# Lecture 11: In-Network Computing
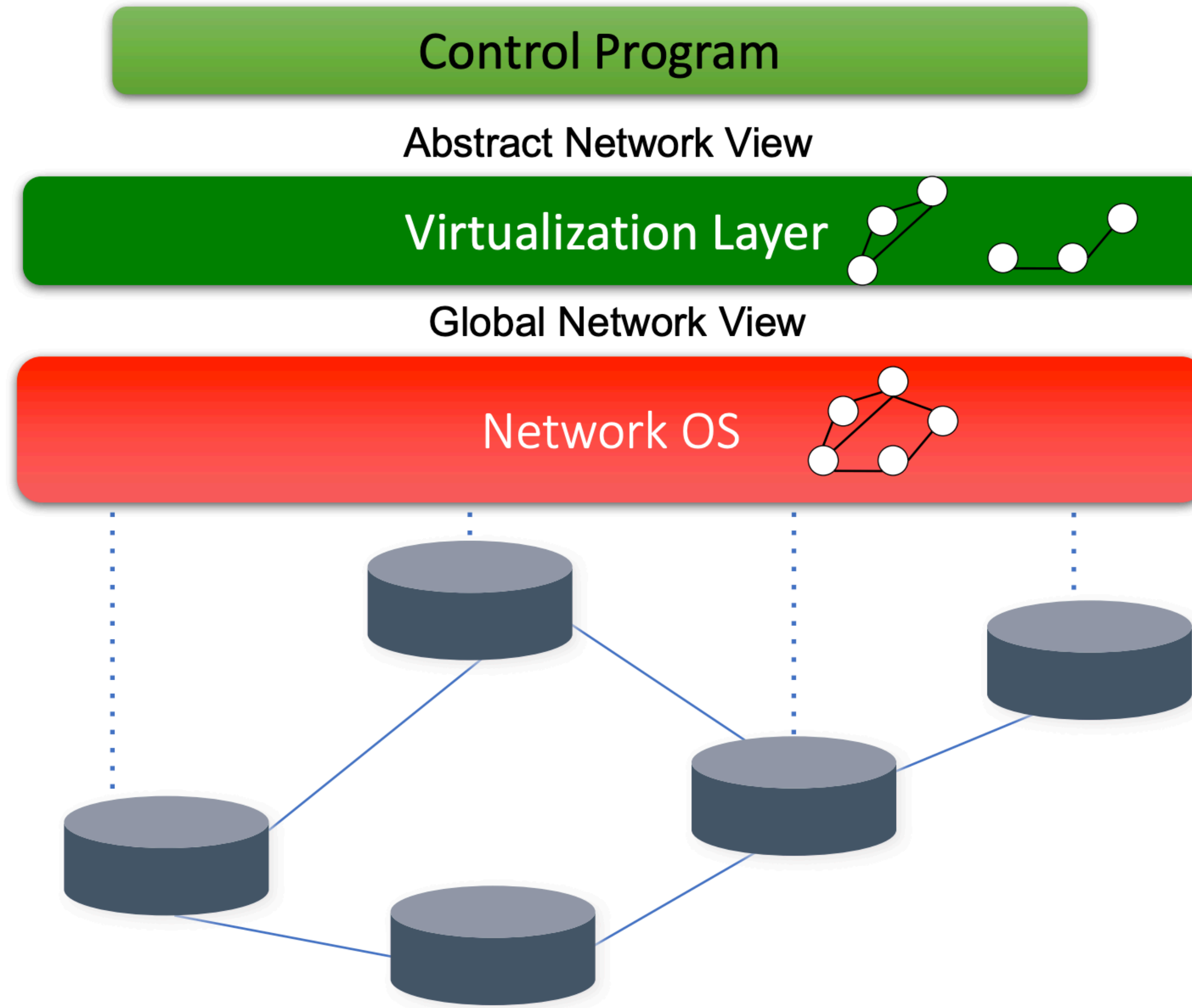
## CS 234 / NetSys 210: Advanced Computer Networks

Sangeetha Abdu Jyothi
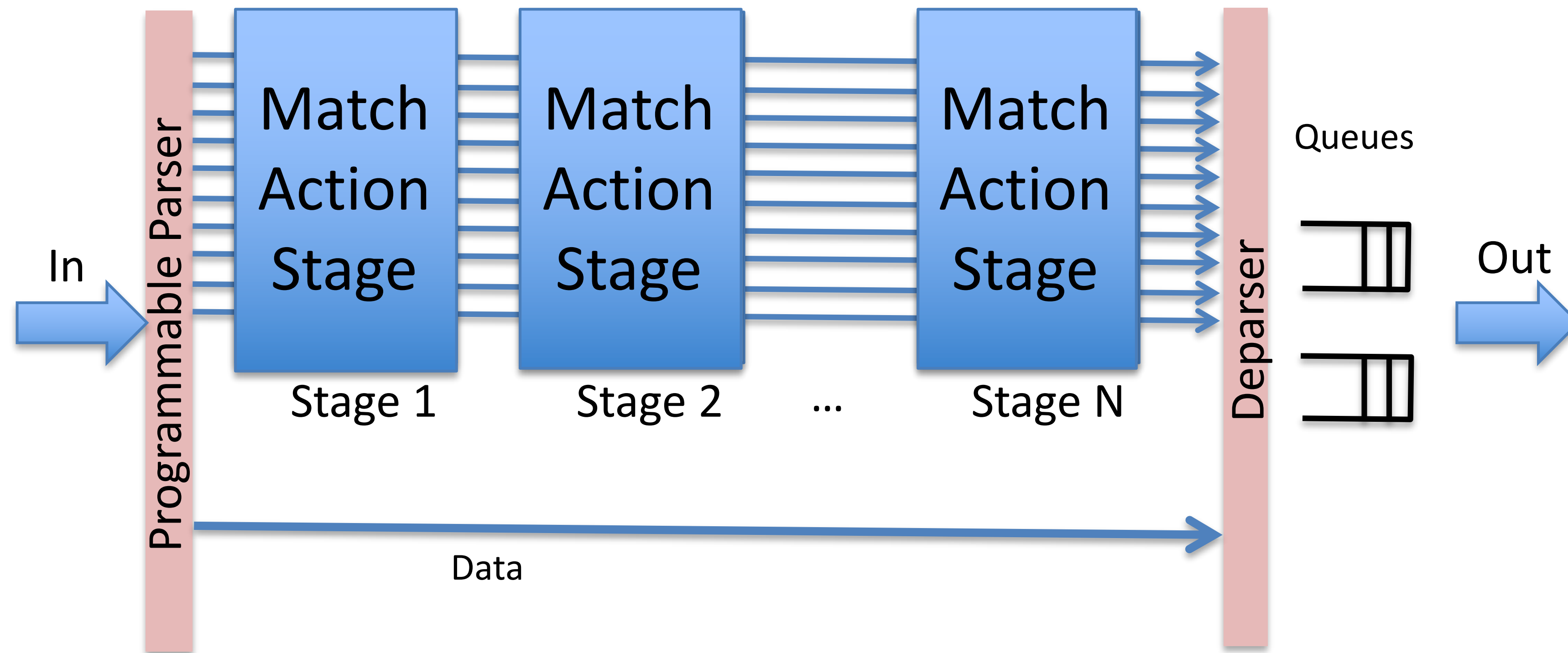
UCIRVINE

# Elmo: Source-Routed Multicast for Public Clouds

Unicast
(one-to-one)

Multicast
(one-to-many)

Broadcast
(one-to-all)

# One to Many Communication Pattern in Cloud

# Limitations of Native Multicast



Processing overhead

Controller

Excessive control churn
due to membership and topology changes

Limited group entries

# Limitations of Unicast-based Alternatives

# Elmo: Source-Routed Multicast for Cloud Services

- Key challenges:

    - How to efficiently encode multicast forwarding policy inside packets?

    - How to process this encoding at line rate?

# 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 Naive Source Routed Multicast

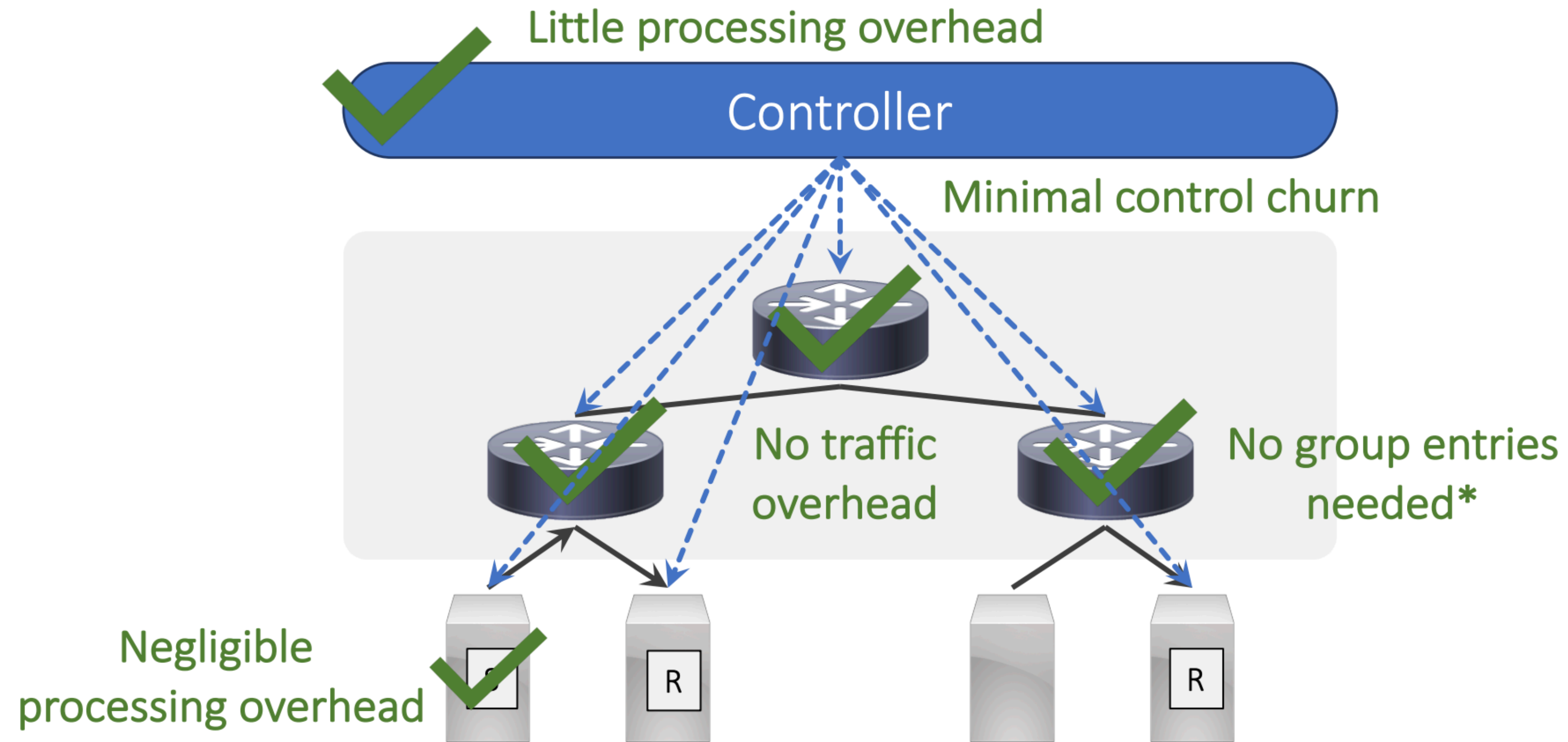A **multicast group** encoded as a list of **(Switch, Ports)** pairs

```
Switch 1: [Ports ]
Switch 2: [.. .. ..]
Switch 3: [.. .. ..]
Switch 4: [.. .. .x ..]
Switch 5: [.x .. .. ..]
```
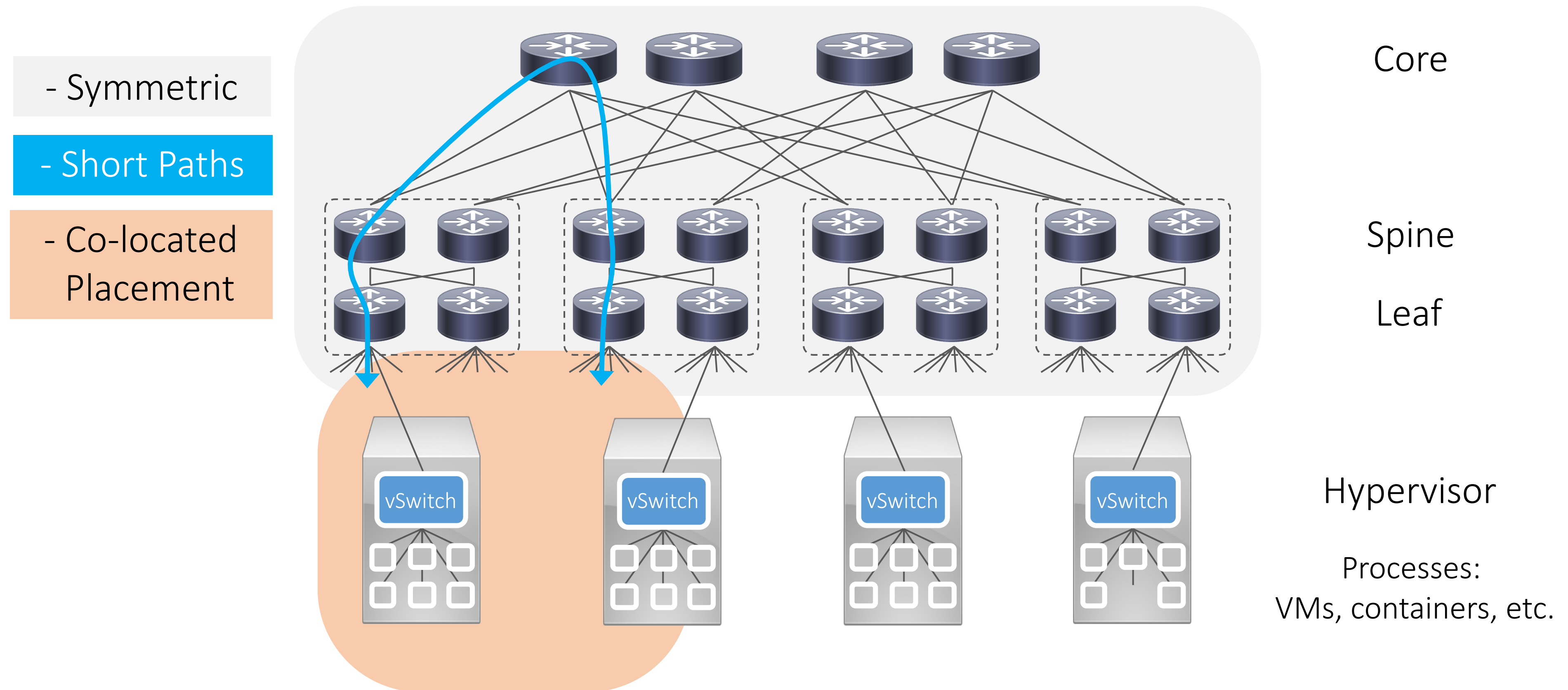
For a data center with:
- **1000** switches
- **48 ports** per switch

**Not Scalable!**

# Exploiting DC Characteristics for efficient encoding

- Symmetric

- Short Paths

- Co-located Placement

Core

Spine

Leaf

Hypervisor

vSwitch

vSwitch

vSwitch

vSwitch

Processes:
VMs, containers, etc.

# Programmable Switches for Line Rate Processing

Barefoot Tofino or
Cavium XPliant

PISCES    PISCES    PISCES    PISCES

15

# 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 .. .. ..]
```

**1** Encode switch ports as a bitmap

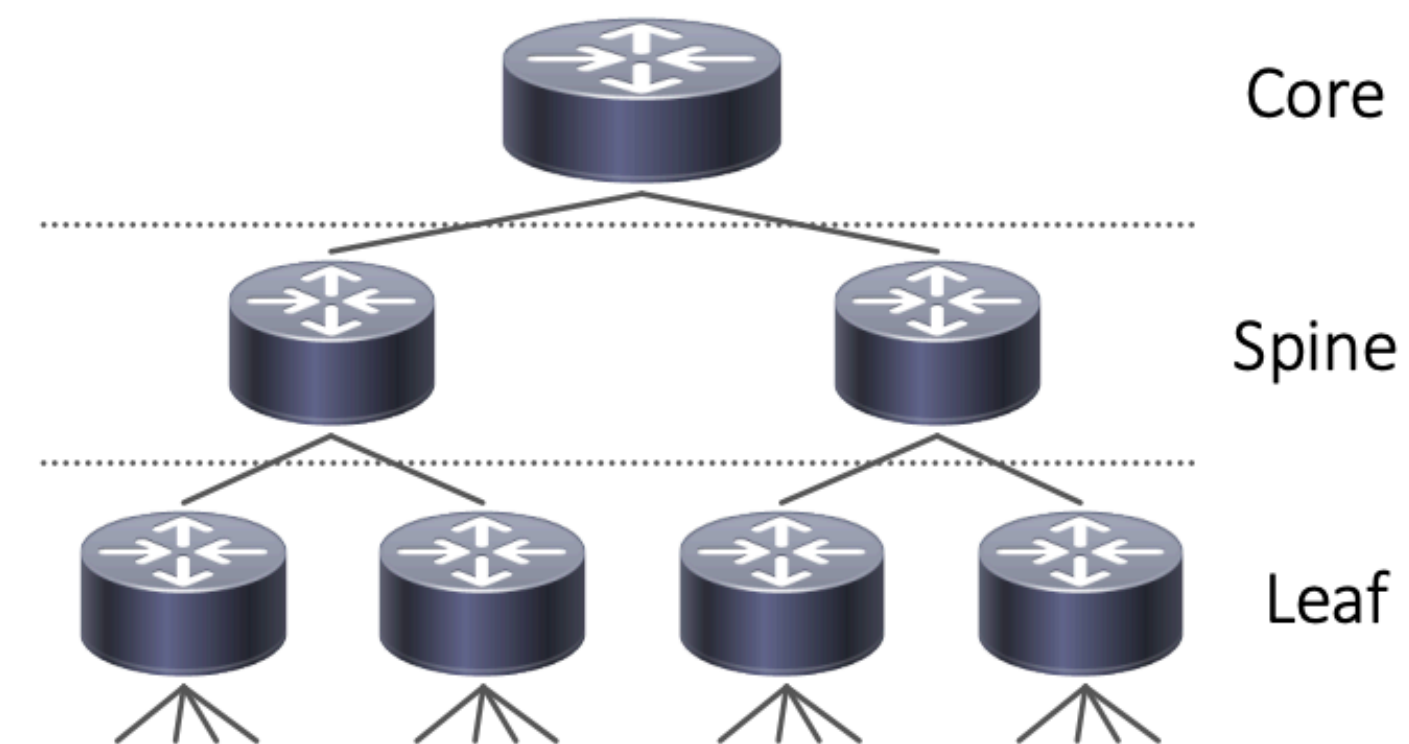**Bitmap** is the internal **data structure** that switches use for **replicating packets**

16

# Encoding a Multicast Policy in Elmo

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

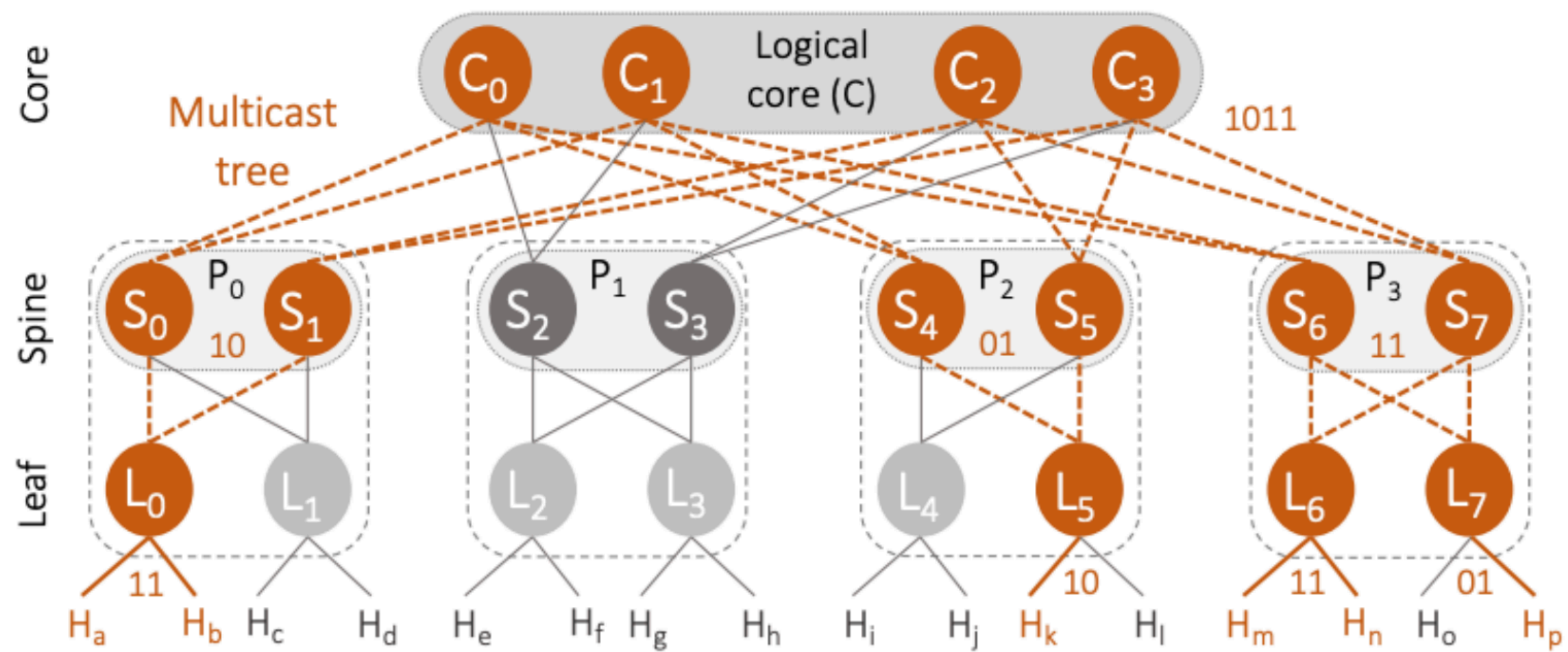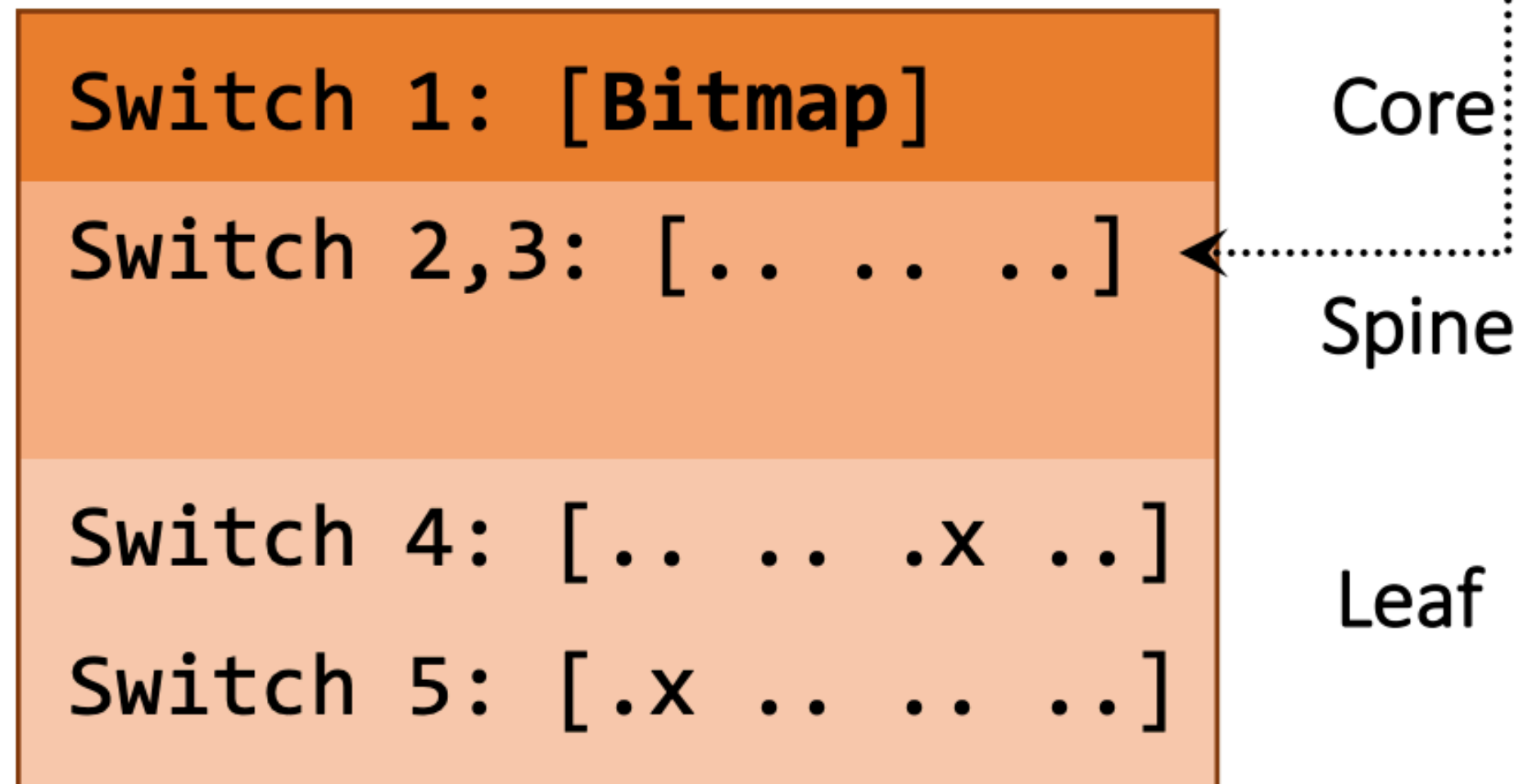| |
|---|
| Switch 1: [Bitmap] |
| Switch 2: [.. .. ..] |
| Switch 3: [.. .. ..] |
| Switch 4: [.. .. .x ..] |
| Switch 5: [.x .. .. ..] |

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

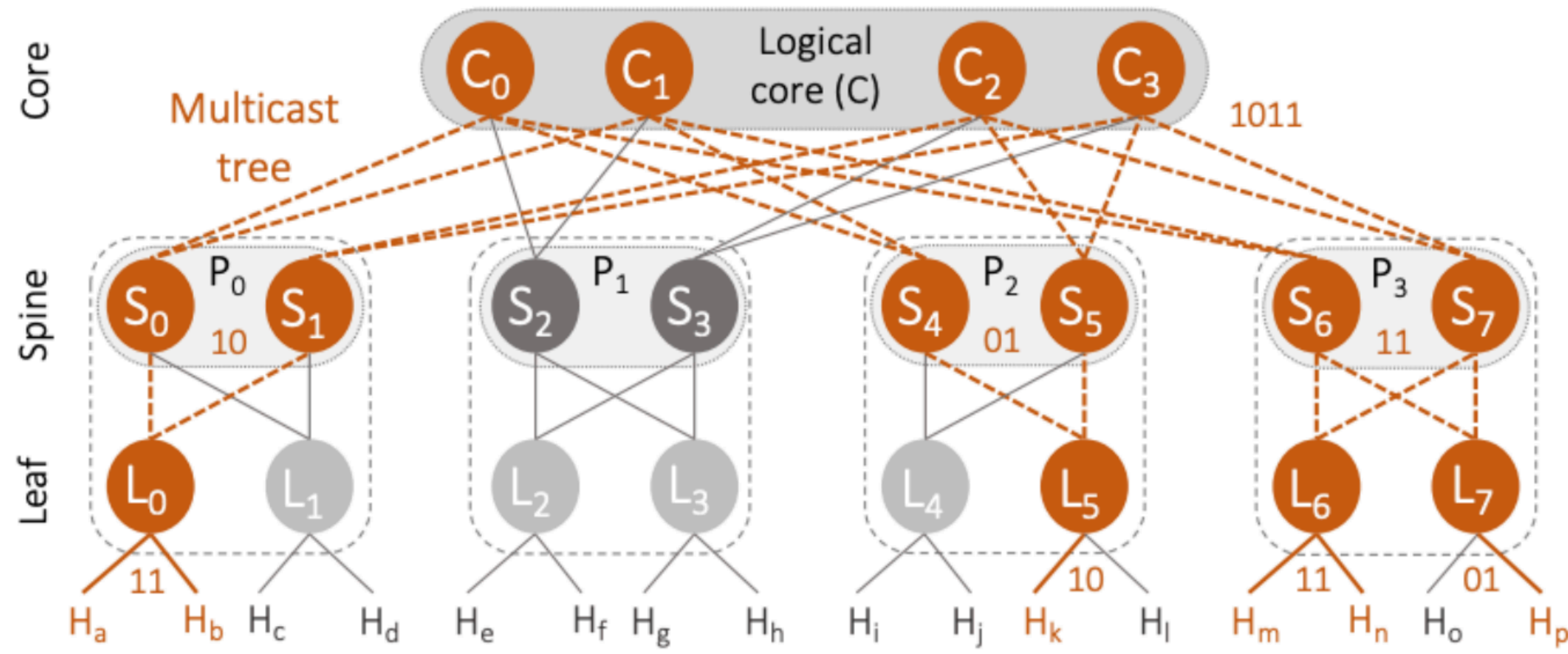A multicast group encoded as a list of (Switch, Ports) pairs

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

Core

Spine

Leaf

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

19

# 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

Default Bitmap
Switch Table Entries

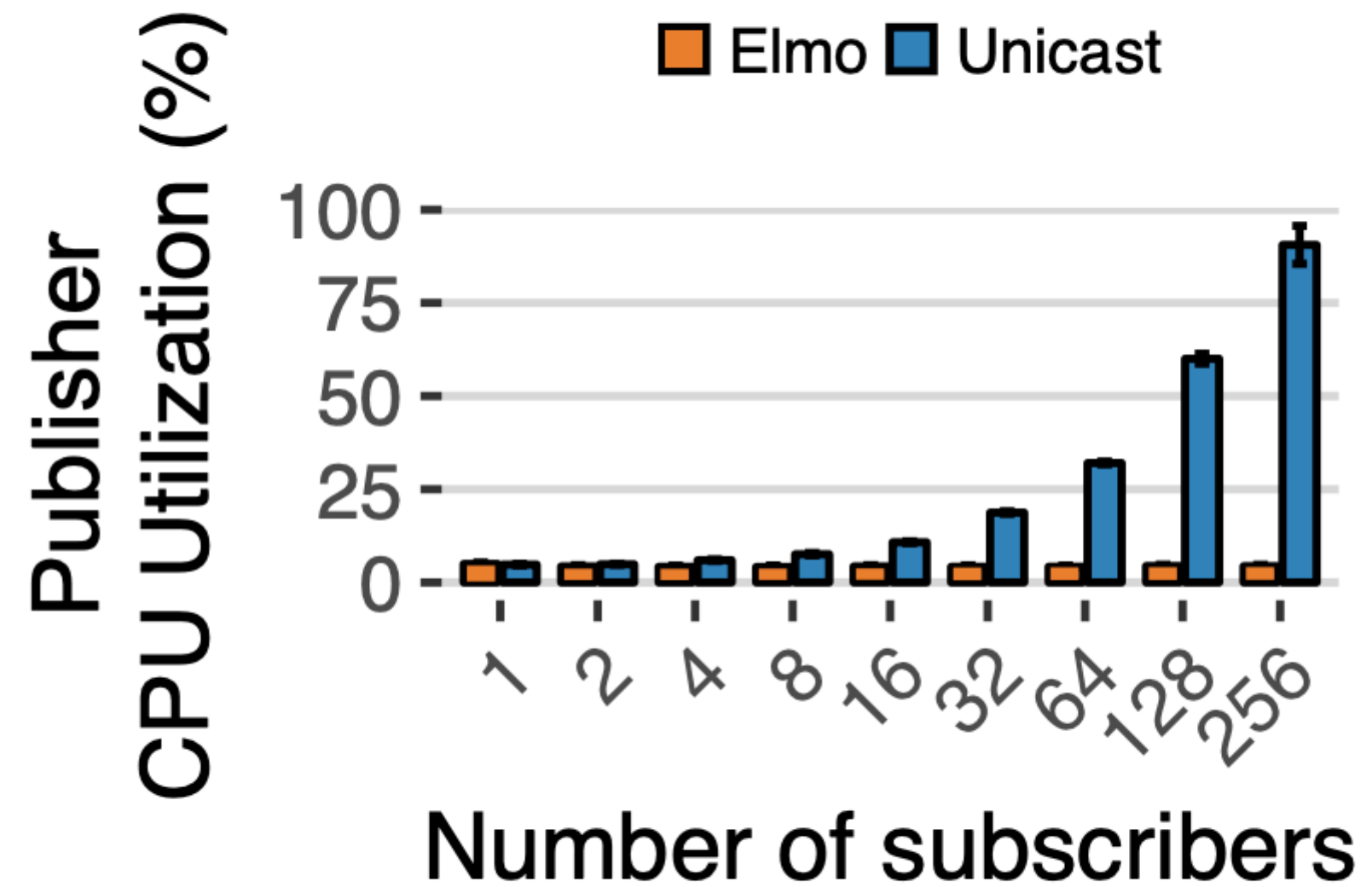**5** Use switch entries and a default bitmap for larger groups
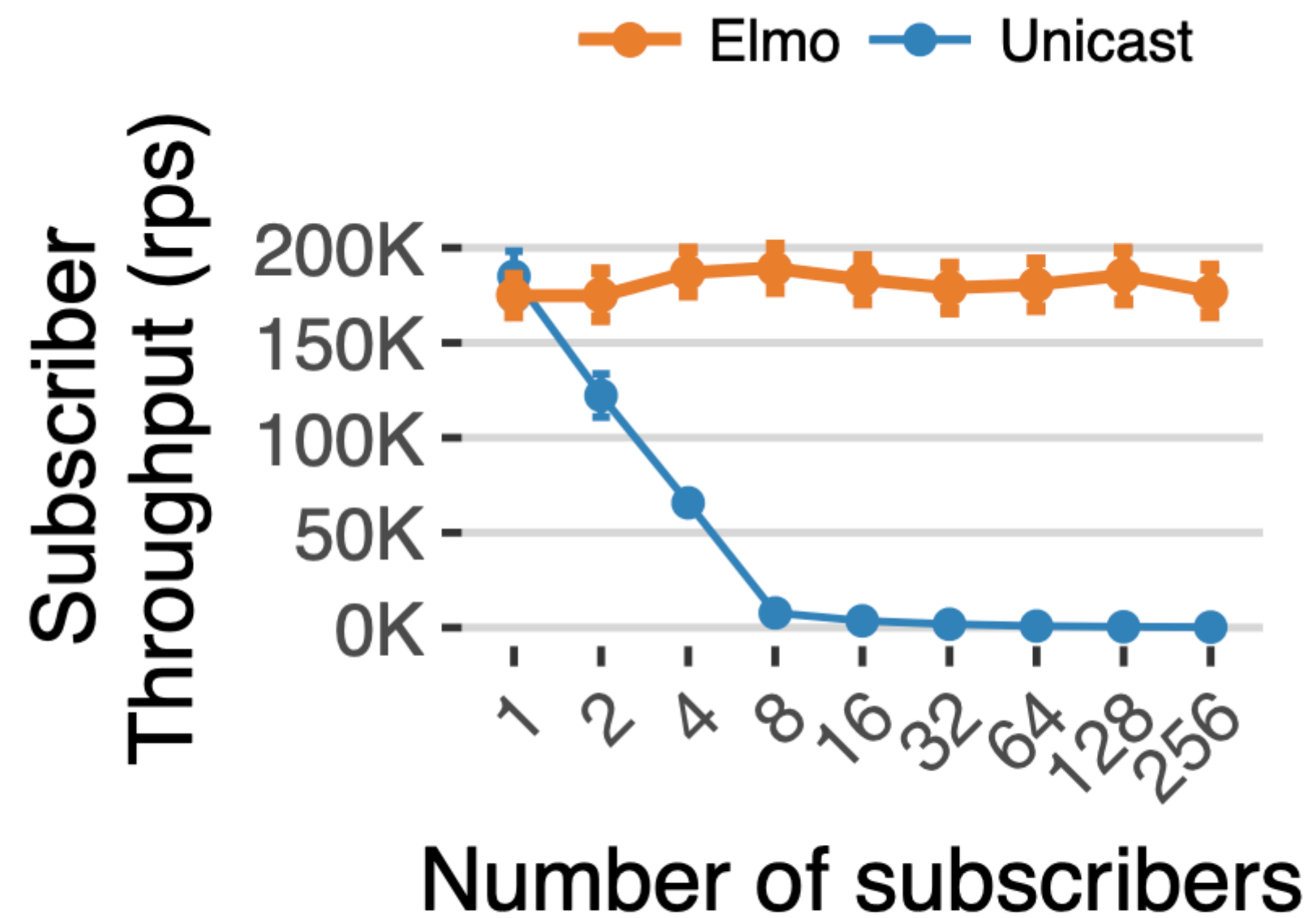
# Encoding a Multicast Tree in Elmo

- Key design decisions:

  - Encoding switch output ports in a bitmap

  - Encoding on the logical topology

  - Sharing bitmap across switches

  - Dealing with limited header space using default p-rules

  - Reducing traffic overhead using s-rules
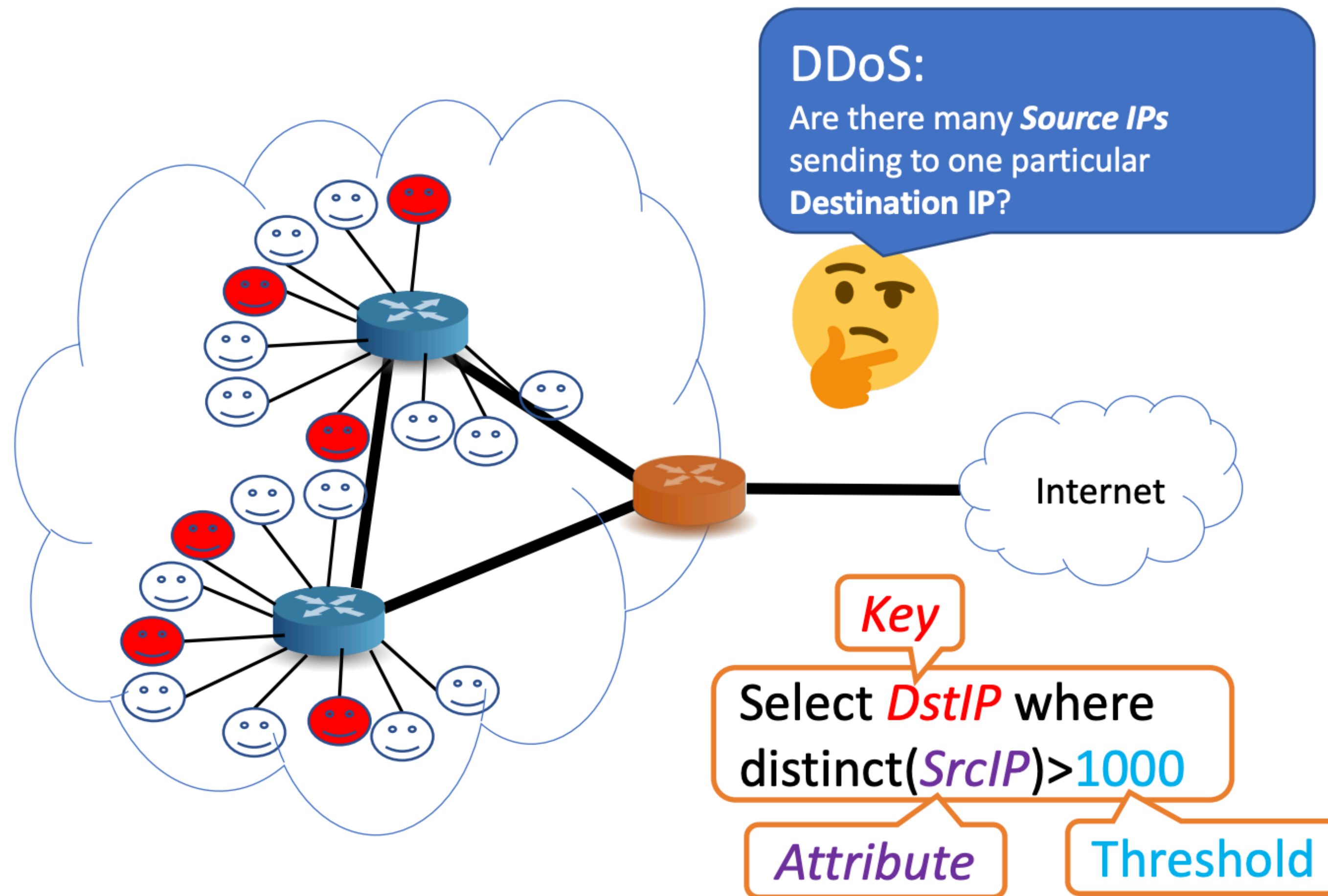
**Supports a Million groups!**

# More In-Network Computing-Based Solutions

Worker 1 updates    Worker 2 updates    ...    Worker N updates

Worker 1    Worker 2    Worker 3    Worker 4

**Aggregate model updates in-network**

Switch

26

## Challenges

</> Limited computation

Limited storage

**No floating points**

Packet loss



**6.5 Tbps**
programmable
data plane

# 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

  - HPCC: High Precision Congestion Control, SIGCOMM'19

- A new protocols for more efficient L2 switching

  - The Deforestation of L2, SIGCOMM'16

# Other app-level use cases

- NetChain [SOSP'17]: in-network key-value store

- NetLock [SIGCOMM'20]: Switching support to manage locks

- NetPaxos [SOSR'15]: implement Paxos on programmable switches

- NoPaxos [OSDI'16]: in-network primitives for distributed protocols

# Thanks!