



ons

NORTH AMERICA

OPEN NETWORKING //

Enabling Collaborative  
Development & Innovation

# VPP Accelerated High Performance & Scalable L3DSR L4 Load Balancer on Top Clos

Yusuke Tatsumi @ Yahoo Japan Corporation & Naoyuki Mori @ Intel

*Acknowledgement:*

*Shunya Kitada, Yuta Kinoshita @ Yahoo Japan Corporation*

*Hongjun Ni, Ray Kinsella @Intel*

*Pierre Pfister, Jerome Tollet @Cisco*

Hosted By

THE LINUX FOUNDATION | LFN NETWORKING



ons  
NORTH AMERICA  
OPEN NETWORKING //  
Enabling Collaborative  
Development & Innovation

# Agenda

- Yahoo! JAPAN introduction, background, and issues
- FD.io VPP Overview
- VPP LB data-plane
- VPP LB control-plane
- Summary
- Call to Action

Hosted By

 THE LINUX FOUNDATION |  LFNETWORKING



# ONS

NORTH AMERICA

OPEN NETWORKING //  
Enabling Collaborative  
Development & Innovation

# Yahoo! JAPAN introduction



Hosted By



\*Other names and brands may be claimed as the property of others.



ons  
NORTH AMERICA  
OPEN NETWORKING //  
Enabling Collaborative  
Development & Innovation

# Yahoo! JAPAN introduction

Monthly Page Views  
**70+ Billion**

Number of services  
**100+**

Daily Unique Browser  
**90+ Million**

Daily Unique Browser  
(Only Smartphone)  
**60+ Million**

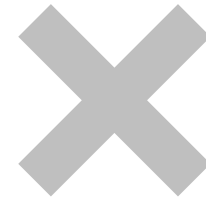
Many Load Balancers support our services!



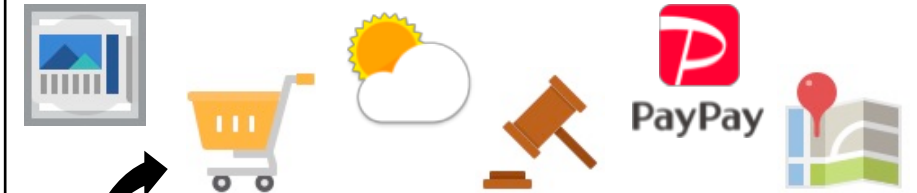
# Background and issues of our LB

## Total responsibility of LB =

Demands for speed:



Explosive demands

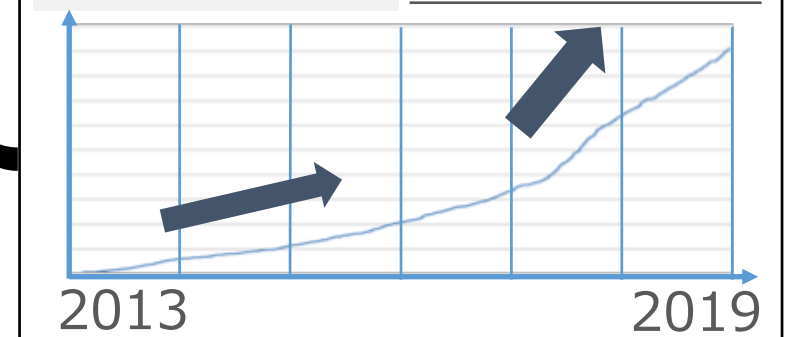


More features are demanded of our LBaaS:

- **Scale-in/out LB capability**
- **Robustness of LB system**
- **Elastic management of VIP**

[3]

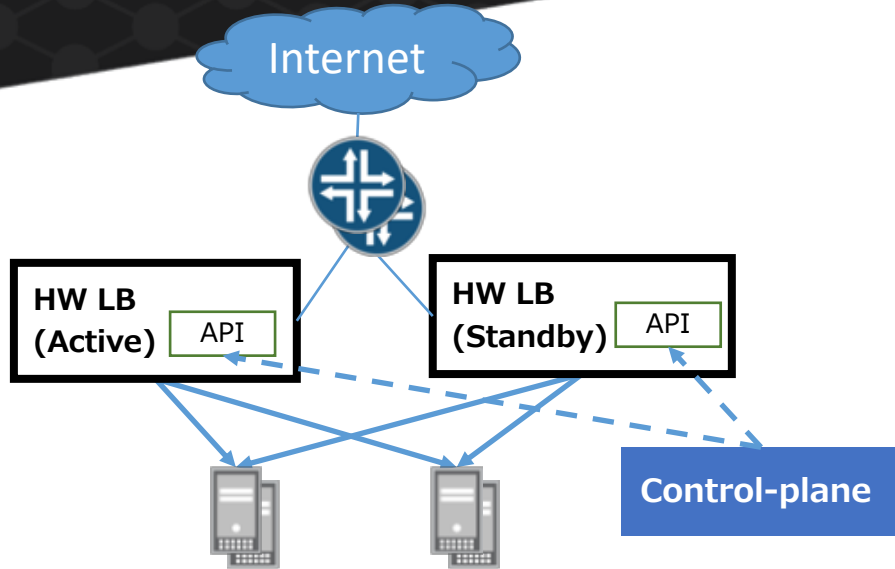
The # of VMs Now: 140k+



Hosted By

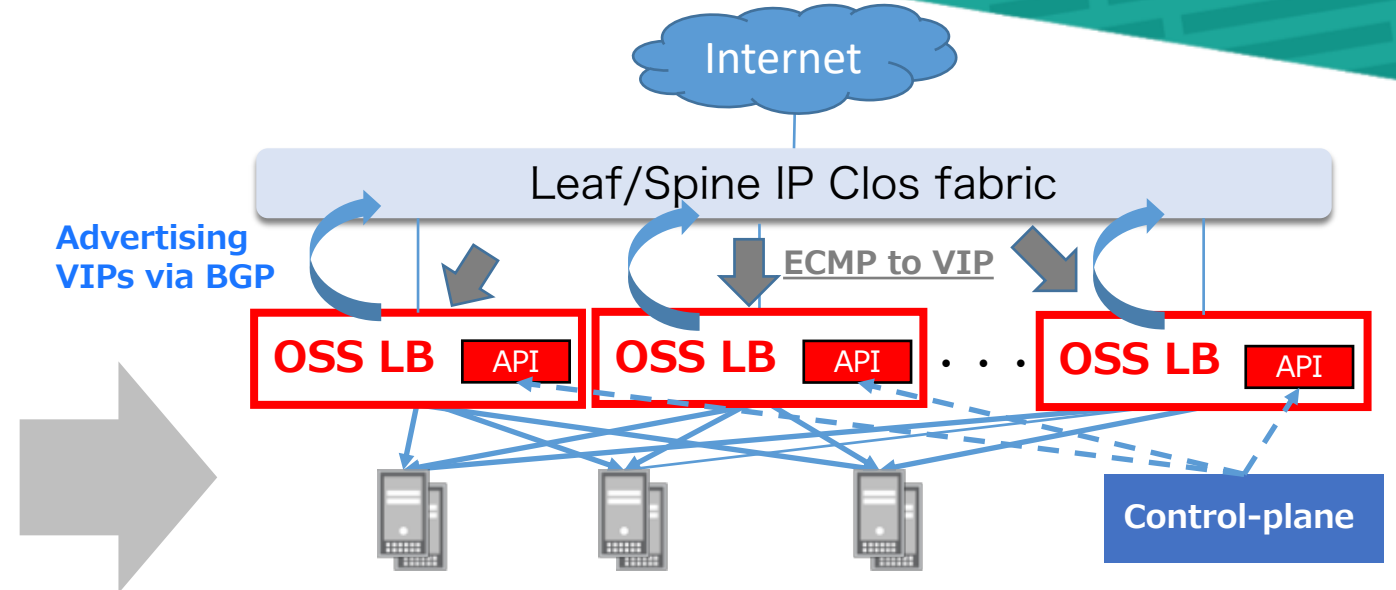


# Limitations of our system and OSS possibilities



## Our current appliance based LBaaS

- Active-backup (2N) LB capability on top MLAG
  - hard to scaling-out
- Vender proprietary (Black box)
  - Lead time & Hardware EoL
- Careful operation (plan/exec) to place LB

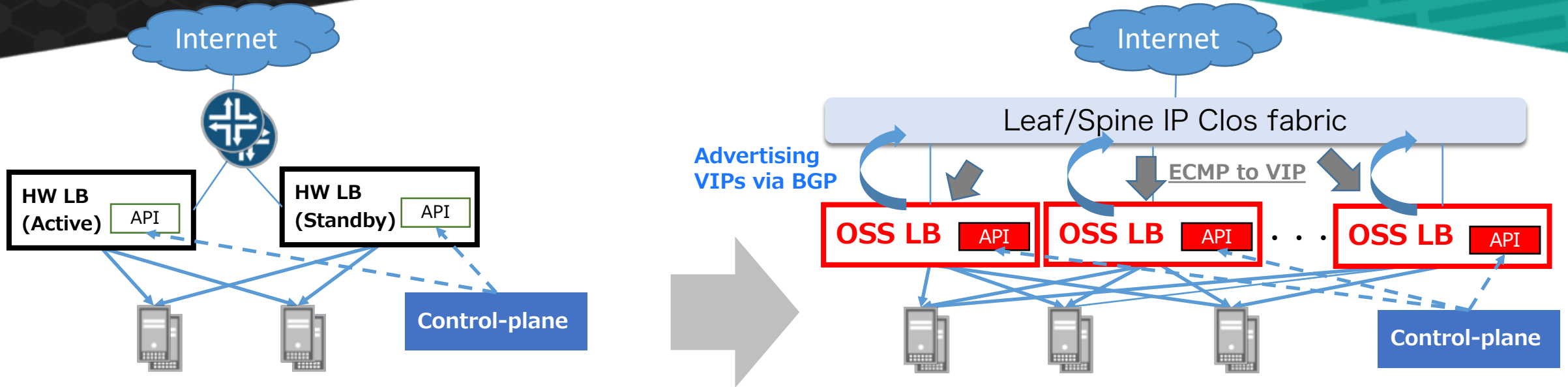


## Scalable LBaaS architecture with OSS

- Scaling-in/out (N+1) LB capability on top Clos
- Commodity server/NIC + OSS (White box)
- Easy operation to place LB with BGP







# Limitations of our system and OSS possibilities

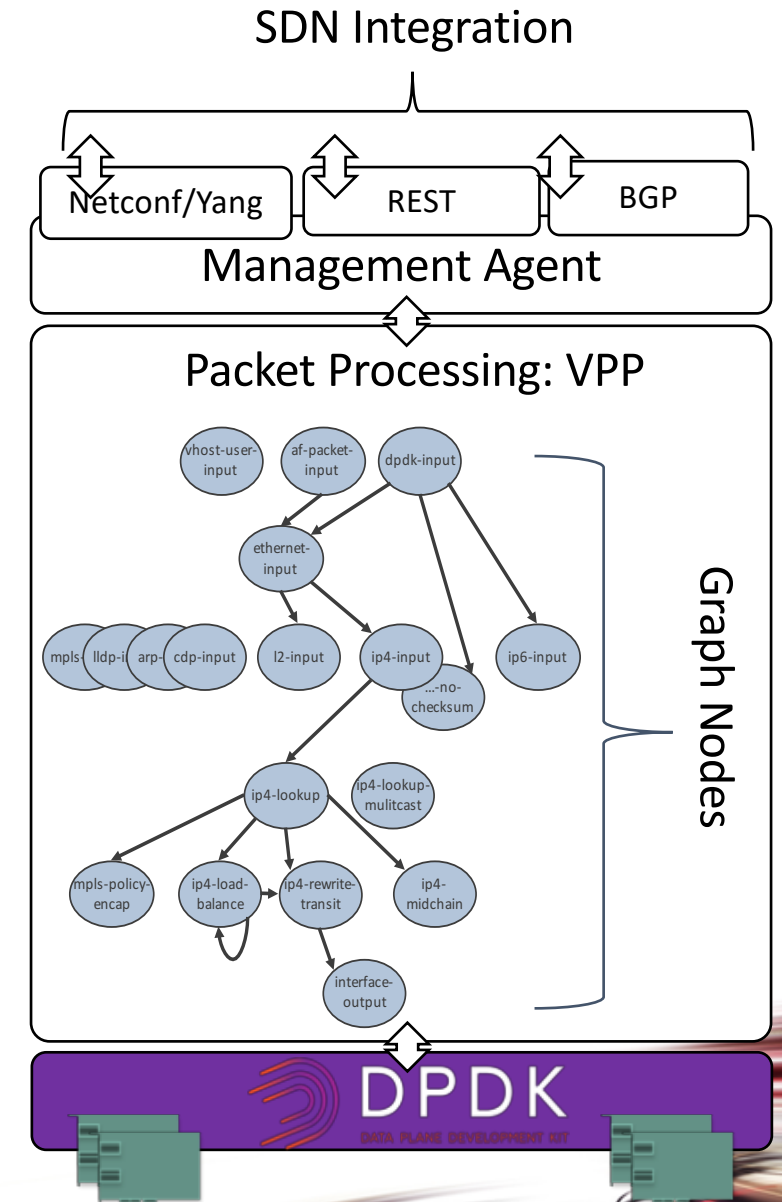


## Strategy to break free of our current LBaaS limitations:

1. Collaborate with existing high performance OSS data-plane FD.io VPP
2. Develop control-plane from scratch based on our operational experiences

# FD.io\* VPP Overview

Universal data plane	Extensible Modular Design
<ul style="list-style-type: none"> <li>Layer 2 – 4 Network Stack</li> <li>CP, TM, Overlays and more ...</li> <li>Linux (and FreeBSD) support</li> <li>Kernel Interfaces (Netmap, FastTap)</li> <li>Container and Virtualization support</li> <li>Appliance, Infrastructure, VNF &amp; CNF</li> </ul> 	<p><b>Architecture</b></p> <ul style="list-style-type: none"> <li>Pluggable, easy to understand &amp; extend.</li> <li>Mature graph node architecture.</li> </ul> <p><b>Plugins</b></p> <ul style="list-style-type: none"> <li>Full control to reorganize the pipeline.</li> <li>Fast, plugins are equal citizens</li> </ul> 
Fast, Scalable and Deterministic	Developer friendly
<ul style="list-style-type: none"> <li>L2XC - 25+ Mpps per core</li> </ul> <p><b>Deterministic</b></p> <ul style="list-style-type: none"> <li>0 packet drops, ~15µs latency</li> <li>Continuous &amp; extensive latency testing.</li> </ul> <p><b>Scalability</b></p> <ul style="list-style-type: none"> <li>Linear scaling with core/thread count</li> <li>Supporting millions of concurrent L[2,3] tables entries.</li> </ul> 	<ul style="list-style-type: none"> <li>Runtime counters for everything. ( throughput, ipc, errors etc )</li> <li>Full pipeline tracing facilities.</li> <li>Multi-language API bindings.</li> <li>VPP command line introspection.</li> </ul> 



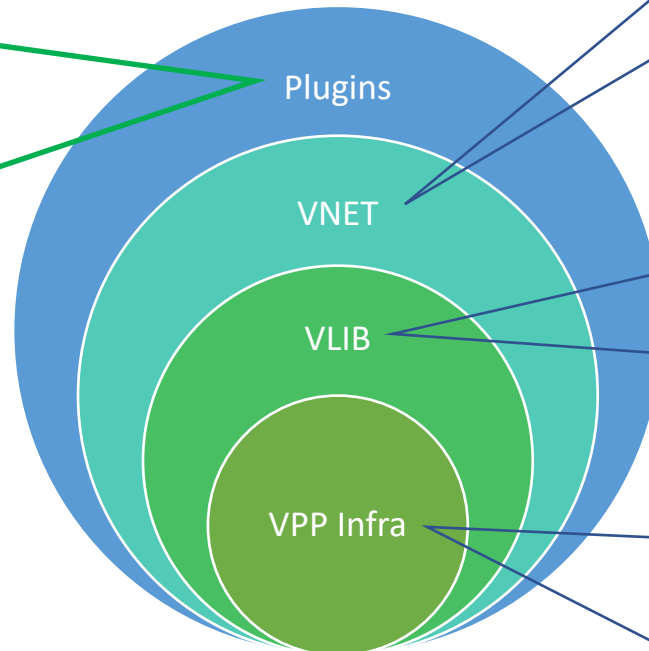
\*Other names and brands may be claimed as the property of others.



# VPP Layering for modular design

## Plugins

- Plugins can be in-tree:  
SNAT,  
Policy ACL,  
Flow Per Packet,  
ILA,  
IOAM,  
**Load Balancer,**  
SIXRD,  
VCGN  
Segment Routing
- Separate fd.io project:  
NSH\_SFC



## VNET

VPP networking source

- Devices
- Layer [2, 3, 4]
- Session Management
- Overlays
- Traffic Management

## VLIB

VPP application management

- buffer, buffer management
- graph node, node management
- tracing, counters
- threading
- CLI
- and most importantly ...
- main()

## VPP Infra

Library of function primitives, for

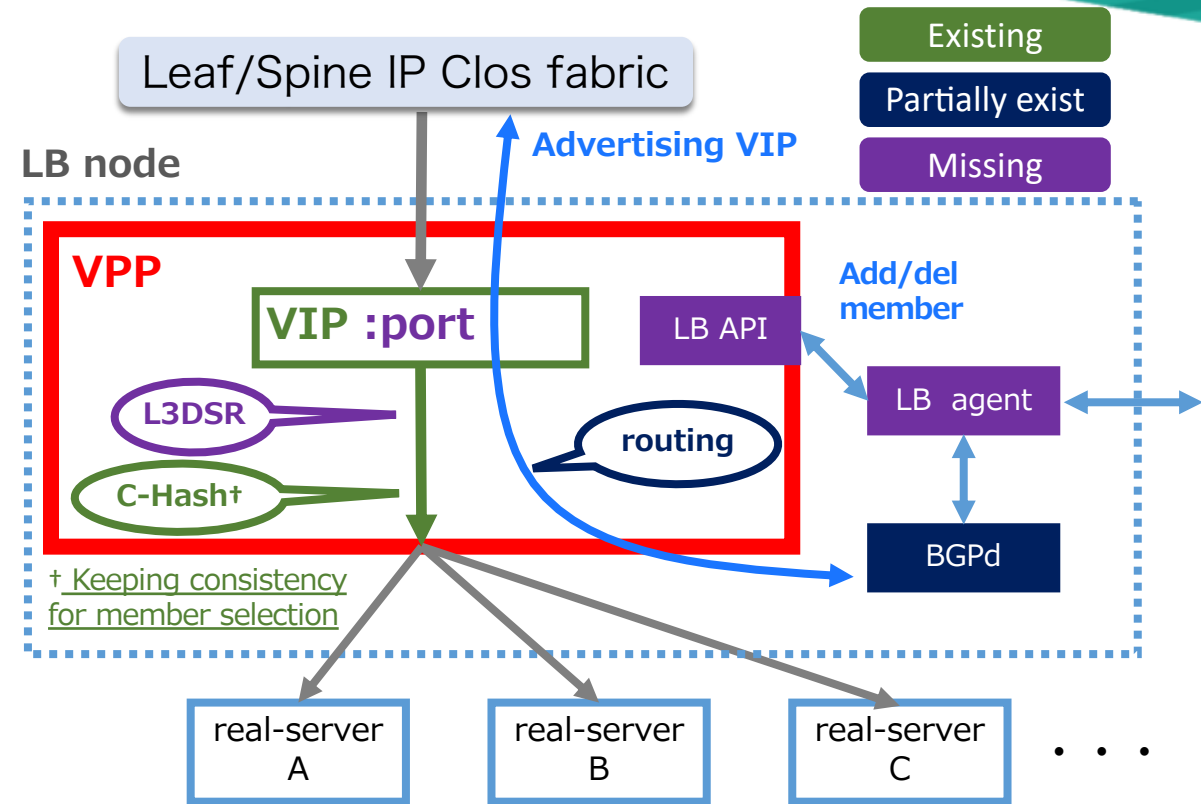
- memory management
- memory operations
- vectors
- rings
- hashing
- timers



# What's available?

# & What's missing?

- OSS based fast data-plane
  - ✓ **VPP**
- Scaling-in/out capability
  - ✓ **Consistent-hash feature in LB**
- Integration with ECMP/BGP on CLOS
  - ✓ **VPP Sandbox Router plugin**
- L3DSR LB (using IPv4 DSCP field) [1]
  - ✗ **VPP LB L3DSR support!**
- Management API
  - ✗ **VPP LB APIs enhancement!**



Hosted By

# VPP LB plugin Enhancements

Category	Feature	Description	First appeared
Framework	Consistent Hash	Brings great scalability and redundancy w/ LB nodes	16.09
	Multi-Threading	Performance scalability w/ CPU cores	16.09
Protocols	GRE4/GRE6	Encap w/ GRE IPv4 and IPv6	16.09
	L3DSR	★ Layer 3 Direct Server Return	18.04
	NAT4/NAT6	NAT (originally for kube-proxy)	18.07
	Port # aware	★ Care TCP/UDP port number for LBaaS integration	18.10
API	Manipulation	★ Set/Get VIP/member over API from controller	18.10
	Statistics	★ Collect more detailed statistics like # of connections	On going

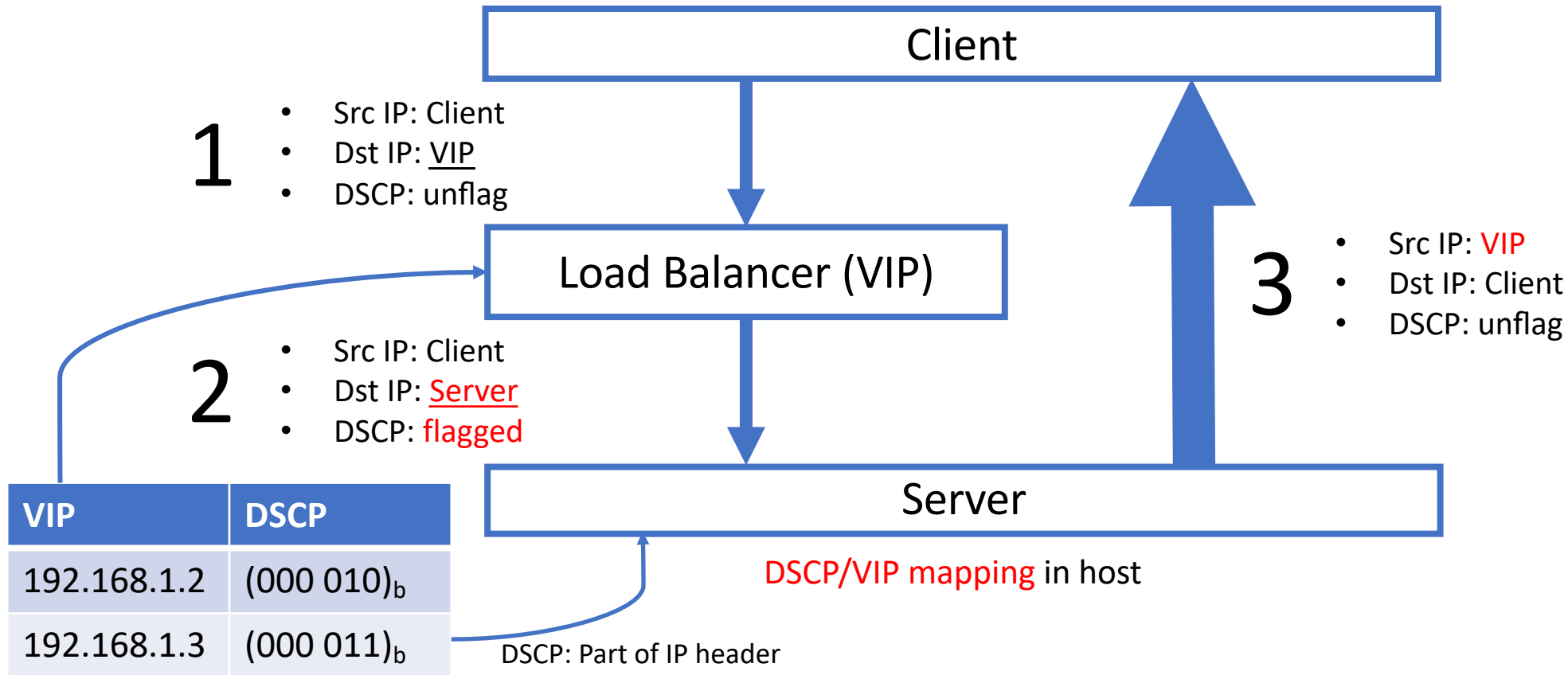
★ Added by this project



\*Other names and brands may be claimed as the property of others.

# Why L3DSR (L3 Direct Server Return) ?

- No tunneling overhead with DSCP
- DSR: Saving LB bandwidth (3. is generally x2-10 larger than 1.)

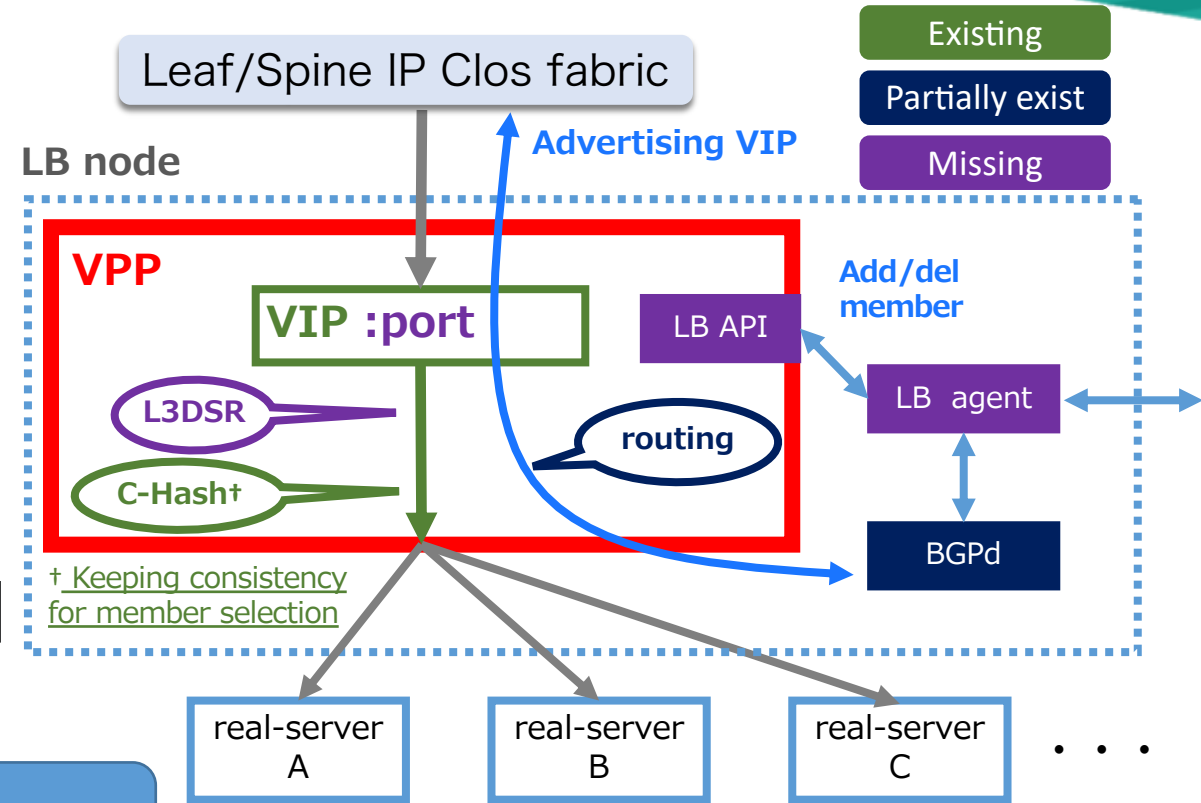




# Filling missing parts, build up system!

- OSS based fast data-plane
- ✓ **VPP**
- Scaling-in/out capability
- ✓ **Consistent-hash feature in LB**
- Integration with ECMP/BGP on CLOS
- ✓ **VPP Sandbox Router plugin**
- L3DSR LB (using IPv4 DSCP field) [1]
- ✓ **VPP LB L3DSR supported**
- Management API
- ✓ **VPP LB APIs**

Successfully Developed!

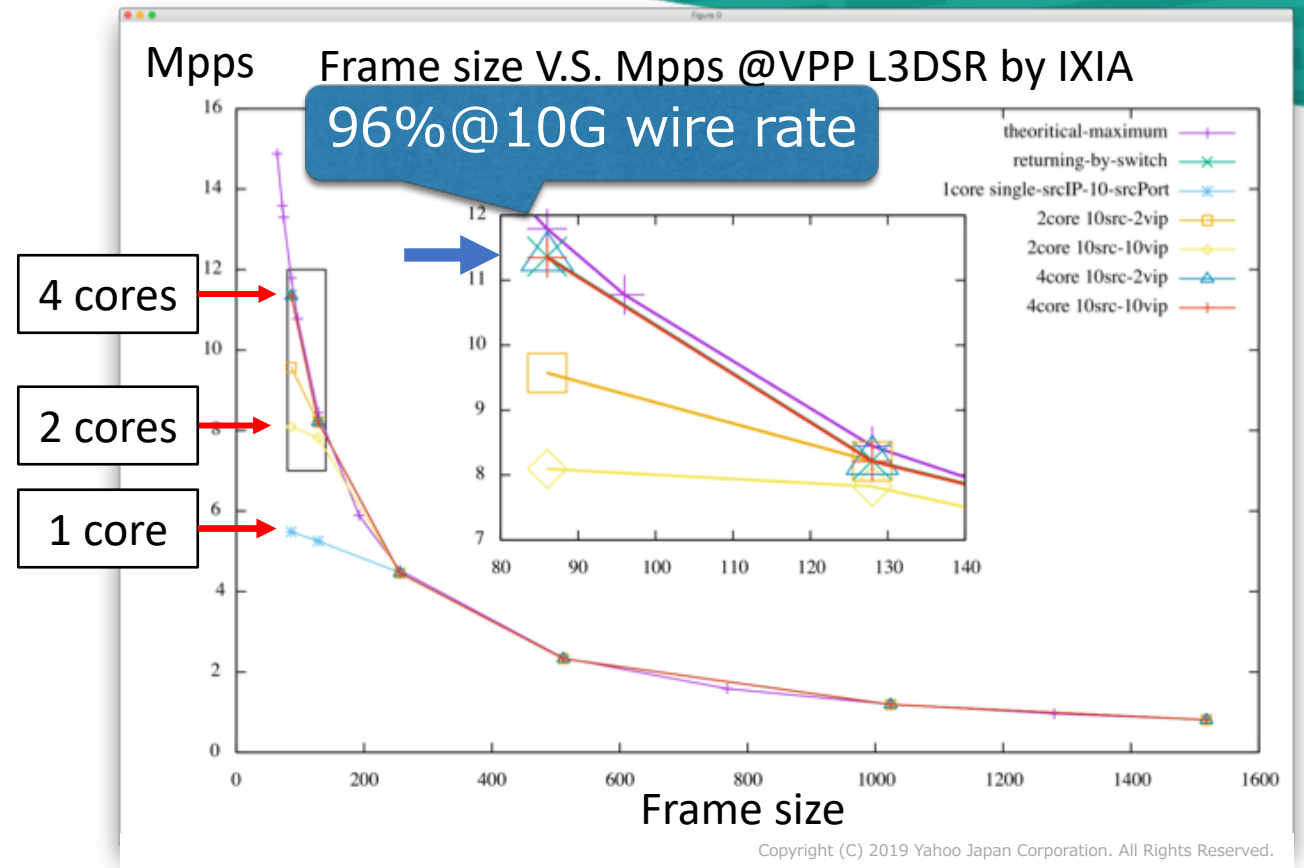
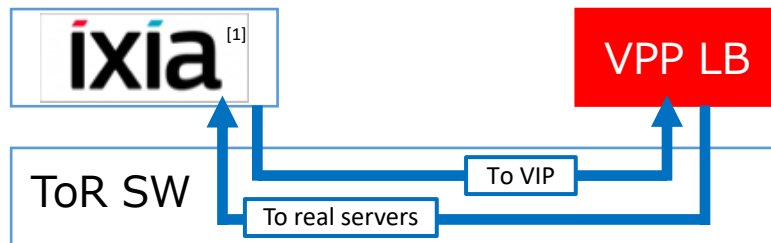


Hosted By



# VPP LB/node performance test @ 10Gbps

- Traffic generator: Ixia IxNetwork
- VPP Server Spec:
  - CPU: Intel® Xeon® Processor E5-2650L v3 \* 2S
  - Memory: 384GB
  - NIC: Intel® Ethernet Network Connection X540-AT2 (10Gbps\*2)

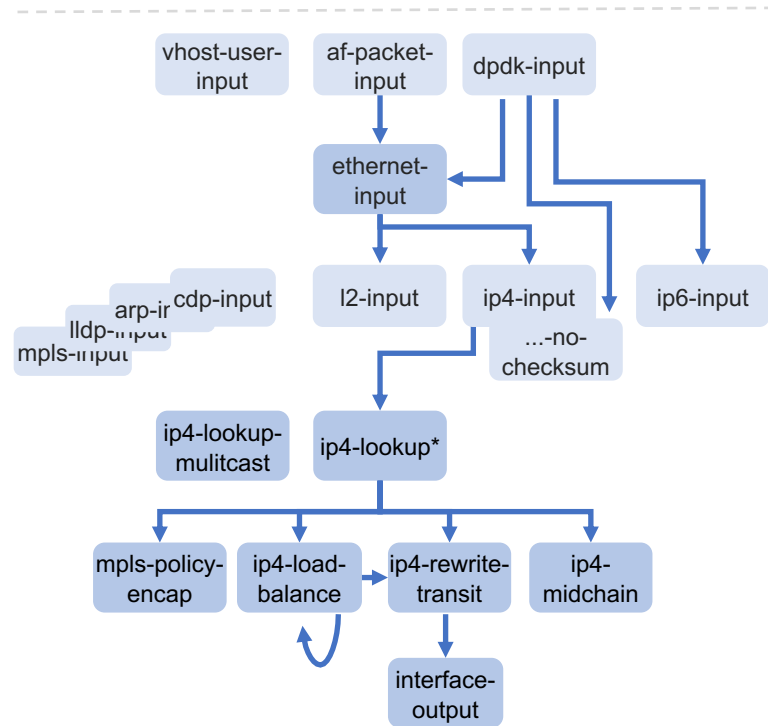


**Close to 10G wire rate performance achieved!**

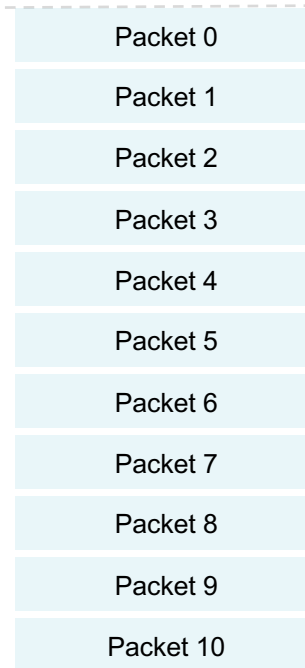
# FD.io VPP – The “Magic” of Vectors

Compute Optimized SW Network Platform

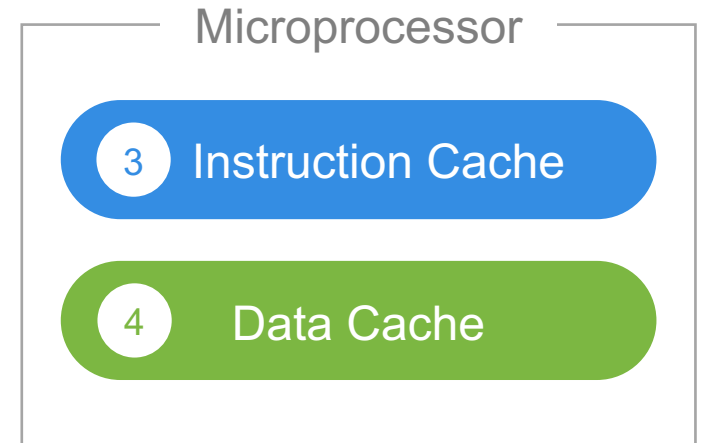
1 Packet processing is decomposed into a directed graph of nodes ...



2 ... packets move through graph nodes in vector ...



3 ... graph nodes are optimized to fit inside the instruction cache ...



4 ... packets are pre-fetched into the data cache.

\* Each graph node implements a “micro-NF”, a “micro-NetworkFunction” processing packets.

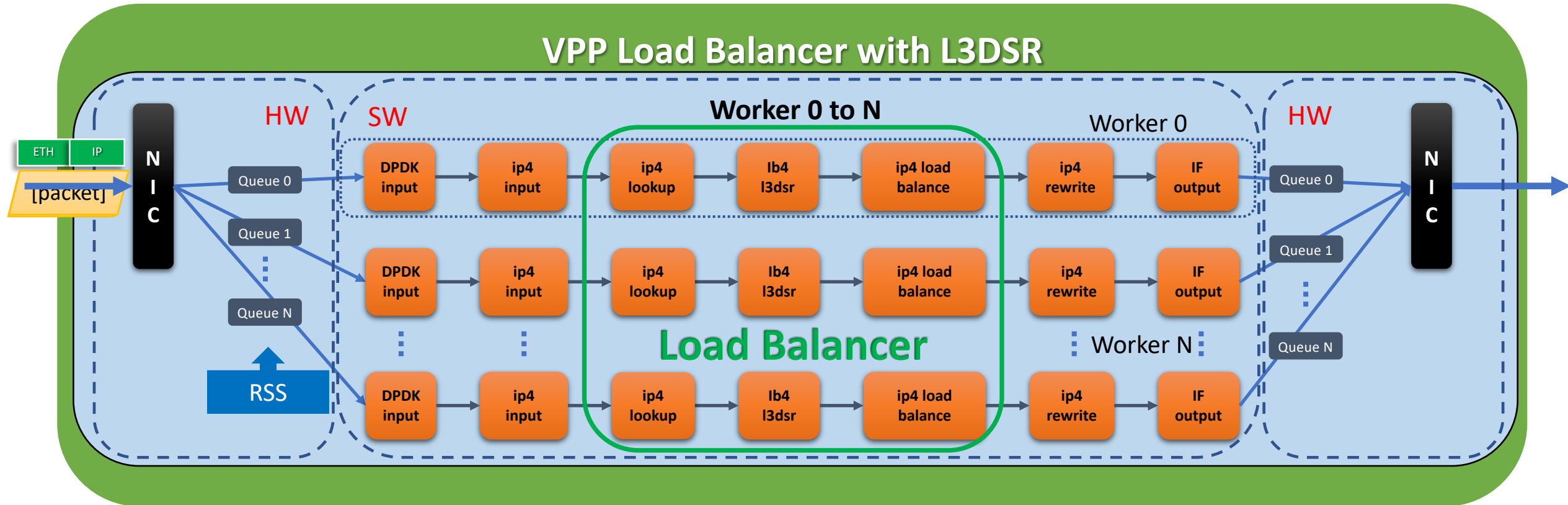
Makes use of modern Intel® Xeon® Processor micro-architectures.

Instruction cache & data cache always hot → Minimized memory latency and usage.



\*Other names and brands may be claimed as the property of others.

# Multi-threading for performance scalability



- RSS (Receive Side Scaling) enables traffic associated with one connection to a given thread.
- Well leveraging HW feature RSS and SW multi-threading for true scalability





ons  
NORTH AMERICA  
OPEN NETWORKING //  
Enabling Collaborative  
Development & Innovation

## Collaboration with VPP resulted in:

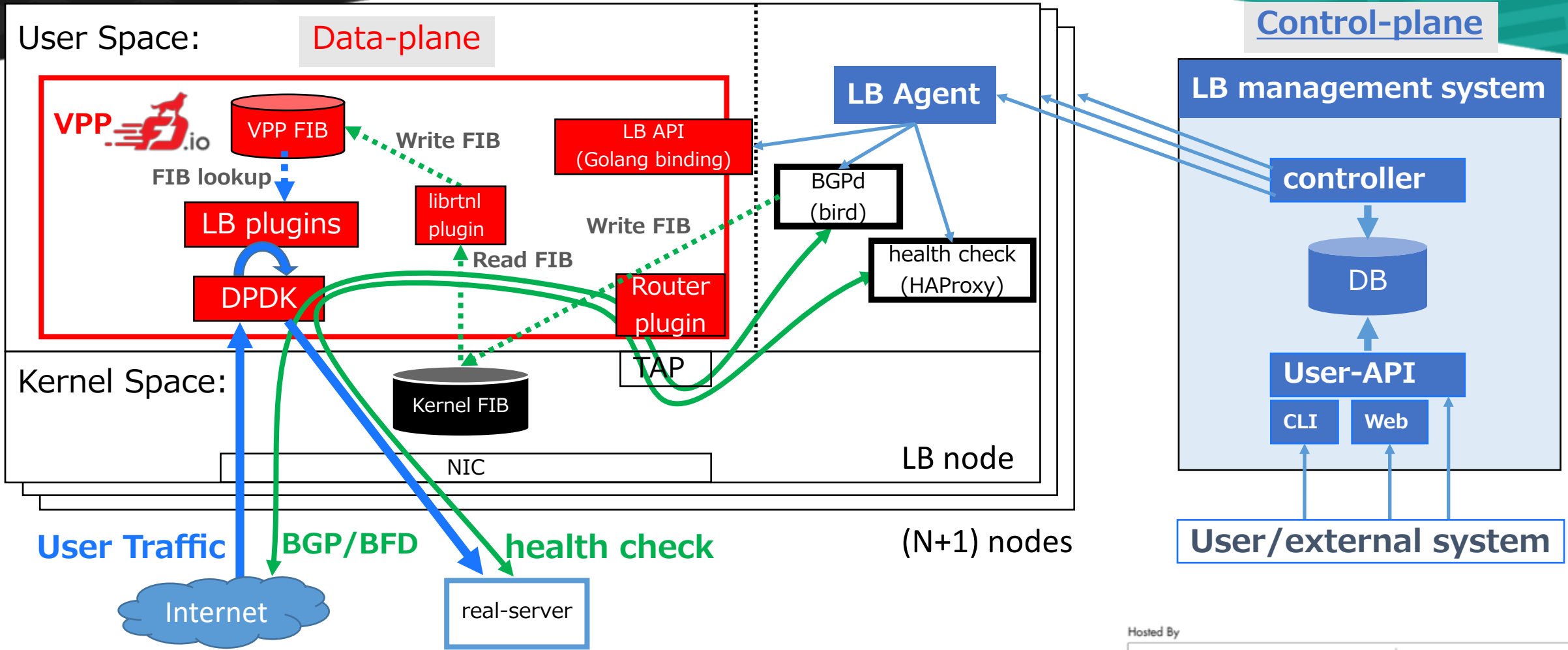
- **OSS based stateless L3DSR L4 Load Balancer**
- **Scaling-in/out under N+1 capability**
- **On top Clos with BGP robustness**
- **Sufficient performance**

Hosted By

 THE LINUX FOUNDATION |  LFNETWORKING



# Internal VPP LB node architecture and control-plane



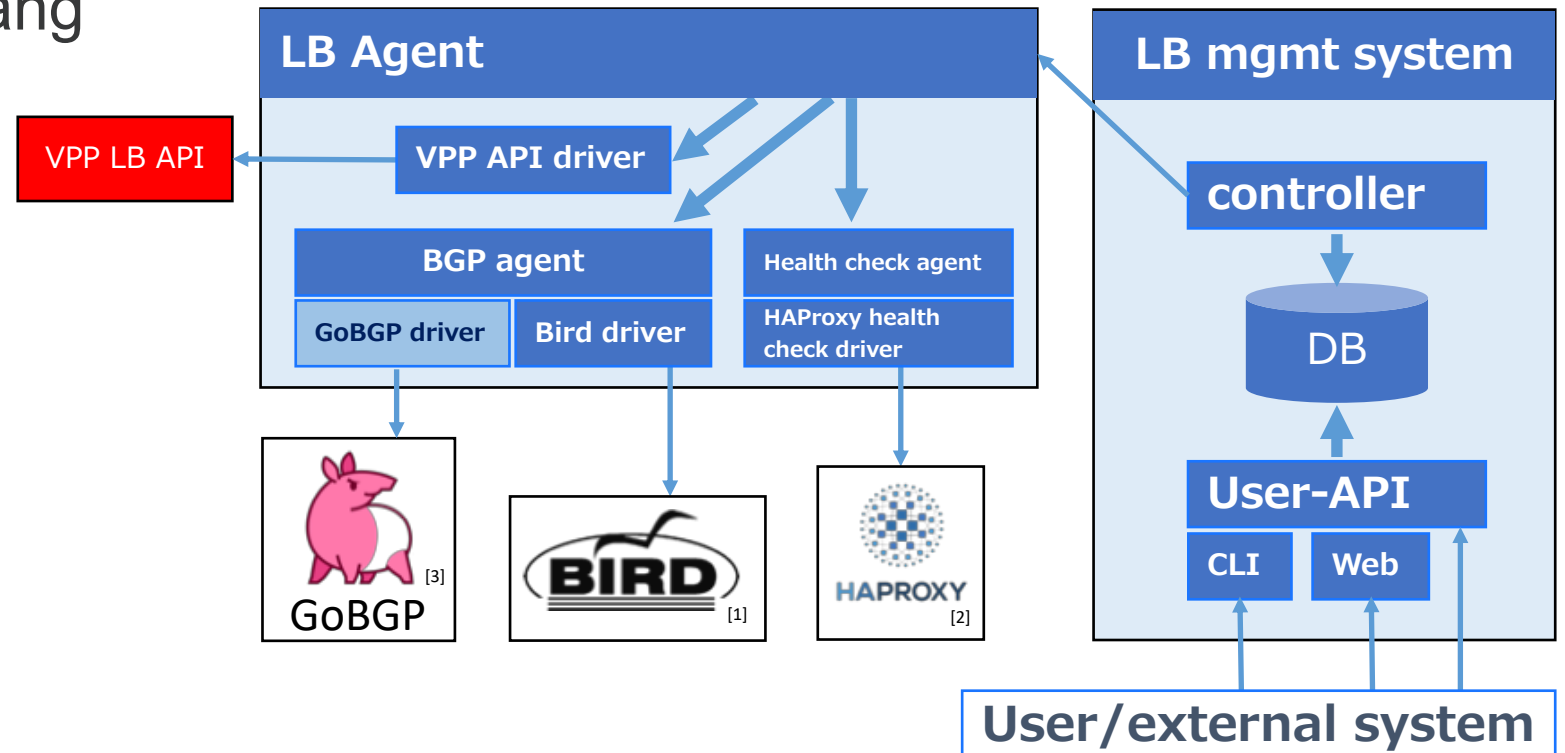
Hosted By

\*Other names and brands may be claimed as the property of others.



# LB control-plane

- Full-scratch written by Golang
  - 22k LoC with gRPC
- Components
  - LB agent and drivers
    - VPP API agent
    - BGP agent
    - Health check agent
  - LB management system
    - Controller & DB
    - User-API & Interface



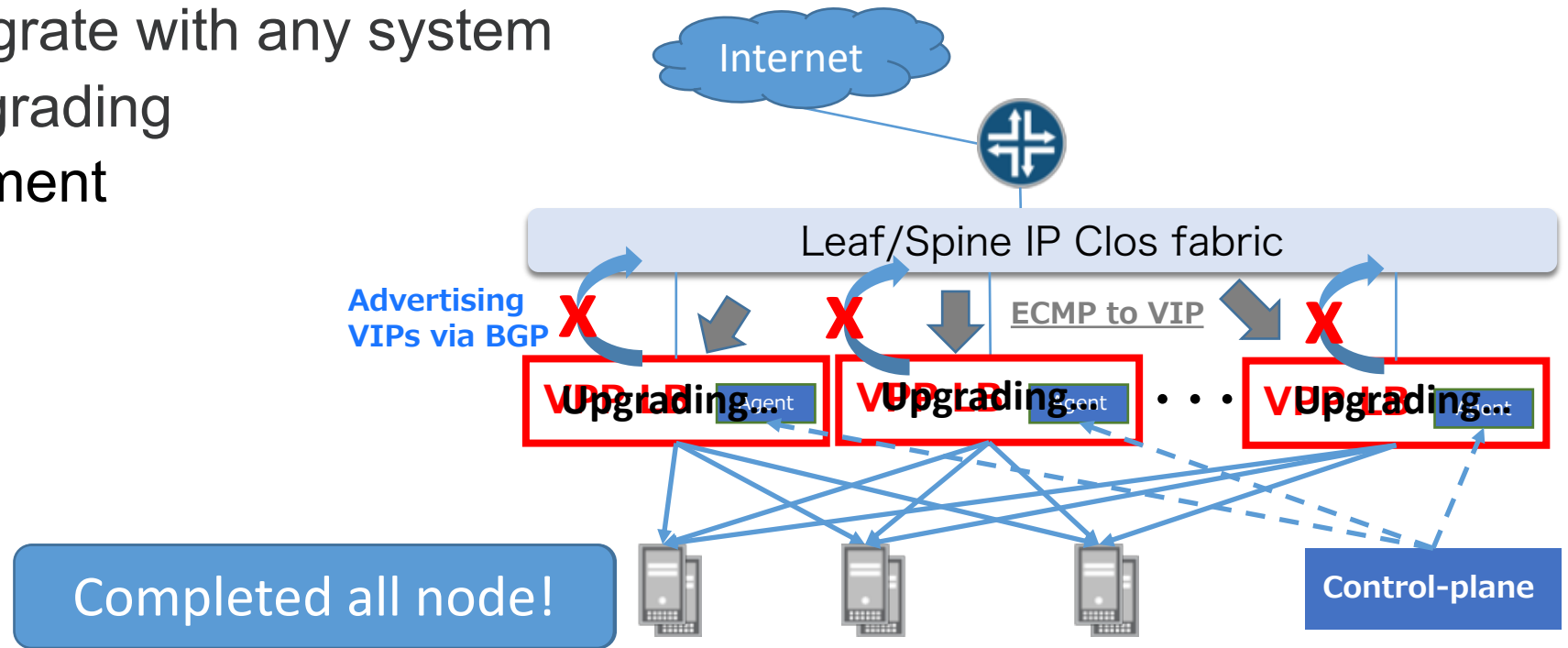
Hosted By



# LB control-plane (Cont.)

- Features

- Serving API to integrate with any system
- Zero-downtime upgrading
- Life cycle management
  - HW EoL migration

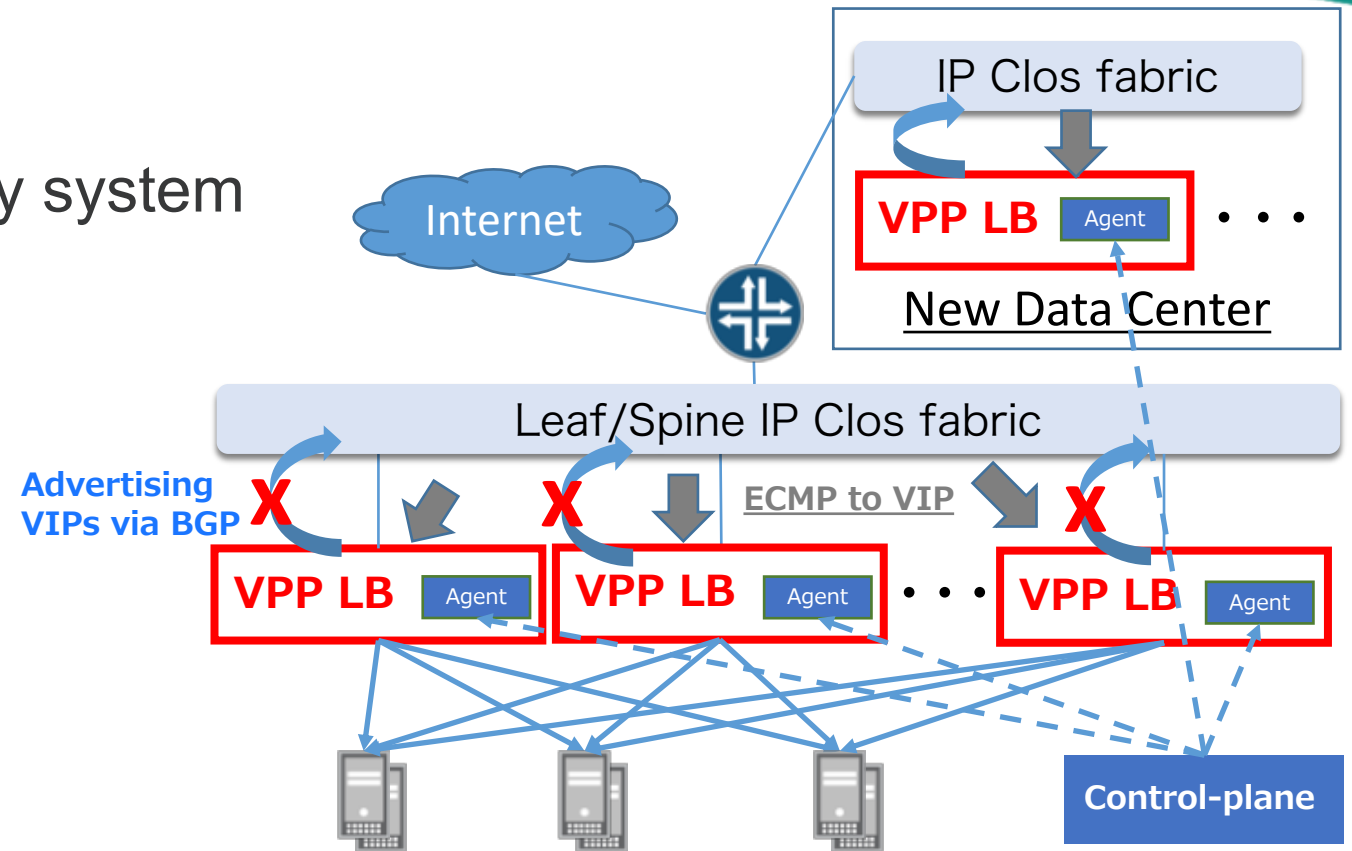


Hosted By



# LB control-plane (Cont.)

- Features
  - Serving API to integrate with any system
  - Zero-downtime upgrading
  - Life cycle management
    - HW EoL migration



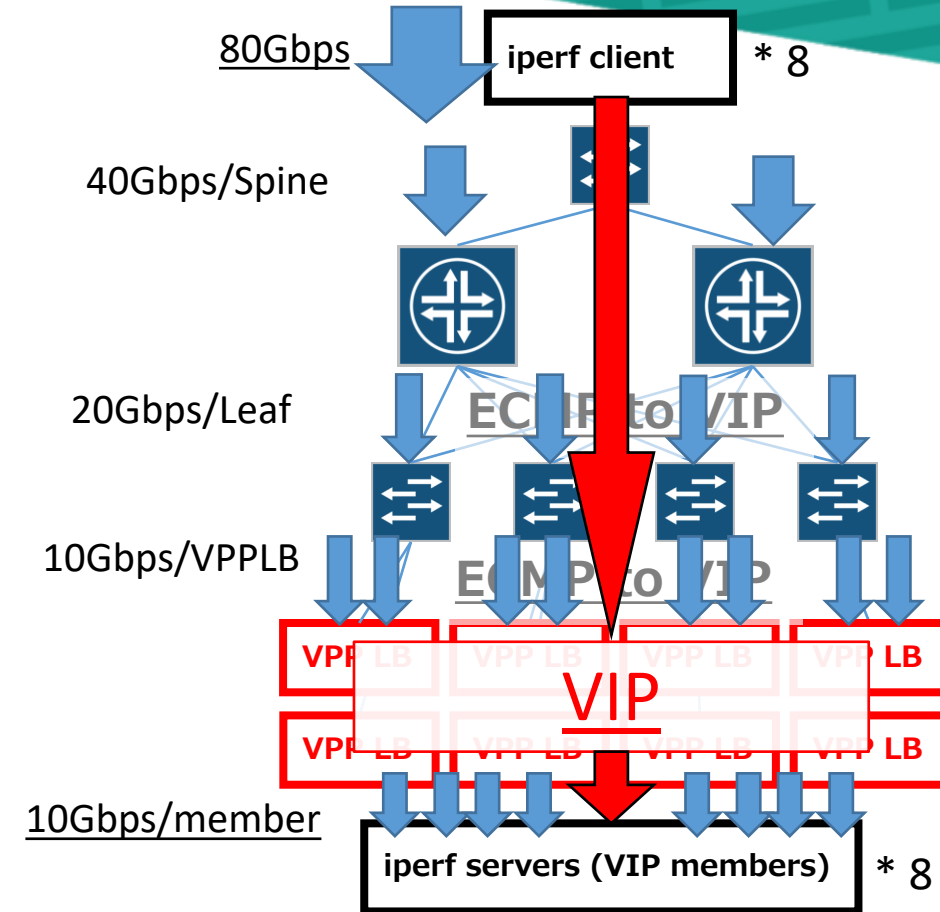
Hosted By



# LB system performance

## @ 80Gbps

- Setting
  - Traffic generator: iperf 10Gbps \* 8
  - 8 VPP-LB nodes: 10Gbps \* 8
  - ECMP
    - Spine
    - Leaf
    - VPP-LB
- 80Gbps to the 1 VIP on 8 VPP-LB



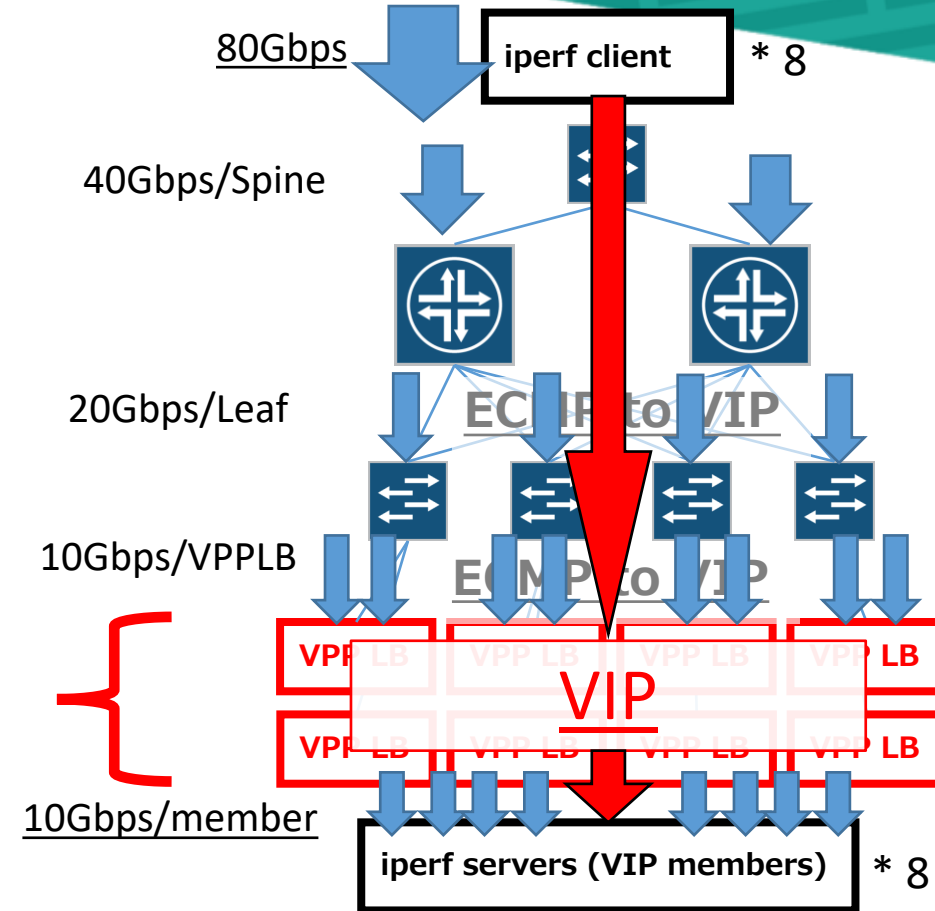
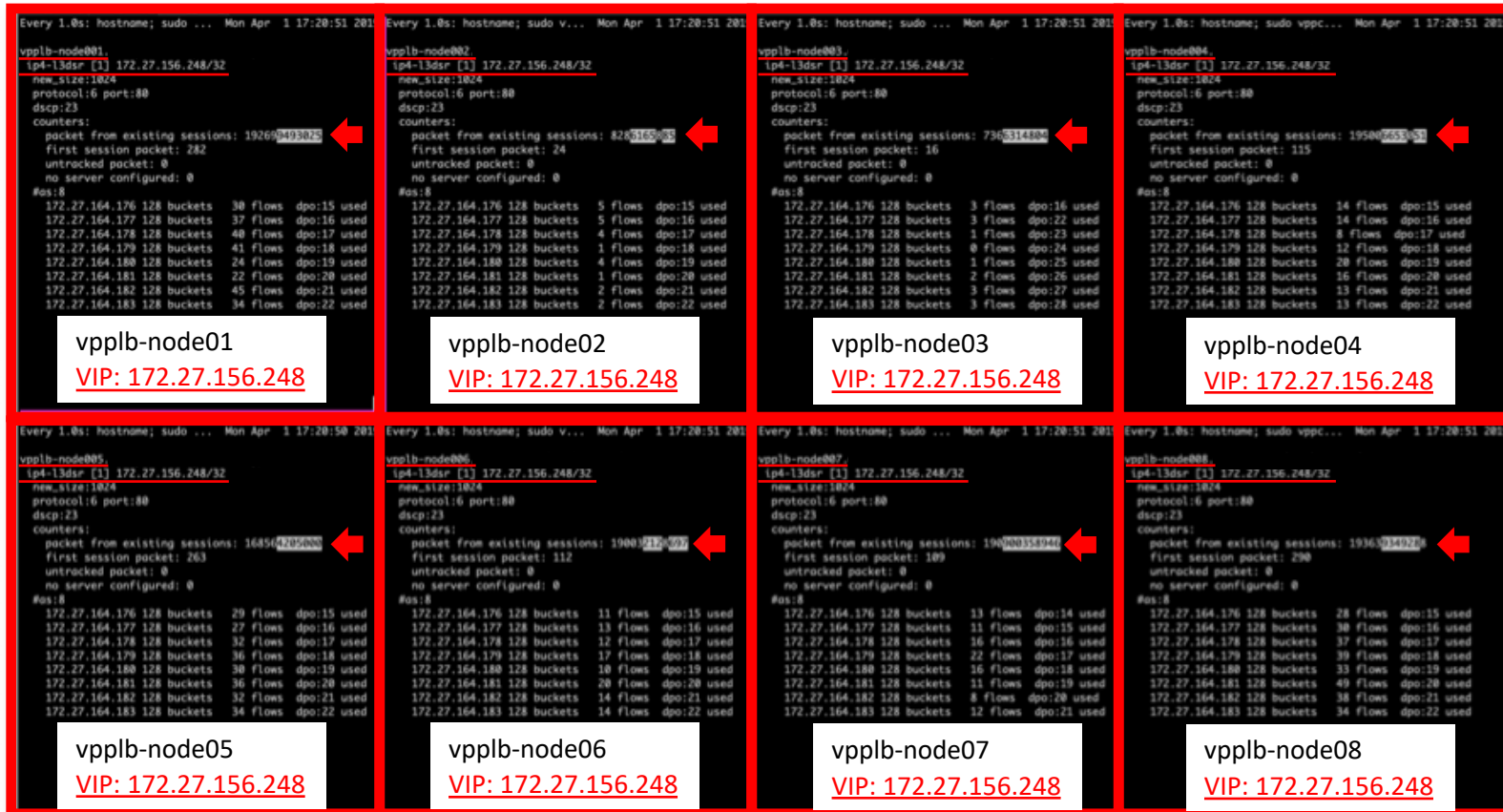
Hosted By



**ONS**  
 NORTH AMERICA  
 OPEN NETWORKING //  
 Enabling Collaborative  
 Development & Innovation

# LB system performance

## @ 80Gbps

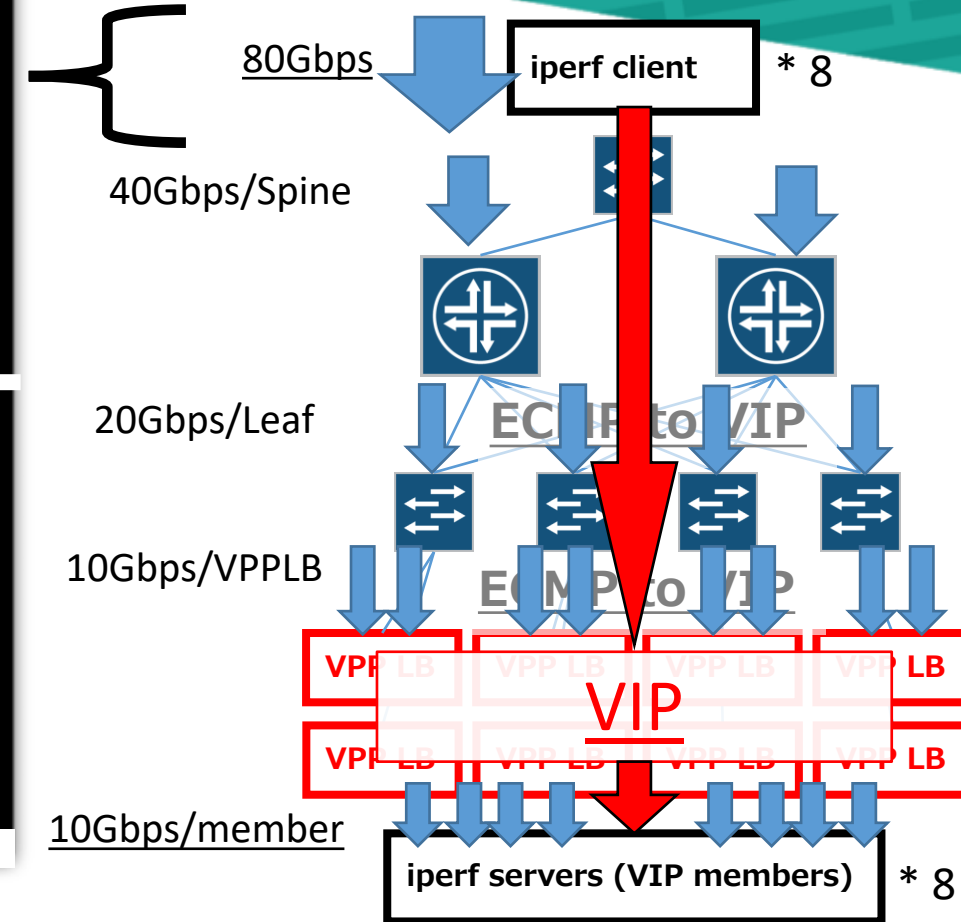
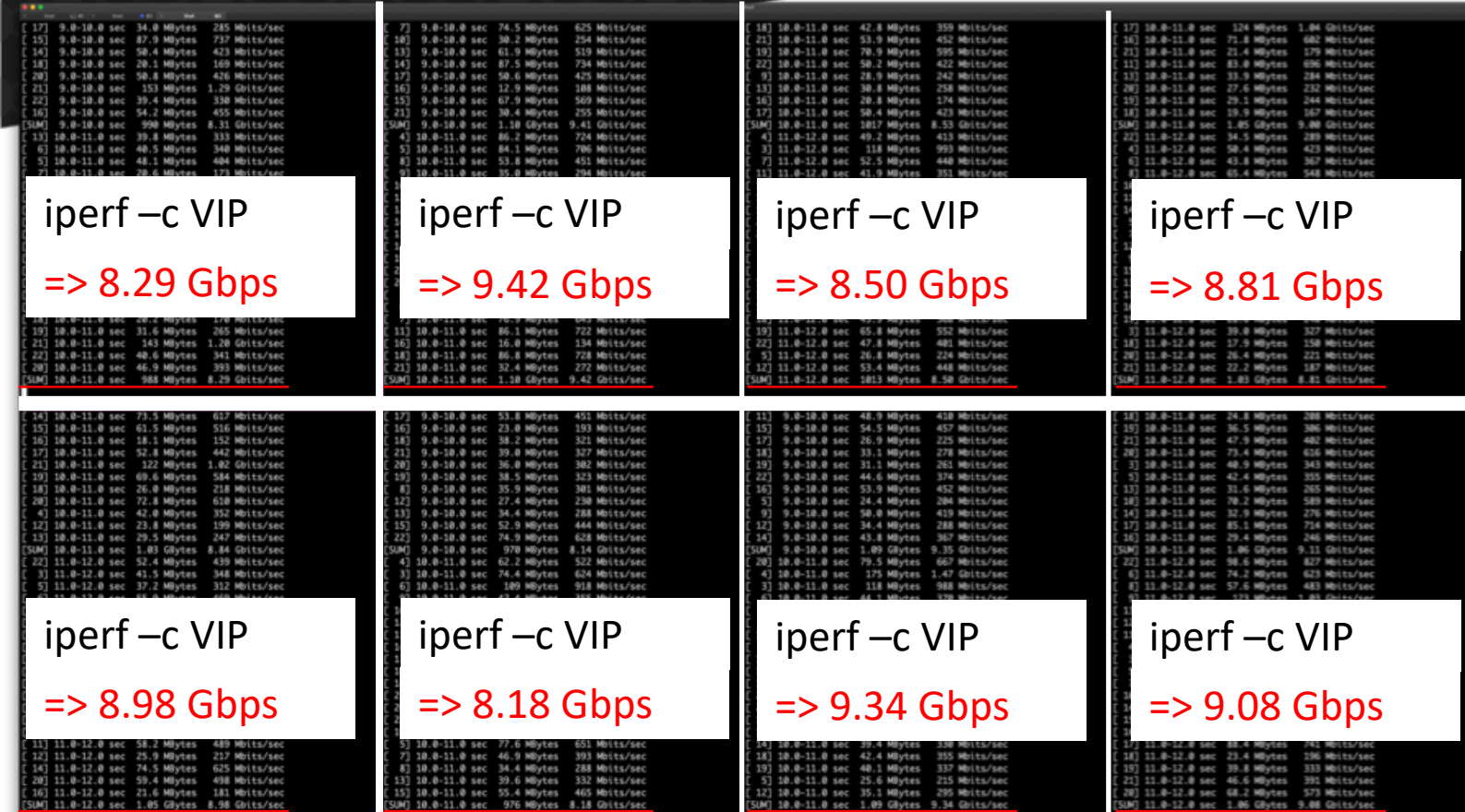




**ONS**  
 NORTH AMERICA  
 OPEN NETWORKING //  
 Enabling Collaborative  
 Development & Innovation

# LB system performance

## @ 80Gbps



Copyright (C) 2019 Yahoo Japan Corporation. All Rights Reserved.

80Gbps balancing via 8 VPP-LB on Top Clos/ECMP

# Scale-out on multiple VPP-LB node & ECMP on top Clos

\*Other names and brands may be claimed as the property of others.

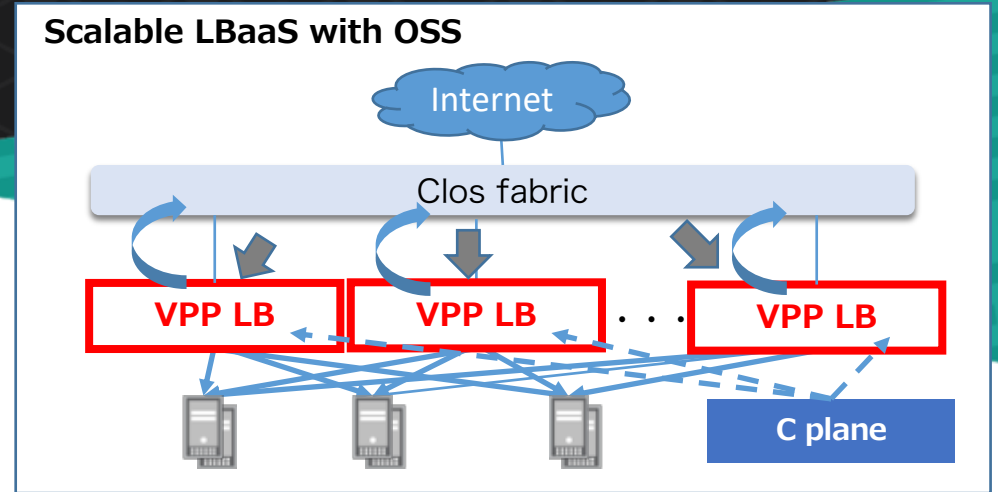




# LBaaS Summary:

## Coordination of our Control-plane with VPP data-plane

- Easy operation to place LB node with BGP
- OSS based stateless L3DSR L4 Load Balancer
- Serving API to integrate with any system
- Scaling-in/out under N+1 capability
- Zero-downtime upgrading
- On top Clos with BGP robustness
- Life cycle management (EoL migration)
- Sufficient performance in operation perspective



## Realization of LBaaS:

- **Scale-in/out LB capability**
- **Robustness of LB system**
- **Elastic management of VIP**

# Call to Action



- **FD.io: A Good open community resolving real-world networking issues and developing new features**
  - Quick response from community
  - Deepening networking knowledge in community
- **Collaboration**
  - Sharing your ideas/issues
  - Coding and verifying from your own perspective
  - Open discussion within the community

**Let's work together!**

- [ytatsumi@yahoo-corp.jp](mailto:ytatsumi@yahoo-corp.jp)
- [naoyuki.mori@intel.com](mailto:naoyuki.mori@intel.com)

# Legal Disclaimers

- Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors.

Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit [www.intel.com/benchmarks](http://www.intel.com/benchmarks).

- Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. No computer system can be absolutely secure. Check with your system manufacturer or retailer or learn more at [www.intel.com](http://www.intel.com).
- Cost reduction scenarios described are intended as examples of how a given Intel- based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.
- No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.
- Intel disclaims all express and implied warranties, including without limitation, the implied warranties of merchantability, fitness for a particular purpose, and non-infringement, as well as any warranty arising from course of performance, course of dealing, or usage in trade.
- All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest Intel product specifications and roadmaps.
- Intel, the Intel logo and Xeon are trademarks of Intel Corporation in the U.S. and/or other countries.
- \*Other names and brands may be claimed as the property of others.
- © 2018 Intel Corporation.





# ons

NORTH AMERICA

**OPEN NETWORKING //**  
Enabling Collaborative  
Development & Innovation

Hosted By

 THE **LINUX** FOUNDATION |  **OLF** NETWORKING



ons

NORTH AMERICA

OPEN NETWORKING //  
Enabling Collaborative  
Development & Innovation

# Backup

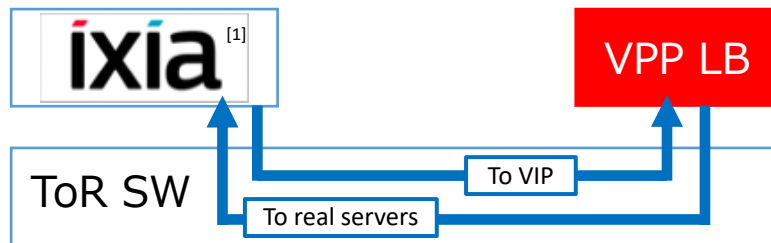
Hosted By

 THE **LINUX** FOUNDATION |  **LF** NETWORKING

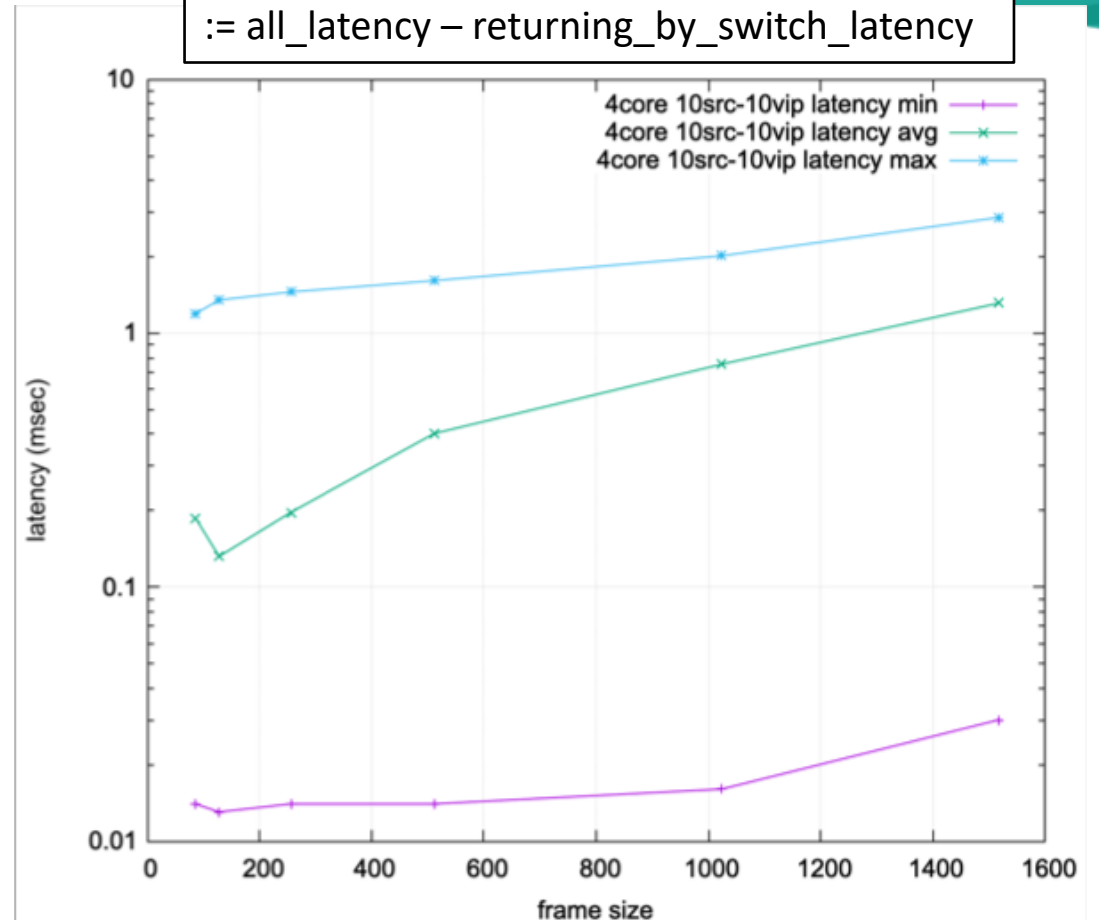


# VPP LB/node performance test @ 10Gbps

- Traffic generator: Ixia IxNetwork
- VPP Server Spec:
  - CPU: Intel® Xeon® Processor E5-2650L v3 \* 2S
  - Memory: 384GB
  - NIC: Intel® Ethernet Network Connection X540-AT2 (10Gbps\*2)



VPPLB Additional latency  
:= all\_latency - returning\_by\_switch\_latency



Copyright (C) 2019 Yahoo Japan Corporation. All Rights Reserved.



# L3DSR LB configuration example

- L3DSR LB (using IPv4 DSCP field) [1]
  - ✓ VPP LB L3DSR supported
- Management API
  - ✓ VPP LB APIs

## L3DSR LB configuration example:

```
lb vip 192.168.1.1/32 protocol tcp port 80 encap l3dsr dscp 23 new_len 1024
lb as 192.168.1.1/32 protocol tcp port 80 172.16.1.1
lb as 192.168.1.1/32 protocol tcp port 80 172.16.1.2
lb as 192.168.1.1/32 protocol tcp port 80 172.16.1.3 * Added by this project
```

