# Use cases of Programmable Dataplane (P4)

## ECE/CS598HPN

*Radhika Mittal*

# Which paper(s) did you read?

- (A) BeauCoup: Network Monitoring

- (B) Elmo: Multicast

- (C) Both

- (D) Neither

# Network Monitoring

- Most popular usecase of programmable dataplanes.

- Lots of recent papers!

- Key challenges:
  - Dealing with small amount of memory.
  - Ensuring high line rate (small processing capability, limited memory access)
  - Supporting a wide variety of queries.
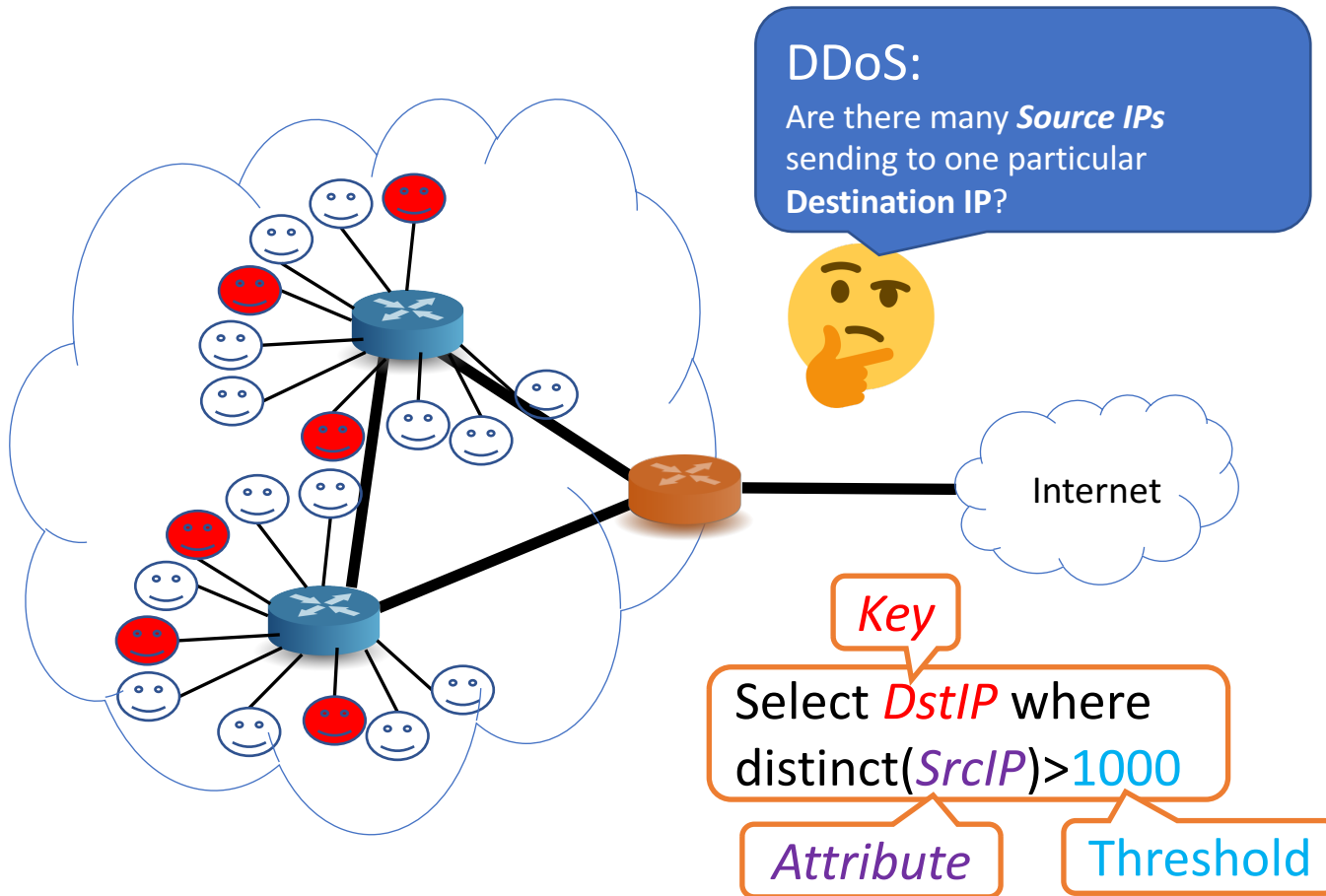
[1][bo'ku] *Adv*. many, a lot.

# *BeauCoup:*[1]

Answering *many* network traffic queries,

*one* memory update at a time!

**Xiaoqi Chen,** Shir Landau-Feibish, Mark Braverman, Jennifer Rexford

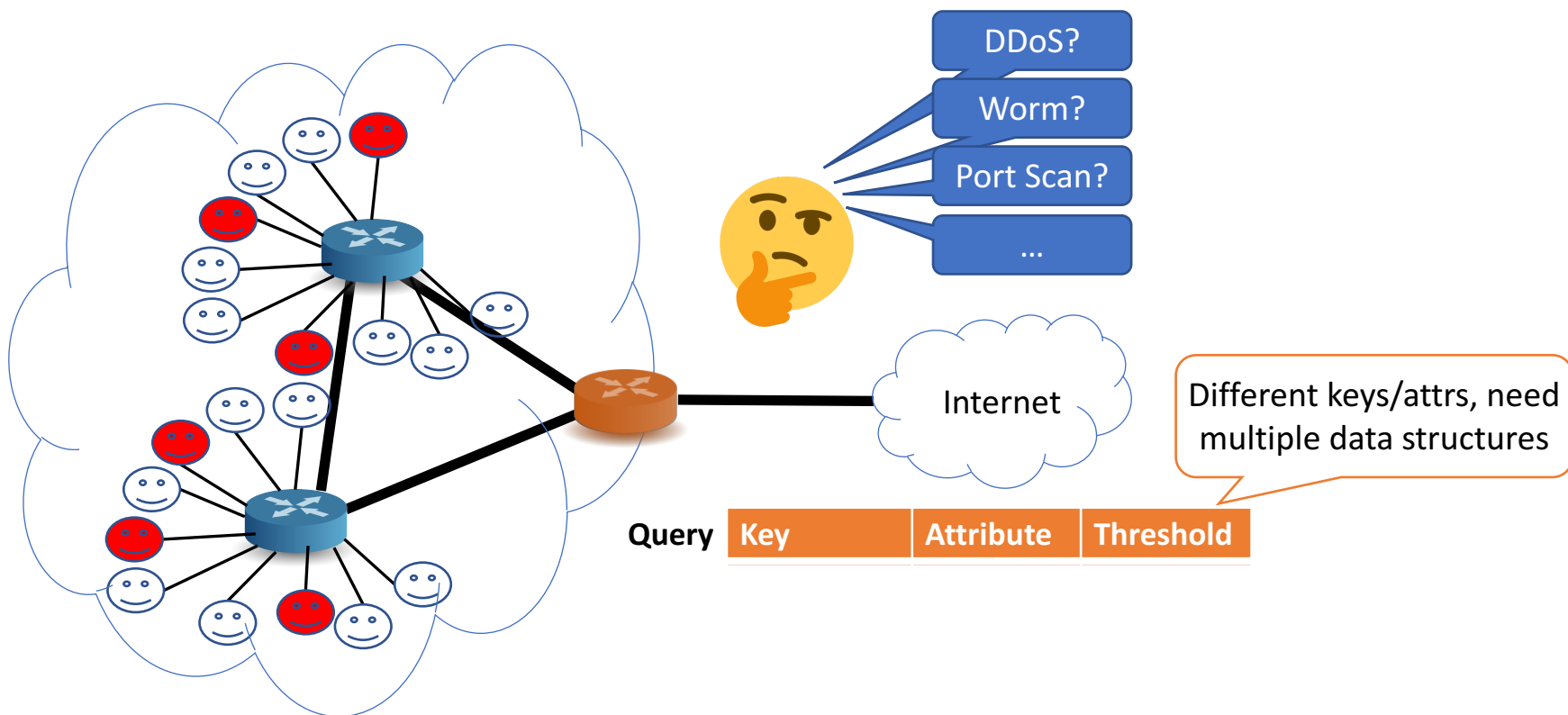SIGCOMM 2020
new york city
USA

PRINCETON
UNIVERSITY

*Slides from Xiaoqi Chen*

# Network traffic query

DDoS:
Are there many **Source IPs** sending to one particular **Destination IP**?

Internet

*Key*

Select *DstIP* where distinct(*SrcIP*)>1000

*Attribute*

Threshold

# *Many* network traffic queries



DDoS?

Worm?

Port Scan?

...

Internet

Different keys/attrs, need multiple data structures

Query | Key | Attribute | Threshold

# *Many* network traffic queries

I have 42 queries

Run 42 data structures?

I can't…

*Spec for today's commodity programmable switch:*

- **XX Tbps** aggregated throughput

- **YY MB** data-plane memory

- **Can only access ZZ bytes of memory per packet**

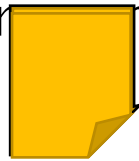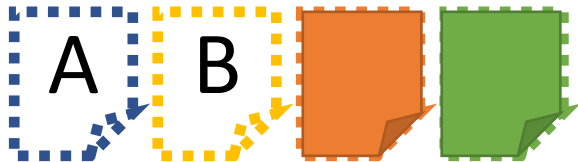# *One* memory update at a time?

- Constant memory update per packet, regardless of the number of queries?

- Game plan:
  1. Each query uses only **o(1)** memory update per packet ***on average***
  2. Combine many different queries, on average uses **O(1)**
  3. Coordinate, ***at most* O(1)** per packet

# BeauCoup's Approach

- Challenge:
  many queries, few memory updates

- Achieving **o(1)** memory access:
  coupon collectors

# The coupon collector problem

- 4 different coupons, collect all of them

- Random draws

- How many total draws are required?

# Naïve Approach

Query: Select *DstIP* where distinct(*SrcIP*)>130

- Map each ScrIP to a coupon
  - How many total coupons?
  - How many do you need to collect?

- Issues with this approach:
  - Too much memory
  - Each packet results in a coupon collection.
    - Exceed O(1) access when multiple such queries are combined.

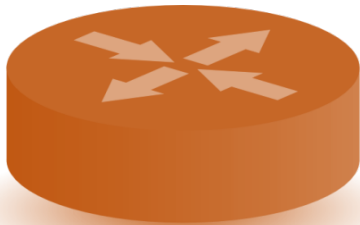# *BeauCoup* coupon collector

$f(SrcIP) \rightarrow$ Coupon ?

Mapping

Select *DstIP* where distinct($SrcIP$)>100

Collect different coupons

Key: 162.249.4.107
Coupons: A 🟨 🟧 D

- $f(10.0.1.15) \rightarrow$ Coupon C
- $f(10.0.1.33) \rightarrow$ Coupon B
- $f(10.0.1.15) \rightarrow$ Coupon C
- $f(10.0.1.42) \rightarrow$ No Coupon

# *BeauCoup* coupon collector

$f(\textit{SrcIP})$ -> Coupon ?

- Generalization: $(\textbf{\textit{m}}, \textbf{\textit{p}}, \textbf{\textit{n}})$-coupon collector
- $\textbf{\textit{m*p}} < 1$, most packets collect no coupon

Example:
(*m*=8, *p*=1%, *n*=4)

Given a new SrcIP, each coupon is drawn with probability *1%*

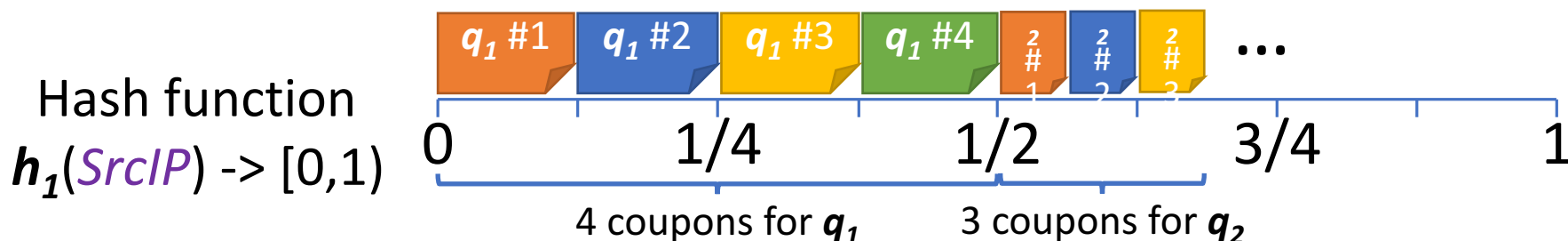*m=8* coupons in total

stop at *n=4* different coupons

# Stacking queries: same attribute

$q_1$: $f$(*SrcIP*) -> Coupon

$m_1$=4, $p_1$=1/8

$q_2$: $f$(*SrcIP*) -> Coupon

$m_2$=3, $p_2$=1/16

⋮

Hash function
$h_1$(*SrcIP*) -> [0,1)



| $q_1$ #1 | $q_1$ #2 | $q_1$ #3 | $q_1$ #4 | $q_2$ #1 | $q_2$ #2 | $q_2$ #3 | ... |

0          1/4          1/2          3/4          1

4 coupons for $q_1$          3 coupons for $q_2$

# One hash function for each attribute

$q_1$: $f$(SrcIP) -> Coupon

$m_1$=4, $p_1$=1/8

$q_6$: $g$(DstIP) -> Coupon

$m_6$=3, $p_6$=1/8

$h_1$(SrcIP) ->

| $q_1$ #1 | $q_1$ #2 | $q_1$ #3 | $q_1$ #4 | ... |

0    1/4    1/2    3/4    1

$h_2$(DstIP) ->

| $q_6$ #1 | $q_6$ #2 | $q_6$ #3 | ... |

0    1/4    1/2    3/4    1

Randomly break ties if a coupon needs to be collected for two different attributes

# System design

- Query compiler: finds coupon collector configurations
  - Stops near query thresholds, minimize error
  - Hardware limits (e.g., memory access limit)
  - Fairness across queries
- Data plane program: collect coupons into in-memory table
  - Simultaneously run *many* queries
  - At most *one* coupon per packet
  - Update queries on-the-fly

# Query compiler

Query set
$Q = \{q_1, q_2, ...\}$

Total memory update
limit: $\Gamma$ per packet

Per-query limit:
$\gamma_q$ per packet

Query $q_i$
Key, Attribute,
Threshold

Compiler

$\gamma_q = \Gamma / |Q|$
(fair allocation)

$q_i$'s Collector Configuration
Total coupons:       $m$
Each probability:    $p$
Coupons to collect:  $n$
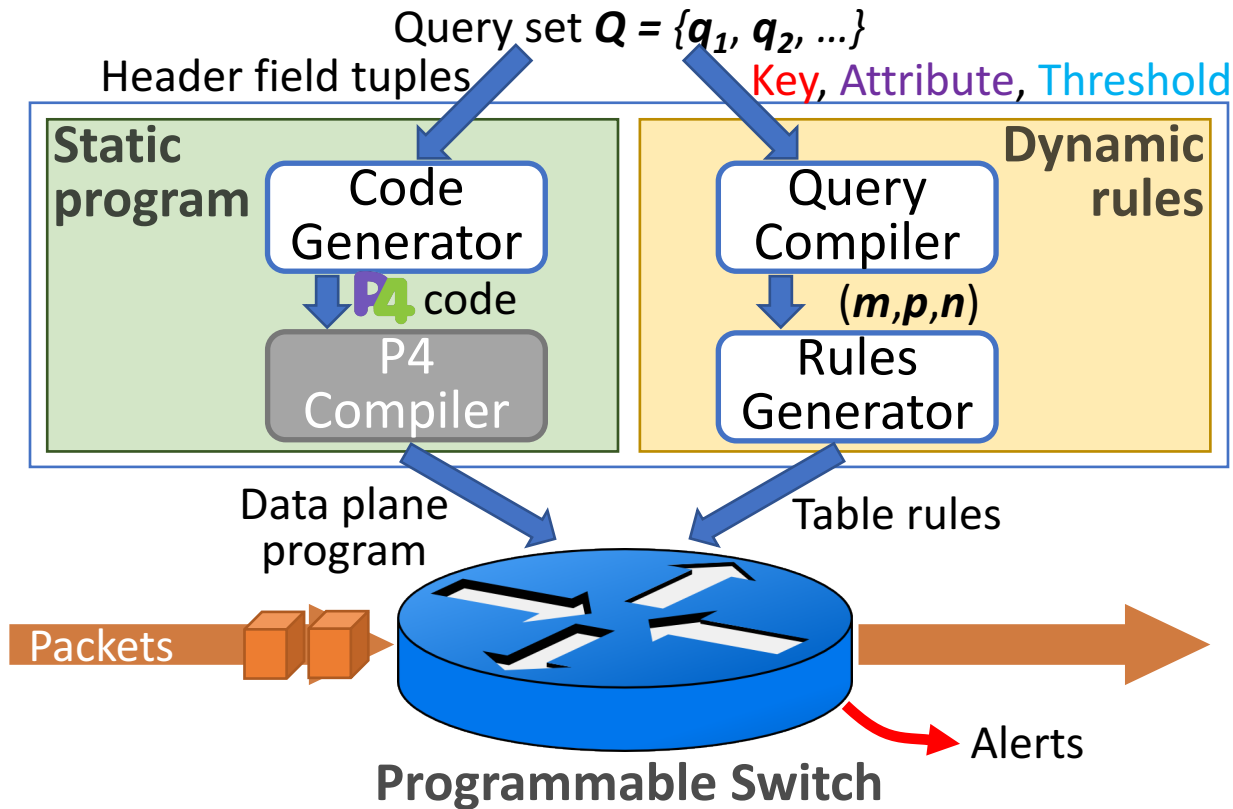
Goal:
I.   Stop near Threshold
II.  Update limit $m*p \leq \gamma_q$
III. HW limit, e.g., $m \leq 32$

# Query compiler

Query set
$Q = \{q_1, q_2, ...\}$

Query Compiler

Total coupons: $m$
Each probability: $p$
Coupons to collect: $n$

P4 Program

Switch Data Plane

Threshold=1000, $\gamma_q$ =0.01
($m$=20, $p$=1/2048, $n$=8)

Threshold=1000

0        2000

# Installing queries into switches

Query set $Q = \{q_1, q_2, ...\}$

Header field tuples

Key, Attribute, Threshold

**Static program**

Code Generator

P4 code

P4 Compiler

**Dynamic rules**

Query Compiler

$(m,p,n)$

Rules Generator

Data plane program
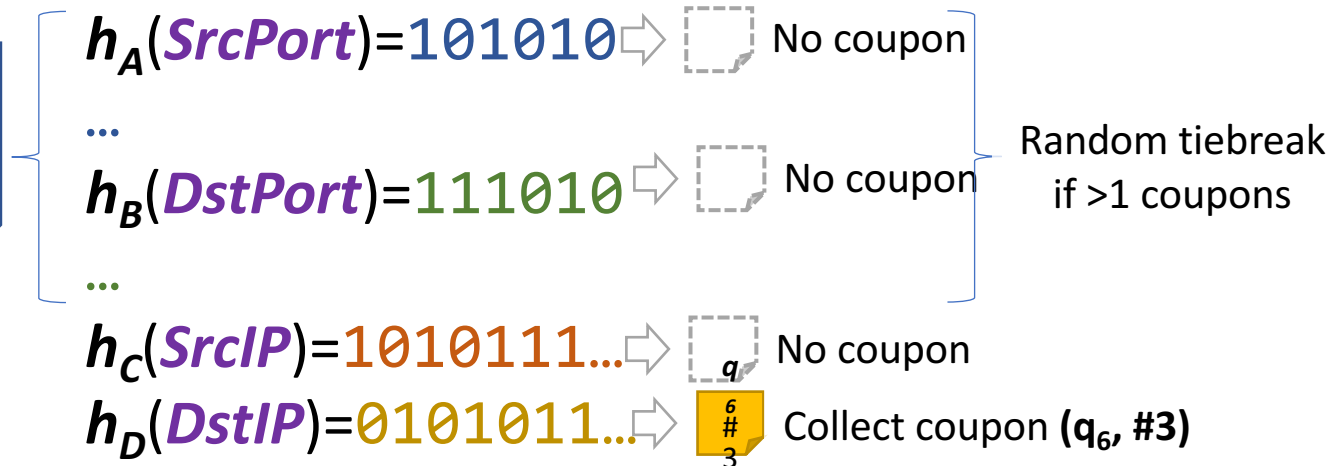
Table rules

Packets

**Programmable Switch**

Alerts

- The installed rules represent query set $Q$
- Update queries *on the fly*, without recompiling P4

# TCAM for selecting a coupon

| Match h_A(SrcPort) | | Query#,Coupon# | |
|---|---|---|---|

| Match h_B(DstPort) | | Query#,Coupon# | |
|---|---|---|---|

| Match h_C(SrcIP) | | Query#,Coupon# | |
|---|---|---|---|

| Match h_D(DstIP) | Query#,Coupon# |
|---|---|
| 000***** | (6,1) |
| 001***** | (6,2) |
| **010******* | **(6,3)** |
| 01101*** | (8,1) |
| … | … |

**Packet**
SrcPort: 25012
DstPort: 443
SrcIP: 10.0.1.15
DstIP: 162.249.4.107

$h_A(SrcPort) = 101010$ ⇒ No coupon

…

$h_B(DstPort) = 111010$ ⇒ No coupon

Random tiebreak
if >1 coupons

…

$h_C(SrcIP) = 1010111…$ ⇒ $q$ No coupon

$h_D(DstIP) = 0101011…$ ⇒ $\frac{6}{\#3}$ Collect coupon $(q_6, \#3)$

# Coupon collector table in SRAM



**Packet**
SrcPort: 27000
DstPort: 443
SrcIP: 10.0.1.33
DstIP: 4.3.2.1

$q_6$ #1

$q_6$ Coupon #1
Key:
10.0.1.33

**Packet**
SrcPort: 25012
DstPort: 443
SrcIP: 10.0.1.15
DstIP:
162.249.4.107

$q_6$ #3

$q_6$ Coupon #3
Key:
10.0.1.15

$q_6$:
*SrcIP*

| Q , Key | Coupons |
|---|---|
| $q_4$: 8.8.8.8:53 | 1  2  3  4 |
| $q_4$: 1.1.1.1:53 | 1  2  3  4 |
| $q_5$: 10.0.0.1 | 1  2  3 |
| $q_6$: 10.0.1.33 | 1  2  3 |
|  | 3 |

Query $q_6$
Key
10. .1.33

Space efficiency:
- Keys from all queries multiplexed into one table
- Only keep rows for "active keys" (at least one coupon)
- Clear rows after timeout

*SrcIP* 10.0.1.33 is sending
to >1000 distinct *DstIP*s.

# Evaluation highlights

- How efficient is BeauCoup?

Uses **4x~10x fewer memory access** than the state-of-the-art to achieve the same accuracy.

- How much hardware resource?

On the Barefoot Tofino programmable switch, BeauCoup occupies **<50% of each resource**

# *BeauCoup:*

Answering *many* network traffic queries,
*one* memory update at a time!

- **Scalable**: built upon *coupon collectors,* runs many queries simultaneously

- **Versatile**: change queries on the fly, without recompiling P4 program

- **Efficient**: achieve the same accuracy using 4x-10x fewer memory accesses

Is this a good usecase of programmable dataplanes?

# What are the limitations?

# Elmo: Source Routed Multicast for Public Clouds

## Muhammad Shahbaz

Lalith Suresh, Jennifer Rexford, Nick Feamster,

Ori Rottenstreich, and Mukesh Hira

Stanford University

PRINCETON UNIVERSITY

Technion — Israel Institute of Technology

vmware

*Slides from Muhammad Shahbaz*

# 1-to-Many Communication in Cloud

# 1-to-Many Communication in Cloud

# 1-to-Many Communication in Cloud

10,000s of tenants

�’ 100s of workloads

➘ Millions of groups



amazon  Google  ◼ Microsoft

# 1-to-Many Communication in Cloud

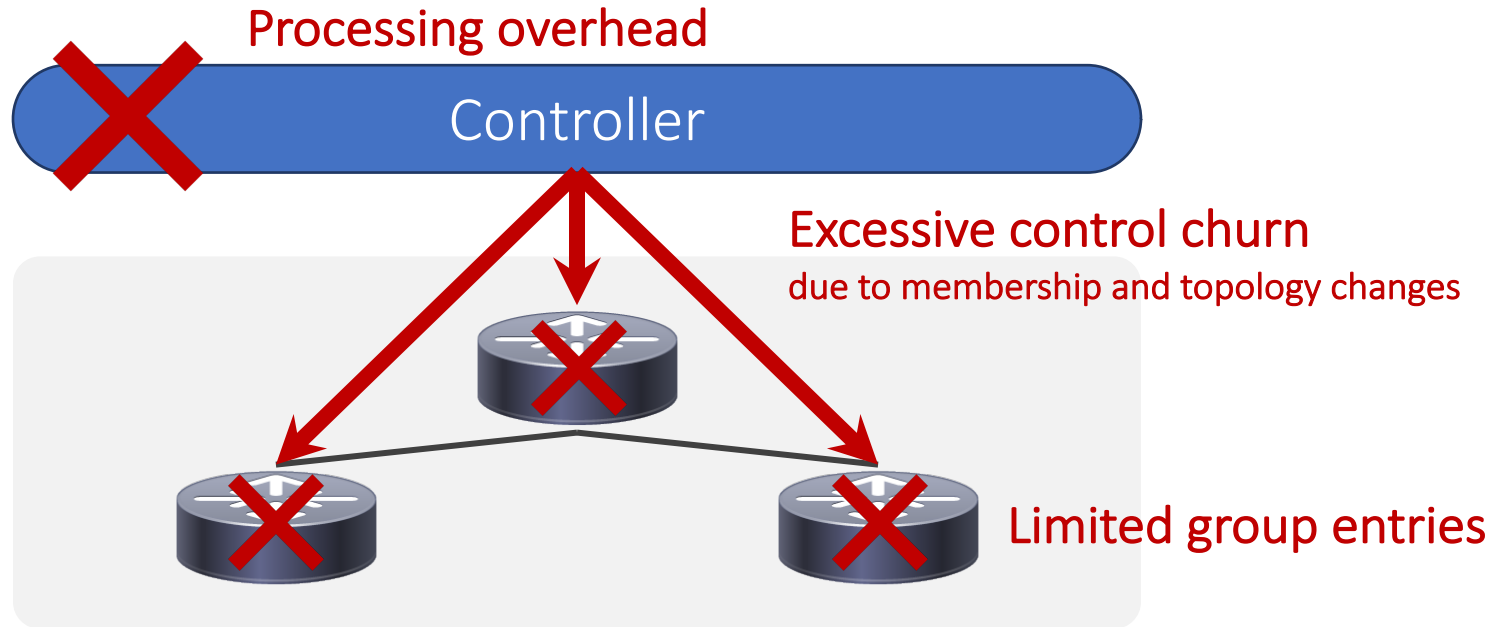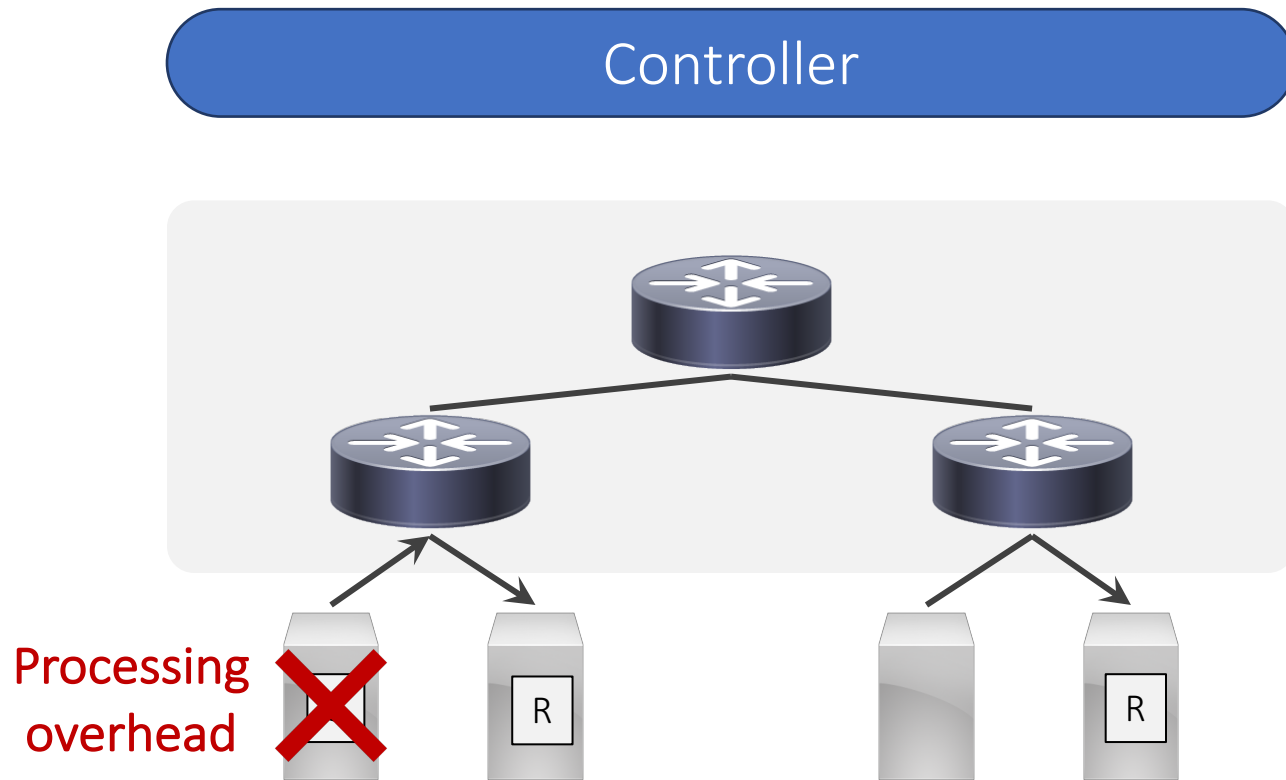10,000s of tenants

➘ 100s of workloads

➘ Millions of groups



Multicast

# 1-to-Many Communication in Cloud

10,000s of tenants

�especially 100s of workloads

➥ Millions of groups

# Limitations of **Native** Multicast

Controller

# Limitations of <u>Native</u> Multicast

Processing overhead

Controller

Excessive control churn
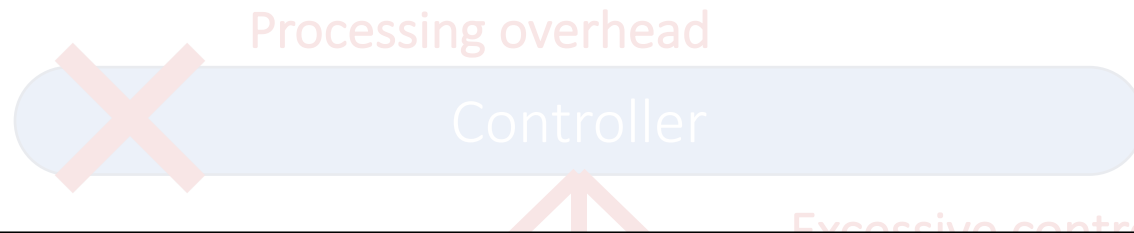due to membership and topology changes

Limited group entries

# Restricted to <u>Unicast-based</u> Alternatives

# Restricted to <u>Unicast-based</u> Alternatives

# 1-to-Many Communication in the Cloud

Processing overhead

Controller

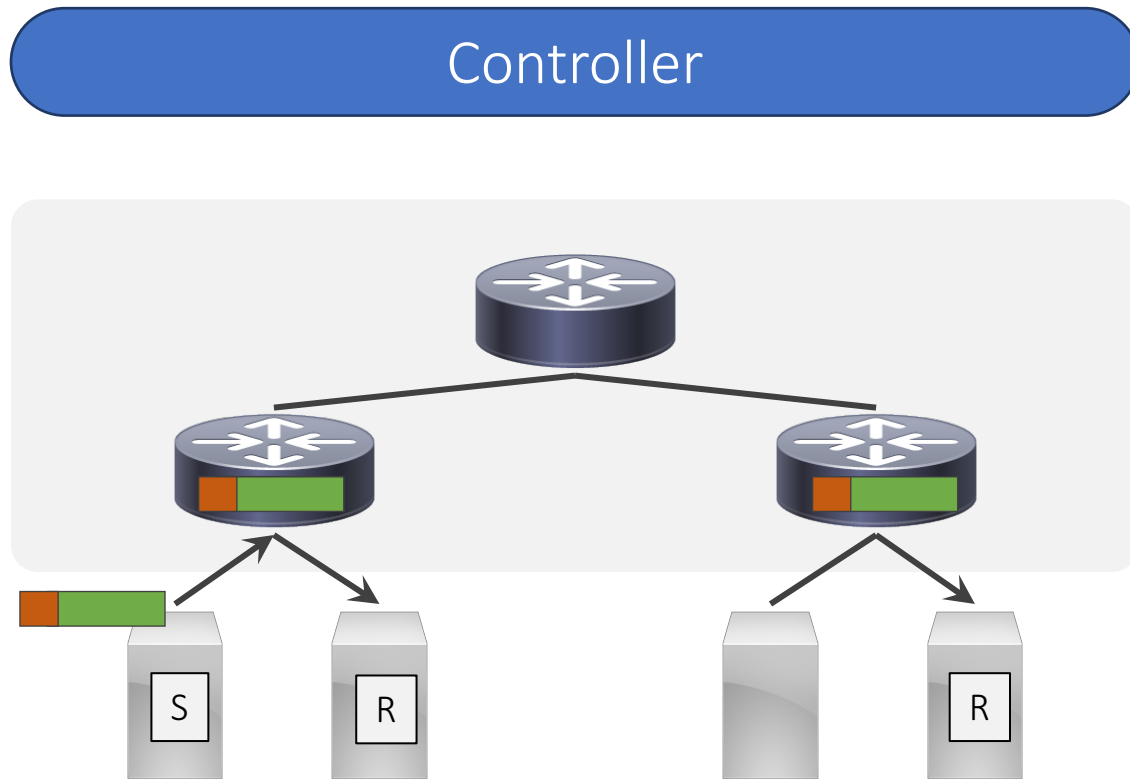Excessive control churn

Need a scheme that **scales** to millions of groups
**without**
excessive control, end-host CPU, and traffic **overheads!**

Processing
overhead

overhead

# Proposal: Source Routed Multicast

# Proposal: Source Routed Multicast

Controller

# Proposal: Source Routed Multicast

Little processing overhead

Controller

Minimal control churn

No traffic overhead

No group entries needed*

Negligible processing overhead

S

R

R

# A Naïve Source Routed Multicast

A **multicast group** encoded as a list of **(Switch, Ports)** pairs

```
Switch 1: [Ports ]
Switch 2: [.. .. ..]
Switch 3: [.. .. ..]
Switch 4: [.. .. .x ..]
Switch 5: [.x .. .. ..]
```

For a data center with:
- **1000** switches
- **48 ports** per switch

O(30) bytes per switch

O(30,000) bytes header

for a group spanning 1000 switches

**20x** the packet size!

# Enabling Source Routed Multicast in Public Clouds

Key attributes:

- <u>Efficiently encode</u> multicast forwarding policy inside packets

- <u>Process</u> this encoding at <u>hardware speed</u> in the switches

- <u>Execute</u> tenants' applications <u>without modification</u>

# Encoding a Multicast Policy in <u>Elmo</u>

A multicast group encoded as
a list of **(Switch, Ports)** pairs

```
Switch 1: [Ports ]
Switch 2: [.. .. ..]
Switch 3: [.. .. ..]
Switch 4: [.. .. .x ..]
Switch 5: [.x .. .. ..]
```

# Encoding a Multicast Policy in <u>Elmo</u>

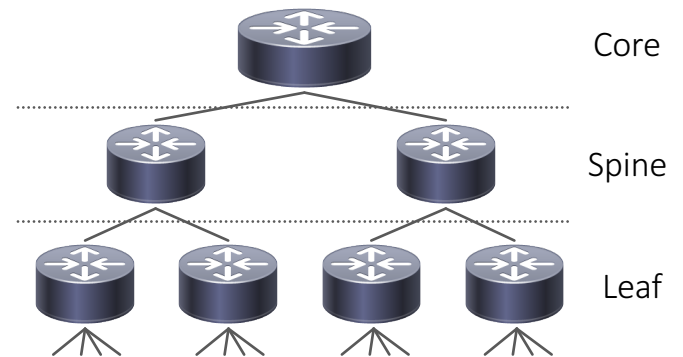A multicast group encoded as a list of **(Switch, Ports)** pairs

```
Switch 1: [Bitmap] ◄
Switch 2: [.. .. ..]
Switch 3: [.. .. ..]
Switch 4: [.. .. .x ..]
Switch 5: [.x .. .. ..]
```

**①** Encode switch ports as a **bitmap**

**Bitmap** is the internal **data structure** that switches use for **replicating packets**

# Encoding a Multicast Policy in Elmo

A multicast group encoded as
a list of **(Switch, Ports)** pairs

```
Switch 1: [Bitmap]
Switch 2: [.. .. ..]
Switch 3: [.. .. ..]
Switch 4: [.. .. .x ..]
Switch 5: [.x .. .. ..]
```

**2** Group switches into **layers**

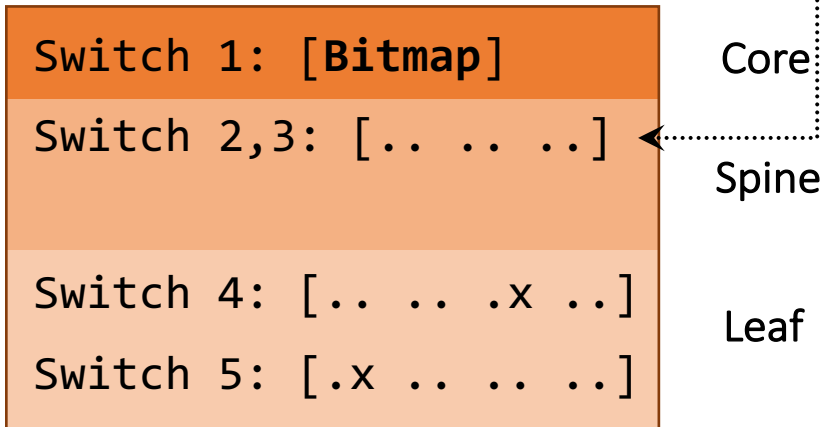# **Encoding a Multicast Policy in <u>Elmo</u>**

A multicast group encoded as
a list of **(Switch, Ports)** pairs

```
Switch 1: [Bitmap]
Switch 2: [.. .. ..]
Switch 3: [.. .. ..]
Switch 4: [.. .. .x ..]
Switch 5: [.x .. .. ..]
```

Core
Spine
Leaf

② Group switches into **layers**



Core

Spine

Leaf

More precisely: *upstream leaf, upstream spine, core, downstream spine, downstream leaf*

# Encoding a Multicast Policy in **Elmo**

A multicast group encoded as
a list of **(Switch, Ports)** pairs

```
Switch 1: [Bitmap]
Switch 2: [.. .. ..]
Switch 3: [.. .. ..]
Switch 4: [.. .. .x ..]
Switch 5: [.x .. .. ..]
```
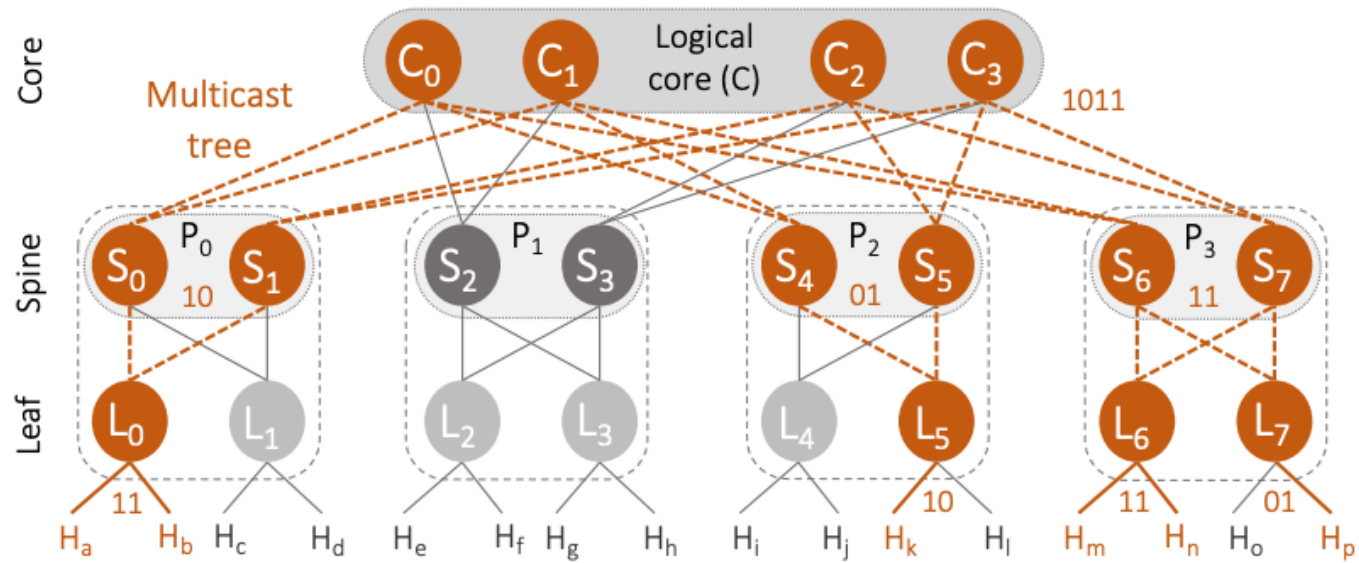
**3** Switches within a layer with **same** ports **share a bitmap**

# Encoding a Multicast Policy in Elmo

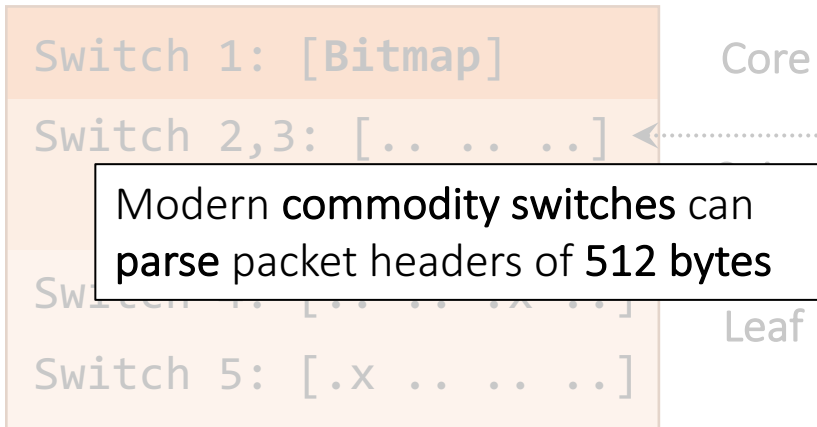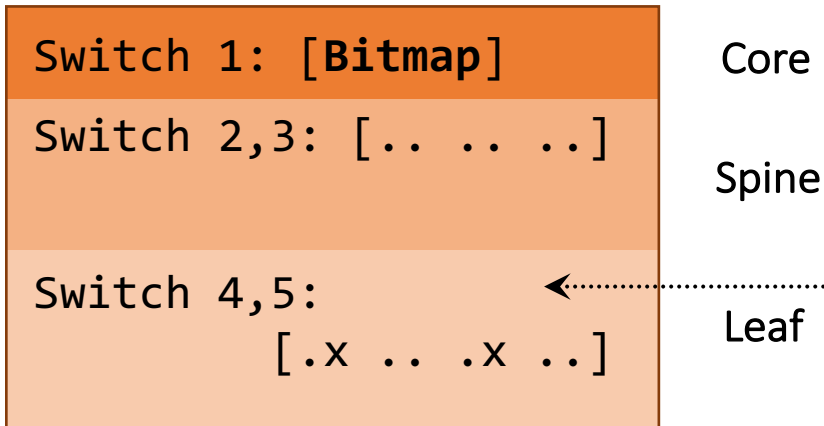A multicast group encoded as a list of (Switch, Ports) pairs

```
Switch 1: [Bitmap]
Switch 2,3: [.. .. ..]

Switch 4: [.. .. .x ..]
Switch 5: [.x .. .. ..]
```

Core

Spine

Leaf

③ Switches within a layer with **same** ports **share a bitmap**

47

# Encoding a Multicast Policy in Elmo

# Encoding a Multicast Policy in **Elmo**

A multicast group encoded as
a list of **(Switch, Ports)** pairs

Switch 1: [**Bitmap**]
Switch 2,3: [.. .. ..]

Core

③ Switches within a layer with **same** ports **share a bitmap**

Modern **commodity switches** can **parse** packet headers of **512 bytes**

Switch 5: [.x .. .. ..]

Leaf

Li et al.

For a data center with:
- 628 switches
- 325 bytes header space

Supports **890,000** groups!

No. of groups

1M
750K
500K
250K
0

# Encoding a Multicast Policy in <u>Elmo</u>

A multicast group encoded as
a list of **(Switch, Ports)** pairs

| |
|---|
| Switch 1: [**Bitmap**] |
| Switch 2,3: [.. .. ..] |
| |
| Switch 4: [.. .. .x ..] |
| Switch 5: [.x .. .. ..] |

Core

Spine

Leaf

④ Switches within a layer with N **different** ports **share a bitmap**

# Encoding a Multicast Policy in <u>Elmo</u>

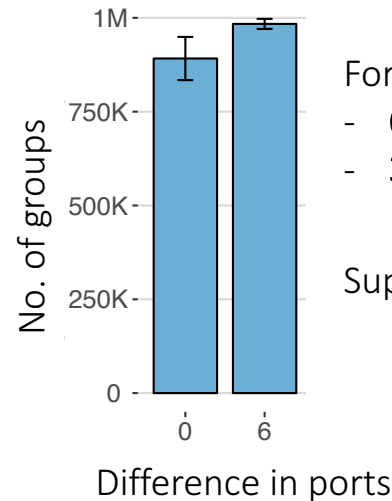A multicast group encoded as a list of **(Switch, Ports)** pairs

```
Switch 1: [Bitmap]
Switch 2,3: [.. .. ..]

Switch 4,5:
         [.x .. .x ..]
```

Core

Spine

Leaf

④ Switches within a layer with **N different** ports **share a bitmap**

51

# Encoding a Multicast Policy in <u>Elmo</u>

A multicast group encoded as
a list of **(Switch, Ports)** pairs

```
Switch 1: [Bitmap]
Switch 2,3: [.. .. ..]

Switch 4,5:
          [.x .. .x ..]
```

Core

Spine

Leaf

**4** Switches within a layer with **N different** ports **share a bitmap**



No. of groups

1M
750K
500K
250K
0

0    6

Difference in ports

Li et al.

For a data center with:
- **628** switches
- **325** bytes header space

Supports **980,000** groups!

# Encoding a Multicast Policy in **Elmo**

A multicast group encoded as
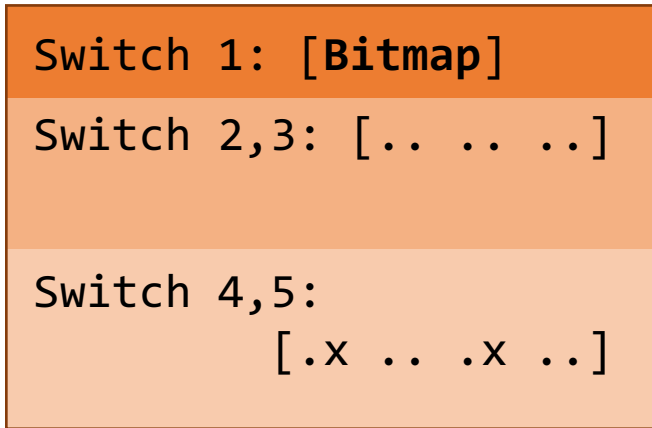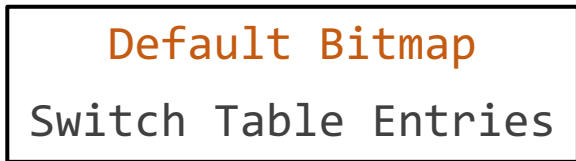a list of **(Switch, Ports)** pairs

Fixed Header Size

```
Switch 1: [Bitmap]        Core

Switch 2,3: [.. .. ..]    Spine

Switch 4,5:
        [.x .. .x ..]     Leaf
```

# **Encoding a Multicast Policy in <u>Elmo</u>**

A multicast group encoded as
a list of **(Switch, Ports)** pairs

Fixed Header Size

| |
|---|
| Switch 1: [**Bitmap**] |
| Switch 2,3: [.. .. ..] |
| Switch 4,5:<br>　　　　[.x .. .x ..] |

Core

Spine

Leaf

Default Bitmap

Switch Table Entries

# Encoding a Multicast Policy in <u>Elmo</u>

A multicast group enco
a list of **(Switch, Ports)**

**Fixed Header Size**

Switch [Bitmap]

Switch

Switch

Switch Table Entries

⑤ Use **switch entries** and a **default bitmap** for **larger groups**

Li et al.

**Switch entries**

5.0K

2.5K

0.0

0   6

**Traffic overhead**

1

0

0   6

Difference in ports

For a data center with:
- **628** switches
- **325** bytes header space

**Groups covered with p-rules**

750K

500K

250K

0

0   6   12

Redundancy limit (R

55

# Encoding a Multicast Policy in **Elmo**

A multicast group encoded as a list of **(Switch, Ports)** pairs

Fixed Header Size

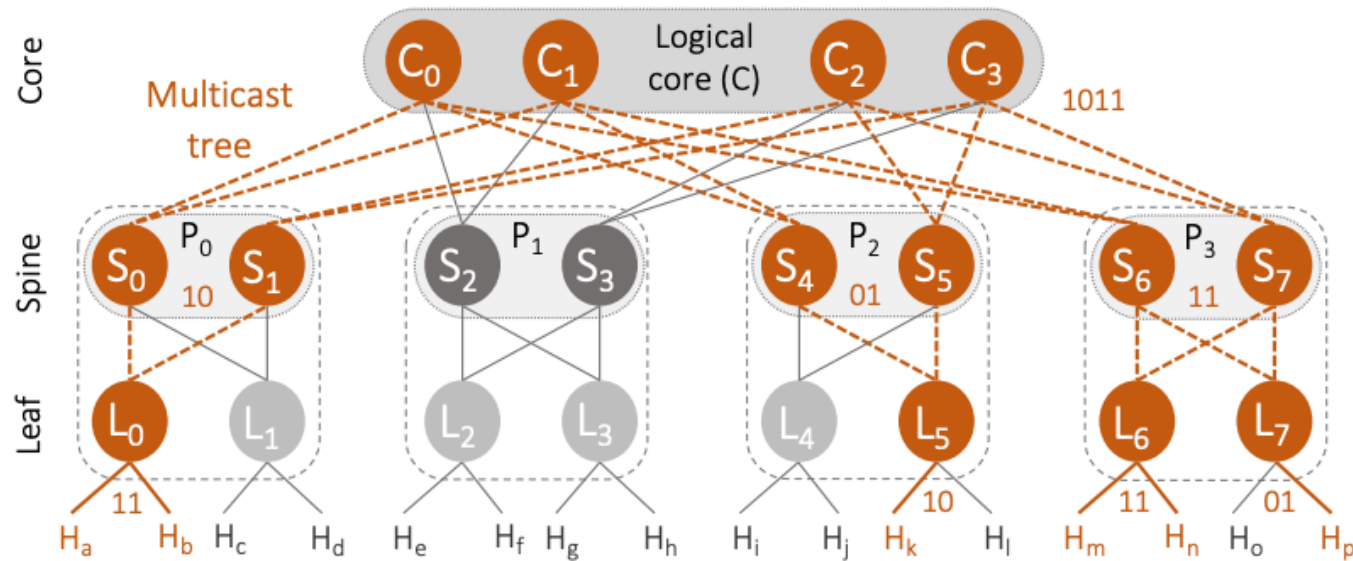| | |
|---|---|
| Switch 1: [**Bitmap**] | Core |
| Switch 2,3: [.. .. ..] | Spine |
| Switch 4,5:<br>        [.x .. .x ..] | Leaf |

**Default Bitmap**
Switch Table Entries

① Encode switch ports as a **bitmap**

② Group switches into **layers**

③ Switches within a layer with:
- **same** ports **share a bitmap**

④ - **N different** ports **share a bitmap**

⑤ Use **switch entries** and a **default bitmap** for **larger groups**

For a data center with:
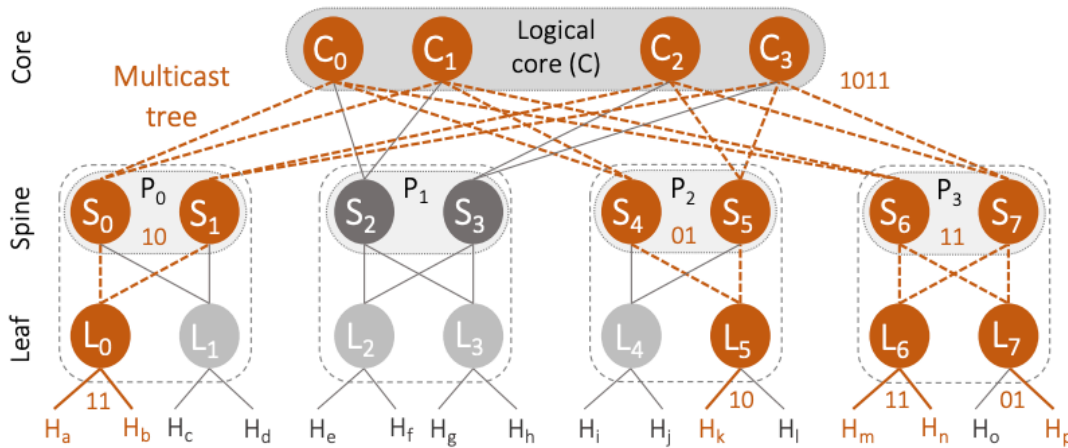- **628** switches
- **325** bytes header space

Supports **a Million** groups!

56

# Encoding a Multicast Policy in Elmo

# Encoding a Multicast Policy in Elmo



Downstream spine and leaf assignments with varying degree of redundancy (R) and s-rules.

| R = 0 | | R = 2 |
|---|---|---|
| #s-rules = 0 | #s-rules = 1 | #s-rules = 0, 1 |
| $(p)$ 10:$[P_0]$<br>$(p)$ 01:$[P_2]$<br>$(d)$ 11:$[P_3]$ | $(p)$ 10:$[P_0]$<br>$(p)$ 01:$[P_2]$<br>$(s)$ 11:$[P_3]$ | $(p)$ 10:$[P_0]$<br>$(p)$ 11:$[P_2,P_3]$ |
| $(p)$ 11:$[L_0,L_6]$<br>$(p)$ 10:$[L_5]$<br>$(d)$ 01:$[L_7]$ | $(p)$ 11:$[L_0,L_6]$<br>$(p)$ 10:$[L_5]$<br>$(s)$ 01:$[L_7]$ | $(p)$ 11:$[L_0,L_6]$<br>$(p)$ 11:$[L_5,L_7]$ |

#p-rules = 2 (and max two switches per p-rule)

legends: p-rule (p), s-rule (s), and default p-rule (d)

| Sender $H_a$ | type | u-leaf | u-spine | d-core | Common downstream spine and leaf p-rules | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | d-spine | | d-leaf | | |
| Outer header(s) VXLAN | u | 01\|M | 00\|M | 0011 | 10:$[P_0]$  11:$[P_3]$ | 11:$[L_0,L_6]$  01:$[L_7]$ | | | Packet body |
| | | At $L_0$: forward to $H_b$ and multipath to $P_0$ | $P_0$: multipath to C | C: forward to $P_2$, $P_3$ | 01:$[P_2]$  Default | 10:$[L_5]$  Default | | | |
| Sender $H_k$ | | | | | | | | | |
| Outer header(s) VXLAN | u | 00\|M | 00\|M | 1001 | | | | | Packet body |
| | | At $L_5$: multipath to $P_2$ | $P_2$: multipath to C | C: forward to $P_0$, $P_3$ | $P_0$: forward to $L_0$<br>$P_2$: forward to $L_5$<br>$P_3$: forward to $L_6$, $L_7$ | $L_0$: forward to $H_a$, $H_b$<br>$L_5$: forward to $H_k$<br>$L_6$: forward to $H_m$, $H_n$<br>$L_7$: forward to $H_p$ | | | |

# Processing a Multicast Policy in <u>Elmo</u>

2. Computes the multicast policy

1. API

Controller

3. Installs **entries** in programmable
- **virtual switches** to push Elmo headers on packets
- **hardware switches**

- More **flow entries** and higher **update rates** than hardware switches
- **No changes** to the tenant application

Virtual Switch
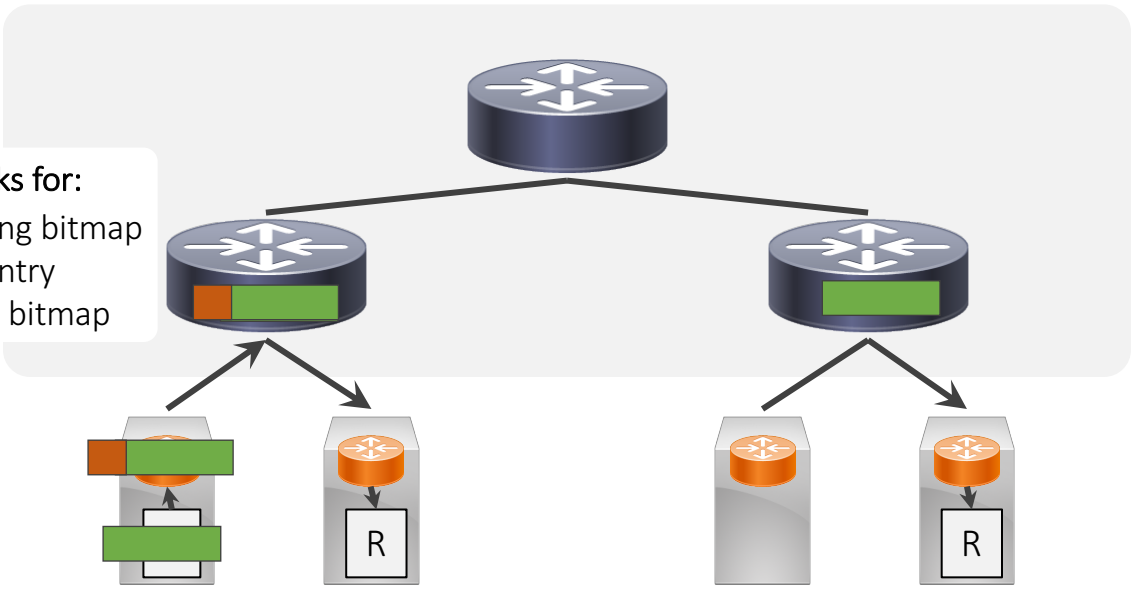
S

R

R

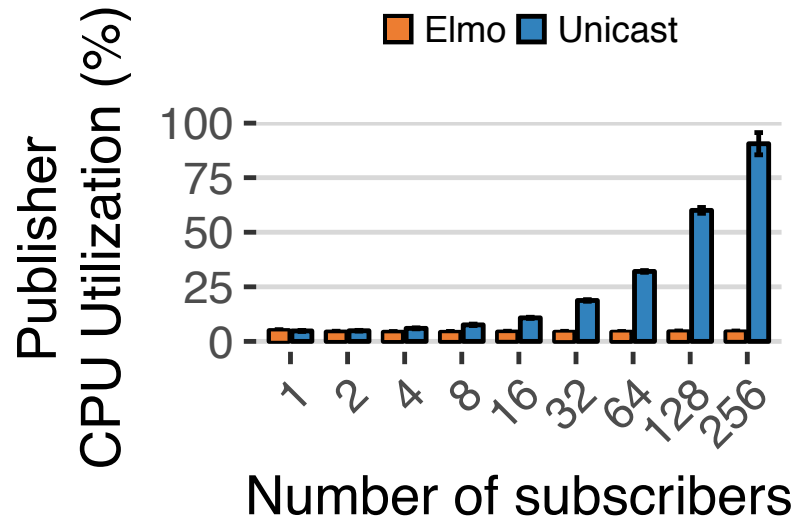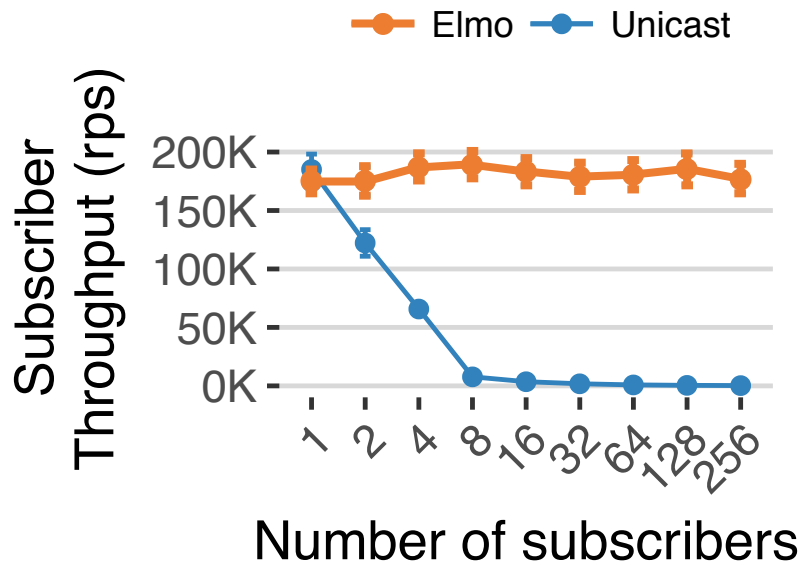# Processing a Multicast Policy in <u>Elmo</u>

Controller

Switch looks for:
    Matching bitmap
or  Table entry
or  Default bitmap

R

R

Implemented using P4 on a Barefoot Tofino Switch

# Applications Run Without Performance Overhead

# Conclusion

### Elmo

Source Routed
Multicast
for Public Clouds

- <u>Designed</u> for multi-tenant data centers

- <u>Compactly encodes multicast policy</u> inside packets

- <u>Operates at hardware speed</u> using programmable data planes

Is this a good usecase of programmable dataplanes?

What are the limitations?

# Other networking usecases

- Load balancing:
  - HULA: Scalable Load Balancing Using Programmable Data Planes, SOSR'16

- Congestion control:
  - Evaluating the Power of Flexible Packet Processing for Network Resource Allocation, NSDI'17
    - *Support RCP and XCP on programmable switches*
  - HPCC: High Precision Congestion Control, SIGCOMM'19
    - *Obtain precise link information for congestion control*

- A new protocols for more efficient L2 switching
  - The Deforestation of L2, SIGCOMM'16

- And others…