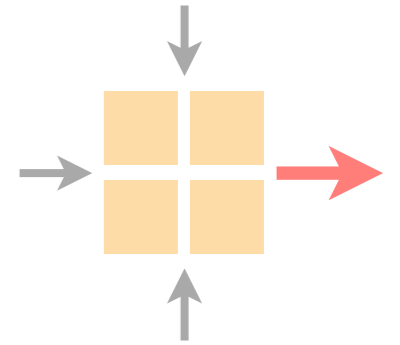


Advanced Topics in Communication Networks

Programming Network Data Planes



Laurent Vanbever

nsg.ee.ethz.ch

ETH Zürich


Nov 1 2018

Last week on

Advanced Topics in Communication Networks

We looked at the Tofino architecture together with two (key, value) store applications: Net/{Cache, Chain}





Programmable Data Plane at Terabit Speeds


 Vladimir Gurevich

 May 16, 2017

BAREFOOT NETWORKS
Copyright © 2017 - Barefoot Networks


**NetCache: Balancing Key-Value Stores
with Fast In-Network Caching**

Xin Jin, Xiaozhou Li, Haoyu Zhang, Robert Soulé,
 Jeongkeun Lee, Nate Foster, Changhoon Kim, Ion Stoica



**NetChain: Scale-Free Sub-RTT
Coordination**

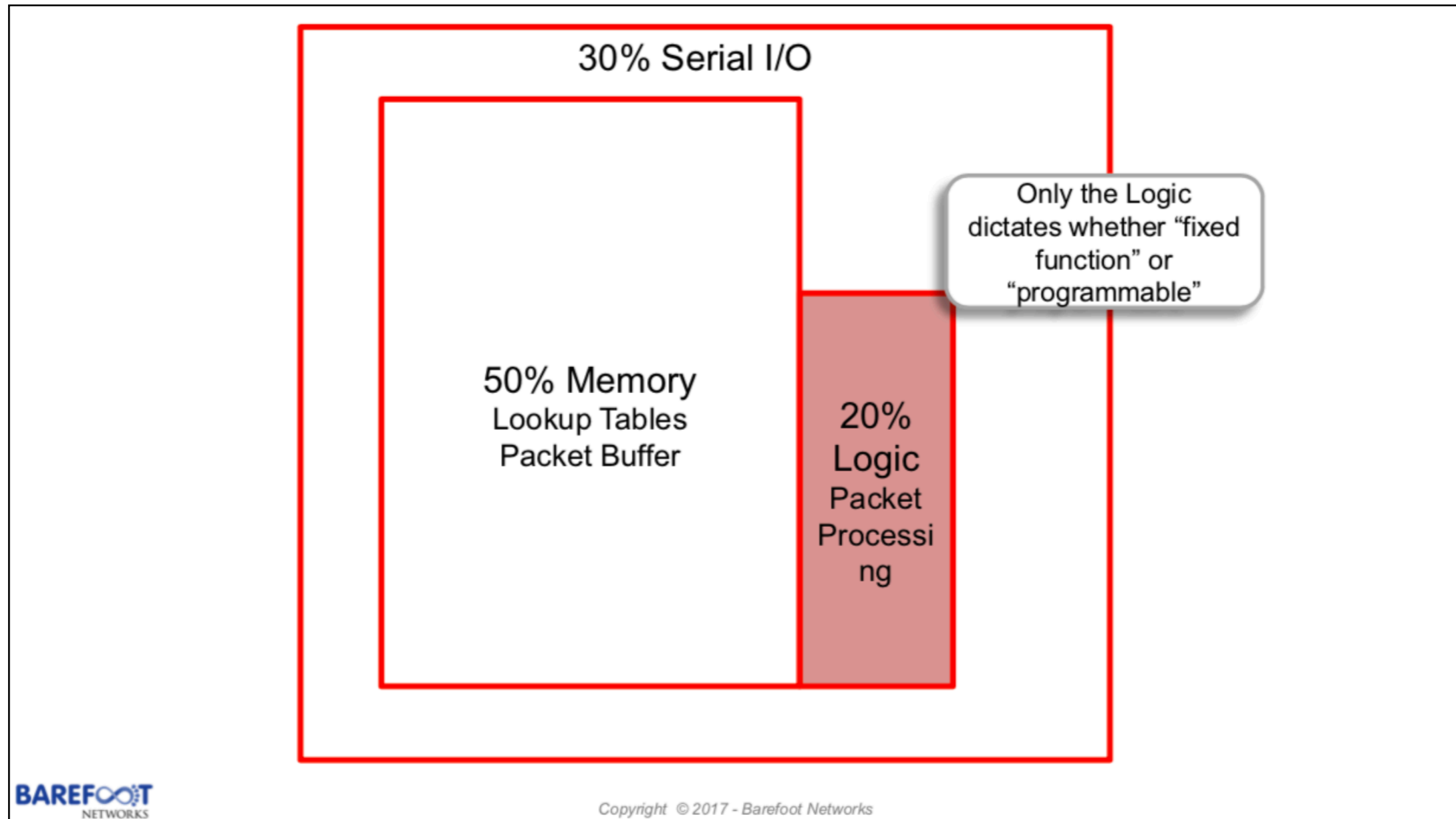
Xin Jin
 Xiaozhou Li, Haoyu Zhang, Robert Soulé, Jeongkeun Lee,
 Nate Foster, Changhoon Kim, Ion Stoica



“Programmable switches are 10-100x slower than fixed-function switches. They cost more and consume more power.”

Conventional wisdom in networking

One of the main enabler for data-plane programmability is the shrinking size of the packet processing logic chip.



Source: Programmable Data Planes at Terabit Speeds, Vladimir Gurevich, 2017

Barefoot Tofino processes packets in parallel,
even though the semantic of a P4 program is sequential

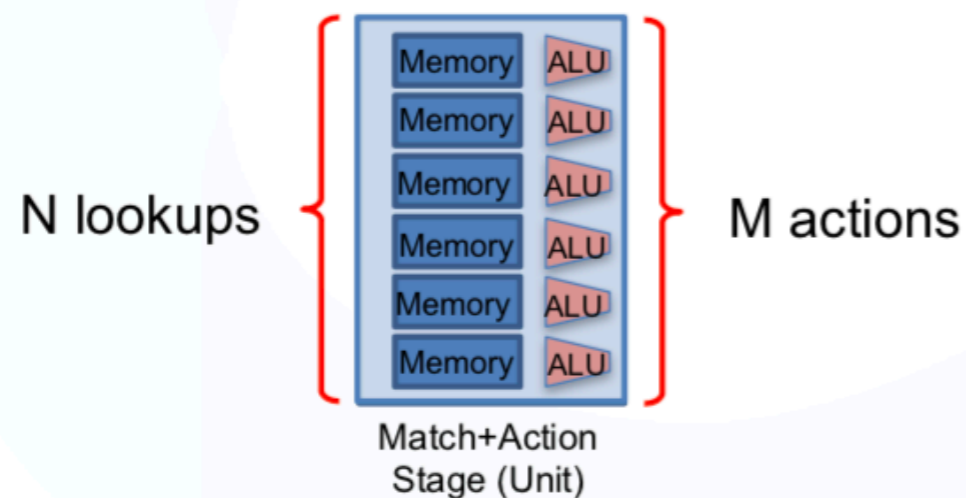
Parallelism and alternatives

- **Sequential semantics does not prohibit parallelism**
- **Doing everything does not mean doing everything all the time**

Barefoot Tofino processes packets in parallel,
even though the semantic of a P4 program is sequential

PISA: Important Details

- Multiple simultaneous lookups and actions can be supported

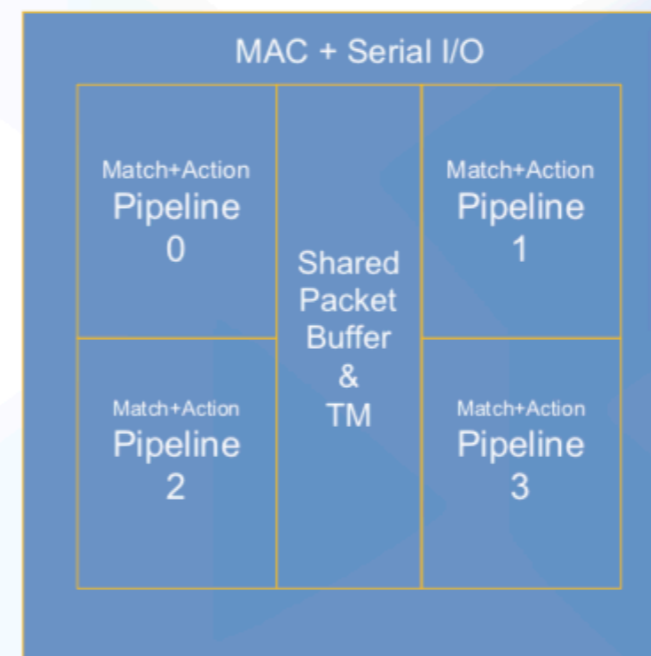


Barefoot Tofino 6.5 Tbps backplane

several billion packets per second at line rate

6.5Tb/s Tofino™ Summary

- **State of the art design**
 - Single Shared Packet Buffer
 - TSMC 16nm FinFET+
- **Four Match+Action Pipelines**
 - Fully programmable PISA Embodiment
 - All compiled programs run at line-rate.
 - Up to 1.3 million IPv4 routes
- **Port Configurations**
 - 65 x 100GE/40GE
 - 130 x 50GE
 - 260 x 25GE/10GE
- **CPU Interfaces**
 - PCIe: Gen3 x4/x2/x1
 - Dedicated 100GE port



Tofino relies on Packet Header Vector (PHV) to pass states between stages—this is one of the limiting factor

Packet Header Vector (PHV)

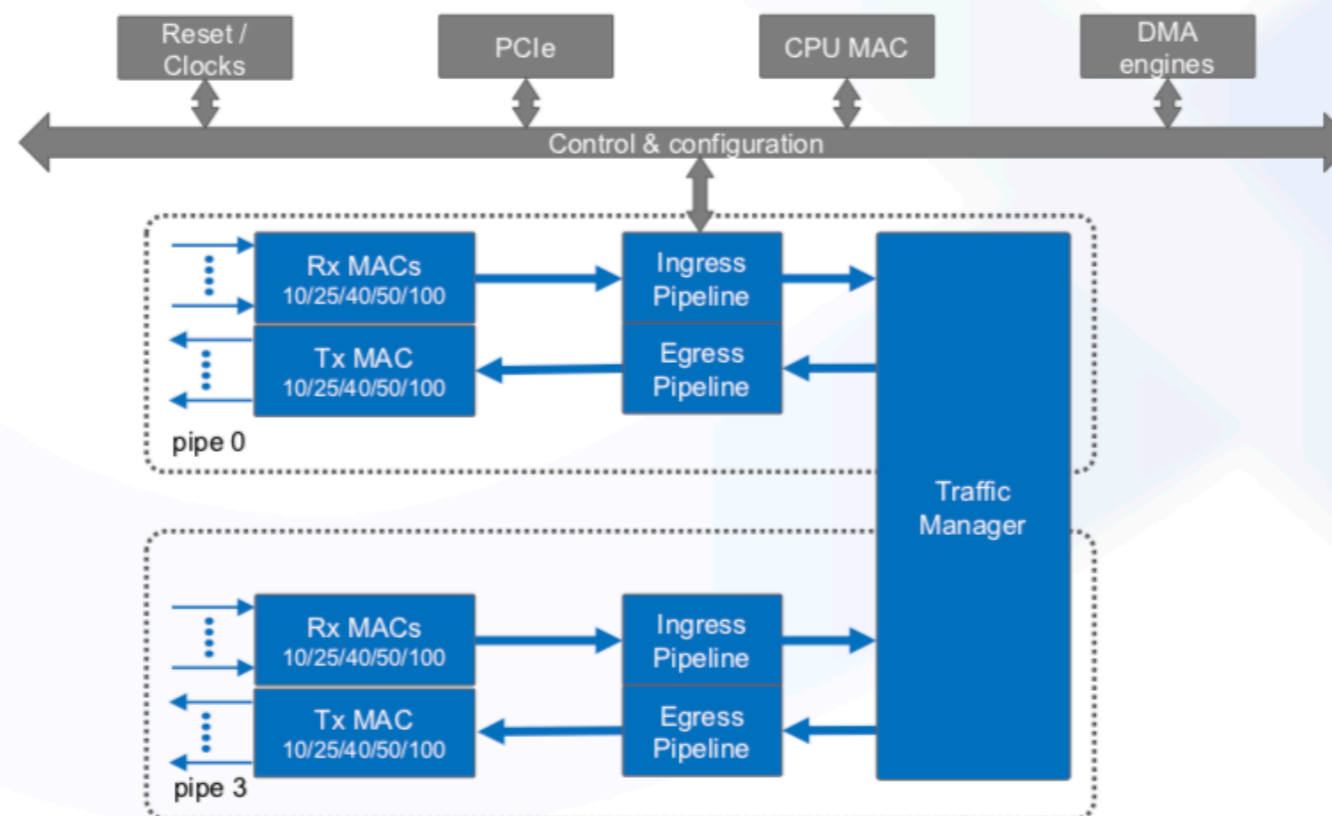
- A set of uniform containers that carry the headers and metadata along the pipeline
- Fields can be packed into any container or their combination
- PHV Allocation step in the compiler decides the actual packing



Tofino uses a folded pipeline in which the *same* stages are used for both the ingress and the egress pipeline

Unified Pipeline

- **There is no difference between ingress and egress processing**
 - The same blocks can be efficiently shared



PROGRAMMABLE DATA PLANE AT TERABIT SPEEDS

BAREFOOT NETWORKS

Programmable Data Plane at Terabit Speeds

Vladimir Gurevich
May 16, 2017

BAREFOOT NETWORKS Copyright © 2017 - Barefoot Networks

NETCACHE: BALANCING KEY-VALUE STORES WITH FAST IN-NETWORK CACHING

NetCache: Balancing Key-Value Stores with Fast In-Network Caching

Xin Jin, Xiaozhou Li, Haoyu Zhang, Robert Soulé
Jeongkeun Lee, Nate Foster, Changhoon Kim, Ion Stoica

JOHNS HOPKINS UNIVERSITY BAREFOOT NETWORKS PRINCETON UNIVERSITY Cornell University Berkeley UNIVERSITY OF CALIFORNIA

NETCHAIN: SCALE-FREE SUB-RTT COORDINATION

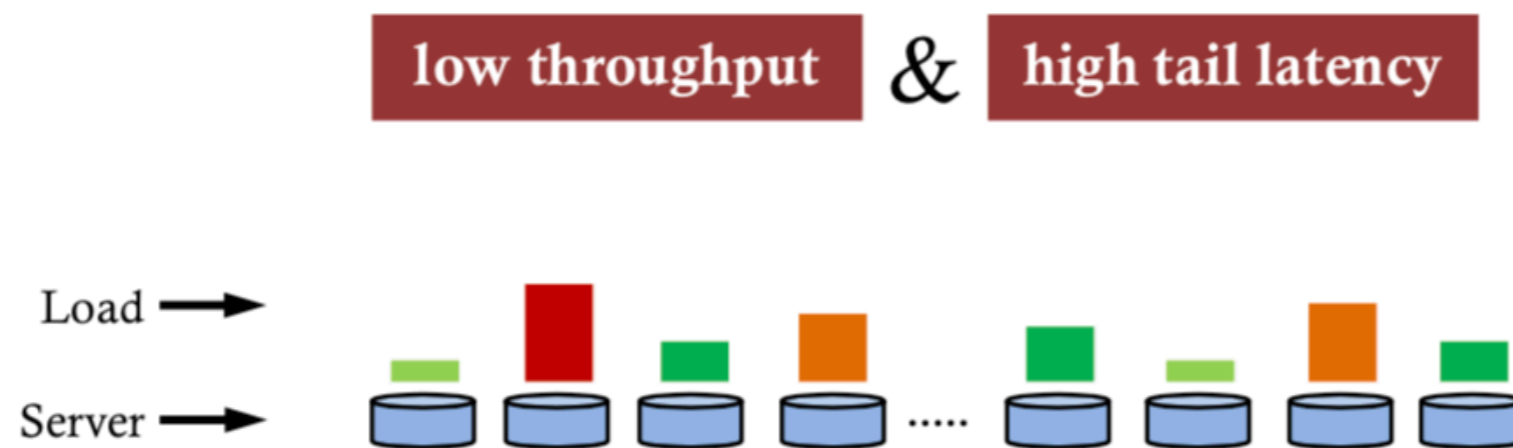
NetChain: Scale-Free Sub-RTT Coordination

Xin Jin
Xiaozhou Li, Haoyu Zhang, Robert Soulé, Jeongkeun Lee,
Nate Foster, Changhoon Kim, Ion Stoica

JOHNS HOPKINS UNIVERSITY BAREFOOT NETWORKS PRINCETON UNIVERSITY Università della Svizzera Italiana Cornell University Berkeley UNIVERSITY OF CALIFORNIA

NetCache solves the problem of load-balancing in key-values stores observing *dynamic, skewed* workload

Key challenge: *highly-skewed* and rapidly-changing workloads



It leverages that a small but very fast cache can provide perfect load-balancing... in theory

Opportunity: fast, small cache can ensure load balancing

[B. Fan et al. SoCC'11, X. Li et al. NSDI'16]

Cache $O(N \log N)$ hottest items

E.g., 10,000 hot objects

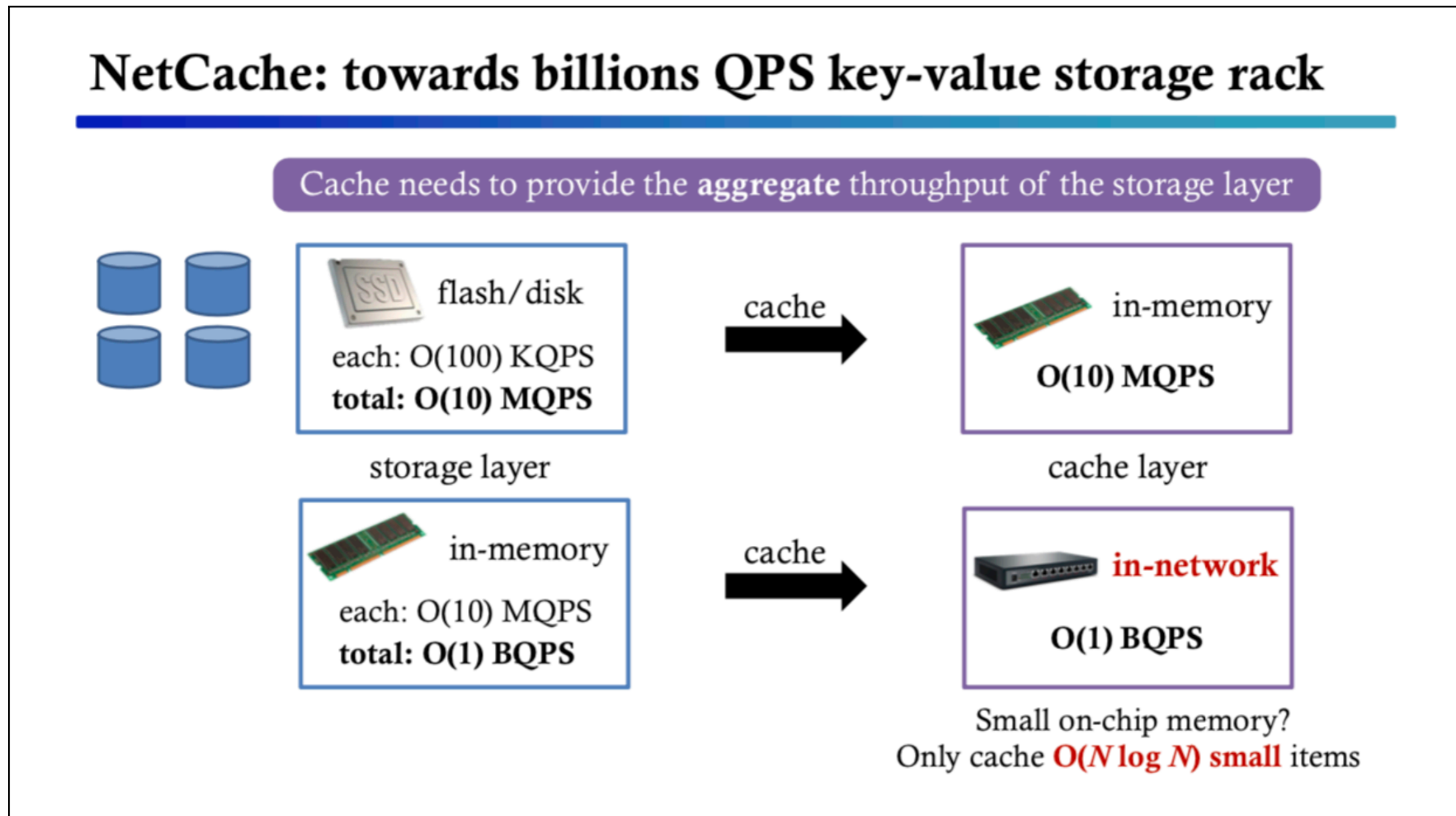
N: # of servers



E.g., 100 backends with 100 billions items

Requirement: cache throughput \geq backend aggregate throughput

NetCache relies on the O(billion) throughput of programmable network devices to achieve it in practice



It relies on a tailored UDP-based protocol, an de/encoding scheme for storing variable length values, and sketches

Key-value caching in network ASIC at line rate ?!

- ❑ How to identify application-level packet fields ?
- ❑ How to store and serve variable-length data ?
- ❑ How to efficiently keep the cache up-to-date ?

~ p4_02_2017_programmable_data_plane_at_terabit_speeds.pdf (page 1 of 40)



BAREFOOT
NETWORKS

Programmable Data Plane at Terabit Speeds


Vladimir Gurevich
May 16, 2017

BAREFOOT NETWORKS Copyright © 2017 - Barefoot Networks

~ NSDI17_NetCache_slides.pdf (page 1 of 30)

NetCache: Balancing Key-Value Stores with Fast In-Network Caching

Xin Jin, Xiaozhou Li, Haoyu Zhang, Robert Soulé
Jeongkeun Lee, Nate Foster, Changhoon Kim, Ion Stoica



~ NSDI18_NetChain_slides.pdf (page 1 of 37)

NetChain: Scale-Free Sub-RTT Coordination

Xin Jin
Xiaozhou Li, Haoyu Zhang, Robert Soulé, Jeongkeun Lee,
Nate Foster, Changhoon Kim, Ion Stoica



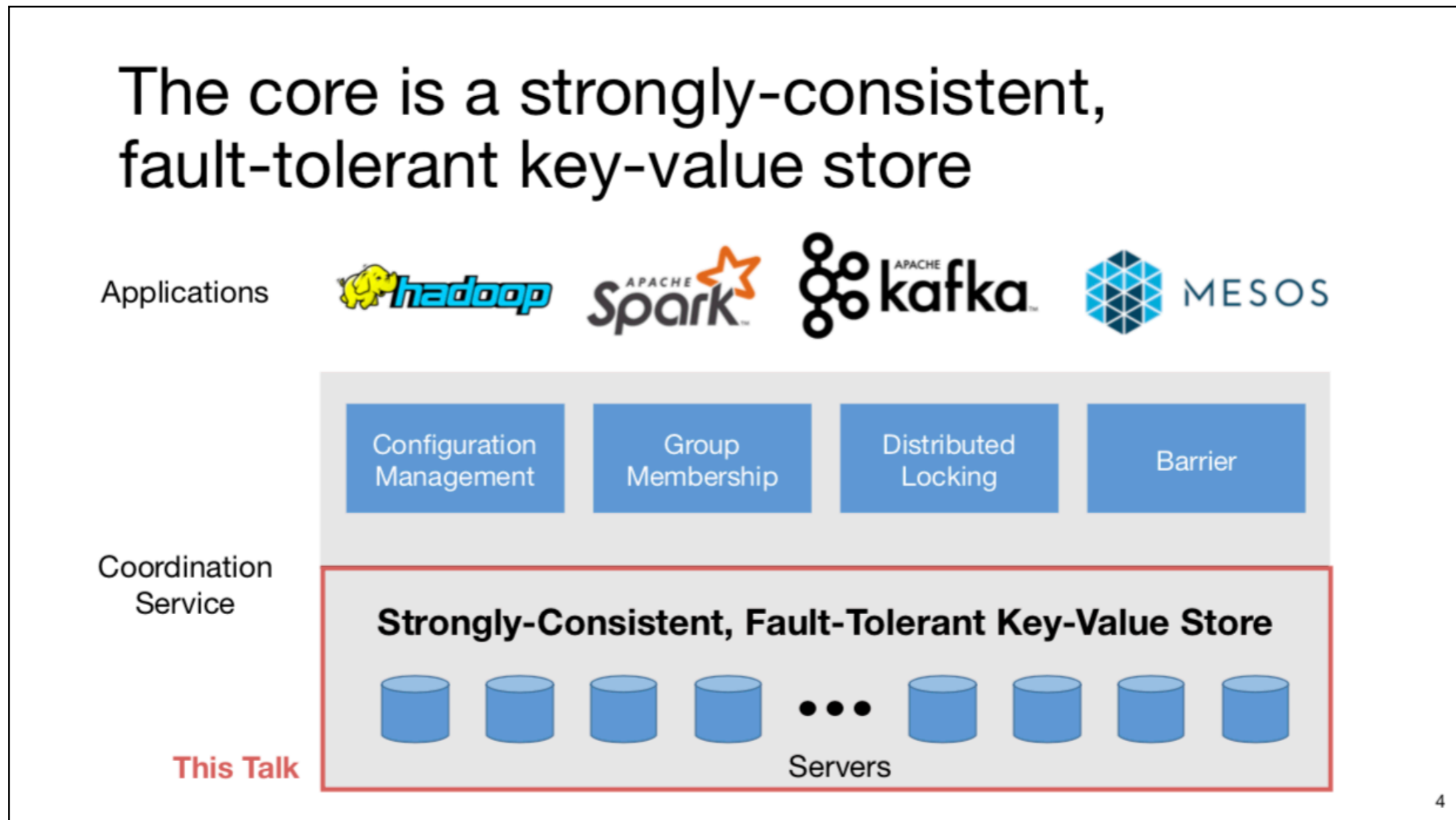
NetChain builds upon NetCache to scale coordination services, a key building block of distributed systems

Conventional wisdom: **avoid coordination**

NetChain: **lightning fast coordination**
enabled by programmable switches

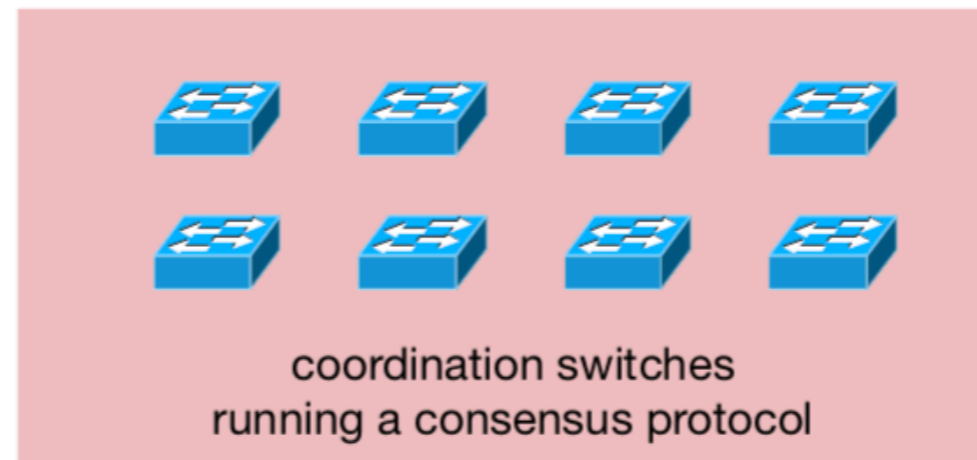
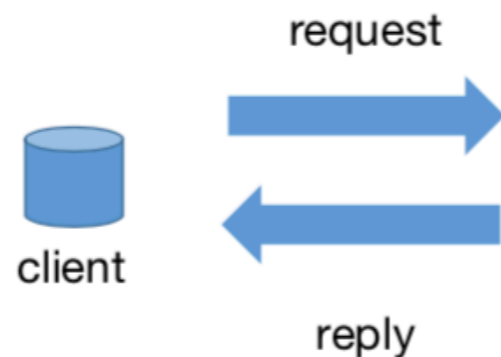
Open the door to rethink distributed systems design

Coordination services typically rely on a replicated key-value store for consistency and fault-tolerance



State of the art server-based coordination services struggle to provide high-throughput and low-latency

Opportunity: **in-network** coordination



- Throughput: **switch throughput**
- Latency: **half of an RTT**

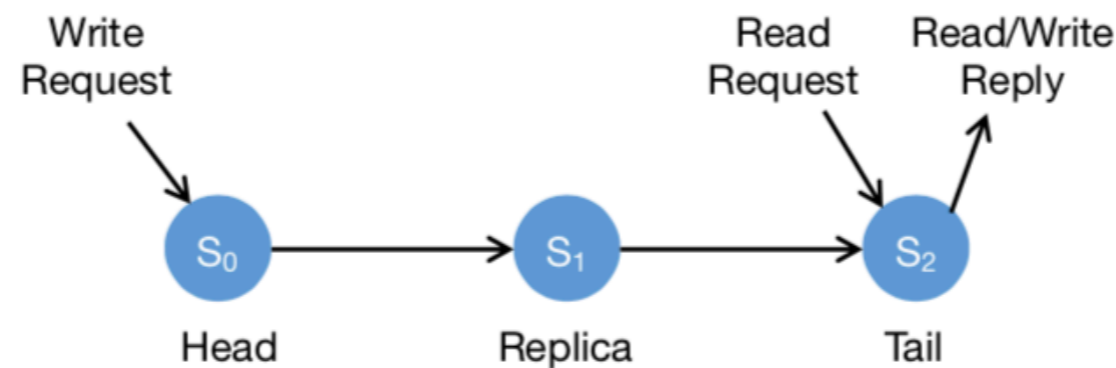
Key challenge is to ensure consistency and fault-tolerance

Design goals for coordination services

- High throughput
 - Low latency
 - Strong consistency
 - Fault tolerance
- } Directly from high-performance switches
- } Chain replication in the network

NetChain does so using chain replication, building upon NetCache for storing values in each switch

What is chain replication



- Storage nodes are organized in a **chain** structure
- Handle operations
 - **Read** from the **tail**
 - **Write** from **head** to **tail**
- Provide strong consistency and fault tolerance
 - Tolerate **f failures** with **f+1 nodes**

NetChain relies on a tailored UDP-based protocol, source-routing mechanisms and message serialization

How to build a strongly-consistent, fault-tolerant, in-network key-value store

- How to store and serve key-value items?
 - How to route queries according to chain structure?
 - How to handle out-of-order delivery in network?
 - How to handle switch failures?
-
- The diagram illustrates the mapping of the four questions to network planes. A large right-facing curly bracket groups the first three questions (storage, routing, and delivery) under the label 'Data Plane'. A horizontal arrow points from the fourth question (switch failures) to the label 'Control Plane'.

This week on

Advanced Topics in Communication Networks

A high-level, **non-exhaustive** overview of the research surrounding data plane programmability

A high-level, non-exhaustive overview of the research surrounding data plane programmability

Data plane programmability for

Performance
Monitoring
Applications offloading

Platforms for Data plane programmability
Correctness
Management

Data plane
programmability for

Performance

Monitoring

Applications offloading

Platforms

for

Data plane

Correctness

programmability

Management

A large set of papers on programmable data planes aim at improving performance, esp. load balancing

HULA [SOSR'16]

DRILL [SIGCOMM'17]



CONGA [SIGCOMM'14]

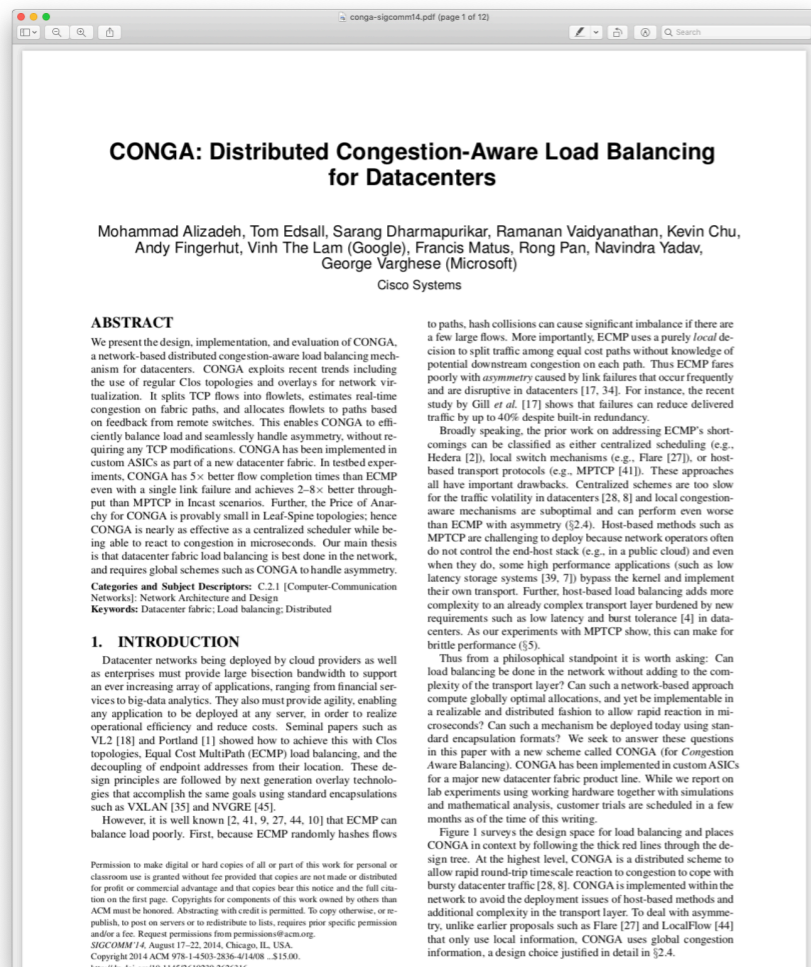


LetFlow [NSDI'17]

A large set of papers on programmable data planes aim at improving performance, esp. load balancing

HULA [SOSR'16]

DRILL [SIGCOMM'17]



CONGA [SIGCOMM'14]



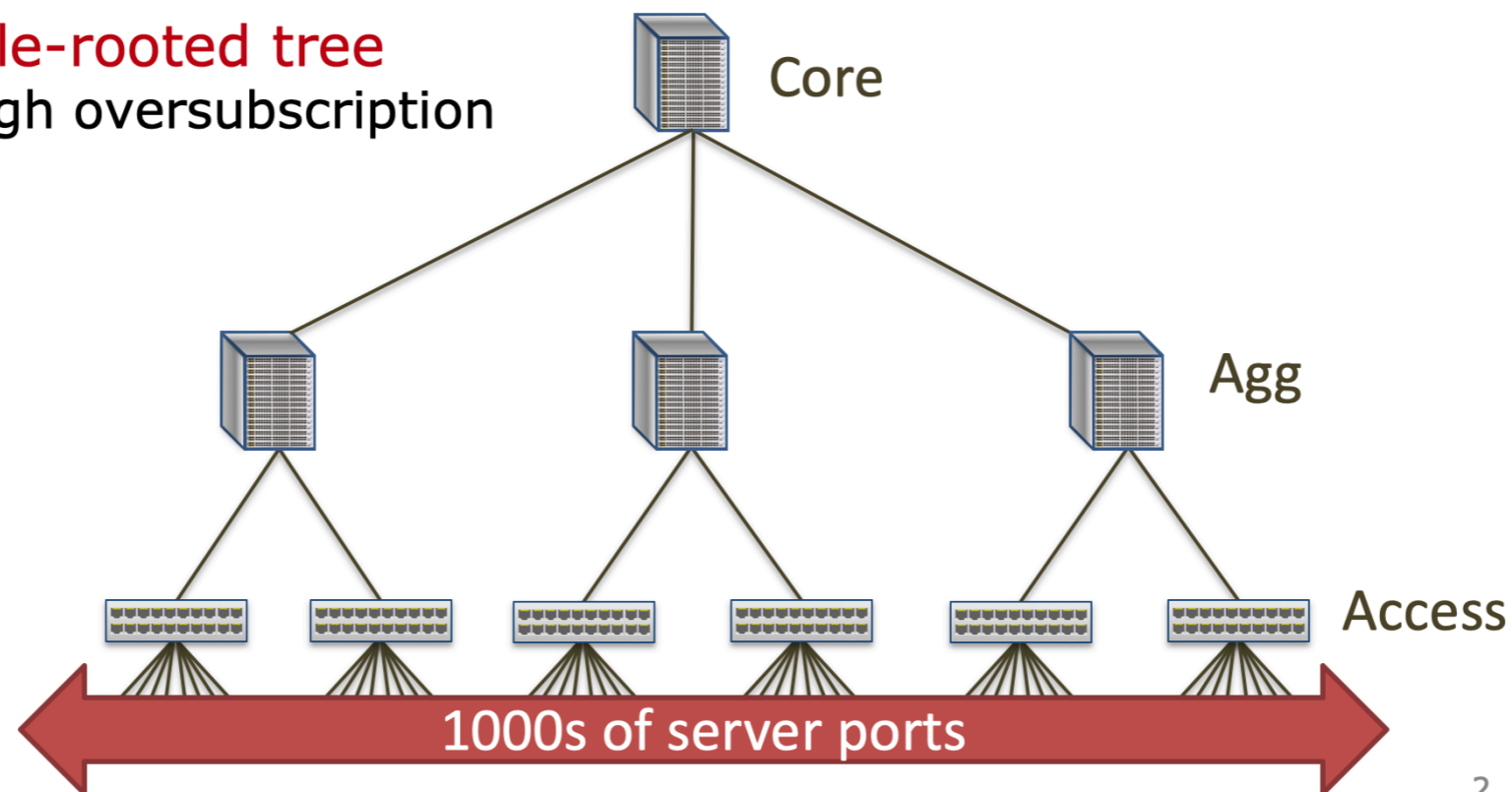
LetFlow [NSDI'17]

Motivation

DC networks need large bisection bandwidth for **distributed** apps (big data, HPC, web services, etc)

Single-rooted tree

- High oversubscription



2

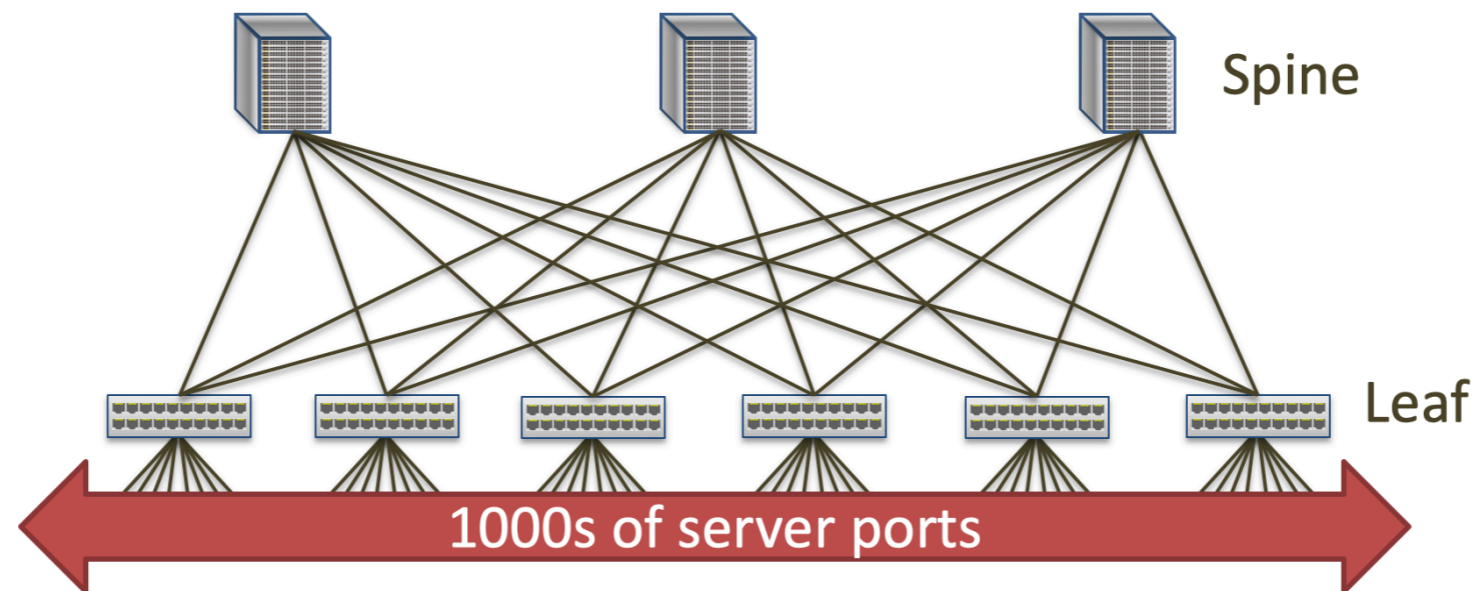
Source: CONGA: Distributed Congestion-Aware Load Balancing for Datacenters, Mohammad Alizadeh et al., 2014

Motivation

DC networks need large bisection bandwidth for **distributed** apps (big data, HPC, web services, etc)

Multi-rooted tree [Fat-tree, Leaf-Spine, ...]

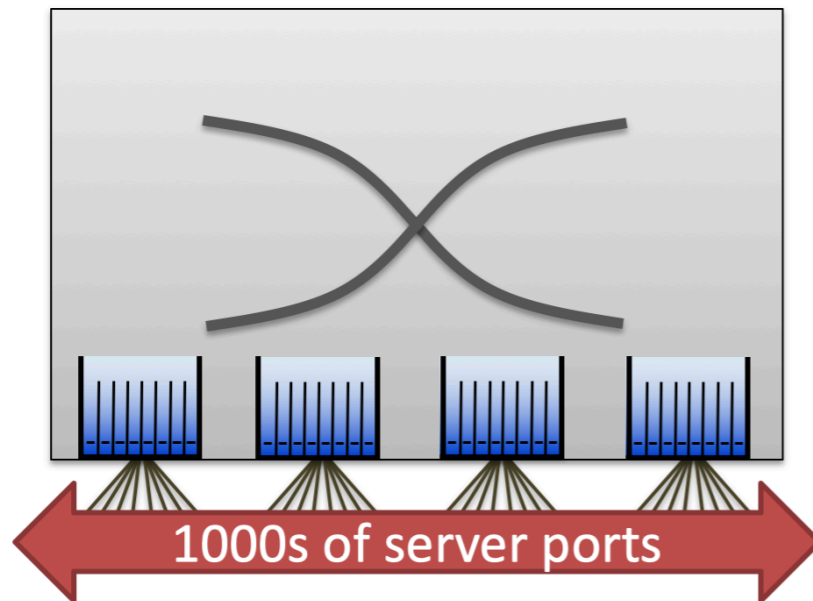
- Full bisection bandwidth, achieved via multipathing



2

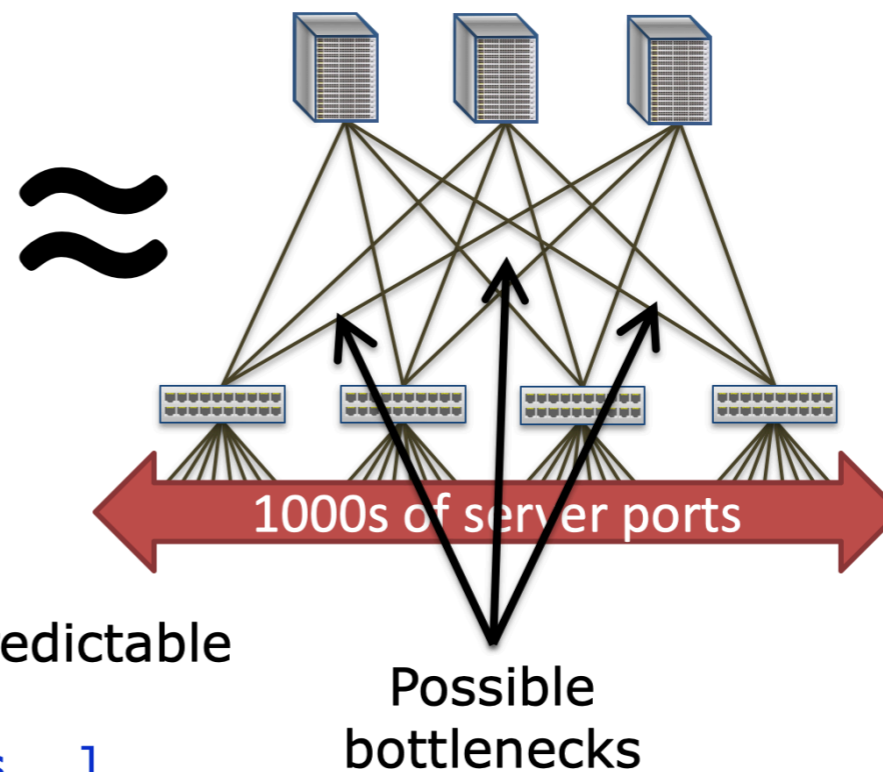
Multi-rooted != Ideal DC Network

Ideal DC network:
Big output-queued switch



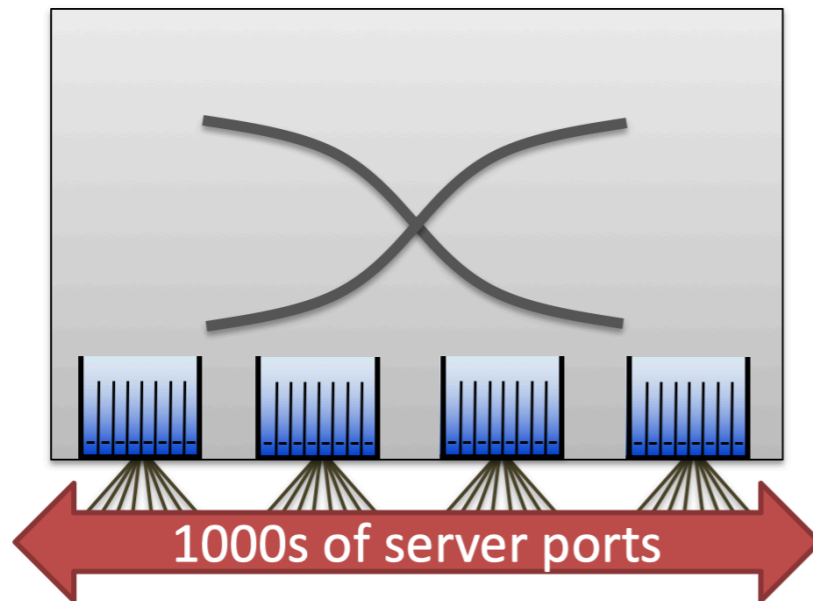
- No internal bottlenecks → predictable
- Simplifies BW management
[EyeQ, FairCloud, pFabric, Varys, ...]

Multi-rooted tree

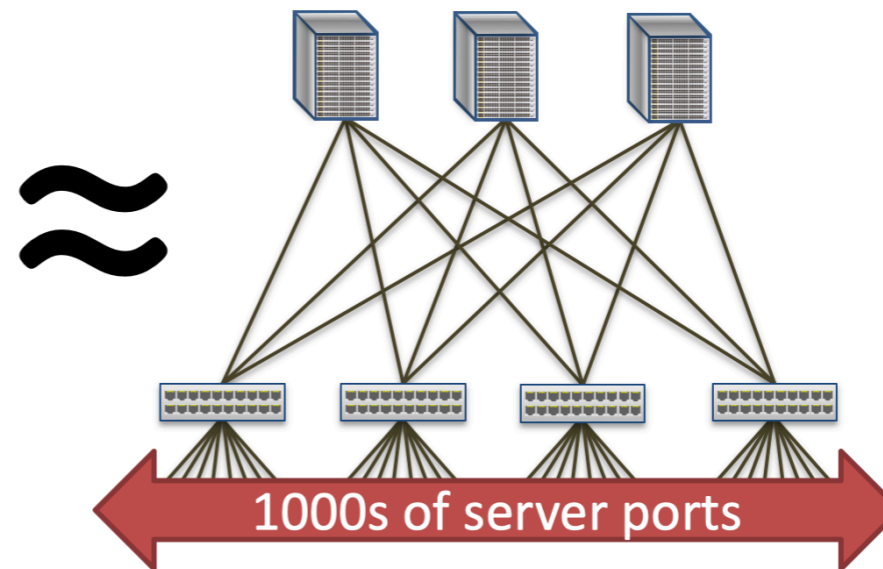


Multi-rooted \neq Ideal DC Network

Ideal DC network:
Big output-queued switch



Multi-rooted tree



Need precise load balancing

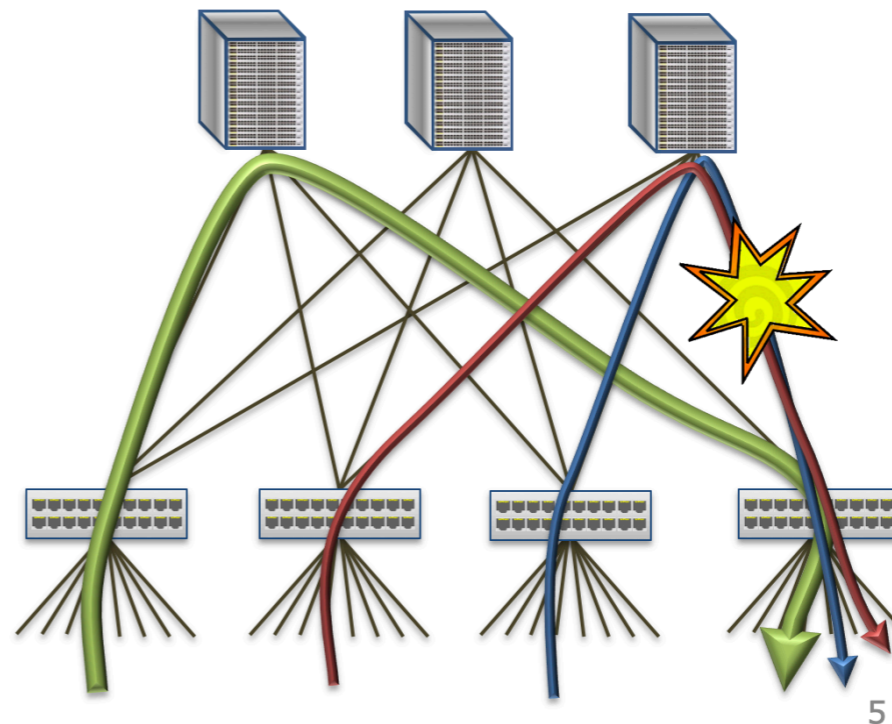
Today: ECMP Load Balancing

Pick among equal-cost paths by a **hash** of 5-tuple

- Approximates Valiant load balancing
- Preserves packet order

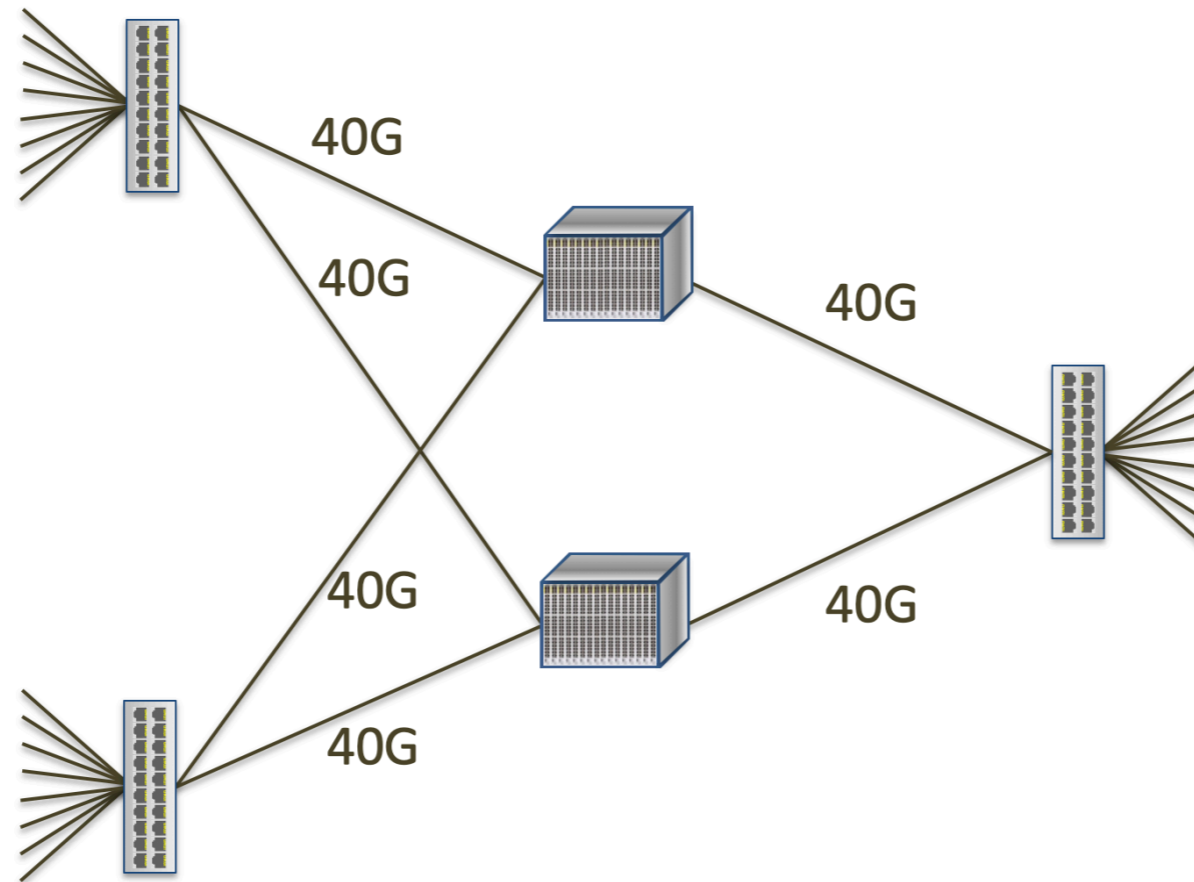
Problems:

- Hash collisions
(coarse granularity)
- Local & stateless
(v. bad with asymmetry
due to link failures)



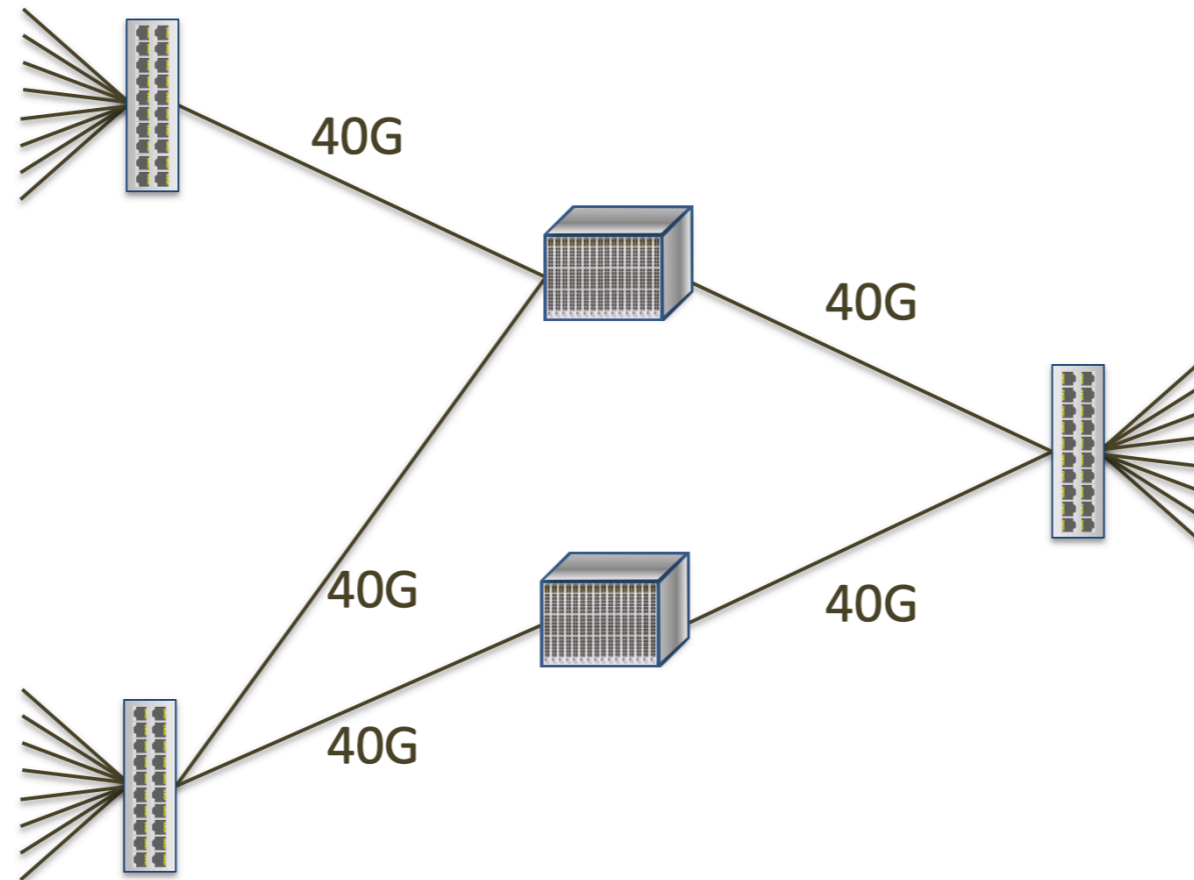
Dealing with Asymmetry

Handling asymmetry needs non-local knowledge

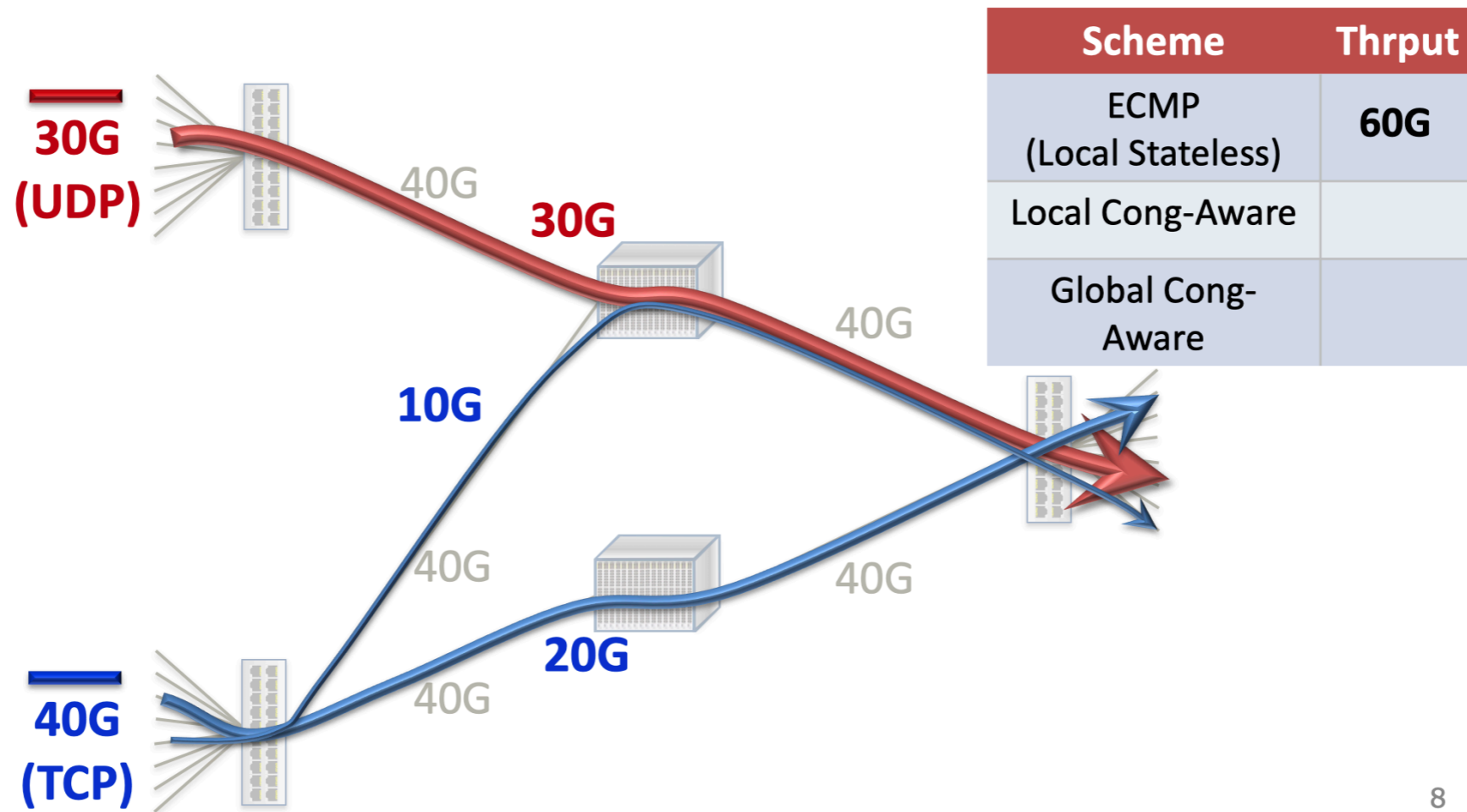


Dealing with Asymmetry

Handling asymmetry needs non-local knowledge

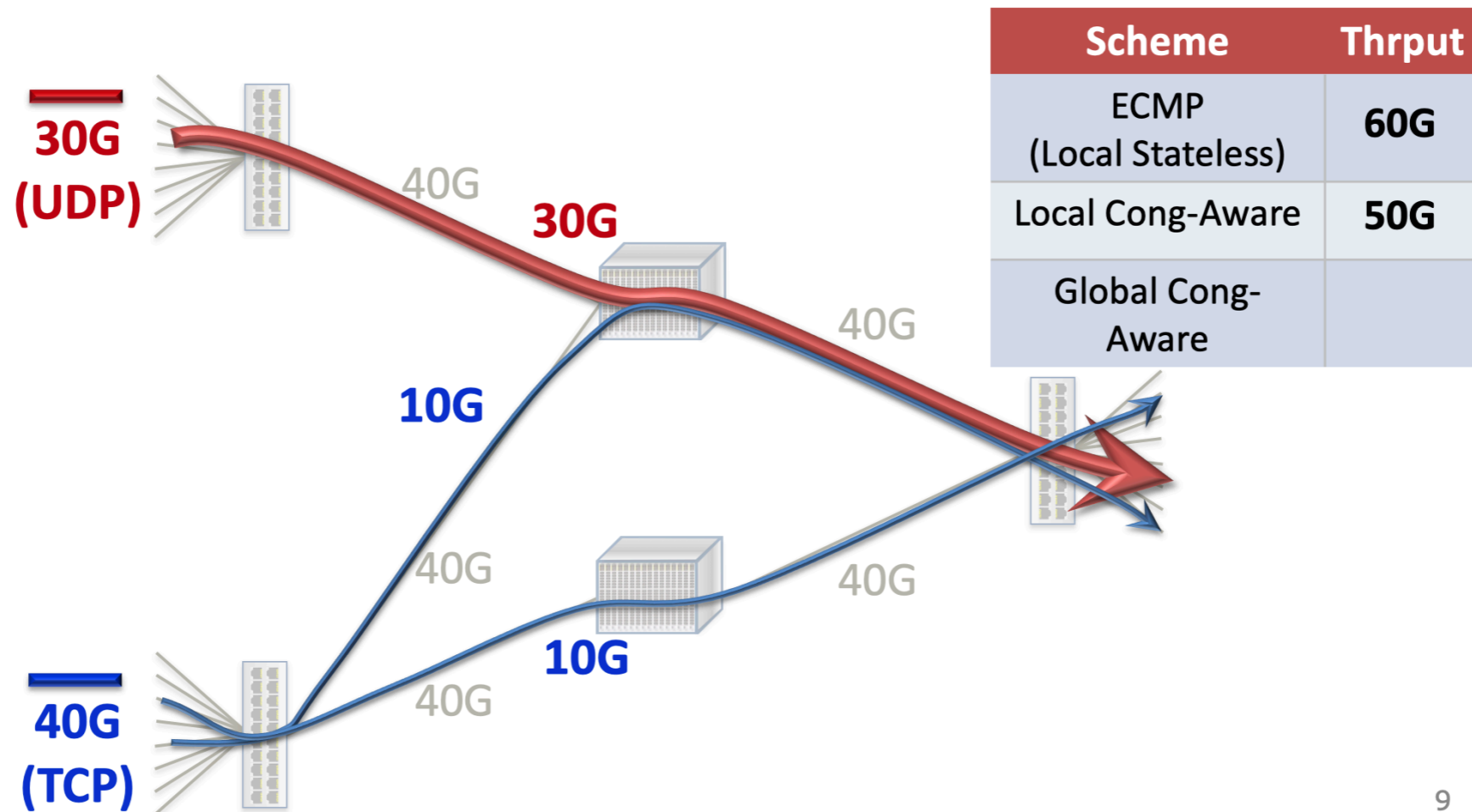


Dealing with Asymmetry: ECMP



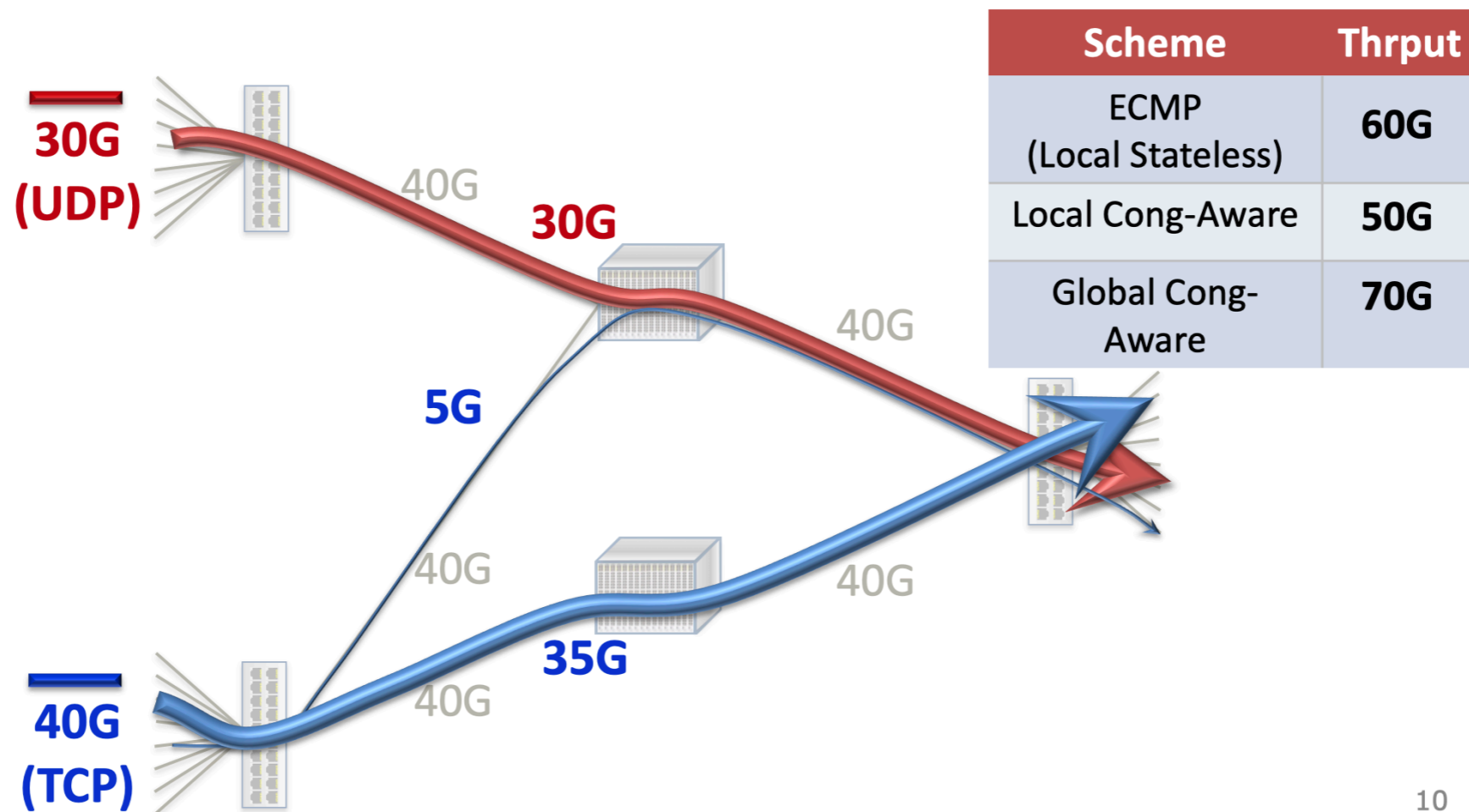
Source: CONGA: Distributed Congestion-Aware Load Balancing for Datacenters,
Mohammad Alizadeh et al., 2014

Dealing with Asymmetry: Local Congestion-Aware



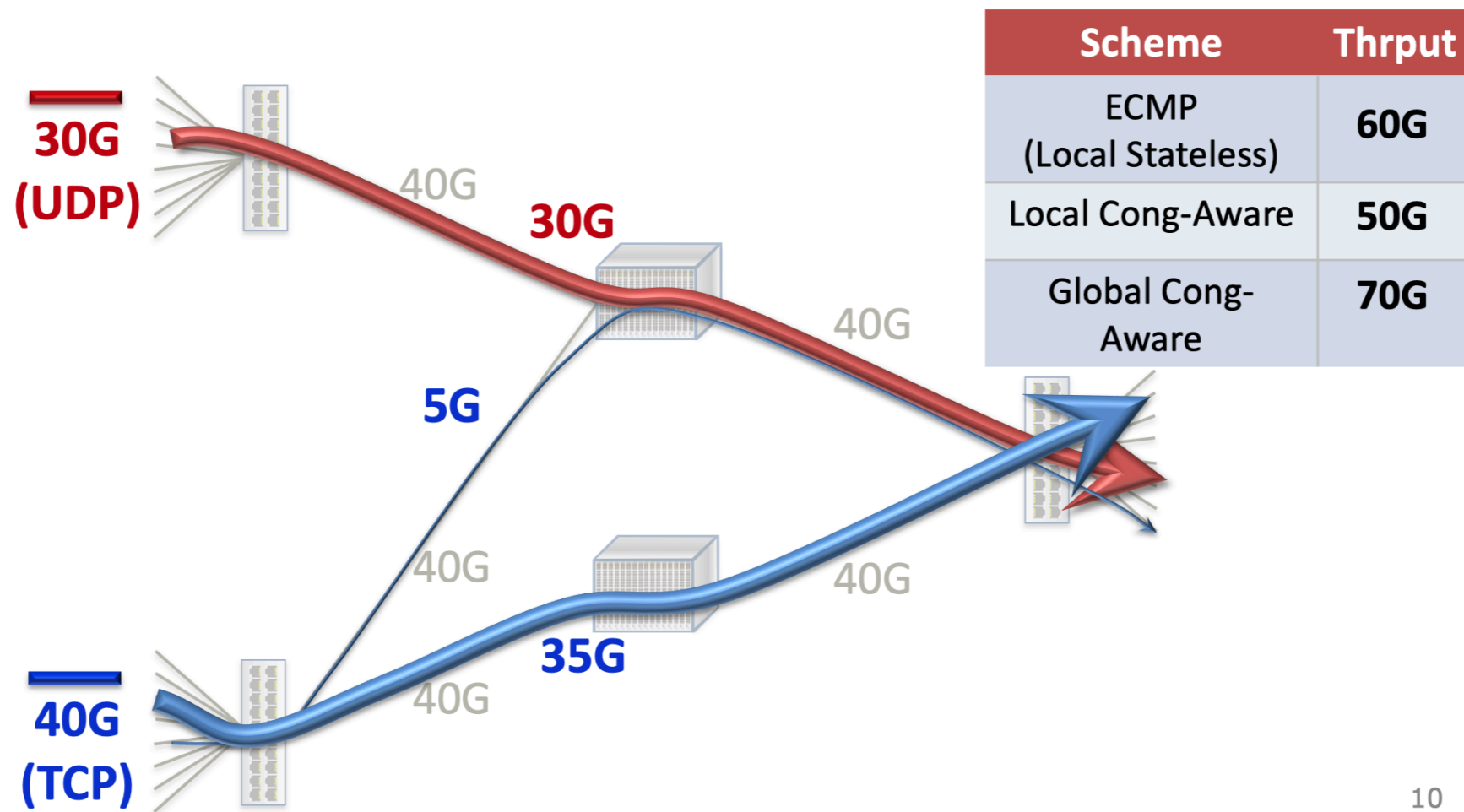
Source: CONGA: Distributed Congestion-Aware Load Balancing for Datacenters,
Mohammad Alizadeh et al., 2014

Dealing with Asymmetry: Global Congestion-Aware



Source: CONGA: Distributed Congestion-Aware Load Balancing for Datacenters,
Mohammad Alizadeh et al., 2014

Dealing with Asymmetry: Global Congestion-Aware



Source: CONGA: Distributed Congestion-Aware Load Balancing for Datacenters,
Mohammad Alizadeh et al., 2014

Global Congestion-Awareness (in Datacenters)

		Datacenter
Opportunity →	{ Latency Topology	microseconds simple, regular
Challenge →	Traffic	volatile, bursty

Global Congestion-Awareness (in Datacenters)

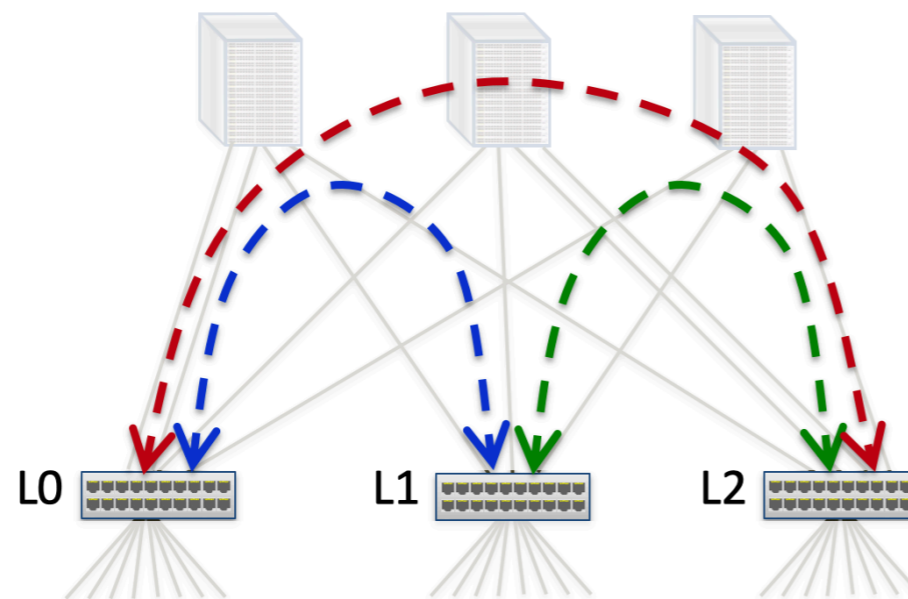
		Datacenter
Opportunity →	{ Latency Topology	microseconds simple, regular
Challenge →	Traffic	volatile, bursty

Key Insight:

Use *extremely* fast, low latency distributed control

CONGA in 1 Slide

1. Leaf switches (top-of-rack) track congestion to other leaves on different paths **in near real-time**
1. Use greedy decisions to minimize bottleneck util



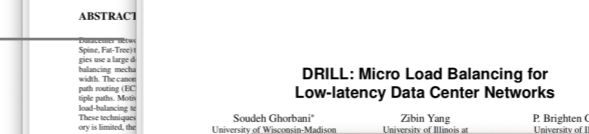
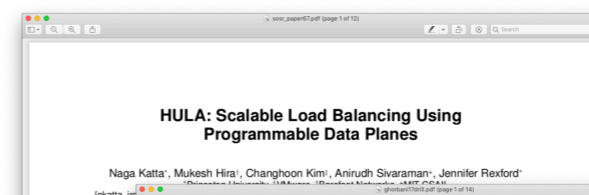
Fast feedback loops
between leaf switches,
directly in dataplane

A large set of papers on programmable data planes aim at improving performance, esp. load balancing

P4-based data-plane load-balancing with better scalability than CONGA

"micro" load balancing, packet-by-packet, can deal with micro-bursts

HULA [SOSR'16]



DRILL [SIGCOMM'17]



stateless, yet congestion-aware load-balancing decision

LetFlow [NSDI'17]

Data plane
programmability for

Performance

Monitoring

Applications offloading

Platforms

for

Data plane

Correctness

programmability

Management

Language-Directed Hardware Design for Network Performance Monitoring

Sonata: Query-Driven Streaming Network Telemetry

Arpit Gupta, Rob Harrison, Marco Canini, Nick Feamster, Jennifer Rexford, Walter Willinger

ABSTRACT
Managing and securing networks requires collecting and analyzing network traffic data in real time. Existing telemetry systems do not allow operators to express the range of queries needed to perform management or scale to large traffic volumes and rates. We present Sonata, an expressive and scalable telemetry system that coordinates joint collection and analysis of network traffic. Sonata provides a declarative interface to express queries for a wide range of common telemetry tasks, to enable real-time execution. Sonata partitions each query across the stream processor and the data plane, running as much of the query as it can on the network switch, at line rate. To optimize the use of limited switch memory, Sonata dynamically refines each query to ensure that available resources focus only on traffic that satisfies the query. Our evaluation shows that Sonata can support a wide range of telemetry tasks while reducing the workload for the stream processor by as much as seven orders of magnitude compared to existing telemetry systems.

CCS CONCEPTS
• Networks → Network monitoring.

KEYWORDS
analytics, programmable switches, stream processing

ACM Reference Format:
Arpit Gupta, Rob Harrison, Marco Canini, Nick Feamster, Jennifer Rexford, and Walter Willinger. 2018. Sonata: Query-Driven Streaming Network Telemetry. In *SIGCOMM '18*. ACM SIGCOMM, 2018. Conference, August 20–25, 2018, Budapest, Hungary. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/3265432.326555>

LossRadar: Fast Detection of Lost Packets in Data Center Networks

FlowRadar: A Better NetFlow for Data Centers

Yuliang Li*, Rui Miao*, Changhoon Kim*, Minlan Yu*

ABSTRACT
Packet loss caused by congestion and network configuration errors is a common problem in data center networks. We also need information about these losses to troubleshoot network issues. Existing solutions are generic and do not take into account the specific characteristics of data center networks. We propose LossRadar, a new way to detect packet loss in data center networks. LossRadar is designed to be easy to integrate into existing network monitoring tools. It uses a small amount of memory and bandwidth to detect packet loss in data center networks. LossRadar is designed to be easy to integrate into existing network monitoring tools. It uses a small amount of memory and bandwidth to detect packet loss in data center networks.

Abstract
NetFlow has been a widely used monitoring tool with a variety of applications. NetFlow maintains an active working set of flows in a hash table that supports flow insertion, collision resolution, and flow removing. This is hard to implement in merchant silicon at data center switches, which has limited per-packet processing time. Therefore, many NetFlow implementations and other monitoring solutions have to sample or select a subset of packets to monitor. In this paper, we observe the need to monitor all the flows without sampling in short time scales. Thus, we design FlowRadar, a new way to maintain flows and their counters that scales to a large number of flows with small memory and bandwidth overhead. The key idea of FlowRadar is to encode per-flow counters with a small memory and constant insertion time at switches, and then to leverage the computing power at the remote collector to perform network-wide decoding and analysis of the flow counters. Our evaluation shows that the memory usage of FlowRadar is close to traditional NetFlow with perfect hashing. With FlowRadar, operators can get better views into their networks as demonstrated by two new monitoring applications we build on top of FlowRadar.

1 Introduction
NetFlow [41] is a widely used monitoring tool for over 20 years, which records the flows (e.g., source IP, destination IP, source port, destination port, and protocol) and their properties (e.g., packet counters, and the flow start and finish times). When a flow finishes after the in-flight timeout, NetFlow exports the corresponding flow records to a remote collector. NetFlow has been used for a variety of monitoring applications such as accounting network usage, capacity planning, troubleshooting, and attack detection. Despite its wide applications, the key problem to im-

Dapper: Data Plane Performance Diagnosis of TCP

Network-Wide Heavy Hitter Detection with Commodity Switches

Rob Harrison, Qizhe Cai, Arpit Gupta, and Jennifer Rexford

ABSTRACT
With more applications running in the data plane, network operators often need to diagnose performance problems in the data plane. We present Dapper, a new way to monitor network-wide heavy hitters. Dapper is designed to be easy to integrate into existing network monitoring tools. It uses a small amount of memory and bandwidth to detect heavy hitters in data center networks. Dapper is designed to be easy to integrate into existing network monitoring tools. It uses a small amount of memory and bandwidth to detect heavy hitters in data center networks.

ABSTRACT
Network operators often need to identify outliers in network traffic, to detect attacks or diagnose performance problems. A common way to detect unusual traffic is to perform “heavy hitter” detection that identifies the top-k flows (or flows exceeding a pre-determined threshold) according to some metric. For example, network operators often track destinations receiving traffic from a large number of distinct sources or TCP incast [4] in real time. In traditional networks, this heavy-hitter detection relies on analyzing packet samples or flow logs [5, 6]. Programmable switches can offer new possibilities for aggregating traffic statistics and identifying large flows directly in the data plane [17, 18, 24, 27]. These

Figure 1: The graph shows the small number of heavy hitters between two major ISPs [12] with different monitoring intervals. Even under high sampling rates, recall quickly diminishes and worsens as the monitoring interval grows.

In-band Network Telemetry (INT)

June 2016

Changhoon Kim, Parag Bhidé, Ed Doe: *Barefoot Networks*
Hugh Holbrook: *Arista*
Anoop Ghanwani: *Dell*
Dan Daly: *Intel*
Mukesh Hira, Bruce Davie: *VMware*

Introduction
Terms
What To Monitor
Switch-level Information
Ingress Information
Egress Information
Buffer Information
Processing INT Headers
INT Header Types
Handling INT Packets
Header Format and Location
INT over any encapsulation
On-the-fly Header Creation
Header Format
Header Location and Format -- INT over Geneve

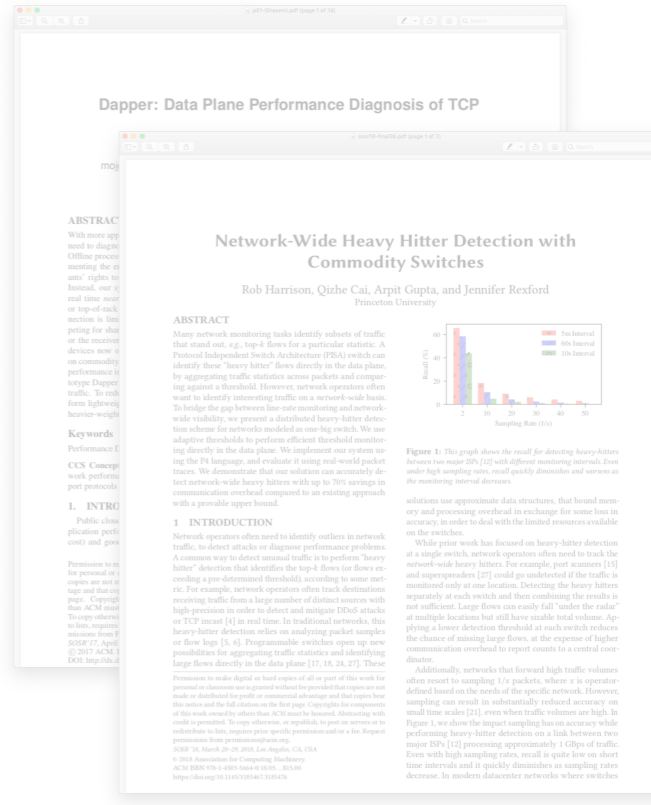
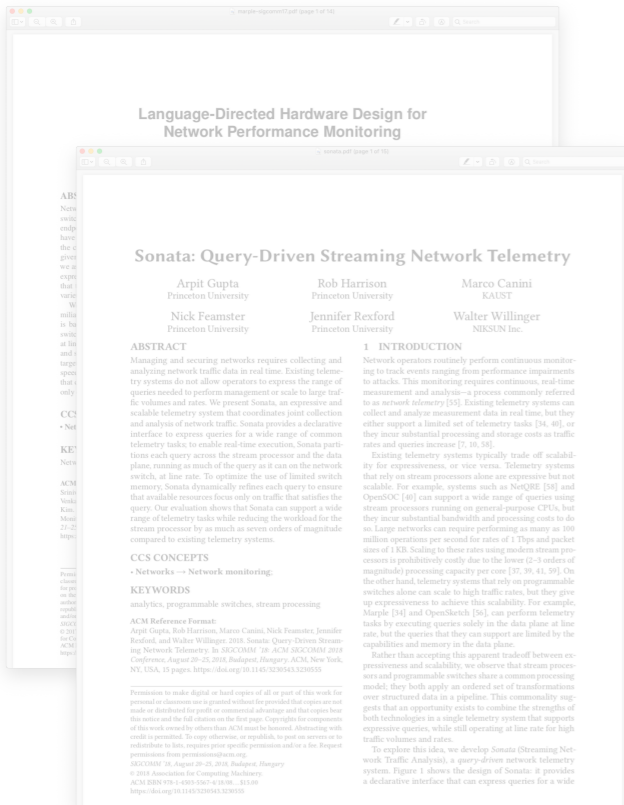
SketchLearn: Relieving User Burdens in Approximate Measurement with Automated Statistical Inference

Elastic Sketch: Adaptive and Fast Network-wide Measurements

One Sketch to Rule Them All: Rethinking Network Flow Monitoring with UnivMon

Zaoxing Liu*, Antonis Manousis*, Gregory Vorsanger*, Vyas Sekar*, Vladimir Braverman*

ABSTRACT
Network management requires accurate estimates of metrics for many applications including traffic engineering (e.g., congestion control, anomaly detection [9]), and forensic analysis [46]. Each such management task requires accurate and timely statistics on different application-level metrics of interest (e.g., the flow size distribution [27], heavy hitters [10], entropy measures [38, 50], or detecting changes in traffic patterns [44]). At a high level, there are two classes of techniques to estimate these metrics of interest. The first class of approaches relies on *general flow monitoring*, typically with some form of packet sampling (e.g., NetFlow [25]). While general flow monitoring is good for coarse-grained visibility, prior work has shown that it provides low accuracy for more fine-grained metrics [30, 31, 43]. These well-known limitations of sampling motivated an alternative class of techniques based on *sketching* or *streaming algorithms*. Here, custom online algorithms and data structures are designed for specific metrics of interest that can yield provable resource-accuracy trade-offs (e.g., [17, 18, 20, 31, 36, 38, 51]). While the body of work in data streaming and sketching has made significant contributions, we argue that this machinery of crafting special-purpose algorithms is unscalable in the long term. As the number of monitoring tasks grows, this entails significant investment in algorithm design and hardware support for new metrics of interest. While recent tools like OpenSketch [47] and SCREAM [41] provide libraries to reduce the implementation effort and offer efficient resource allocation, they do not address the fundamental need to design and operate new custom sketches for each task. Furthermore, as any given point in time the data plane resources have to be committed (a priori) to a specific set of metrics to monitor and will have fundamental blind spots for other metrics that are not currently being tracked. Ideally, we want a monitoring framework that offers both generality by delaying the binding to specific applications of interest but at the same time provides the required fidelity for estimating these metrics. Achieving generality and high fidelity simultaneously has been an elusive goal both in theory [23] (Question 2.4) as well as in practice [45]. In this paper, we present the *UnivMon* (short for Universal Monitoring) framework that can simultaneously achieve both generality and high fidelity across a broad spectrum of monitoring tasks [31, 36, 38, 51]. UnivMon builds on and



In-band Network Telemetry (INT)

June 2016

Changhoon Kim, Parag Bhide, Ed Doe: *Barefoot Networks*

Hugh Holbrook: *Arista*

Anoop Ghanwani: *Dell*

Dan Daly: *Intel*

Mukesh Hira, Bruce Davie: *VMware*

- [Introduction](#)
- [Terms](#)
- [What To Monitor](#)
- [Switch-level Information](#)
- [Ingress Information](#)
- [Egress Information](#)
- [Buffer Information](#)
- [Processing INT Headers](#)
- [INT Header Types](#)
- [Handling INT Packets](#)
- [Header Format and Location](#)
- [INT over any encapsulation](#)
- [On-the-fly Header Creation](#)
- [Header Format](#)
- [Header Location and Format -- INT over Geneve](#)



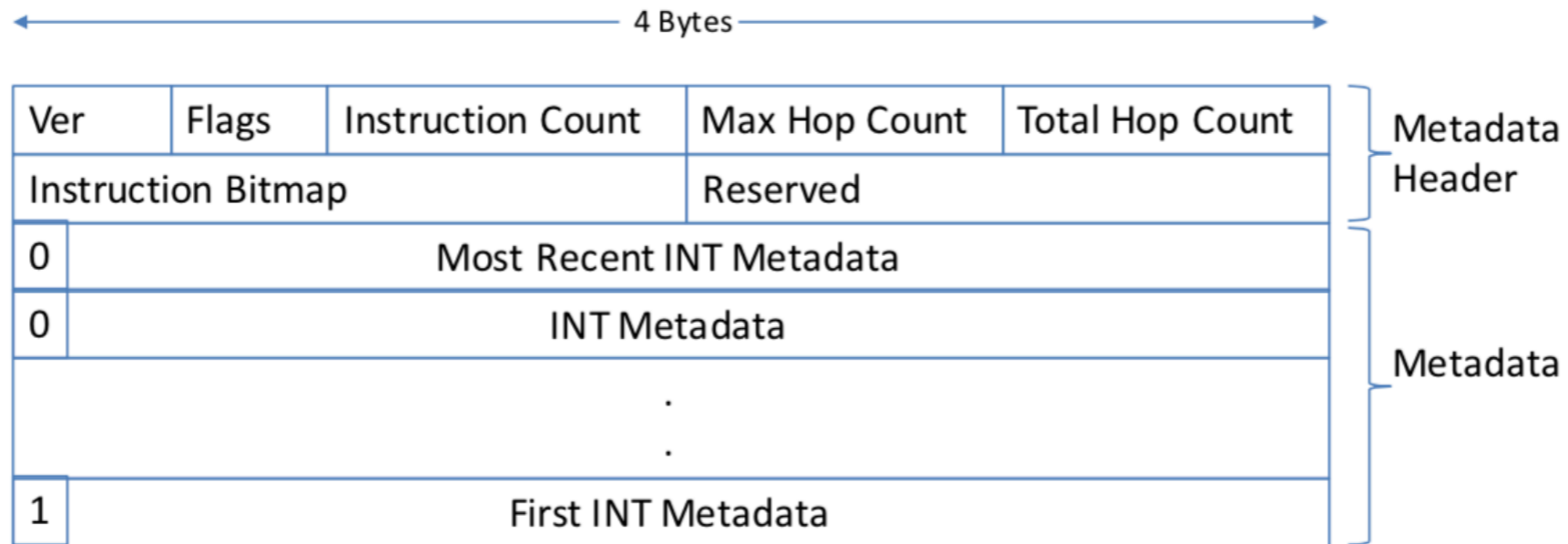
Current monitoring methods are inadequate

- Not fast enough
 - Involve CPU and control planes
 - Network state changes rapidly
- Do not provide end-to-end state
 - Difficult to correlate per-element state with the actual path of a flow

INT : In-band Network Telemetry

- Mechanism for collecting network state in the dataplane
 - As close to **realtime** as possible
 - At current and future **line rates**
 - With a framework that can **adapt** over time
- Examples of network state
 - Switch ID, Ingress Port ID, Egress Port ID
 - Egress Link Utilization
 - Hop Latency
 - Egress Queue Occupancy
 - Egress Queue Congestion Status
 -

INT Header Format



Source: In-band Network Telemetry, Mukesh Hira and Naga Katta, 2015

INT using P4

- P4 enables flexible packet parsing and modification for INT
- P4 allows INT to adapt to
 - Any Encapsulation format
 - Any State required to be collected
 - Any feature, protocol – current and future

INT : P4 Code Snippet

Exact-match Table Definition

```
table int_inst {
  reads {
    int_header.instruction_mask : exact;
  }
  actions {
    int_set_header_i0;
    int_set_header_i1;
    int_set_header_i2;
    int_set_header_i3;
    .....
  }
}
```

Action Definitions

```
action int_set_header_i0() {
}
action int_set_header_i1() {
  int_set_header_3();
}
action int_set_header_i2() {
  int_set_header_2();
}
action int_set_header_i3() {
  int_set_header_3();
  int_set_header_2();
}
.....
```

HULA: INT + Flowlet routing

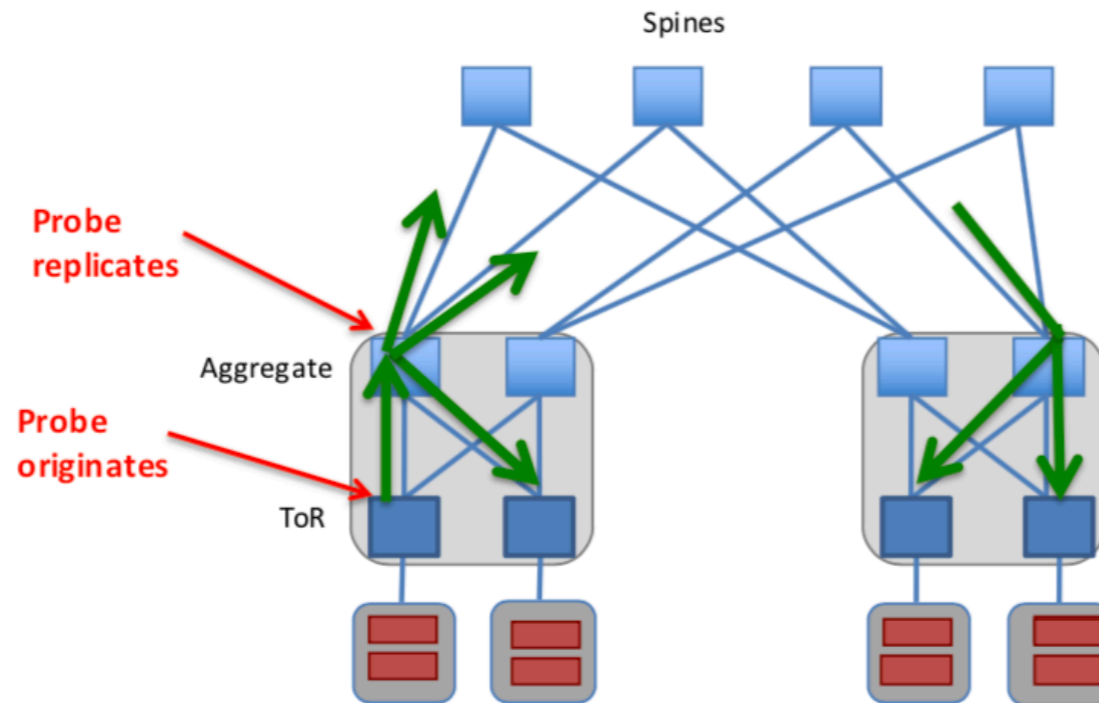
1. Periodic INT probes

- disseminate path utilization to switches

2. Flowlet detection and path selection

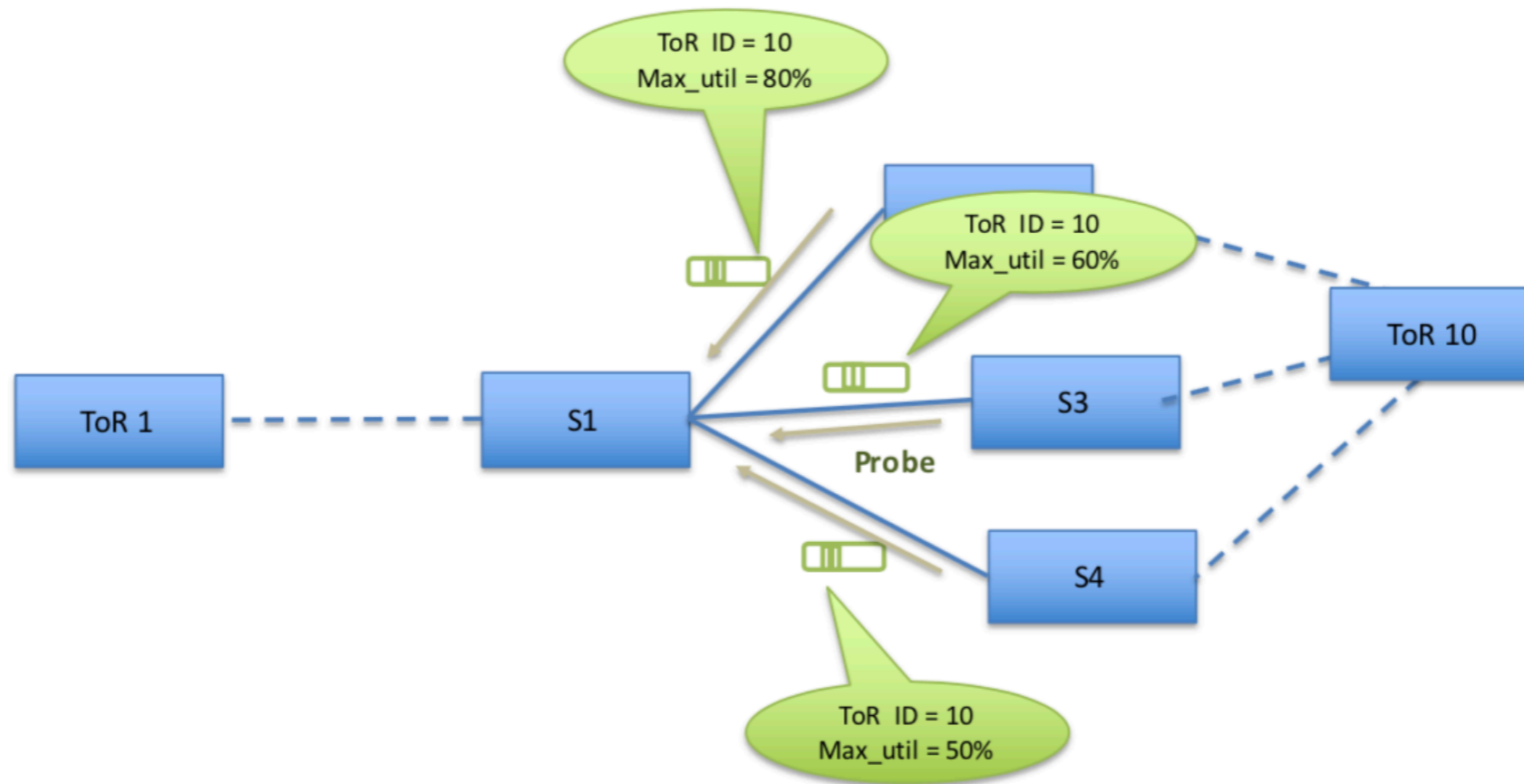
- happens at **all** switches
- hop-by-hop adaptive routing

INT probes traverse multiple paths



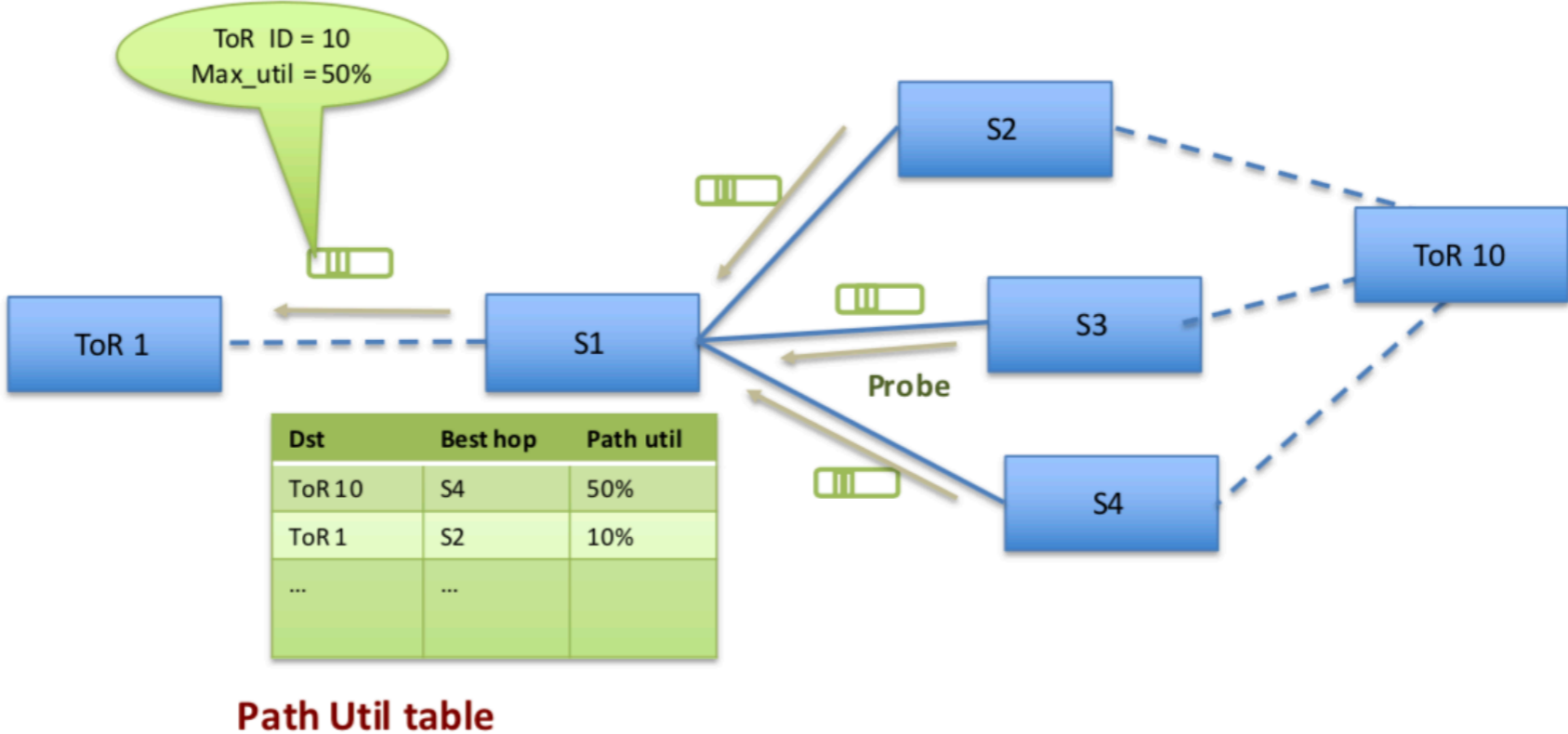
Source: In-band Network Telemetry, Mukesh Hira and Naga Katta, 2015

Probes carry path utilization



Source: In-band Network Telemetry, Mukesh Hira and Naga Katta, 2015

Probes update switch state



Source: In-band Network Telemetry, Mukesh Hira and Naga Katta, 2015

Summary

- INT provides real-time network state directly in the dataplane
 - Scales to arbitrarily large networks
 - Scales to current and future link speeds
 - Can adapt to any network, any encap, any application
- Knowledge of real-time network state opens up new possibilities
 - Enhanced monitoring and troubleshooting
 - Network-state aware routing
 - ...

Language-Directed Hardware Design for Network Performance Monitoring

Sonata: Query-Driven Streaming Network Telemetry

Arpit Gupta, Rob Harrison, Marco Canini, Nick Feamster, Jennifer Rexford, Walter Willinger

ABSTRACT
Managing and securing networks requires collecting and analyzing network traffic data in real time. Existing telemetry systems do not allow operators to express the range of queries needed to perform management or scale to large traffic volumes and rates. We present Sonata, an expressive and scalable telemetry system that coordinates joint collection and analysis of network traffic. Sonata provides a declarative interface to express queries for a wide range of common telemetry tasks, to enable real-time execution. Sonata partitions each query across the stream processor and the data plane, running as much of the query as it can on the network switch, at line rate. To optimize the use of limited switch memory, Sonata dynamically refines each query to ensure that available resources focus only on traffic that satisfies the query. Our evaluation shows that Sonata can support a wide range of telemetry tasks while reducing the workload for the stream processor by as much as seven orders of magnitude compared to existing telemetry systems.

CCS CONCEPTS
• Networks → Network monitoring.

KEYWORDS
analytics, programmable switches, stream processing

ACM Reference Format:
Arpit Gupta, Rob Harrison, Marco Canini, Nick Feamster, Jennifer Rexford, and Walter Willinger. 2018. Sonata: Query-Driven Streaming Network Telemetry. In *SIGCOMM '18*. ACM, SIGCOMM 2018 Conference, August 20–25, 2018, Budapest, Hungary. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/3266453.3266555>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions.acm.org.
SIGCOMM '18, August 20–25, 2018, Budapest, Hungary.
© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-6649-0/18/08...\$15.00
<https://doi.org/10.1145/3266453.3266555>

LossRadar: Fast Detection of Lost Packets in Data Center Networks

FlowRadar: A Better NetFlow for Data Centers

Yuliang Li*, Rui Miao*, Changhoon Kim*, Minlan Yu*
*University of Southern California *Barfoot Networks

ABSTRACT
Packet loss caused by congestion and hardware failures in data center networks is a common problem. Existing network-wide detection of loss is not accurate enough. We propose FlowRadar, a new way to detect loss in data center networks. FlowRadar is easy to implement and does not require any changes to the network. It is designed to be used in data center networks. FlowRadar is easy to implement and does not require any changes to the network. It is designed to be used in data center networks.

Abstract
NetFlow has been a widely used monitoring tool with a variety of applications. NetFlow maintains an active working set of flows in a hash table that supports flow insertion, collision resolution, and flow removing. This is hard to implement in merchant silicon at data center switches, which has limited per-packet processing time. Therefore, many NetFlow implementations and other monitoring solutions have to sample or select a subset of packets to monitor. In this paper, we observe the need to monitor all the flows without sampling in short time scales. Thus, we design FlowRadar, a new way to maintain flows and their counters that scales to a large number of flows with small memory and bandwidth overhead. The key idea of FlowRadar is to encode per-flow counters with a small memory and constant insertion time at switches, and then to leverage the computing power at the remote collector to perform network-wide decoding and analysis of the flow counters. Our evaluation shows that the memory usage of FlowRadar is close to traditional NetFlow with perfect hashing. With FlowRadar, operators can get better views into their networks as demonstrated by two new monitoring applications we build on top of FlowRadar.

1 Introduction
NetFlow [4] is a widely used monitoring tool for over 20 years, which records the flows (e.g., source IP, destination IP, source port, destination port, and protocol) and their properties (e.g., packet counters, and the flow start- and finish times). When a flow finishes after the in-flight timeout, NetFlow exports the corresponding flow records to a remote collector. NetFlow has been used for a variety of monitoring applications such as accounting network usage, capacity planning, troubleshooting, and attack detection.

Despite its wide applications, the key problem to im-

Dapper: Data Plane Performance Diagnosis of TCP

Network-Wide Heavy Hitter Detection with Commodity Switches

Rob Harrison, Qizhe Cai, Arpit Gupta, and Jennifer Rexford
Princeton University

ABSTRACT
With more applications running in the data plane, network operators often need to track the network-wide heavy hitters to identify these "heavy hitter" flows directly in the data plane, by aggregating traffic statistics across packets and comparing against a threshold. However, network operators often want to identify interesting traffic on a network-wide basis. To bridge the gap between line-rate monitoring and network-wide visibility, we present a distributed heavy-hitter detection scheme for network-wide one-hop switches. We use adaptive thresholds to perform efficient threshold monitoring directly in the data plane. We implement our system using the P4 language, and evaluate it using real-world packet traces. We demonstrate that our solution can accurately detect network-wide heavy hitters with up to 70% savings in communication overhead compared to an existing approach with a provable upper bound.

1 INTRODUCTION
Network operators often need to identify outliers in network traffic, to detect attacks or diagnose performance problems. A common way to detect unusual traffic is to perform "heavy hitter" detection that identifies the top-k flows (or flows exceeding a pre-determined threshold) according to some metric. For example, network operators often track destinations receiving traffic from a large number of distinct sources with high precision in order to detect and mitigate DDoS attacks or TCP incursions [4] in real time. In traditional networks, this heavy-hitter detection relies on analyzing packet samples or flow logs [5, 6]. Programmable switches offer new possibilities for aggregating traffic statistics and identifying large flows directly in the data plane [17, 18, 24, 27]. These

Figure 1: The graph shows the small number of detecting heavy hitters between two major ISPs [12] with different monitoring intervals. Even under high sampling rates, recall quickly diminishes and worsens as the monitoring interval grows.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions.acm.org.
SIGCOMM '18, August 20–25, 2018, Budapest, Hungary.
© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-6649-0/18/08...\$15.00
<https://doi.org/10.1145/3266453.3266555>

In-band Network Telemetry (INT)

June 2016

Changhoon Kim, Parag Bhidé, Ed Doe: *Barfoot Networks*
Hugh Holbrook: *Arista*
Anoop Ghanwani: *Dell*
Dan Daly: *Intel*
Mukesh Hira, Bruce Davie: *VMware*

Introduction
Terms
What To Monitor
Switch-level Information
Ingress Information
Egress Information
Buffer Information
Processing INT Headers
INT Header Types
Handling INT Packets
Header Format and Location
INT over any encapsulation
On-the-fly Header Creation
Header Format
Header Location and Format – INT over Geneve

SketchLearn: Relieving User Burdens in Approximate Measurement with Automated Statistical Inference

Elastic Sketch: Adaptive and Fast Network-wide Measurements

State Key Laboratory of Information Security, Institute of Software, Chinese Academy of Sciences, Beijing, China

ABSTRACT
Network management requires accurate estimates of metrics for many applications including traffic engineering (e.g., heavy hitters, anomaly detection [9], and forensic analysis [46]). Each such management task requires accurate and timely statistics on different application-level metrics of interest (e.g., the flow size distribution [27], heavy hitters [10], entropy measures [38, 50], or detecting changes in traffic patterns [44]). At a high level, there are two classes of techniques to estimate these metrics of interest. The first class of approaches relies on *general flow monitoring*, typically with some form of packet sampling (e.g., NetFlow [25]). While general flow monitoring is good for coarse-grained visibility, prior work has shown that it provides low accuracy for more fine-grained metrics [30, 31, 43]. These well-known limitations of sampling motivated an alternative class of techniques based on *sketching* or *streaming algorithms*. Here, custom online algorithms and data structures are designed for specific metrics of interest that can yield provable resource-accuracy trade-offs (e.g., [17, 18, 20, 31, 36, 38, 45]).

One Sketch to Rule Them All: Rethinking Network Flow Monitoring with UnivMon

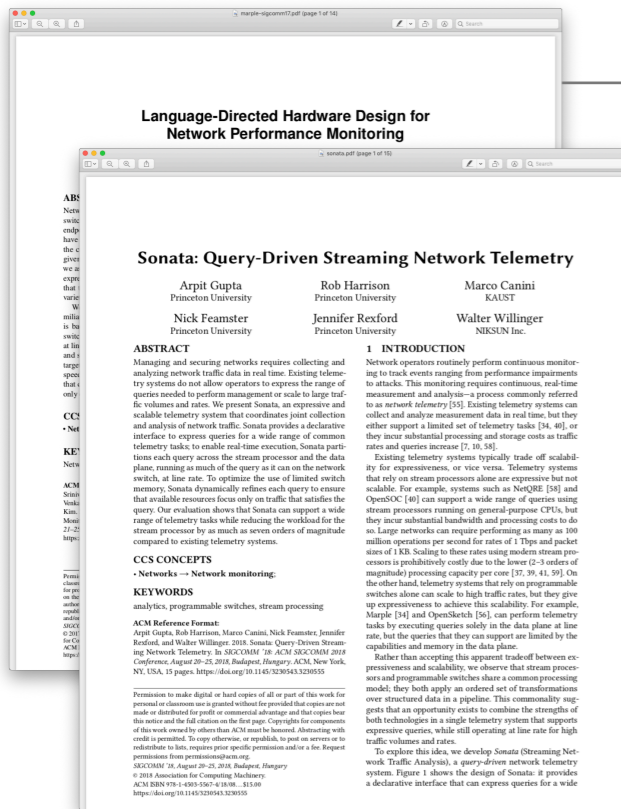
Zaoxing Liu*, Antonis Manousis*, Gregory Vorsanger*, Vyas Sekar*, Vladimir Braverman*
*Johns Hopkins University *Carnegie Mellon University

1 Introduction
Network management is multi-faceted and encompasses a range of tasks including traffic engineering [11, 32], attack and anomaly detection [9], and forensic analysis [46]. Each such management task requires accurate and timely statistics on different application-level metrics of interest (e.g., the flow size distribution [27], heavy hitters [10], entropy measures [38, 50], or detecting changes in traffic patterns [44]). At a high level, there are two classes of techniques to estimate these metrics of interest. The first class of approaches relies on *general flow monitoring*, typically with some form of packet sampling (e.g., NetFlow [25]). While general flow monitoring is good for coarse-grained visibility, prior work has shown that it provides low accuracy for more fine-grained metrics [30, 31, 43]. These well-known limitations of sampling motivated an alternative class of techniques based on *sketching* or *streaming algorithms*. Here, custom online algorithms and data structures are designed for specific metrics of interest that can yield provable resource-accuracy trade-offs (e.g., [17, 18, 20, 31, 36, 38, 45]).

CCS CONCEPTS
• Networks → Network monitoring; Network measurement.

Keywords
Flow Monitoring, Sketching, Streaming Algorithms

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions.acm.org.
SIGCOMM '18, August 22–26, 2018, Philadelphia, PA, USA.
© 2018 ACM. New York, NY, USA, 12 pages.
DOI: <https://doi.org/10.1145/3292487.3292496>



MARPLE [SIGCOMM'17]

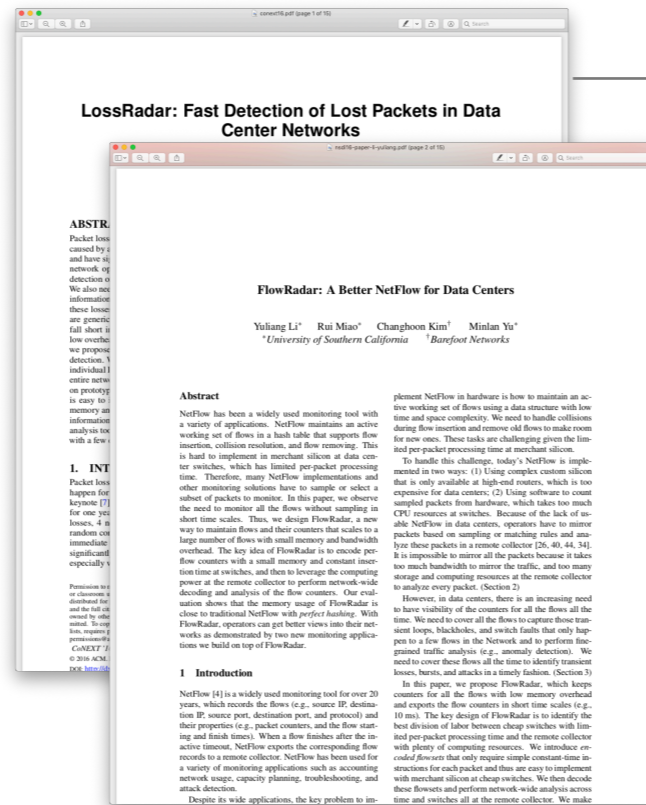
SONATA [SIGCOMM'18]

Both papers enable operators to express **monitoring queries**

```
result = filter(pktstream, qid == Q and switch == S
               and t_out - t_in > 1ms)
returns a stream of packets experiencing high queuing latencies
```

A compiler then compiles these queries to: switch programs + control code

The two papers differ among others in the types of queries they support



LossRadar [CoNEXT'16]

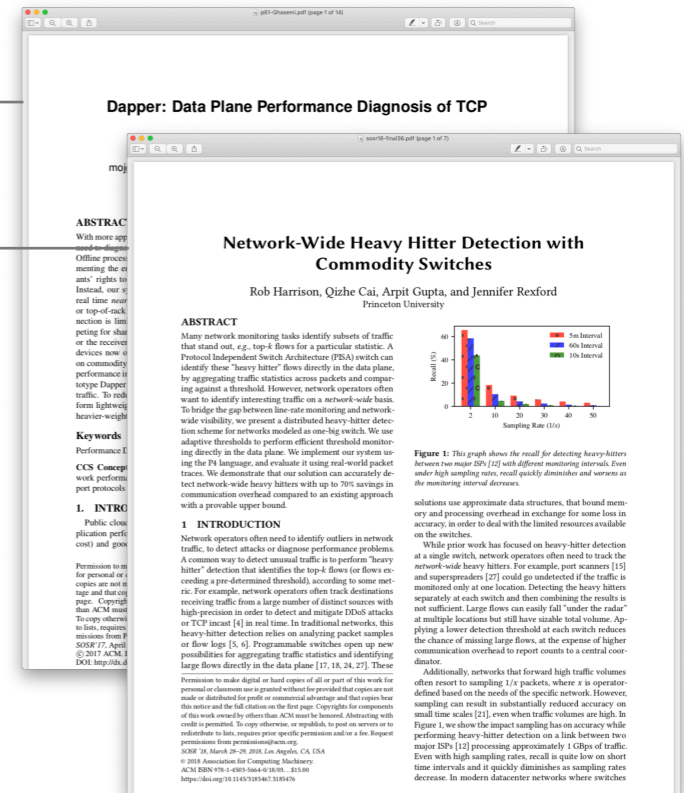
FlowRadar [NSDI'16]

Develop techniques and tools to monitor *all flows* by

- relying on in-switch data structures (Bloom Filters) and
- decoding them at the controller-level

DAPPER [SOSR'17]

Network-Wide HH [SOSR'18]



Develop P4-based detection mechanisms to

- diagnose TCP performance issue (e.g. small receiver buffers)
- heavy-hitter (e.g. port scanners, superspreader, DDoS)

Introduce techniques to make sketch-based monitoring more practical (by making sketches adaptive or "universal")

SketchLearn [SIGCOMM'18]

Elastic Sketch [SIGCOMM'18]

UnivMon [SIGCOMM'16]



Data plane
programmability

for

Performance
Monitoring

Applications offloading

Platforms

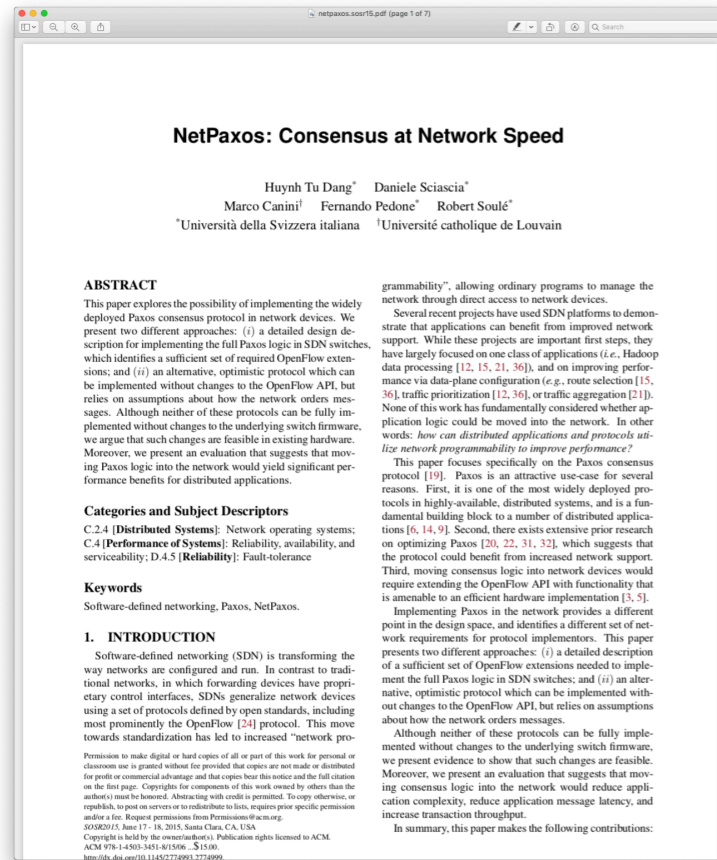
for

Data plane
programmability

Correctness

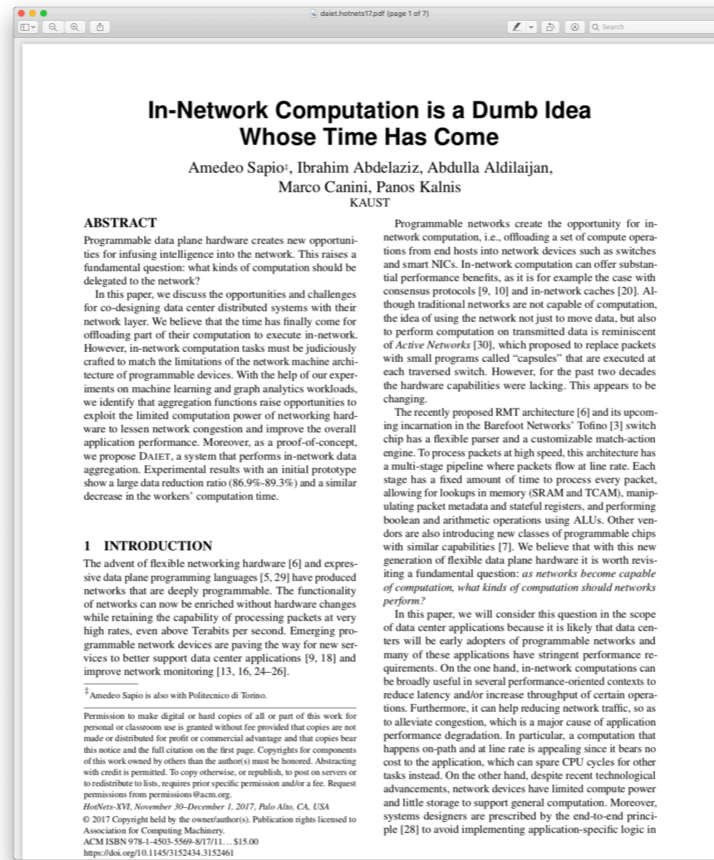
Management

[SOSR'15]



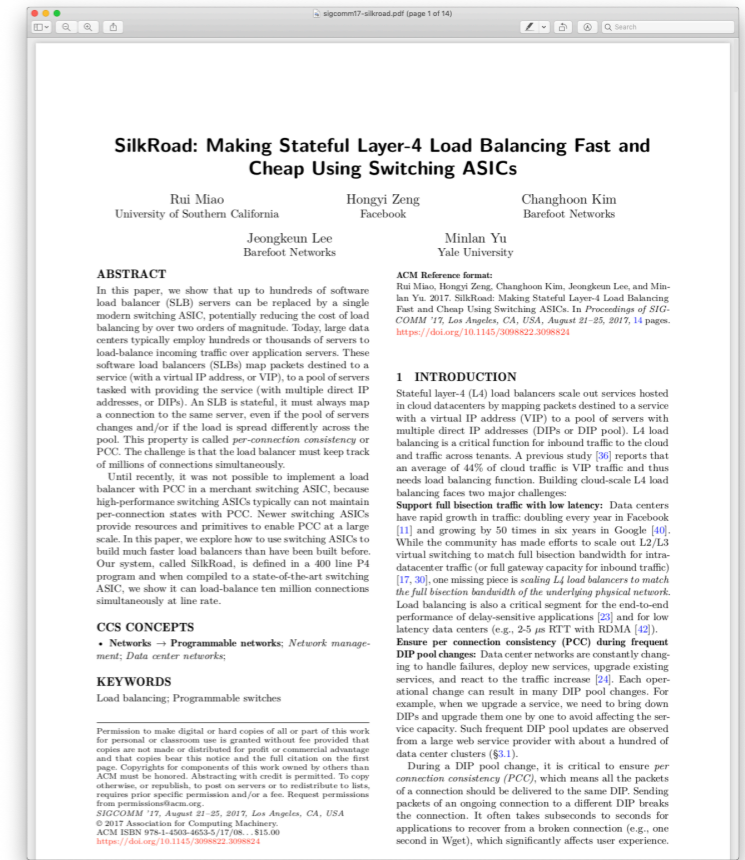
Consensus at network speed

[HotNets'17]



In-Network Aggregation
(e.g., for MapReduce, graph analytics, ML)

[SIGCOMM'17]



Stateful layer-4 load balancers

+ NetCache [SOSP'17], NetChain [NSDI'18]

Data plane
programmability

for

Performance
Monitoring
Applications offloading

Platforms

for

Data plane
programmability

Correctness
Management

"Data-plane" programmability goes beyond
switch programmability (or P4 for that matter)

Offloading...

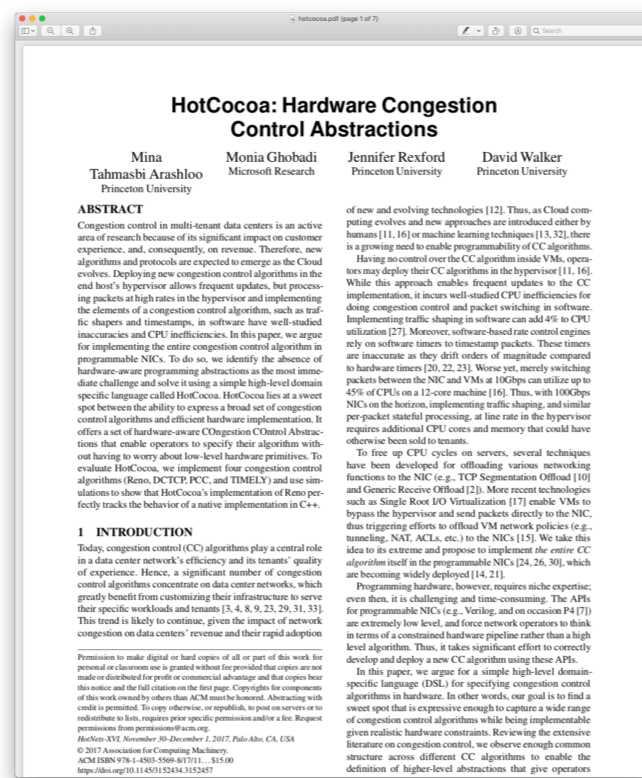
... to FPGA-based SmartNICs

host networking

congestion control



[NSDI'18]



HotCocoa: Hardware Congestion Control Abstractions

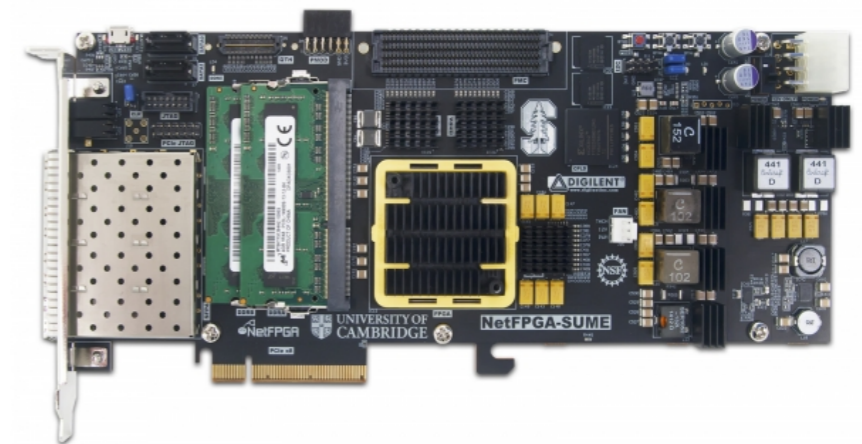
Mina Tahmasbi Arashloo, Monia Ghobadi, Jennifer Rexford, David Walker

ABSTRACT Congestion control in multi-tenant data centers is an active area of research because of its significant impact on customer experience, and, consequently, on revenue. Therefore, new algorithms and protocols are expected to emerge as the Cloud evolves. Deploying new congestion control algorithms in the end host's hypervisor allows frequent updates, but processing packets at high rates in the hypervisor and implementing the elements of a congestion control algorithm, such as traffic shapers and timestamps, in software have well-studied inaccuracies and CPU inefficiencies. In this paper, we argue for implementing the entire congestion control algorithm in programmable NICs. To do so, we identify the absence of hardware-aware programming abstractions as the most immediate challenge and solve it using a simple high-level domain specific language called HotCocoa. HotCocoa lies at a sweet spot between the ability to express a broad set of congestion control algorithms and efficient hardware implementation. It offers a set of hardware-aware Congestion Control Abstractions that enable operators to specify their algorithm without having to worry about low-level hardware primitives. To evaluate HotCocoa, we implement four congestion control algorithms (Reno, DCTCP, PCC, and TIMELY) and use simulations to show that HotCocoa's implementation of Reno perfectly tracks the behavior of a native implementation in C++.

1 INTRODUCTION Today, congestion control (CC) algorithms play a central role in a data center network's efficiency and its tenants' quality of experience. Hence, a significant number of congestion control algorithms concentrate on data center networks, which greatly benefit from customizing their infrastructure to serve their specific workloads and tenants [3, 4, 8, 9, 23, 29, 31, 33]. This trend is likely to continue, given the impact of network congestion on data centers' revenue and their rapid adoption

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org. HotNets '17, November 29-December 1, 2017, Palo Alto, CA, USA. © 2017 Association for Computing Machinery. ACM ISBN 978-1-4503-5969-8/17/11...\$15.00 http://doi.org/10.1145/3142434.3142487

[HotNets'17]

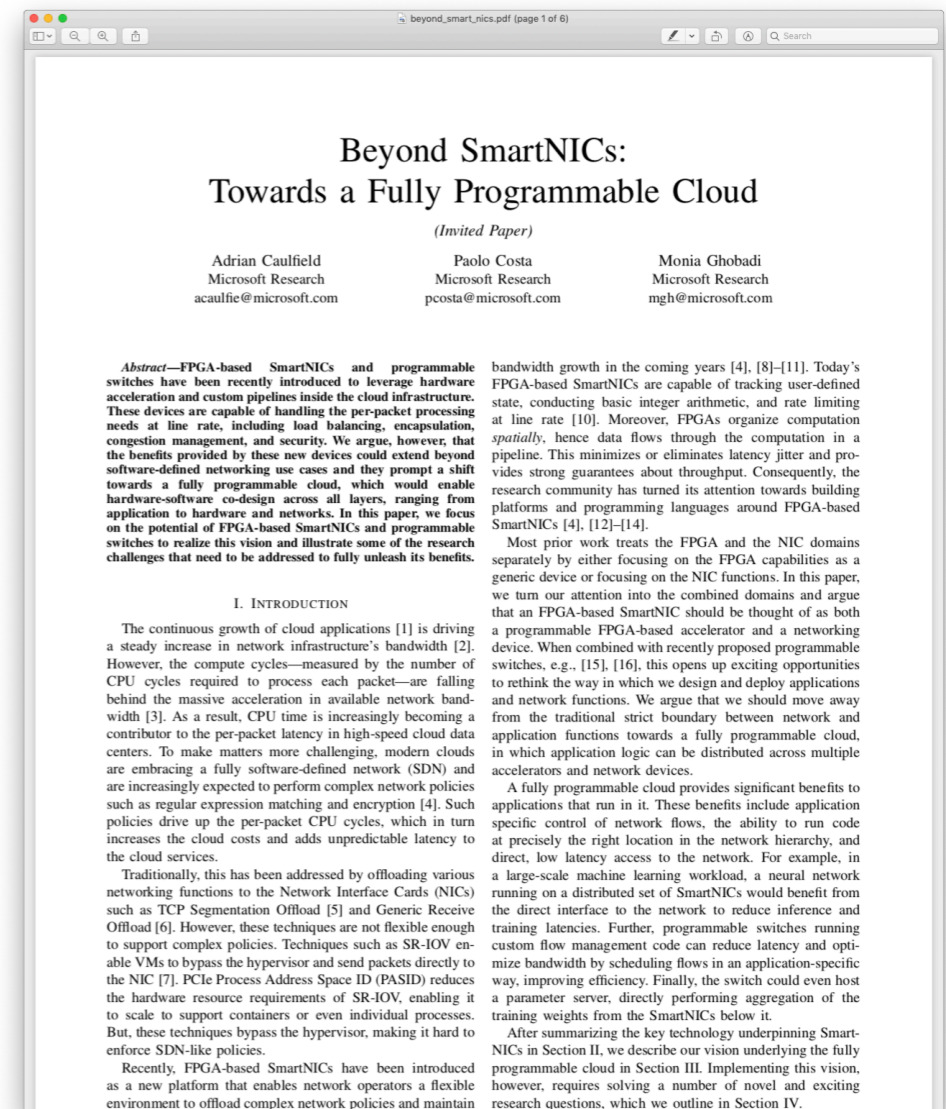


NetFPGA SUME board

Host-based programmability + SmartNICs + programmable switches = fully programmable platforms

Big question is

How to combine them best?



IEEE International Conference on High Performance Switching and Routing, 2018

Data plane
programmability

for

Performance
Monitoring
Applications offloading

Platforms

for

Data plane
programmability

Correctness

Management

So you've a programmable networks... How do you make sure that it works as it should?!

p4v: Practical Verification for Programmable Data Planes*

Jed Liu, William Hallahan, Cole Schlesinger, Milad Sharif, Jeongkeun Lee
Barefoot Networks, Yale University, Barefoot Networks, Barefoot Networks, Barefoot Networks
Ithaca, NY, USA, New Haven, CT, USA, Santa Clara, CA, USA, Santa Clara, CA, USA, Santa Clara, CA, USA

Robert Soulé, Han Wang, Călin Cașcaval, Nick McKeown, Nate Foster
University of Lugano, Barefoot Networks, Barefoot Networks, Stanford University, Cornell University
Lugano, Switzerland, Santa Clara, CA, USA, Santa Clara, CA, USA, Stanford, CA, USA, Ithaca, NY, USA

ABSTRACT
We present the design and implementation of p4v, a practical tool for verifying data planes described using the P4 programming language. The design of p4v is based on classic verification techniques but adds several key innovations including a novel mechanism for incorporating assumptions about the control plane and domain-specific optimizations which are needed to scale to large programs. We present case studies showing that p4v verifies important properties and finds bugs in real-world programs. We conduct experiments to quantify the scalability of p4v on a wide range of additional examples. We show that with just a few hundred lines of control-plane annotations, p4v is able to verify critical safety properties for switch.p4, a program that implements the functionality of on a modern data center switch, in under three minutes.

CCS CONCEPTS
• **Networks** → *Programming interfaces*; • **Software and its engineering** → *Software verification*;

KEYWORDS
Programmable data planes, P4, verification.

ACM Reference Format:
Jed Liu, William Hallahan, Cole Schlesinger, Milad Sharif, Jeongkeun Lee, Robert Soulé, Han Wang, Călin Cașcaval, Nick McKeown, and Nate Foster. 2018. p4v: Practical Verification for Programmable Data Planes. In *SIGCOMM '18: SIGCOMM 2018, August 20–25, 2018, Budapest, Hungary*. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3230543.3230582>

*This version fixes small typographic errors in the official published version.

SIGCOMM '18, August 20–25, 2018, Budapest, Hungary
© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.
This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *SIGCOMM '18: SIGCOMM 2018, August 20–25, 2018, Budapest, Hungary*, <https://doi.org/10.1145/3230543.3230582>.

[SIGCOMM'18]

Debugging P4 programs with Vera
Radu Stoescu, Dragos Dumitrescu, Matei Popovici, Lorina Negreanu, Costin Raiciu
University Politehnica of Bucharest
firstname.lastname@cs.pub.ro

ABSTRACT
We present Vera, a tool that verifies P4 programs using symbolic execution. Vera automatically uncovers a number of common bugs including parsing/deparsing errors, invalid memory accesses, loops and tunneling errors, among others. Vera can also be used to verify user-specified properties in a novel language we call NetCTL.
To enable scalable, exhaustive verification of P4 program snapshots, Vera automatically generates all valid header layouts and uses a novel data-structure for match-action processing optimized for verification. These techniques allow Vera to scale very well: it only takes between 5s-15s to track the execution of a purely symbolic packet in the latest P4 program currently available (kELCC).
model updates in milliseconds. Vera can update a single concrete dataplane at once by all to insert symbolic table entries; the highlights possible control plane errors.
We have used Vera to analyze many the P4 tutorials, P4 programs in the the switch code from <https://p4.org/>, bugs in each of them in seconds/min

CCS CONCEPTS
• **General and reference** → *Verify Network reliability*; • **Programmab**

1 INTRODUCTION
Programmable network dataplanes by P4 [2] promise to help networks application demands. On the downsi

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permission from permissions@acm.org.
© 2018 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-5267-4/18/08...\$15.00
<https://doi.org/10.1145/3230543.3230582>

Verification of P4 Programs in Feasible Time using Assertions
Miguel Neves, Lucas Freire, Alberto Schaeffer-Filho, Marinho Barcellos
Institute of Informatics
UFRRG

ABSTRACT
Recent trends in software-defined networking have extended network programmability to the data plane. Unfortunately, the chance of introducing bugs increases significantly. Verification can help prevent bugs by assuring that the program does not violate its requirements. Although research on the verification of P4 programs is very active, we still need tools to make easier for programmers to express properties and to rapidly verify complex invariants. In this paper, we leverage assertions and symbolic execution to propose a more general P4 verification approach. Developers annotate P4 programs with assertions expressing general network correctness properties; the result is transformed into C models and all possible paths symbolically executed. We implement a prototype, and use it to show the feasibility of the verification approach. Because symbolic execution does not scale well, we investigate a set of techniques to speed up the process for the specific case of P4 programs. We use the prototype implemented to show the gains provided by three speed up techniques (use of constraints, program slicing, parallelization), and experiment with different compiler optimization choices. We show our tool can uncover a broad range of bugs, and can do it in less than a minute considering various P4 applications.

CCS CONCEPTS
• **Networks** → *Programmable networks*; • **Software and its engineering** → *Software verification and validation*;

KEYWORDS
P4; Verification; Programmable Data Planes

ACM Reference Format:
Miguel Neves, Lucas Freire, Alberto Schaeffer-Filho, Marinho Barcellos. 2018. Verification of P4 Programs in Feasible Time using Assertions. In *The 18th International Conference on Emerging Networking Experiments and Technologies (CoNEXT '18), December 4–7, 2018, Heraklion, Greece*. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3281411.3281421>

1 INTRODUCTION
Data plane programmability allows operators to quickly deploy new protocols and develop network services. Through programming languages such as P4 [2], it is possible to specify a few instructions

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permission from permissions@acm.org.
© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-5267-4/18/08...\$15.00
<https://doi.org/10.1145/3281411.3281421>

To network functionality can introduce bugs that may cause great damage. Recently, faulty routers in two airline networks have grounded airplanes for days (for both Delta and Southwest Airlines), showing just how disruptive the effects of incorrect network behavior can be. Given the momentum behind programmable networks, we expect such faults and many others will cripple programmable networks.
In this paper, we argue that dataplane programs should be verified before deployment to enable safe operation. We present Vera, a verification tool that enables debugging of P4 programs both before deployment and at runtime. At its core, Vera translates P4 to SEFL, a network language designed for verification, and relies on symbolic execution with Symmet

which packet headers should be manipulated, and how, by different forwarding devices in the infrastructure. Despite the flexibility, this paradigm also increases the chance of introducing bugs into the network due to incorrect implementations.
Testing/debugging, verification and enforcement are complementary approaches that can help solve this problem. During development, data plane programs can be debugged and tested, providing a wide range of inputs and checking if the corresponding outputs match the expected behavior. Verification, on its turn, can be used on programs to find bugs that would violate any of the properties stated by their requirements, including bugs that are hard to reproduce in testing. Lastly, with enforcement, the data plane can be monitored during execution to trap and block actions that would result in property violations.
In this paper, we focus on *verification*: we propose an approach to model and check (at compile time) general security and correctness properties of P4 programs, and implement it in a tool that provides network verification in feasible time. Several approaches have been developed to check if a given fixed-function (non-P4) data plane satisfies a set of intended properties [8, 25, 29, 32]. Moreover, verifying P4-programmed data planes is an active area of research, with recent papers proposing verification techniques based on SMT solving [24, 27] and custom symbolic execution [33]¹. In contrast, this work shows how to efficiently verify P4 programs leveraging a popular, off-the-shelf symbolic execution engine [4].
We propose an expressive assertion language (highly influenced by P4) that enables programmers to specify their intended properties by annotating their P4 code. Once annotated, a program is symbolically executed, with assertions being checked while all its paths are traversed. Given that the time taken to perform the symbolic execution grows exponentially with the program complexity, we show how a variety of speed up techniques can be employed to reduce the verification time and number of executed instructions. These techniques consist of using annotations in code to constrain the paths to be traversed according to properties and/or protocols of interest, program slicing to reduce the complexity of the model under verification, and parallelization of symbolic execution. Besides, we experiment with code optimization features offered by current compilers.
To evaluate our approach, we built a prototype using KLEE [4] and the P4 Reference Compiler [20] for the current language version, P4₁₆. We applied it to four real P4 applications collected from the literature: Switch [21], NetPass [5], Dapper [11], and DCp4 [31]. Our results show that the proposed verification process can uncover a broad range of bugs either in the data plane program itself or in its control plane configuration. A detailed performance analysis also shows that, although the verification time grows exponentially with

¹[24, 33] were independently developed at the same time as this work.

[SIGCOMM'18]

[CoNEXT'18]

So you've a programmable networks... How do you make sure that it works as it should?!

p4v.pdf (page 1 of 14)

p4v: Practical Verification for Programmable Data Planes*

Jed Liu
Barefoot Networks
Ithaca, NY, USA

William Hallahan
Yale University
New Haven, CT, USA

Cole Schlesinger
Barefoot Networks
Santa Clara, CA, USA

Milad Sharif
Barefoot Networks
Santa Clara, CA, USA

Jeongkeun Lee
Barefoot Networks
Santa Clara, CA, USA

Robert Soulé
University of Lugano
Lugano, Switzerland

Han Wang
Barefoot Networks
Santa Clara, CA, USA

Călin Cașcaval
Barefoot Networks
Santa Clara, CA, USA

Nick McKeown
Stanford University
Stanford, CA, USA

Nate Foster
Cornell University
Ithaca, NY, USA

ABSTRACT

We present the design and implementation of p4v, a practical tool for verifying data planes described using the P4 programming language. The design of p4v is based on classic verification techniques but adds several key innovations including a novel mechanism for incorporating assumptions about the control plane and domain-specific optimizations which are needed to scale to large programs. We present case studies showing that p4v verifies important properties and finds bugs in real-world programs. We conduct experiments to quantify the scalability of p4v on a wide range of additional examples. We show that with just a few hundred lines of control-plane annotations, p4v is able to verify critical safety properties for swi tch.p4, a program that implements the functionality of on a modern data center switch, in under three minutes.

CCS CONCEPTS

• Networks → Programming interfaces; • Software and its engineering → Software verification;

KEYWORDS

Programmable data planes, P4, verification.

ACM Reference Format:

Jed Liu, William Hallahan, Cole Schlesinger, Milad Sharif, Jeongkeun Lee, Robert Soulé, Han Wang, Călin Cașcaval, Nick McKeown, and Nate Foster. 2018. p4v: Practical Verification for Programmable Data Planes. In SIGCOMM '18: SIGCOMM 2018, August 20–25, 2018, Budapest, Hungary. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3230543.3230582>

*This version fixes small typographic errors in the official published version.

SIGCOMM '18, August 20–25, 2018, Budapest, Hungary
© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.
This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in SIGCOMM '18: SIGCOMM 2018, August 20–25, 2018, Budapest, Hungary, <https://doi.org/10.1145/3230543.3230582>.

[SIGCOMM'18]

vera.pdf (page 1 of 15)

Debugging P4 programs with Vera

Radu Stoescu, Dragos Dumitrescu, Matei Popovici, Lorina Negreanu, Costin Raiciu
University Politehnica of Bucharest
firstname.lastname@cs.pub.ro

ABSTRACT

We present Vera, a tool that verifies P4 programs using symbolic execution. Vera automatically uncovers a number of common bugs including parsing/deparsing errors, invalid memory accesses, loops and tunneling errors, among others. Vera can also be used to verify user-specified properties in a novel language we call NetCTL.

To enable scalable, exhaustive verification of P4 program snapshots, Vera automatically generates all valid header layouts and uses a novel data-structure for match-action processing optimized for verification. These techniques allow Vera to scale very well: it only takes between 5s-15s to track the execution of a purely symbolic packet in the latest P4 program currently available (sKLOC), model updates in milliseconds, Vera re-compile dataplanes at once by all to insert symbolic table entries; the highlights possible control plane errors.

We have used Vera to analyze many the P4 tutorials, P4 programs in the switch code from <https://p4.org>, bugs in each of them in seconds/min.

CCS CONCEPTS

• General and reference → Verify Network reliability; Programmab

1 INTRODUCTION

Programmable network dataplanes by P4 [2] promise to help networks application demands. On the downsi

Permission to make digital or hard copies of personal or classroom use is granted within are not made or distributed for profit or com copies bear this notice and the full citation o components of this work owned by othe be honored. Abstracting with credit is perm republic, to post on servers or to redistribu permission and/or a fee. Request permission: SIGCOMM '18, August 20–25, 2018, Budapest, I © 2018 Copyright held by the owner/authors; to the Association for Computing Machinery. ACM ISBN 978-1-4503-5667-9/18/08... \$15.00 <https://doi.org/10.1145/3230543.3230582>

to network functionality can introduce bugs that may cause great damage. Recently, faulty routers in two airline networks have grounded airplanes for days (for both Delta and Southwest Airlines), showing just how disruptive the effects of incorrect network behavior can be. Given the momentum behind programmable networks, we expect such faults and many others will cripple programmable networks.

In this paper, we argue that dataplane programs should be verified before deployment to enable safe operation. We present Vera, a verification tool that enables debugging of P4 programs both before deployment and at runtime. At its core, Vera translates P4 to SEFL, a network language designed for verification, and relies on symbolic execution with Symex

ABSTRACT

Recent trends in software-defined networking have extended network programmability to the data plane. Unfortunately, the chance of introducing bugs increases significantly. Verification can help prevent bugs by ensuring that the program does not violate its requirements. Although research on the verification of P4 programs is very active, we still need tools to make easier for programmers to express properties and to rapidly verify complex invariants. In this paper, we leverage assertions and symbolic execution to propose a more general P4 verification approach. Developers annotate P4 programs with assertions expressing general network correctness properties; the result is transformed into C models and all possible paths symbolically executed. We implement a prototype, and use it to show the feasibility of the verification approach. Because symbolic execution does not scale well, we investigate a set of techniques to speed up the process for the specific case of P4 programs. We use the prototype implemented to show the gains provided by three speed up techniques (use of constraints, program slicing, parallelization), and experiment with different compiler optimization choices. We show our tool can uncover a broad range of bugs, and can do it in less than a minute considering various P4 applications.

CCS CONCEPTS

• Networks → Programmable networks; • Software and its engineering → Software verification and validation;

KEYWORDS

P4; Verification; Programmable Data Planes

ACM Reference Format:

Miguel Neves, Lucas Freire, Alberto Schaeffer-Filho, Marinho Barcellos. 2018. Verification of P4 Programs in Feasible Time using Assertions. In The 16th International Conference on Emerging Networking Experiments and Technologies (CoNEXT '18), December 4–7, 2018, Heraklion, Greece. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3281141.3281421>

1 INTRODUCTION

Data plane programmability allows operators to quickly deploy new protocols and develop network services. Through programming languages such as P4 [2], it is possible to specify new instructions

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permission from permissions@acm.org.

CoNEXT '18, December 4–7, 2018, Heraklion, Greece
© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-5667-9/18/12... \$15.00 <https://doi.org/10.1145/3281141.3281421>

which packet headers should be manipulated, and how, by different forwarding devices in the infrastructure. Despite the flexibility, this paradigm also increases the chance of introducing bugs into the network due to incorrect implementations.

Testing, debugging, verification and enforcement are complementary approaches that can help solve this problem. During development, data plane programs can be debugged and tested, providing a wide range of inputs and checking if the corresponding outputs match the expected behavior. Verification, on its turn, can be used on programs to find bugs that would violate any of the properties stated by their requirements, including bugs that are hard to reproduce by testing. Lastly, with enforcement, the data plane can be monitored during execution to trap and block actions that would result in property violations.

In this paper, we focus on verification: we propose an approach to model and check (at compile time) general security and correctness properties of P4 programs, and implement it in a tool that provides network verification in feasible time. Several approaches have been developed to check if a given fixed-function (non-P4) data plane satisfies a set of intended properties [8, 25, 29, 32]. Moreover, verifying P4-programmed data planes is an active area of research, with recent projects proposing verification techniques based on SMT solving [24, 27] and custom symbolic execution [31]. In contrast, this work shows how to efficiently verify P4 programs leveraging a popular, off-the-shelf symbolic execution engine [4].

We propose an expressive assertion language (highly influenced by P4) that enables programmers to specify their intended properties by annotating their P4 code. Once annotated, a program is symbolically executed, with assertions being checked while all its paths are traversed. Given that the time taken to perform the symbolic execution grows exponentially with the program complexity, we show how a variety of speed up techniques can be employed to reduce the verification time and number of executed instructions. These techniques consist of using annotations in code to constrain the paths to be traversed according to properties and/or protocols of interest, program slicing to reduce the complexity of the model under verification, and parallelization of symbolic execution. Besides, we experiment with code optimization features offered by current compilers.

To evaluate our approach, we built a prototype using KLEE [4] and the P4 Reference Compiler [20] for the current language version, P4v. We applied it to four real P4 applications collected from the literature: Switch [21], NetPass [5], Duplex [11], and DCp4 [31]. Our results show that the proposed verification process can uncover a broad range of bugs either in the data plane program itself or in its control plane configuration. A detailed performance analysis also shows that, although the verification time grows exponentially with

[24, 31] were independently developed at the same time as this work.

[SIGCOMM'18]

[CoNEXT'18]

Programmable routers...

(specifically, programmable data planes)

...how do they work?

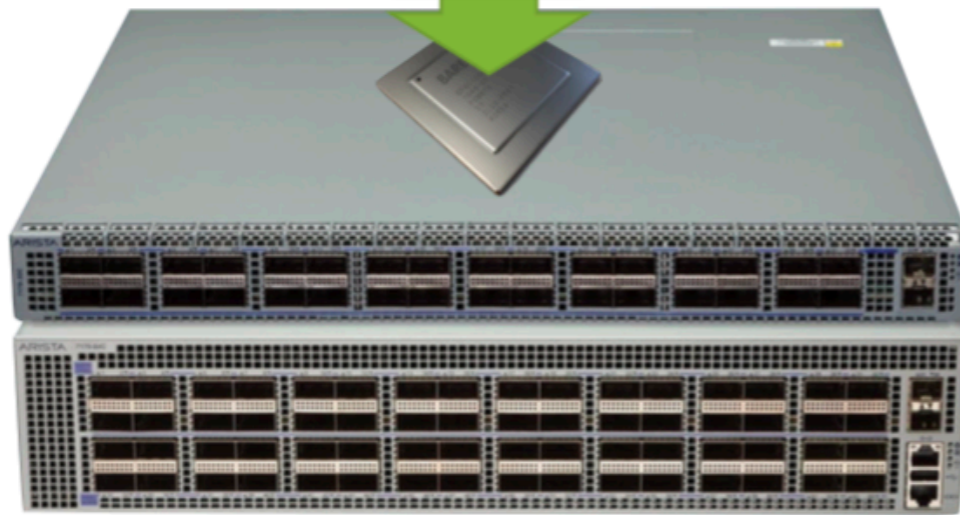


Arista 7170 series switches

Let's verify!



Bit-level description
of data-plane behaviour



Arista 7170 series switches

Give programmers language-based
verification tools

P4 also used as HDL for fixed-function
devices

P4 by example

- P4 is a low-level language → many gotchas
- Let's explore by example!
 - IPv6 router w/ access control list (ACL)

```
control ingress { apply(acl); }  
  
table acl {  
  reads { ipv6.dstAddr: lpm; }  
  actions { allow; deny; }  
}  
  
action allow() {  
  modify_field(std_meta.egress_spec, 1);  
}  
  
action deny() { drop(); }
```

What could *possibly* go wrong?

What if we didn't receive an IPv6 packet?

ipv6 header will be **invalid**

What goes wrong

Table reads arbitrary values

→ Intended ACL policy violated

Can read values from a previous packet

→ Side channel vulnerability!

Real programs are complicated:
hard to keep validity in your head

```
control ingress { apply(acl); }

table acl {
  reads { ipv6.dstAddr: lpm; }
  actions { allow; deny; }
}

action allow() {
  modify_field(std_meta.egress_spec, 1);
}

action deny() { drop(); }
```

Property #1: header validity

What if acl table misses (no rule matches)?

Forwarding decision is unspecified

What goes wrong

Forwarding behaviour depends on hardware

- May not do what you expect!
- Code not portable

```
control ingress { apply(acl); }

table acl {
  reads { ipv6.dstAddr: lpm; }
  actions { allow; deny; }
}

action allow() {
  modify_field(std_meta.egress_spec, 1);
}

action deny() { drop(); }
```

Property #2: unambiguous forwarding

Types of properties

General safety

- **Header validity**
- Arithmetic-overflow checking
- Index bounds checking (header stacks, registers, meters, ...)

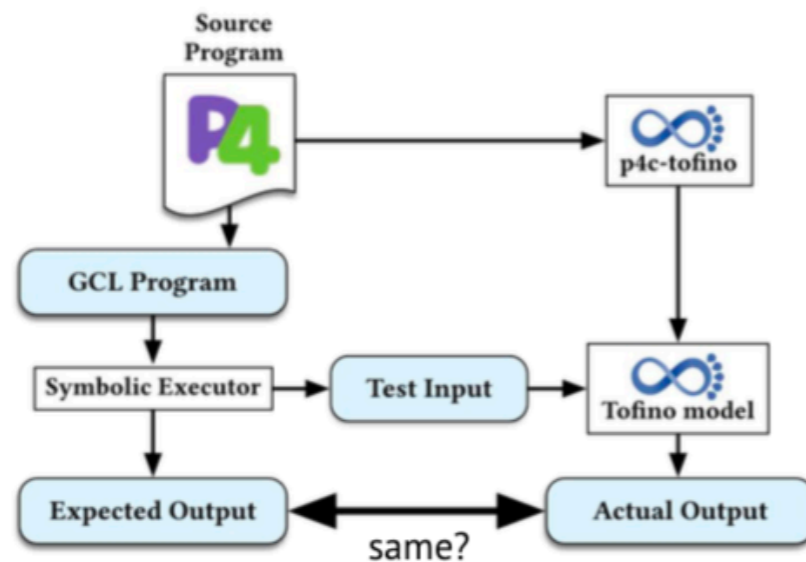
Architectural

- **Unambiguous forwarding**
- **Reparseability**
- **Mutual exclusion of headers**
- Correct metadata usage (e.g., read-only metadata)

Program-specific

- Custom assertions in P4 program — e.g., IPv4 ttl correctly decremented

Challenge #1: imprecise semantics



- P4 language spec doesn't give precise semantics
- Defined semantics by translation to GCL (a simple imperative language)
- Tested semantics
 - Symbolically executed GCL to generate input-output tests for several programs
 - Ran w/ Barefoot P4 compiler & Tofino simulator

Challenge #2: modelling the control plane

- A P4 program is just half the program
 - Table rules are not statically known
 - Populated by the control plane at run time
- Control planes are carefully programmed
 - Tables rarely take arbitrary actions
- To rule out false positives, need to model behaviour of control plane



```
table acl {
  reads {
    ipv6.dstAddr: lpm;
  }
  actions { allow; deny; }
}
```



```
( @[ Action ] acl <hit> (allow);
  std_meta.egress_spec := 1)

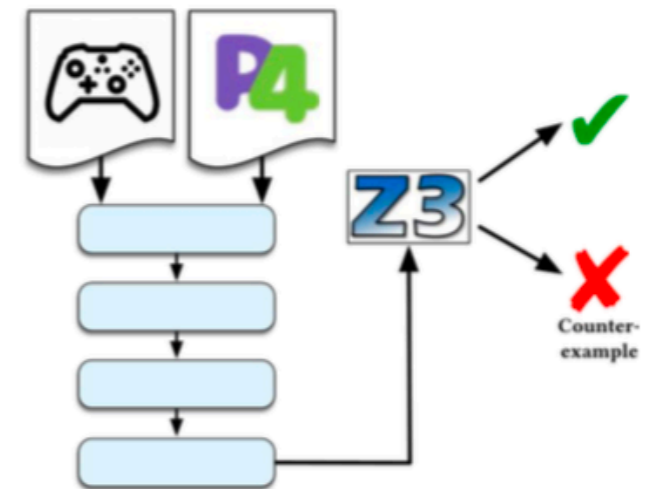
[] ( @[ Action ] acl <hit> (deny);
    std_meta.egress_spec := 511)

[] @[ Action ] acl <miss>
```

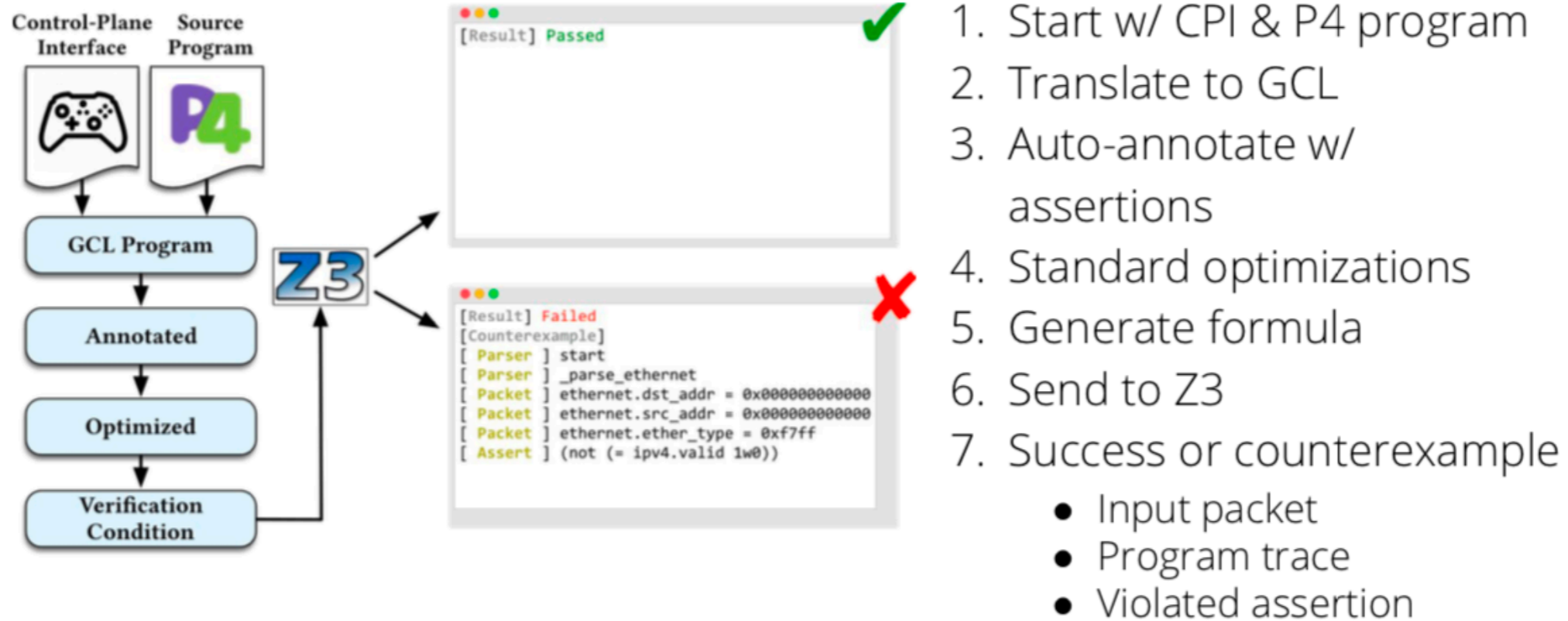
Tables translated into *unconstrained* nondeterministic choice

p4v overview

- **Automated** tool for verifying P4 programs
- Considers **all paths**
 - But also practical for **large programs**
- Includes basic safety properties for any program
- **Extensible** framework
 - Verify custom, program-specific properties
 - Assert-style debugging



p4v architecture



Source: p4v, Practical Verification for Programmable Data Planes, Liu et al., 2018

Data plane
programmability for

Performance
Monitoring
Applications offloading

Platforms for Data plane
Correctness programmability
Management

So you've a *verified* programmable networks...

How do you manage it?!

How do you perform planned maintenance?

now that you've state in your switches...

How do you run multiple applications in your switches?

monitoring, forwarding, load-balancing, etc.

How do you share resources amongst applications?

especially memory and # packet operations

We need an Operating System for the data plane

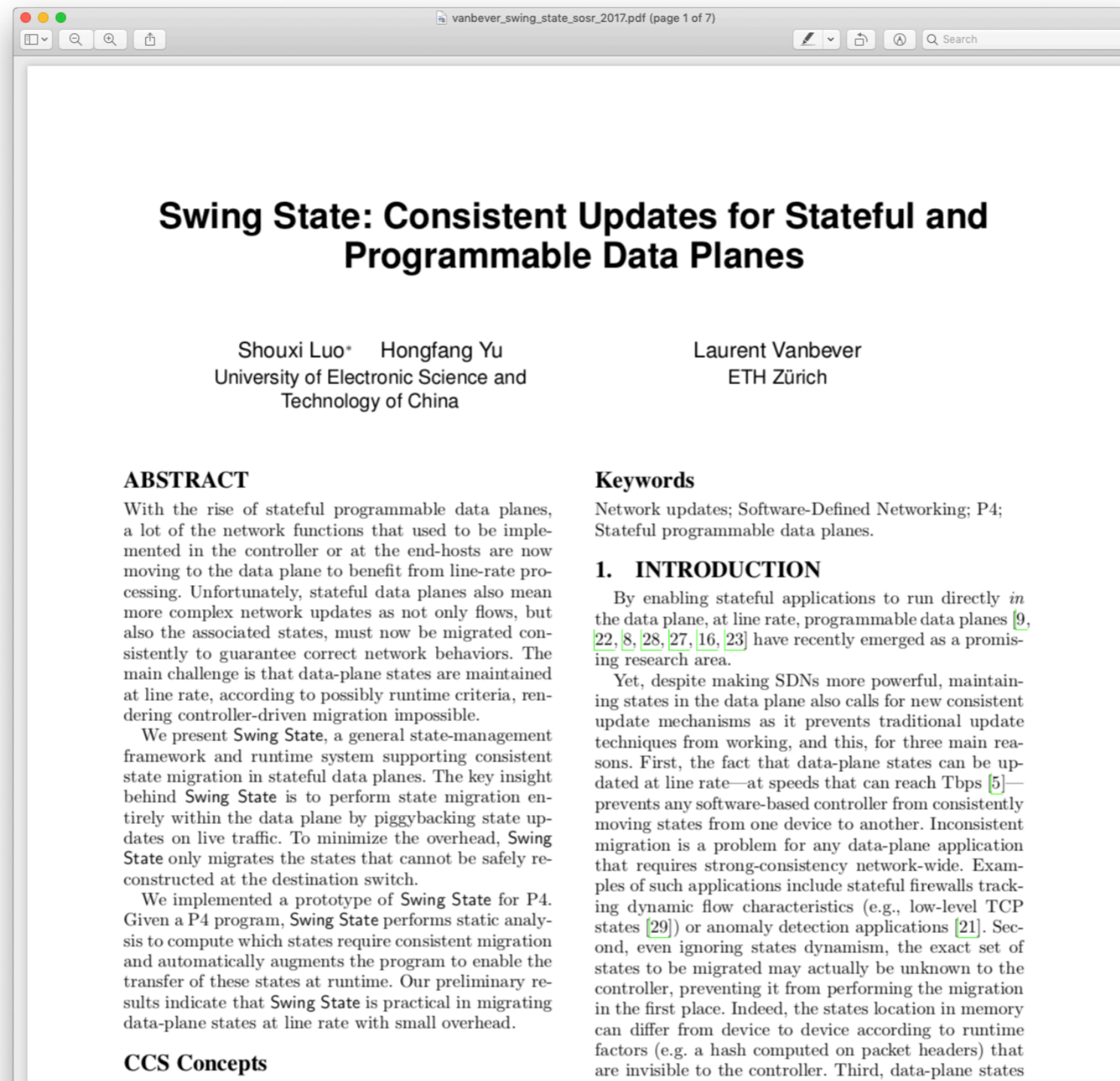
Definition
Wikipedia

An operating system is a system software that manages computer hardware and software resources and provides common services for computer programs.

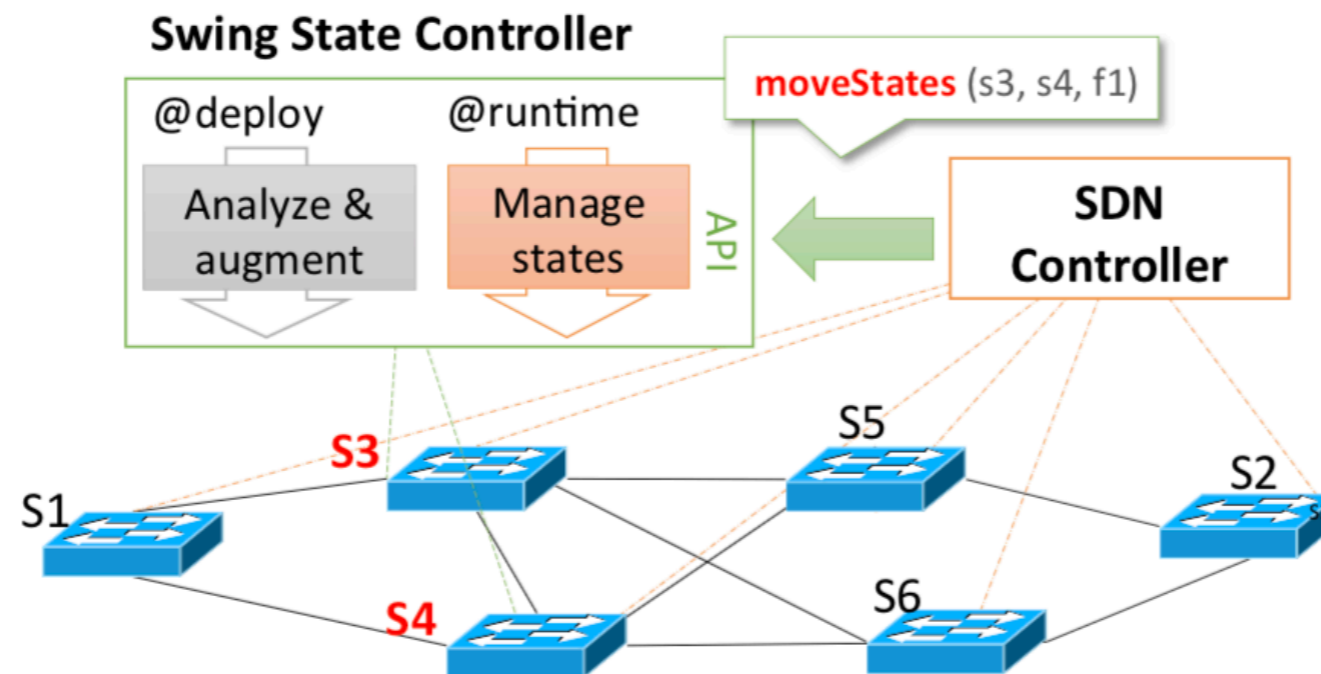
Do we have that? **Nope.** Not yet at least.

We're working on it...

[SOSR'17]




Swing State is a state management framework with 1 primitive: **moveStates**



Source: Swing State: Consistent Updates for Stateful and Programmable Data Planes
Luo et al., SOSR 2017

Advanced Topics in Communication Networks



Lectures/Exercices

~7 weeks


how to program in P4

Group project

≥ 7 weeks

in teams of 2—3

Advanced Topics in Communication Networks



Lectures/Exercices

~7 weeks

how to program in P4

Group project

≥ 7 weeks

in teams of 2—3

The group project starts this week

It accounts for 50% of your final grade

The evaluation of your project will depend on your implementation, report, and presentation

The evaluation of your project will depend on your implementation, report, and presentation

implementation

70%

achieves the basic goals

is properly documented

runs...

The evaluation of your project will depend on your implementation, report, and presentation

implementation

70%

achieves the basic goals

is properly documented

runs...

report

15%, 10 pages max

describes the main building blocks

evaluates the solution

describes what each group member did

The evaluation of your project will depend on your implementation, report, and presentation

implementation

70%

achieves the basic goals

is properly documented

runs...

report

15%, 10 pages max

describes the main building blocks

evaluates the solution

describes what each group member did

presentation

15%, 12 min. +questions

summarizes the problem and the solution

contains a *live* demo

involves all group members

The final deadline for the project is

Wed Dec 19 at 23.59pm

This week

Select a proposal from the list (see Doodle)
or send us your own proposal by email

Every week

Meet with the responsible assistant
schedule a recurring slot in [10.15am; noon]

Wed Dec 19
11.59pm

Send us an archive with report, code, slides

Thu Dec 20
8.15am—

Groups presentation + course/exam debrief
attendance is mandatory

The project has to be done in groups of 3 students
"Matching" process for incomplete groups via Slack

Project grade is shared by each group member
provided that each collaborated (roughly equally)

- Let us know in advance if that's *not* the case
- Briefly describe in the report the contribution of each group member
- Each group member should be involved in the presentation and be able to answer questions

Details about each proposal is available on our website

Advanced Topics in Communication Networks Project Proposals

Proposal #1: Hardware-Based RSVP

Responsible: Albert Gran Alcoz

Resource Reservation Protocol (RSVP) [1] is a signaling protocol that allows connections in a network to perform bandwidth requests throughout a given path. It is a protocol that has been included in different solutions both in the traffic engineering field and in quality of service. Integrated Services (IntServ) was the first in adopting it, in the late 1990s, as a means to provide guaranteed quality of service in multimedia networks. Some years later, and with higher success, RSVP was extended for traffic engineering purposes in the RSVP-TE protocol [2] to be used for the establishment of virtual circuits in MPLS. RSVP suggests users in a network to perform bandwidth reservations before starting data transmissions. For that, packet probes are forwarded from source to destination, letting routers in between identify the amount of resources requested by the new connection. Routers will receive those requests and reply to them by annotating in the same packet their resource availability. Flows will only be admitted if all routers along the path have agreed on having enough resources for hosting the new request. Although achieving notable and robust performance, being able to provide 100% resource guarantees, the high price that RSVP requires in terms of scalability and complexity, has made from it a not very successful solution in multiple scenarios until nowadays. Among the main drawbacks, the most remarkable ones are the time required to set up a new connection (too high especially for real-time flows), the amount of state to be stored in each switch along the path (to keep track of reservations), and the periodic overheads needed to refresh reservation requests.

In this project, we propose the design and implementation of an evolved version of RSVP, based on P4, to be run directly on hardware. We strongly believe that a signaling protocol executed at line rate in the data-plane can be quicker in deploying configurations and faster in reacting to updates.

Students are expected to come up with a data-plane implementation, aiming to overcome RSVP original

Register your proposal (one per group)
from **Friday 3pm until Sunday 11.59pm**

The screenshot shows a Doodle poll interface. At the top, there is a navigation bar with the Doodle logo, links for 'Plans', 'Help', and 'English', and buttons for 'Sign up', 'Log in', and 'Create a Doodle'. The main heading is 'Adv-Net Group Projects' by Roland Meier, posted 18 hours ago. A central instruction box states: 'Please register as teams of 3 people (write the names of all team members). If a team has only 2 members, we might add another person. Reservations with only 1 name or with more than 3 names will be removed.' Below this is a table with 6 columns representing different proposals. The first column is for the overall poll status, showing '2 participants'. The other columns show the proposal name, the number of participants (e.g., '1/1' for the first two proposals, '0/1' for the others), and a registration button with a circular indicator.

	Proposal #1: Hardware-Based RSVP	Proposal #2: Data-plane driven network convergence	Proposal #3: Intra-domain routing in the data-plane	Proposal #4: Delay-based routing entirely in the data-plane	Proposal #5: Advanced stateful firewall
2 participants	✓ 1/1	✓ 1/1	✓ 0/1	✓ 0/1	✓ 0/1
<input type="text" value="Enter your name"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

If you want to propose your own project,
send me an email describing it by **Friday (Nov 2) 3pm**

Ivanbever@ethz.ch

Quick overview of the proposals



Albert



Thomas



Roland



Alexander



Edgar

Quick overview of the proposals



Albert



Thomas



Roland



Alexander



Edgar

Proposal #1

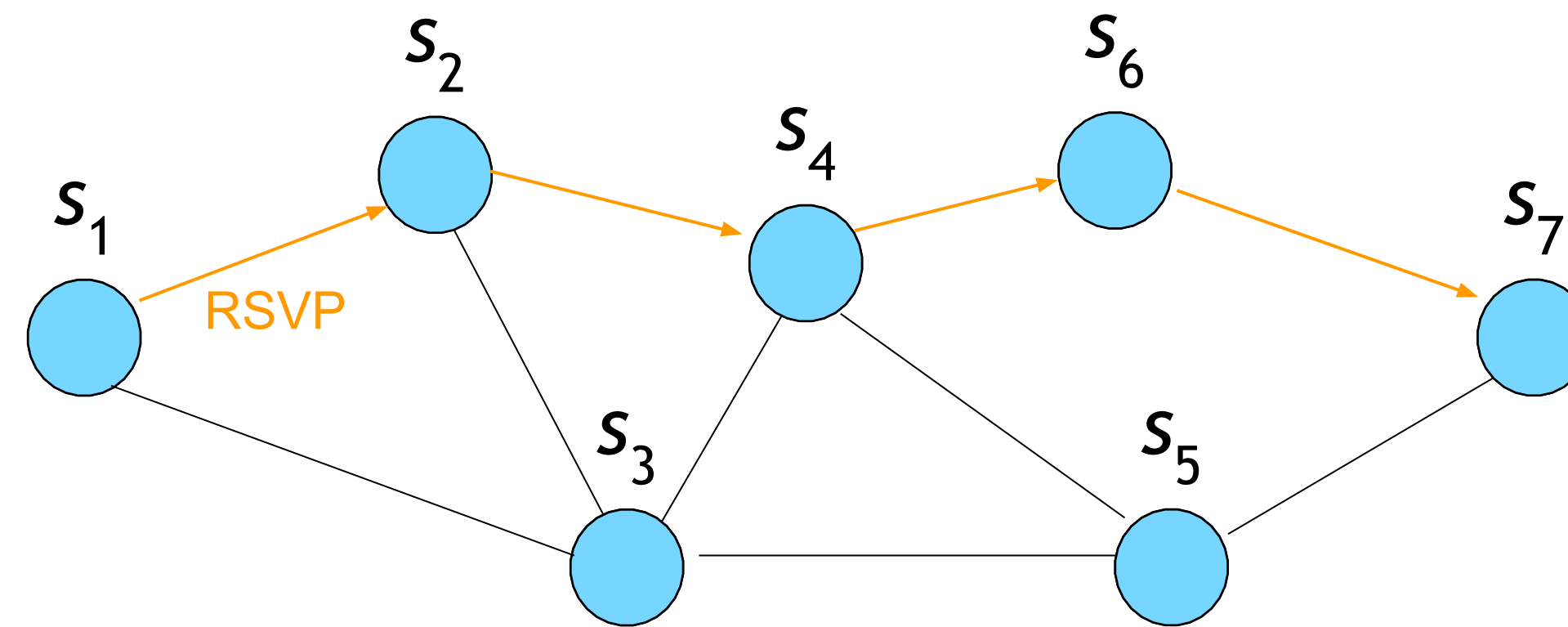
Hardware-Based RSVP

Bandwidth reservations throughout a given path:

- Quality of Service guarantees (IntServ)
- Establishment of virtual circuits (MPLS)

Exclusive **data plane** implementation:

- Personalized headers
- Header stacks
- Registers
- Bloom filters



Faster and more scalable than traditionally

Quick overview of the proposals



Albert



Thomas



Roland



Alexander

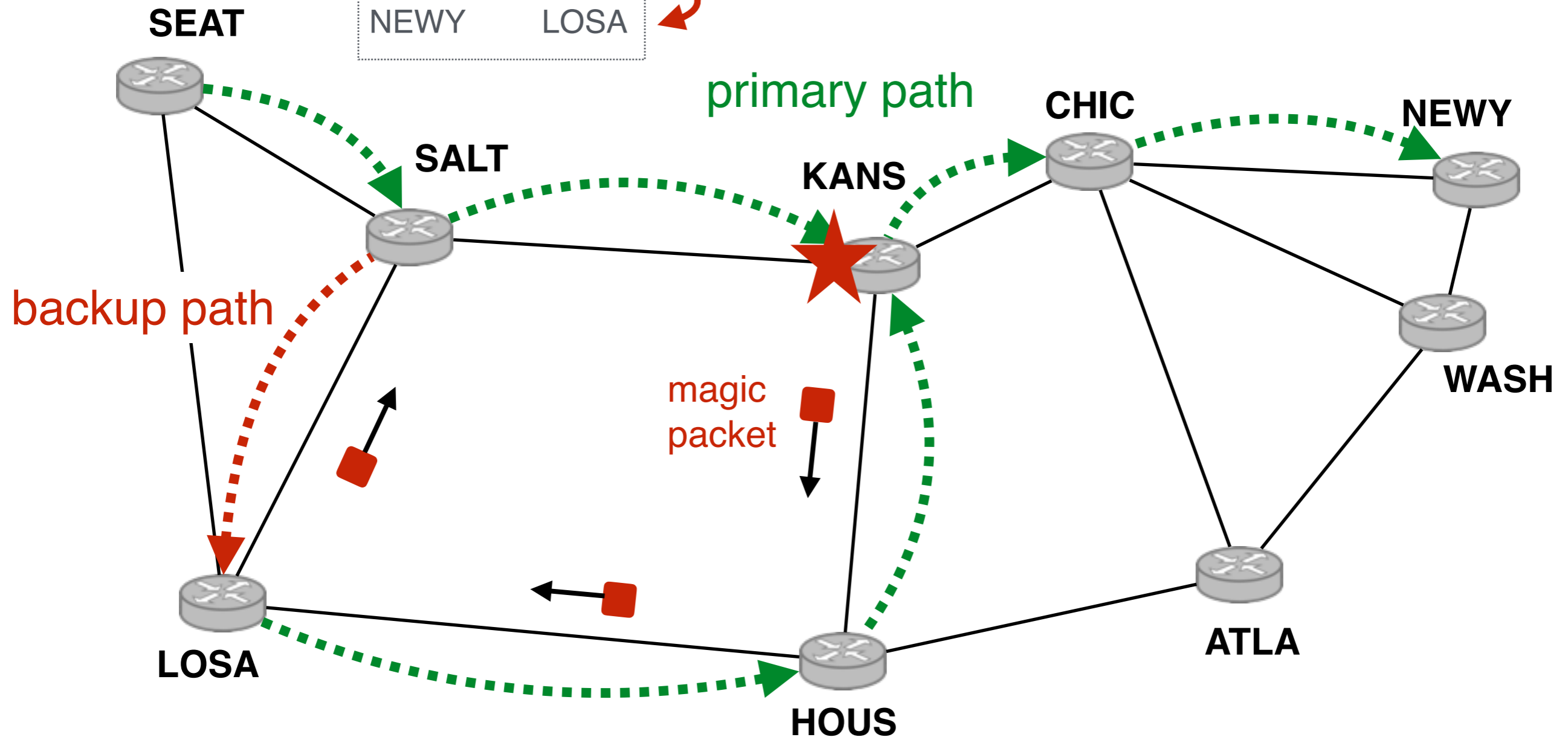


Edgar

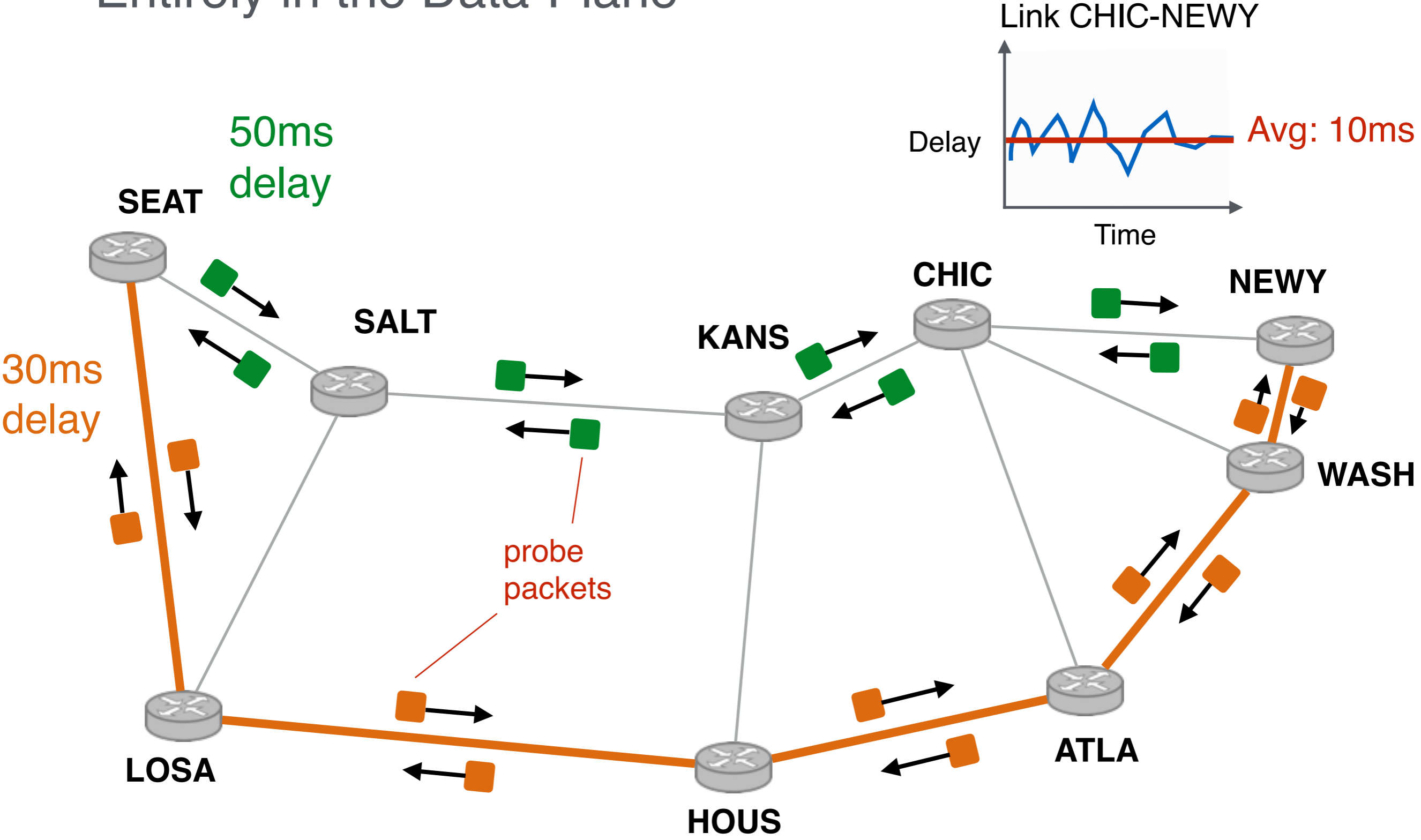
Proposal #2: Data-plane Driven Network Convergence

fw table at SALT

dest.	next-hop
NEWY	KANS
NEWY	LOSA



Proposal #3: Delay-based Routing Entirely in the Data-Plane



Quick overview of the proposals



Albert



Thomas



Roland



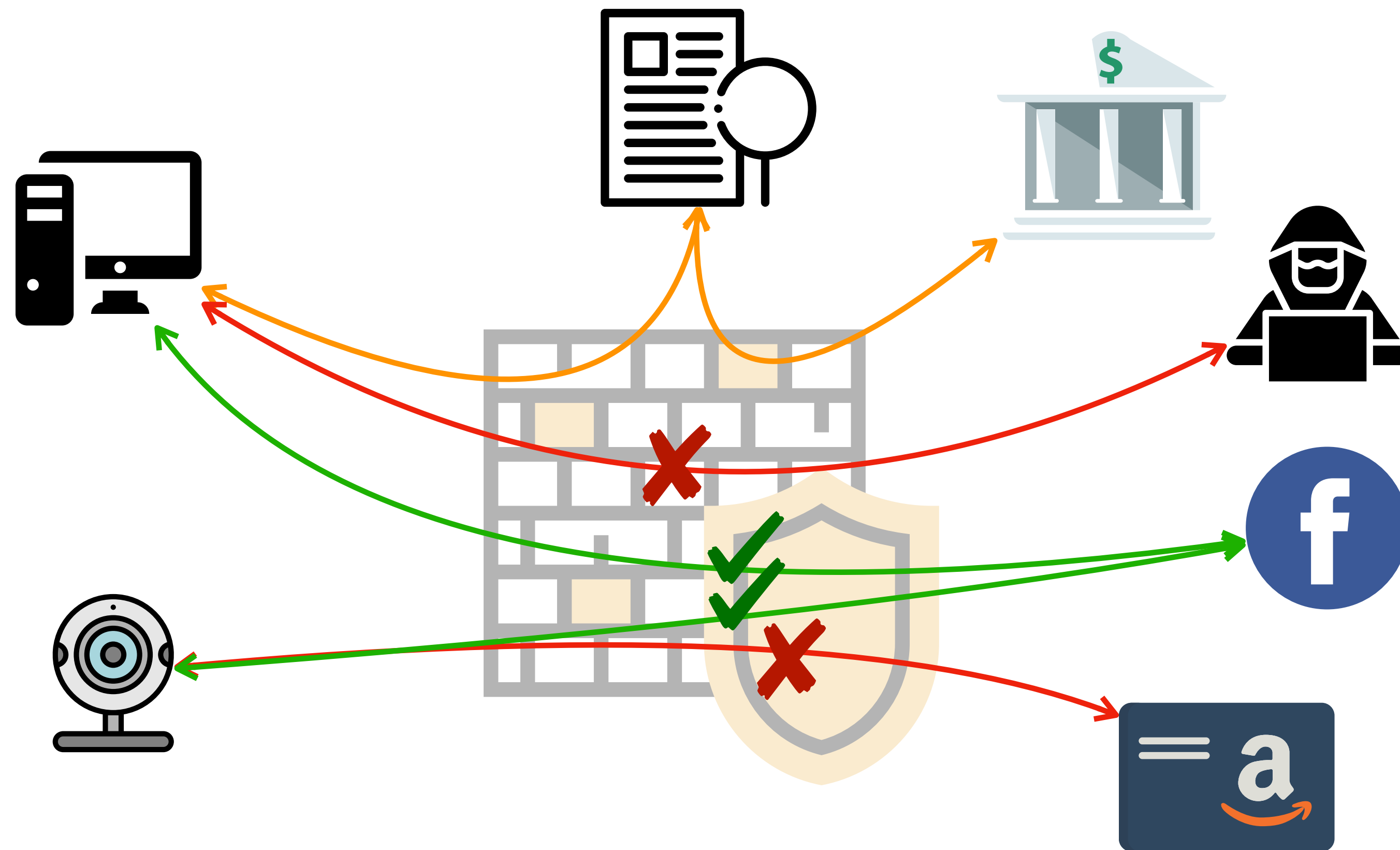
Alexander



Edgar

Proposal #4

Advanced stateful firewall

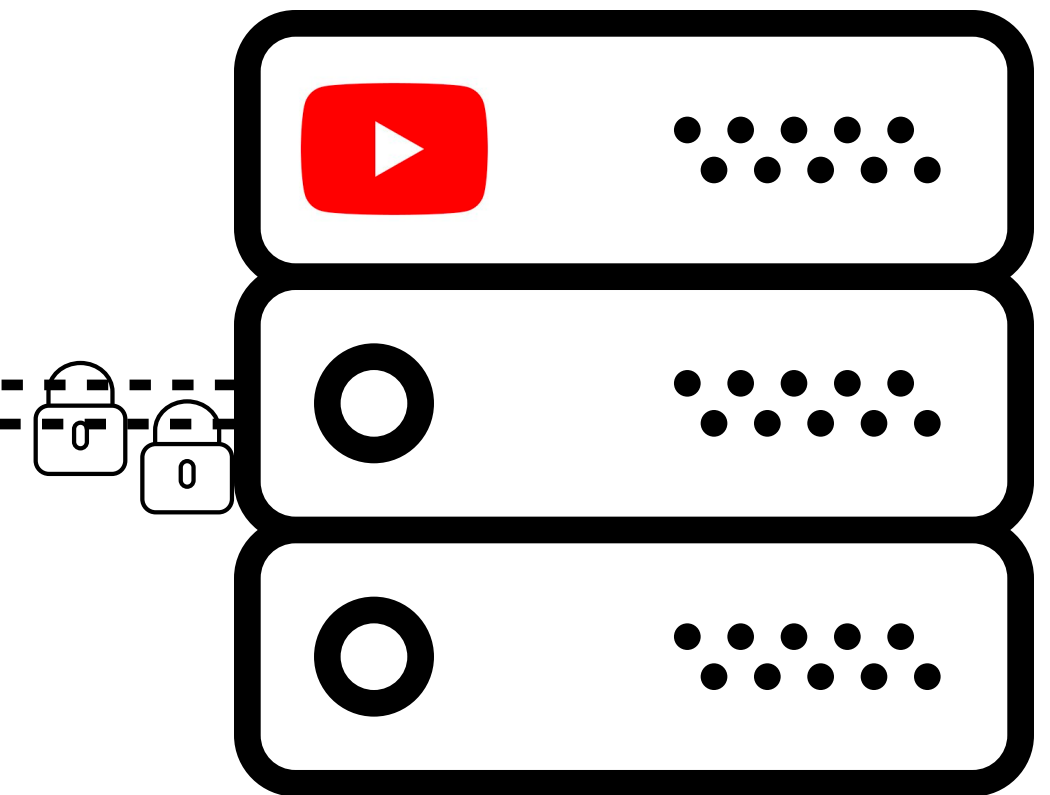
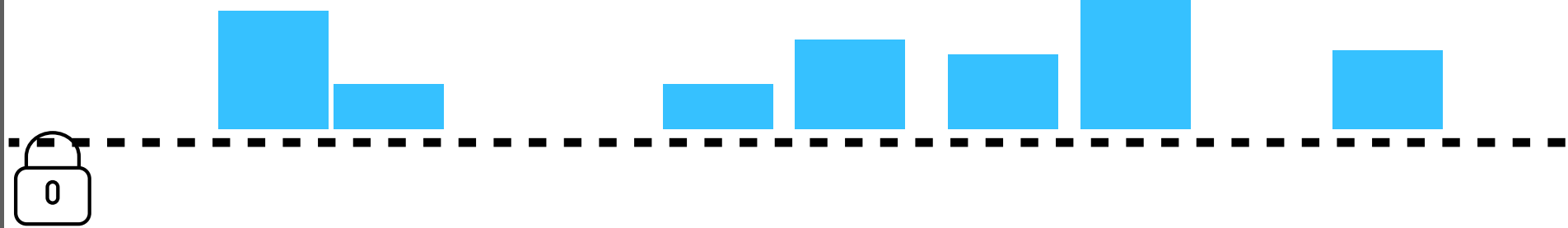
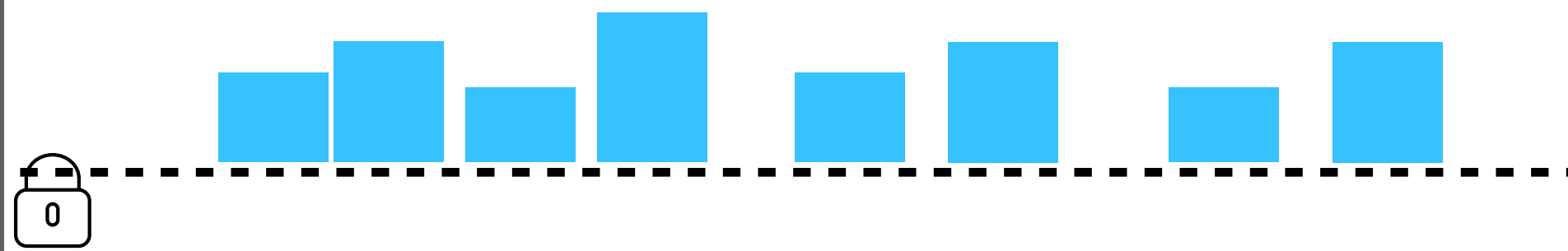


- ✓ Fine-grained access policies
- ✓ Deep packet inspection (DPI)
- ✓ VPN
- ✓ Attack detection
- ✓ Spoofing detection
- ✓ (add your idea here)

...

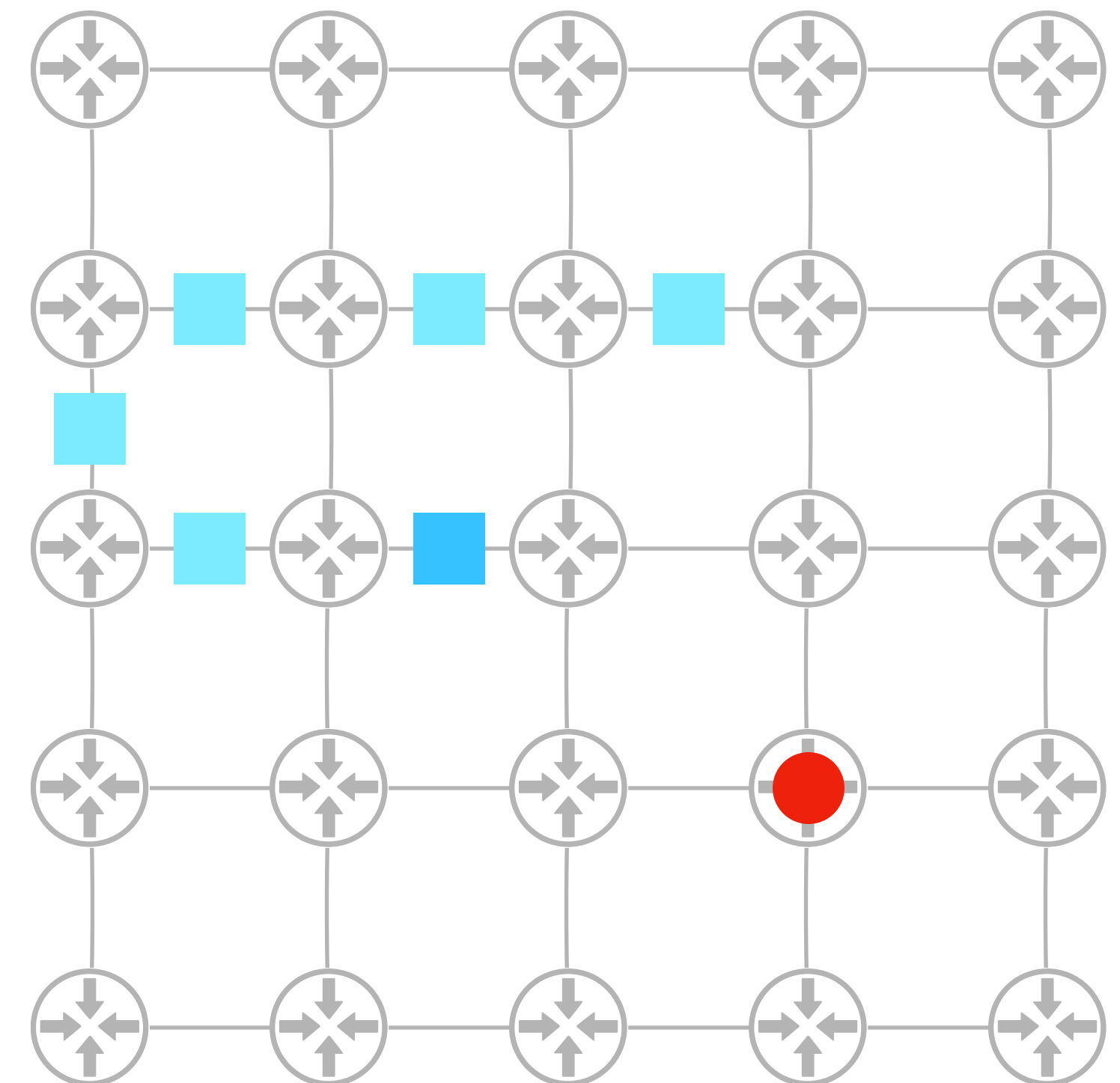
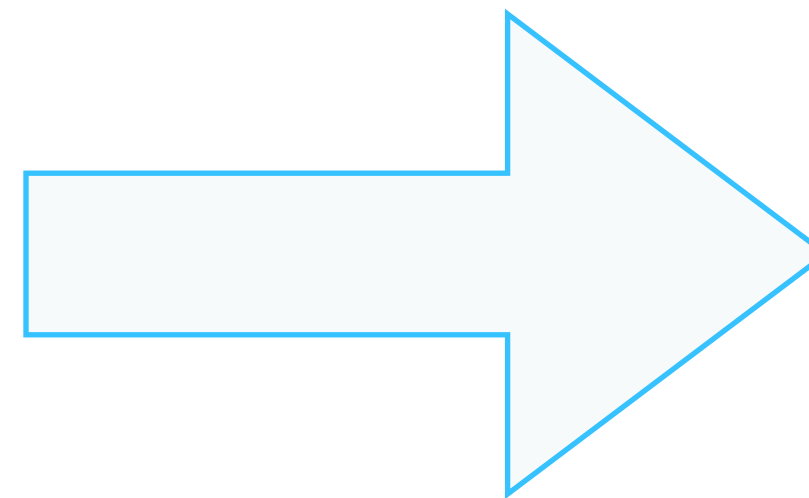
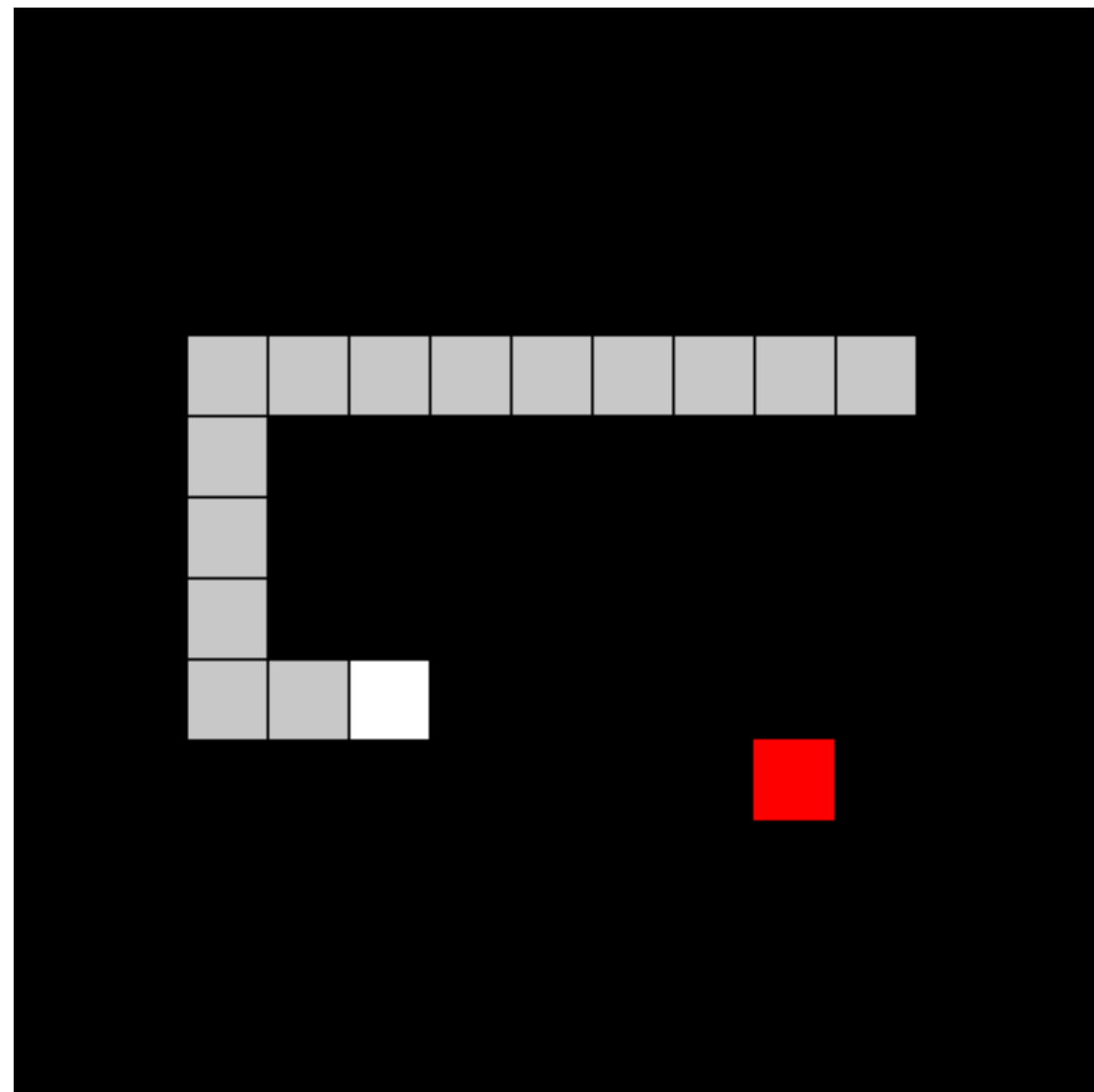
Proposal #5

I know what you're seeing now



Proposal #6

Playing snake in the data plane



Quick overview of the proposals



Albert



Thomas



Roland



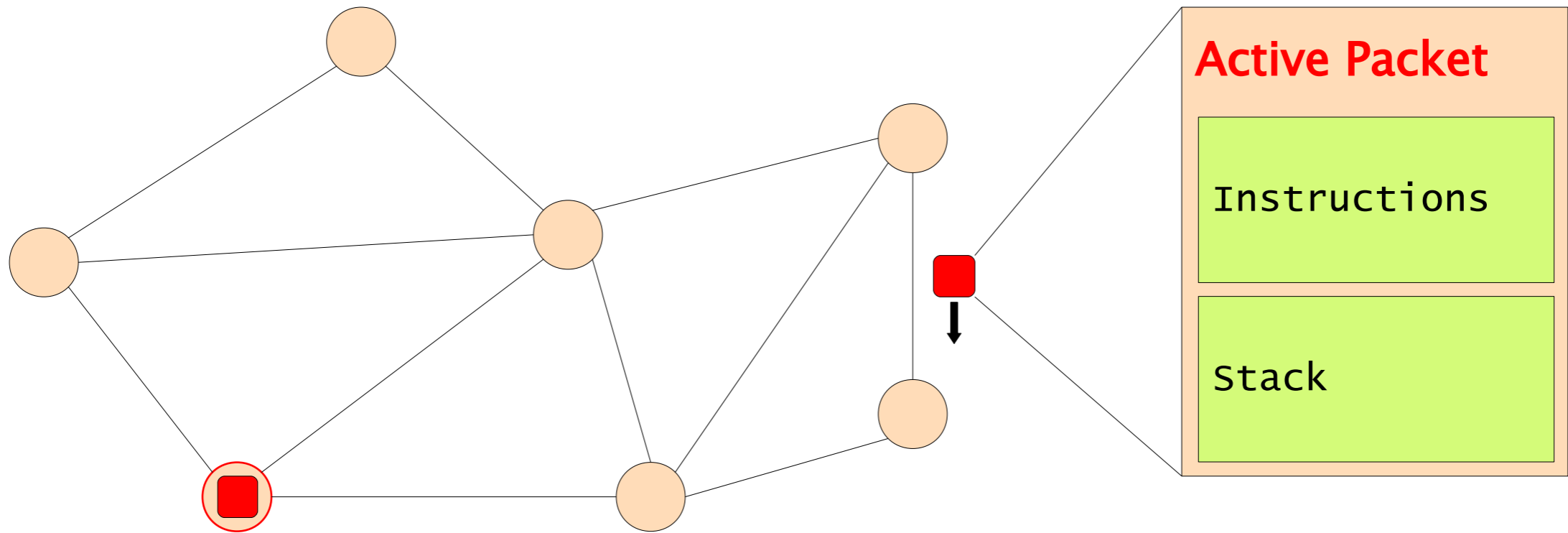
Alexander



Edgar

Proposal #7

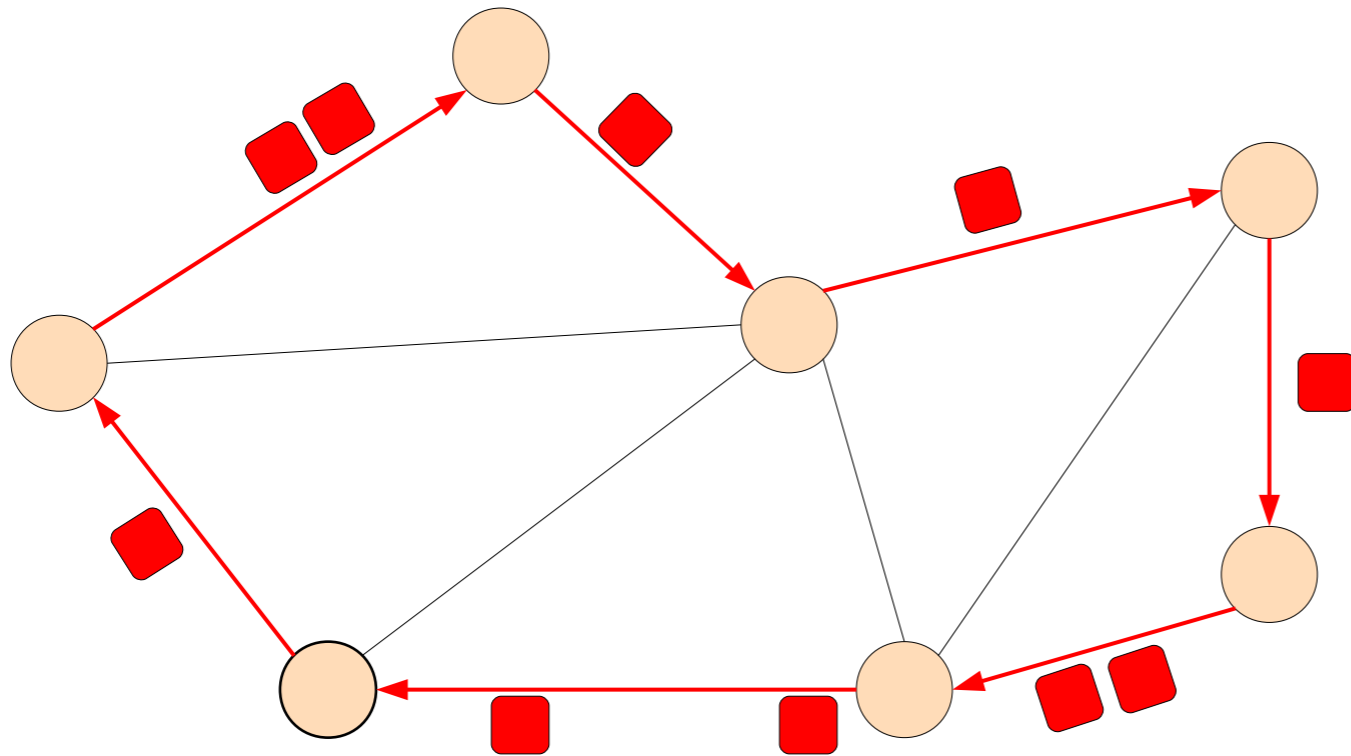
In **Active Networks**, packets carry **programs**.



The programs **are executed**
on each switch along the path

Proposal #8

Storing data in the cloud ~~the right way!~~



Store data in a
forwarding loop

Quick overview of the proposals



Albert



Thomas



Roland



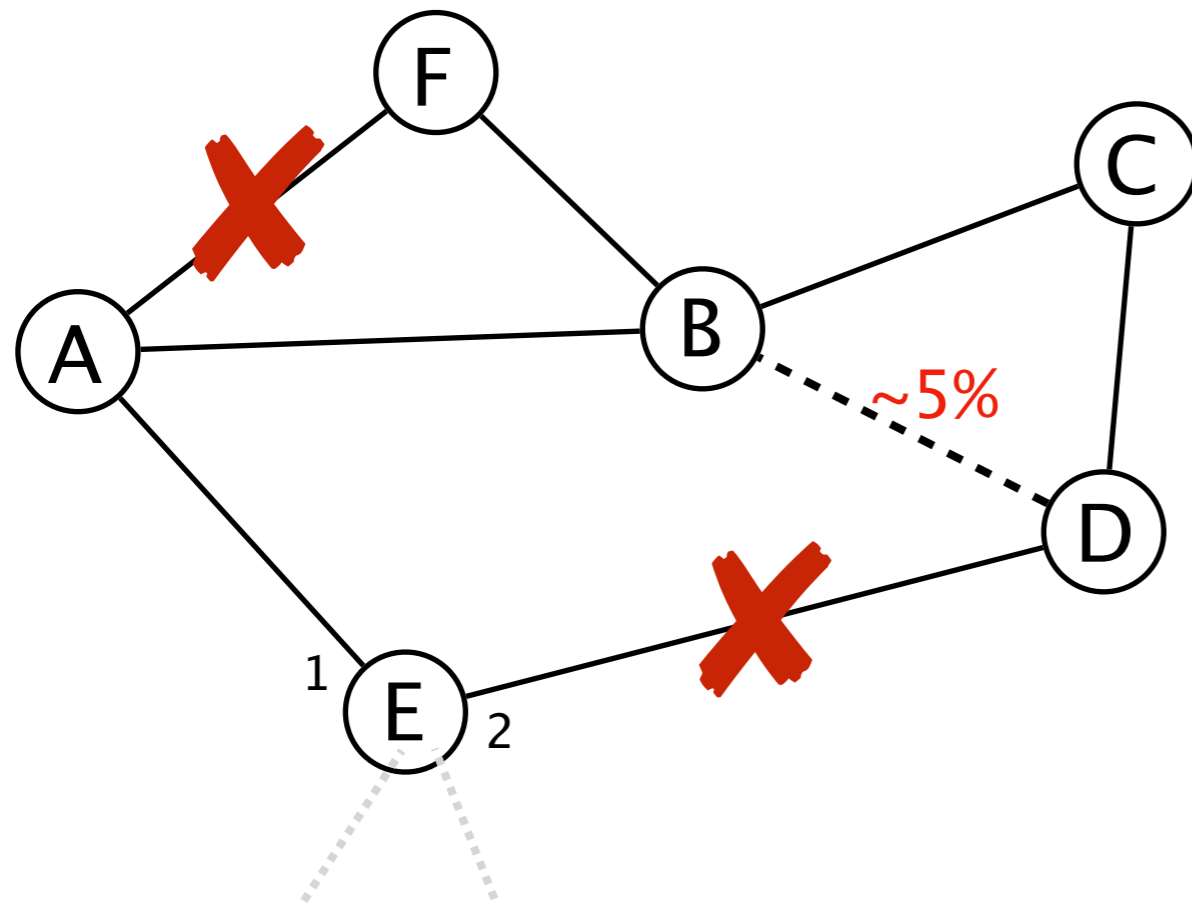
Alexander



Edgar

Proposal #9

Data Plane Failure Detection



prefix	port
10.1.1.0/24	01
10.1.2.0/24	10

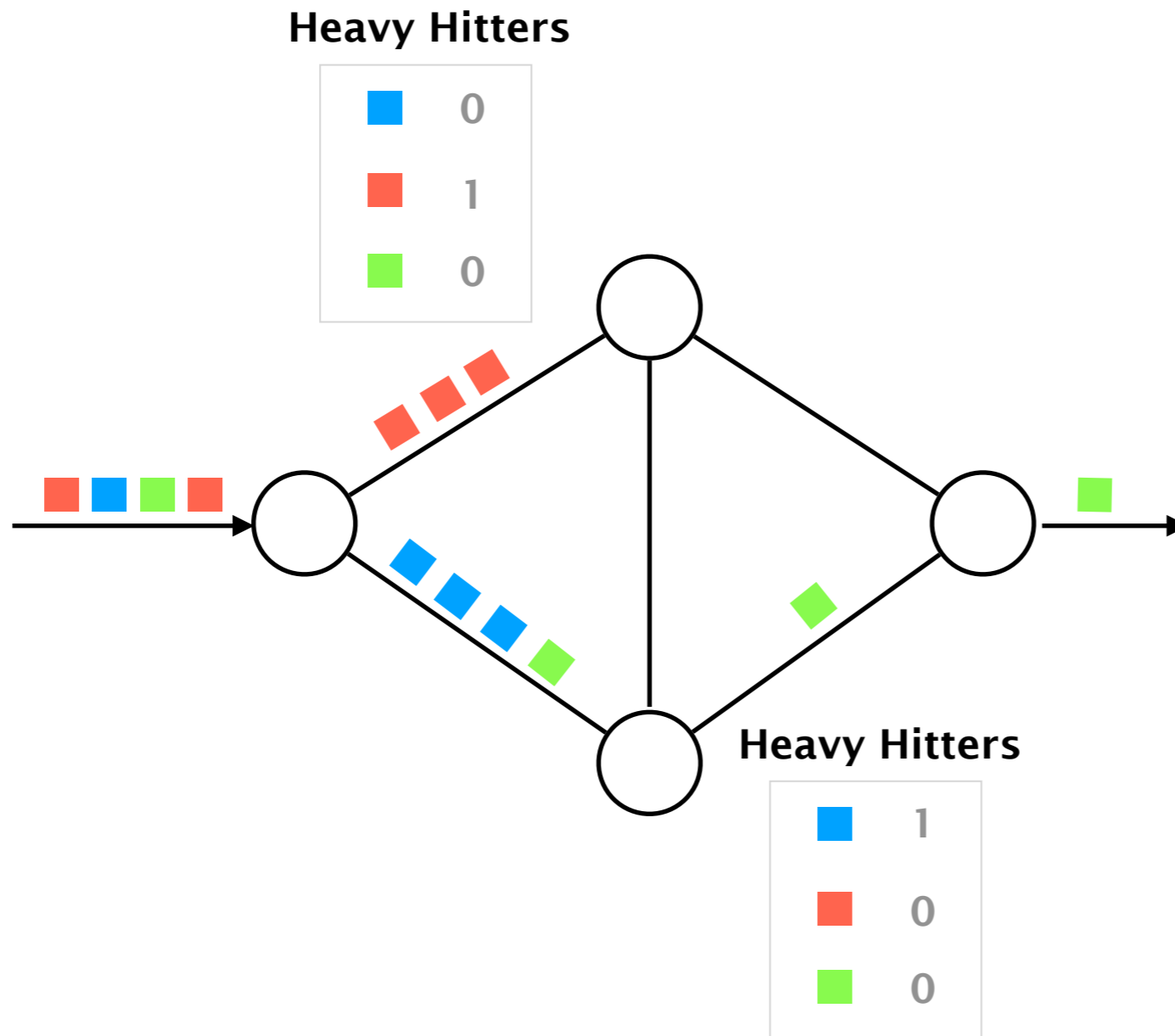
Detect **local** and **remote** link failures (A-C)

Detect **random** packet drops (B-C)

Detect **corrupted** table entries (E)

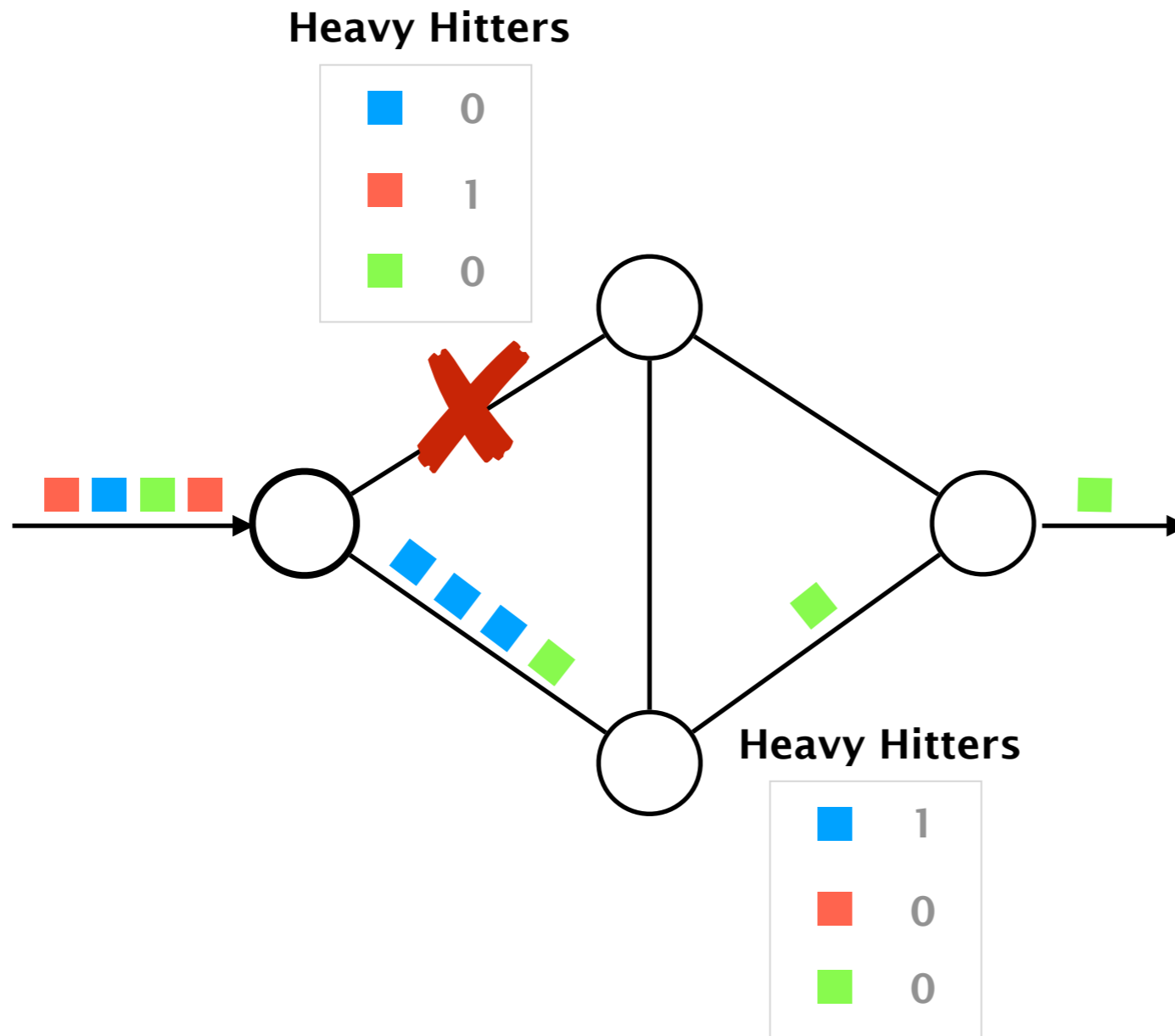
Proposal #10

Stateful Application Migration



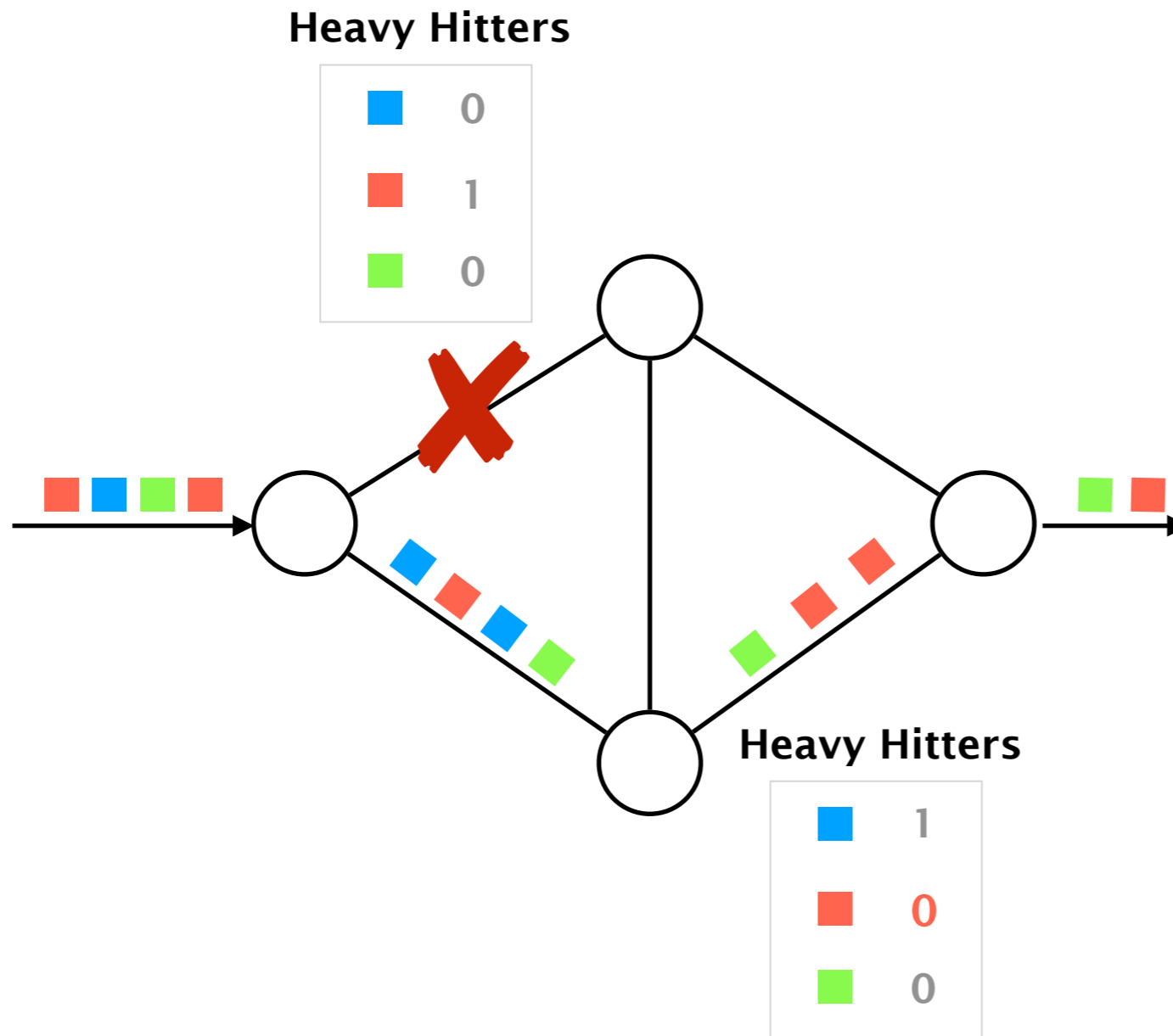
Proposal #10

Stateful Application Migration



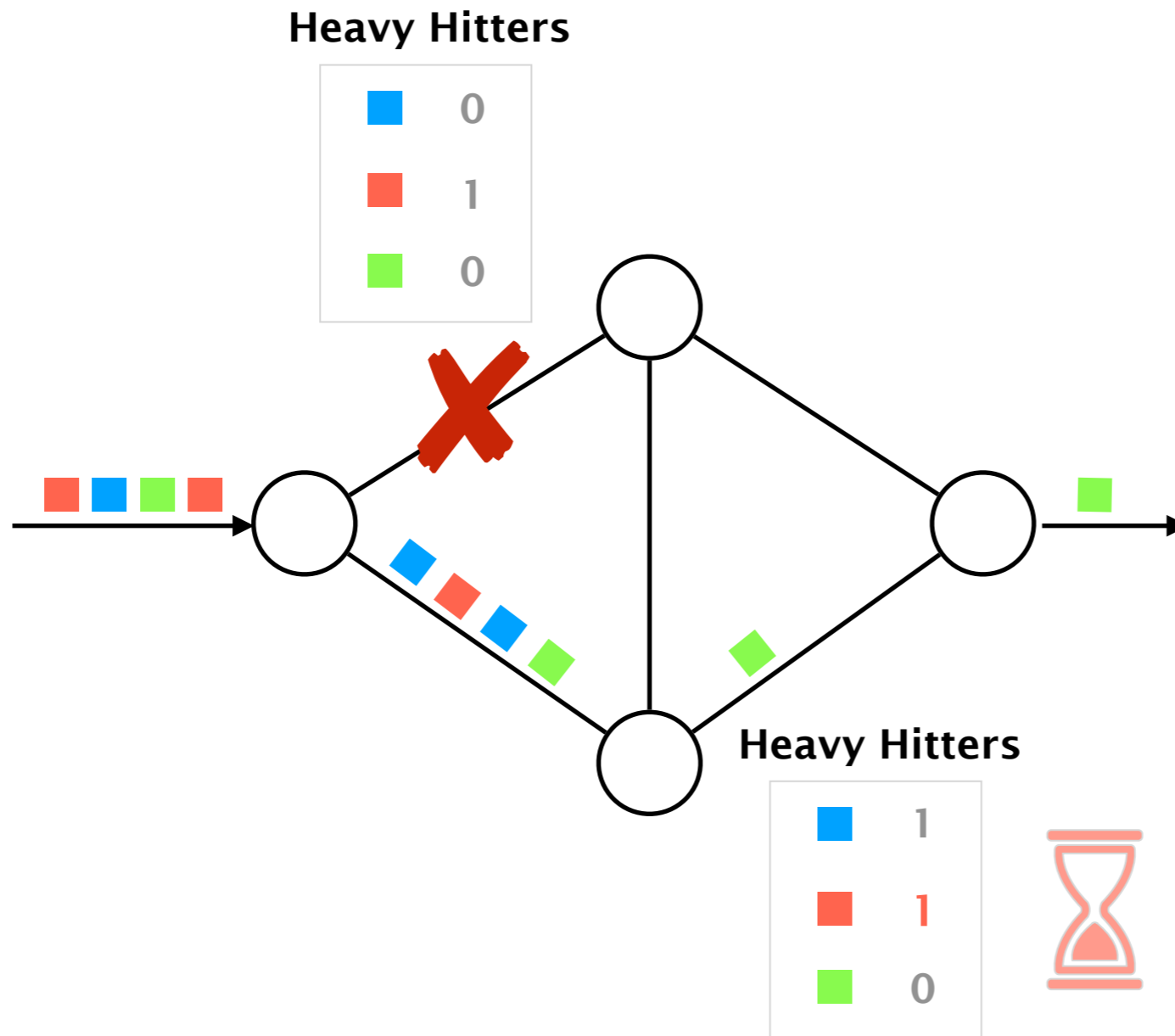
Proposal #10

Stateful Application Migration



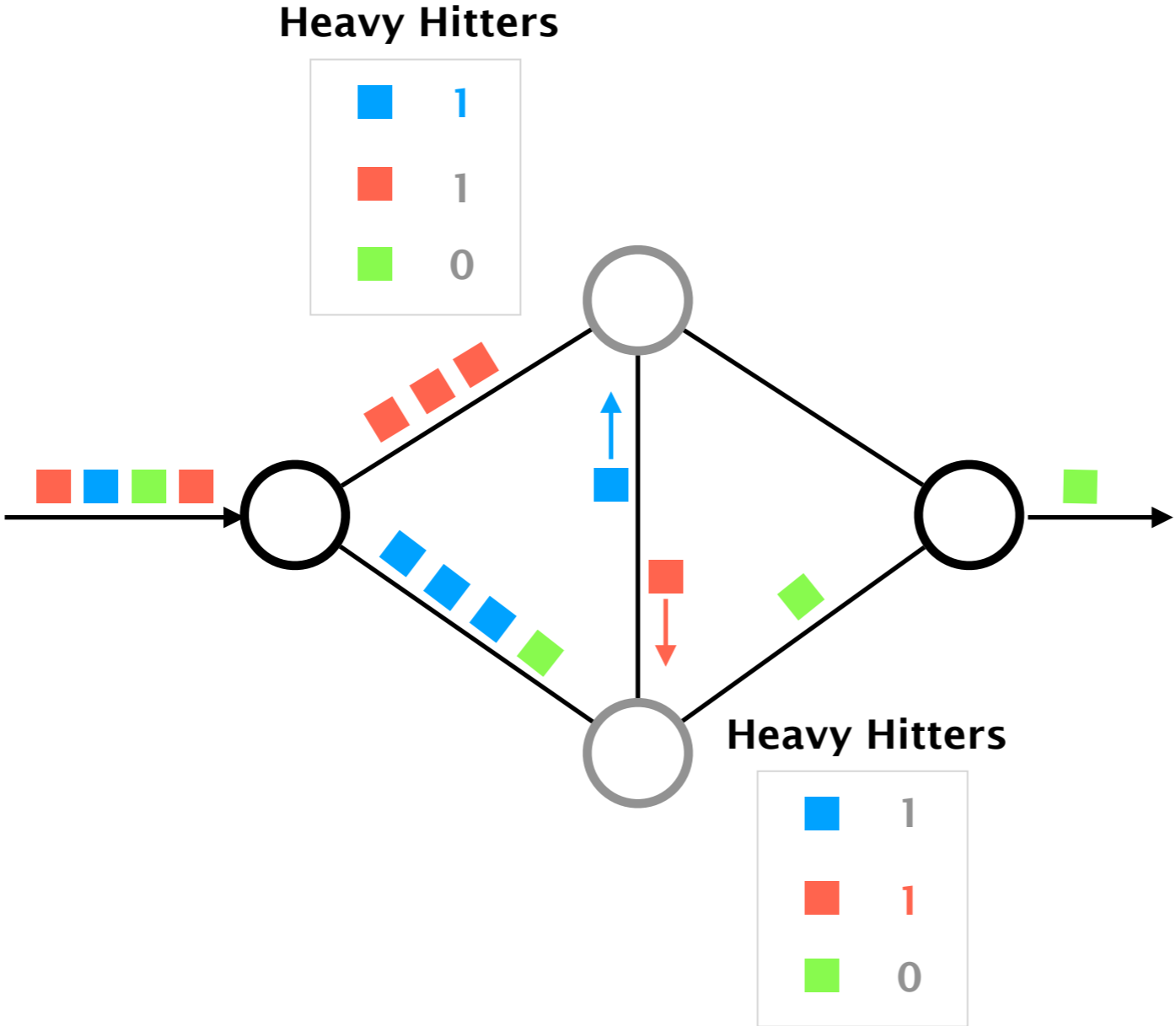
Proposal #10

Stateful Application Migration



Proposal #10

Stateful Application Migration



Proposal #11

P4 Switch

Management and Configuration API

Control Plane

Basic Features

I2 forwarding, learning, multicast

ipv4, ipv6, I3 multicast

ECMP, Weighted ECMP

ICMP

ARP

ECN

Simple QoS

Advanced Features

Spanning Tree Protocol

netflow, sFlow or similar

VXLAN, MPLS, Gre

DHCP Server

DNS Cache

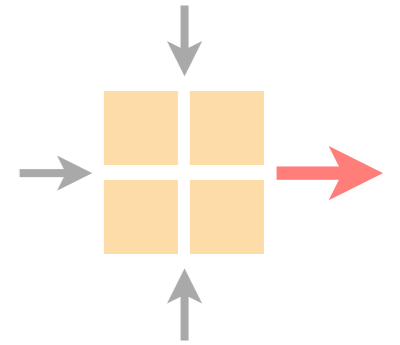
Simple Firewall

NAT

Data Plane

Advanced Topics in Communication Networks

Programming Network Data Planes



Laurent Vanbever

nsg.ee.ethz.ch

ETH Zürich

Nov 1 2018