# Duet: Cloud Scale Load Balancing with Hardware and Software

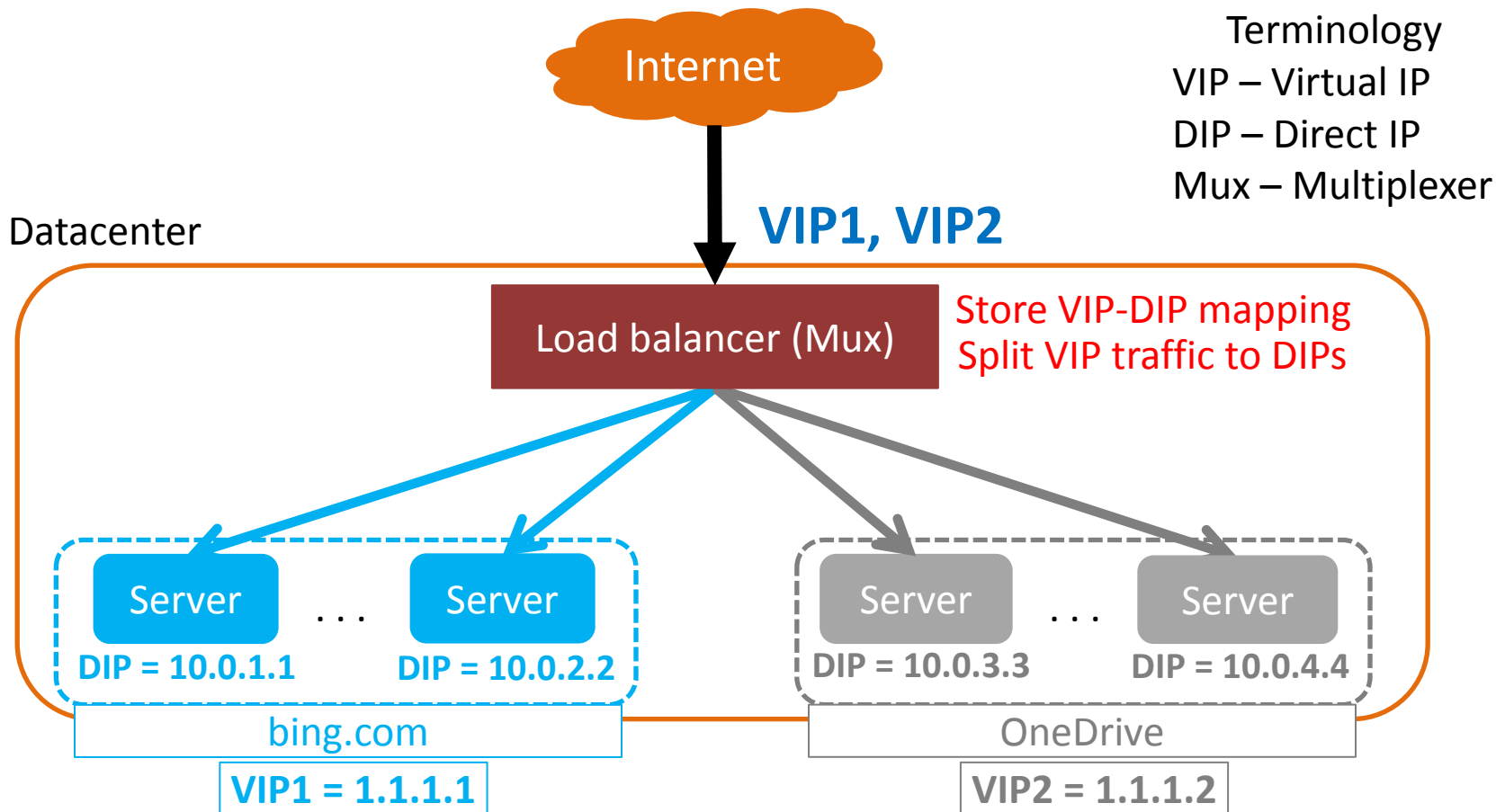**Rohan Gandhi**    Hongqiang Harry Liu   Guohan Lu    Jitu Padhye

Y. Charlie Hu                Lihua Yuan        Ming Zhang

PURDUE UNIVERSITY

Microsoft

# Load Balancer is Critical
# For Online Services

Internet

Terminology
VIP – Virtual IP
DIP – Direct IP
Mux – Multiplexer

Datacenter

**VIP1, VIP2**

Load balancer (Mux)

Store VIP-DIP mapping
Split VIP traffic to DIPs

Server ... Server          Server ... Server

**DIP = 10.0.1.1**   **DIP = 10.0.2.2**          **DIP = 10.0.3.3**   **DIP = 10.0.4.4**

bing.com                          OneDrive

**VIP1 = 1.1.1.1**                **VIP2 = 1.1.1.2**

Load balancer provides high availability and scalability

# Existing LBs Have Limitations

## Specialized Hardware LBs

**Too costly**

$100+ million for 15 Tbps

**Poor robustness**

1+1 redundancy
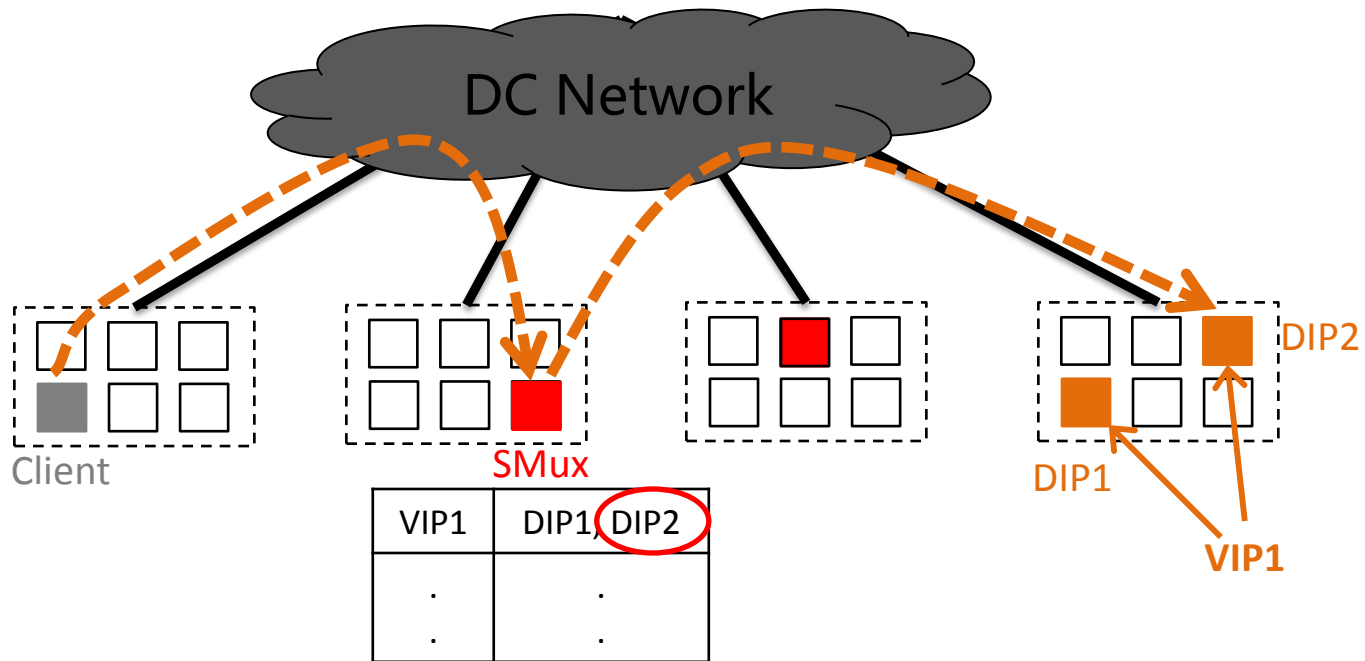
## Software LBs

**Scale with demand**

Scale up/down according to VIP traffic

**High robustness**

n+1 redundancy
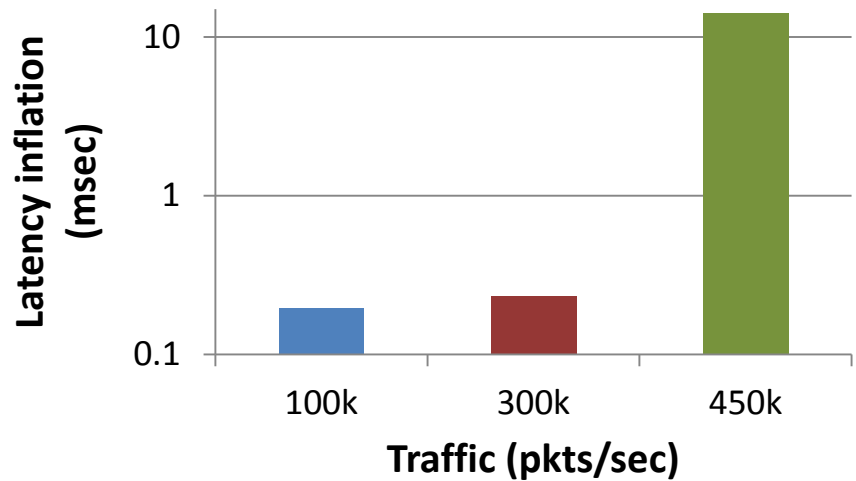
Software LBs have
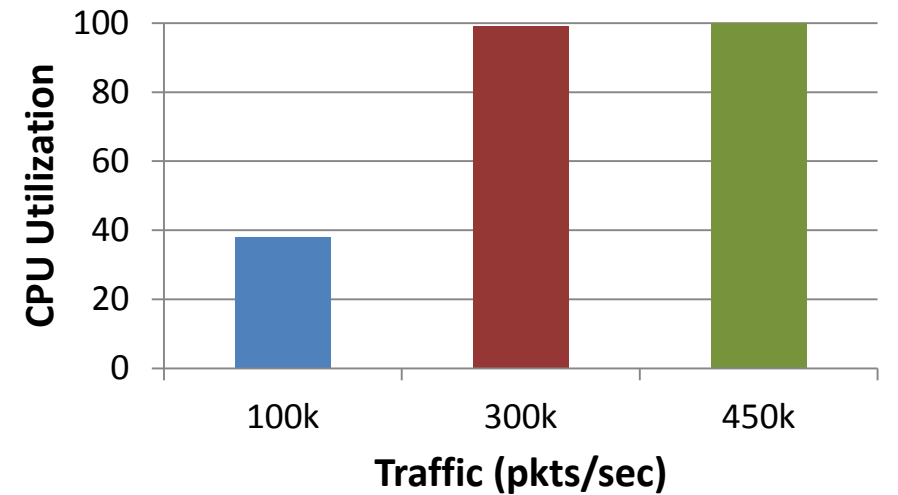cost and performance limitations

# Software LB Design



Software LB Benefits:

- High robustness
- Scale with demand

# Software LB has Limitations



High latency inflation: 200 usec

Low capacity: 300k pkts/sec

5k SMuxes needed at 15Tbps traffic in 50k server DC

# How can we build high performance, low cost and robust load balancer?

## Duet ideas

- Use commodity switches as hardware Muxes

- Use software Muxes as a backstop

# Can Switch Act As a Mux

## Switches offer:

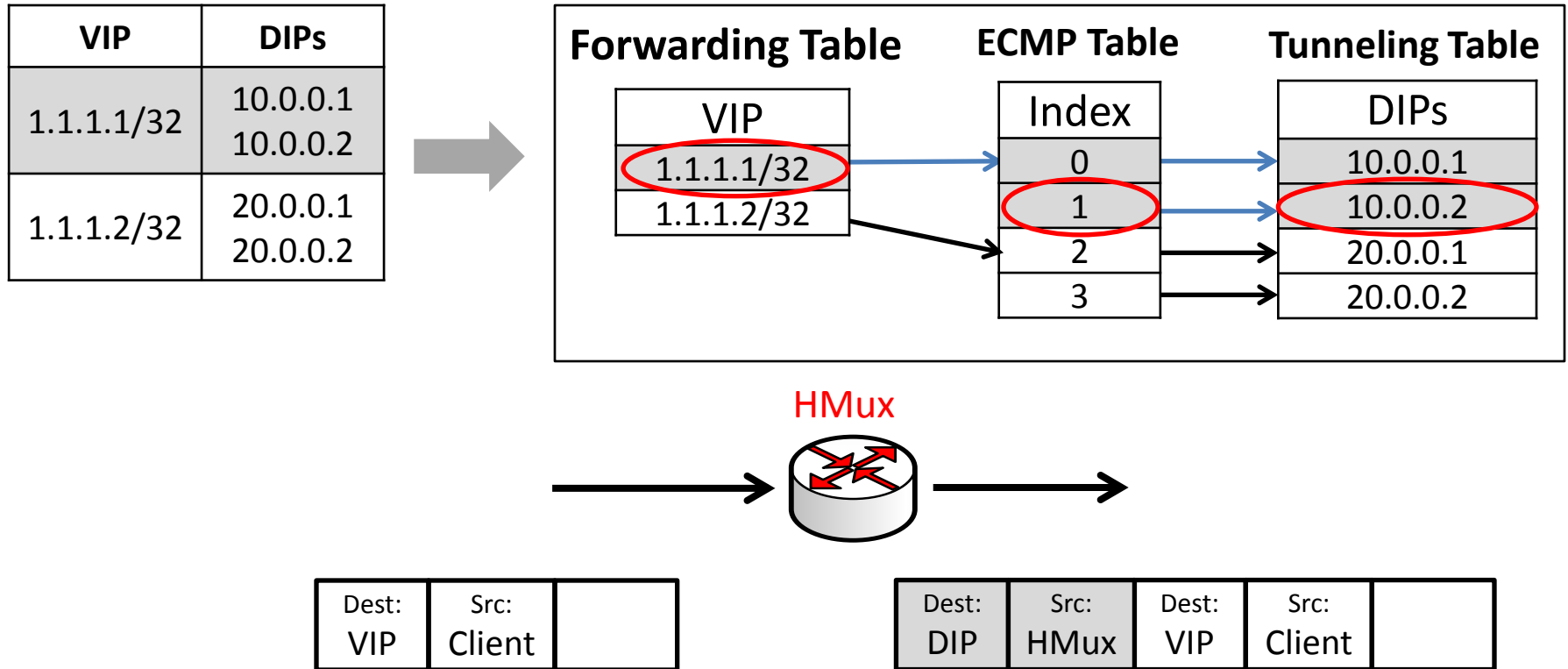- High capacity (500+ million pkts/sec)

- Low latency inflation (1 usec)

## Mux functionalities

- Split VIP traffic across DIPs

- Forward VIP traffic to DIPs

## Switch resources

- ECMP

- Tunneling

# Implementing HMux on Switch

| VIP | DIPs |
|---|---|
| 1.1.1.1/32 | 10.0.0.1<br>10.0.0.2 |
| 1.1.1.2/32 | 20.0.0.1<br>20.0.0.2 |

**Forwarding Table**

| VIP |
|---|
| 1.1.1.1/32 |
| 1.1.1.2/32 |

**ECMP Table**

| Index |
|---|
| 0 |
| 1 |
| 2 |
| 3 |

**Tunneling Table**

| DIPs |
|---|
| 10.0.0.1 |
| 10.0.0.2 |
| 20.0.0.1 |
| 20.0.0.2 |

HMux

| Dest:<br>VIP | Src:<br>Client | |
|---|---|---|

| Dest:<br>DIP | Src:<br>HMux | Dest:<br>VIP | Src:<br>Client | |
|---|---|---|---|---|

# Key Design Challenges

- Limited switch memory

- High failure robustness

- VIP assignment

- VIP migration

# Challenge 1: Switches have Limited Memory

Workload: 100k+ VIPs and 1+ millions DIPs
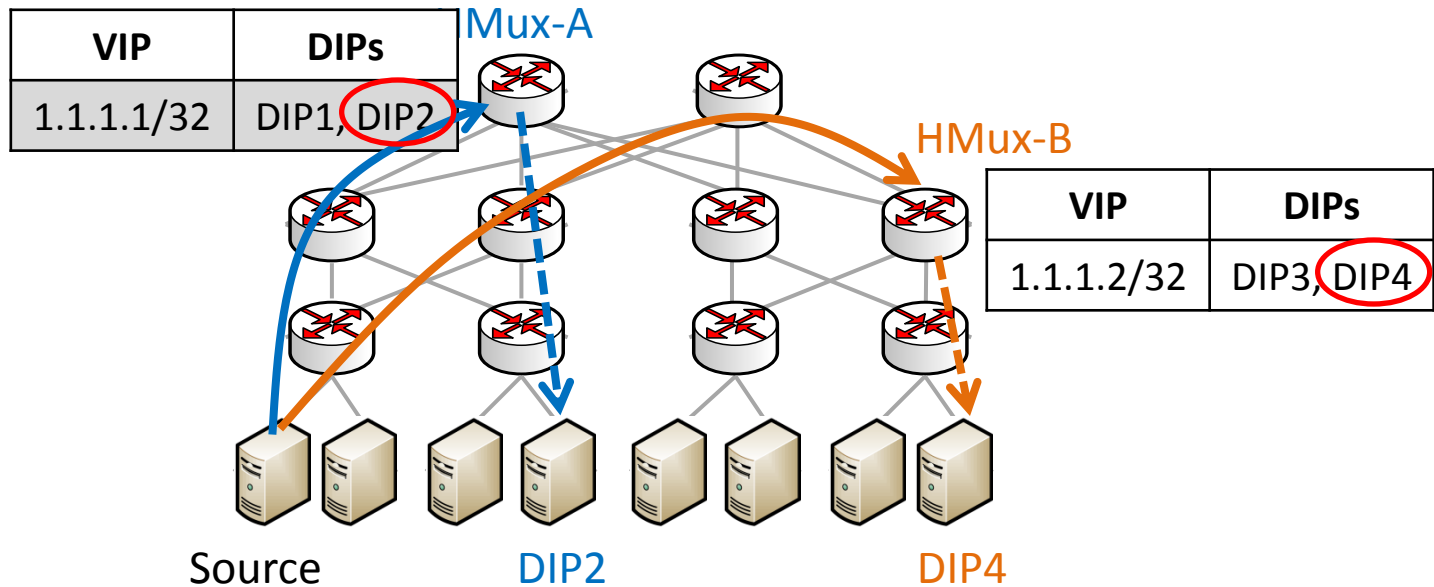
Single HMux cannot store all VIPs and DIPs

| Table | Forwarding | ECMP | Tunneling |
|---|---|---|---|
| Max. size | 16k | 4k | 512 |

Max VIPs                    Max DIPs

# Solution: Partitioning VIPs across HMuxes



| VIP | DIPs |
|---|---|
| 1.1.1.1/32 | DIP1, DIP2 |

Mux-A

HMux-B

| VIP | DIPs |
|---|---|
| 1.1.1.2/32 | DIP3, DIP4 |

Source          DIP2          DIP4

| Capacity | VIPs | DIPs |
|---|---|---|
| Single HMux | 16k | 512 |
| All HMuxes | 16k | 512 * 2k = 1M |

Fixed          Scales with #DIPs

# Challenge 2: High Robustness



HMux-A 1.1.1.1/32 — FAIL

HMux-B 1.1.1.2/32

Source    DIP1    DIP2

- Availability during failure?
- Large number of VIPs?

# Idea: Integrate SMux with HMux

|  | **HMuxes** | **SMuxes** | **Duet** |
|---|---|---|---|
| Low latency | ✓ | ✗ | ✓ |
| High capacity | ✓ | ✗ | ✓ **?** |
| High availability | ✗ | ✓ | ✓ |
| Scale to large #VIPs | ✗ | ✓ | ✓ |

# Solution: Use SMuxes As a Backstop



HMux-A

| VIP | DIP |
|-----|-----|
| 1.1.1.1/32 | DIP1, DIP2 |

HMux-B

| VIP | DIP |
|-----|-----|
| 1.1.1.2/32 | DIP3, DIP4 |

SMux-A

| VIP | DIP |
|-----|-----|
| 1.1.1.1/32 | DIP1, DIP2 |
| 1.1.1.2/32 | DIP3, DIP4 |

SMux-B

| VIP | DIP |
|-----|-----|
| 1.1.1.1/32 | DIP1, DIP2 |
| 1.1.1.2/32 | DIP3, DIP4 |

# Solution: Use SMuxes As a Backstop



HMux-A

VIP = 1.1.1.1/32

HMux-B

VIP = 1.1.1.2/32

Source

SMux-A

VIP = 1.1.1.0/24

SMux-B

VIP = 1.1.1.0/24

- High availability during failure
- Scale to large #VIPs

# VIP Traffic Distribution is Highly Skewed



Top 10% VIPs carry 99% traffic

Duet handles 86-99.9% traffic using HMuxes

# Challenge 3: How to Assign VIPs?

VIP-1

VIP-2

.
.
.

VIP-k

Objective: Maximize traffic handled by HMuxes

Input:
VIP traffic, DIP locations
Topology

Constraints:
Switch memory
Link capacity

SMux          SMux

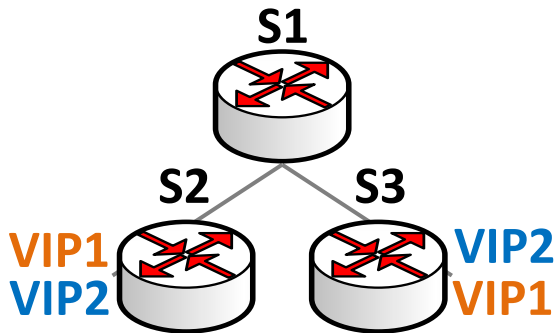# Challenge 4: How to Migrate VIPs?

Current

New

Maintain VIP availability

S1

S2    S3

VIP1    VIP2
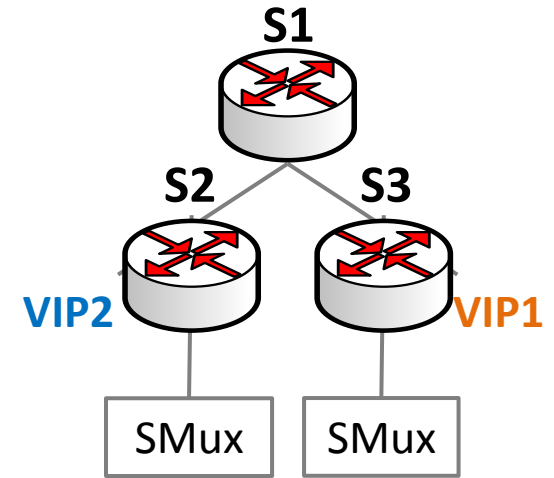
S1

S2    S3

VIP2    VIP1

Withdraw and announce

Announce before withdraw

S1

S2    S3

VIP1    VIP2
VIP2    VIP1

S1

S2    S3

VIP2    VIP1    Limited Memory

VIP1    VIP2

# Solution: Migrate VIPs through SMuxes

Current

S1

S2    S3

**VIP1**           **VIP2**

SMux    SMux

Withdraw VIPs from
old location

S1

S2    S3

**VIP1, VIP2**  **VIP1, VIP2**

SMux    SMux

New

S1

S2    S3

**VIP2**           **VIP1**

SMux    SMux

Announce VIPs from
new location

Fast and maintains VIP availability

# Duet Extensions

☑ SNAT

☑ Support VIPs with 512+ DIPs

☑ Port based load balancing

☑ Load balancing in virtualized networks
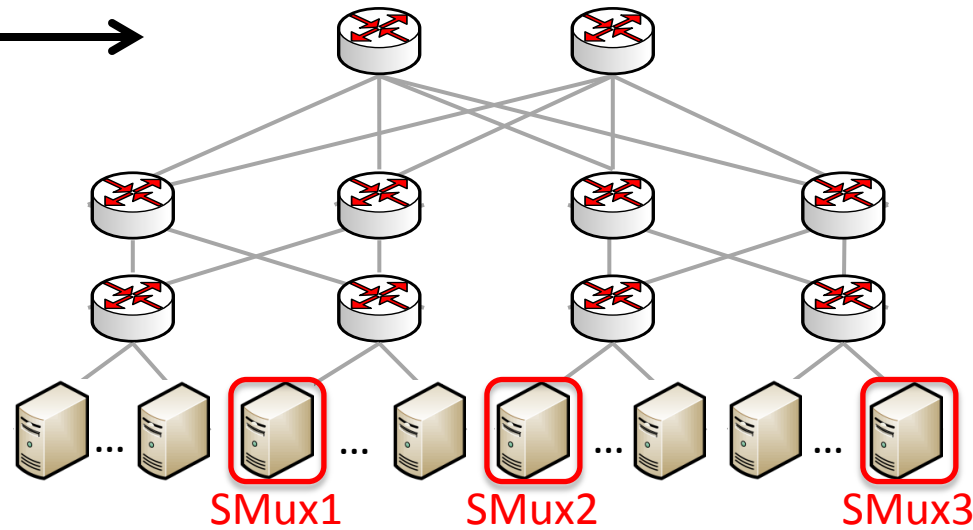
# Experimental Setup
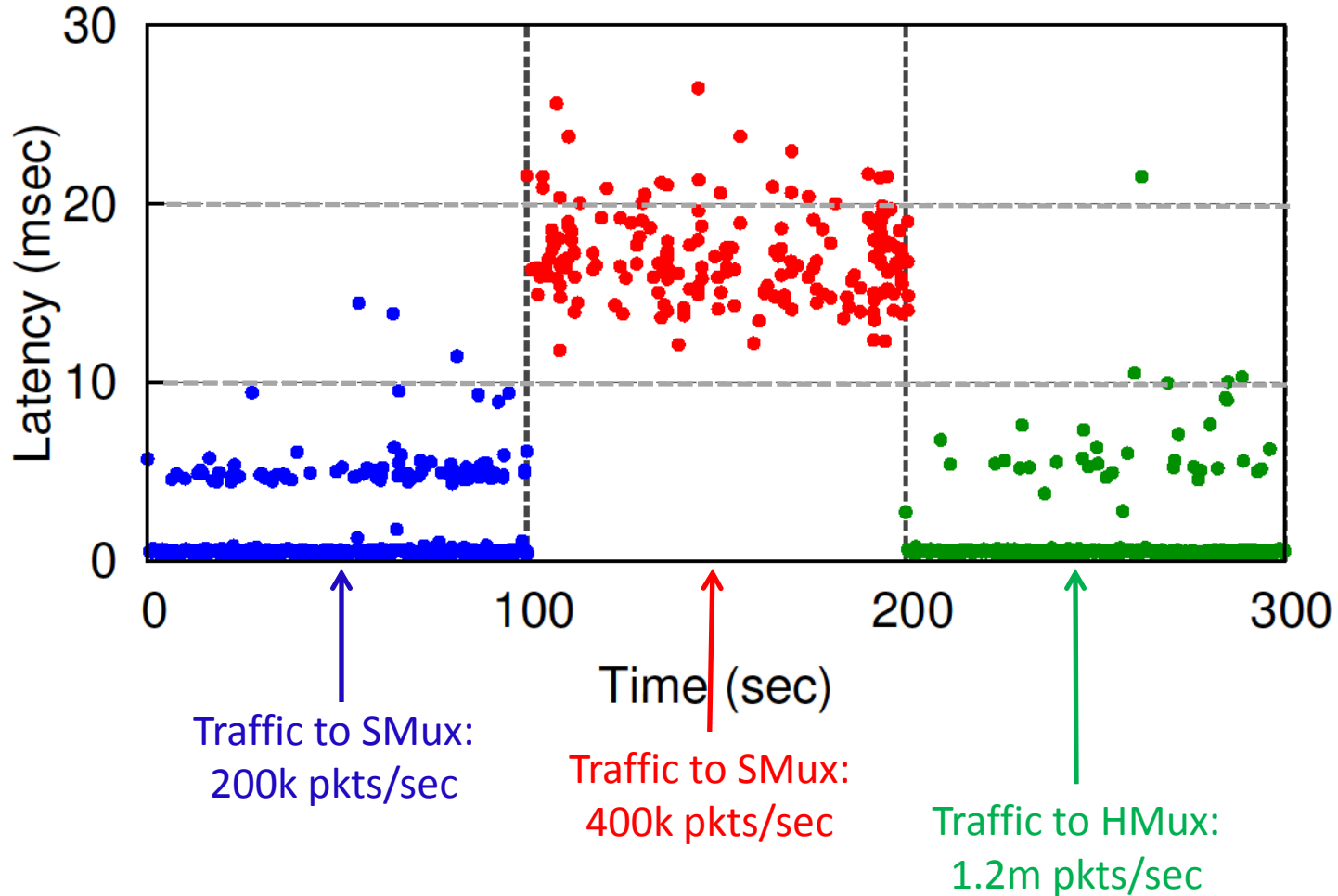
## Testbed

- 10 switches, 3 SMuxes
- 10 VIPs, 34 DIPs

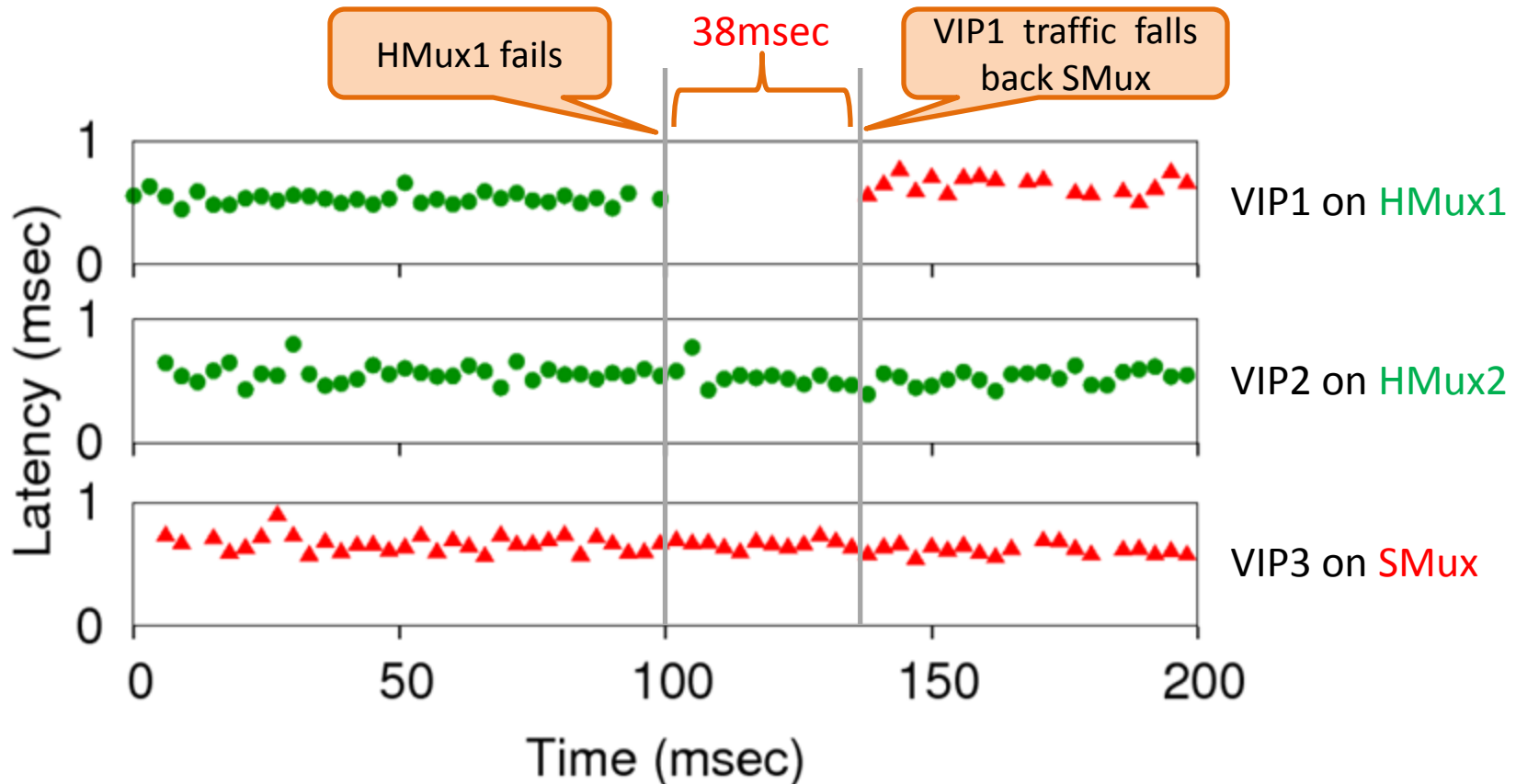## Simulation

- Topology and traffic trace from Azure DC

- High capacity
- High availability
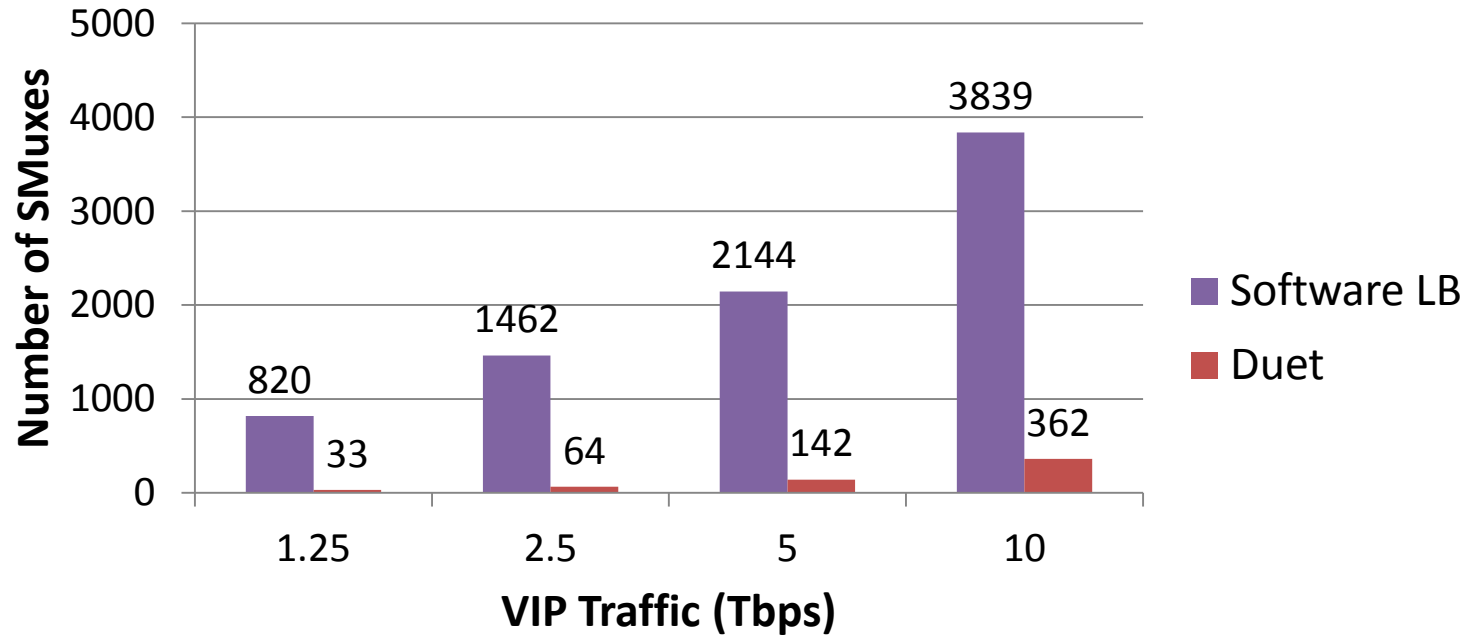- Low cost

# Duet Provides High Capacity



Traffic to SMux:
200k pkts/sec

Traffic to SMux:
400k pkts/sec

Traffic to HMux:
1.2m pkts/sec

HMux has larger capacity and lower latency than SMux

# Duet Provides High Availability

# Duet Reduces Cost



Duet reduces cost by 10-24x

# Summary

- Specialized and software LBs have cost and performance problems

- Duet key ideas:
    - Use commodity switches as HMuxes
    - Use small number of SMuxes as backstop

- Benefits:
    - Low latency
    - High capacity
    - High robustness
    - Low cost