

Harp: Improving VPC Network Availability via Efficient Failure Detection and Rerouting in Tencent Cloud

Jiayu Hu^{*}, Feng Jin^{*}, Xianping Zhou^{*}, Kai Zhang[^], Zhen Shen^{*}, Yongkang Luo^{*}

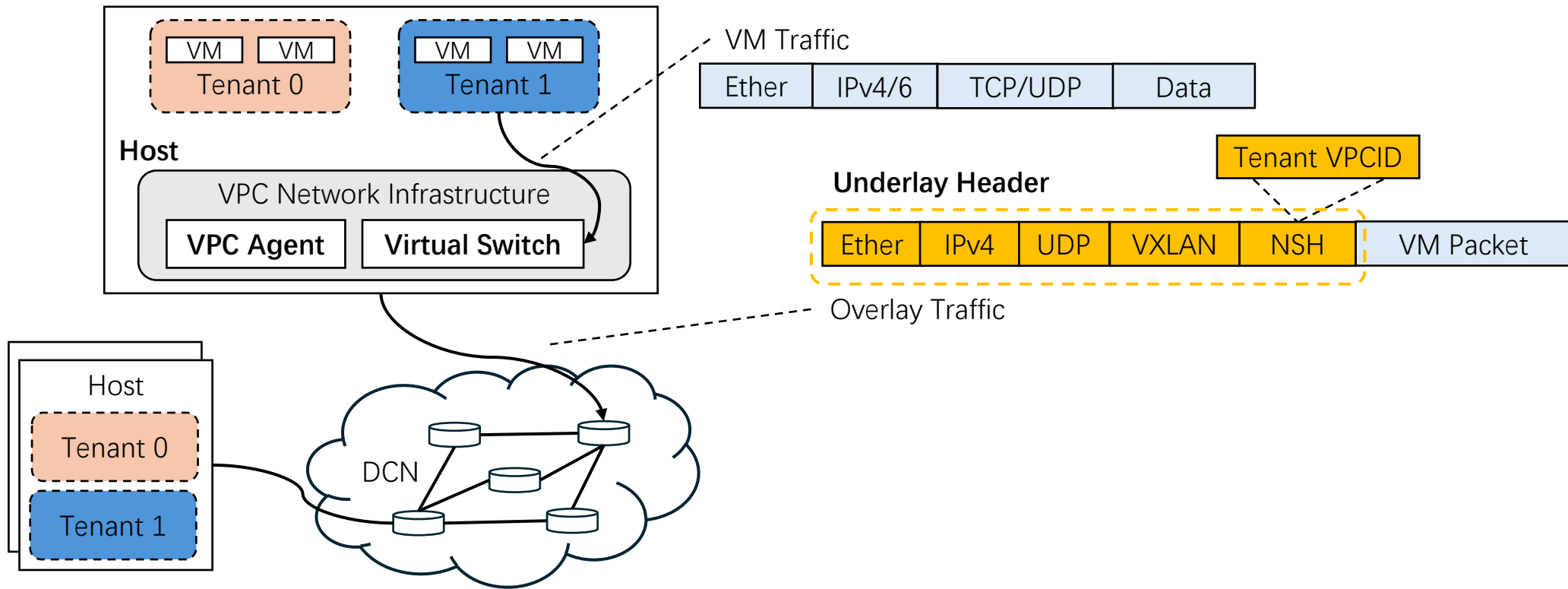
^{*}Tencent, [^]Fudan University

Tencent 腾讯



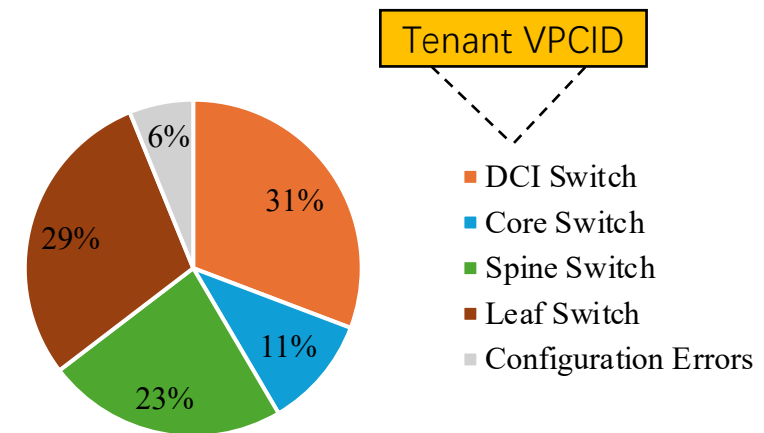
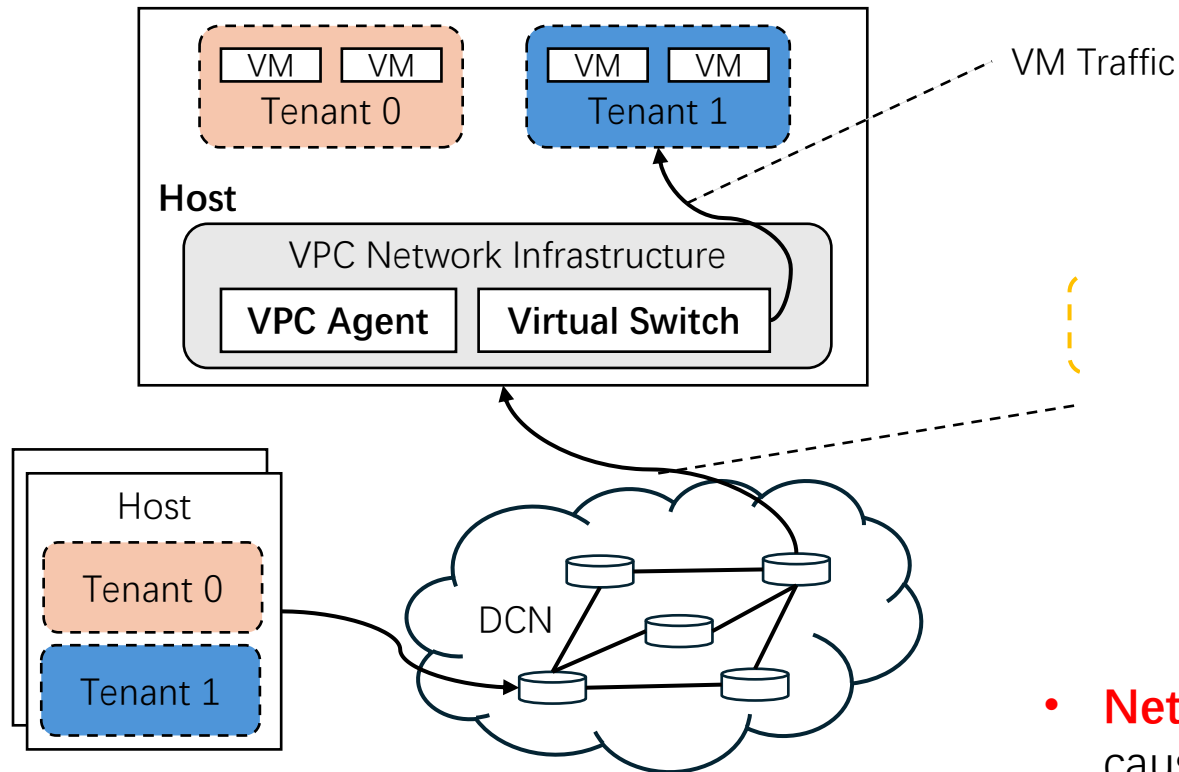
Virtual Private Cloud

- Virtual Private Cloud (VPC) is a network technique that enables tenants to customize networks in Cloud.



Virtual Private Cloud

- Virtual Private Cloud (VPC) enables tenants to customize their networks in Cloud.



Network Failure Distribution in Tencent Cloud

- Network failures** inevitably occur in data centers, causing **VPC packet drops** and even **breaking SLAs**.

Improving VPC Network Availability

1. Low Failure Recovery Time

- There are many **real-time applications** in Tencent.
- For example, Redis considers an operation failed if replies are not returned **within 200ms**.
- Localization-then-mitigation approaches, such as Pingmesh 😞

Harp's solution

In-Band Failure
Detection Mechanism

2. Deploy on Legacy Network Devices

- Many network devices in Tencent data centers only **support IPv4**.
- IPv6 based approaches, such as PRR 😞

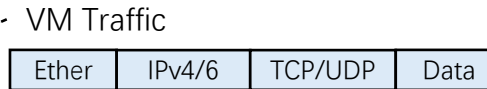
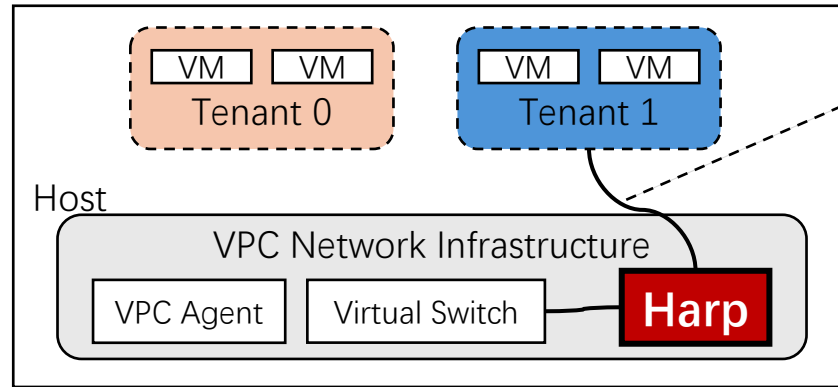
ECMP-based Path
Control Approach

3. Compatible with Existing Software

- Cloud applications are based on various **protocols**, which **cannot be changed by cloud providers**.
- Transport layer approaches, such as MPTCP 😞

VPC Network
Infrastructure Layer
Solution

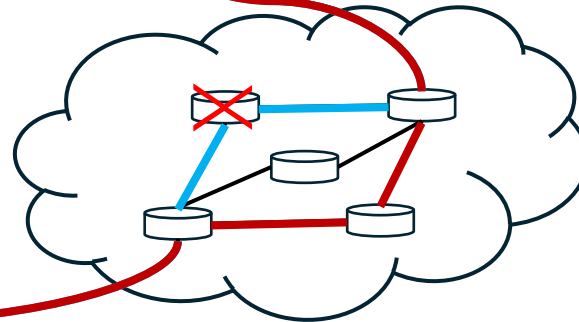
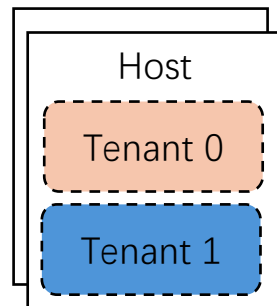
Harp: System Overview



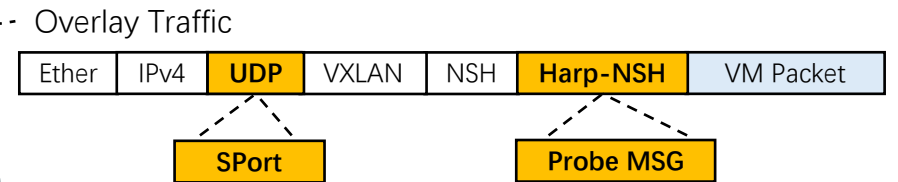
Monitor the health of physical paths

Work for all VM applications

change the ECMP selection



Repath by modifying UDP SRC ports



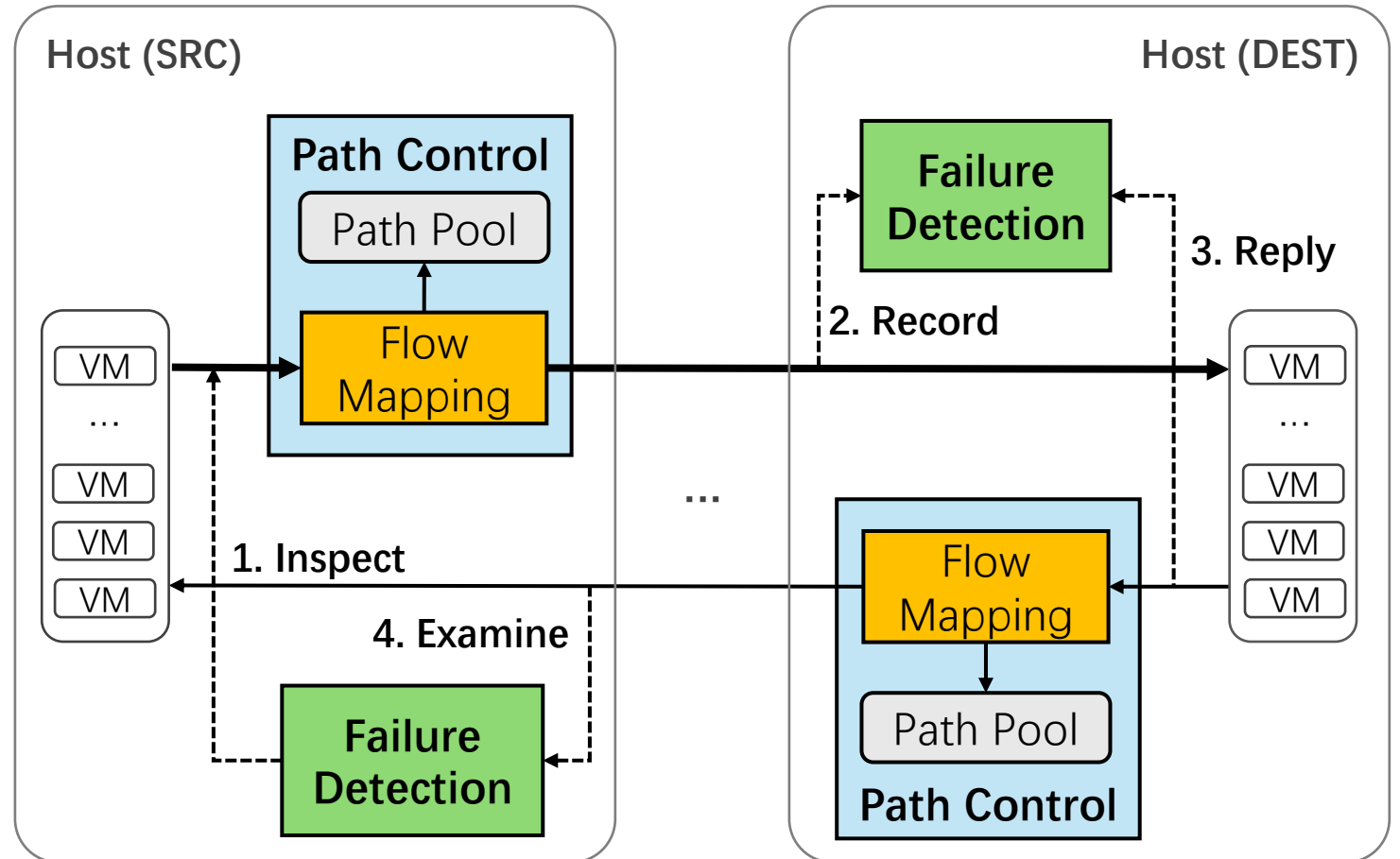
Harp: System Overview

- **Path Control**

- Build a physical path pool for each communicating host pair
- Map flows to paths
- Repath after failure

- **Failure Detection**

- Insert inspection messages for used paths
- Record active paths and return that information
- Examine the health of used paths

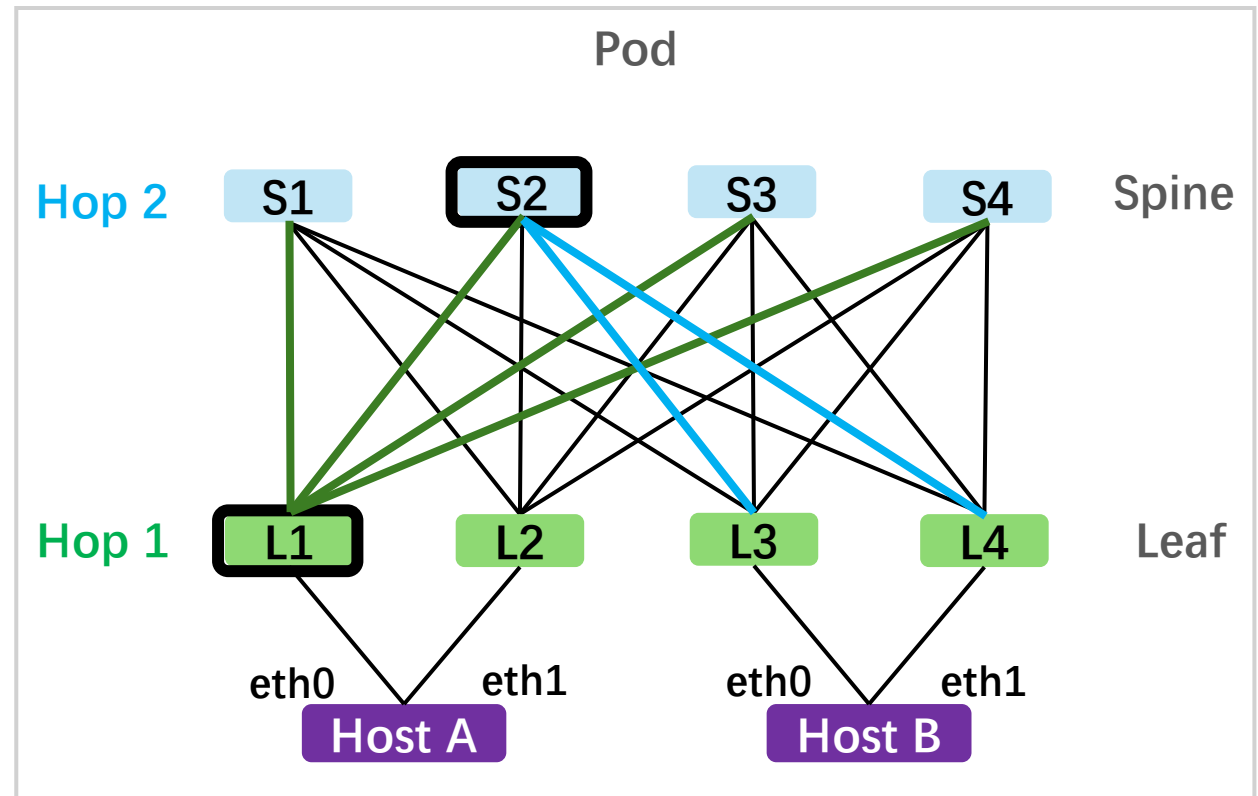


ECMP-based Path Control

K	SPORT Group	rPath	SPORT collection
0	11,27	0	11
1	9,65	1	27
2	78,79	2	9
3	98,99	3	65
		4	78
		5	79
		6	98
		7	99

K	SPORT Group
0	11,9,78,98
1	27,65,79,99

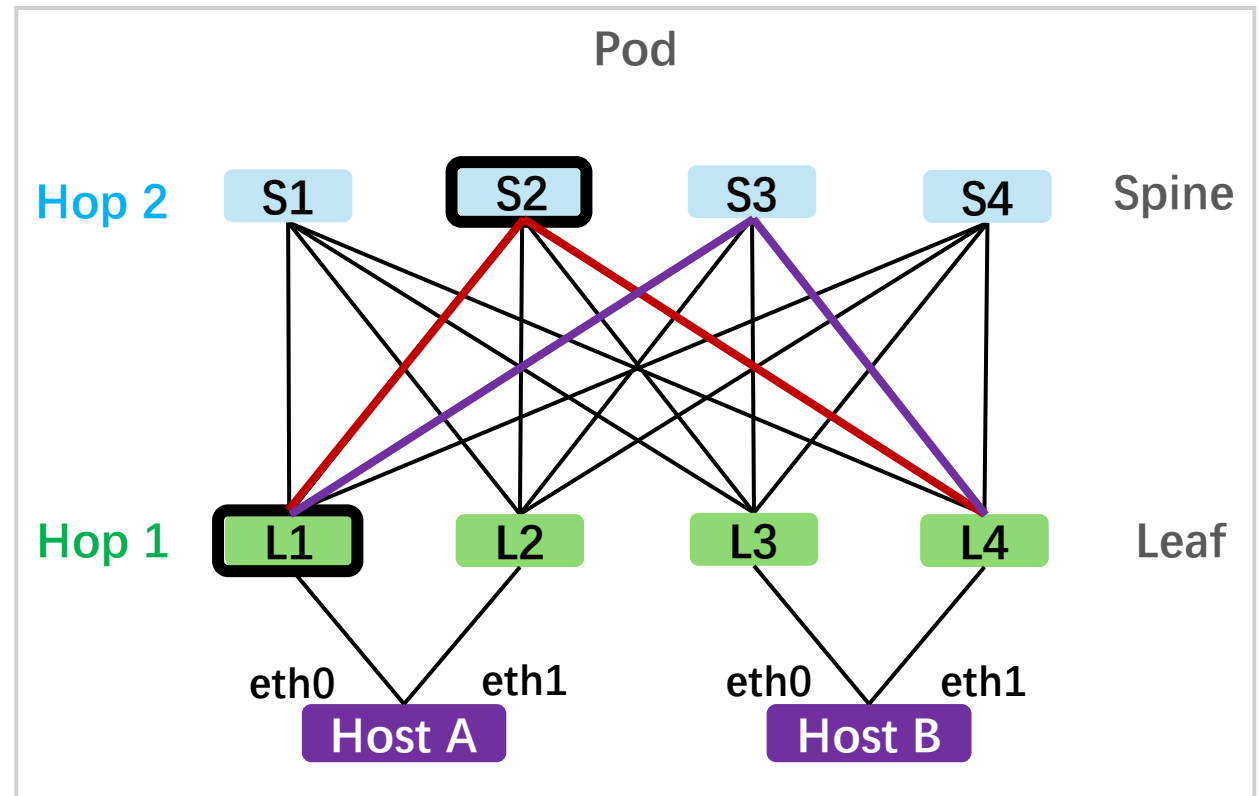
In the 5-tuple, only the SPORT changes for traffic between Host A and Host B.



- ECMP(SIP, DIP, **SPORT**, DPORT, Proto) = **K**
- SPORT Groups of Hop 1 \cap SPORT Groups of Hop 2 \Rightarrow 8 SPORT collections (8 *rPaths*)

ECMP-based Path Control

K	SPORT Group	rPath	SPORT collection
0	11,27	0	11
1	9,65	1	27
2	78,79	2	9
3	98,99	3	65
		4	78
K	SPORT Group	5	79
0	11,9,78,98	6	98
1	27,65,79,99	7	99

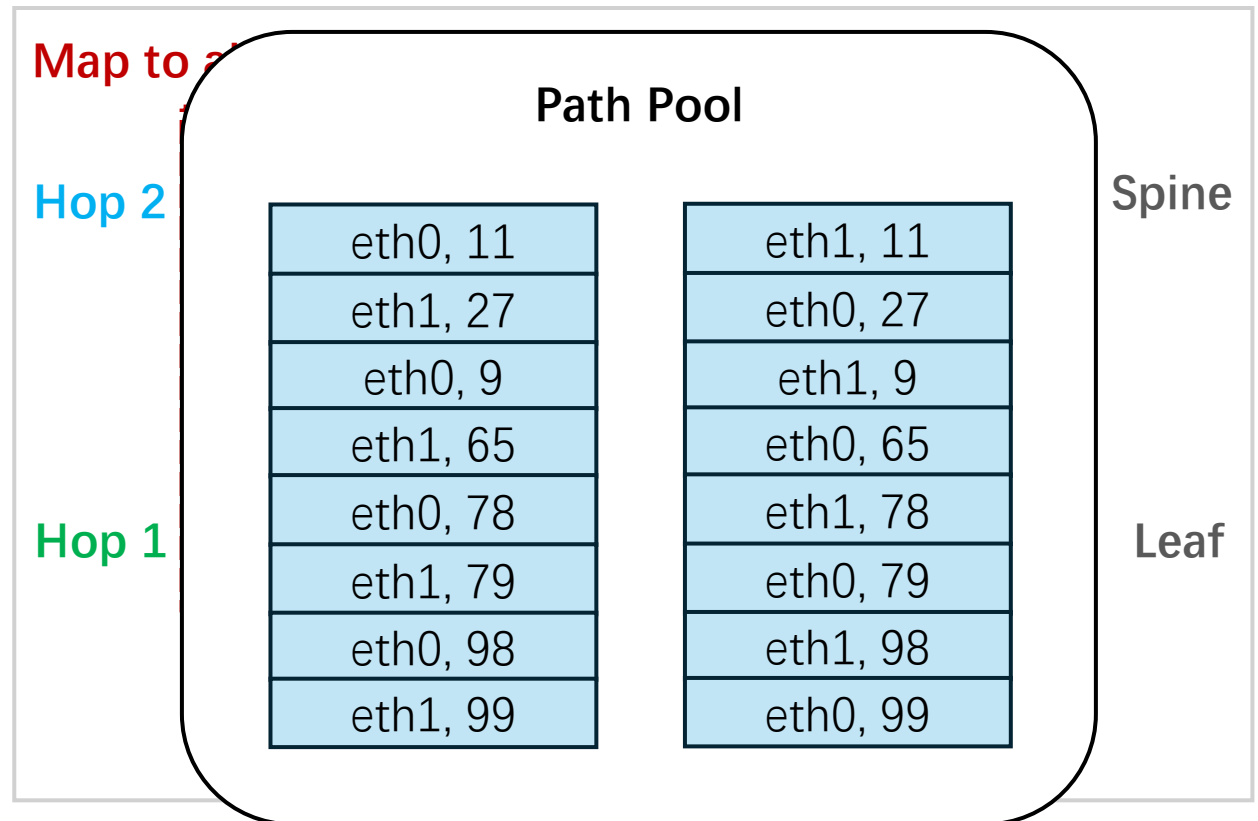


- SPORTs from **different rPaths** lead flows to **distinct paths**.

- ECMP(SIP, DIP, **SPORT**, DPORT, Proto) = **K**
- SPORT Groups of Hop 1 \cap SPORT Groups of Hop 2
=> 8 SPORT collections (8 *rPaths*)

ECMP-based Path Control

K	SPORT Group	rPath	SPORT collection
0	11,27	0	11
1	9,65	1	27
2	78,79	2	9
3	98,99	3	65
		4	78
K	SPORT Group	5	79
0	11,9,78,98	6	98
1	27,65,79,99	7	99



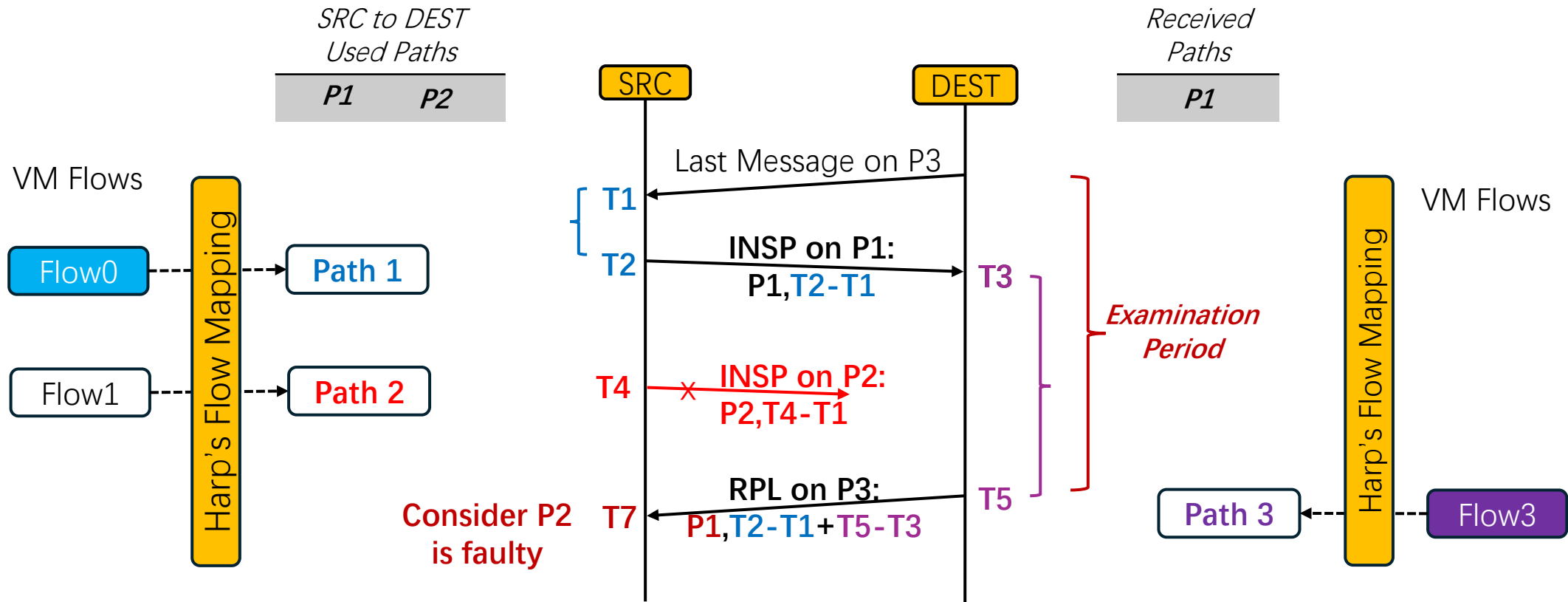
- SPORTs from different rPaths lead flows to distinct paths.
- **(NIC interface, SPORT from rPath)** controls a physical path.

- ECMP(SIP, DIP, **SPORT**, DPORT, Proto)=K
- SPORT Groups of Hop 1 \cap SPORT Groups of Hop 2 \Rightarrow 8 SPORT collections (8 rPaths)

In-Band Failure Detection

Inspection Period — (blue line)
Reply Processing Delay — (purple line)

1. Send an inspection message (**INSP**) for each used path.
2. Send a reply message (**RPL**) for all received paths during the past period.
3. A path is **faulty** if it is **used** but **not replied**.

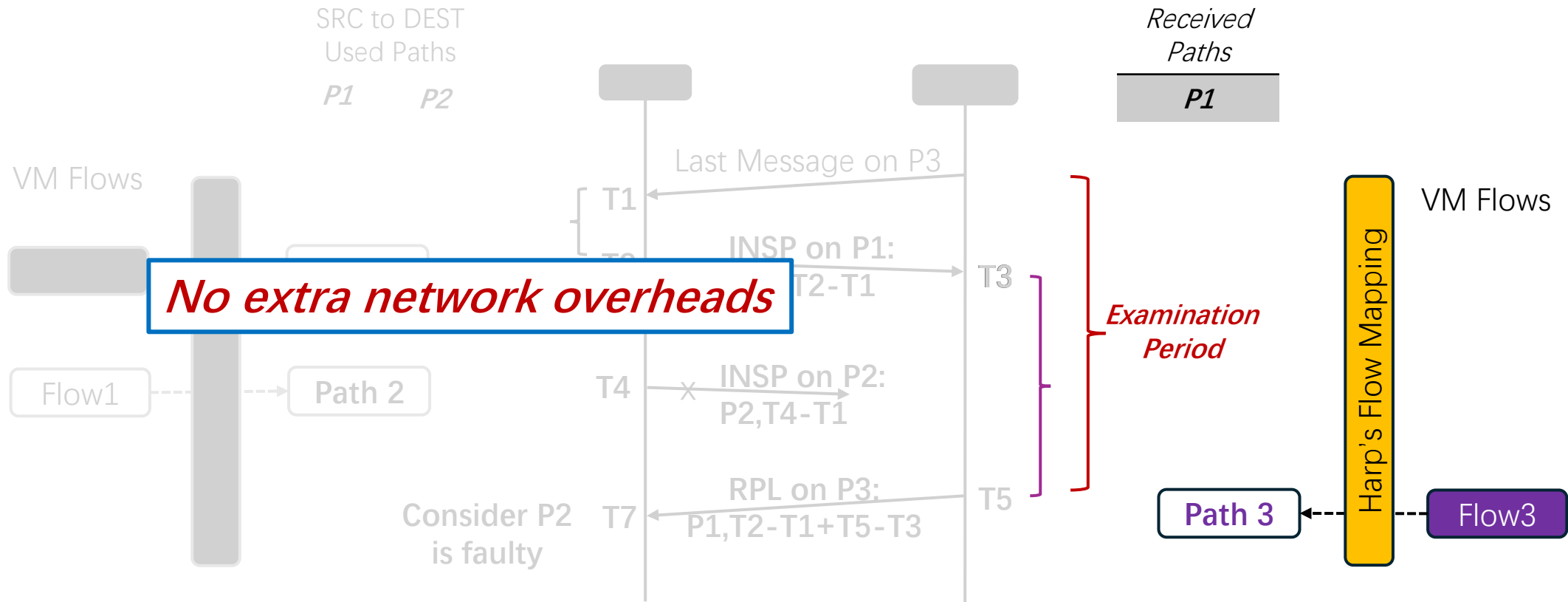


In-Band Failure Detection

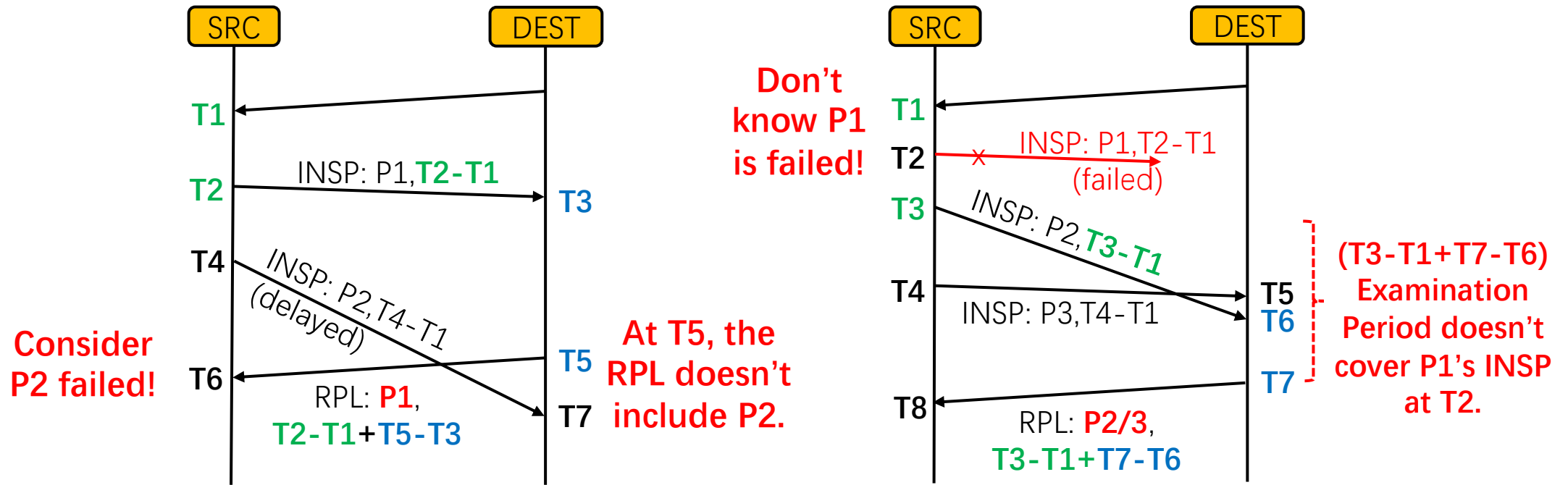
Inspection Period — (blue line)
Reply Processing Delay — (purple line)

1. Send an inspection message (**INSP**) for each used path.
3. A path is **faulty** if it is **used** but **not replied**.

2. Send a reply message (**RPL**) for all received paths during the past period.



Improving Fault Tolerance

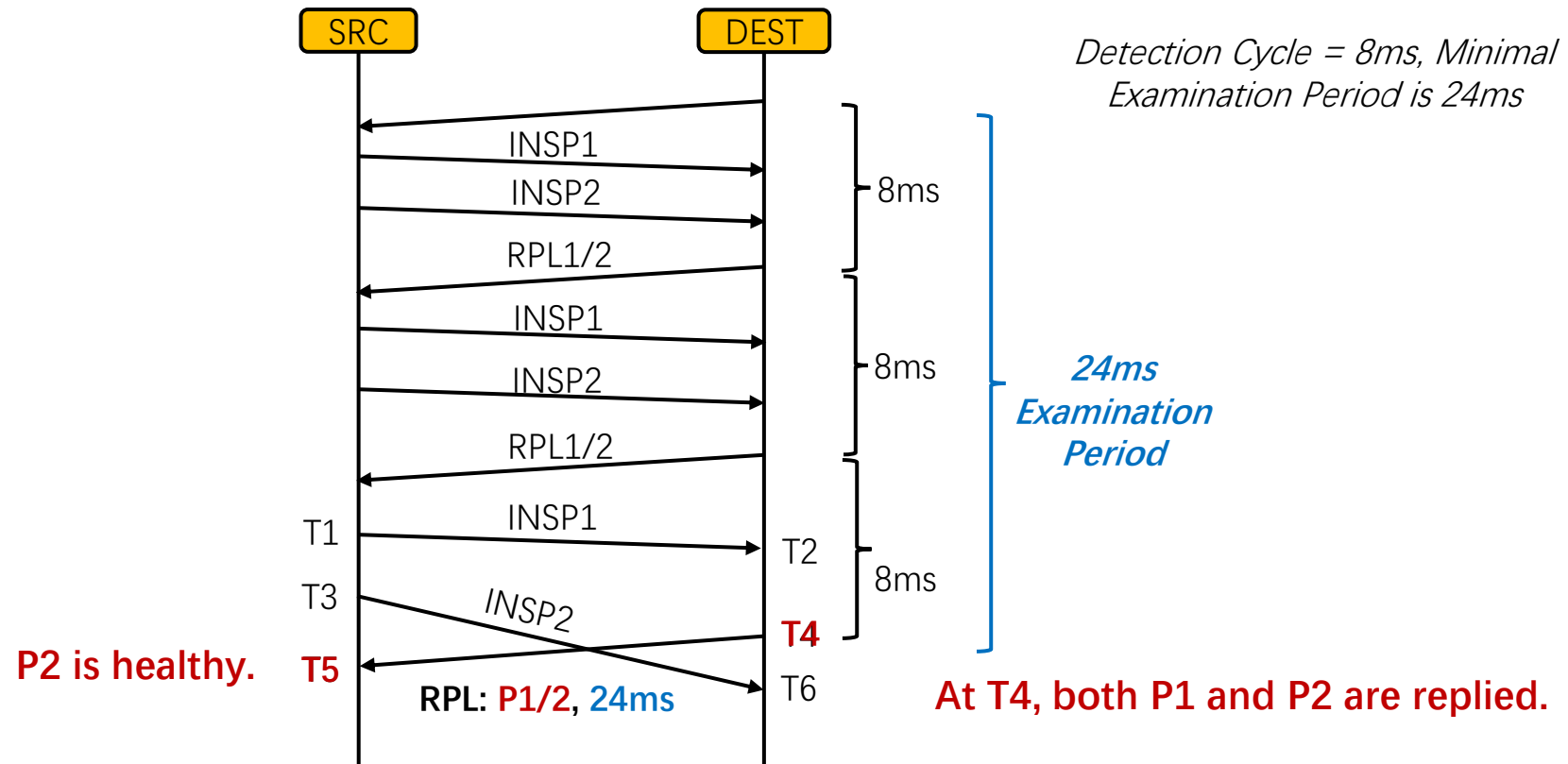


Case 1: Misjudging a healthy path as faulty

Case 2: Failing to find a faulty Path

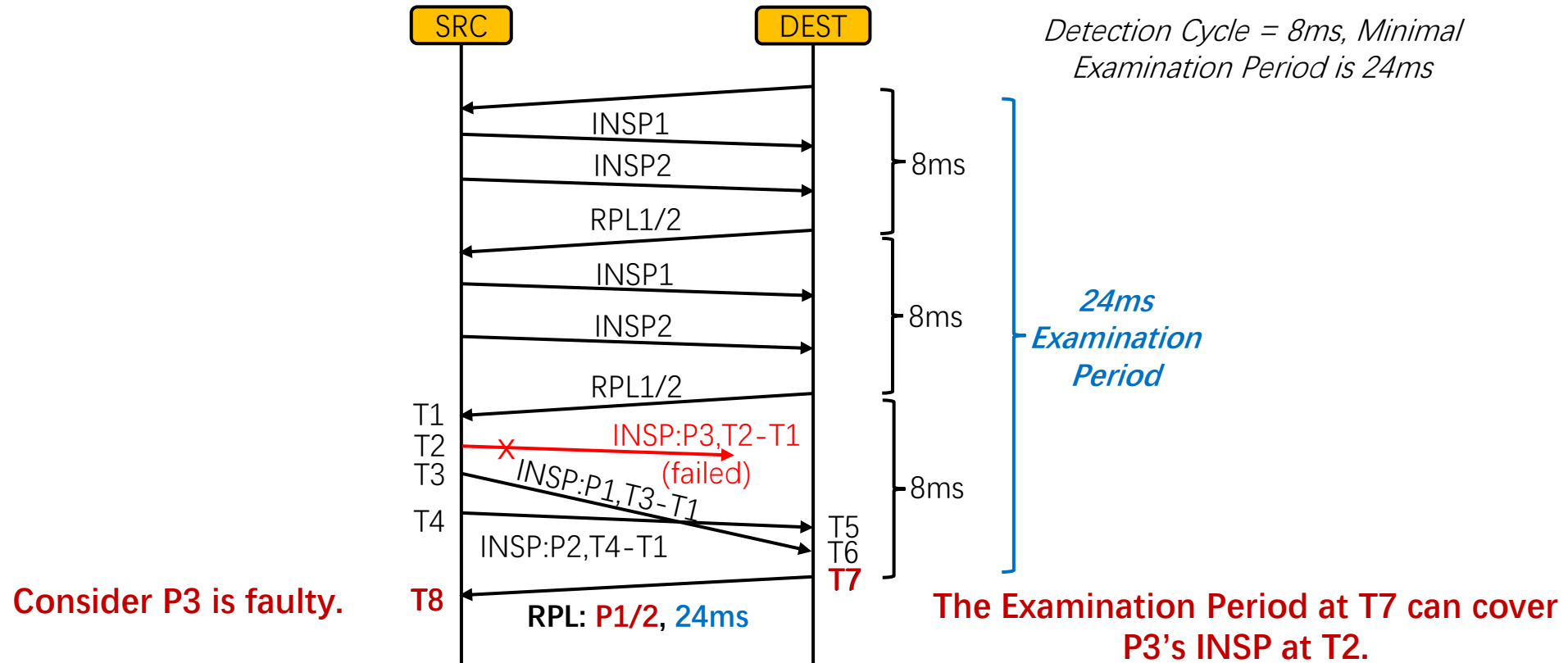
Improving Fault Tolerance

- Harp requires the Examination Period to span **multiple detection cycles**.



Improving Fault Tolerance

- Harp requires the Examination Period to span **multiple detection cycles**.



Improving Fault Tolerance

- Harp requires the Examination Period to span **multiple detection cycles**.



Detection Cycle = 8ms, Minimal is 24ms

Advantages

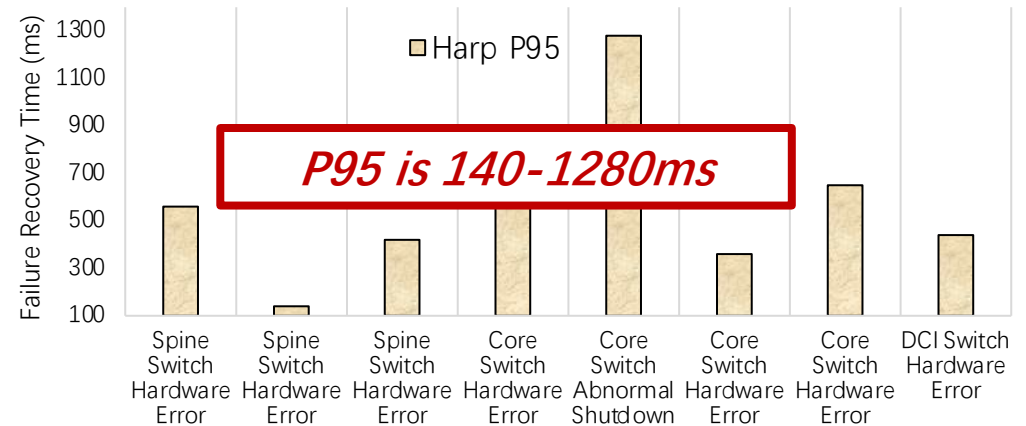
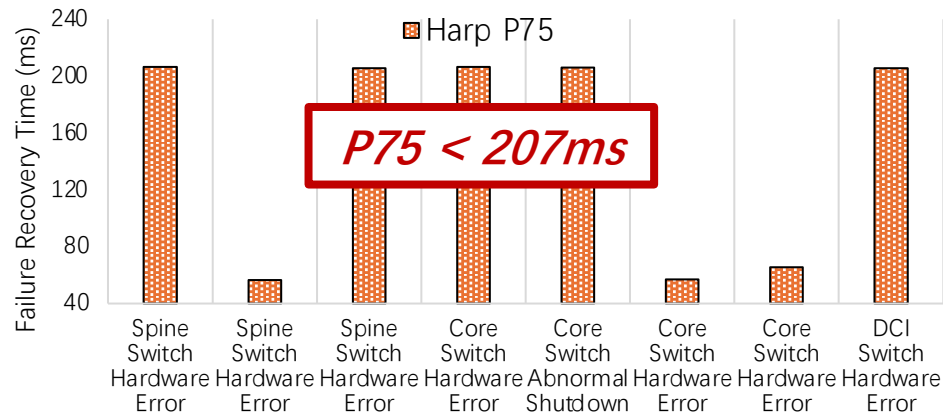
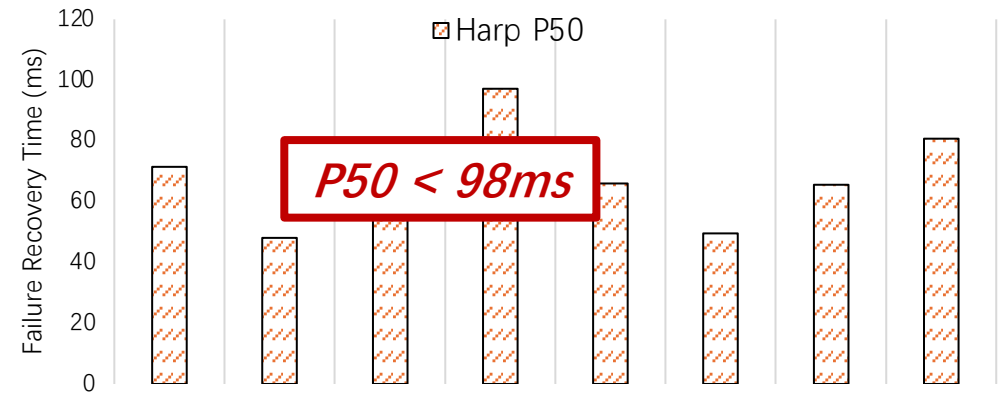
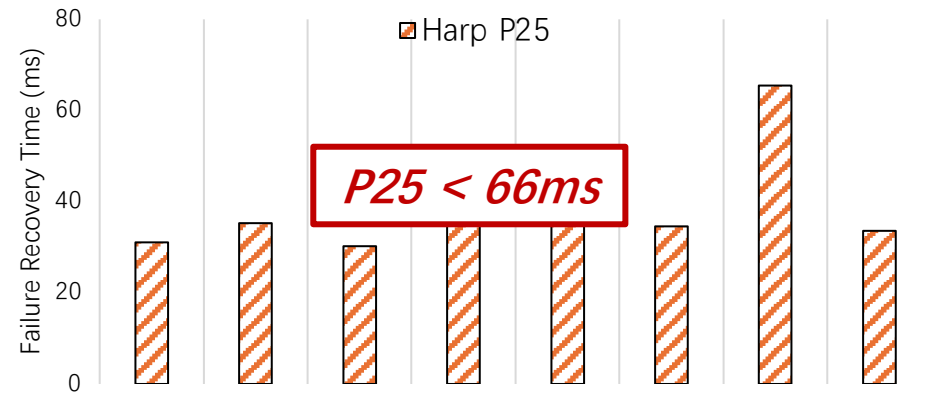
- Knowing the **states of all paths** increases path switching **success probability**.
- Run the mechanism at **millisecond intervals**, ensuring a low failure recovery latency.
- Probe the paths used by VM traffic, enabling **high path coverage**.

Consider P3 is faulty.

16

RPL: P2/3, 24ms

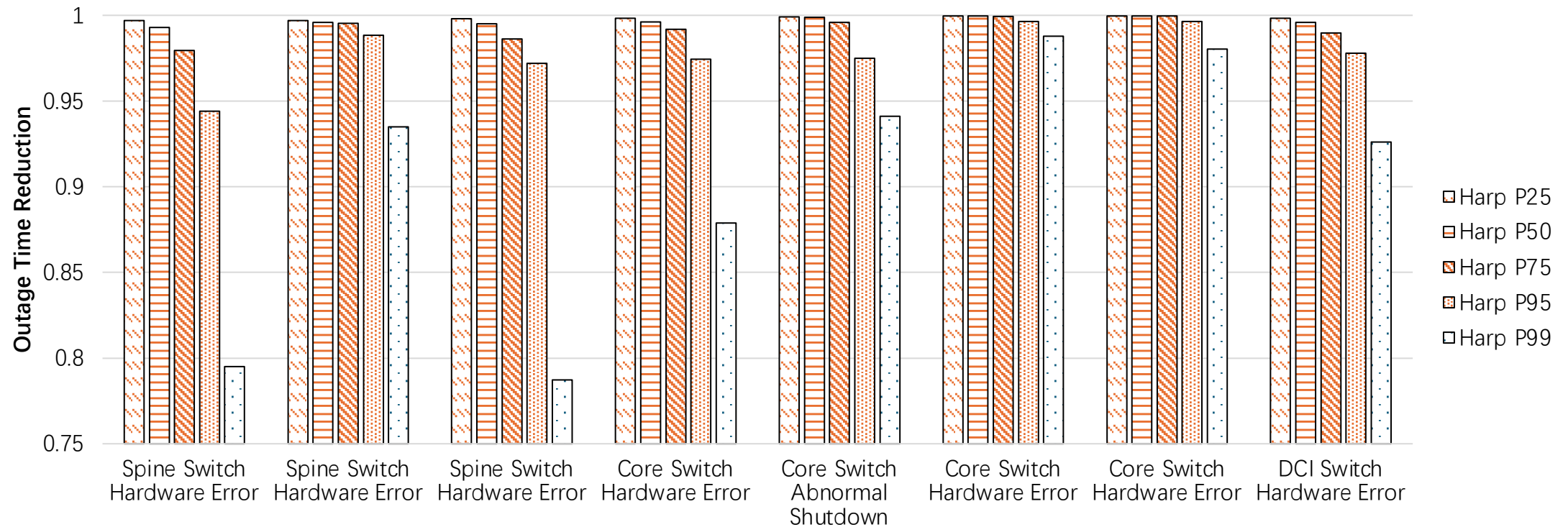
Evaluation of Failure Recovery Time



Eight real-world network failures in Tencent Cloud.

Evaluation of Failure Recovery Time

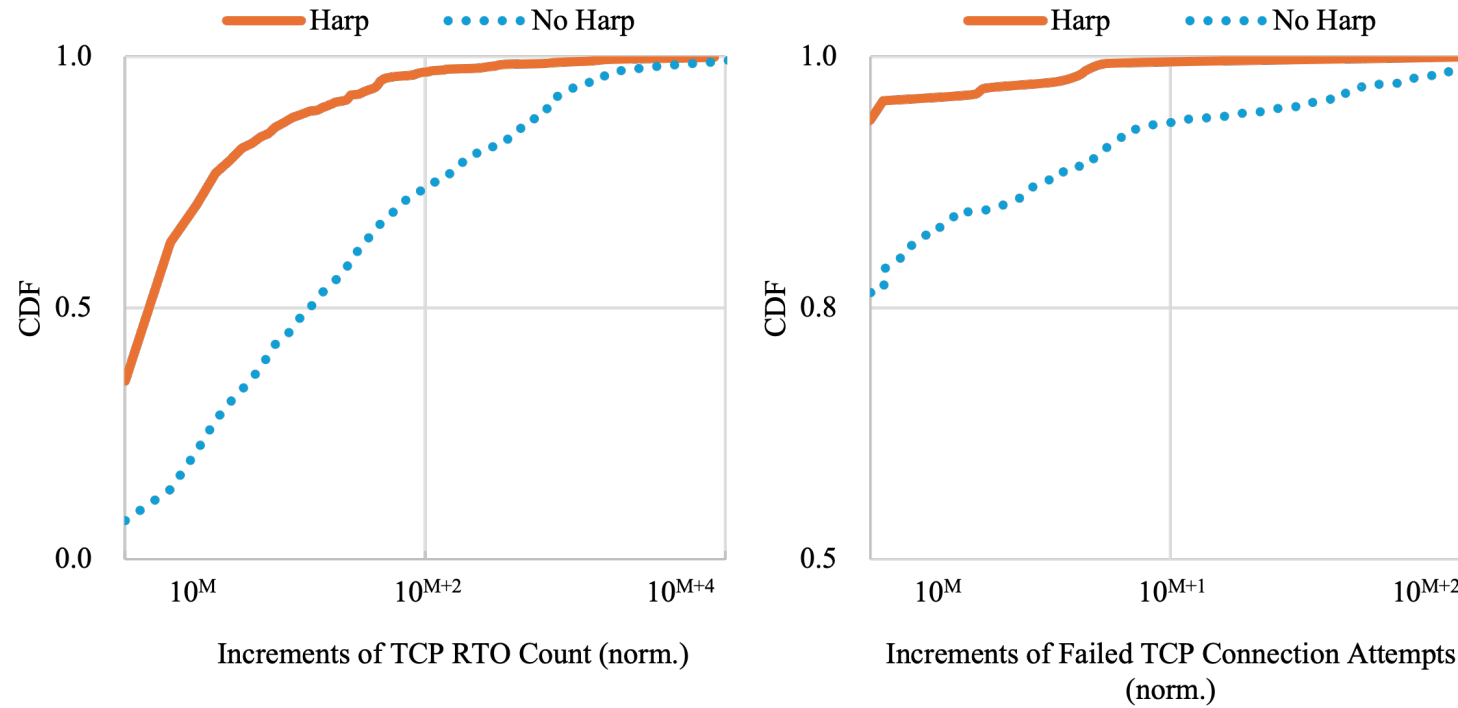
- Comparison with the **Pingmesh-based Approach** used in Tencent Cloud



- Harp significantly **reduces** the VPC network's **outage time** by **78.71%-99.97%**.

Eight real-world network failures in Tencent Cloud.

Evaluation of VM Performance



- Numbers of TCP RTOs and Failed TCP Connection Attempts are **significantly reduced with Harp.**

Conclusion

- Harp has been deployed in Tencent Cloud for over two years across almost all regions.
- Harp can detect failures and restore VPC network connectivity on a sub-second timescale.
- Harp does not require special hardware supports.
- Harp is transparent to cloud applications.