

# In-Network Address Caching for Virtual Networks

Lior Zeno, Ang Chen, Mark Silberstein



Virtual networks *enable* the public cloud,  
but they are also *difficult* to manage.

*Virtual-to-physical IP translation* is *challenging*,  
but *necessary* for packet forwarding.



**Sailfish: Accelerating Cloud-Scale Multi-Tenant Multi-Service Gateways with Programmable Switches**

**Bluebird: High-performance SDN for Bare-metal Cloud Services**

Manikandan Arumugam<sup>1</sup>, Deepak Bansal<sup>3</sup>, Navdeep Bhatia<sup>1</sup>, James Roemer<sup>3</sup>, Simon Canner<sup>1</sup>

ABSTRACT  
The cloud  
hub of  
gateways

**Andromeda: Performance, Isolation, and Velocity at Scale in Cloud Network Virtualization**

Mic  
Dima  
The bare  
as a Ser  
tomers a  
Nil

**Achelous: Enabling Programmability, Elasticity, and Reliability in Hyperscale Cloud Networks**

Chengkun Wei<sup>†\*</sup>, Xing Li<sup>†§\*</sup>, Ye Yang<sup>§†\*</sup>, Xiaochong Jiang<sup>†</sup>, Tianyu Xu<sup>†</sup>, Bowen Yang<sup>§</sup>, Taotao Wu<sup>§</sup>, Chao Xu<sup>§</sup>, Yilong Lv<sup>§</sup>, Haifeng Gao<sup>§</sup>, Zhenhao Zhang<sup>§</sup>, Zikang Chen<sup>§</sup>, Zeke Wang<sup>†</sup>, Zihui Zhang<sup>†</sup>, Shunmin Zhu<sup>§§□</sup>, Wenzhi Chen<sup>†□</sup>

<sup>†</sup>Zhejiang University <sup>§</sup>Alibaba Cloud <sup>□</sup>Tsinghua University

ABSTRACT Abstract

system together end to end. We developed *Andromeda*,

# What

In-network IP address caching

# Why

Need efficient updates & fast routing

# How

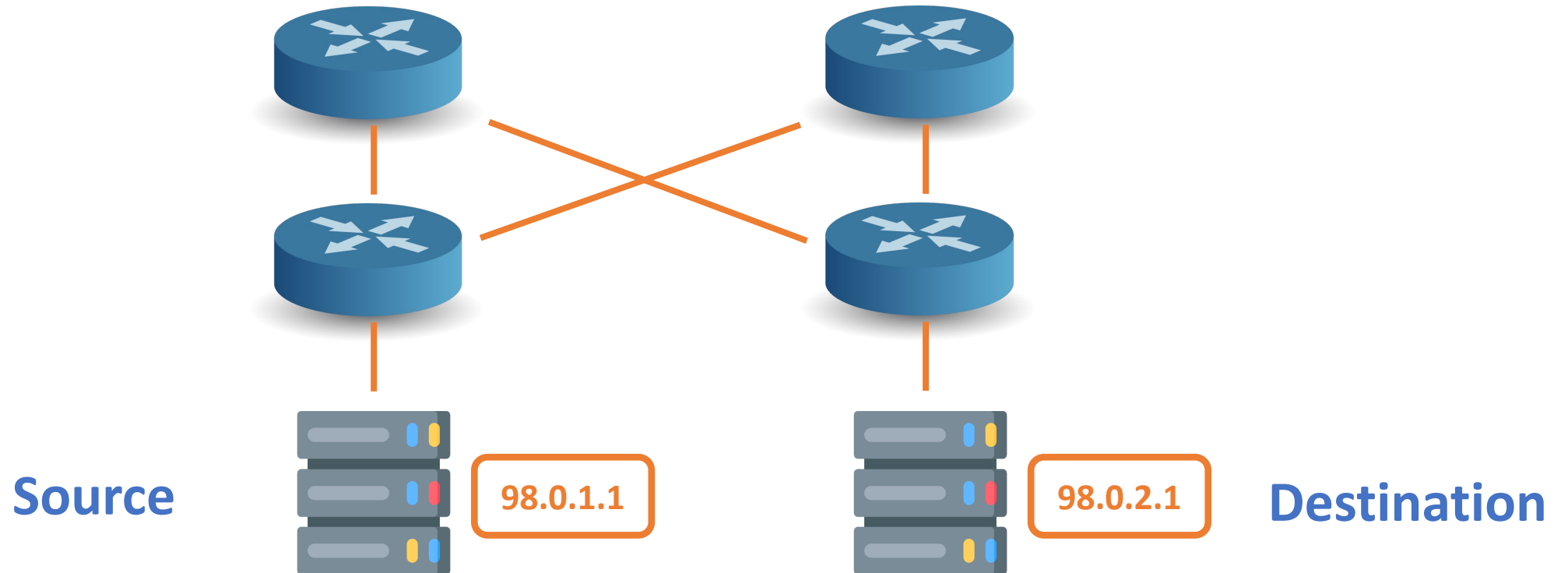
Systematic approach backed by distributed computing results

# Results

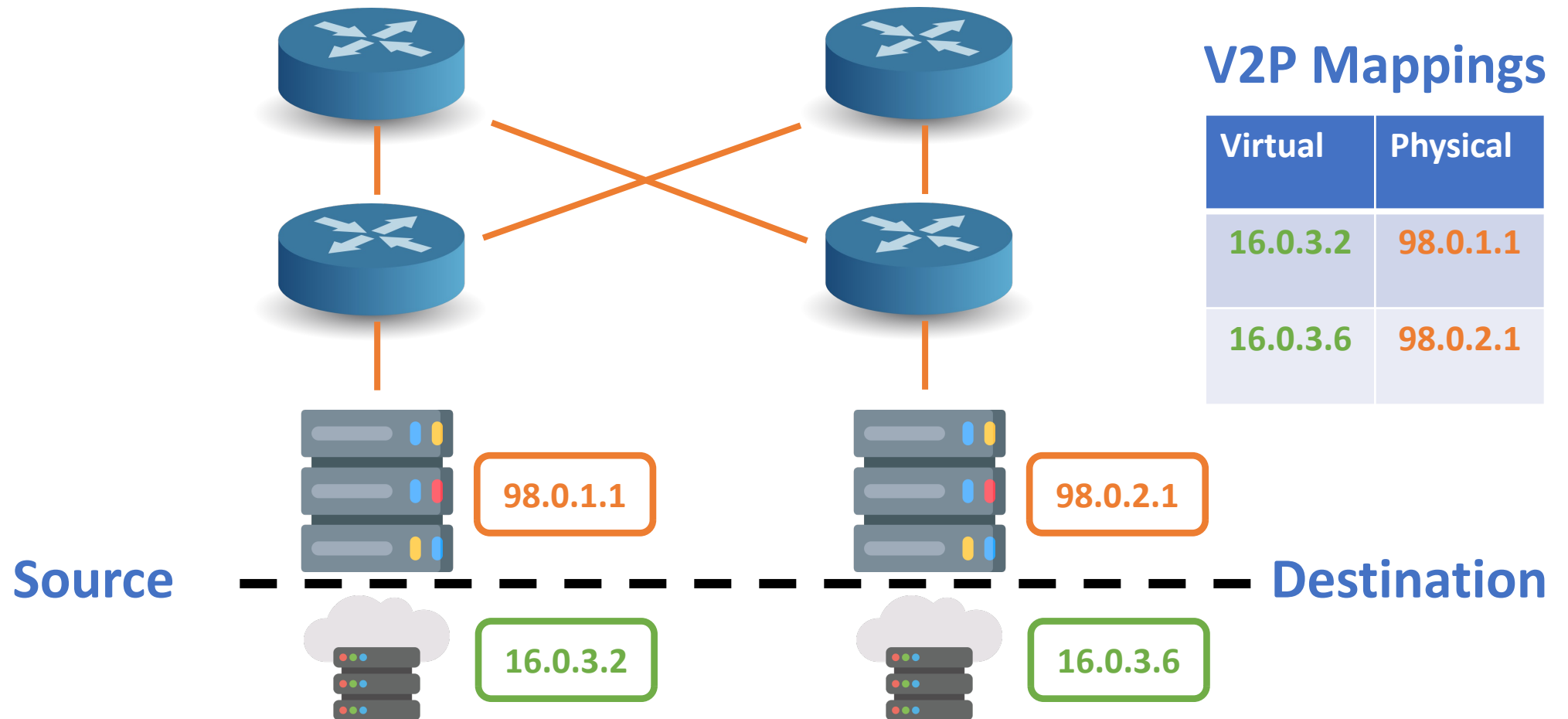
<b>Better</b>	<b>7.8×</b>	<b>4.3×</b>	<b>6.1×</b>
	<b>FCT</b>	<b>First-packet</b>	<b>Bandwidth overhead</b>

# Background:

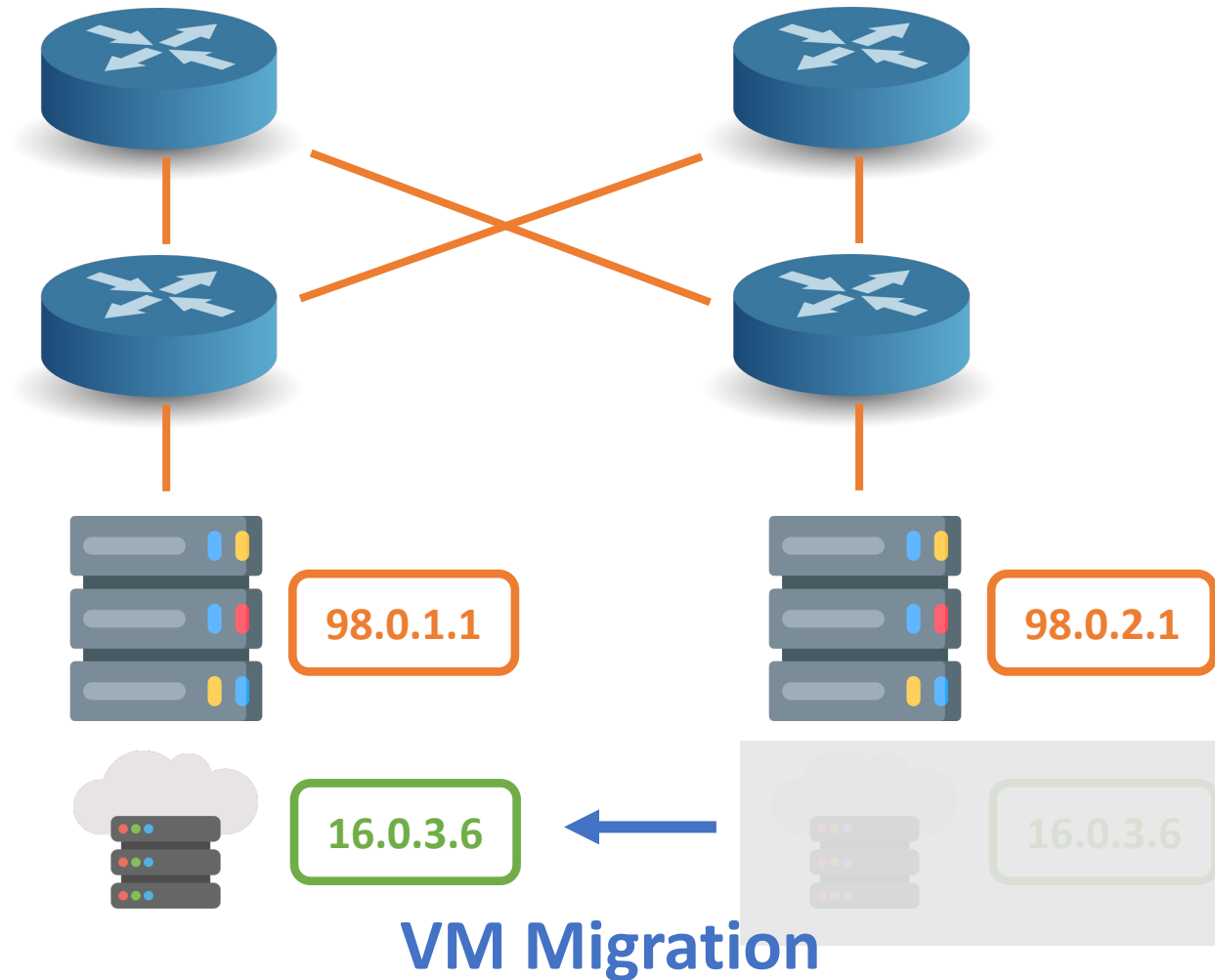
## Packet Forwarding in Physical Networks



# Background: Packet Forwarding in Virtual Networks



# Background: VM Roaming



## V2P Mappings

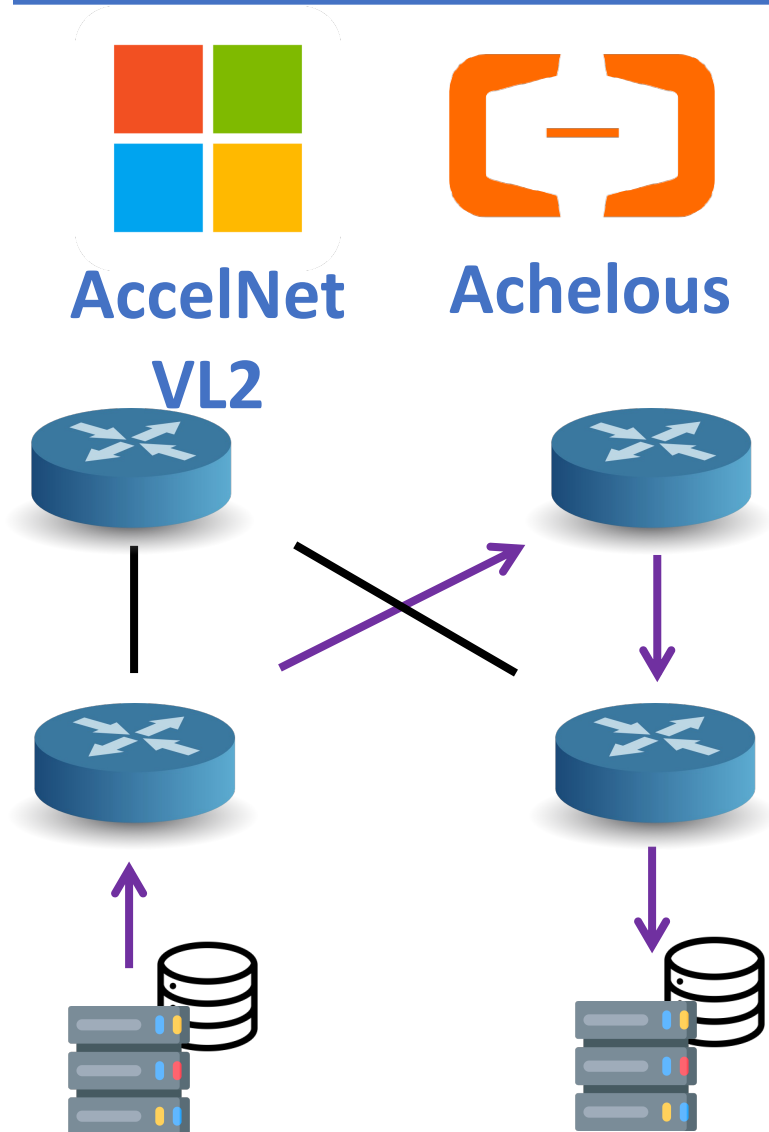
Virtual	Physical
16.0.3.6	98.0.2.1



Virtual	Physical
16.0.3.6	98.0.1.1

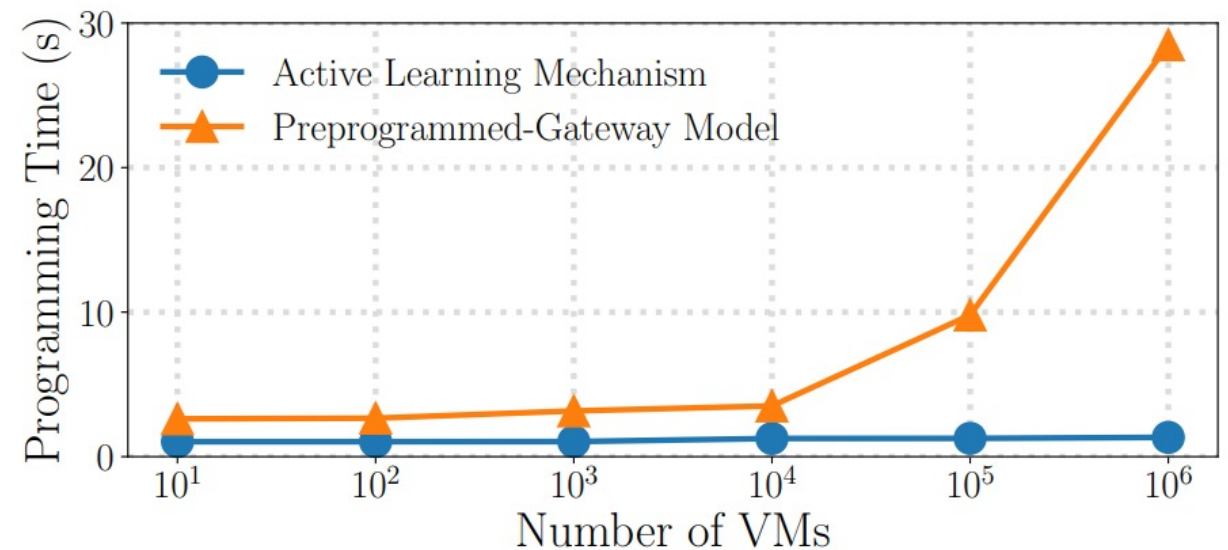
# Where to Perform Address Translation?

# Host-Driven Design



Fast packet forwarding

Slow updates: large networks – seconds, single server – 10s ms





# Gateway-Driven Designs

Bluebird



Sailfish



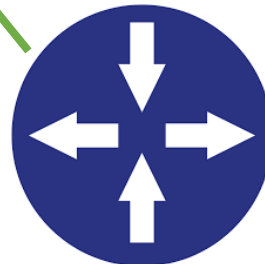
Zeta

HUAWEI

Andromeda

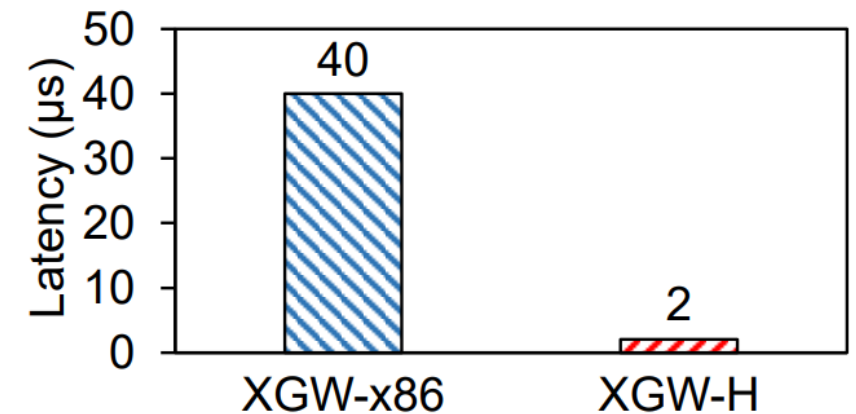


Gateway



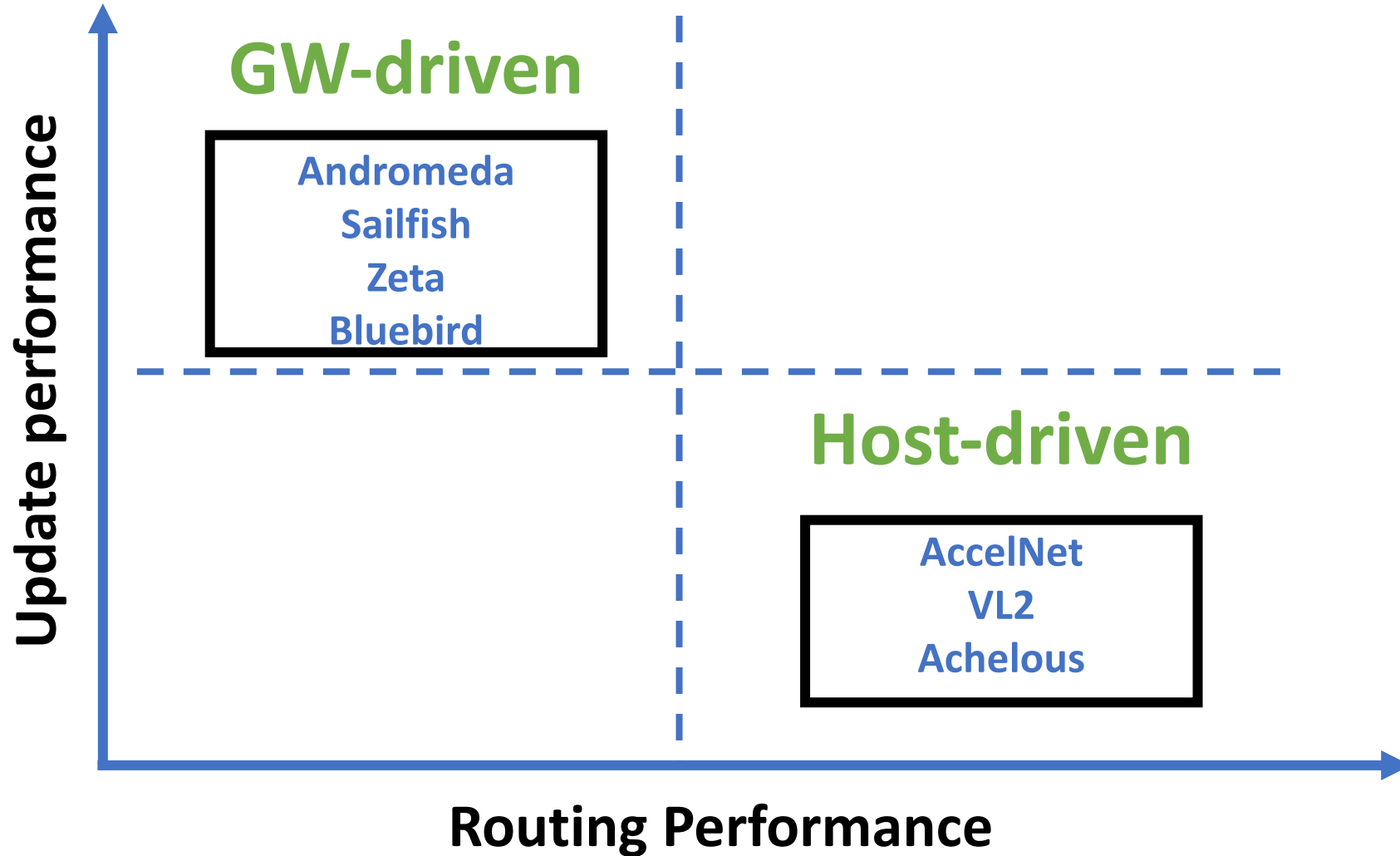
Fast updates

Slow packet forwarding (40 $\mu$ s processing overhead), network bandwidth overhead



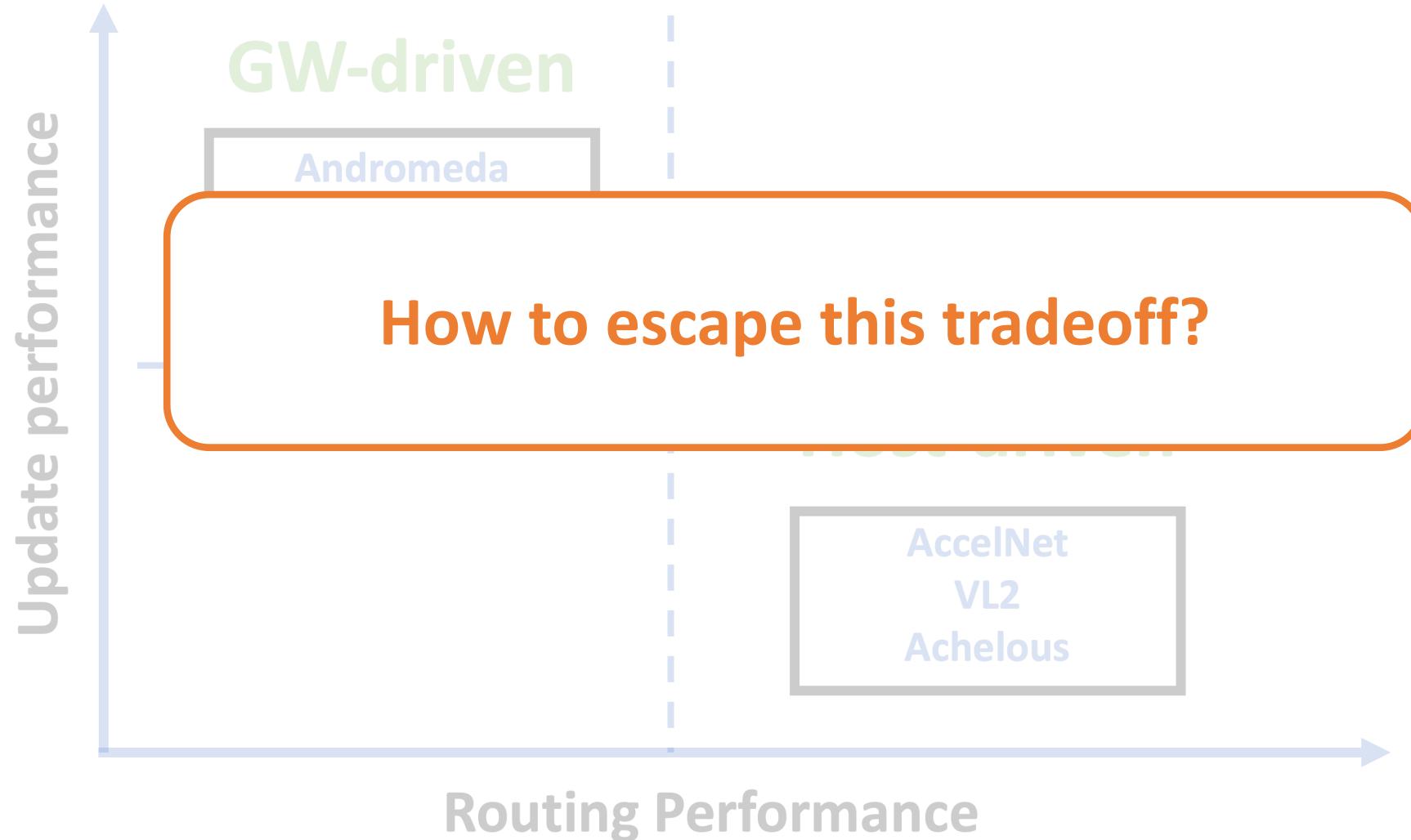
# The Tradeoff

---



# The Lookup-Update Tradeoff

---



# The Read (Lookup)-Write (Update) Tradeoff



Limiting the replication factor balances read and write costs

## The Dangers of Replication and a Solution

Jim Gray (Gray@Microsoft.com)

THE COST OF DATA REPLICATION +

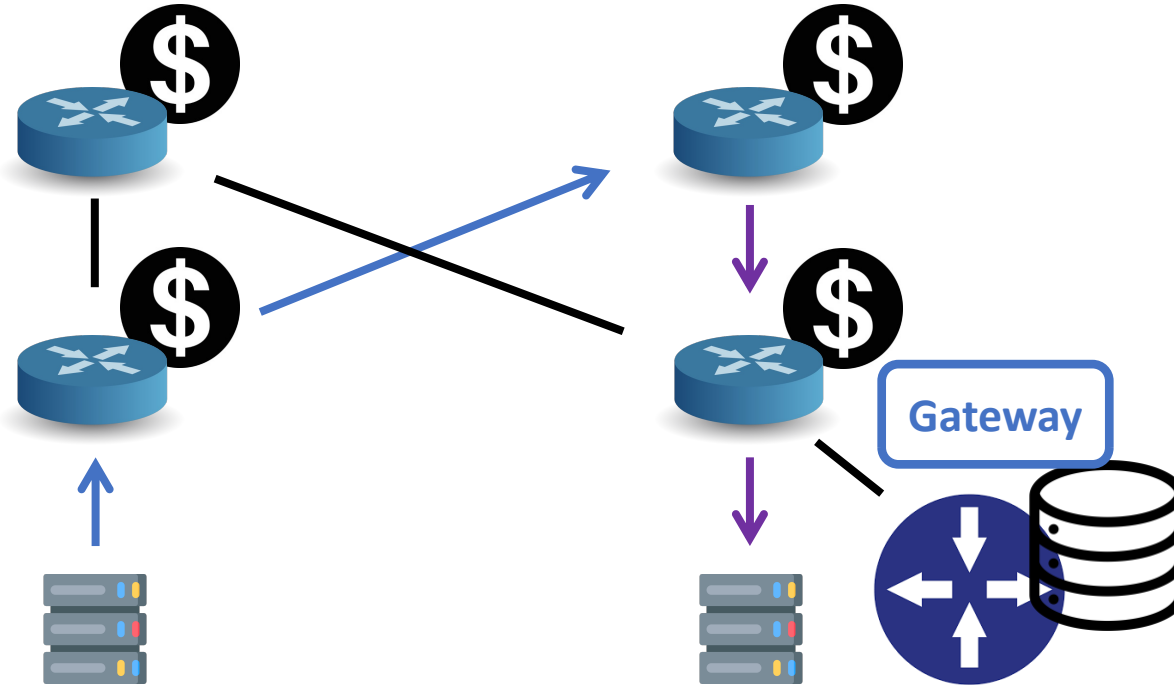
Hector Garcia-Molina  
Daniel Barbara

Department of Electrical Engineering and Computer Science  
Princeton University  
Princeton, New Jersey 08544

# Our Idea: In-Network Address Caching



Switches learn from traffic and *cache* V2P mappings within the data plane



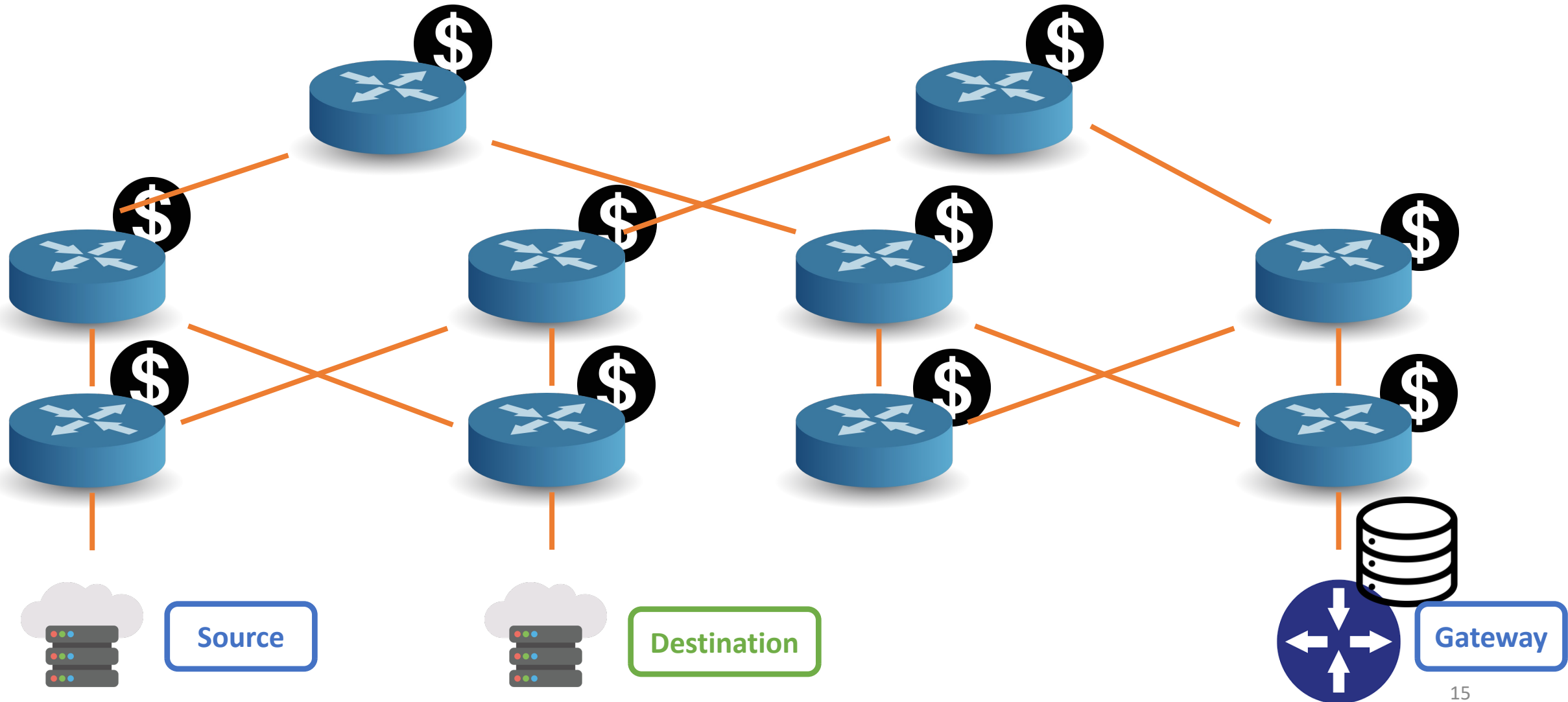
Fast routing  
Fast updates  
Reduced gateway costs  
No routing changes  
Incrementally deployable  
Compatible with existing switches

# Agenda

---

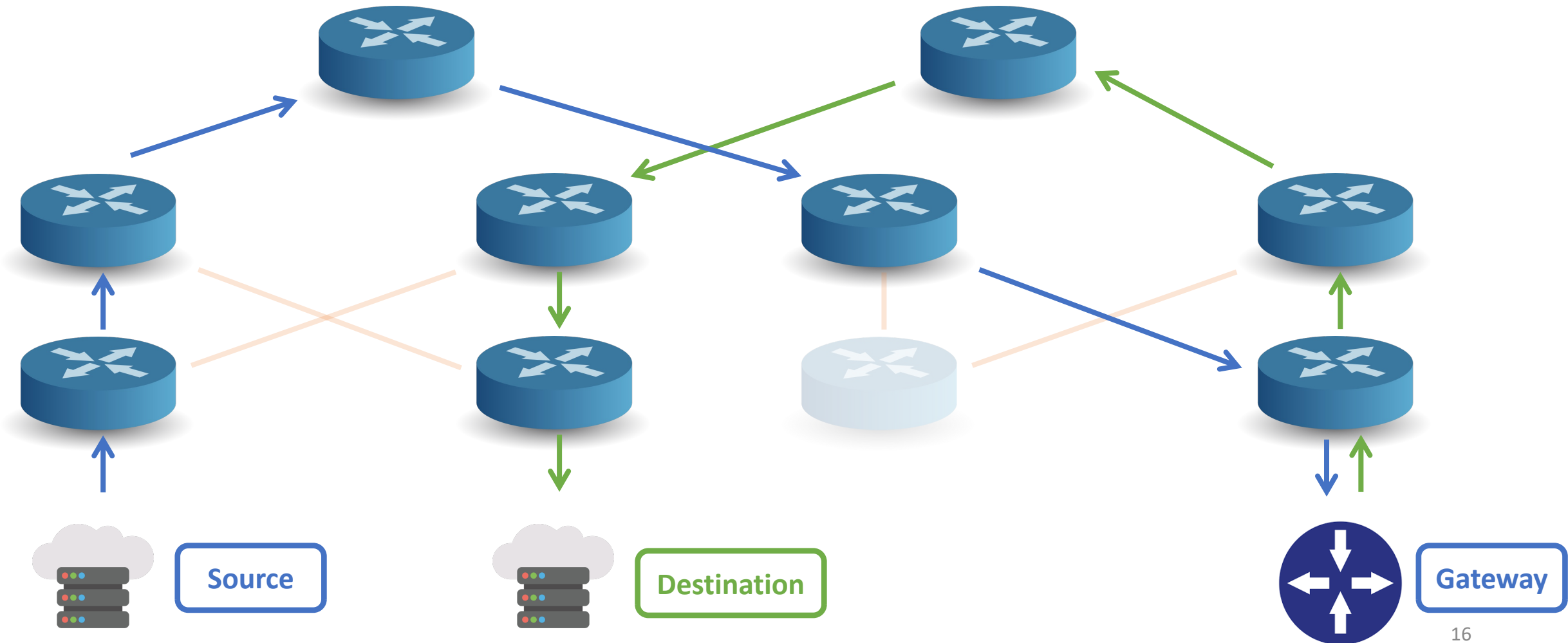
- Background
- SwitchV2P: overview
- Challenges
- Design
- Experimental results

# In-Network Address Translation



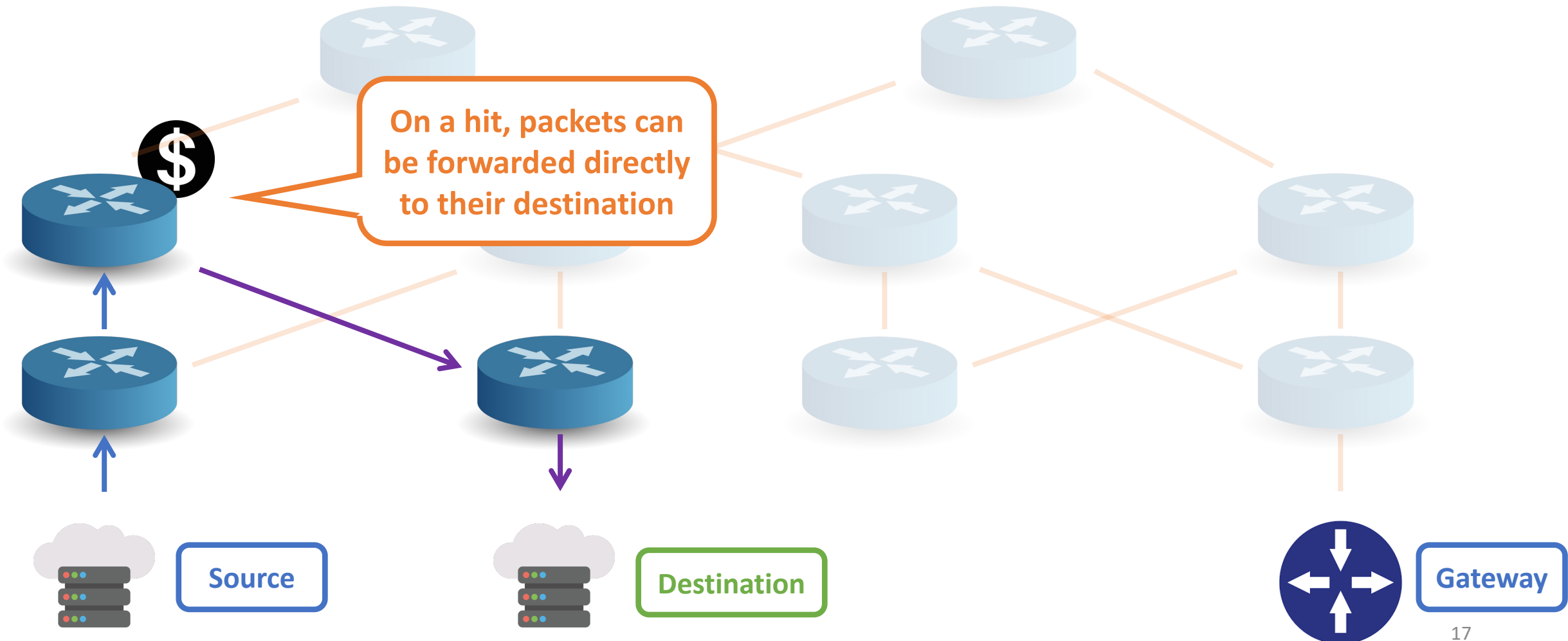
# Cache Miss

---



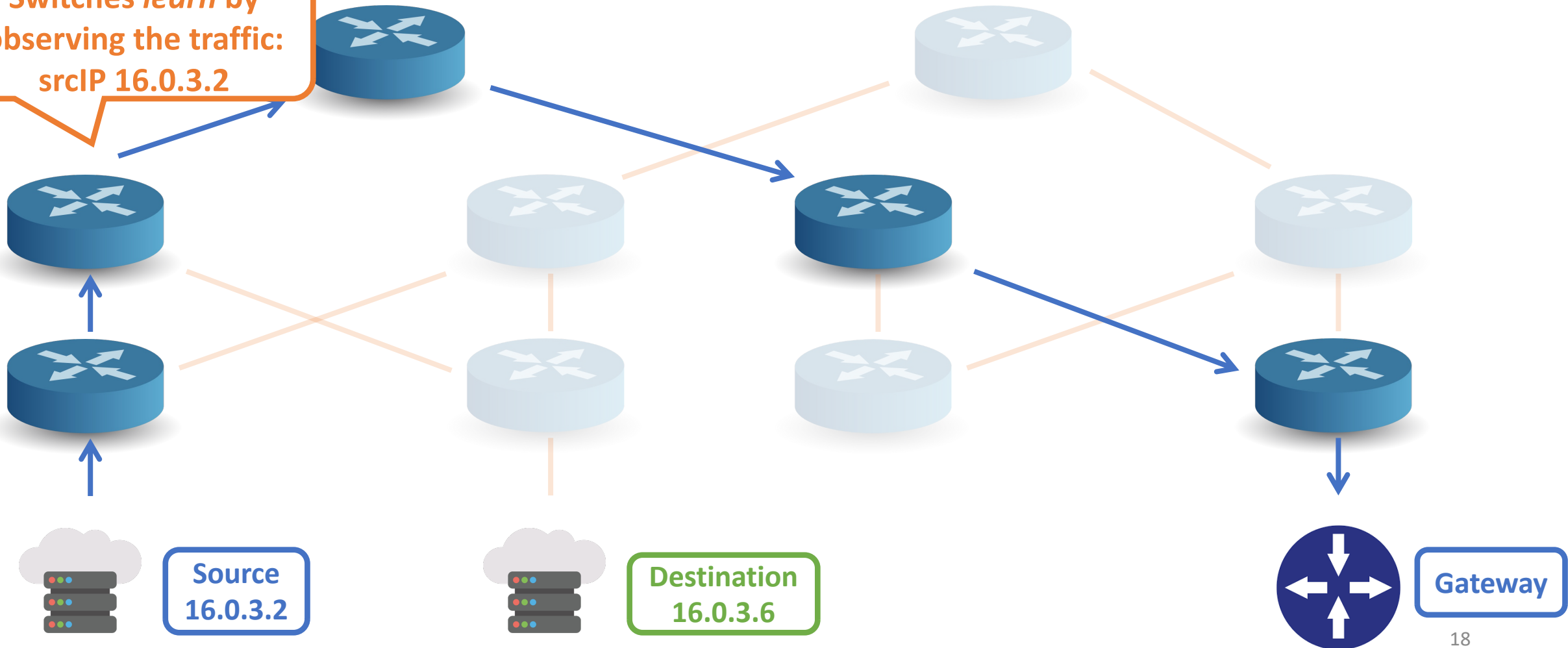


# Cache Hit: Source → Destination

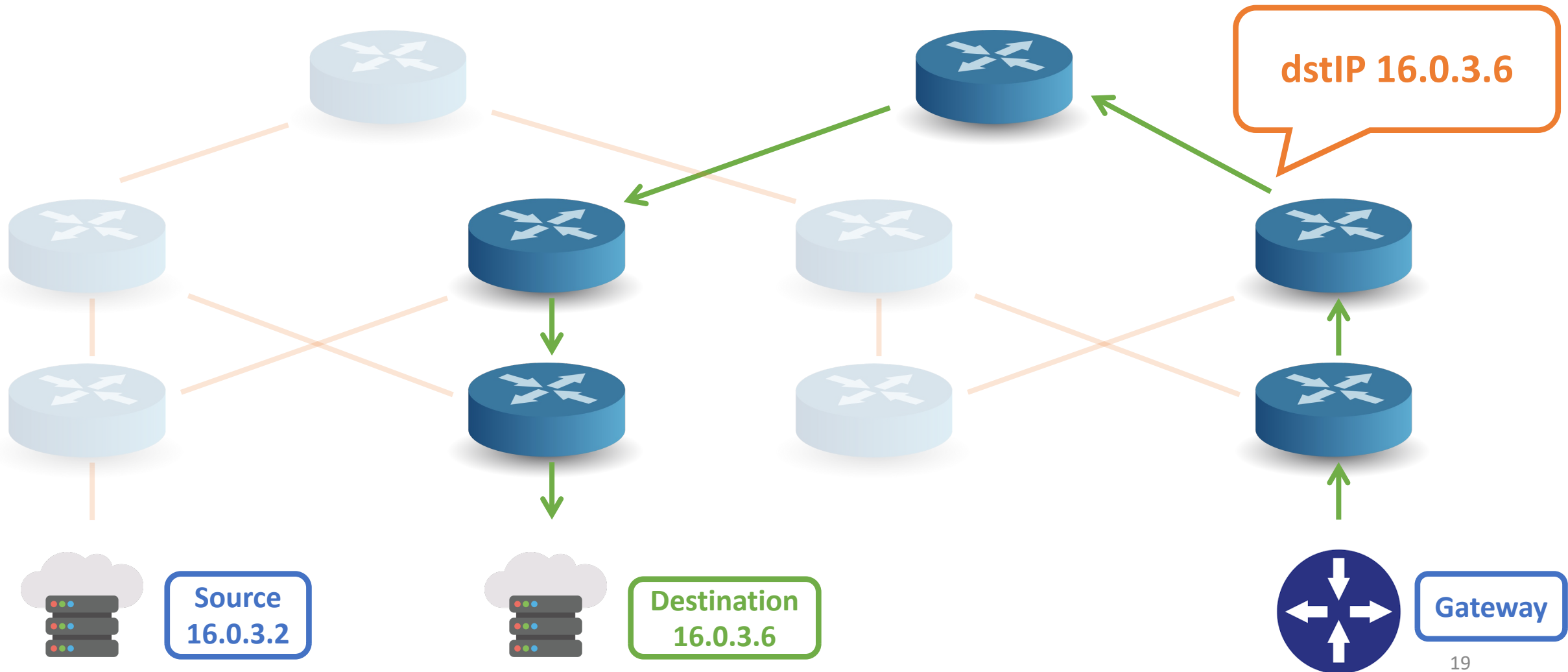


# Caching by Learning

Switches learn by observing the traffic:  
srcIP 16.0.3.2

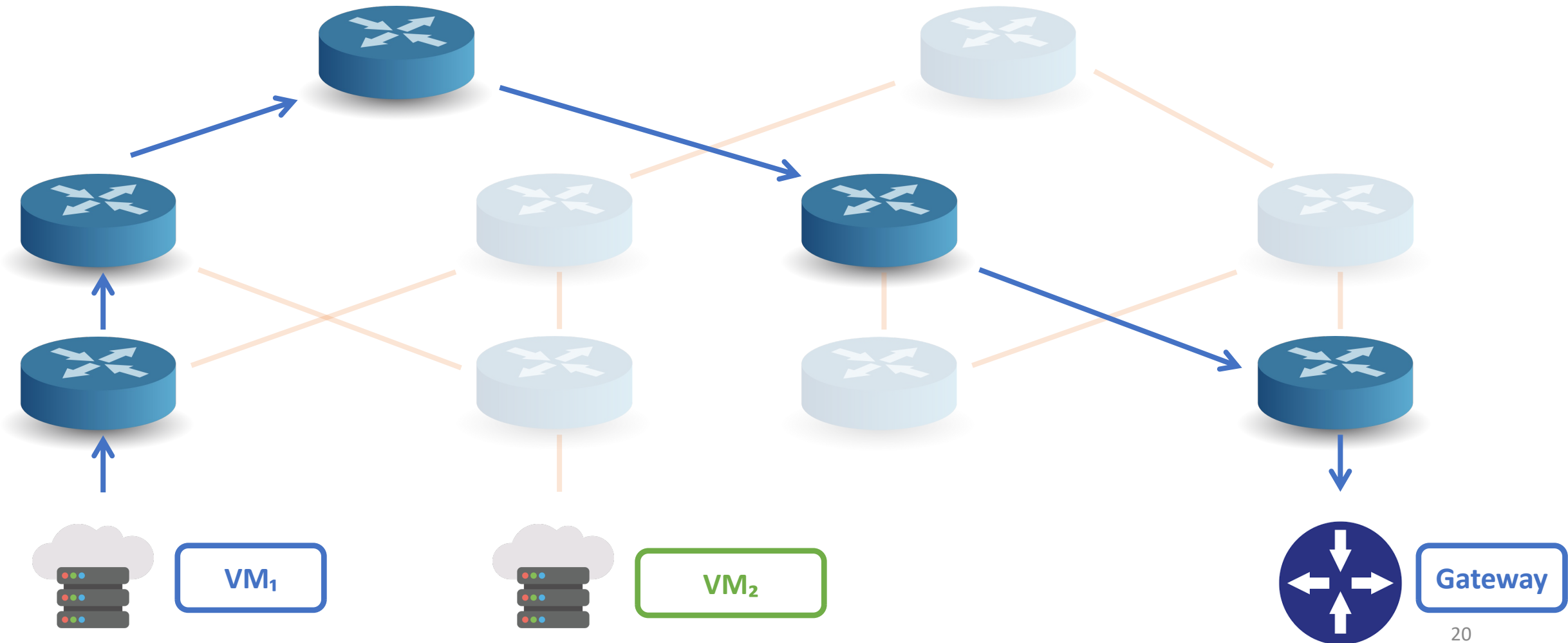


# Caching by Learning



# The Greedy Approach

---





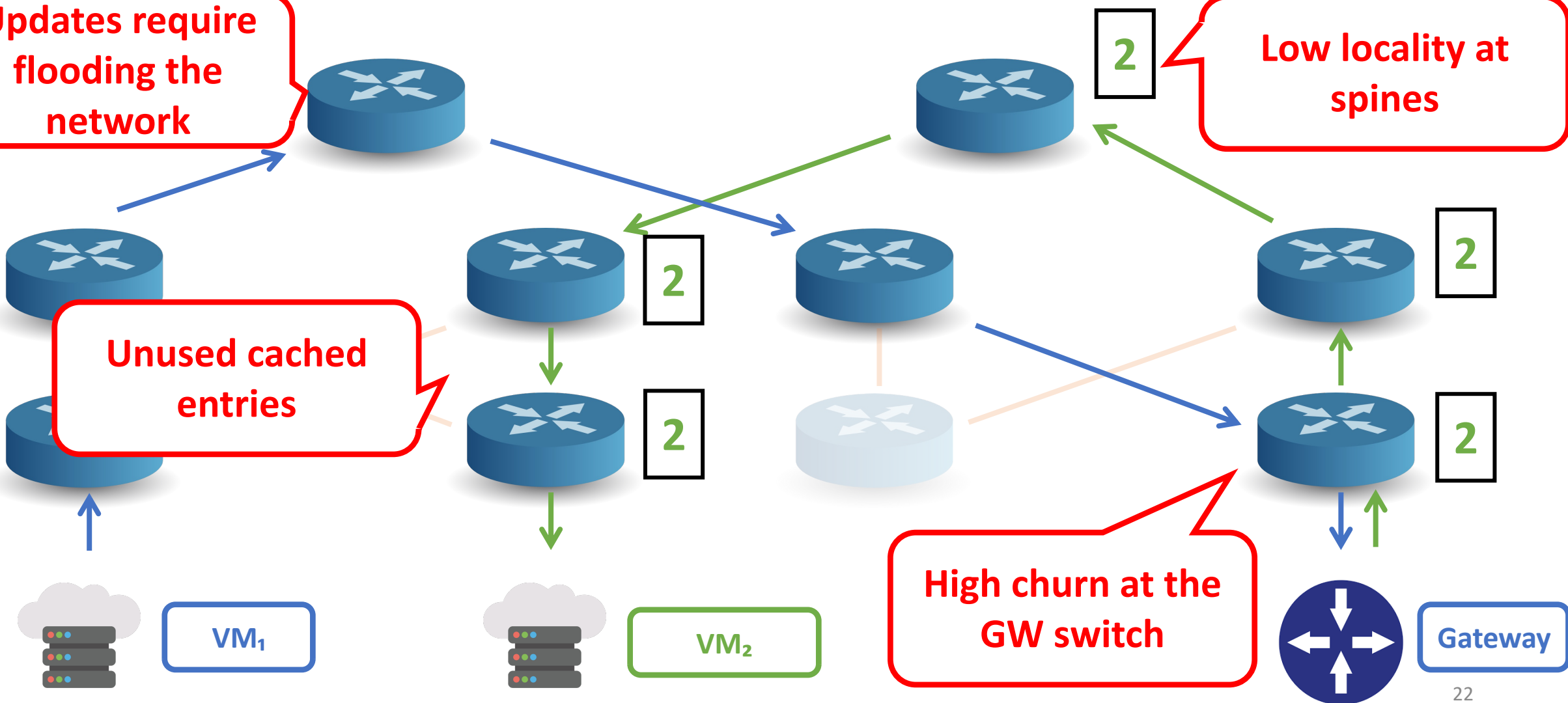
# Local Decisions are not Enough

Updates require flooding the network

Low locality at spines

Unused cached entries

High churn at the GW switch



# Topology-Aware Caching

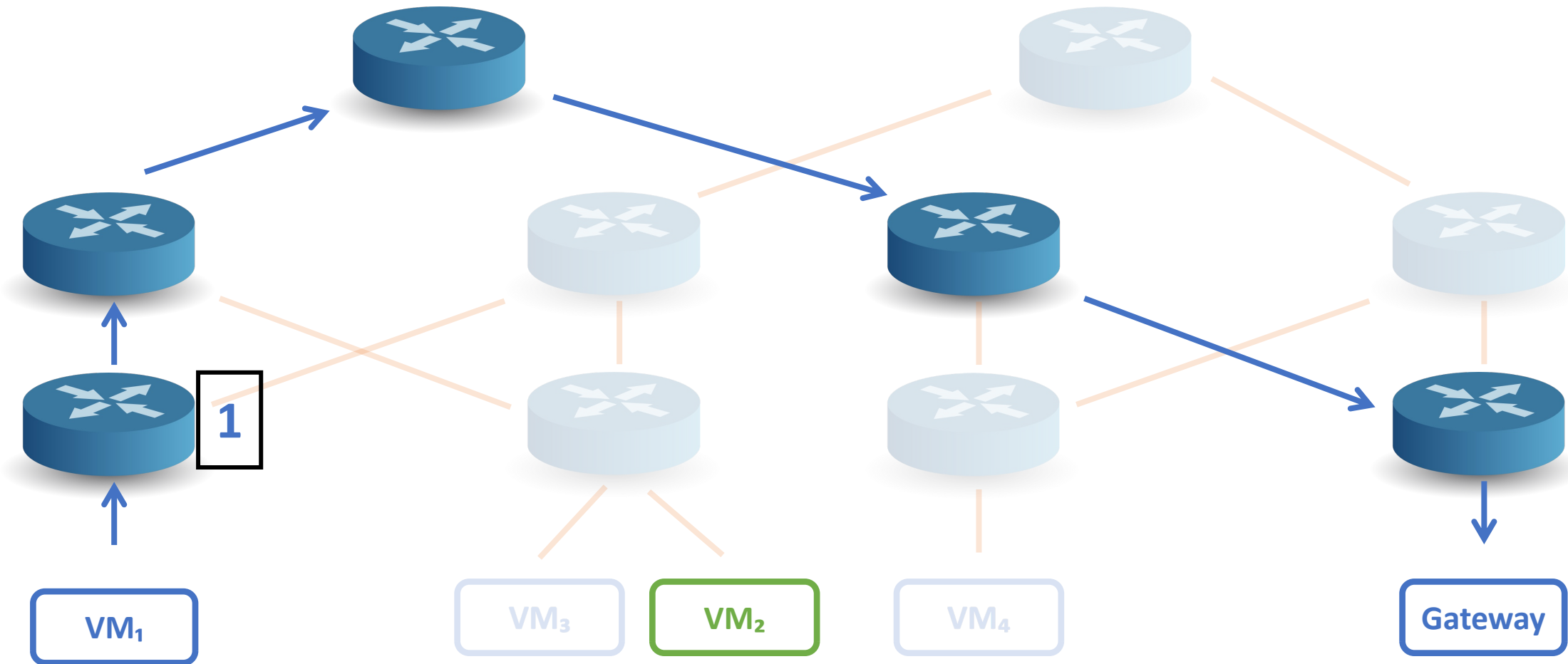
---

- Direct-mapped cache with small metadata (1 bit)
- ToRs learn source addresses
- GW ToRs learn destination addresses
- Evicted entries are spilled to other switches
- Popular entries are promoted to upper levels
- Move mappings to the traffic

**Please see the paper for more details!**

# Example: $VM_1 \rightarrow VM_2$

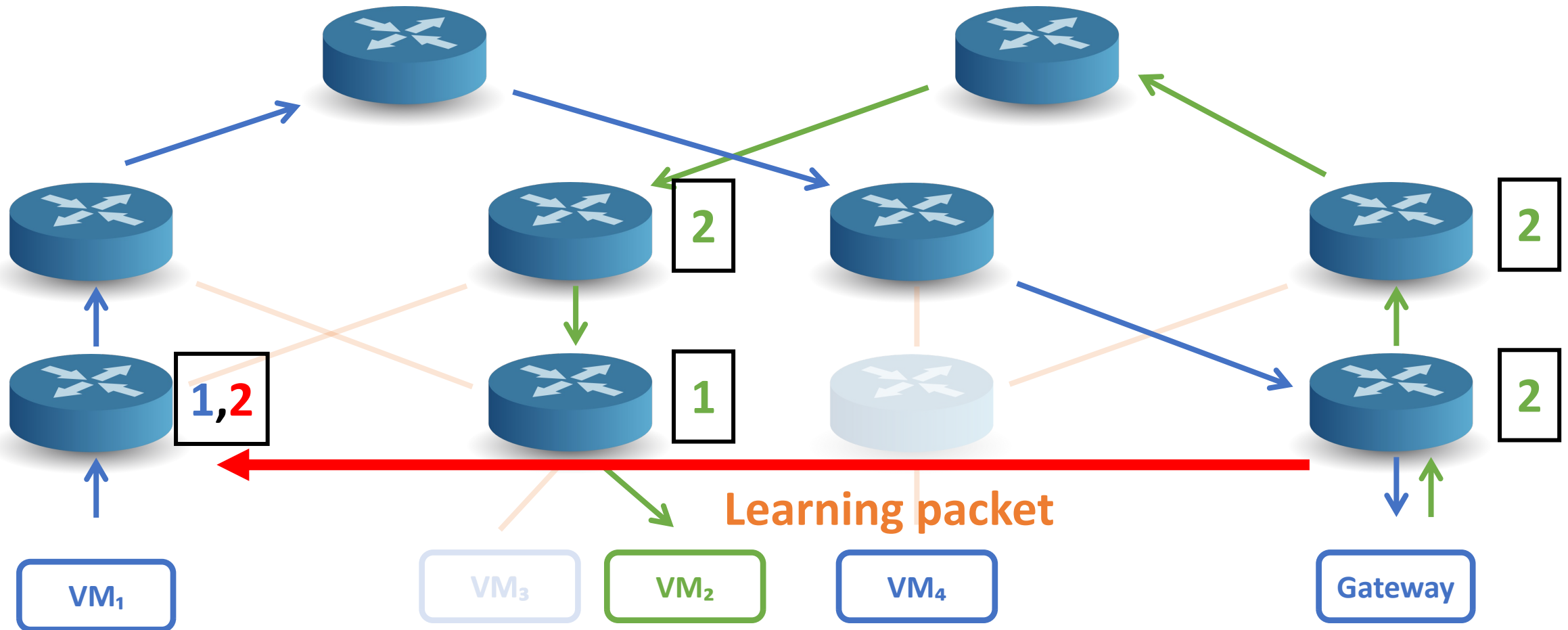
---



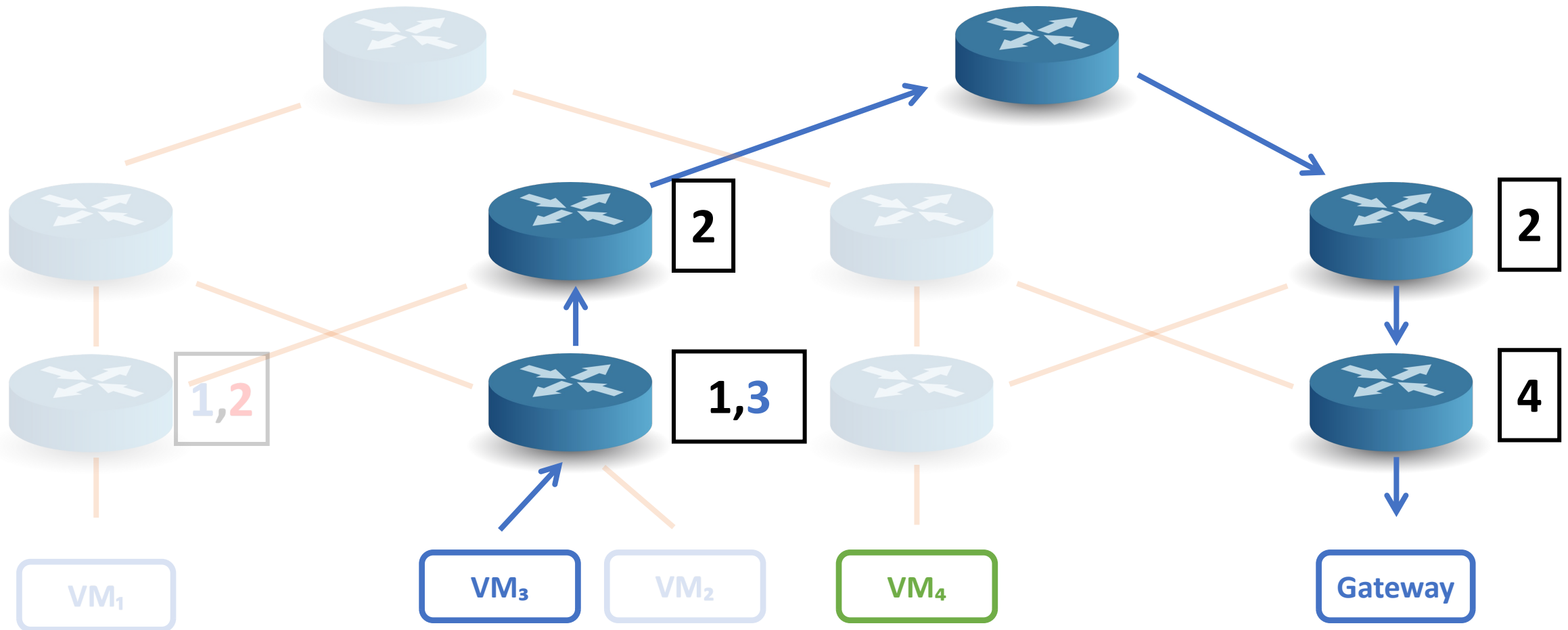




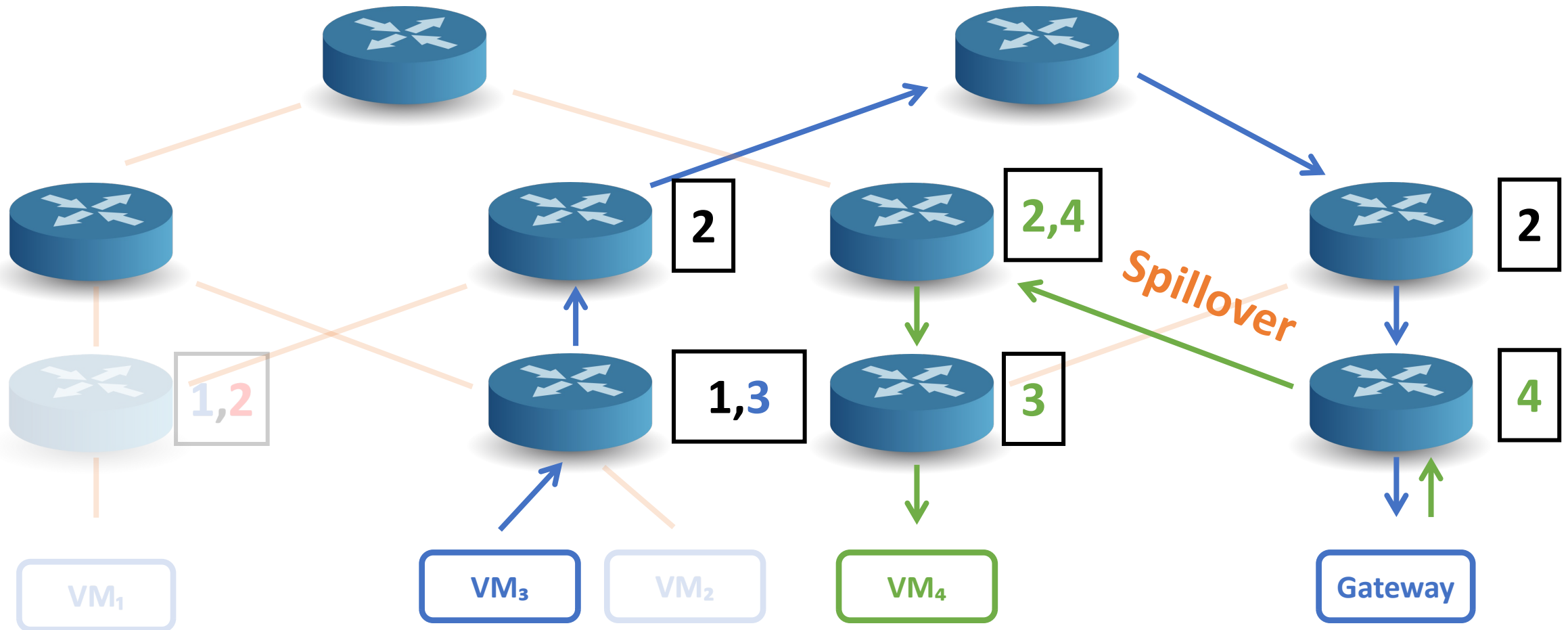
# Example: $VM_1 \rightarrow VM_2$



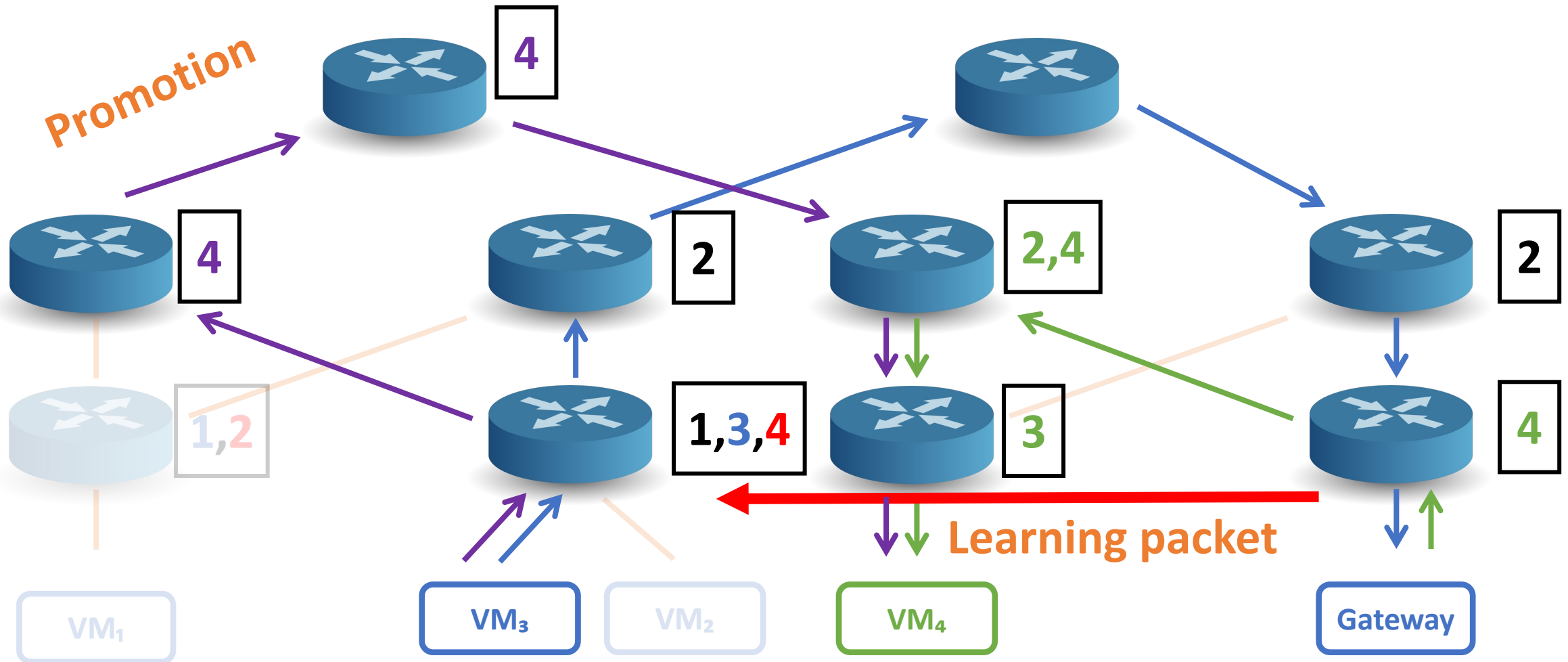
# Example: $VM_3 \rightarrow VM_4$



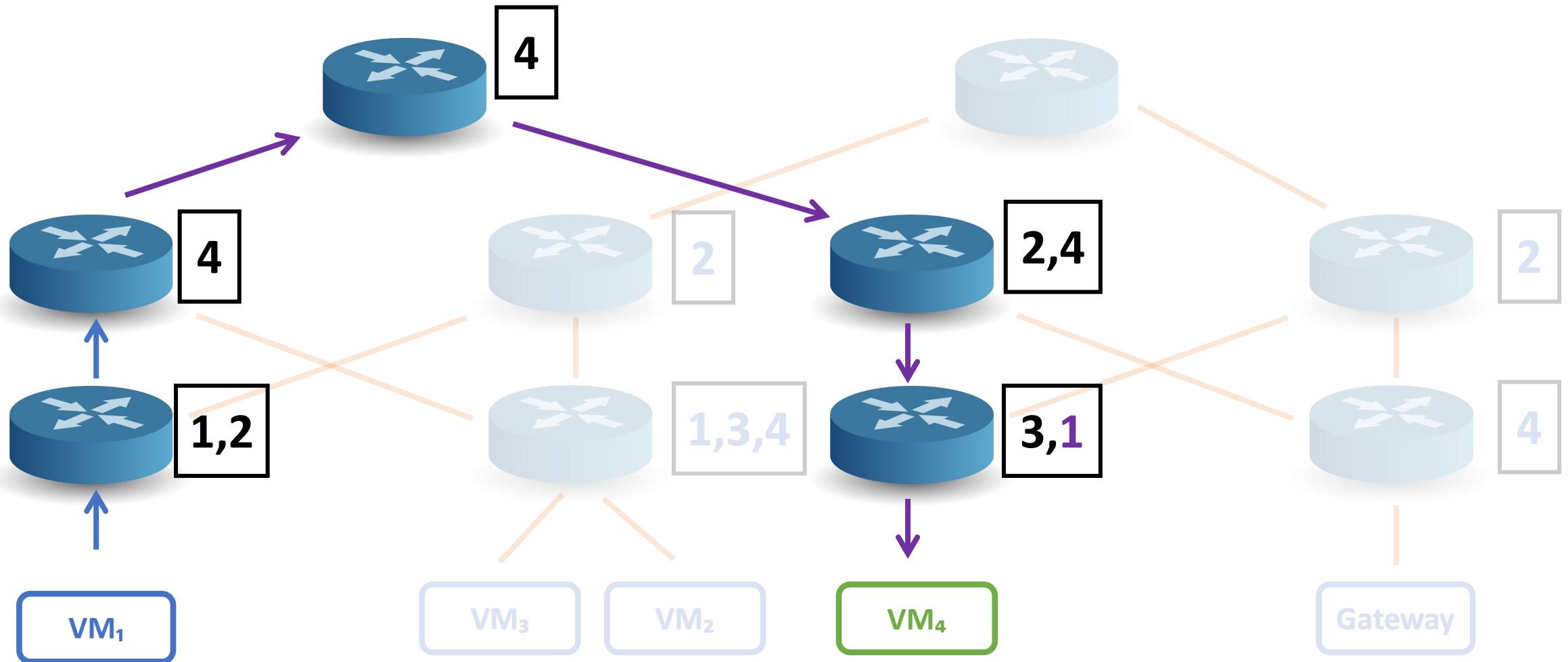
# Example: VM<sub>3</sub> → VM<sub>4</sub>



# Example: VM<sub>3</sub> → VM<sub>4</sub>

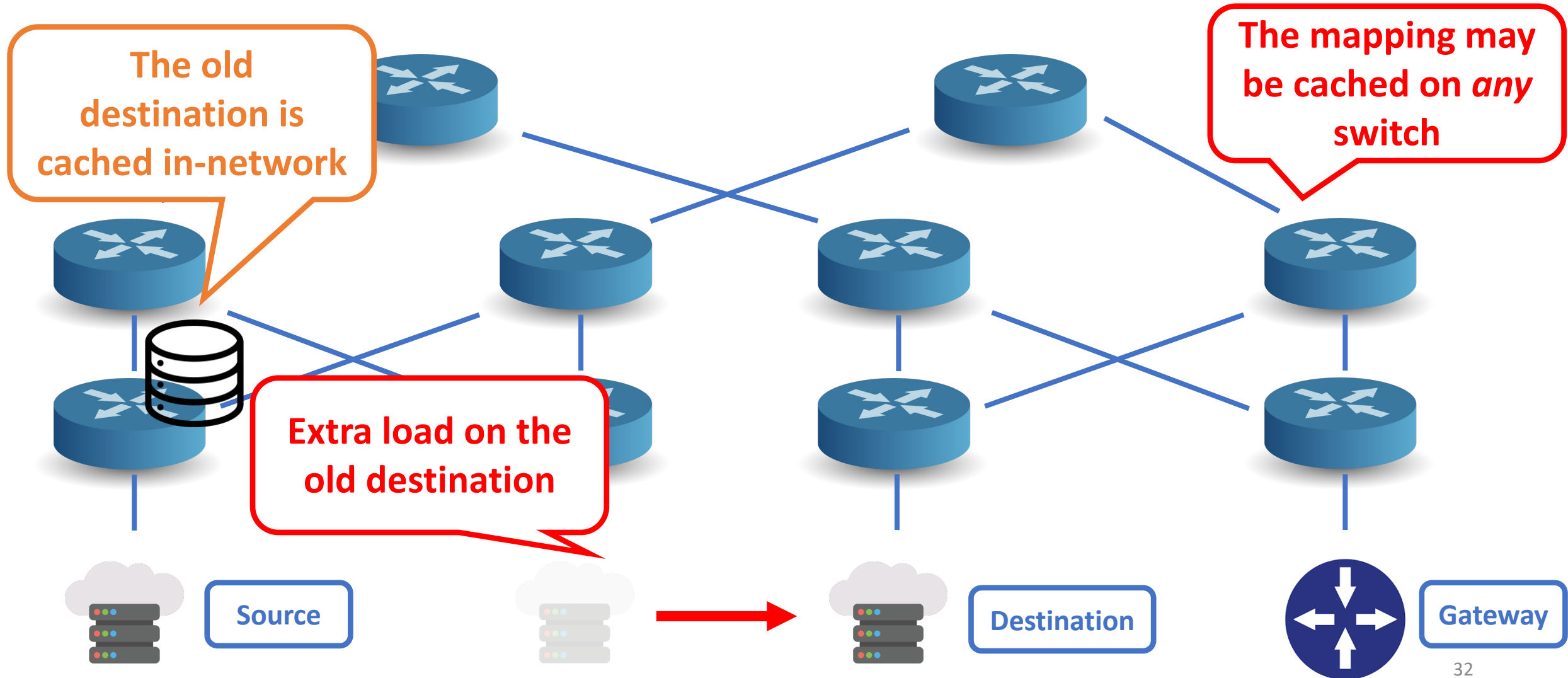


# Example: $VM_1 \rightarrow VM_4$





# Updates





# Updates

---

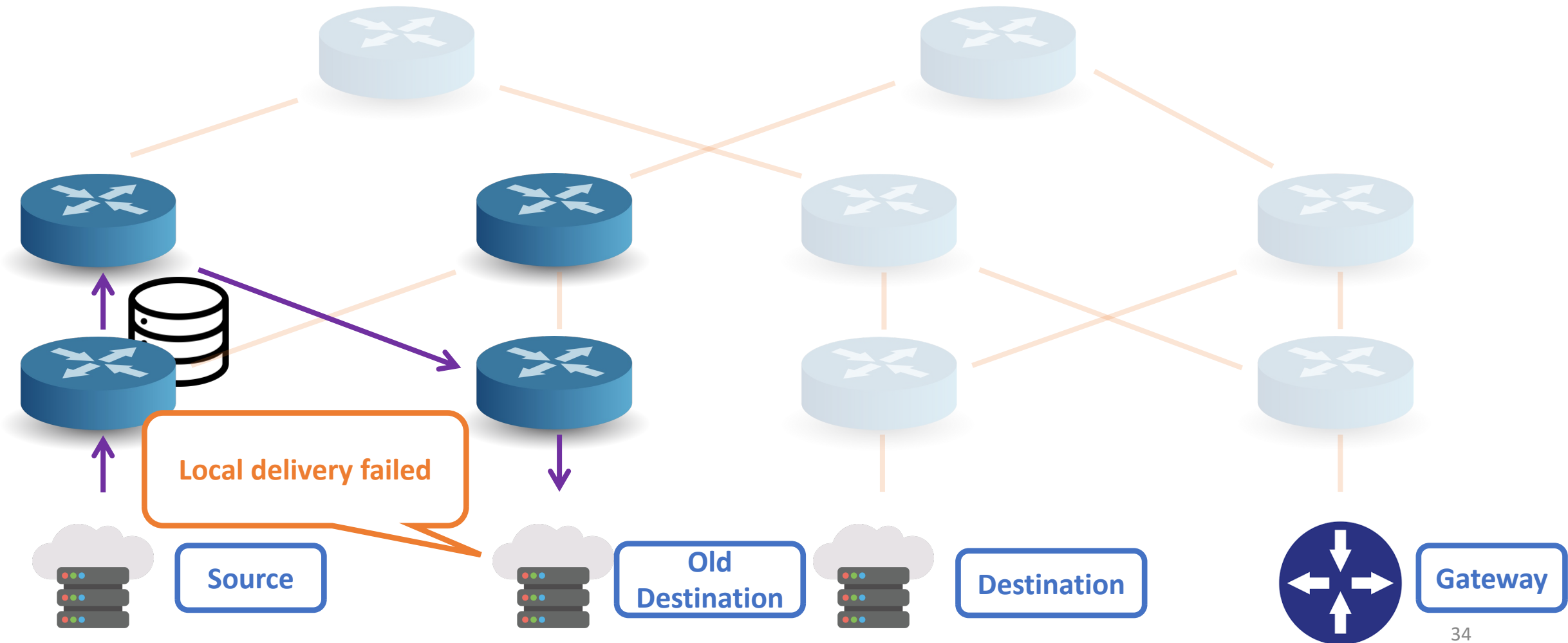
## Goal

**Minimize the number of misdelivered and  
invalidation messages**

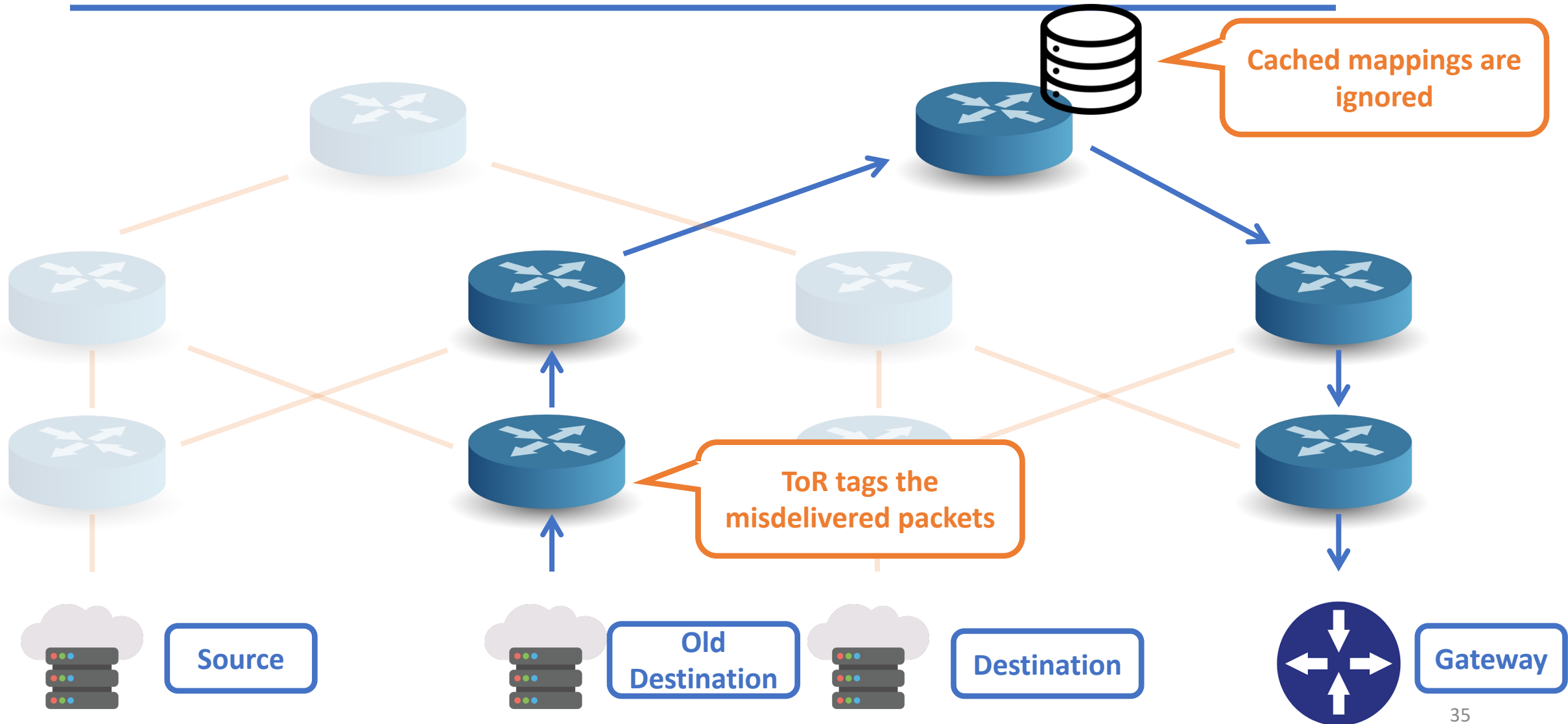
## Key idea

**Cache coherence is not necessary**

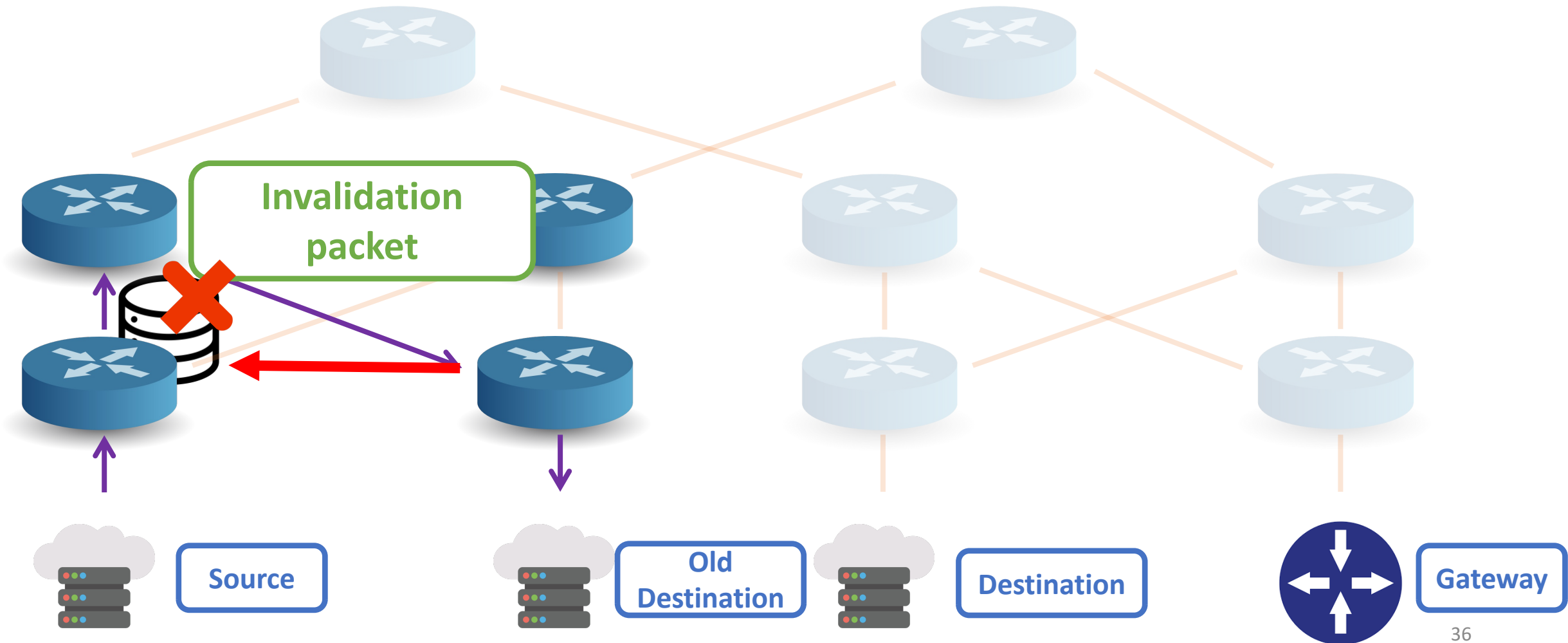
# Keeping Correct Forwarding



# Keeping Correct Forwarding



# Lazy Invalidation



# Simulations

---

- Large network topologies: 10K VMs, 128 servers, 80 switches (>800 switches for Alibaba)
- Traces: Hadoop, WebSearch, Alibaba RPC, Microbursts, Video
- Network- and application-level metrics
- Main results:
  - Up to 7.8× reduction in FCT and 4.3× reduction in first packet latency
  - Low miss rates (below 1%) - same performance with an order-of-magnitude fewer gateways
  - Reduced network load
  - Low migration costs

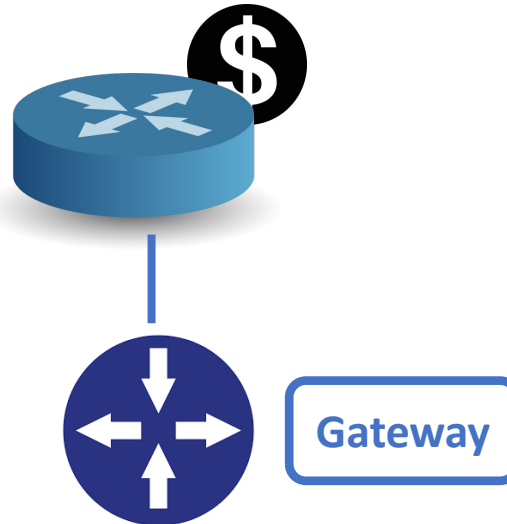
# Baselines

---

Andromeda  
(w/o offloading)



Sailfish



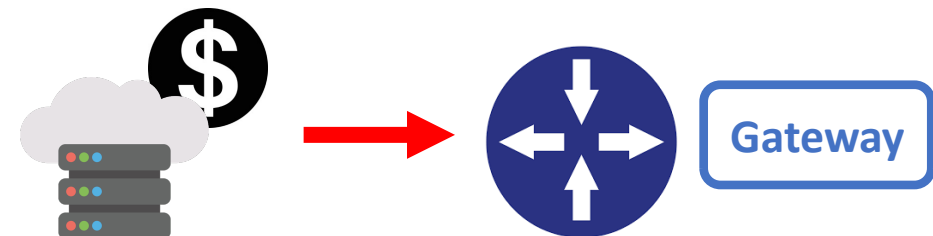
LocalLearning



AccelNet

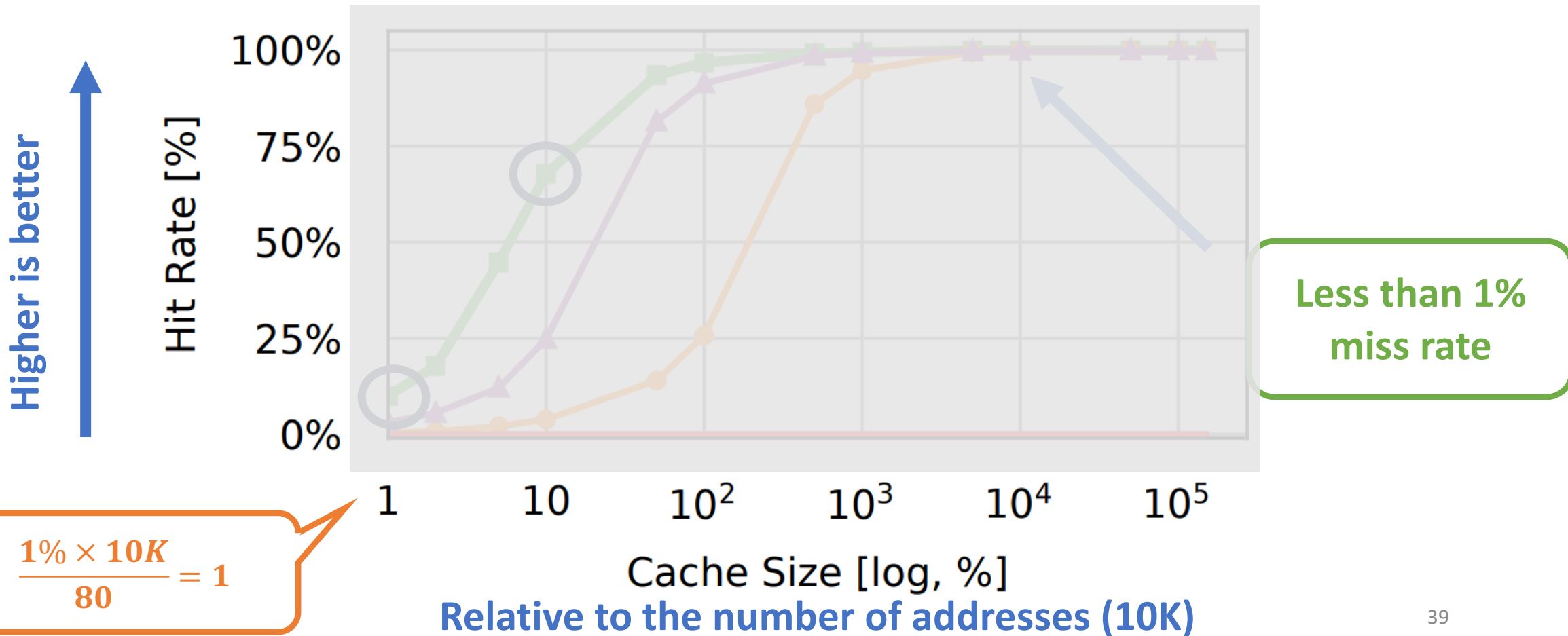


VL2



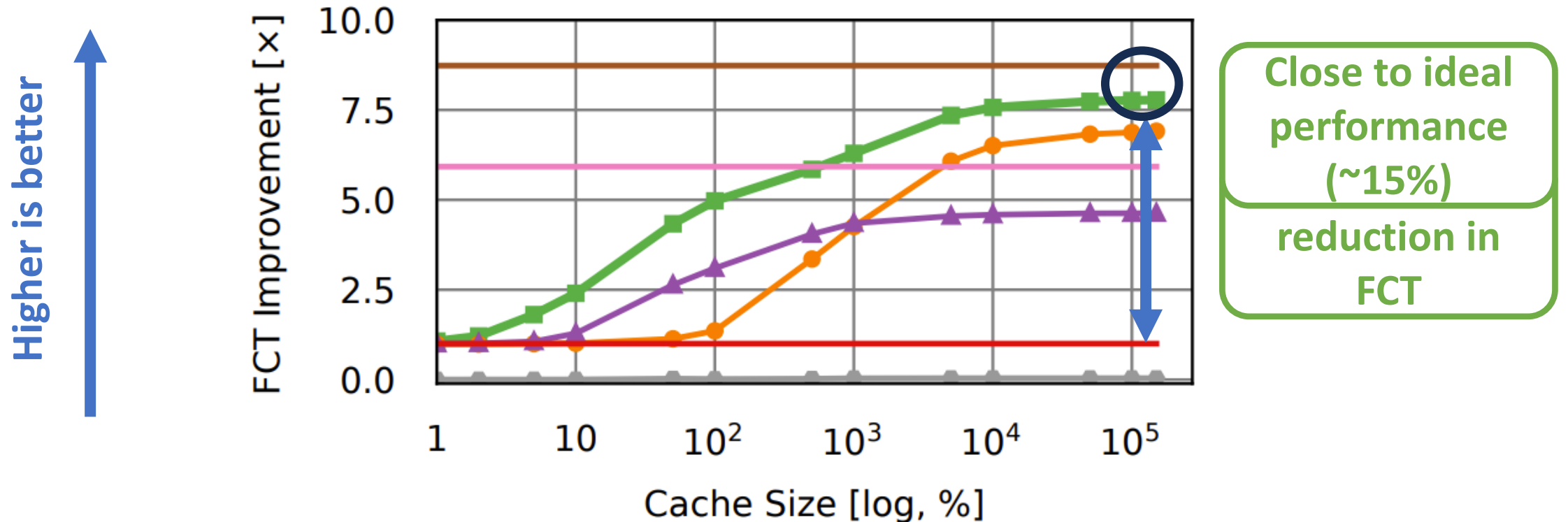
# Hadoop: Hit Rate

SwitchV2P LocalLearning Sailfish Andromeda



# Hadoop: FCT

SwitchV2P LocalLearning Sailfish Andromeda  
AccelNet VL2



Relative to the number of addresses (10K)



# Updates: Results

---

	Avg. Packet Latency	Gateway Packets	Misdelivered Packets	Total Invalidation Packets
<b>NoCache</b>	1×	100%	1×	
<b>OnDemand</b>	0.25×	0%	11×	
<b>SwitchV2P</b>	0.25×	8.7%	1.2×	24

**SwitchV2P reduces the load on the stale destination with a small number of invalidation packets**

# Conclusions

---

- Give the power to the switches!
- In-network address translation is practical and efficient
- Key ideas: topology-aware caching, move mappings to the traffic, and lazy invalidations
- Up to 7.8× reduction in FCT and 4.3× reduction in first packet latency
- Up to 6.1× reduction in bandwidth overheads

**Thank you!**

**Questions?**



**liorz@campus.technion.ac.il**