

Supporting Bottleneck Structure Graphs: Use Cases and Requirements

Jordi Ros-Giralt, Sruthi Yellamraju, Qin Wu, Richard Yang, Luis Contreras, Kai Gao, Jensen Zhang

I-Draft: draft-giraltyellamraju-alto-bsg-multidomain

<https://datatracker.ietf.org/doc/draft-giraltyellamraju-alto-bsg-requirements/>

IETF Plenary 115

PANRG Session

11/10/2022

Table of Contents

- Brief Introduction to Bottleneck Structures
- Bottleneck Structure Graphs (BSGs): Use Cases
- Production Deployments
- Discussion Q&A

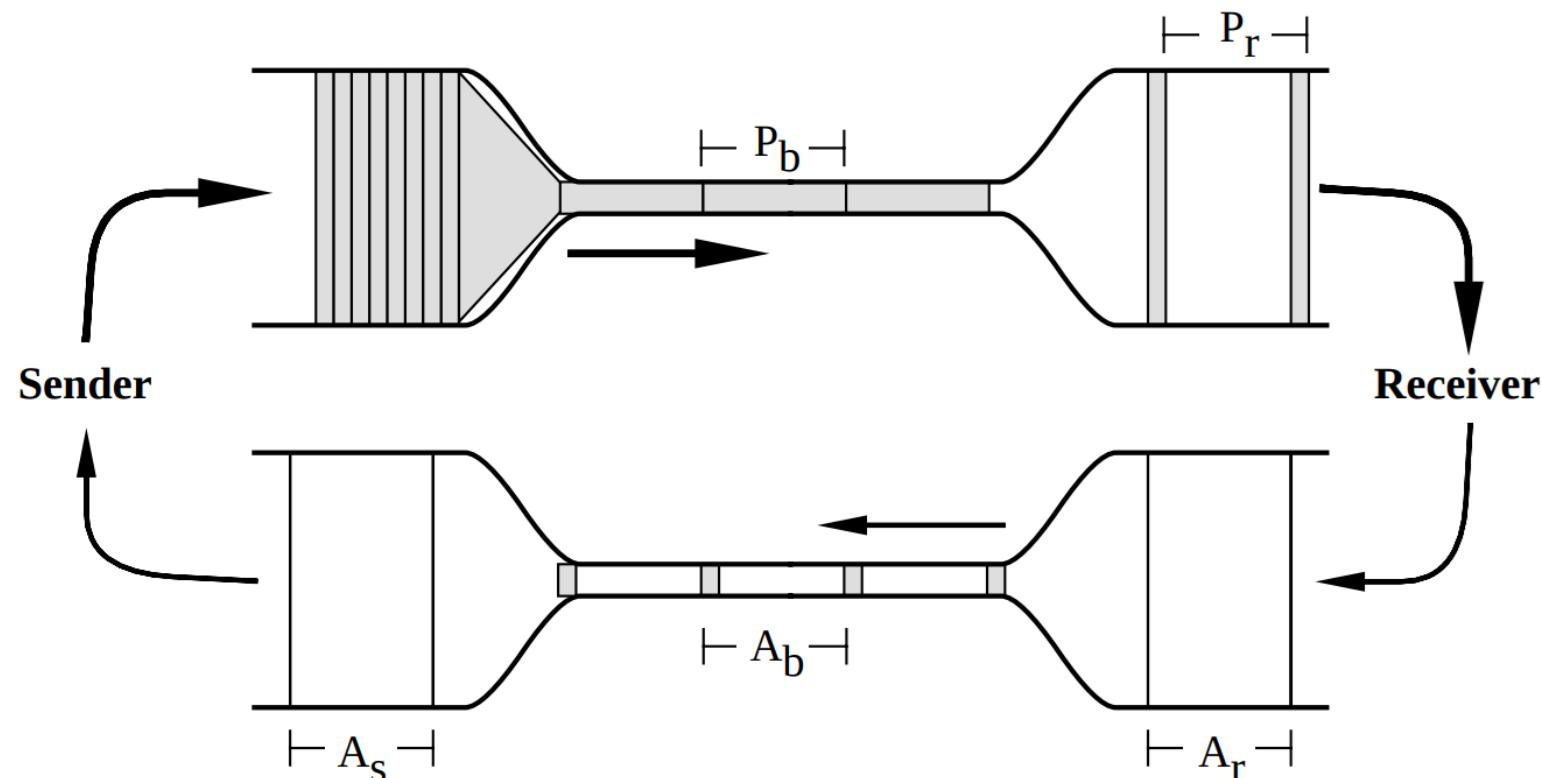
Introduction to Bottleneck Structures

Framework and Implementation Details in the Following I-Drafts and Papers

- [1] IETF Draft: "Supporting Bottleneck Structure Graphs in ALTO: Use Cases and Requirements", <https://datatracker.ietf.org/doc/draft-giraltyellamraju-alto-bsg-requirements/>
- [2] IETF Draft: "Bottleneck Structure Graphs in Multidomain Networks: Introduction and Requirements for ALTO", <https://datatracker.ietf.org/doc/draft-giraltyellamraju-alto-bsg-multidomain/>
- [3] "On the Bottleneck Structure of Congestion-Controlled Networks," ACM SIGMETRICS, Boston, June 2020 [<https://bit.ly/3Urng9M>].
- [4] "Designing Data Center Networks Using Bottleneck Structures," accepted for publication at ACM SIGCOMM 2021 [<https://bit.ly/3TaJpZ5>].
- [5] "A Quantitative Theory of Bottleneck Structures for Data Networks", Qualcomm Technologies, Inc. Technical Report, 2022 [<https://bit.ly/3DG4u7U>].

Conventional View: Single Bottleneck Model

Figure 1: Window Flow Control 'Self-clocking'

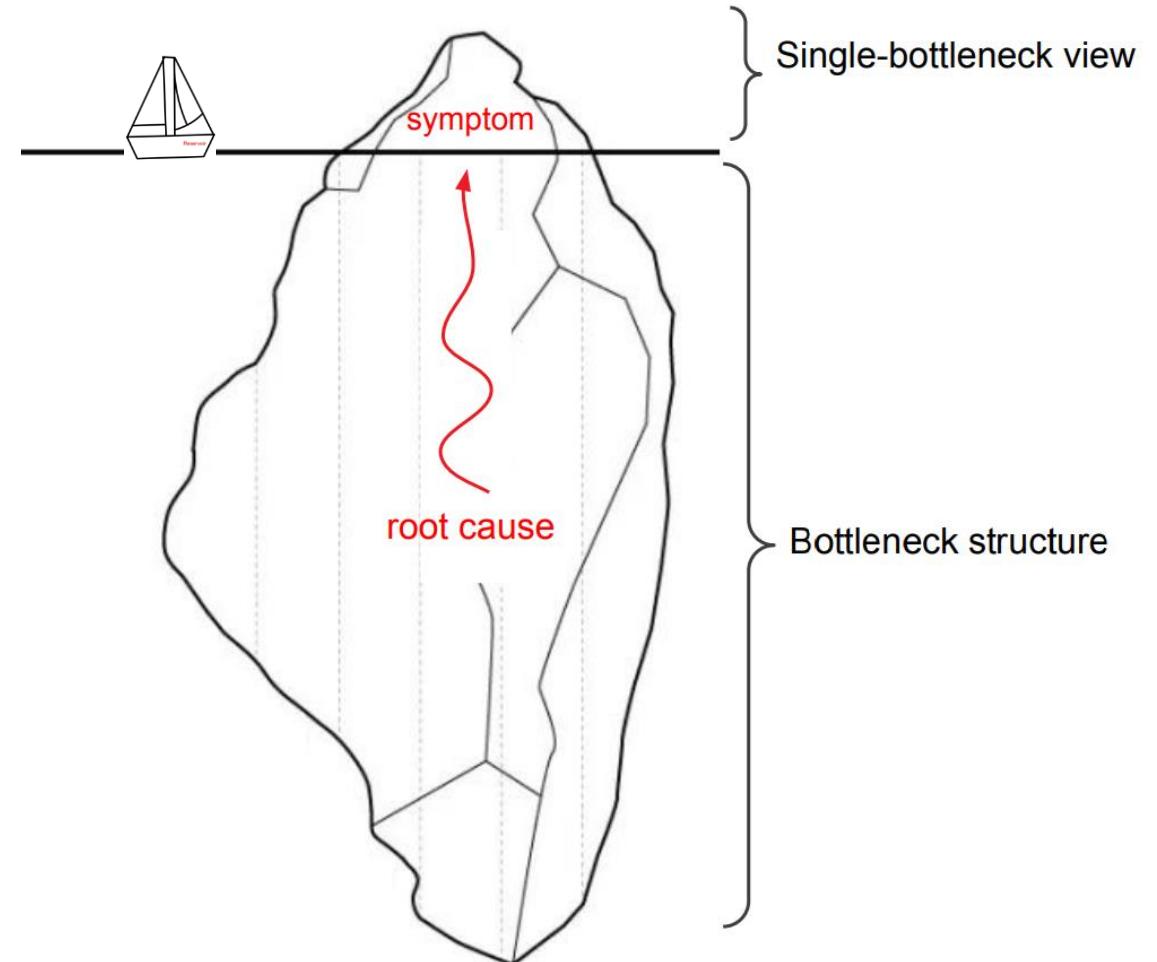


[*] Van Jacobson, "Congestion Avoidance and Control," SIGCOMM, 1988 [<https://bit.ly/3FQouFf>]

Problem Positioning: The Hiding Root Cause of System-Wide Performance

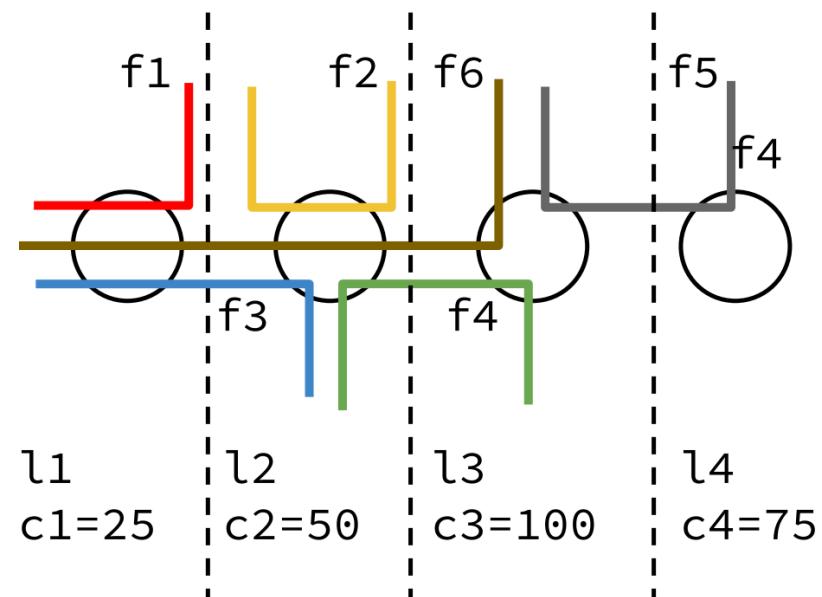
Analogy:

- Structure of the congestion problem in data networks:
 - The single-bottleneck problem is the tip of the iceberg (the symptom)
 - The bottleneck structure is the submerged portion (determines system-wide performance)



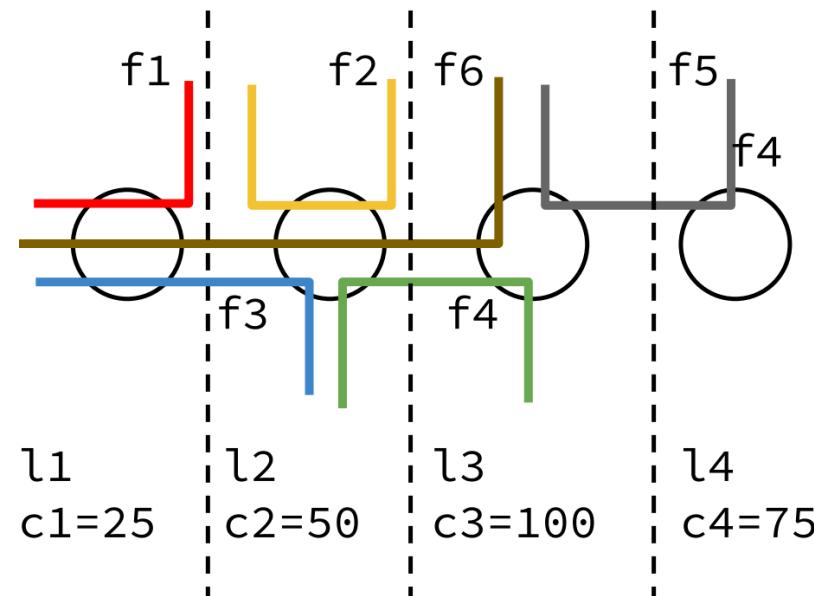
Simple Example of Bottleneck Structure

Communication Network:

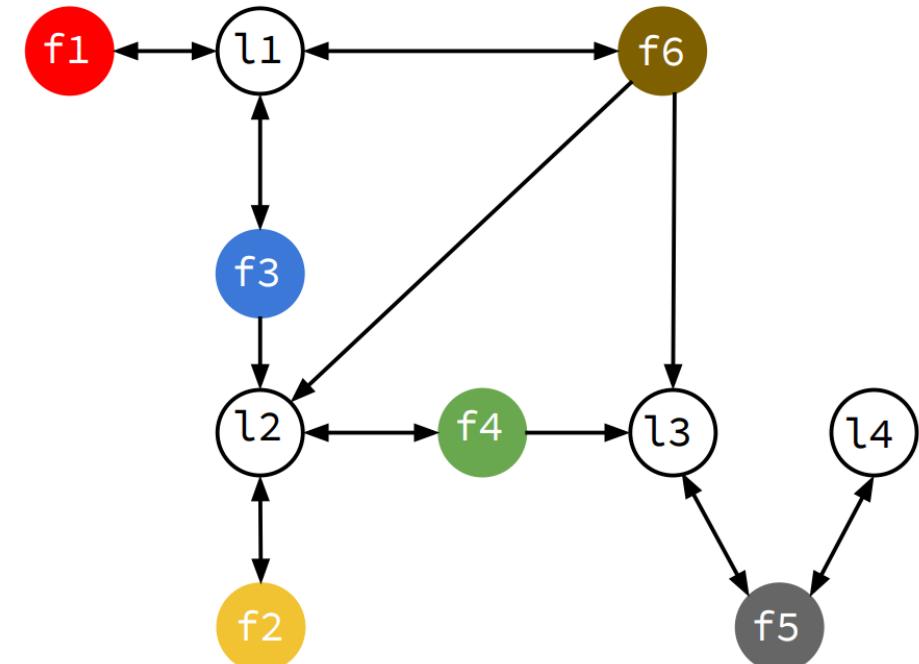


Simple Example of Bottleneck Structure

Communication Network:

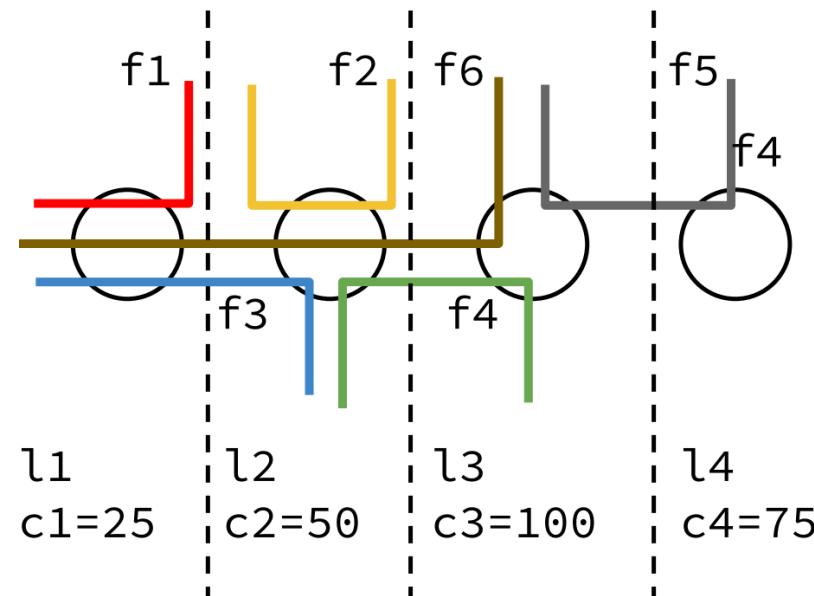


Bottleneck Structure:

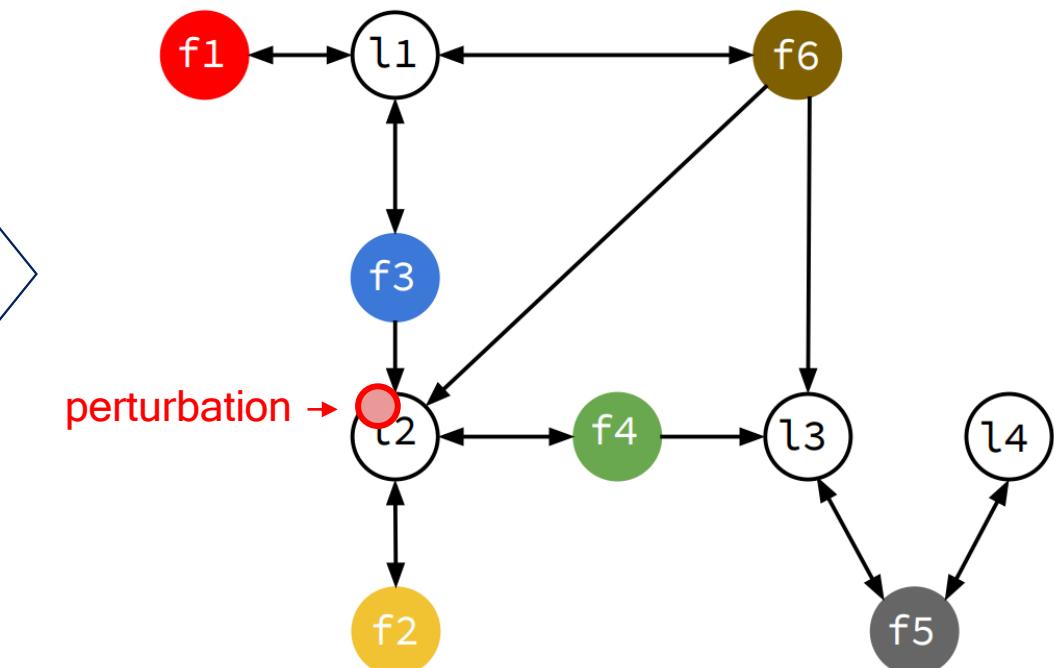


Simple Example of Bottleneck Structure

Communication Network:

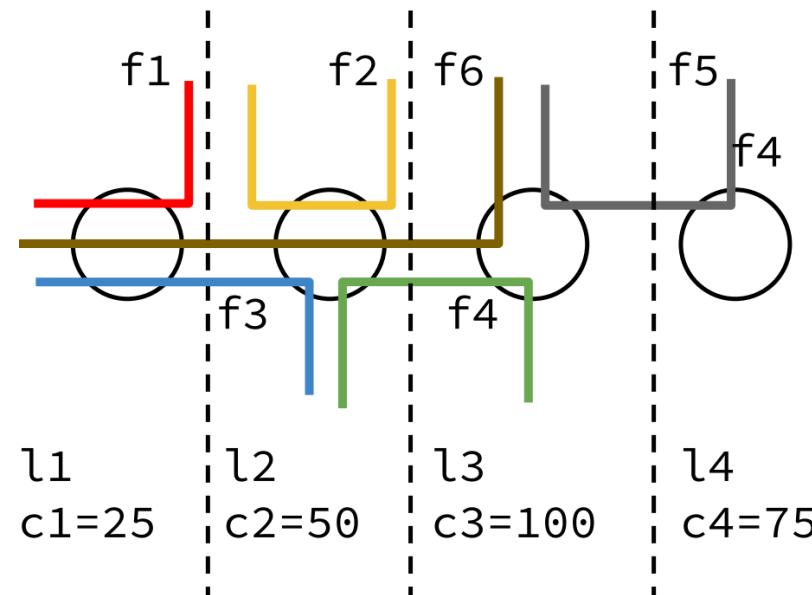


Bottleneck Structure:

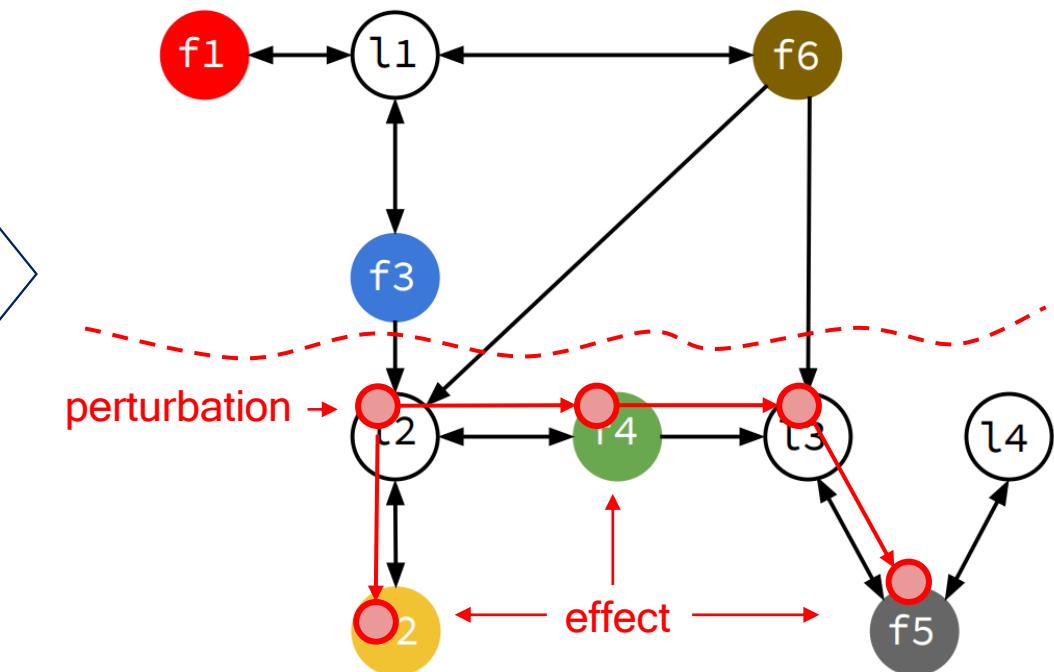


Simple Example of Bottleneck Structure

Communication Network:

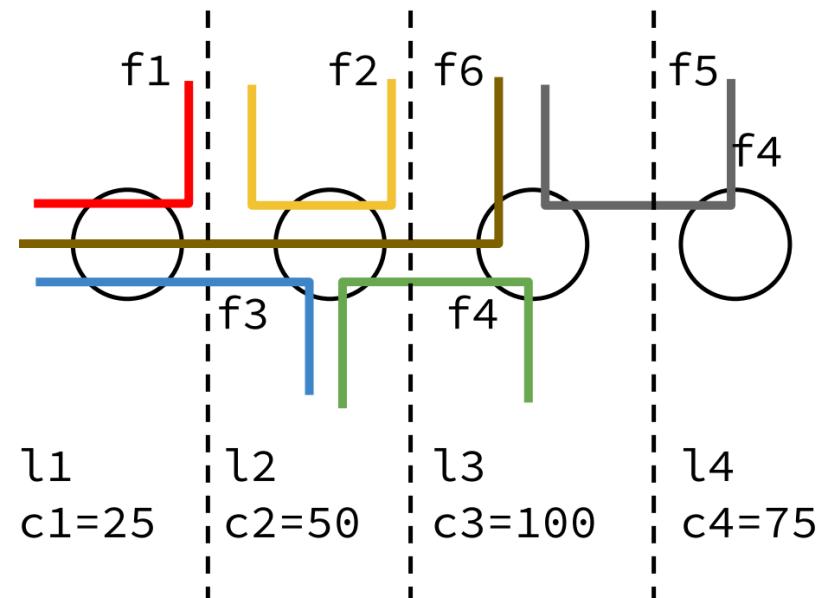


Bottleneck Structure:

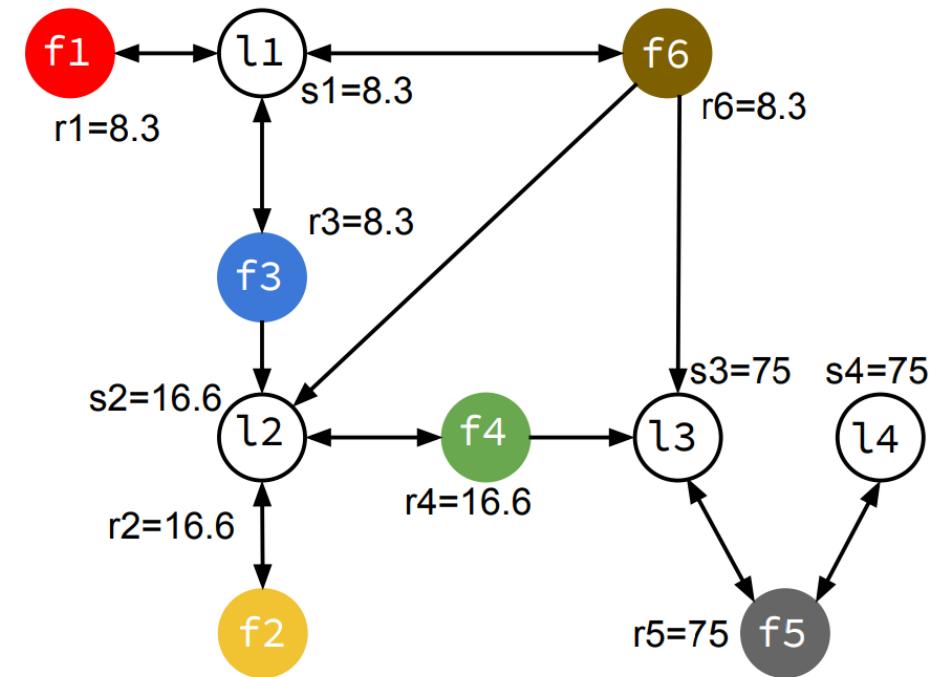


Simple Example of Bottleneck Structure

Communication Network:



Bottleneck Structure:



Flow bandwidth allocation: $\mathbf{r} = [8.3, 16.6, 8.3, 16.6, 75, 8.3]$

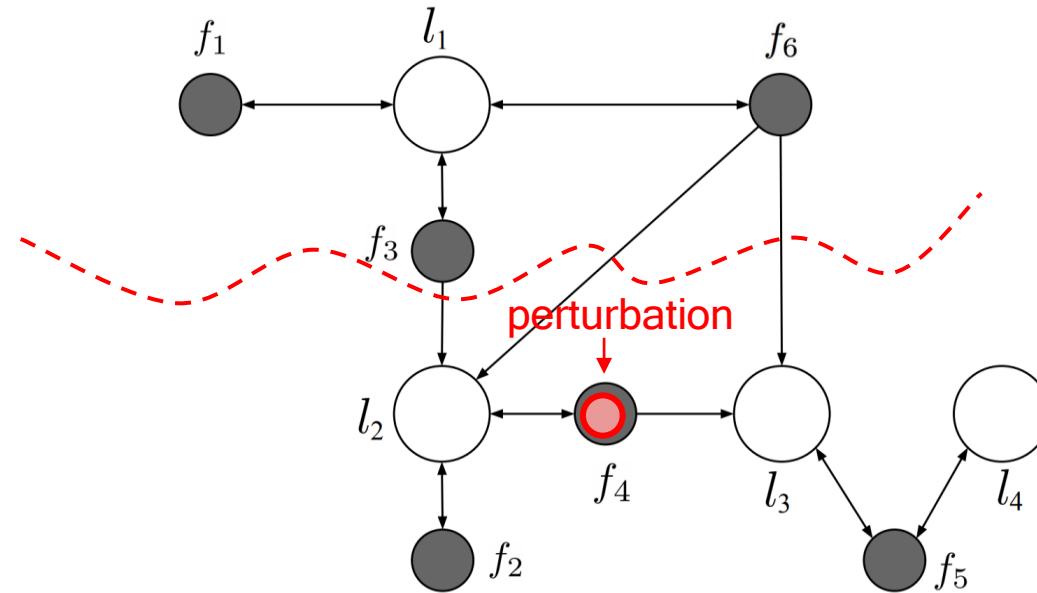
Propagation Lemmas

Can the flap of a butterfly's wings in America set off a tornado in Asia?



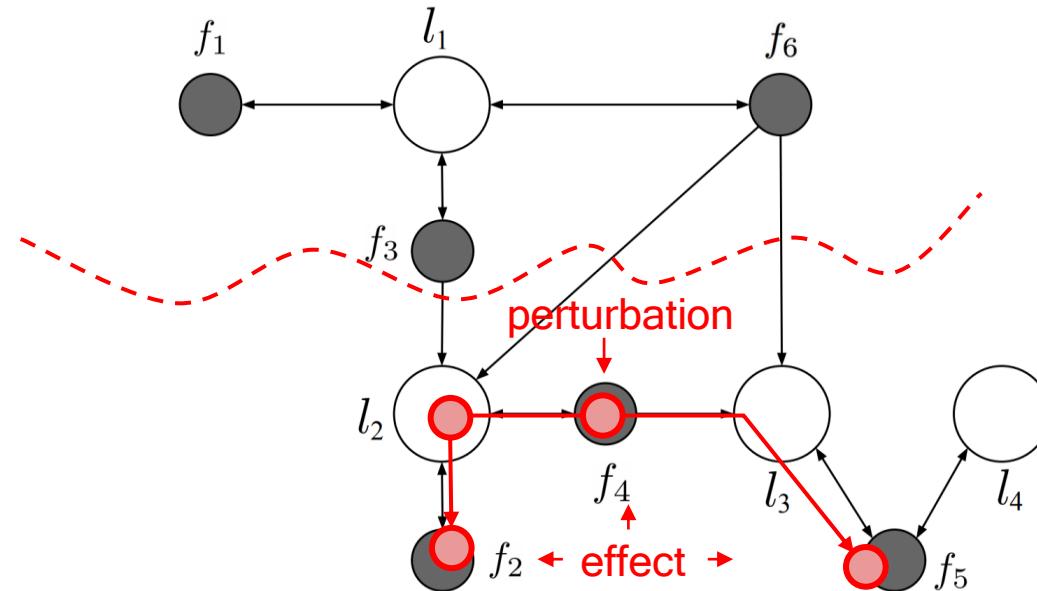
Propagation Lemmas

A flow f can influence the performance of another flow f' iff the bottleneck structure has a directed path from f to f' . (Same lemma exists for bottleneck links.)



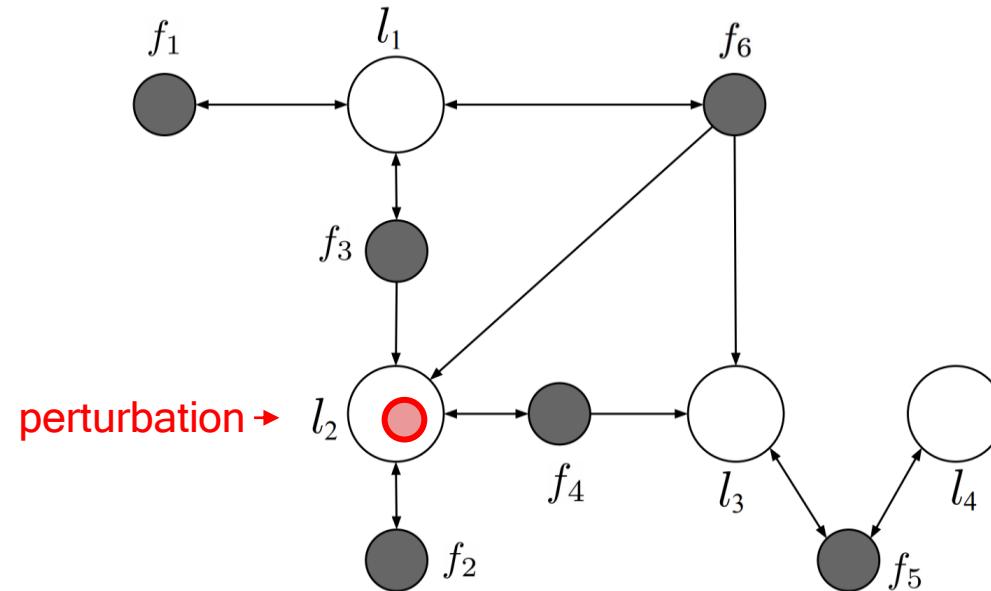
Propagation Lemmas

A flow f can influence the performance of another flow f' iff the bottleneck structure has a directed path from f to f' . (Same lemma exists for bottleneck links.)



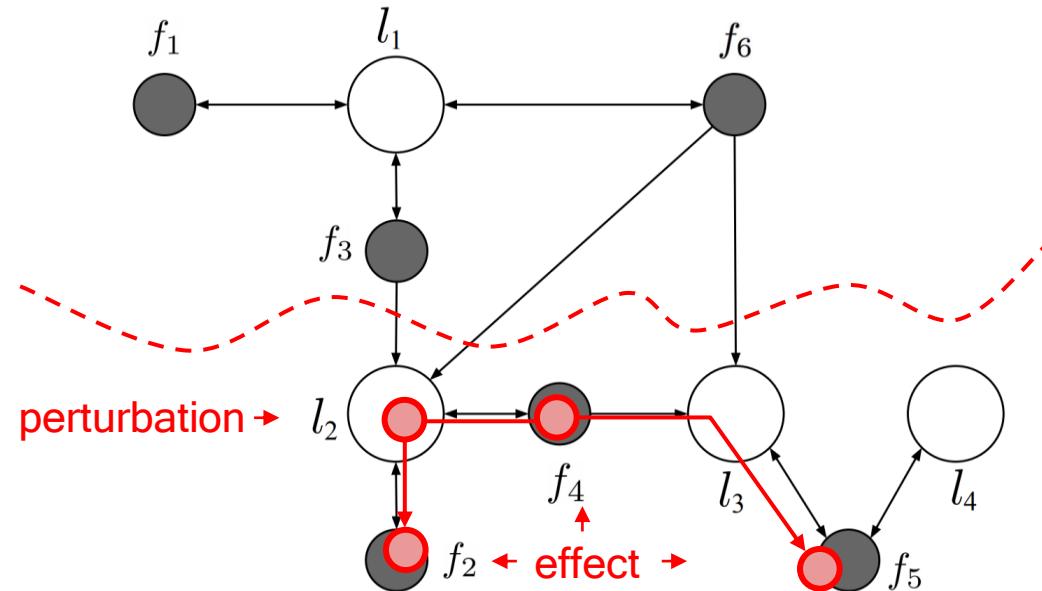
Propagation Lemmas

A flow f can influence the performance of another flow f' iff the bottleneck structure has a directed path from f to f' . (Same lemma exists for bottleneck links.)



Propagation Lemmas

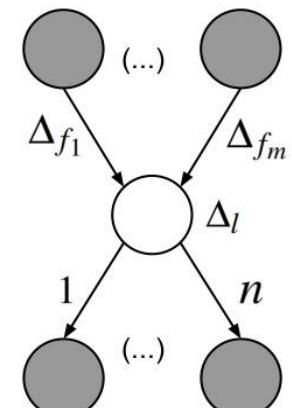
A flow f can influence the performance of another flow f' iff the bottleneck structure has a directed path from f to f' . (Same lemma exists for bottleneck links.)



Propagation Lemmas

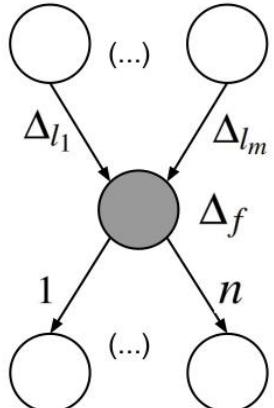
Link and flow equations:

(a) Link equation:



$$\Delta_l = - \sum_{1 \leq i \leq m} \Delta_{f_i} / n$$

(b) Flow equation:

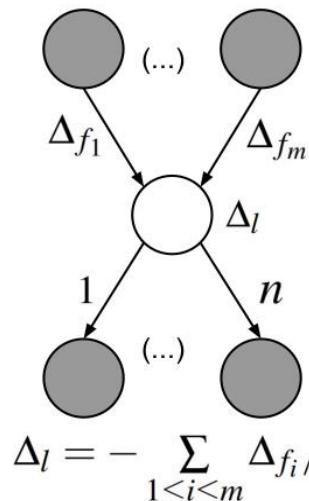


$$\Delta_f = \min\{\Delta_{l_i}, 1 \leq i \leq m\}$$

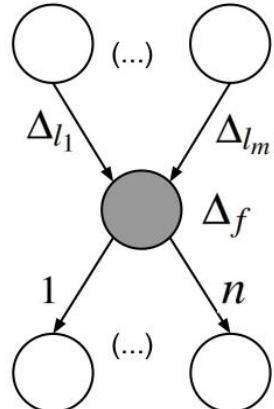
Propagation Lemmas

Link and flow equations:

(a) Link equation:

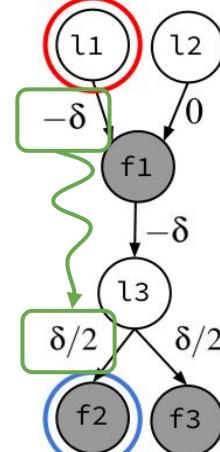


(b) Flow equation:

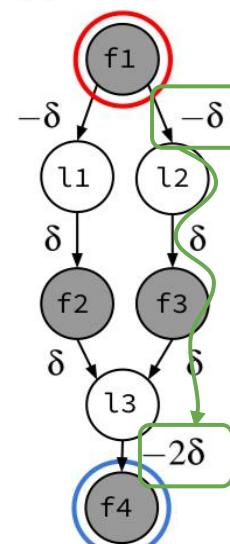


Example:

(c) Link gradient:



(d) Flow gradient:

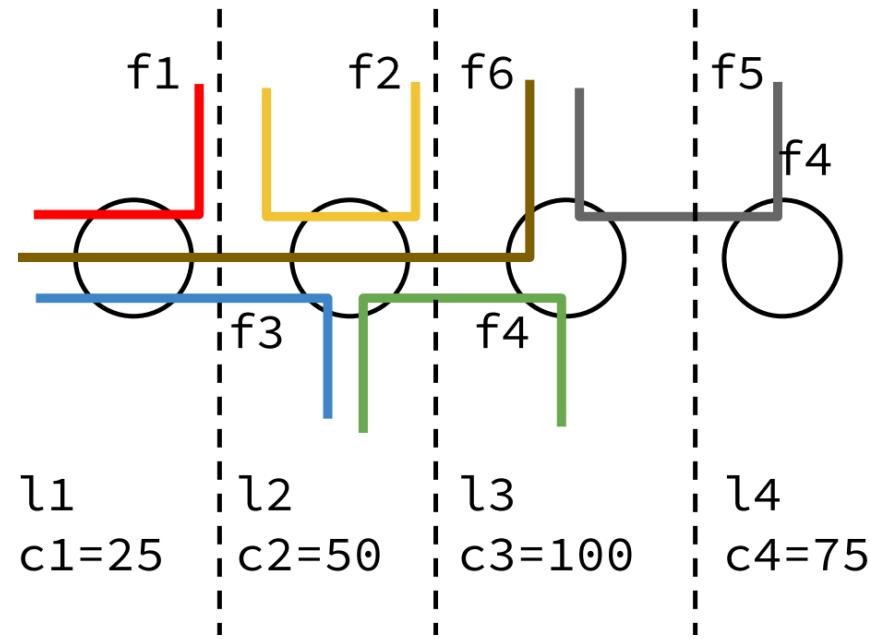


$\nabla_{l_1}(f_2) = \partial r_{f_2} / \partial c_{l_1} = \Delta_{f_2} / (-\delta) = \frac{\delta/2}{-\delta} = -1/2$

$\nabla_{f_1}(f_4) = \partial r_{f_4} / \partial r_{f_1} = \Delta_{f_4} / (-\delta) = -2\delta / (-\delta) = 2$

Optimal Flow Throughput Reduction

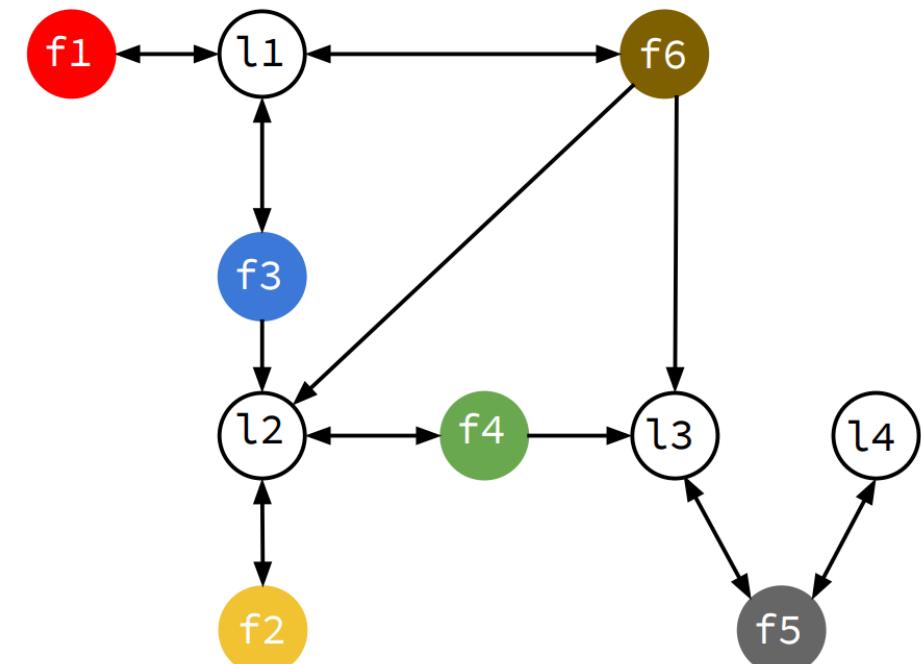
Communication Network:



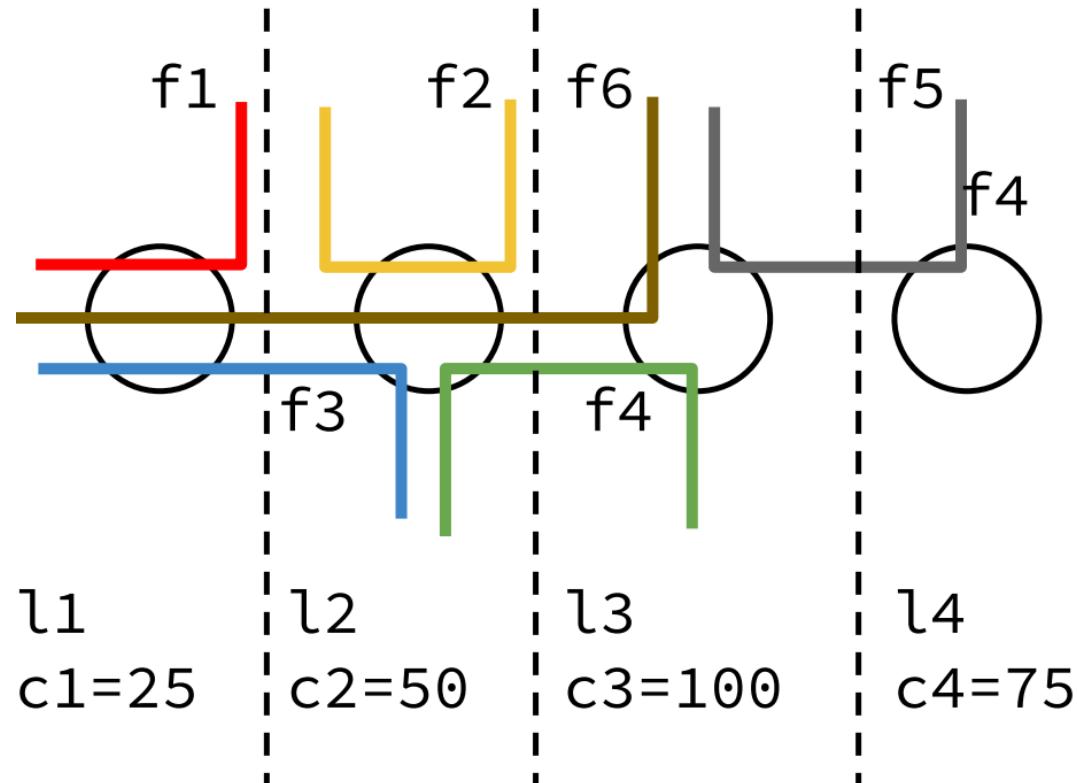
$$\mathbf{r} = [8.\overline{3}, 16.\overline{6}, 8.\overline{3}, 16.\overline{6}, 75, 8.\overline{3}]$$

$f_1 \quad f_2 \quad f_3 \quad f_4 \quad f_5 \quad f_6$

Bottleneck Structure:



Optimal Flow Throughput Reduction



$$f_1 \quad f_2 \quad f_3 \quad f_4 \quad f_5 \quad f_6$$

$$\mathbf{r} = [8.3, 16.6, 8.3, 16.6, 75, 8.3]$$

F : Total network flow

$$\partial F / \partial r_1^- = 1$$

$$\partial F / \partial r_2^- = 1$$

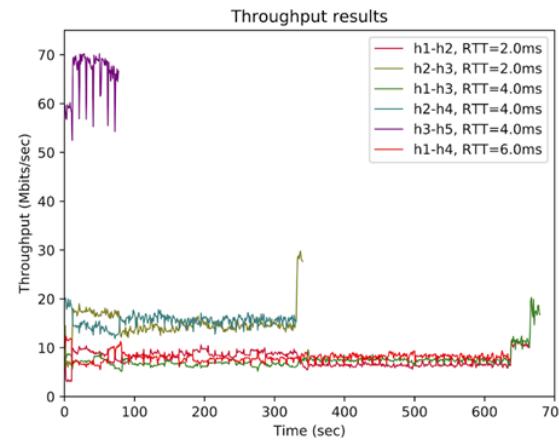
$$\partial F / \partial r_3^- = 1/4$$

$$\partial F / \partial r_4^- = 0$$

$$\partial F / \partial r_5^- = 1$$

$$\boxed{\partial F / \partial r_6^- = -1/2}$$

Optimal Flow Throughput Reduction



(a) Without removing any flow.

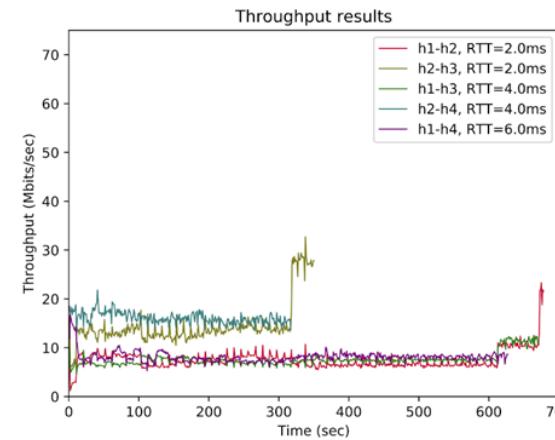
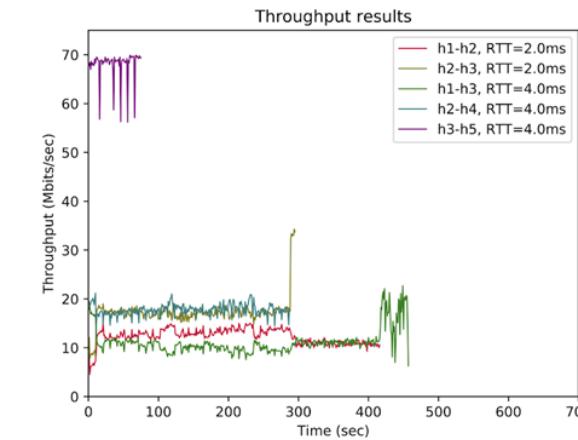
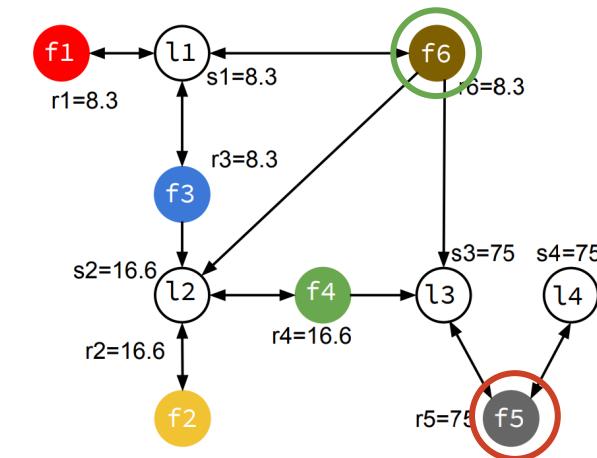
(b) Removing the heavy-hitter flow f_5 .(c) Removing a low-hitter flow f_6 .

Table 3: As predicted by the theory of bottleneck ordering,
flow f_6 is a significantly higher impact flow than flow f_5 .

Comp. time (secs)	f_1	f_2	f_3	f_4	f_5	f_6	Slowest
With all flows	664	340	679	331	77	636	679
Without flow f_5	678	350	671	317	—	611	678
Without flow f_6	416	295	457	288	75	—	457
Avg rate (Mbps)	f_1	f_2	f_3	f_4	f_5	f_6	Total
With all flows	7.7	15.1	7.5	15.4	65.8	8.1	119.6
Without flow f_5	7.5	14.5	7.6	16.1	—	8.3	54
Without flow f_6	12.2	17.2	11.1	17.7	68.1	—	126.3

Bottleneck Structure:

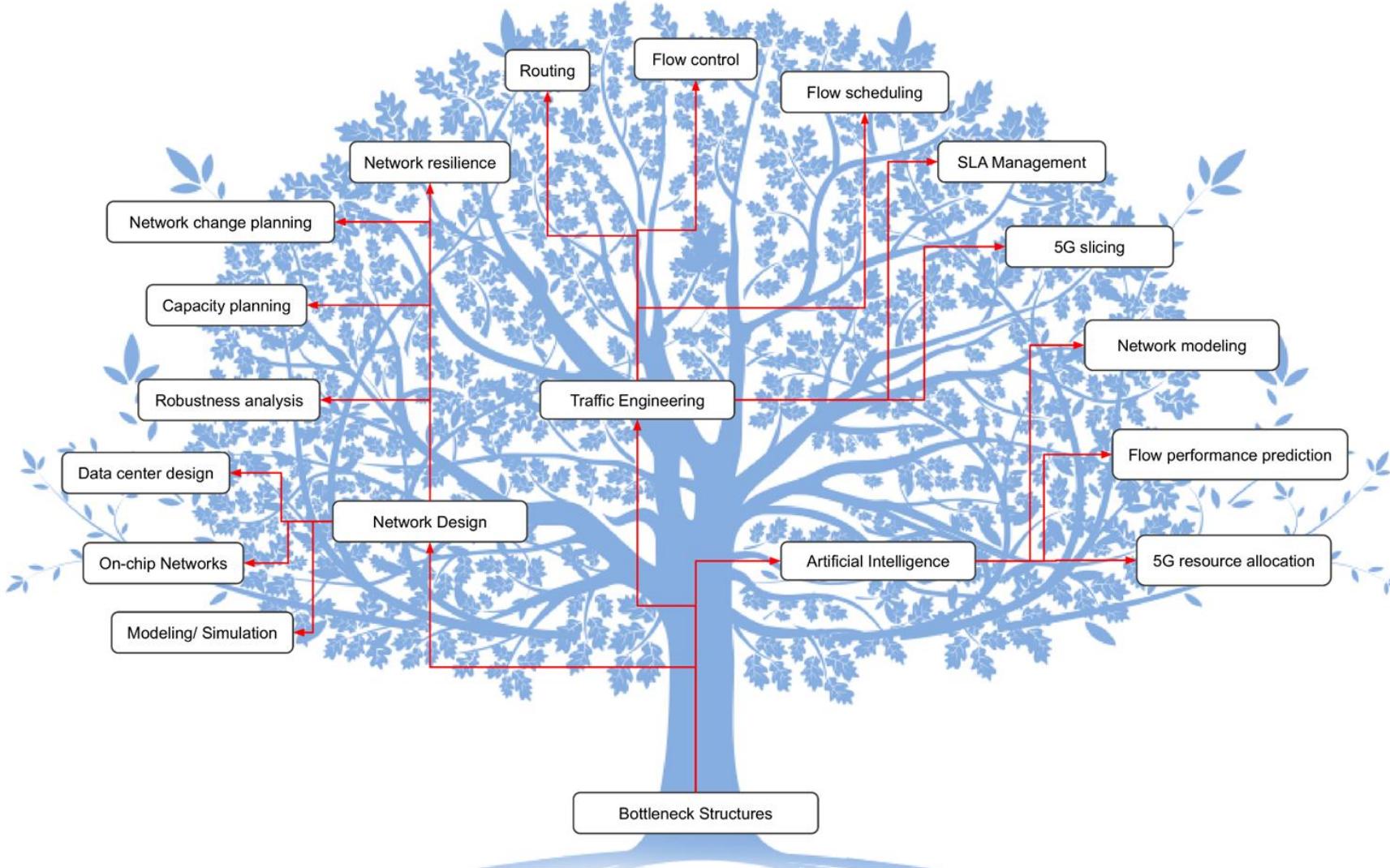


Types of Perturbations (Derivatives) Supported by the Bottleneck Structure Graph

- Flow routing
- Traffic shaping (BW enforcement)
- Link capacity upgrades
- Link capacity fluctuations (e.g., SNR in a wireless channel)
- Path shortcuts
- Flow scheduling
- Flow completion
- Job mapping
- Multi-job scheduling

Bottleneck Structure Graphs (BSGs): Use Cases

Bottleneck Structure Graphs (BSG): Use Cases



Potential WGs collaborations with ALTO

PANRG

PCE

TEAS

CDNI

COINRG

NETMOD

DETNET

NMRG / digital twins

CAN (BOF)

IAB / Path Signals

Others...

Bottleneck Structure Graphs (BSG): Use Cases Documented in the I-Draft

- Application Rate Limiting for CDN and Edge Cloud Applications
- Time-bound Constrained Flow Acceleration for Large Data Set Transfers
- Application Performance Optimization Through AI Modeling
- Optimized Joint Routing and Congestion Control
- Service Placement for Edge Computing
- Training Neural Networks and AI Inference for Edge Clouds, Data Centers and Planet-Scale Networks
- 5G Network Slicing

Bottleneck Structure Graphs (BSG): Use Cases Documented in the I-Draft

- Application Rate Limiting for CDN and Edge Cloud Applications
- Time-bound Constrained Flow Acceleration for Large Data Set Transfers
- Application Performance Optimization Through AI Modeling
- Optimized Joint Routing and Congestion Control
- Service Placement for Edge Computing
- Training Neural Networks and AI Inference for Edge Clouds, Data Centers and Planet-Scale Networks
- 5G Network Slicing

We will focus on “Optimized Joint Routing and Congestion Control”. For details on all other use cases, see the I-Draft:

<https://datatracker.ietf.org/doc/draft-giraltyellamraju-alto-bsg-requirements/>

BSG Use Cases: Optimizing Joint Routing and Congestion Control

Assume Google's B4 Network from [B4-SIGCOM]:

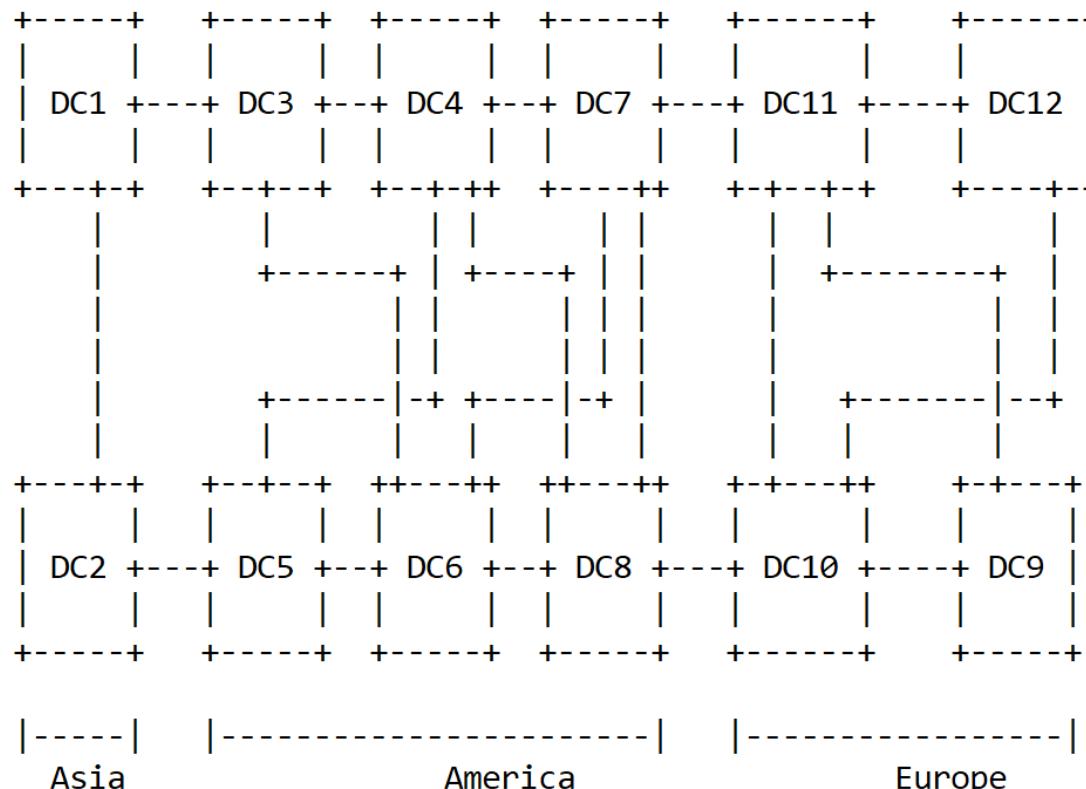


Figure 4: Google's B4 network introduced in [B4-SIGCOMM].

	Link	Adjacent data centers	Link	Adjacent data centers
11	DC1, DC2		111	DC10, DC12
12	DC1, DC3		112	DC4, DC5
13	DC3, DC4		113	DC5, DC6
14	DC2, DC5		114	DC11, DC12
15	DC3, DC6		115	DC4, DC7
16	DC6, DC7		116	DC4, DC8
17	DC7, DC8		117	DC7, DC8
18	DC8, DC10		118	DC9, DC11
19	DC9, DC10		119	DC10, DC11
110	DC7, DC11			

Table 1: Link connectivity (adjacency matrix) in the B4 network.

BSG Use Cases: Optimizing Joint Routing and Congestion Control

Assume Google's B4 Network from [B4-SIGCOM] (a bit more human friendly view):



BSG Use Cases: Optimizing Joint Routing and Congestion Control

- Assume a simple configuration with a pair of flows (one for each direction) connecting every data center in the US with every data center in Europe.
- All links are assumed to have a capacity of 10 Gbps except for the transatlantic links (DC7-DC11 and DC8-DC10), which are configured at 25 Gbps.
- Then the bottleneck structure is the graph shown on the right.

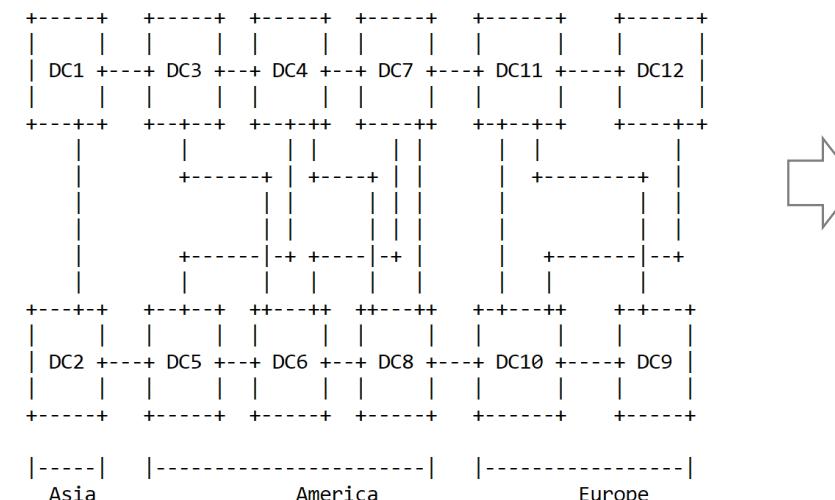


Figure 4: Google's B4 network introduced in [B4-SIGCOMM].

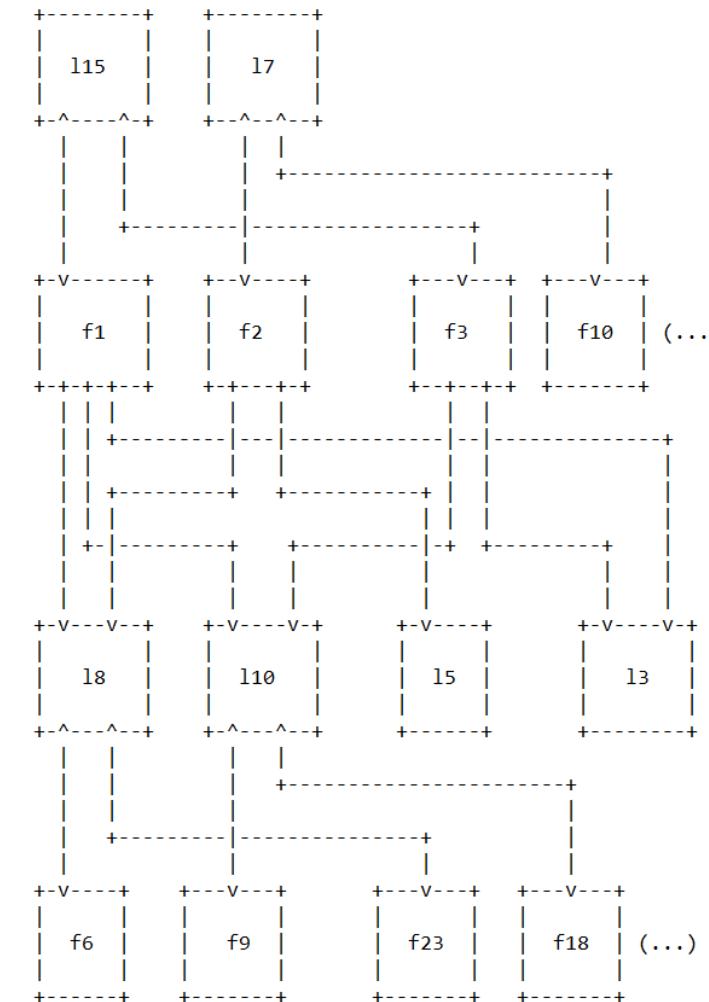
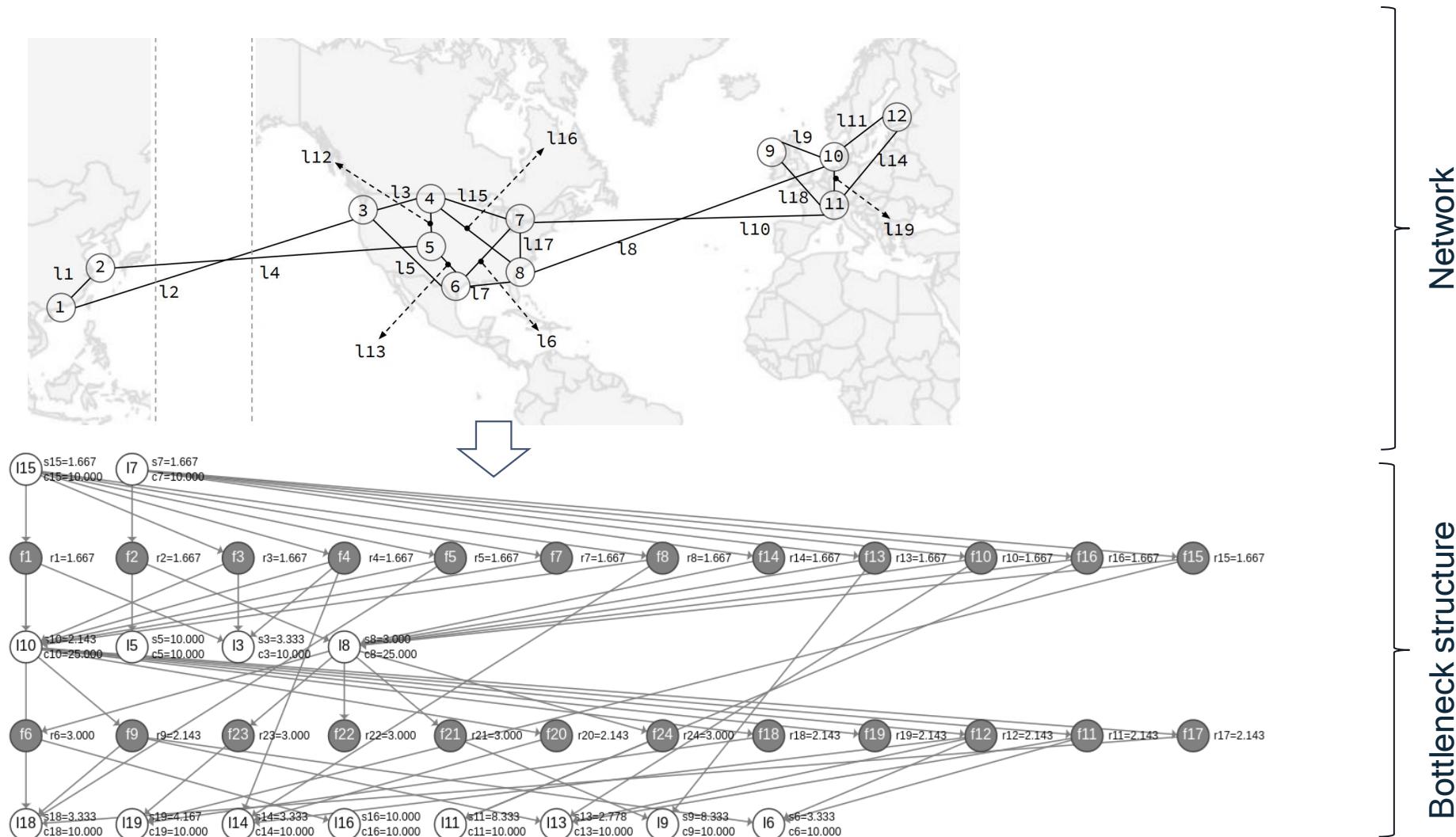
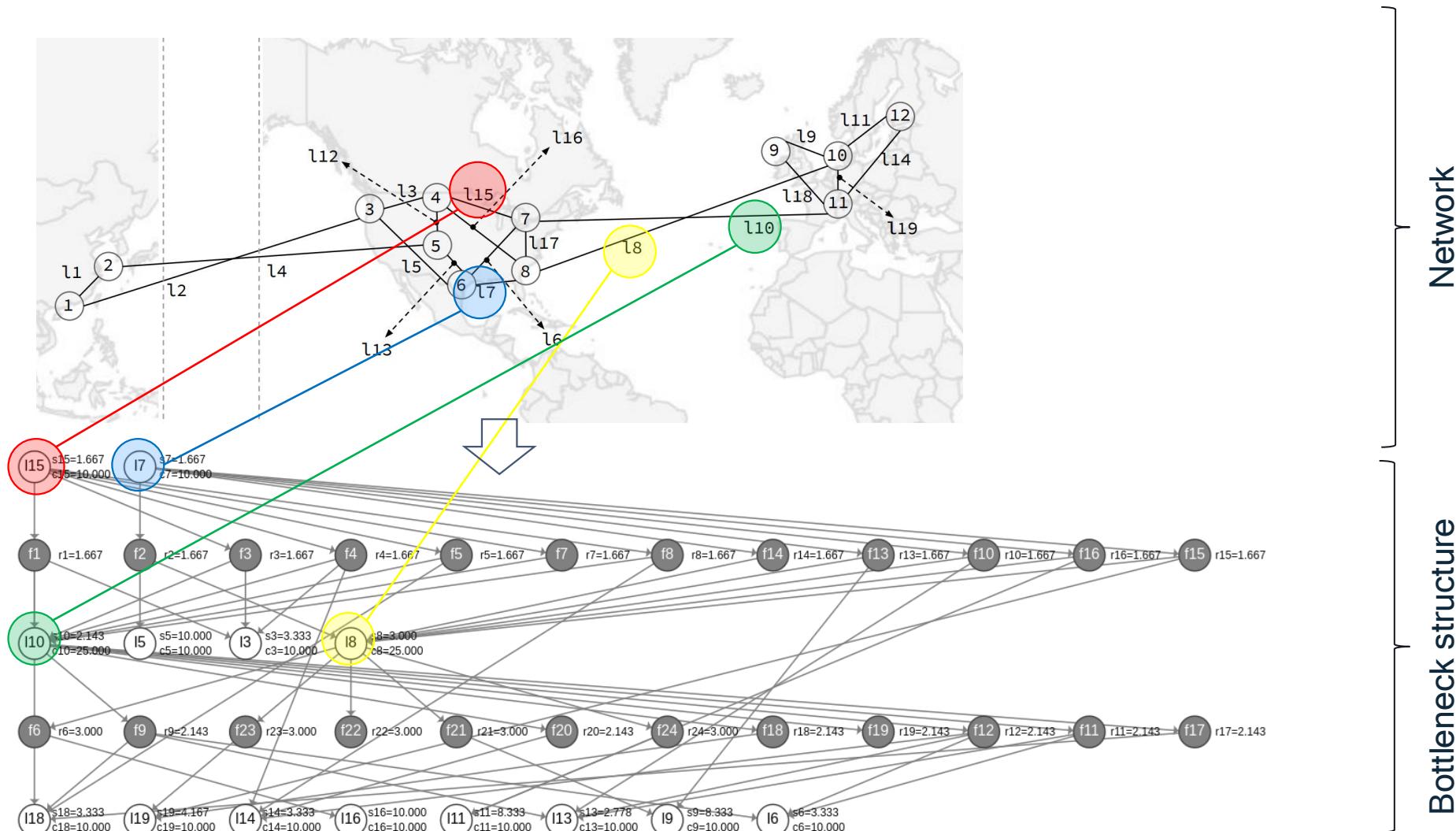


Figure 5: Bottleneck structure of Google's B4 network example.

BSG Use Cases: Optimizing Joint Routing and Congestion Control

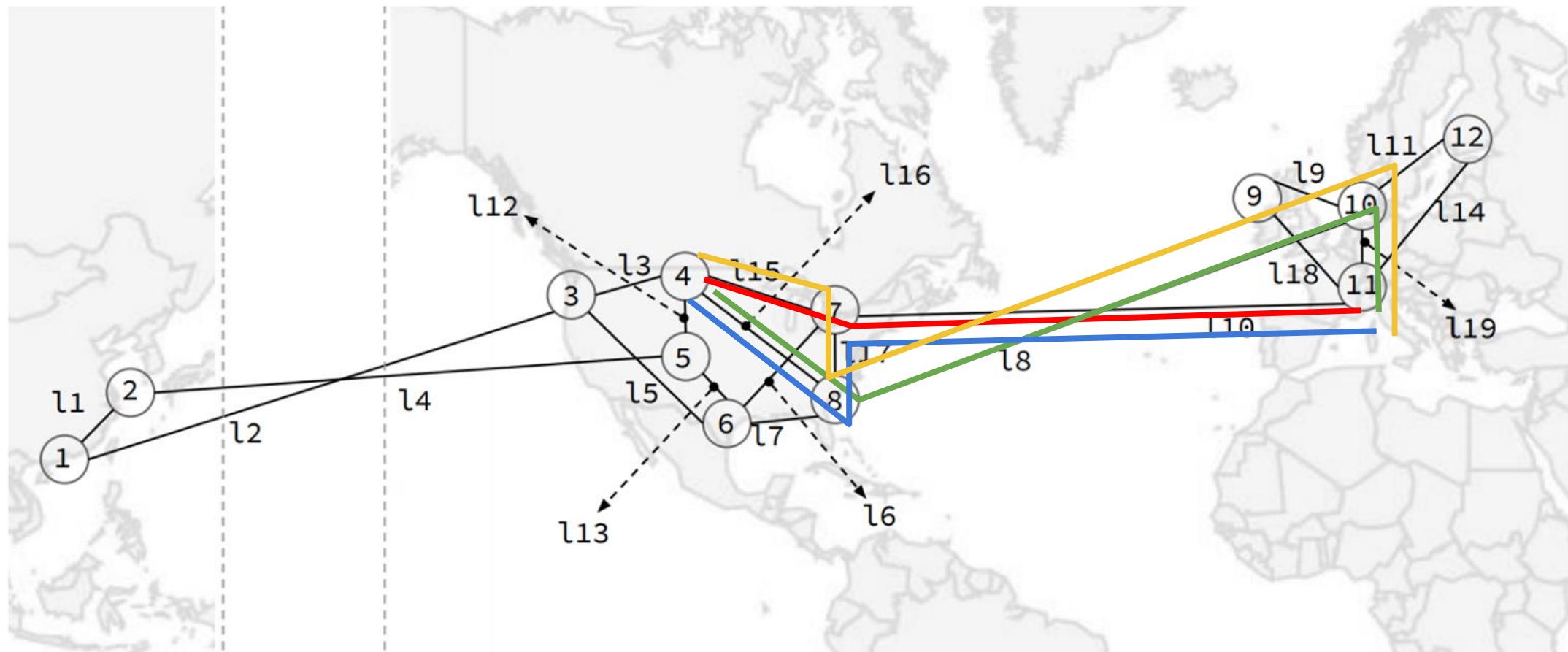


BSG Use Cases: Optimizing Joint Routing and Congestion Control



BSG Use Cases: Optimizing Joint Routing and Congestion Control

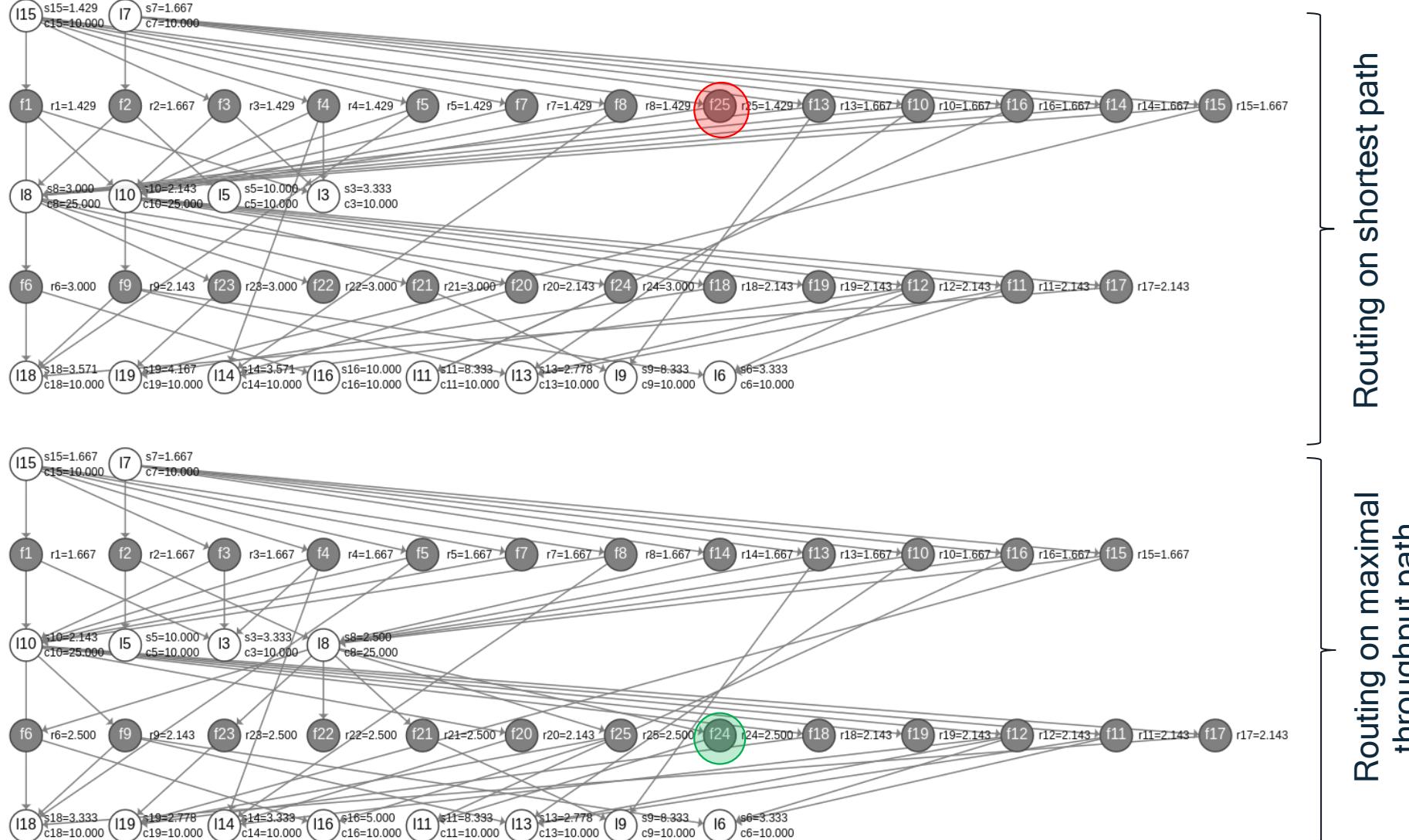
Suppose that an application needs to place a new flow on Google's B4 network to transfer a large data set from data center 11 (DC11) to data center 4 (DC4). There are multiple path choices:



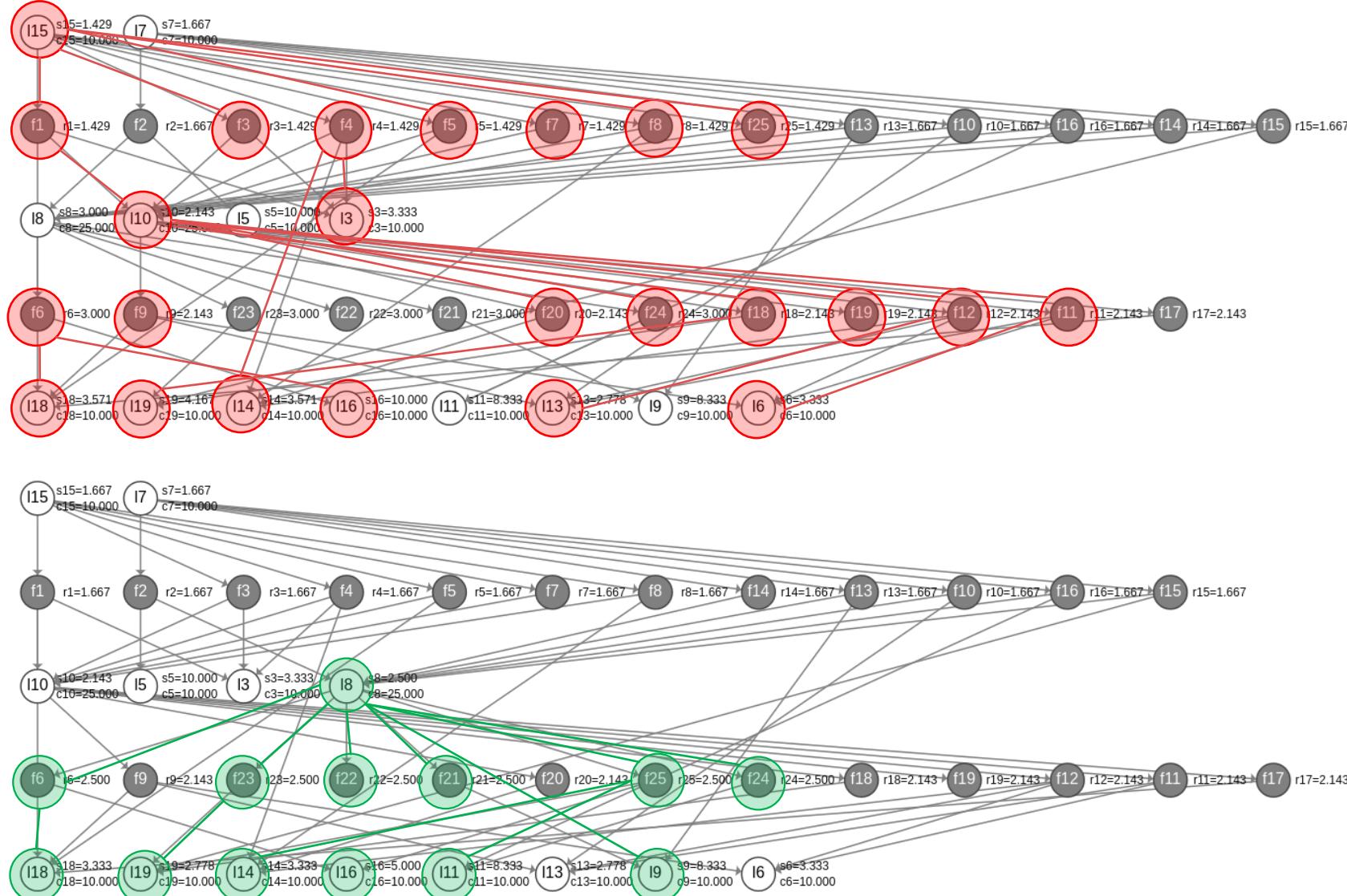
BSG Use Cases: Optimizing Joint Routing and Congestion Control

- Using bottleneck structures, we can compute in $O(V+E*\log(V))$ the path that will yield maximal throughput while considering the reaction of the congestion control algorithm.
- The optimal path corresponds to DC11 → I19 → DC10 → I8 → DC8 → I16 → DC4 yielding a throughput of 2.5 Gbps.
- Note that this is higher than the shortest path DC11 -> I10 -> DC7 -> I15 -> DC4, which yields a throughput of 1.429 Gbps.
- SLA management: Bottleneck structures can also be used to qualify and quantify the ripple effects produced on all other flows when placing the new flow to ensure their SLAs are preserved. See next slide.

BSG Use Cases: Optimizing Joint Routing and Congestion Control



BSG Use Cases: Optimizing Joint Routing and Congestion Control



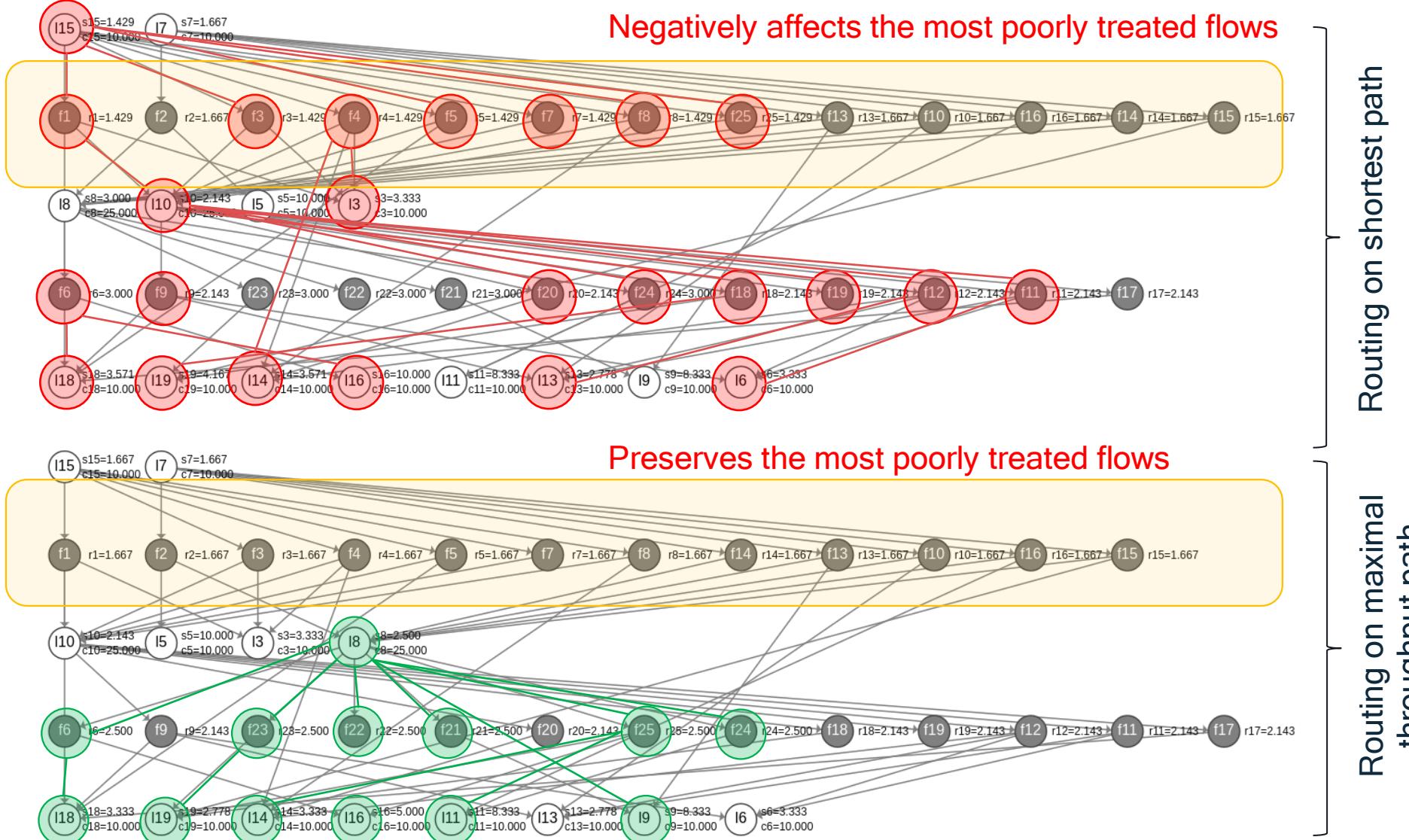
Routing on shortest path

Flows affected when placing flow on shortest path

Routing on maximal throughput path

Flows affected when placing flow on maximal throughput path

BSG Use Cases: Optimizing Joint Routing and Congestion Control



Production Deployments

Production Deployments

A full-stack software (southbound/northbound APIs and kernel) implementation of Bottleneck Structure Graphs called GradientGraph is deployed in two US-wide production networks:

- National Research Platform: <https://nationalresearchplatform.org/>
- DOE/Esnets: <https://www.es.net/>
- For lack of time in this session, we plan to report deployment experience in the next IETF Meetings. We can also discuss offline if you are interested.

Discussion

Discussion

Guidance on what RGs and WGs would be suitable to discuss the proposed distributed protocol:

- PANRG
- COINRG
- ALTO
- ICCRG
- PCE
- TEAS
- NMRG
- Others