



# Realizing Source Routed Multicast Using Mellanox's Programmable Hardware Switches

Matty Kadosh<sup>1</sup>, Yonatan Piasezky<sup>1</sup>, Barak Gafni<sup>1</sup>, Lalith Suresh<sup>2</sup>, Muhammad Shahbaz<sup>3</sup>, Sujata Banerjee<sup>2</sup>

<sup>1</sup>Mellanox, <sup>2</sup>VMware, <sup>3</sup>Stanford

Sponsored By



# Group Communication in Public Clouds

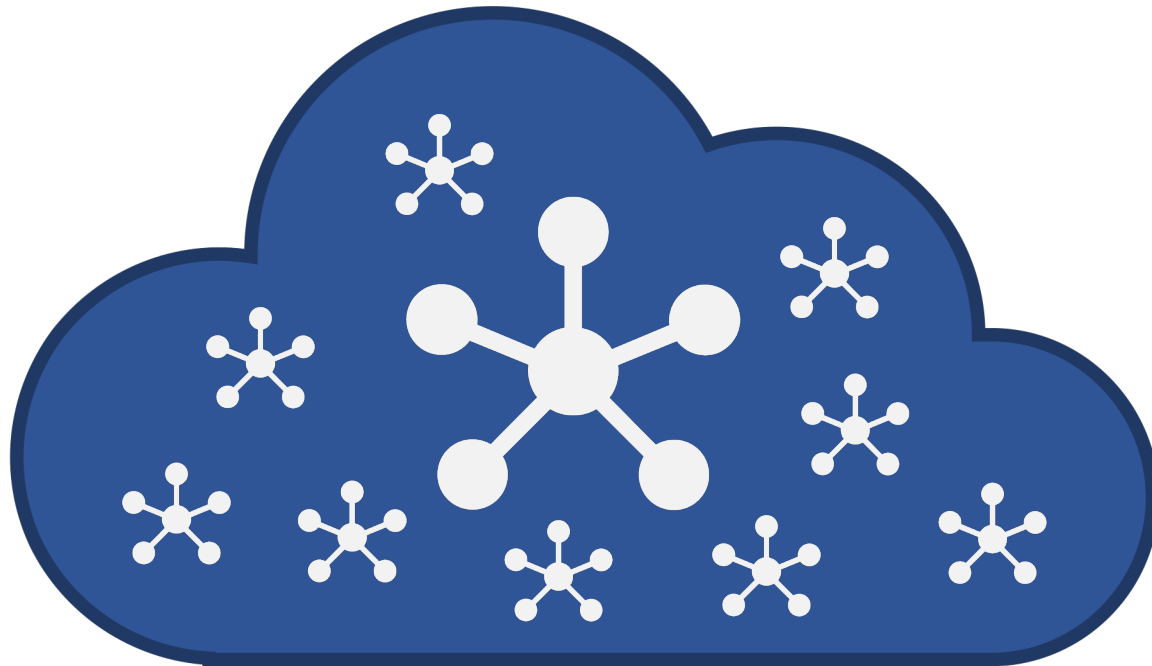
10,000s of tenants

→ 100s of workloads

Replication  
for databases and state machines

Publish-Subscribe  
like ZeroMQ and RabbitMQ

Infrastructure Apps  
like VMware NSX and OpenStack



Millions of distinct  
**group**  
communications

amazon Google Microsoft

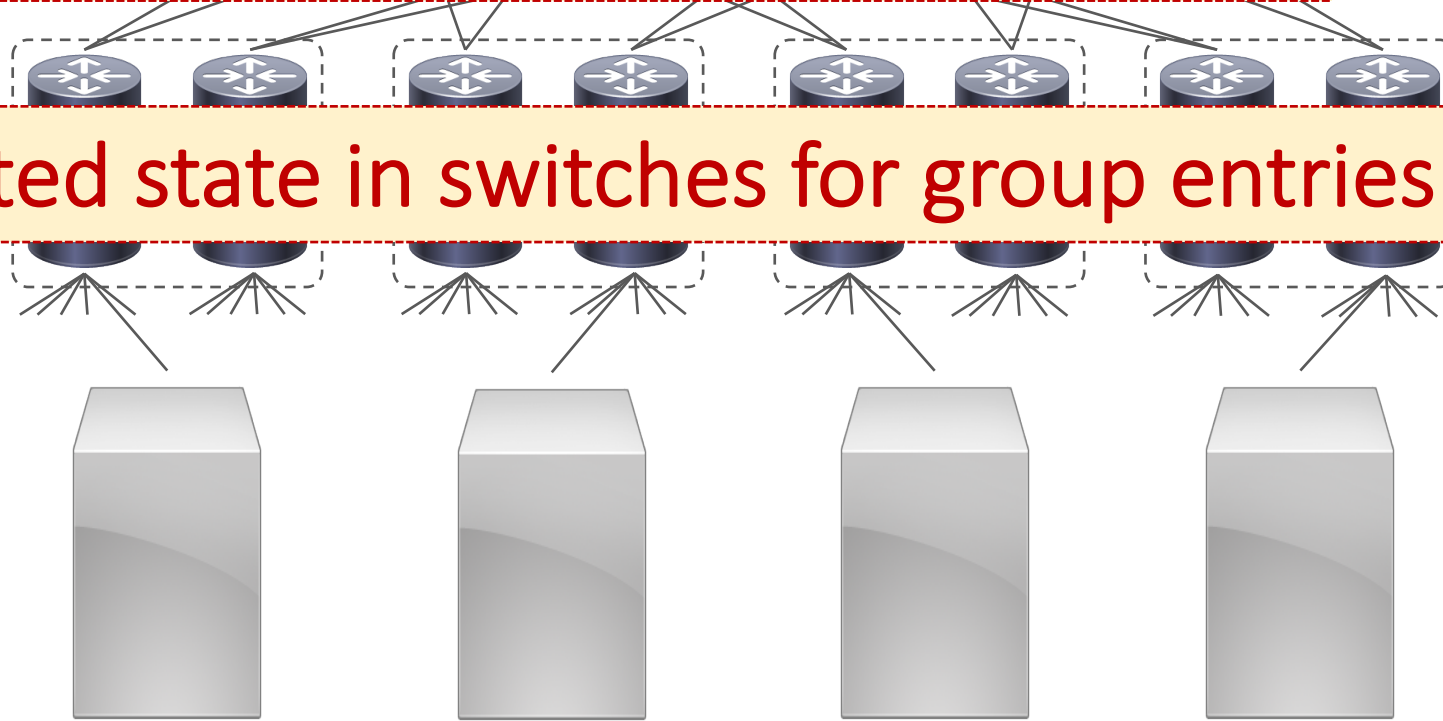
# Limitations of Native Multicast

Data Center Controller

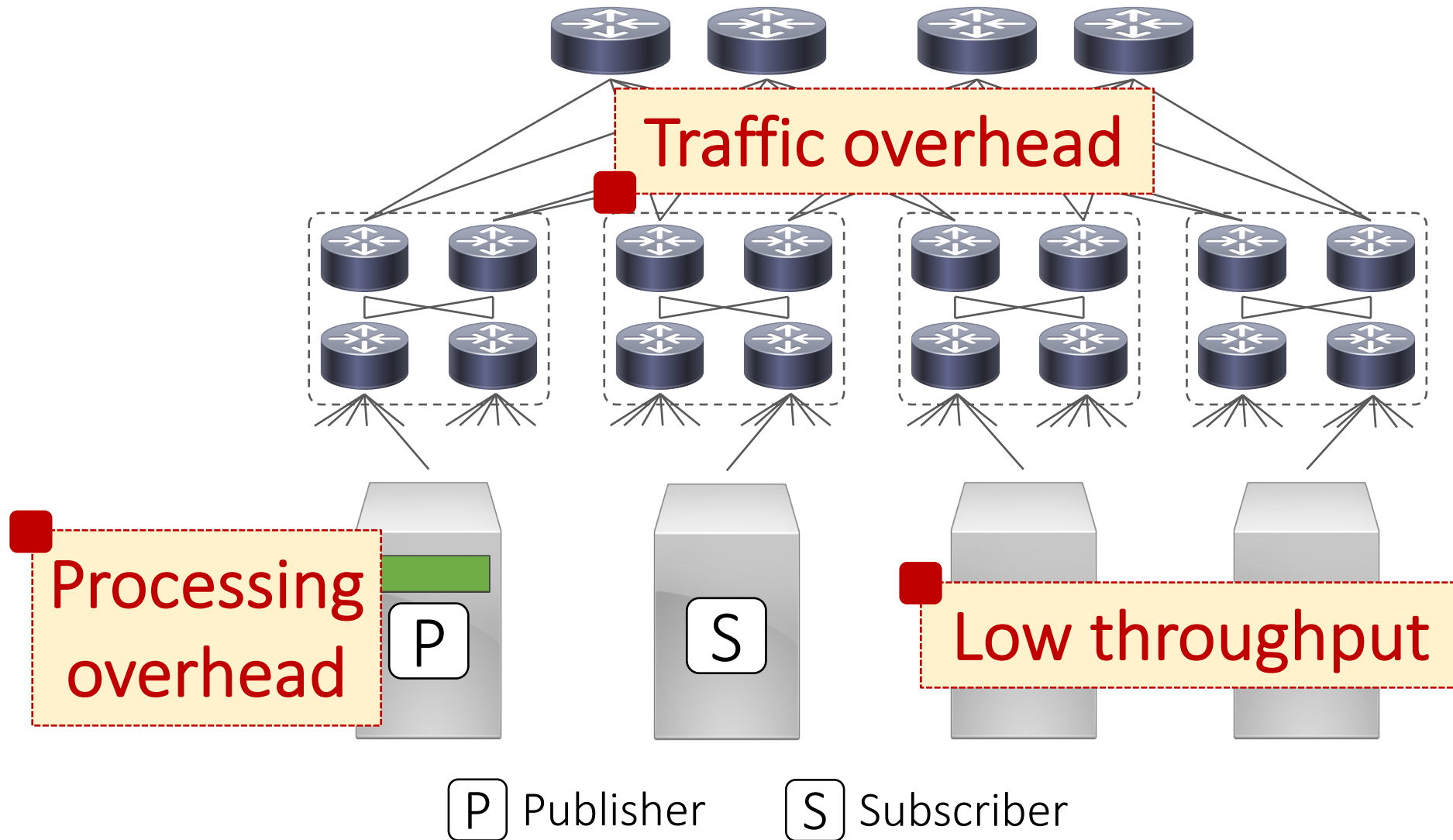
Processing overhead

Excessive control churn  
due to membership and topology changes

Limited state in switches for group entries < 10K

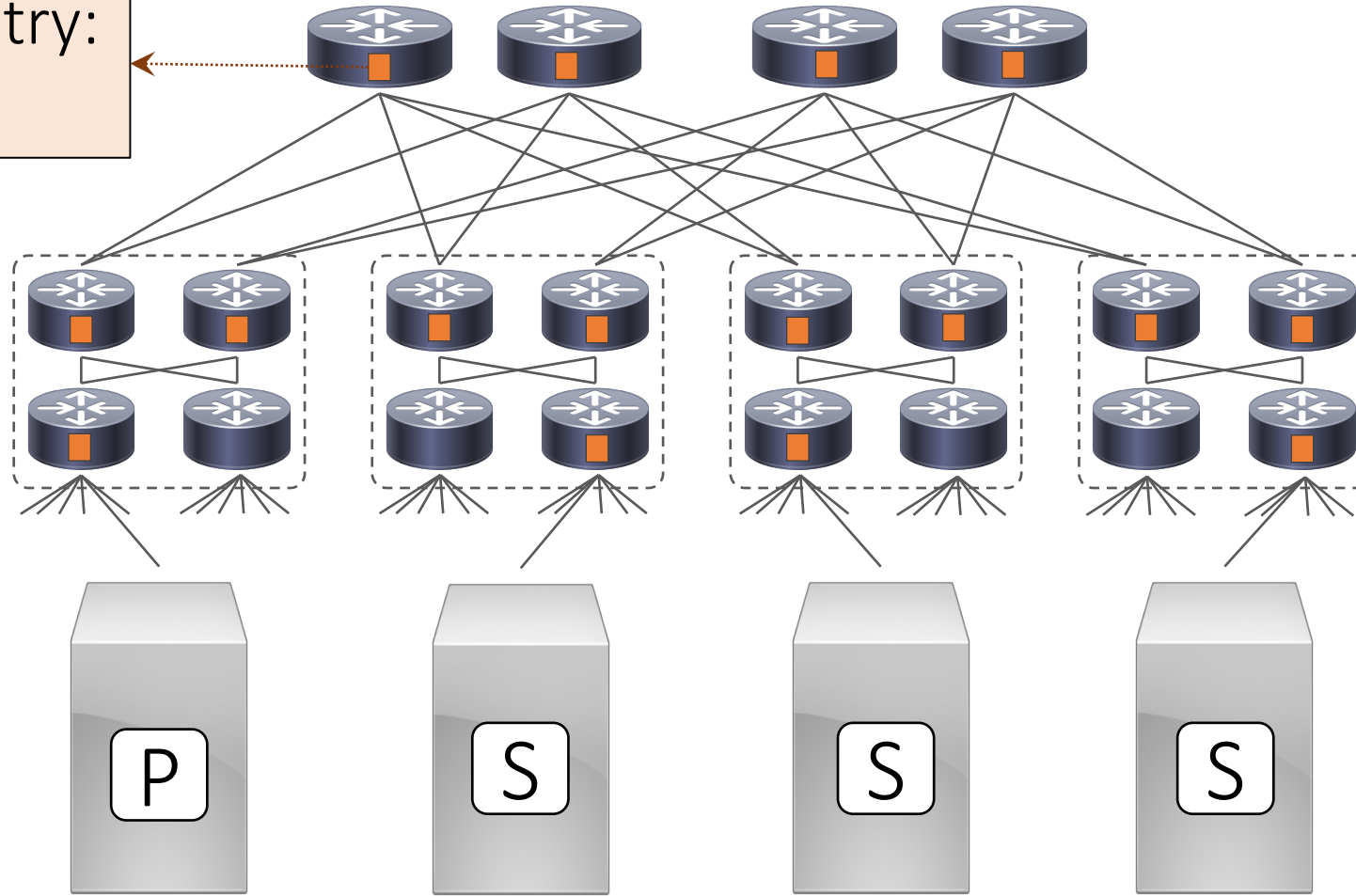


# Restricted to Application-Level Multicast



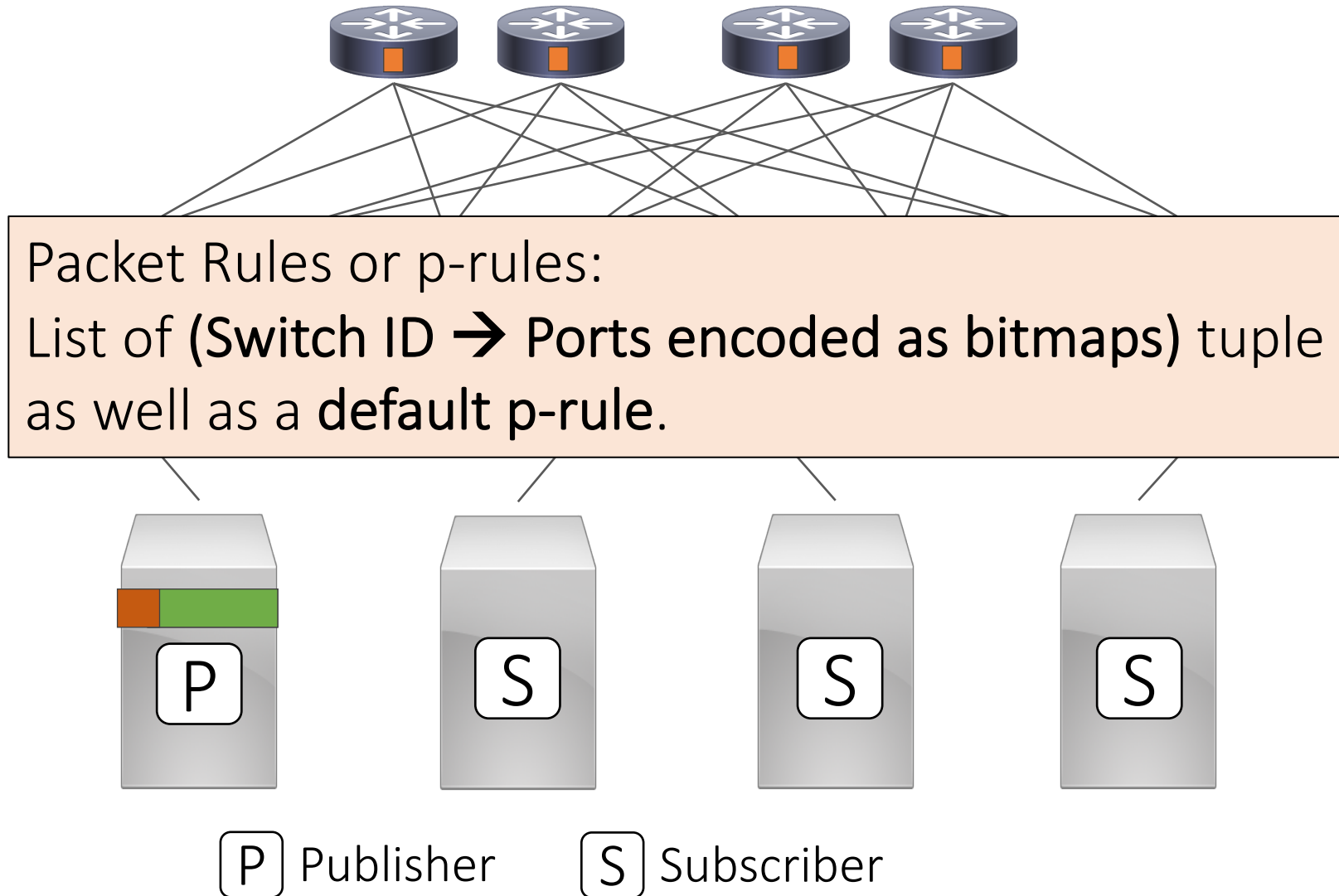
# Elmo: Source Routed Multicast → SIGCOMM'19

Group Table Entry:  
ID → Ports

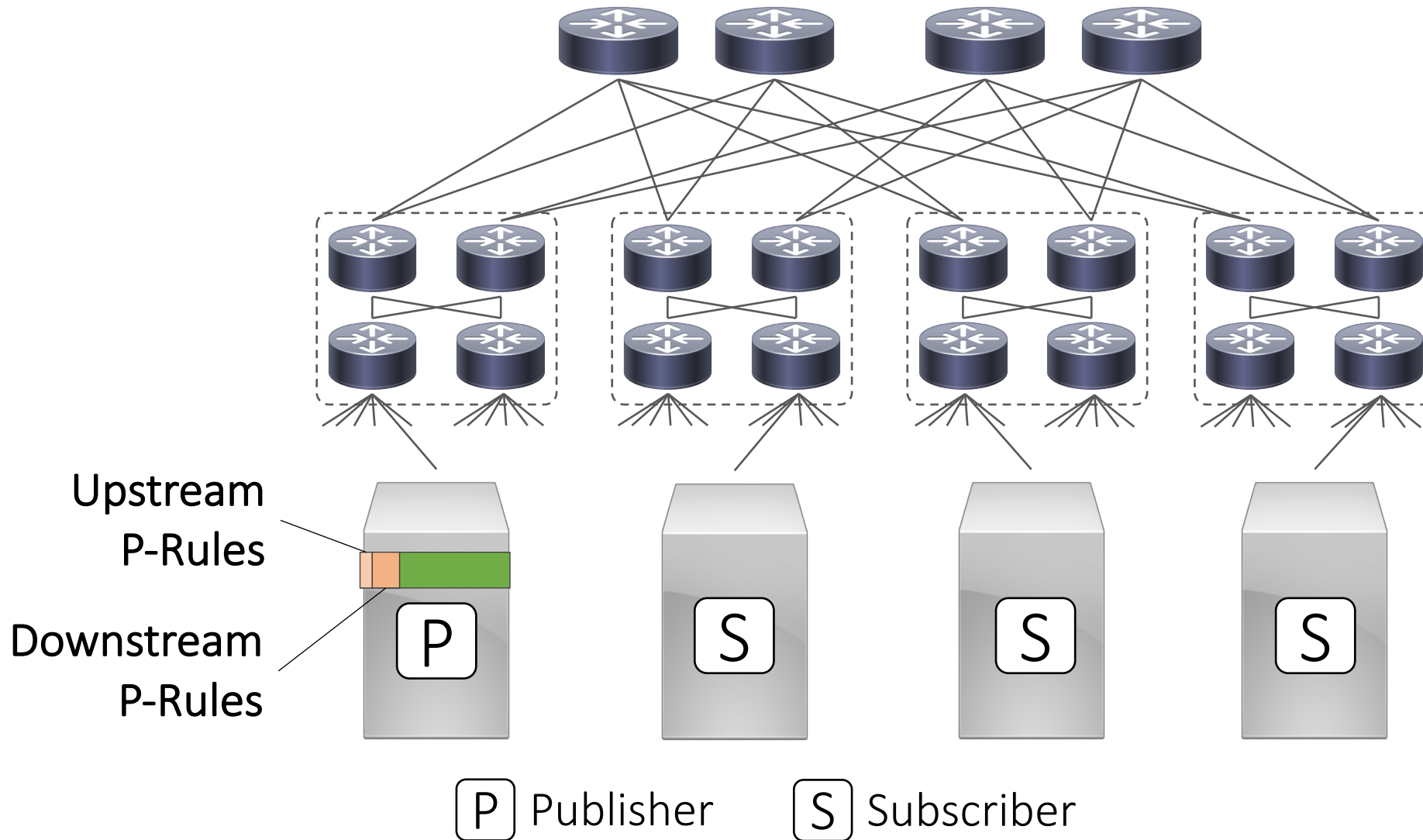


**P** Publisher    **S** Subscriber

# Elmo: Source Routed Multicast

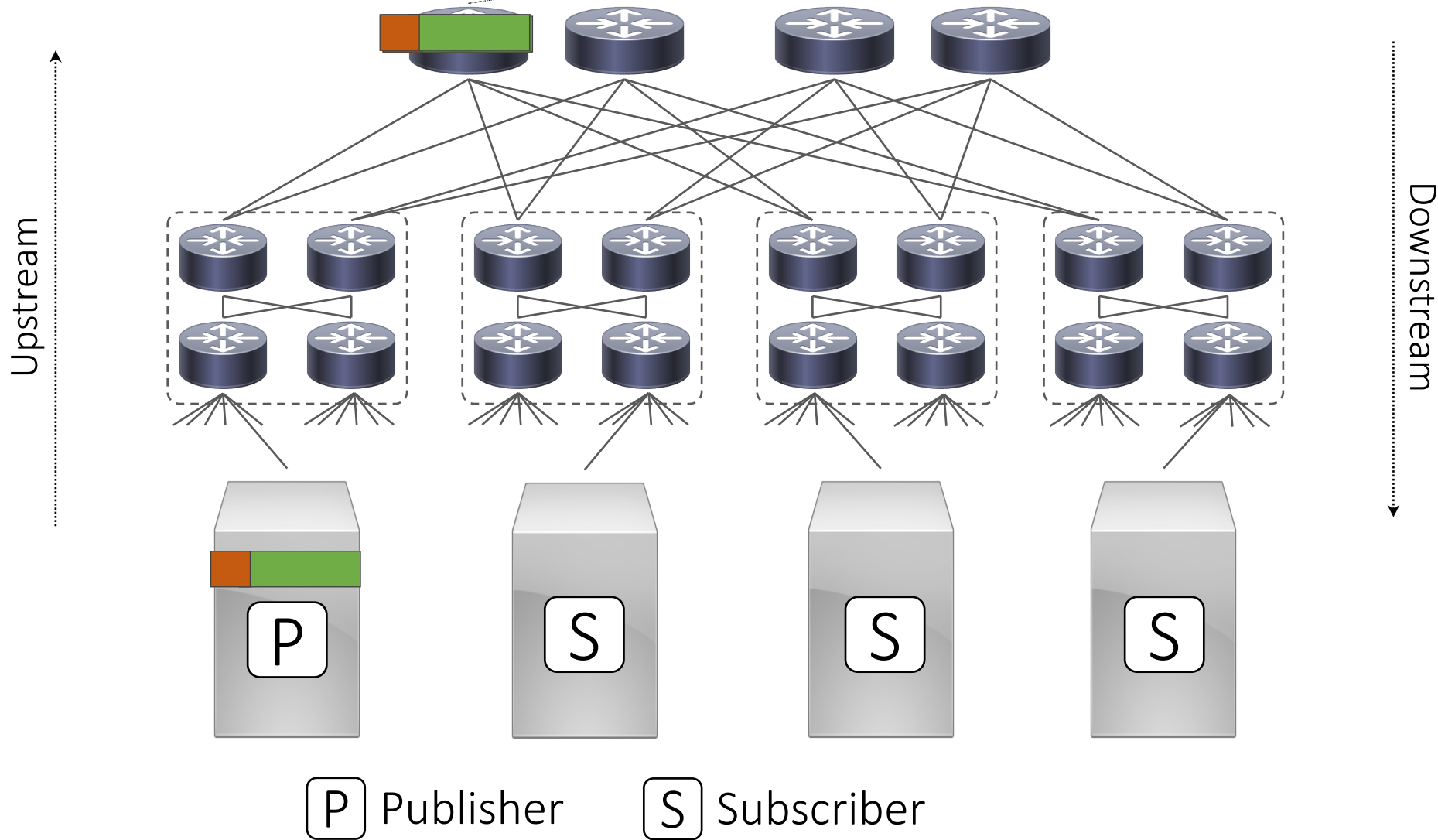


# Elmo: Source Routed Multicast



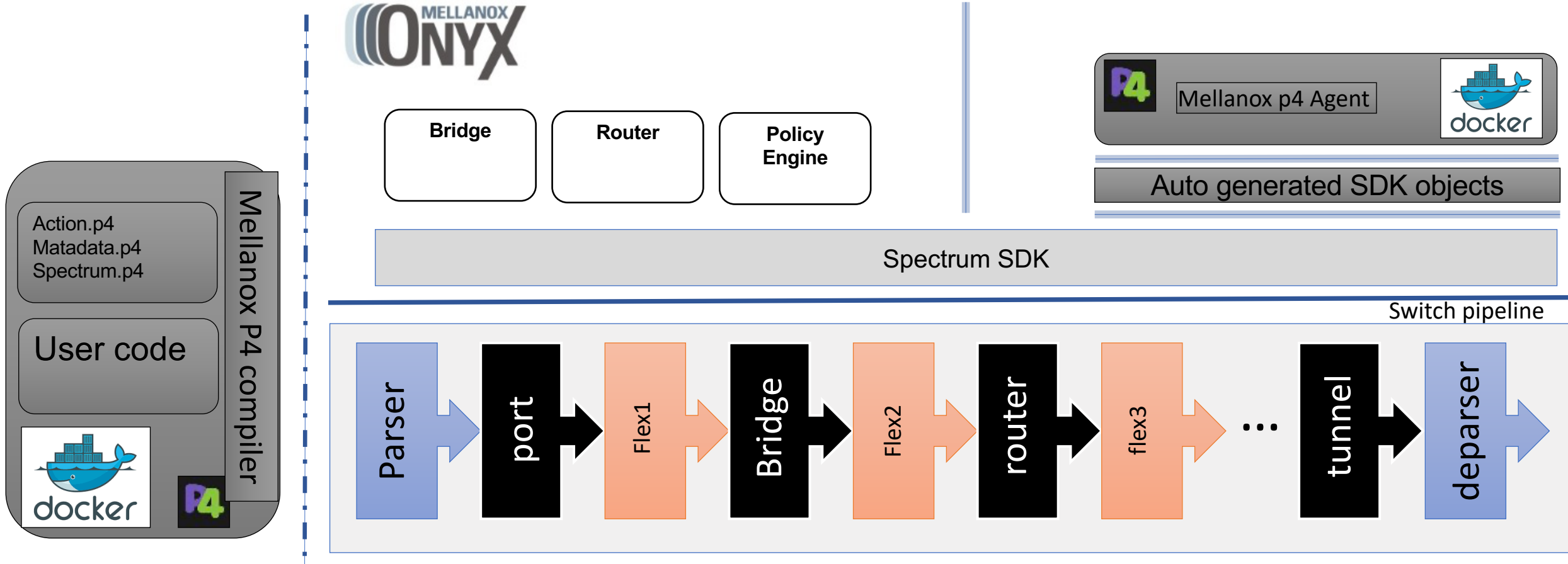
# Elmo: Source Routed Multicast

- Read a p-rule
- Read the default p-rule



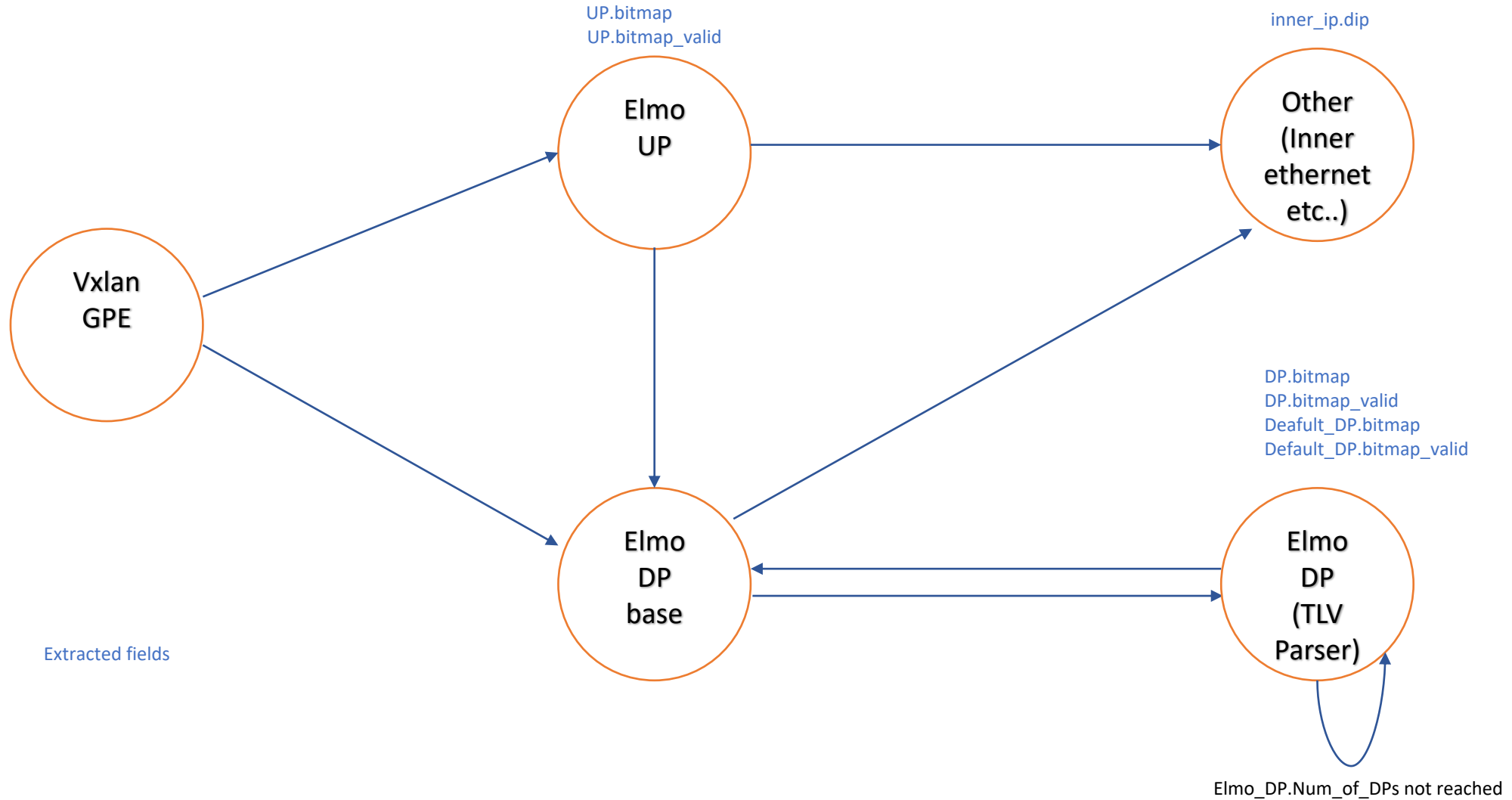


# Mellanox Programmable switch model



- Hybrid – Integration between legacy (switch router) and programmable pipeline
- NOS (ONVX / SONiC) and user applications run in parallel

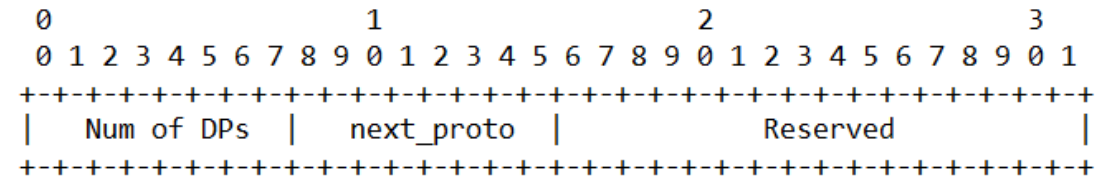
# Parser state machine



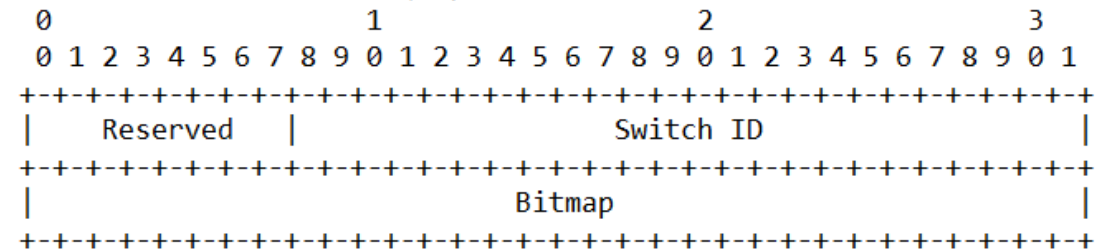
# Options parser

- Options are common among network protocols (IPv4, TCP, etc.. )
- Options follows some common structure
  - Base header has a known length
  - Total header length (computed)
  - Total options length
  - Options are built in a TLV fashion:
    - Type (self-indicator)
    - Length (some granularity)
    - Type and Length fields are fixed
  - This structure mainly exists to support unknown options
  - State transition is defined in the base header
- In Elmo:
  - Downstream P-rules are options
  - Unknown switch ID
  - Default p-rule – common Switch ID

Elmo downstream header (DP base):



Elmo downstream P-Rule (DP):



# Single switch Functionality

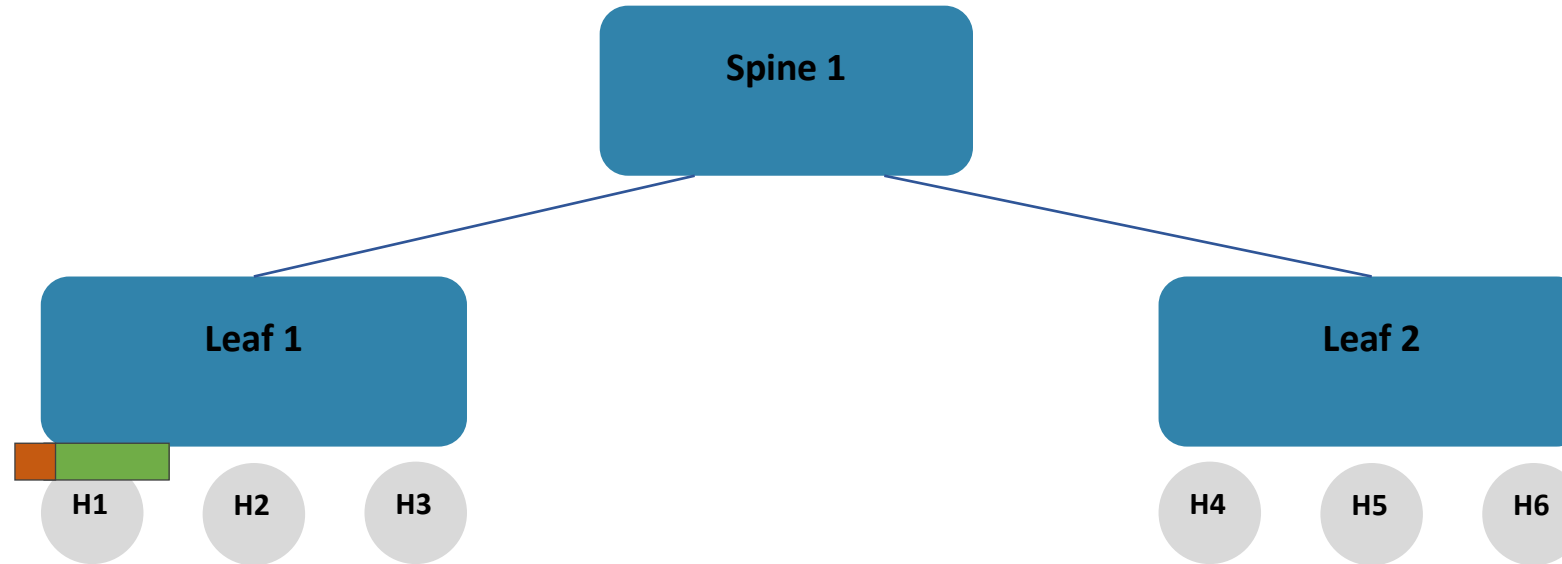
- Upstream bitmap
- Downstream bitmap
- Default p-rule
  - Increase scale on the expense of excess traffic
- Normal forwarding by the legacy pipeline

```
action set_egress_ports_dp() {  
    metadata.egress_ports = headers.elmo_downstream_p_rule.bitmap;  
}
```

```
table initial_bitmap_table {  
    key = {  
        headers.elmo_upstream_p_rule.isValid() : exact;  
        headers.elmo_downstream_p_rule.isValid() : ternary;  
    }  
}
```

```
apply  
{  
    if (initial_bitmap_table.apply().miss) {  
        if (s_rule_table.apply().miss) {  
            default_bitmap_table.apply();  
        }  
    }  
}
```

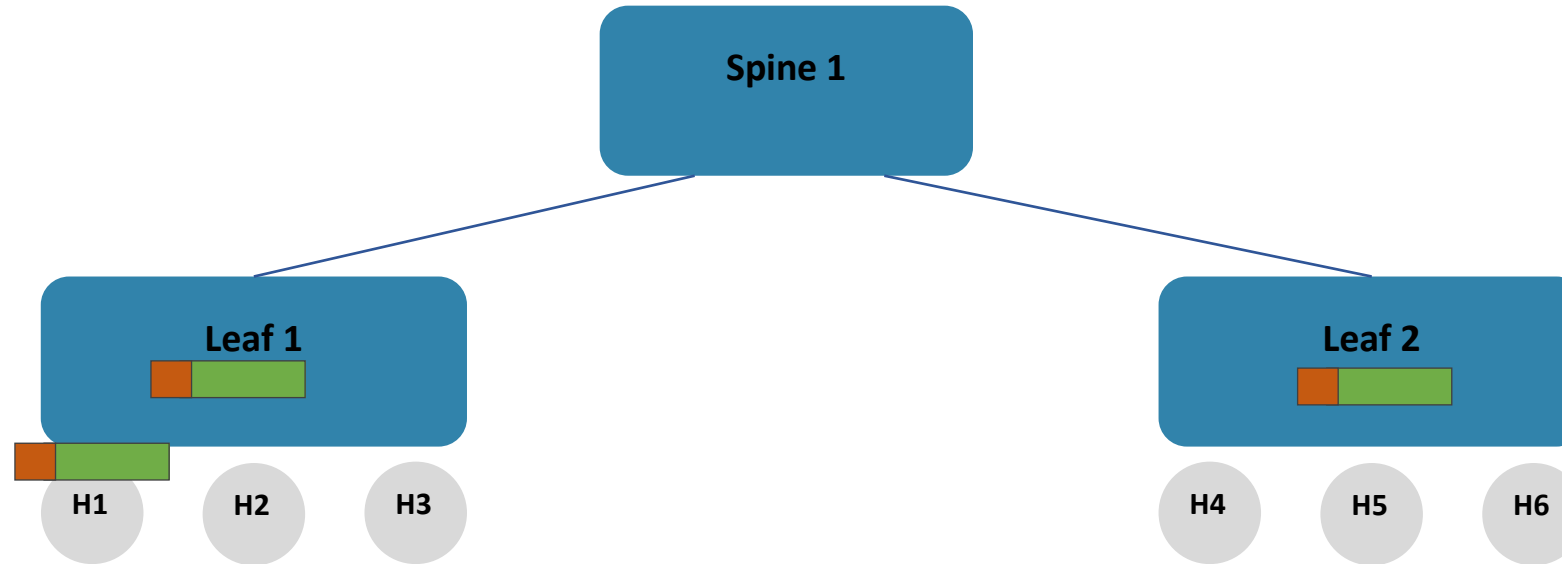
# Demo



## ■ Mcast groups:

1. H1 transmits to H2:  
UP bitmap : [000, 010]  
no DP header
2. H1 transmits to H3, H4 and H6:  
UP bitmap: [001, 001]  
DP: [S1: 001, L2: 101]
3. H1 transmits to H2, H4, H5:  
UP bitmap: [001, 010]  
DP:[S1: 001, L2: 110]
4. H1 transmits to H5, H6:  
UP bitmap: [001, 000]  
DP: [S1: 001, default: 011]

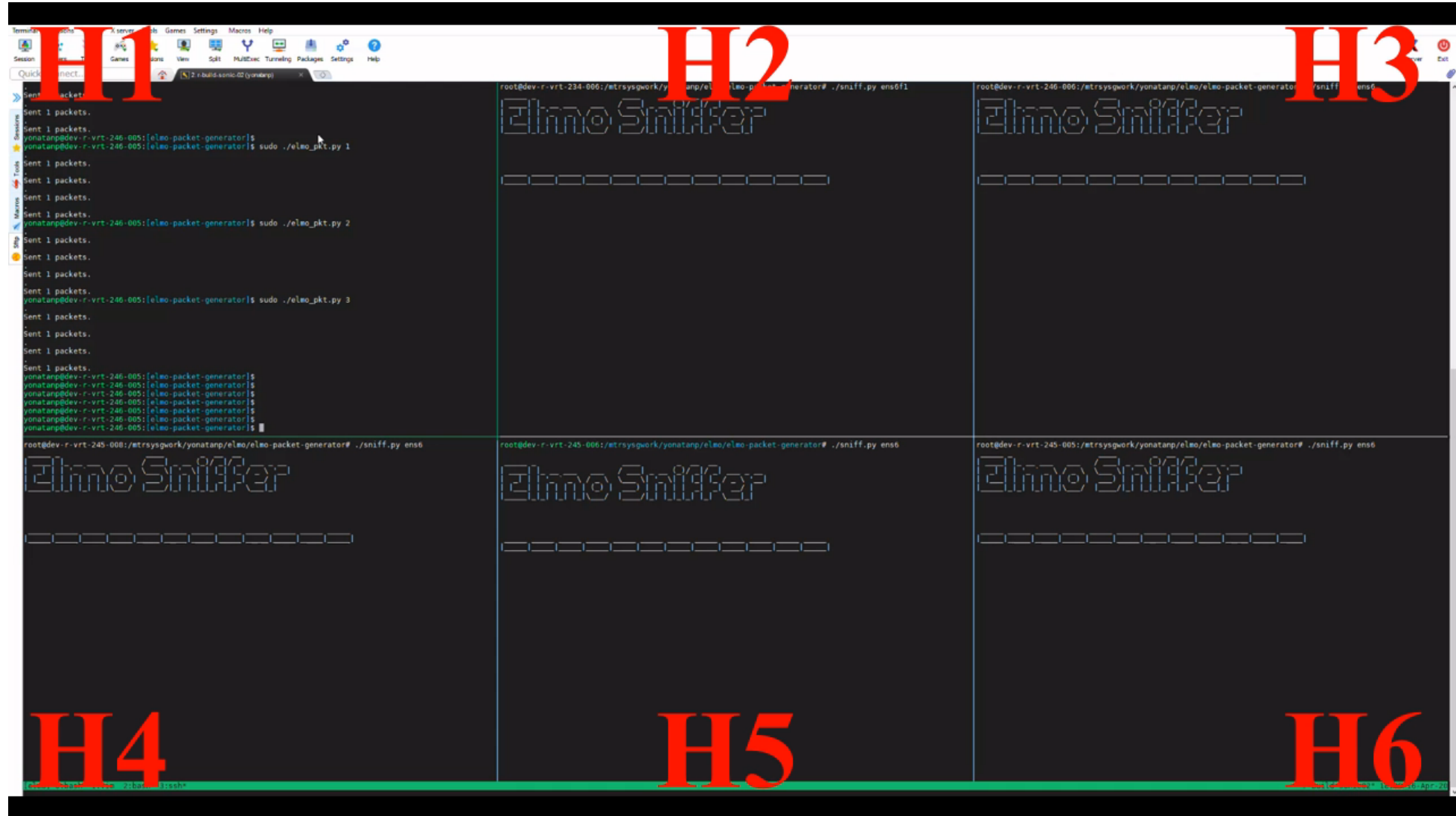
# Demo



## ■ Mcast groups:

1. H1 transmits to H2:  
UP bitmap : [000, 010]  
no DP header
2. H1 transmits to H3, H4 and H6:  
UP bitmap: [001, 001]  
DP: [S1: 001, L2: 101]
3. H1 transmits to H2, H4, H5:  
UP bitmap: [001, 010]  
DP:[S1: 001, L2: 110]
4. H1 transmits to H5, H6:  
UP bitmap: [001, 000]  
DP: [S1: 001, default: 011]

# Demo



# Challenges

- In the following slides, we'll share our experience from this work.
- Challenges encountered during this work:
  - Multicast
  - Options parsing
  - Extraction



# Multicast

- Multicast is not handled by the PSA model (extern)
- Hard for stateless switch multicast
- This work - directly expose MC bitmap to the dataplane:
  - `metadata.egress_ports = headers.elmo_downstream_default_p_rule.bitmap;`
- Multicast group table can be easily supported
- Hybrid architecture - support non-physical ports as well (e.g., router interface etc..)

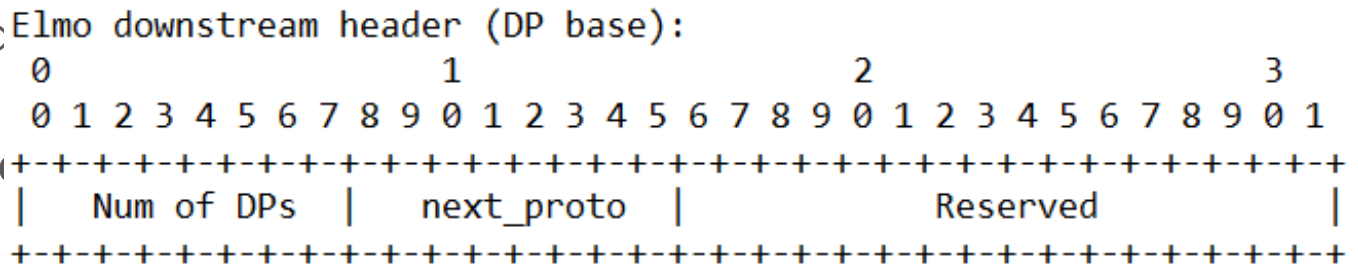
# Options parsing

- Options current implementation in P4
  - possible but not trivial
  - Not easily offloadable
- Common use case
  - Worthwhile to have standard fashion of defining
  - Easily HW offloaded by the different vendors.
- Build a sub-parser prototype which follows the observed structure

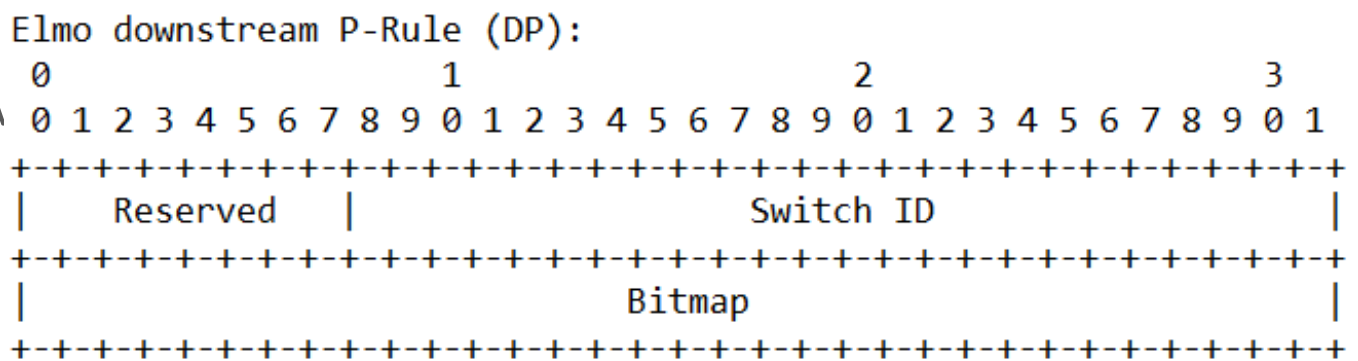
# Extractions

- Extract methods act on entire headers –consumes them and advances the cursor
- It is further assumed that HW will extract all the fields of the accepted header.
  - Could be costly

- What if you require a subset of the fields



- May prevent HW optimizing by selectively
  - Dynamically loaded control



- Advanced field extraction features like

# Extraction - proposal

- We implemented field extraction primitives on our architecture:
  - `Void extract_field<T>(out T headerLvalue.field);`
  - `void extract_field<T>(out T headerLvalue.field, in bit<32> variableFieldOffset);`
  - `void extract_field<T>(out T headerLvalue.field, in bit<32> variableFieldSizeInBits, in bit<32> variableFieldOffset);`
- Extract a single field and advance the cursor,
  - Adds to current header primitives (not replace)
- Useful also for:
  - variable offset fields
  - more than one variable length field in a header
- Another option - Usage analysis in the compiler backend
  - Sufficient for monolithic P4 executables
  - Problematic for target architectures which allow dynamic insertion of control pipelines (which share the same parser)

# Conclusions

- Elmo compactly encodes multicast policy inside packets
- Designed for multi-tenant data centers scales
- Demonstrated, for the first time, Elmo implementation with wire speed performance using hybrid programmable dataplane
- All legacy forwarding and control plane is intact



Sponsored By



## Thank You

Matty Kadosh (Mellanox), [mattyk@mellanox.com](mailto:mattyk@mellanox.com)  
Lalith Suresh (VMware), [lsuresh@vmware.com](mailto:lsuresh@vmware.com)

## Acknowledgment

Jen Rexford<sup>1</sup>

Nick Feamster<sup>2</sup>

Ori Rottenstreich<sup>3, 5</sup>

Mukesh Hira<sup>4</sup>, Mihai Budiu<sup>4</sup>, Ben Pfaff<sup>4</sup>

Alan Lo<sup>5</sup>, Aviv Kfir<sup>5</sup>, Jose Yallouz<sup>5</sup>

[1] Princeton [2] U. Chicago [3] Technion [4] VMware [5] Mellanox