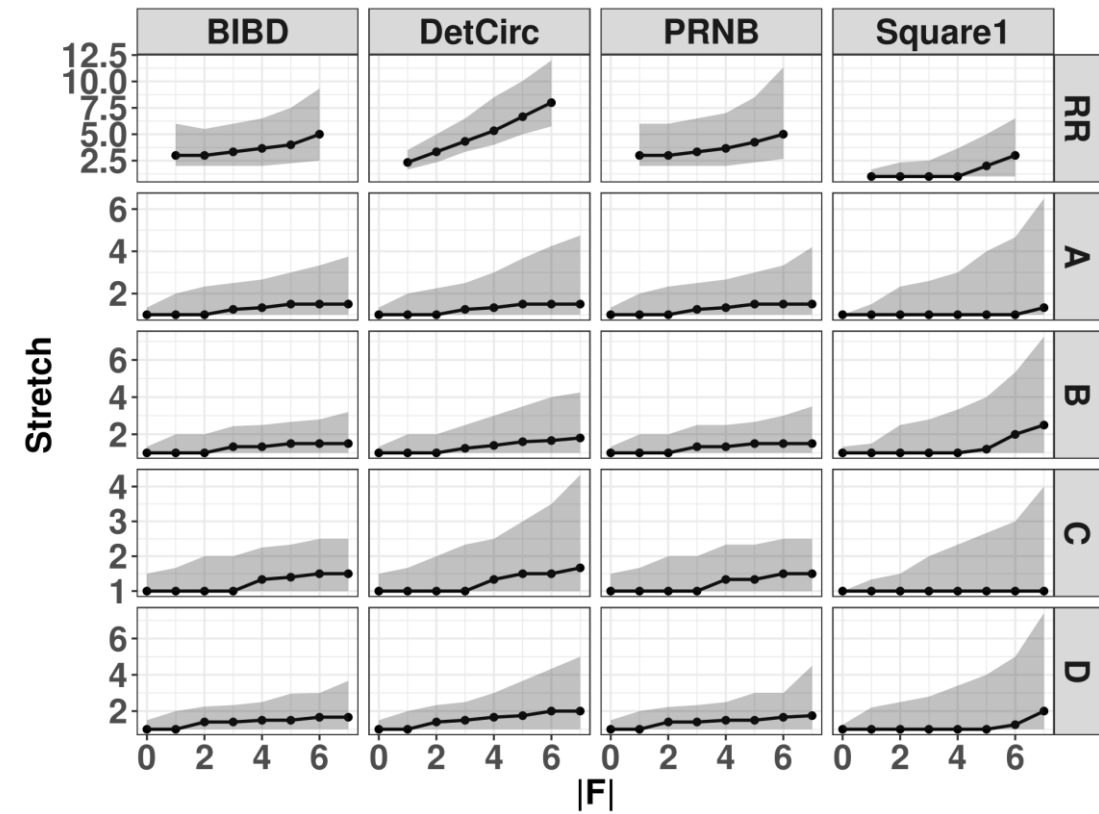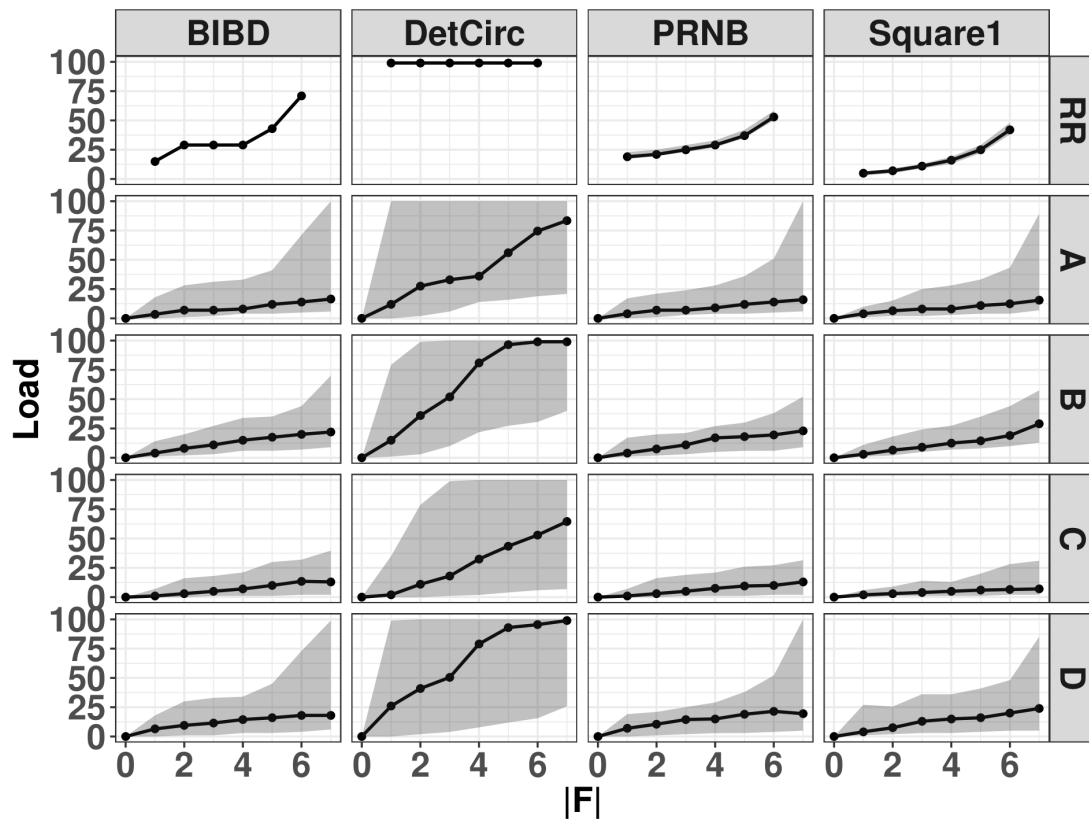
# Rocketfuel ASes

| AS | 1239 A | 2914 B | 3356 C | 7018 D |
|---|---|---|---|---|
| Number of nodes | 389 | 225 | 377 | 204 |
| Number of links | 3621 | 1696 | 4736 | 1667 |
| Eccentricity | 6 | 6 | 6 | 6 |
| Avg shortest path length | 3.06 | 2.48 | 3.14 | 3.17 |

TABLE I: Properties of 8-connected cores of various ASes

# Evaluation: Resiliency

# Evaluation: Deterministic Worst-Case Failures

# Evaluation: Random Failures