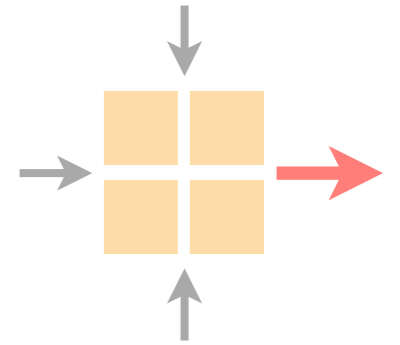Advanced Topics in Communication Networks

# Programming Network Data Planes

Laurent Vanbever

nsg.ee.ethz.ch

ETH Zürich

Oct 29 2019

Last week on

Advanced Topics in Communication Networks

P4 hardware
target

P4-based
applications

How do we build a *fast*
reprogrammable switch?

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

How can we allow network programmability in the field, at reasonable cost, and without <mark>sacrificing speed</mark>

supporting Tbps of
backplane throughput

# Let's look at a concrete design:

# Reconfigurable Match Tables (RMT)

sdn-chip-sigcomm-2013.pdf (page 1 of 12)

Search

## Forwarding Metamorphosis: Fast Programmable Match-Action Processing in Hardware for SDN

Pat Bosshart[†], Glen Gibb[‡], Hun-Seok Kim[†], George Varghese[§], Nick McKeown[‡],
Martin Izzard[†], Fernando Mujica[†], Mark Horowitz[‡]
[†]Texas Instruments   [‡]Stanford University   [§]Microsoft Research
pat.bosshart@gmail.com  {grg, nickm, horowitz}@stanford.edu
varghese@microsoft.com  {hkim, izzard, fmujica}@ti.com

## ABSTRACT

In Software Defined Networking (SDN) the control plane is physically separate from the forwarding plane. Control software programs the forwarding plane (e.g., switches and routers) using an open interface, such as OpenFlow. This paper aims to overcomes two limitations in current switching chips and the OpenFlow protocol: i) current hardware switches are quite rigid, allowing "Match-Action" processing on only a fixed set of fields, and ii) the OpenFlow specification only defines a limited repertoire of packet processing actions. We propose the RMT (reconfigurable match tables) model, a new RISC-inspired pipelined architecture for switching chips, and we identify the essential minimal set of action primitives to specify how headers are processed in hardware. RMT allows the forwarding plane to be changed in the field without modifying hardware. As in OpenFlow, the programmer can specify multiple match tables of arbitrary width and depth, subject only to an overall resource limit, with each table configurable for matching on arbitrary fields. However, RMT allows the programmer to modify *all* header fields much more comprehensively than in OpenFlow. Our paper describes the design of a 64 port by 10 Gb/s switch chip implementing the RMT model. Our concrete design demonstrates, contrary to concerns within the community, that flexible OpenFlow hardware switch implementations are feasible at almost no additional cost or power.

## Categories and Subject Descriptors

## 1. INTRODUCTION

*To improve is to change; to be perfect is to change often.*
— Churchill

Good abstractions—such as virtual memory and time-sharing—are paramount in computer systems because they allow systems to deal with change and allow simplicity of programming at the next higher layer. Networking has progressed because of key abstractions: TCP provides the abstraction of connected queues between endpoints, and IP provides a simple datagram abstraction from an endpoint to the network edge. However, routing and forwarding *within* the network remain a confusing conglomerate of routing protocols (e.g., BGP, ICMP, MPLS) and forwarding behaviors (e.g., routers, bridges, firewalls), and the control and forwarding planes remain intertwined inside closed, vertically integrated boxes.

Software-defined networking (SDN) took a key step in abstracting network functions by separating the roles of the control and forwarding planes via an *open* interface between them (e.g., OpenFlow [27]). The control plane is lifted up and out of the switch, placing it in external software. This programmatic control of the forwarding plane allows network owners to add new functionality to their network, while replicating the behavior of existing protocols. OpenFlow has become quite well-known as an interface between the control plane and the forwarding plane based on the approach known as "Match-Action". Roughly, a subset of packet bytes

[SIGCOMM'13]

The paper argues that flexibility does not come at the price of performance or cost

## Outline

- Conventional switch chips are inflexible
- SDN demands flexibility…sounds expensive…
- How do we do it: The RMT switch model
- Flexibility costs less than 15%

# Enter…

# Reconfigurable Match Tables (RMT)

## Outline

- Conventional switch chip are inflexible
- SDN demands flexibility…sounds expensive…
- <span style="color:red">How do we do it: The RMT switch model</span>
- Flexibility costs less than 15%

13

# What kind of switch architecture could support flexibility and yet run at Terabits per second?

| | |
|---|---|
| **Throughput** aggregate | 1 Tbps |
| **Packet size** average | 1000 bits |
| **# operations** per packet (avg.) | 10 |
| Requirements | 10 billion op./second |

# Pipelined architectures organize processing through a sequence of processing units and local memory

# For flexibility,
# each processing unit/memory can be made generic

# Each CPU can process distinct packets, with up to 10 packets going through the pipeline simultaneously

The runtime behavior of the parser & the match stages is defined through the RMT abstract model

# The RMT Abstract Model

- Parse graph
- Table graph

How do we implement in hardware
a programmable parser and a logical pipeline?

# How do we implement in hardware a programmable parser and a logical pipeline?

ancs48-gibb.pdf (page 1 of 12)

Search

## Design Principles for Packet Parsers

Glen Gibb[†], George Varghese[‡], Mark Horowitz[†], Nick McKeown[†]
[†]Stanford University  [‡]Microsoft Research
{grg, horowitz, nickm}@stanford.edu   varghese@microsoft.com

## ABSTRACT

All network devices must parse packet headers to decide how packets should be processed. A $64 \times 10$ Gb/s Ethernet switch must parse one billion packets per second to extract fields used in forwarding decisions. Although a necessary part of all switch hardware, very little has been written on parser design and the trade-offs between different designs. Is it better to design one fast parser, or several slow parsers? What is the cost of making the parser reconfigurable in the field? What design decisions most impact power and area?

In this paper, we describe trade-offs in parser design, identify design principles for switch and router designers, and describe a parser generator that outputs synthesizable Verilog that is available for download. We show that i) packet parsers today occupy about 1-2% of the chip, and ii) while future packet parsers will need to be programmable, this only doubles the (already small) area needed.

## Categories and Subject Descriptors

C.2.1 [**Computer-Communication Networks**]: Network Architecture and Design—*Network Communications*

## Keywords

Parsing; Design principles; Reconfigurable parsers

## 1.  INTRODUCTION

Despite their variety, *every* network device examines fields

Ethernet        IPv4              TCP            Payload
Len: 14B        Len: ?            Len: ?

Next: IPv4   Len: 20B   Next: TCP          Len: 20B

Figure 1:  A TCP packet.

In practice, packets often contain many more headers. These extra headers carry information about higher level protocols (e.g., HTTP headers) or additional information that existing headers do not provide (e.g., VLANs[1] in a college campus, or MPLS[2] in a public Internet backbone). It is common for a packet to have eight or more different packet headers during its lifetime.

To parse a packet, a network device has to identify the headers in sequence before extracting and processing specific fields. A packet parser seems straightforward since it knows *a priori* which header types to expect.

In practice, designing a parser is quite challenging:

1. **Throughput.**  Most parsers must run at line-rate, supporting continuous minimum-length back-to-back packets. A 10 Gb/s Ethernet link can deliver a new packet every 70 ns; a state-of-the-art Ethernet switch ASIC with $64 \times 40$ Gb/s ports must process a new packet every 270 ps.

[ANCS'13]

# Parsing is the (complex) process of identifying and extracting the appropriate fields in a packet header

Throughput

Parser must run at line-rate

parse 1 packet every 70 ns on a 10 Gbps link

Dependency

Parsing involves sequential processing
as headers typically point to the next one
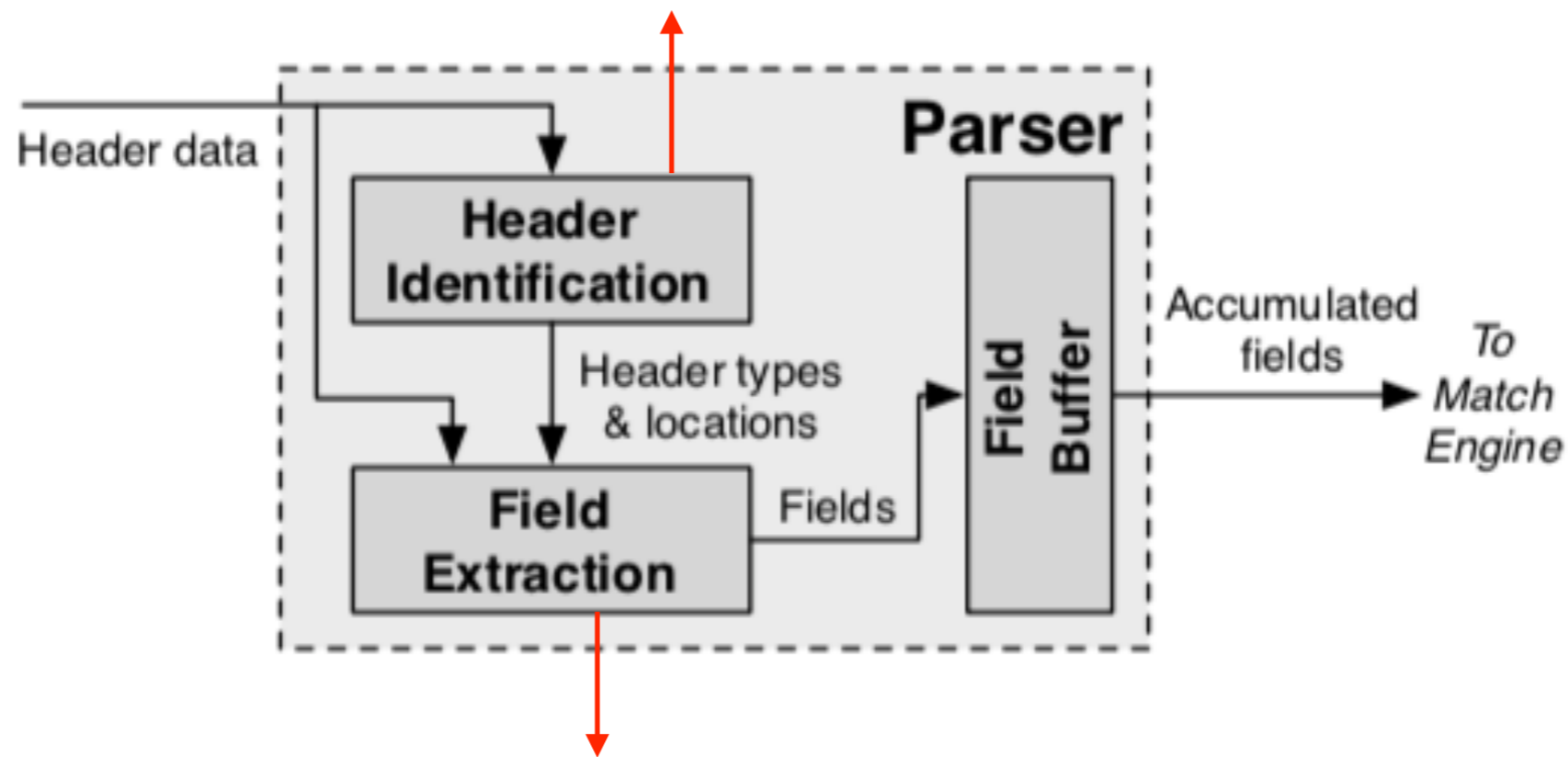
Incompleteness

Some headers do not even identify
the subsequent header

Heterogeneity

Many header formats exist that
can appear in various orders/locations

# A parser can be divided into two separate blocks:
## header identification and field extraction

implements the parse graph's
state machine



extracts the chosen fields
from identified headers

# In a programmable parser, the two modules rely on runtime information instead of hard-coded logic

stored in memory,

e.g. in RAM and/or TCAM



stores the bit sequences
that identify the headers

stores the next state,
the fields to extract,
and any other data (if any)

Source: Design Principles for Packet Parsers, Gibb et al.

How do we implement in hardware
a programmable parser and <span style="color:red">a logical pipeline</span>?

# A compiler translates a given RMT logical pipeline (specified in P4) into a physical one



RMT Logical to Physical Table Mapping

Table Graph

# The compiler maps each individual logical stage to one or more physical stage.



Table Graph

# The RMT pipeline
# in a few statistics

## Our Switch Design

- 64 x 10Gb ports
  - 960M packets/second
  - 1GHz pipeline
- Programmable parser
- 32 Match/action stages

- Huge TCAM: 10x current chips
  - 64K TCAM words x 640b
- SRAM hash tables for exact matches
  - 128K words x 640b
- 224 action processors per stage
- All OpenFlow statistics counters

28

# Building a RMT pipeline is only 15% more expensive than building a fixed-function switching pipeline

## Outline

- Conventional switch chip are inflexible
- SDN demands flexibility…sounds expensive…
- How do I do it: The RMT switch model
- Flexibility costs less than 15%

# The biggest cost is the memory…
# *not* the processing logic

## Cost of Configurability:
## Comparison with Conventional Switch

- Many functions identical:  I/O, data buffer, queueing…
- Make extra functions optional: statistics
- Memory dominates area
  - Compare memory area/bit and bit count
- RMT must use memory bits efficiently to compete on cost
- Techniques for flexibility
  - Match stage unit RAM configurability
  - Ingress/egress resource sharing
  - Table predication allows multiple tables per stage
  - Match memory overhead reduction
  - Match memory multi-word packing

# That was just an academic paper
## Let's look at a real flexible pipeline



A small subset of our lab @ITET with **two Tofino 3.2 Tbps, 32x 100 GbE QSFP28**

# That was just an academic paper
# Let's look at a real flexible pipeline



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

# Barefoot Tofino 6.5 Tbps backplane
## several billion packets per second at line rate



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

# Barefoot Tofino 6.5 Tbps backplane
## several billion packets per second at line rate

**Tofino. Simplified Block Diagram**

Reset / Clocks

PCIe

CPU MAC

DMA engines

Control & configuration

| Rx MACs 10/25/40/50/100 | Ingress Pipeline | | Egress Pipeline | Tx MAC 10/25/40/50/100 | pipe 0 |

| Rx MACs 10/25/40/50/100 | Ingress Pipeline | | Egress Pipeline | Tx MAC 10/25/40/50/100 | pipe 1 |

Traffic Manager

| Rx MACs 10/25/40/50/100 | Ingress Pipeline | | Egress Pipeline | Tx MAC 10/25/40/50/100 | pipe 2 |

| Rx MACs 10/25/40/50/100 | Ingress Pipeline | | Egress Pipeline | Tx MAC 10/25/40/50/100 | pipe 3 |

**Each pipe has 16x100G MACs + a Packet**
**Additional ports for recirculation, Packet Generator, CPU**

**BAREFOOT** NETWORKS

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

# Tofino relies on Packet Header Vector (PHV) to pass states between stages



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

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



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

# What's next?

Tofino 2: 12.8 Tbps (7 nm switching ASIC)

# This week on

# Advanced Topics in Communication Networks

P4 hardware
target

P4-based
applications

What cool things
can we do with it?

Data plane        for        Performance
programmability               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 — Princeton University
Rob Harrison — Princeton University
Marco Canini — KAUST
Nick Feamster — Princeton University
Jennifer Rexford — Princeton University
Walter Willinger — NIKSUN Inc.

# 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 — †Barefoot Networks

# 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

# 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

# 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†
† Johns Hopkins University   ‡ Carnegie Mellon University

# 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*

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

# 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
  - ....

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

# INT Header Format



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

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

# 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();
}
.....
```

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

# 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

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

# INT probes traverse multiple paths



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

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

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
  - ...

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

# Language-Directed Hardware Design for Network Performance Monitoring

# Sonata: Query-Driven Streaming Network Telemetry

Arpit Gupta
Princeton University

Rob Harrison
Princeton University

Marco Canini
KAUST

Nick Feamster
Princeton University

Jennifer Rexford
Princeton University

Walter Willinger
NIKSUN Inc.

### 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

### 1 INTRODUCTION
Network operators routinely perform continuous monitoring to track events ranging from performance impairments to attacks. This monitoring requires continuous, real-time measurement and analysis—a process commonly referred to as network telemetry [35]. Existing telemetry systems can collect and analyze measurement data in real time, but they either support a limited set of telemetry tasks [34, 46], or they incur substantial processing and storage costs as traffic rates and queries increase [7, 10, 58].

---

# 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   †Barefoot 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.

### 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 flow starting and finish times). When a flow finishes after the inactive timeout, NetFlow exports the corresponding flow records to a remote collector.

---

# 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
Many network monitoring tasks identify subsets of traffic that stand out, e.g., top-k flows for a particular statistic. A Protocol Independent Switch Architecture (PISA) switch can 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 networks modeled as one-big switch. 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.

**Figure 1:** This graph shows the recall for 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 decreases.

### 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.

---

# 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*

---

# 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†
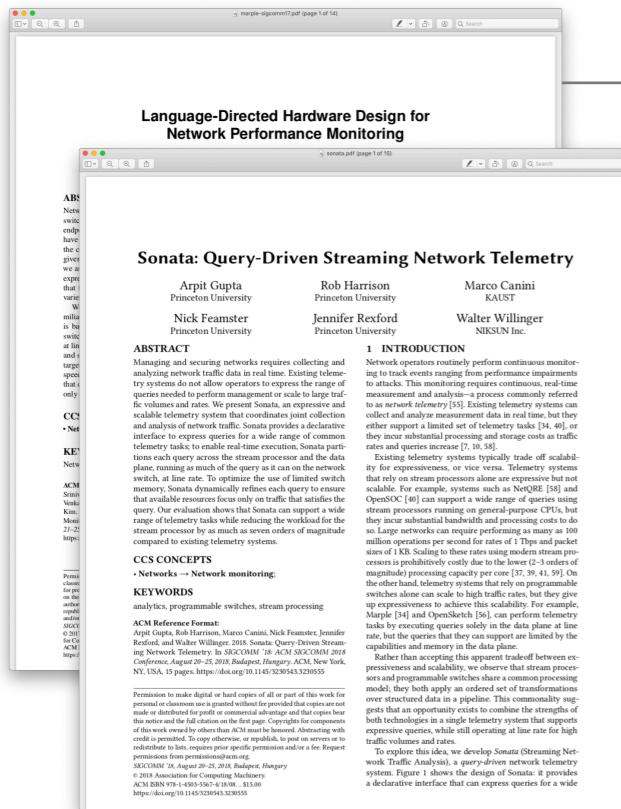† Johns Hopkins University   ‡ Carnegie Mellon University

### ABSTRACT
Network management requires accurate estimates of metrics for many applications including traffic engineering (e.g., heavy hitters), anomaly detection (e.g., entropy of source addresses), and security (e.g., DDoS detection). Obtaining accurate estimates given router CPU and memory constraints is a challenging problem. Existing approaches fall in one of two undesirable extremes: (1) low fidelity general-purpose approaches such as sampling, or (2) high fidelity but complex algorithms customized to specific application-level metrics. Ideally, a solution should be both general (i.e., supports many applications) and provide accuracy comparable to custom algorithms. This paper presents UnivMon, a framework for flow monitoring which leverages recent theoretical advances and demonstrates that it is possible to achieve both generality and high accuracy. UnivMon uses an application-agnostic data plane monitoring primitive; different (and possibly unforeseen) estimation algorithms run in the control plane, and use the statistics from the data plane to compute application-level metrics. We present a proof-of-concept implementation of UnivMon using P4 and develop simple coordination techniques to provide a "one-big-switch" abstraction for network-wide monitoring. We evaluate the effectiveness of UnivMon using a range of trace-driven evaluations and workloads, and show that it offers comparable (and sometimes better) accuracy relative to custom sketching solutions across a range of monitoring tasks.

### CCS Concepts
•Networks → Network monitoring; Network measurement;

### Keywords
Flow Monitoring, Sketching, Streaming Algorithms

### 1 Introduction
Network management is multi-faceted and encompasses a range of tasks including traffic engineering (e.g., heavy hitters), attack and anomaly detection [49], 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 [37], heavy hitters [10], entropy measures [38, 50], or detecting changes in traffic patterns [44].

MARPLE [SIGCOMM'17]

Language-Directed Hardware Design for
Network Performance Monitoring

SONATA [SIGCOMM'18]

Sonata: Query-Driven Streaming Network Telemetry

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]

**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

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

Consensus at network speed



In-Network Aggregation

(e.g., for MapReduce, graph analytics, ML)



Stateful layer-4 load balancers

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

Consensus at network speed



In-Network Aggregation

(e.g., for MapReduce, graph analytics, ML)



Stateful layer-4 load balancers

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

# 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

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

NetCache is a **rack-scale key-value store** that leverages **in-network data plane caching** to achieve

| billions QPS throughput | & | ~10 µs latency |

even under

| highly-skewed | & | rapidly-changing |

workloads.

**New generation of systems enabled by programmable switches** ☺

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



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

low throughput & high tail latency

Load →
Server →

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 ≥ backend aggregate throughput

Source: NetCache: Balancing Key-Value Stores with Fast In-Network Caching, Xin Jin, 2017

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



## NetCache: towards billions QPS key-value storage rack

Cache needs to provide the **aggregate** throughput of the storage layer

flash/disk
each: O(100) KQPS
**total: O(10) MQPS**

storage layer

cache →

in-memory
**O(10) MQPS**

cache layer

in-memory
each: O(10) MQPS
**total: O(1) BQPS**

cache →

in-network
**O(1) BQPS**

Small on-chip memory?
Only cache $O(N \log N)$ small items

Source: NetCache: Balancing Key-Value Stores with Fast In-Network Caching, Xin Jin, 2017

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 ?

# 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 ?

# NetCache Packet Format

**Existing Protocols**  **NetCache Protocol**

| ETH | IP | TCP/UDP | OP | SEQ | KEY | VALUE |

L2/L3 Routing

reserved port #

read, write, delete, etc.

❑ Application-layer protocol: compatible with existing L2-L4 layers

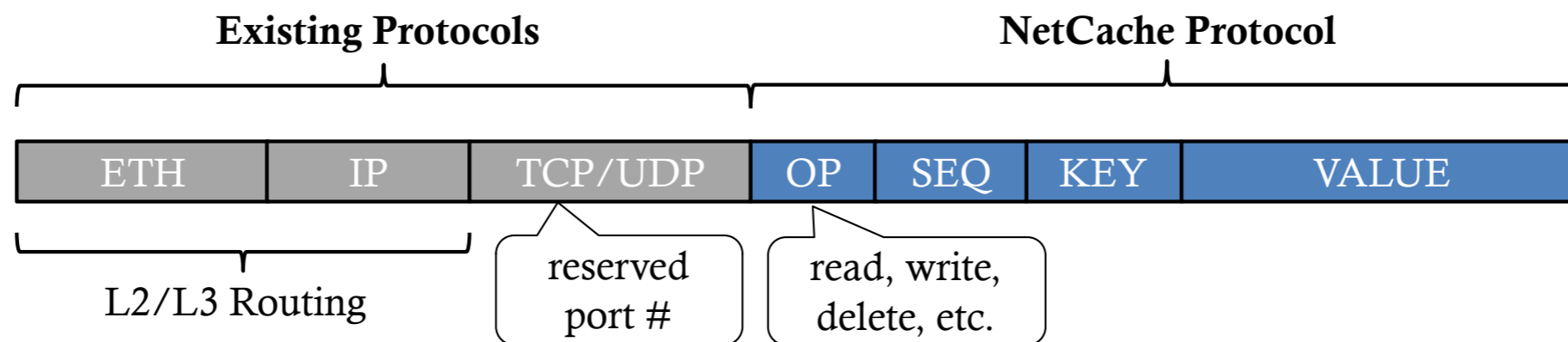❑ Only the top of rack switch needs to parse NetCache fields

# 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 ?

Source: NetCache: Balancing Key-Value Stores with Fast In-Network Caching, Xin Jin, 2017

# Key-value store using register array in network ASIC

| Match | pkt.key == A | pkt.key == B |
|-------|--------------|--------------|
| Action | **process_array**(0) | **process_array**(1) |

pkt.value: A      B

```
action process_array(idx):
  if pkt.op == read:
    pkt.value ◄── array[idx]
  elif pkt.op == cache_update:
    array[idx] ◄── pkt.value
```

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| A | B | | |

Register Array

# Variable-length key-value store in network ASIC?

| Match | pkt.key == A | pkt.key == B |
|---|---|---|
| Action | **process_array**(0) | **process_array**(1) |

pkt.value:  A          B

```
      0   1   2   3
    ┌───┬───┬───┬───┐
    │ A │ B │   │   │
    └───┴───┴───┴───┘
```
Register Array

## Key Challenges:

❑ No loop or string due to strict timing requirements

❑ Need to minimize hardware resources consumption
  - Number of table entries
  - Size of action data from each entry
  - Size of intermediate metadata across tables
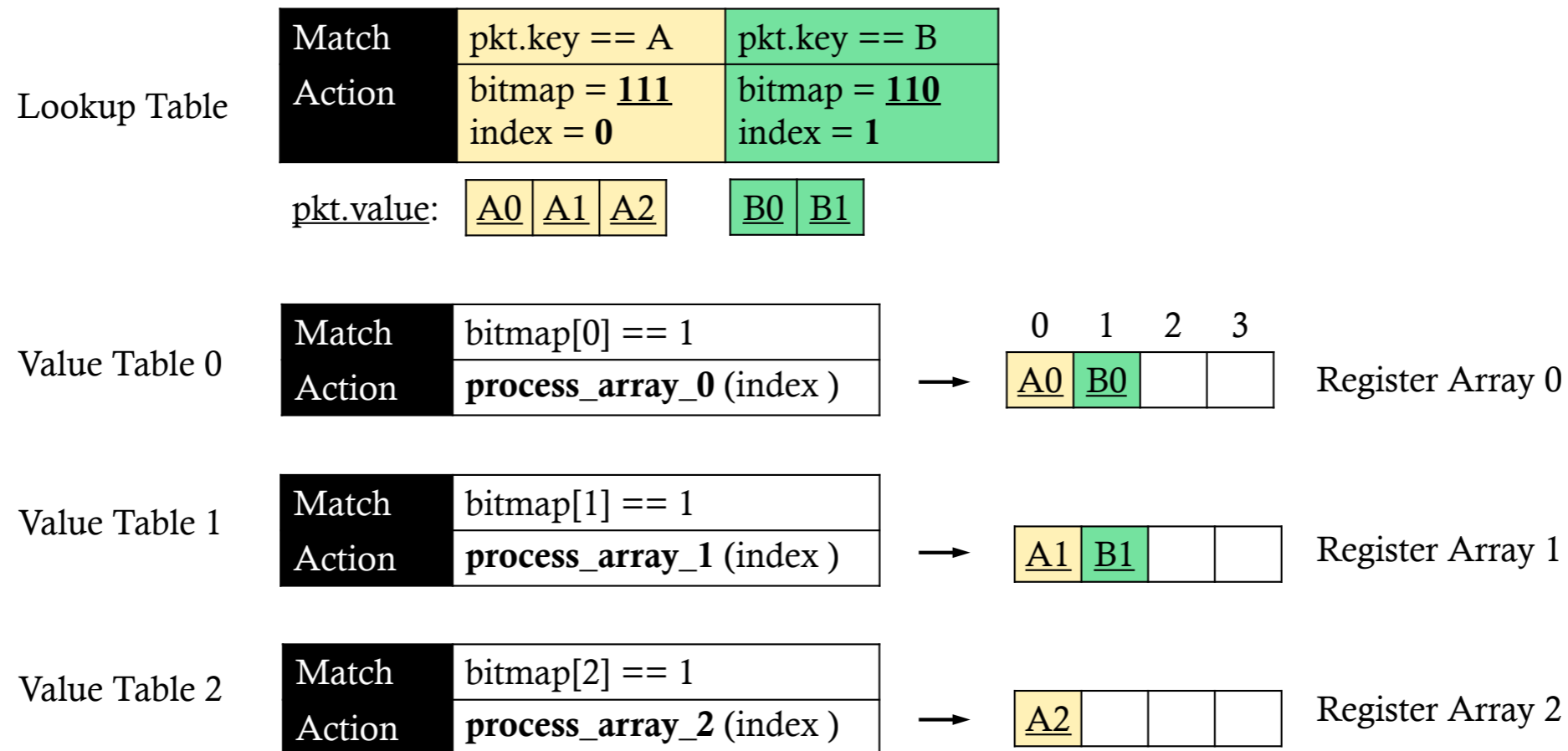
# Combine outputs from multiple arrays

**Lookup Table**

| Match | pkt.key == A |
|-------|--------------|
| Action | bitmap = **111** index = **0** |

pkt.value: | A0 | A1 | A2 |

**Bitmap** indicates arrays that store the key's value

**Index** indicates slots in the arrays to get the value

**Minimal hardware resource overhead**

**Value Table 0**

| Match | bitmap[0] == 1 |
|-------|----------------|
| Action | **process_array_0** (index ) |

0  1  2  3

| A0 | | | | Register Array 0

**Value Table 1**

| Match | bitmap[1] == 1 |
|-------|----------------|
| Action | **process_array_1** (index ) |

| A1 | | | | Register Array 1

**Value Table 2**

| Match | bitmap[2] == 1 |
|-------|----------------|
| Action | **process_array_2** (index ) |

| A2 | | | | Register Array 2

Source: NetCache: Balancing Key-Value Stores with Fast In-Network Caching, Xin Jin, 2017

# Combine outputs from multiple arrays

**Lookup Table**

| Match | pkt.key == A | pkt.key == B |
|---|---|---|
| Action | bitmap = **111** index = **0** | bitmap = **110** index = **1** |

pkt.value:  | A0 | A1 | A2 |   | B0 | B1 |

**Value Table 0**

| Match | bitmap[0] == 1 |
|---|---|
| Action | **process_array_0** (index ) |

  0   1   2   3
→ | A0 | B0 |  |  |   Register Array 0

**Value Table 1**

| Match | bitmap[1] == 1 |
|---|---|
| Action | **process_array_1** (index ) |

→ | A1 | B1 |  |  |   Register Array 1

**Value Table 2**

| Match | bitmap[2] == 1 |
|---|---|
| Action | **process_array_2** (index ) |

→ | A2 |  |  |  |   Register Array 2

# Combine outputs from multiple arrays

| Lookup Table | Match | pkt.key == A | pkt.key == B | pkt.key == C |
|---|---|---|---|---|
| | Action | bitmap = **111** <br> index = **0** | bitmap = **110** <br> index = **1** | bitmap = **010** <br> index = **2** |

pkt.value: | A0 | A1 | A2 |   | B0 | B1 |   | C0 |

| Value Table 0 | Match | bitmap[0] == 1 |
|---|---|---|
| | Action | **process_array_0** (index ) |

0  1  2  3
| A0 | B0 |   |   |   Register Array 0

| Value Table 1 | Match | bitmap[1] == 1 |
|---|---|---|
| | Action | **process_array_1** (index ) |

| A1 | B1 | C0 |   |   Register Array 1

| Value Table 2 | Match | bitmap[2] == 1 |
|---|---|---|
| | Action | **process_array_2** (index ) |

| A2 |   |   |   |   Register Array 2

# Combine outputs from multiple arrays

| Lookup Table | Match | pkt.key == A | pkt.key == B | pkt.key == C | pkt.key == D |
|---|---|---|---|---|---|
| | Action | bitmap = **111** index = **0** | bitmap = **110** index = **1** | bitmap = **010** index = **2** | bitmap = **101** index = **2** |

pkt.value:  | A0 | A1 | A2 |     | B0 | B1 |     | C0 |     | D0 | D1 |

| Value Table 0 | Match | bitmap[0] == 1 |
|---|---|---|
| | Action | **process_array_0** (index ) |

0   1   2   3

| A0 | B0 | D0 |  |   Register Array 0

| Value Table 1 | Match | bitmap[1] == 1 |
|---|---|---|
| | Action | **process_array_1** (index ) |

| A1 | B1 | C0 |  |   Register Array 1

| Value Table 2 | Match | bitmap[2] == 1 |
|---|---|---|
| | Action | **process_array_2** (index ) |

| A2 |  | D1 |  |   Register Array 2

# 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 ?

Source: NetCache: Balancing Key-Value Stores with Fast In-Network Caching, Xin Jin, 2017

# Cache insertion and eviction

❑ Challenge: cache the hottest $O(N \log N)$ items with **limited insertion rate**

❑ Goal: react quickly and effectively to workload changes with **minimal updates**



**1** Data plane reports hot keys

**2** Control plane compares loads of new hot and sampled cached keys

**3** Control plane fetches values for keys to be inserted to the cache

**4** Control plane inserts and evicts keys

Source: NetCache: Balancing Key-Value Stores with Fast In-Network Caching, Xin Jin, 2017

# Query statistics in the data plane



**Count-Min sketch**

**Bloom filter**

**Per-key counters for each cached item**

- ❑ Cached key: per-key counter array
- ❑ Uncached key
  - ▪ Count-Min sketch: report new hot keys
  - ▪ Bloom filter: remove duplicated hot key reports

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



[NSDI'18]

### congestion control



[HotNets'17]



NetFPGA SUME board

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

Big question is

**How to combine them best?**

## Beyond SmartNICs:
## Towards a Fully Programmable Cloud

*(Invited Paper)*

Adrian Caulfield
Microsoft Research
acaulfie@microsoft.com

Paolo Costa
Microsoft Research
pcosta@microsoft.com

Monia Ghobadi
Microsoft Research
mgh@microsoft.com

*Abstract*—FPGA-based SmartNICs and programmable switches have been recently introduced to leverage hardware acceleration and custom pipelines inside the cloud infrastructure. These devices are capable of handling the per-packet processing needs at line rate, including load balancing, encapsulation, congestion management, and security. We argue, however, that the benefits provided by these new devices could extend beyond software-defined networking use cases and they prompt a shift towards a fully programmable cloud, which would enable hardware-software co-design across all layers, ranging from application to hardware and networks. In this paper, we focus on the potential of FPGA-based SmartNICs and programmable switches to realize this vision and illustrate some of the research challenges that need to be addressed to fully unleash its benefits.

## I. INTRODUCTION

The continuous growth of cloud applications [1] is driving a steady increase in network infrastructure's bandwidth [2]. However, the compute cycles—measured by the number of CPU cycles required to process each packet—are falling behind the massive acceleration in available network bandwidth [3]. As a result, CPU time is increasingly becoming a contributor to the per-packet latency in high-speed cloud data centers. To make matters more challenging, modern clouds are embracing a fully software-defined network (SDN) and are increasingly expected to perform complex network policies such as regular expression matching and encryption [4]. Such policies drive up the per-packet CPU cycles, which in turn increases the cloud costs and adds unpredictable latency to the cloud services.

Traditionally, this has been addressed by offloading various networking functions to the Network Interface Cards (NICs) such as TCP Segmentation Offload [5] and Generic Receive Offload [6]. However, these techniques are not flexible enough to support complex policies. Techniques such as SR-IOV enable VMs to bypass the hypervisor and send packets directly to the NIC [7]. PCIe Process Address Space ID (PASID) reduces the hardware resource requirements of SR-IOV, enabling it to scale to support containers or even individual processes. But, these techniques bypass the hypervisor, making it hard to enforce SDN-like policies.

Recently, FPGA-based SmartNICs have been introduced as a new platform that enables network operators a flexible environment to offload complex network policies and maintain

bandwidth growth in the coming years [4], [8]–[11]. Today's FPGA-based SmartNICs are capable of tracking user-defined state, conducting basic integer arithmetic, and rate limiting at line rate [10]. Moreover, FPGAs organize computation *spatially*, hence data flows through the computation in a pipeline. This minimizes or eliminates latency jitter and provides strong guarantees about throughput. Consequently, the research community has turned its attention towards building platforms and programming languages around FPGA-based SmartNICs [4], [12]–[14].

Most prior work treats the FPGA and the NIC domains separately by either focusing on the FPGA capabilities as a generic device or focusing on the NIC functions. In this paper, we turn our attention into the combined domains and argue that an FPGA-based SmartNIC should be thought of as both a programmable FPGA-based accelerator and a networking device. When combined with recently proposed programmable switches, e.g., [15], [16], this opens up exciting opportunities to rethink the way in which we design and deploy applications and network functions. We argue that we should move away from the traditional strict boundary between network and application functions towards a fully programmable cloud, in which application logic can be distributed across multiple accelerators and network devices.

A fully programmable cloud provides significant benefits to applications that run in it. These benefits include application specific control of network flows, the ability to run code at precisely the right location in the network hierarchy, and direct, low latency access to the network. For example, in a large-scale machine learning workload, a neural network running on a distributed set of SmartNICs would benefit from the direct interface to the network to reduce inference and training latencies. Further, programmable switches running custom flow management code can reduce latency and optimize bandwidth by scheduling flows in an application-specific way, improving efficiency. Finally, the switch could even host a parameter server, directly performing aggregation of the training weights from the SmartNICs below it.

After summarizing the key technology underpinning Smart-NICs in Section II, we describe our vision underlying the fully programmable cloud in Section III. Implementing this vision, however, requires solving a number of novel and exciting research questions, which we outline in Section IV.

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?!



[SIGCOMM'18]

[SIGCOMM'18]

[CoNEXT'18]

# So you've a programmable networks…

## How do you make sure that it works as it should?!



[SIGCOMM'18]

[SIGCOMM'18]

[CoNEXT'18]

# 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?

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

**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

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

**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

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

# 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

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

# 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

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

# 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

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

Data plane programmability     for     Performance

Monitoring

Applications offloading


Platforms     for     Data plane programmability

Correctness

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]

vanbever_swing_state_sosr_2017.pdf (page 1 of 7)

## Swing State: Consistent Updates for Stateful and Programmable Data Planes

Shouxi Luo*    Hongfang Yu
University of Electronic Science and
Technology of China

Laurent Vanbever
ETH Zürich

### ABSTRACT

With the rise of stateful programmable data planes, a lot of the network functions that used to be implemented in the controller or at the end-hosts are now moving to the data plane to benefit from line-rate processing. Unfortunately, stateful data planes also mean more complex network updates as not only flows, but also the associated states, must now be migrated consistently to guarantee correct network behaviors. The main challenge is that data-plane states are maintained at line rate, according to possibly runtime criteria, rendering controller-driven migration impossible.

We present Swing State, a general state-management framework and runtime system supporting consistent state migration in stateful data planes. The key insight behind Swing State is to perform state migration entirely within the data plane by piggybacking state updates on live traffic. To minimize the overhead, Swing State only migrates the states that cannot be safely reconstructed at the destination switch.

We implemented a prototype of Swing State for P4. Given a P4 program, Swing State performs static analysis to compute which states require consistent migration and automatically augments the program to enable the transfer of these states at runtime. Our preliminary results indicate that Swing State is practical in migrating data-plane states at line rate with small overhead.

### CCS Concepts

### Keywords

Network updates; Software-Defined Networking; P4; Stateful programmable data planes.

## 1. INTRODUCTION

By enabling stateful applications to run directly *in* the data plane, at line rate, programmable data planes [9, 22, 8, 28, 27, 16, 23] have recently emerged as a promising research area.

Yet, despite making SDNs more powerful, maintaining states in the data plane also calls for new consistent update mechanisms as it prevents traditional update techniques from working, and this, for three main reasons. First, the fact that data-plane states can be updated at line rate—at speeds that can reach Tbps [5]—prevents any software-based controller from consistently moving states from one device to another. Inconsistent migration is a problem for any data-plane application that requires strong-consistency network-wide. Examples of such applications include stateful firewalls tracking dynamic flow characteristics (e.g., low-level TCP states [29]) or anomaly detection applications [21]. Second, even ignoring states dynamism, the exact set of states to be migrated may actually be unknown to the controller, preventing it from performing the migration in the first place. Indeed, the states location in memory can differ from device to device according to runtime factors (e.g. a hash computed on packet headers) that are invisible to the controller. Third, data-plane states

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

# Advanced Topics in Communication Networks

COMPLETED

Lectures/Exercices

Group project

~7 weeks

>= 7 weeks

how to program in P4

in teams of 3

# Advanced Topics in Communication Networks

Lectures/Exercices **COMPLETED**

Group project

~7 weeks

>= 7 weeks

how to program in P4

in teams of 3

The group project starts next 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 + results can be reproduced

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

implementation

70%

achieves the basic goals

is properly documented

runs + results can be reproduced

**You'll have to write
a detailed README** (in Markdown)
We'll provide you with a template

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

implementation

70%

achieves the basic goals

is properly documented

runs + results can be reproduced

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 + results can be reproduced

report

15%, 10 pages max

describes the main building blocks

evaluates the solution

describes what each group member did

presentation

15%, 10 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 16 at 23.59pm

| This week | Select a proposal from the list (adv–net.ethz.ch) |
| | or send us your own proposal by email |

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

| Mon Dec 16 11.59pm | Send us an archive with report, code, slides |

| Tue Dec 17 1.15pm— | 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

If you want to propose your own project,
send us an email describing it by <span style="color:red">Thu Oct 31 11.59am</span>

lvanbever@ethz.ch, cedgar@ethz.ch

# Quick overview of the proposals



Albert          Thomas          Roland          Alexander          Maria          Edgar

# Quick overview of the proposals
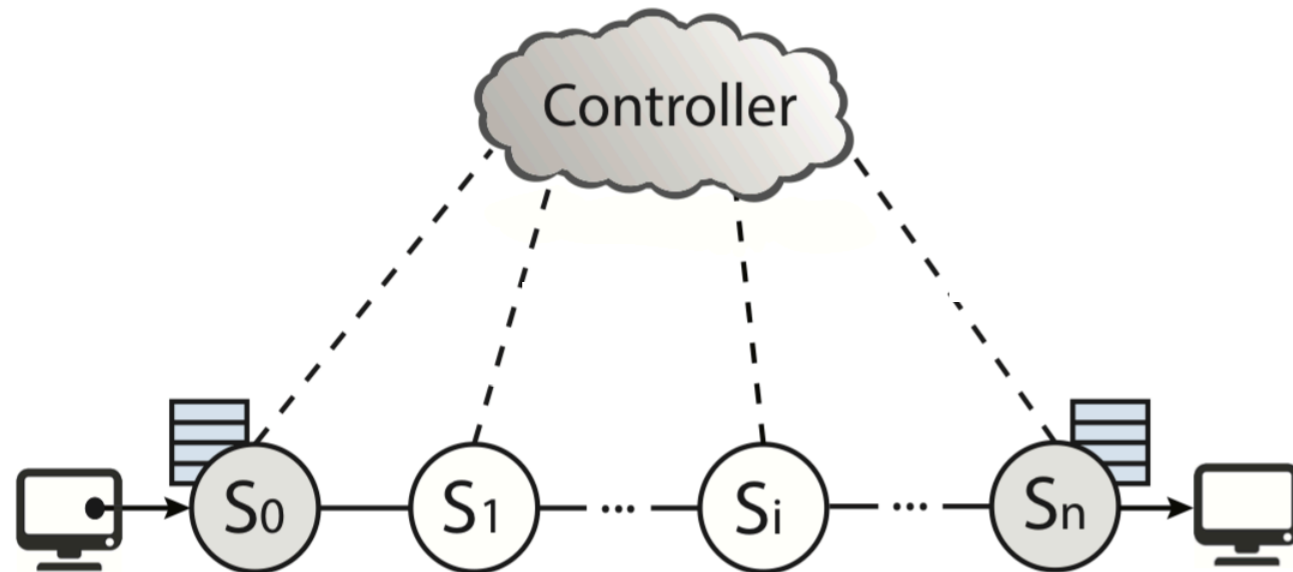


Albert      Thomas      Roland      Alexander      Maria      Edgar

# SDNSec: Forwarding Accountability for SDN Data Plane
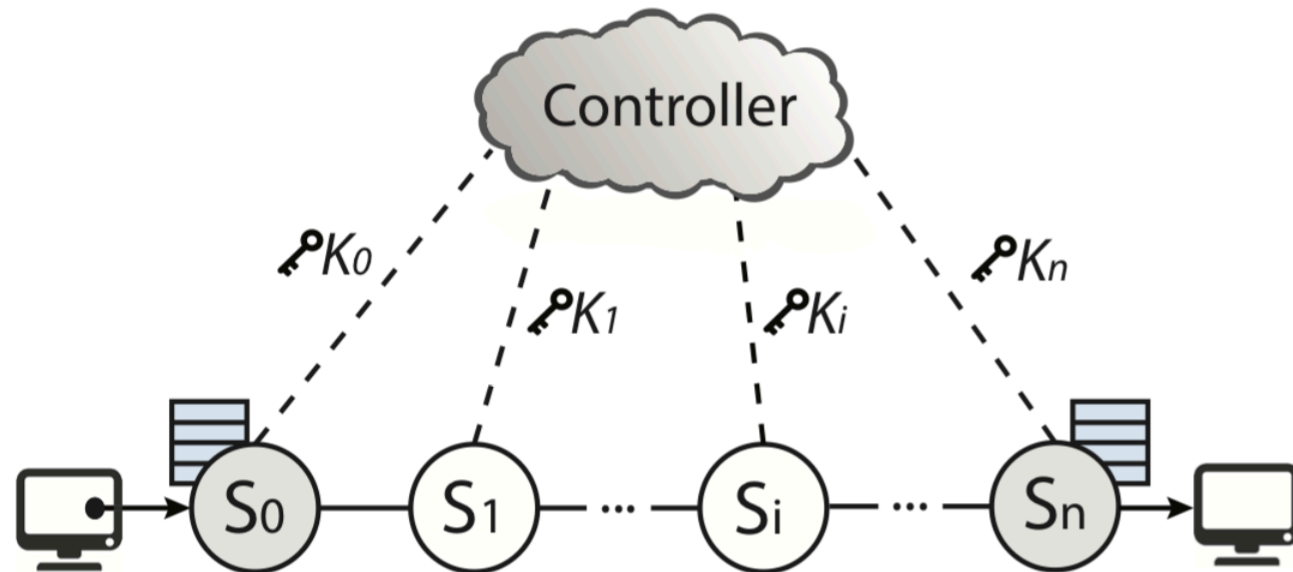


Current data plane lacks accountability:

✗ Enforcing forwarding policies

✗ Validating that policies have not been violated

✗ Consistency guarantees under reconfiguration
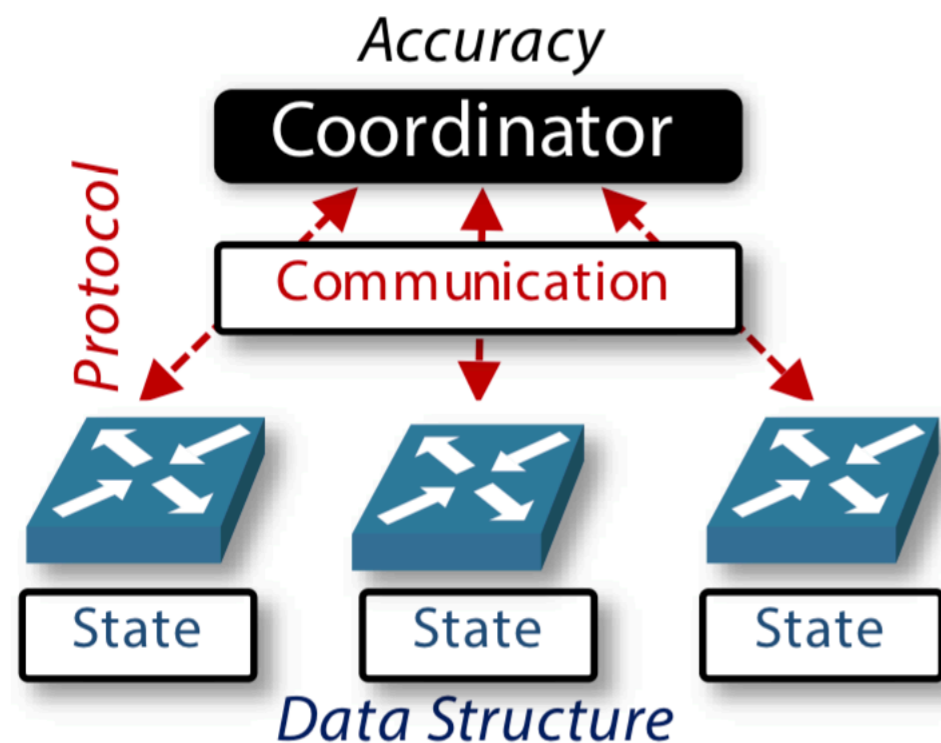
# SDNSec: Forwarding Accountability for SDN Data Plane



With SDNSec:

- ✓ Ingress-switch adds path in header

- ✓ Core-switches extract header, decrypt and forward

- ✓ Controller verifies policy

[ICCCN 2016] NEC Corporation Japan, ETH Zurich (Perrig et. al.)

# Herding the Elephants: Detecting Network–Wide Heavy Hitters with Limited Resources



Accuracy

Coordinator

Protocol
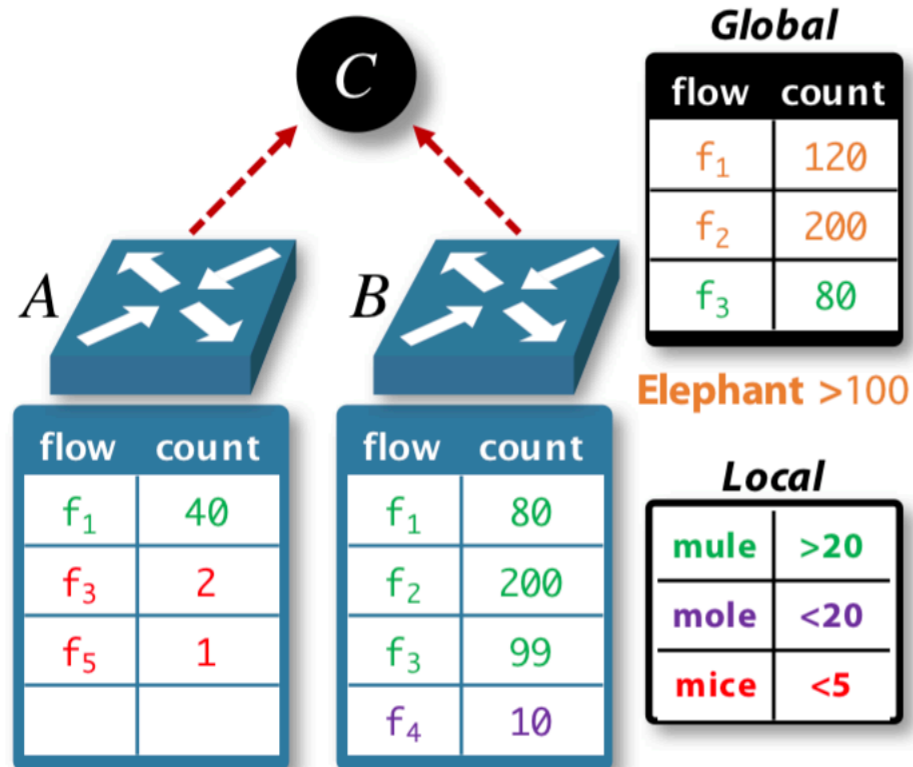
Communication

State    State    State

Data Structure

Separating elephant from mice is key in network management:

✗ Sampling is not accurate and results are delayed

✗ App-specific sketches limit network visibility

[Semantic scholar] Princeton, Walter Robert J. Harrison (Rexford et. al.)

# Proposal #2

# Herding the Elephants: Detecting Network-Wide Heavy Hitters with Limited Resources



Global

| flow | count |
|------|-------|
| f₁ | 120 |
| f₂ | 200 |
| f₃ | 80 |

**Elephant** >100

Local

| mule | >20 |
|------|-----|
| mole | <20 |
| mice | <5 |

| flow | count |
|------|-------|
| f₁ | 40 |
| f₃ | 2 |
| f₅ | 1 |
| | |

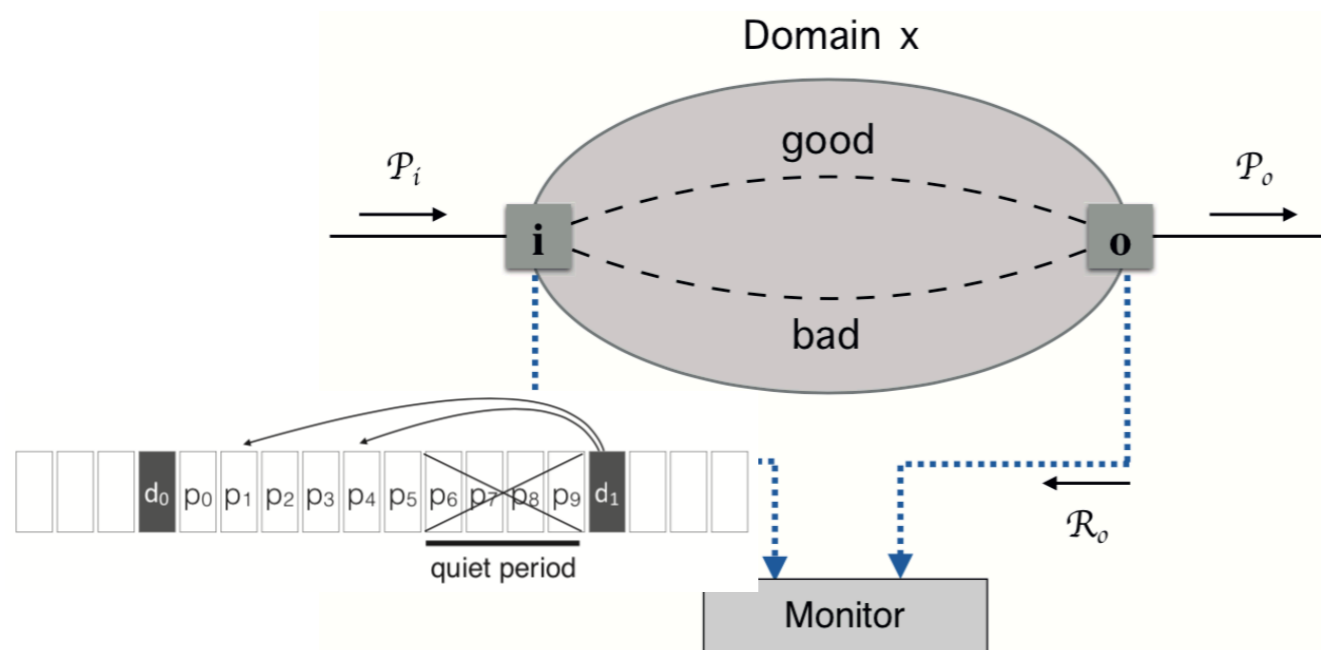| flow | count |
|------|-------|
| f₁ | 80 |
| f₂ | 200 |
| f₃ | 99 |
| f₄ | 10 |

✓ Herd provides accuracy network wide

✓ Switches allocate resources based on flow type

✓ Switches notify controller when local heavy hitter

✓ Controller finds global heavy hitters

Extension: Network-Wide Heavy Hitter Detection with Commodity Switches

[Semantic scholar] Princeton, Walter Robert J. Harrison (Rexford et. al.)

# Proposal #3
## Retroactive Packet Sampling for Traffic Receipts



✗ Network nodes could cheat in monitoring

✗ Performing better for selected samples

✓ Delayed disclosure mechanism prevents it

✓ Estimates loss-rate and delay from controller

Extension: SQR: In-Network Packet Loss Recovery

[SIGMETRICS 2019] EPFL Lausanne, ETH Zurich (Perrig et. al.)

# Quick overview of the proposals

Albert     Thomas     Roland     Alexander     Maria     Edgar

# Blink: Fast Connectivity Recovery Entirely in the Data Plane
*NSDI'19*



**TCP retransmits over time**

**Failure**

**primary**

**backup #2**

**TCP flows**

**backup #1**

# Blink: Fast Connectivity Recovery Entirely in the Data Plane
## *NSDI'19*

**TCP retransmits over time**

**Failure**

**TCP flows**

primary

backup #2

backup #1

## Goal: improving blink

1. Selecting flows with low RTTs

2. Monitoring backup next-hops continuously to reroute faster

3. Monitoring the throughput to improve accuracy

# NetCache: Balancing Key-Value Stores with Fast In-Network Caching
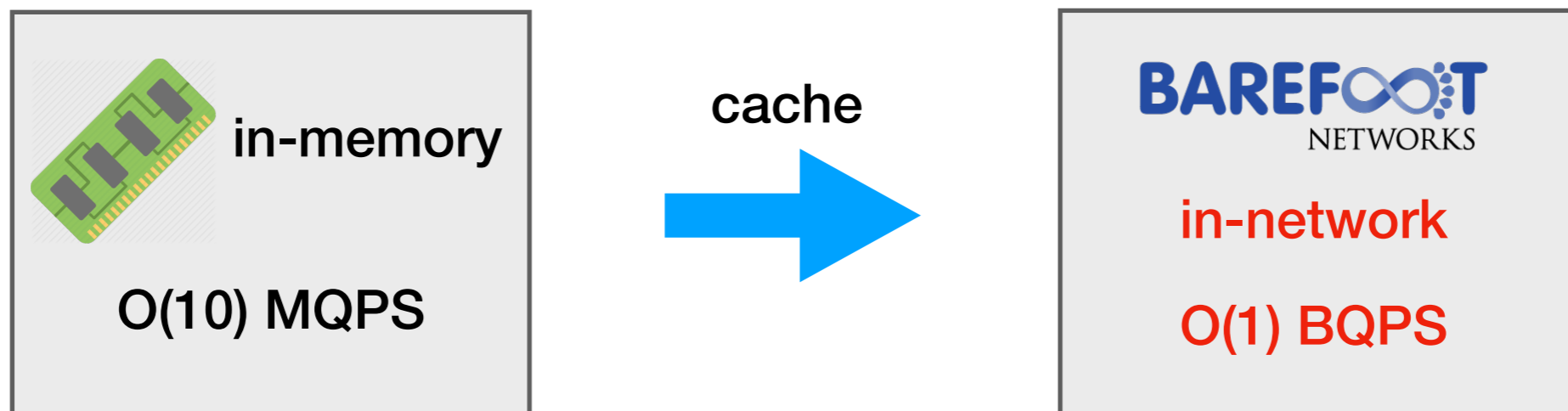*SOSP'17*   (for 2 students only)
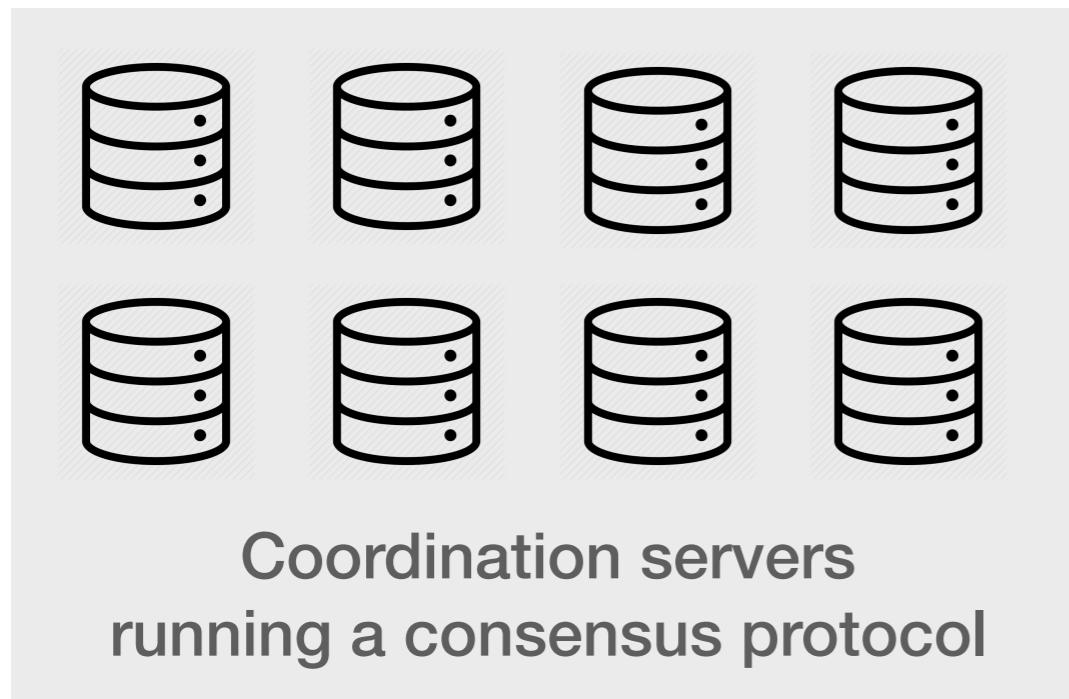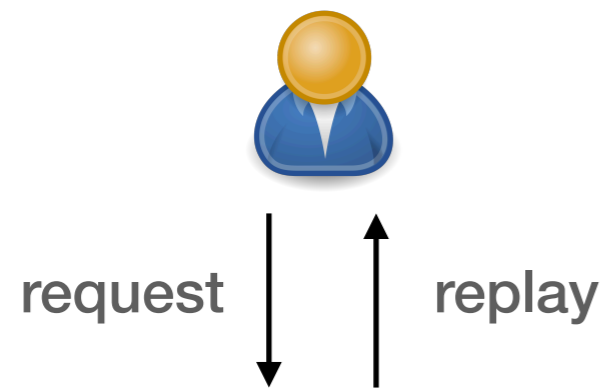
Traditional way to implement a key value store:



flash/disk

O(100) KQPS

cache

in-memory

O(10) MQPS

# NetCache: Balancing Key-Value Stores with Fast In-Network Caching
*SOSP'17*    (for 2 students only)

Traditional way to implement a key value store:

flash/disk

O(100) KQPS

cache →

in-memory

O(10) MQPS

## NetCache:

in-memory

O(10) MQPS

cache →

**BAREFOOT** NETWORKS

in-network

O(1) BQPS

Traditionally, key-value stores are replicated for *fault-tolerance*

request | replay

Coordination servers
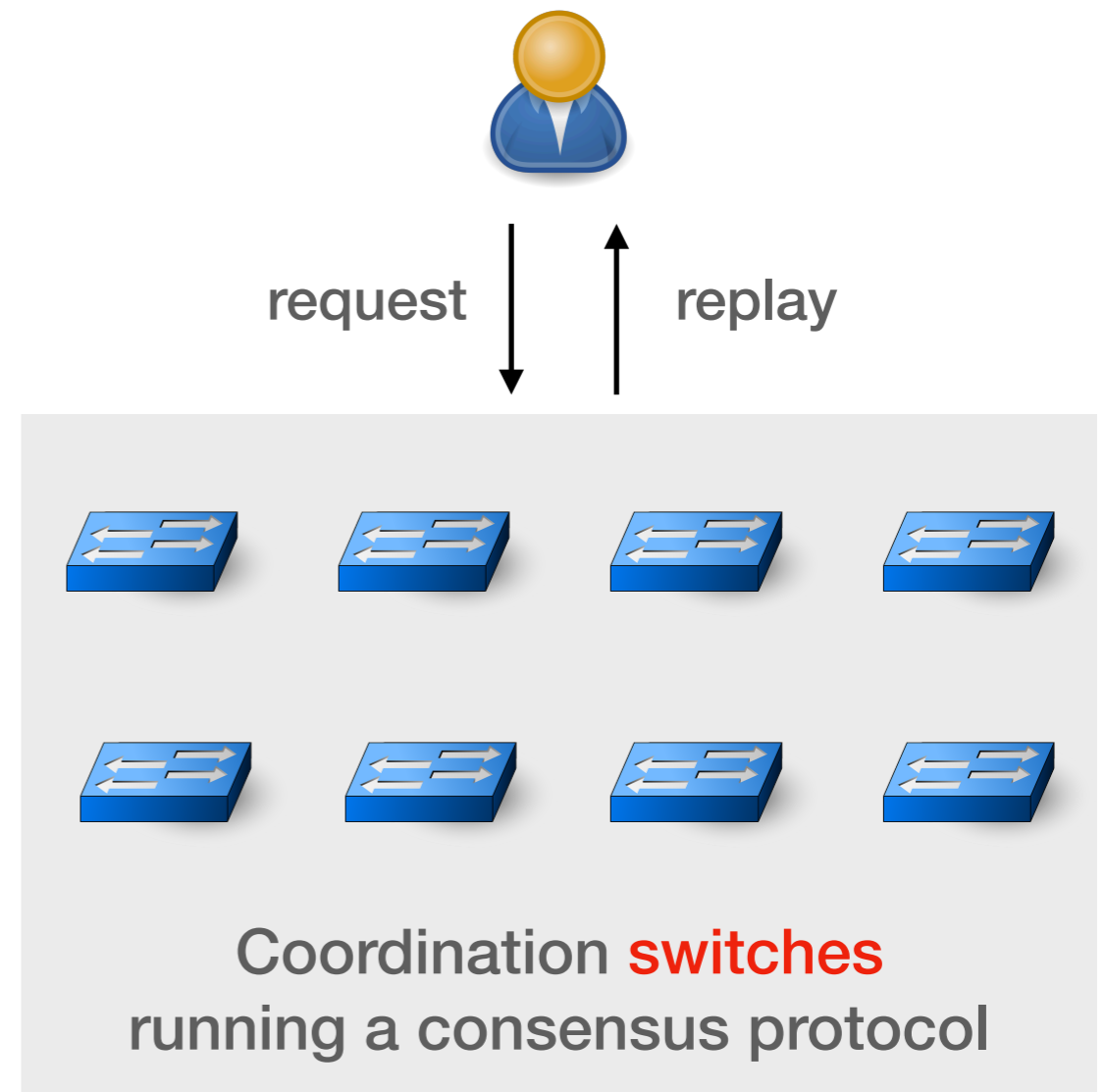running a consensus protocol

# NetChain: Scale-Free sub-RTT Coordination
*NSDI'18*    (for 2 students only)

Traditionally, key-value stores are replicated for *fault-tolerance*

NetChain

request    replay

request    replay

Coordination servers
running a consensus protocol

Coordination switches
running a consensus protocol

# Quick overview of the proposals



Albert    Thomas    **Roland**    Alexander    Maria    Edgar

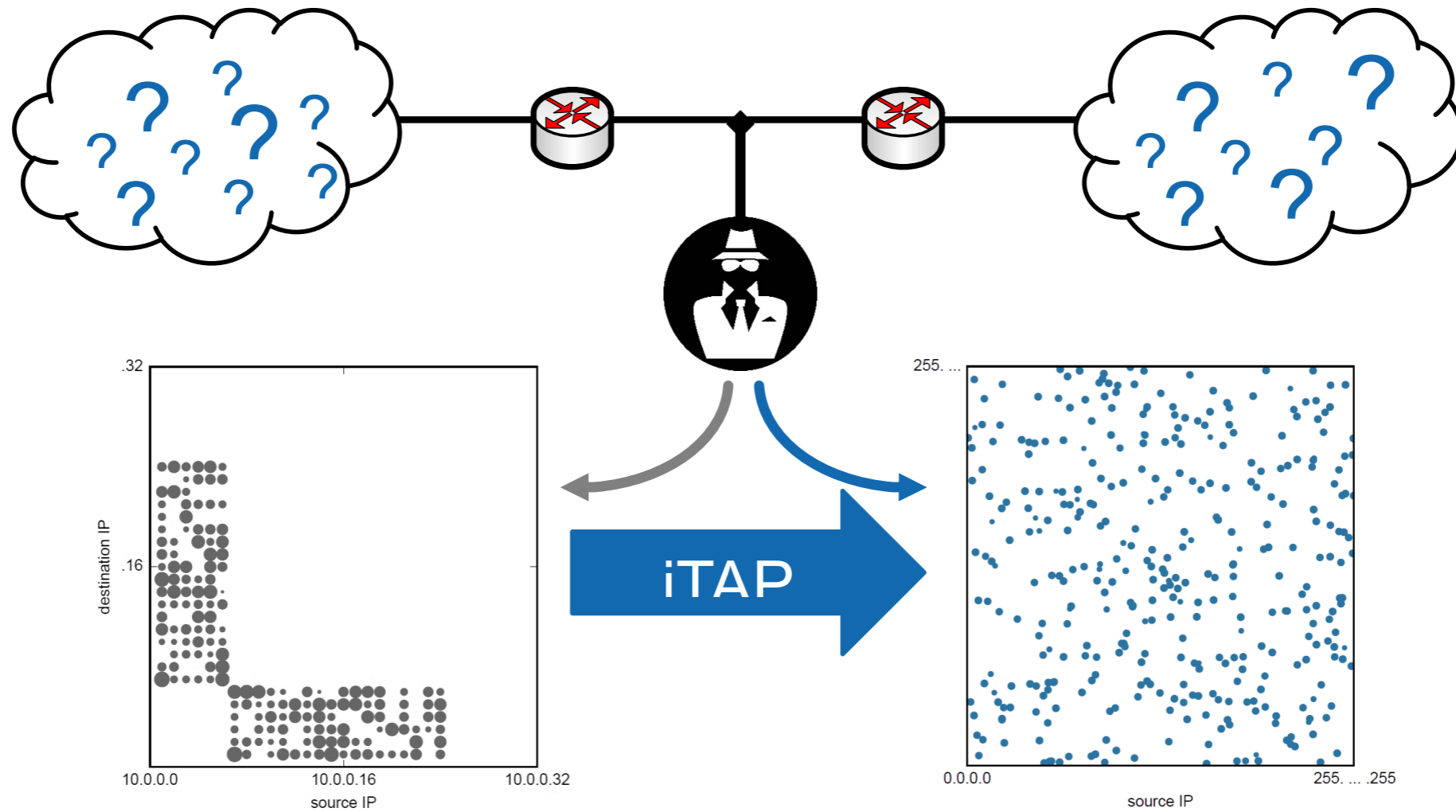# NetHide: Secure and Practical Network Topology Obfuscation



*If I receive a packet to X with TTL = i, I should send it to Y with TTL = j*

# pForest: In-Network Inference
# with Random Forests

# iTAP: In-Network Traffic Analysis Prevention Using Software-Defined Networks

# Quick overview of the proposals
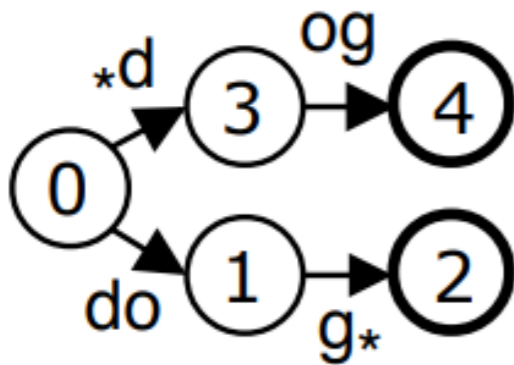
Albert     Thomas     Roland     **Alexander**     Maria     Edgar

# Proposal #4

## Fast String Searching on PISA

P4 is very limited, e.g. it cannot work with strings.
Or can it? It can even handle regular expressions!



| Match | | Action |
|-------|-------|--------|
| state | chars | |
| 0 | do | set_state(1) |
| 3 | og | accept(4) |
| 1 | g* | accept(2) |
| 0 | *d | set_state(3) |

```
$ grep P4 \
    lecture.txt
```

*In the control-plane:*
Translate regex
to automaton.

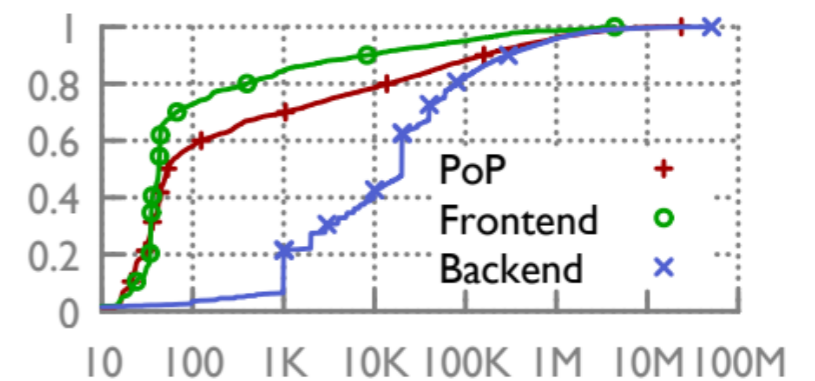*In the data-plane:*
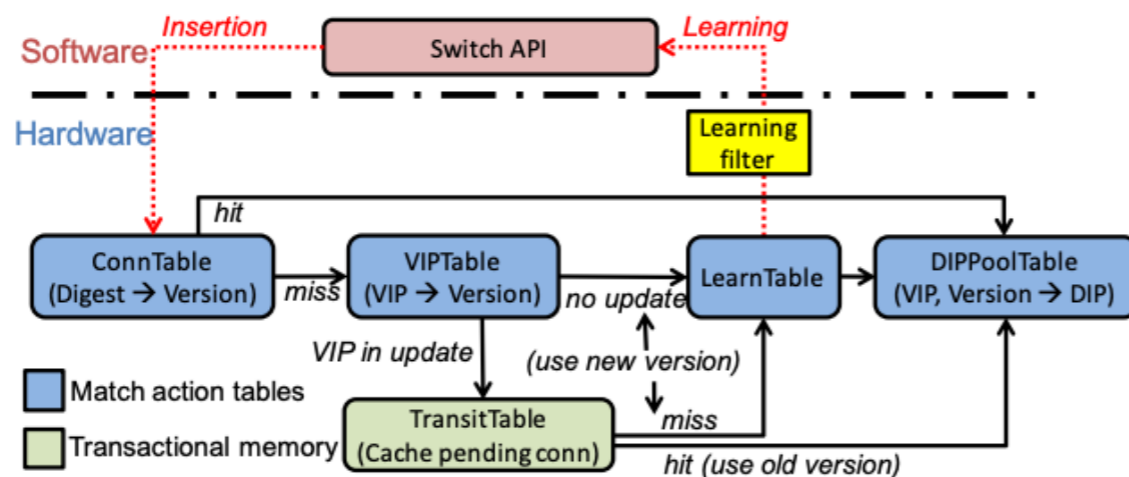Execute automaton
using recirculation.

*Evaluation:*
Compare to grep
and co.

[SOSR 2019] USI, Barefoot (Jepsen et. al.)

# SilkRoad: Making Stateful Layer-4 Load Balancing Fast and Cheap Using Switching ASICs

*SilkRoad* using a P4 switch to replace software load balancers.
It can handle millions of stateful connections using multi-level caching.



*In the control-plane:*
Accept incoming connections.

*In the data-plane:*
Keep track of existing connections.

*Evaluation:*
Test performance at large scale.

# Proposal #6

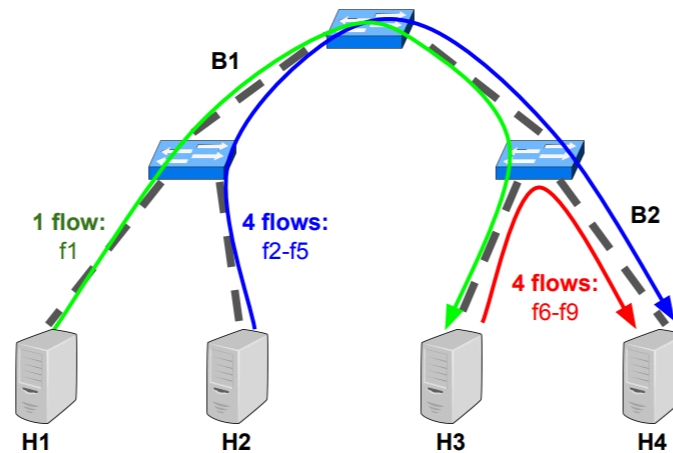# A Distributed Algorithm to Calculate
# Max-Min Fair Rates Without Per-Flow State

*s-Perc* is a congestion control algorithm that proactively assigns
per-flow sending rates without per-flow state on devices.



*In the control-plane:*

Implement the
s-Perc algorithm.
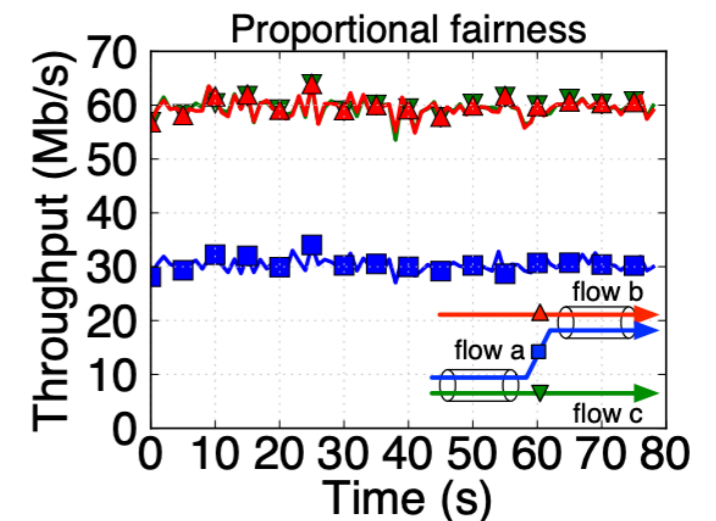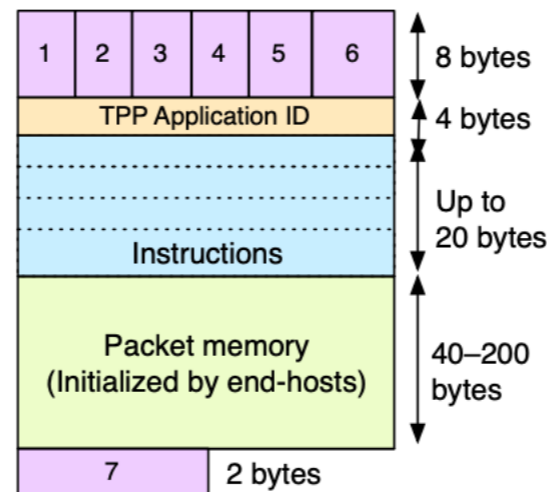
*In the data-plane:*

Create and parse
control messages.

*Evaluation:*

Compare with TCP
and other protocols.

# Millions of Little Minions: Using Packets for Low Latency Network Programming and Visibility

In active networks, packets carry programs, which are run by switches.

| Instruction |
|---|
| LOAD, PUSH |
| STORE, POP |
| CSTORE |
| CEXEC |

| | | | | | | |
|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 8 bytes |
| TPP Application ID | | | | | | 4 bytes |
| Instructions | | | | | | Up to 20 bytes |
| Packet memory (Initialized by end-hosts) | | | | | | 40–200 bytes |
| 7 | | | | | | 2 bytes |

Proportional fairness

*In the control-plane:*
Compile and start packet programs.

*In the data-plane:*
Parse packets and execute instructions.

*Evaluation:*
Test the performance of packet programs.

[SIGCOMM 2014] Stanford University, Cisco, Barefoot (Jeyakumar et. al.)

# Quick overview of the proposals

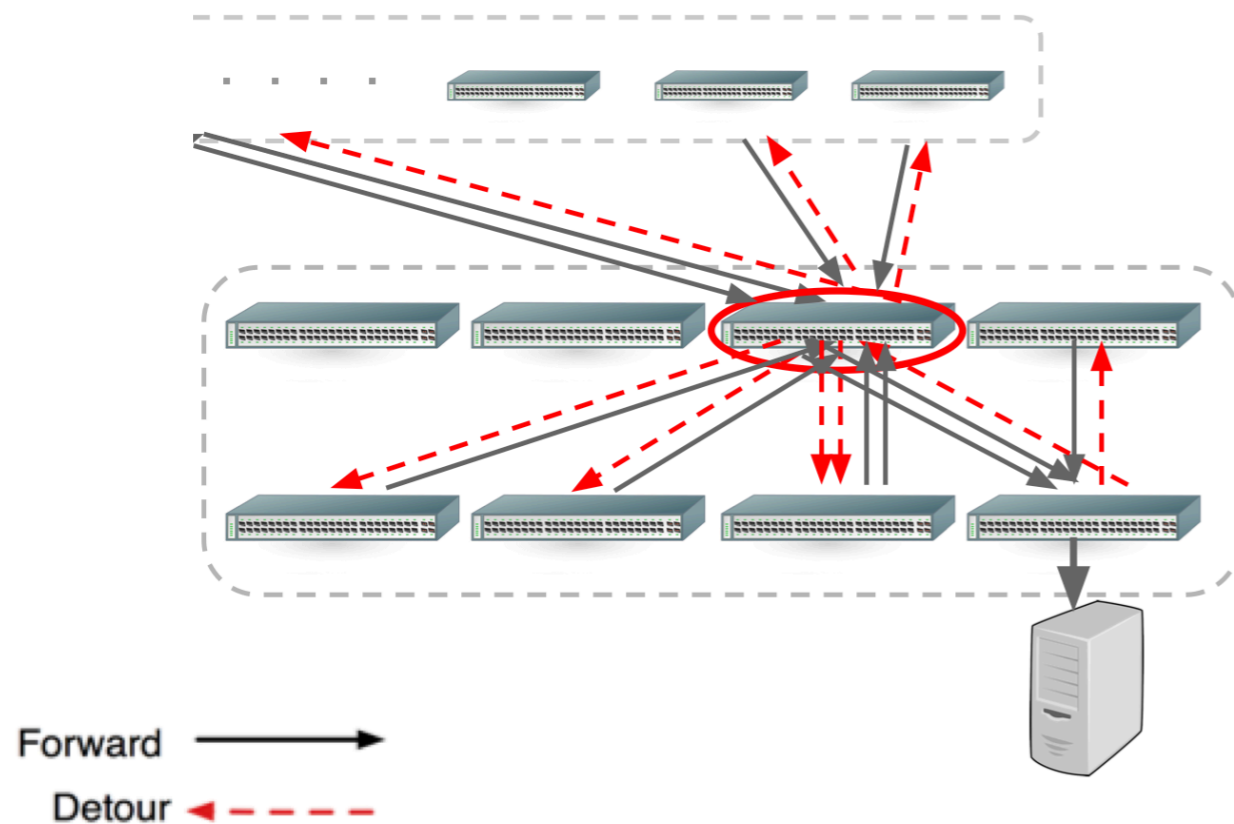Albert   Thomas   Roland   Alexander   Maria   Edgar

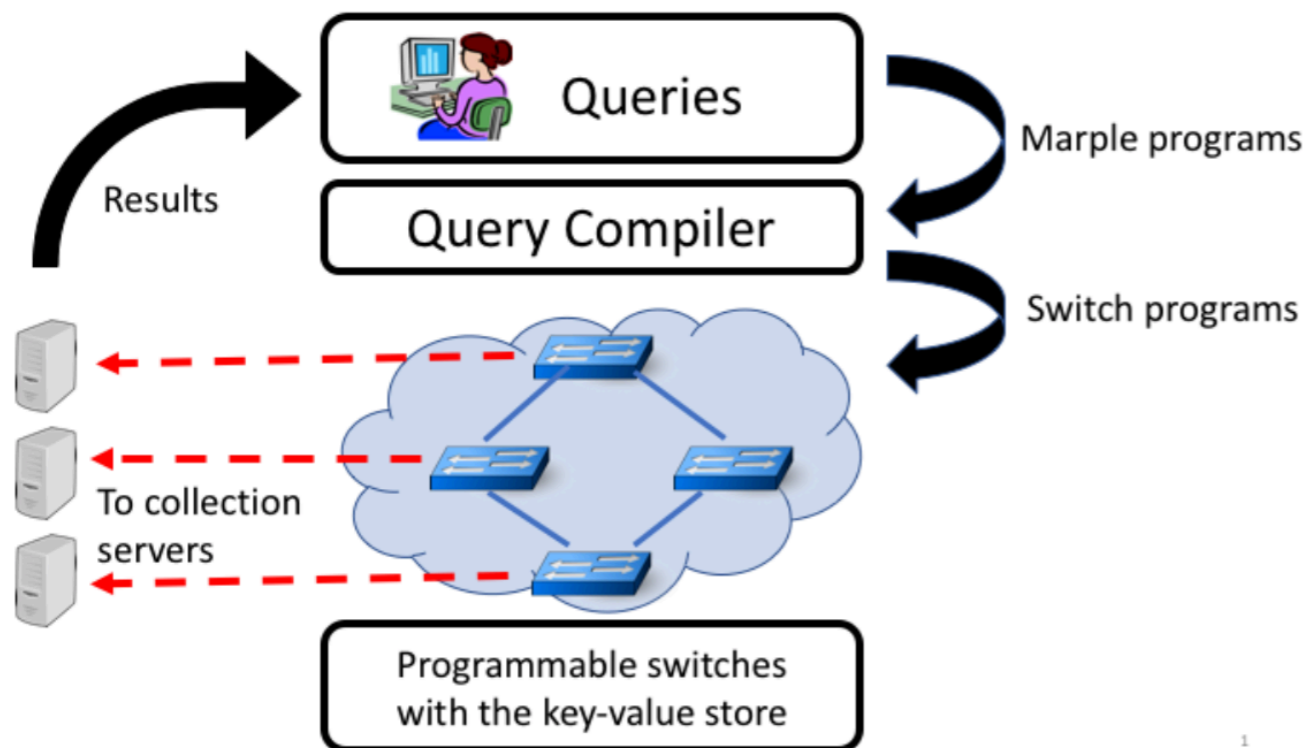# **DIBS**: Just–in–time congestion mitigation for Data Centers

**currently**

- DC patterns can cause **congestion.**
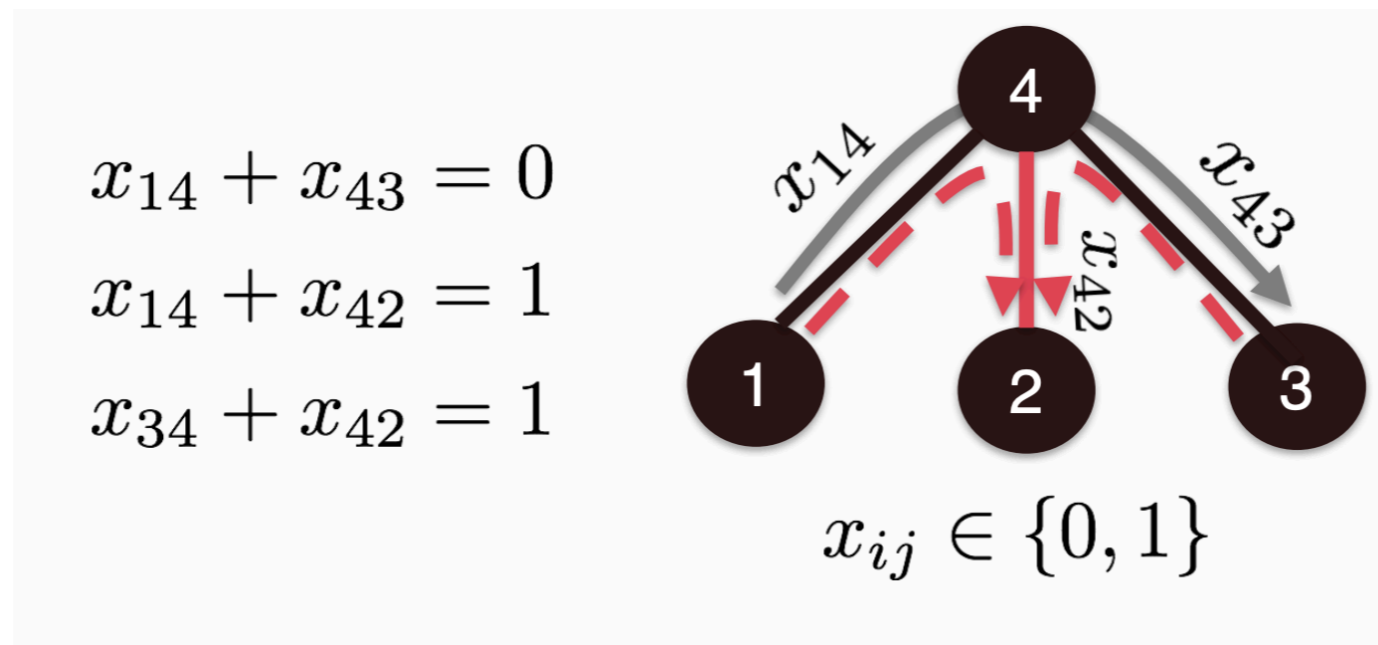
- Switches drop packets they cannot buffer.

**with DIBS**

✱ **detours** to neighboring switches.

✱ minimizes drops, which speeds up **job completion time.**

Forward →

Detour ◄- - -

# **Marple**: Language–Directed Hardware Design for Network Performance Monitoring



Results

Marple programs

Switch programs

Queries

Query Compiler

To collection servers

Programmable switches with the key-value store

* The operator writes a query in a domain–specific language called Marple.

* The query is compiled into a switch program that runs on the network's programmable switches, augmented with new switch hardware primitives that we design in service of Marple.

* The switches stream results out to collection servers, where the operator can retrieve query results.

# 007: Democratically Finding the Cause of Packet Drops



$$x_{14} + x_{43} = 0$$
$$x_{14} + x_{42} = 1$$
$$x_{34} + x_{42} = 1$$

$$x_{ij} \in \{0, 1\}$$

**Need to detect short-lived & concurrent failures despite noise**

✽007 scales by uses traceroute to find paths of flows that had packet drops

✽007 finds faulty links democratically democracy by letting hosts vote

**Implementation with p4 switches.**

✽detect retransmissions in switches

✽issue traceroutes directly from data plane

✽combine traceroutes in control plane

# Quick overview of the proposals
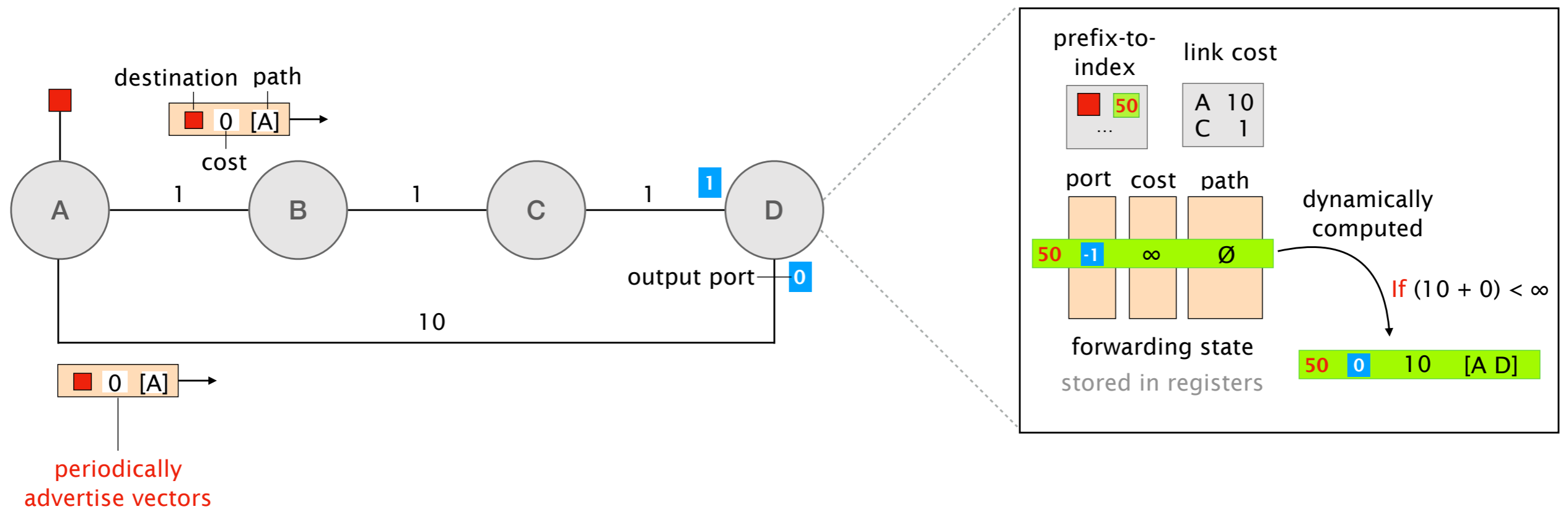


Albert   Thomas   Roland   Alexander   Maria   **Edgar**

# Hardware–Accelerated Network Control Planes

Modern programmable devices can perform small computations on **billions** of small packets per second.
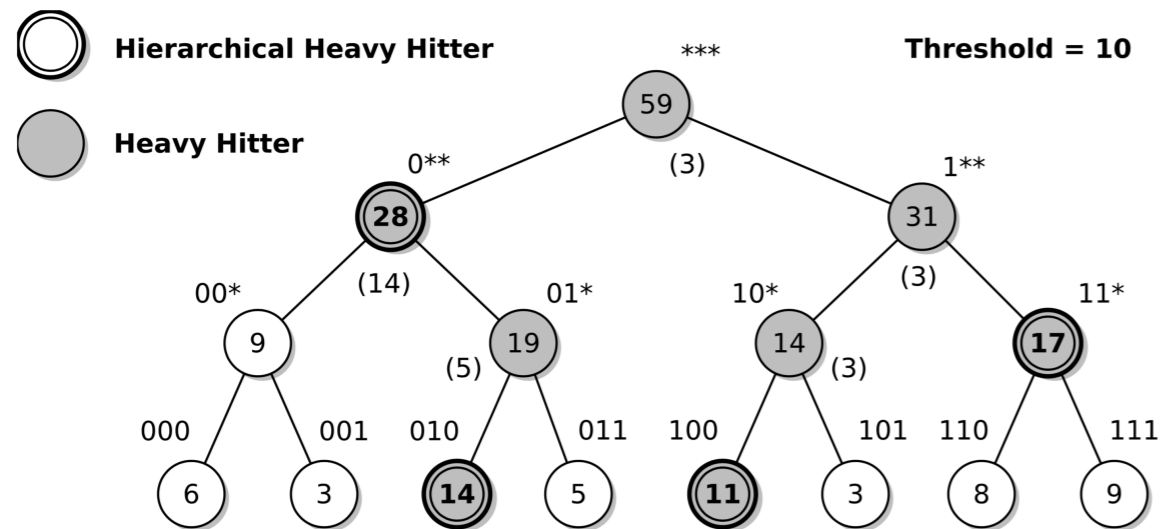
This paper shows how to leverage that to run **control plane** algorithms directly in the **data plane**

# Seek and Push: Detecting Large Traffic Aggregates in the Dataplane
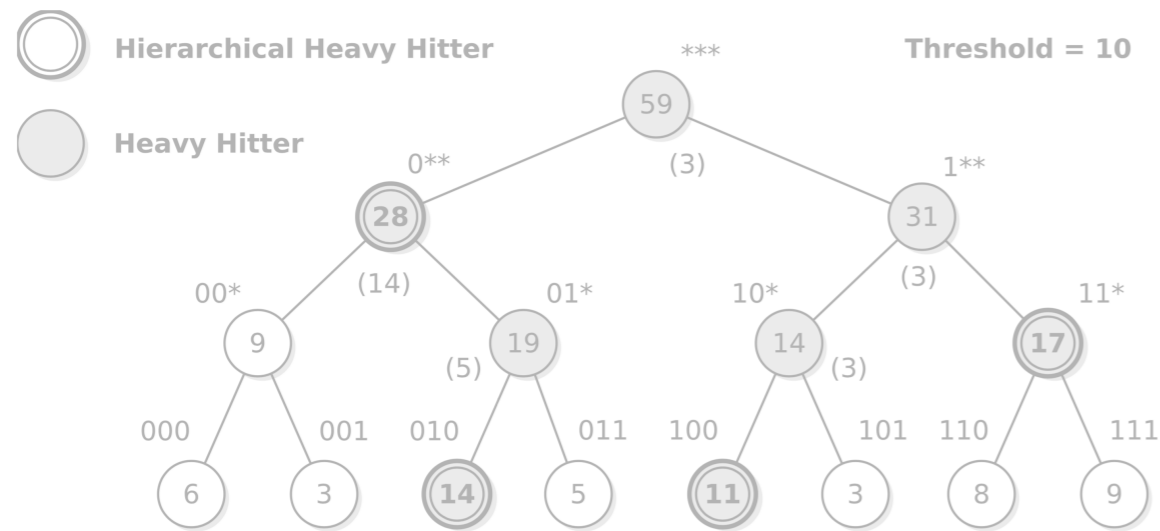
[arXiv 2018] CESNET, Cambridge (Kučera et. al.)

They present a data structure called *Elastic Tire*
that is able to detect: heavy hitters, traffic shifts
and superspreaders.

# Seek and Push: Detecting Large Traffic Aggregates in the Dataplane
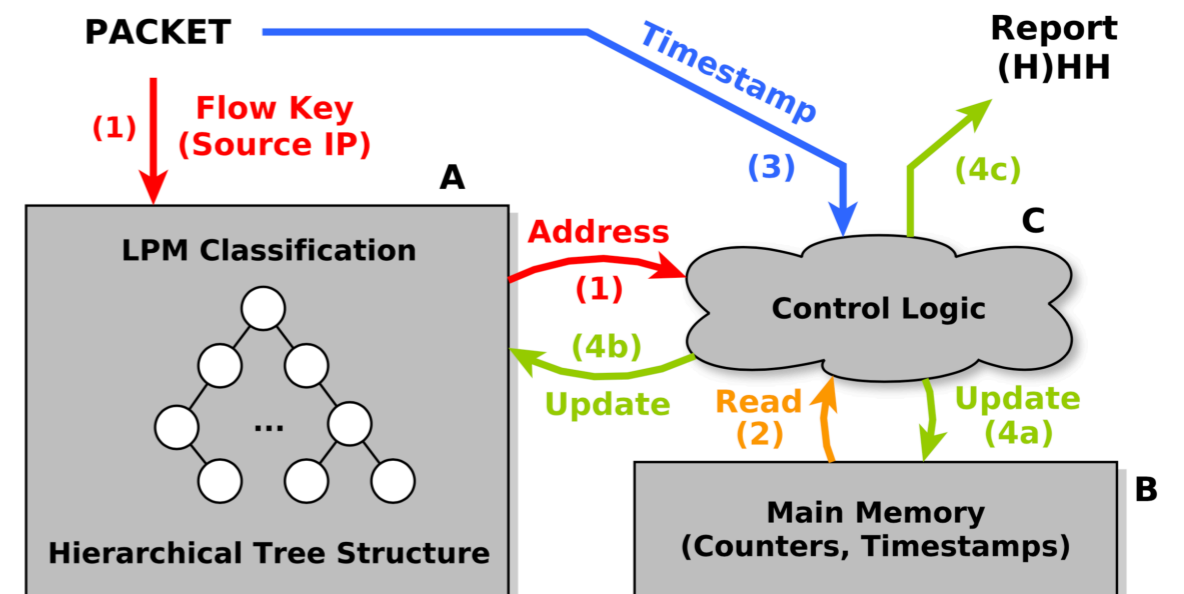
[arXiv 2018] CESNET, Cambridge (Kučera et. al.)

They present a data structure called *Elastic Tire*
that is able to detect: heavy hitters, traffic shifts
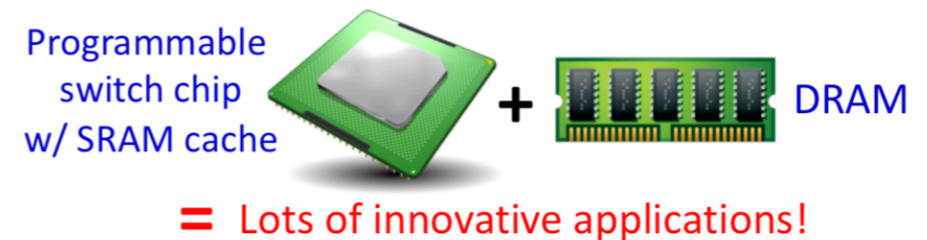and superspreaders.

High-level architecture:
1. Matching the flow using a dynamic LPM tree
2. Update Statistics
3. Control logic to update or report

# Generic External Memory for Switch Data Planes

Programmable switches are flexible but only have a limited on-ship SRAM and TCAMS

Programmable switch chip w/ SRAM cache

**+**

DRAM

**=** Lots of innovative applications!

# Generic External Memory for Switch Data Planes

Programmable switches are flexible but only have a limited on-ship SRAM and TCAMS

Programmable switch chip w/ SRAM cache + DRAM

= Lots of innovative applications!

Leverage RDMA to access remote memories at minimal latency and CPU usage

**Switch Control Plane**
RDMA channel controller

**Switch Data Plane**
On-chip registers

Port ...   ...   ...

Incoming/outgoing packets

DRAM
| Action | Packet |
| Action | Packet |
| Action | Packet |
RDMA NIC

Memory access req./resp.
RDMA initialization

# Generic External Memory for Switch Data Planes

Packet buffer extension

Remote buffer servers

Extending Lookup Tables

Customers' bare-metal servers

Remote table server

Customers' VMs

Extending State for network monitoring

Remote telemetry servers
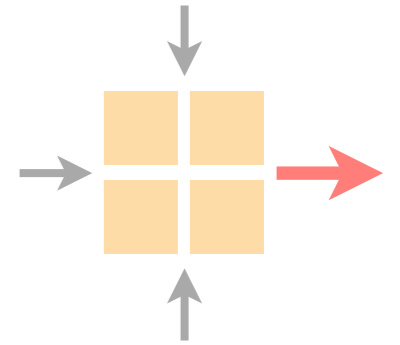
Advanced Topics in Communication Networks

# Programming Network Data Planes

Laurent Vanbever

nsg.ee.ethz.ch

ETH Zürich

Oct 29 2019