# Mellanox Overview



**1999**
Mellanox Founded

**$1.09B**
2018 Revenue

**~2,500**
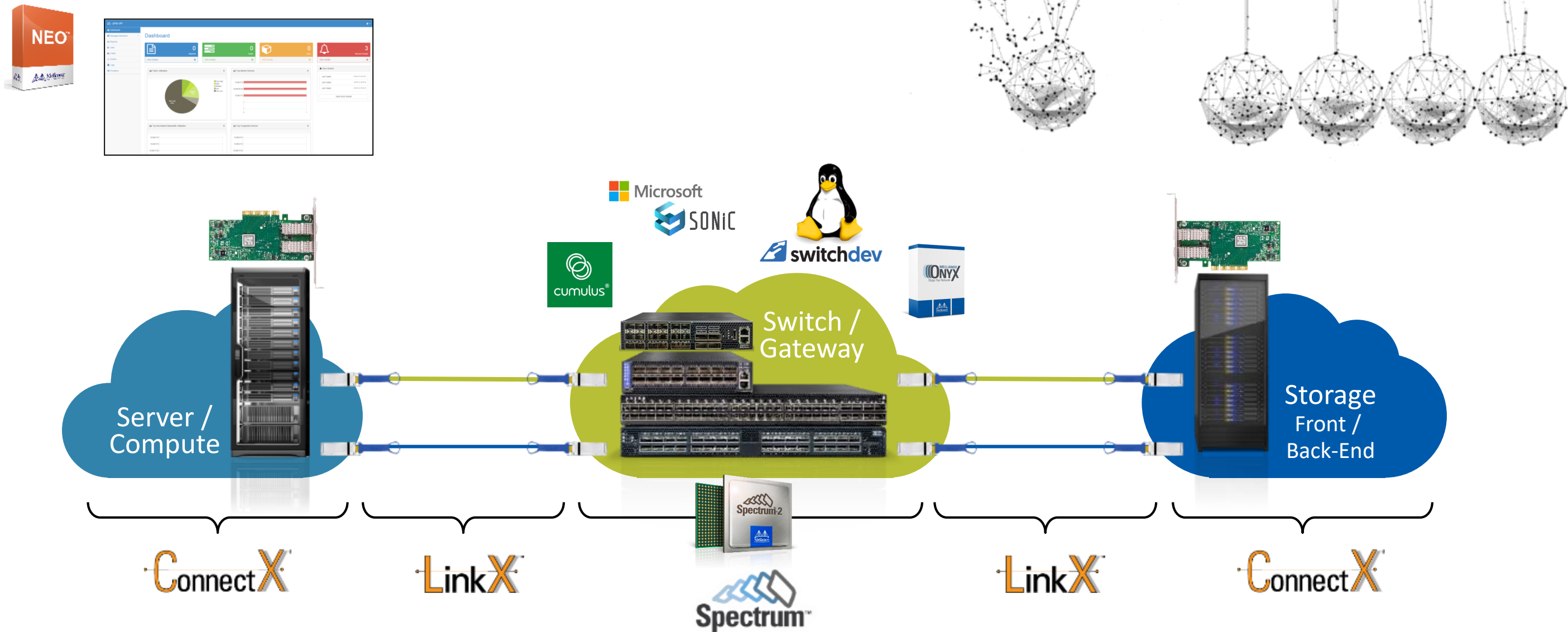Employees worldwide
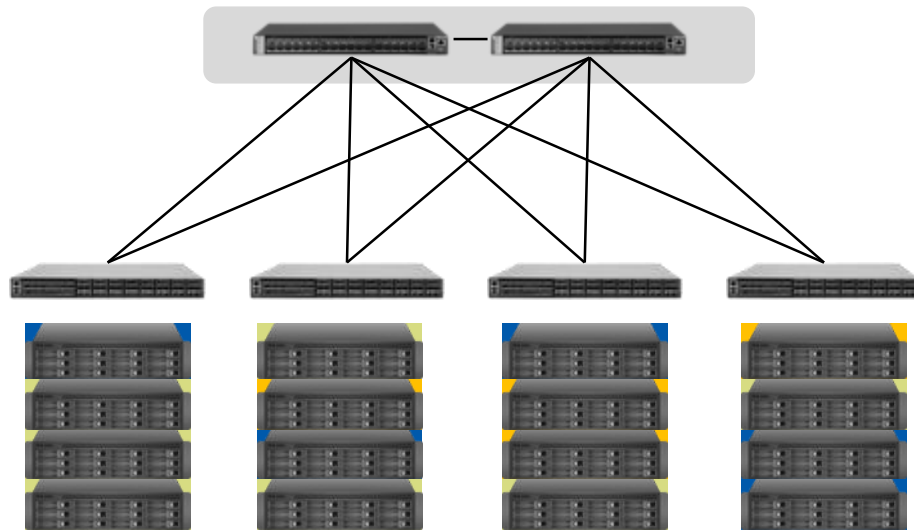
**NASDAQ®**

Ticker: MLNX

Company Headquarters:

- ◼ Yokneam, Israel
- ◼ Sunnyvale, California
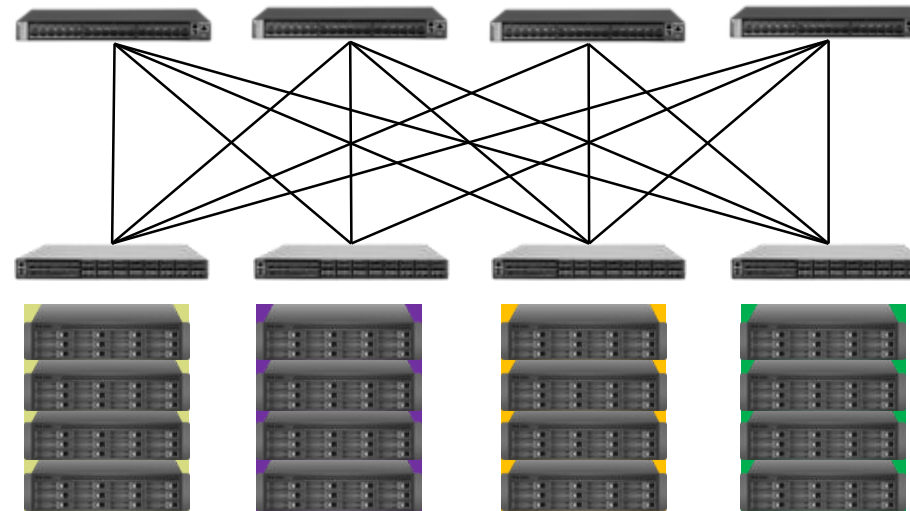- ◼ Worldwide Offices

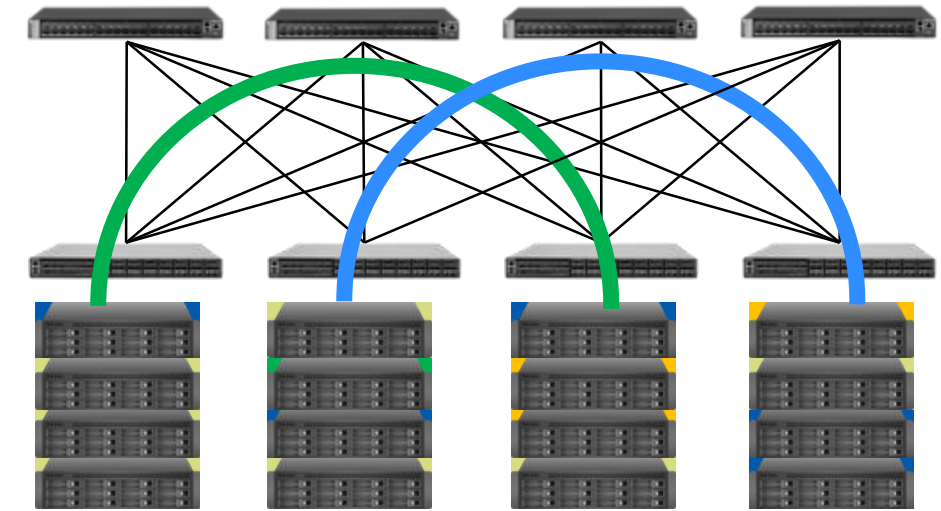# End-to-End Interconnect Solutions

# Leaf/Spine Deployments

**Layer 2 / MLAG**

**Layer 3 / ECMP**

**L2 over Layer 3 VXLAN**



Connect out via spines, L3 GW on spines or above

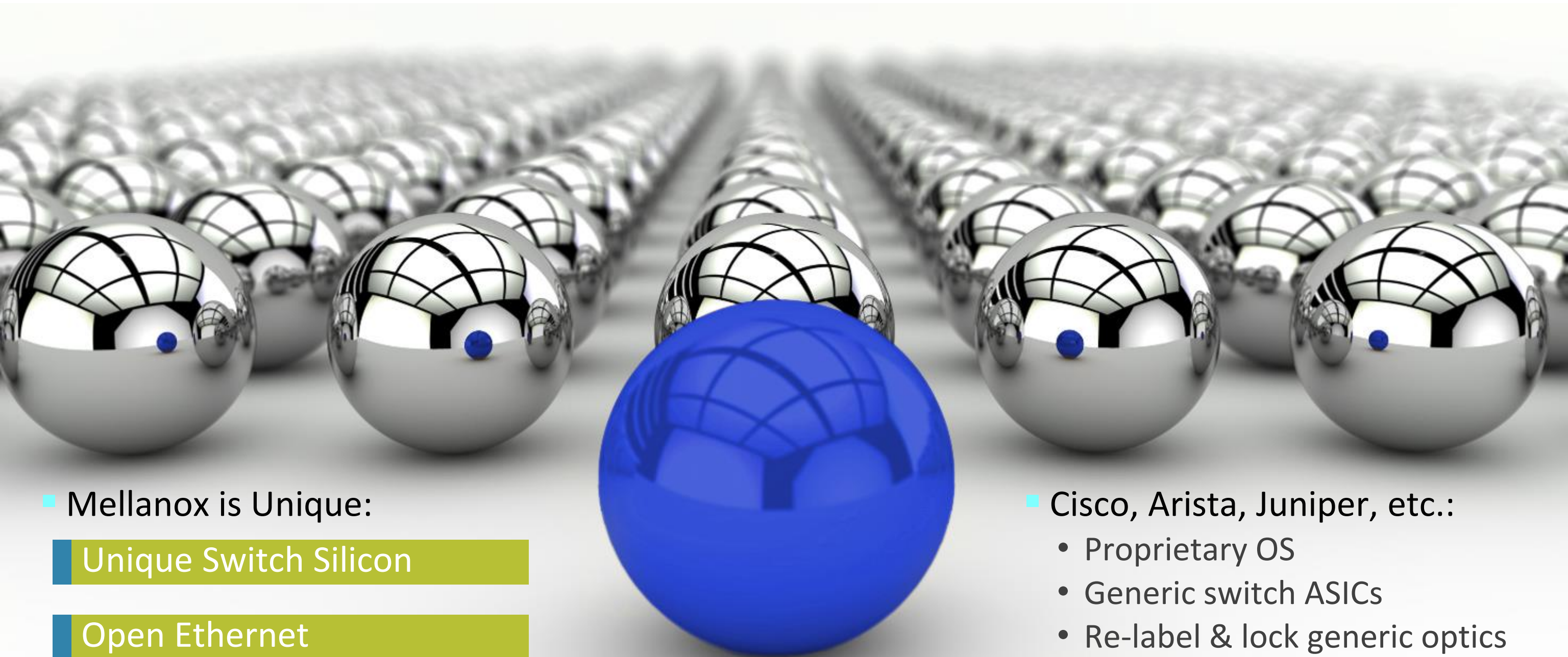**BGP from the Host –** Kuberentes
**VXLAN from the host-** VMware , OpenStack, Kubernetes

Anycast L3 GW on TORs or FW as GW located on the border leaf.
EVPEN Type 5 for Routes out of the fabric.

# Mellanox – Not Like Other Network Vendors

- Mellanox is Unique:

  Unique Switch Silicon
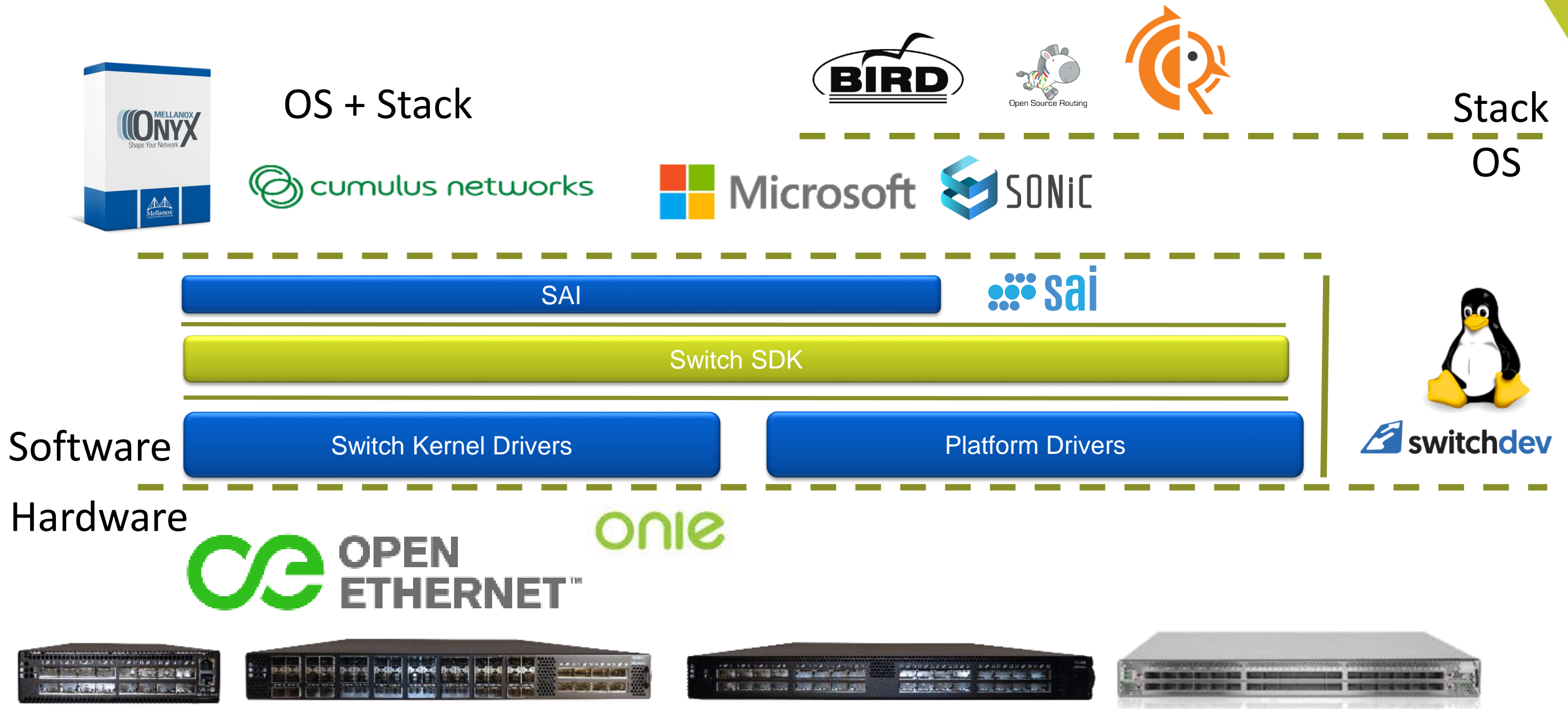
  Open Ethernet

  End to End

- Cisco, Arista, Juniper, etc.:
  - Proprietary OS
  - Generic switch ASICs
  - Re-label & lock generic optics
  - License Features

# We Are Building an Open Ecosystem

# Open Ethernet SN2000 Series

**SN2700 – 32x100GbE** (up to 64 x 50/25/10GbE)
The Ideal 100GbE ToR / Aggregation

**SN2410 – 8x100GbE + 48x25GbE**
25GbE ➜ 100GbE ToR

**SN2100 – 16x100GbE ports (64x25GbE)**
Ideal storage/Database Switch
Highest 25GbE Density per rack unit

**SN2010 – 18x10/25GbE + 4x40/100GbE**
Ideal HCI ToR Switch

- Predictable Performance
- Fair Traffic Distribution for Cloud
- Best-in-Class Throughput, Latency, Power Consumption
- Zero Packet Loss

## 300ns
SN2700 – 169W
SN2410 – 165W
SN2100 – 94W

Energy efficiency

**Spectrum**™

# Spectrum 2 - Open Ethernet SN3000 Series

**SN3700C – 32x100GbE (128x 1-25GbE)**
100GbE Spine/ToR

**SN3700 – 32x200GbE (128x 1-50GbE)**
200GbE Spine

**SN3800 – 64x100GbE**
Spine/Super Spine

**SN3510 – 48x25/50GbE + 6x400GbE**
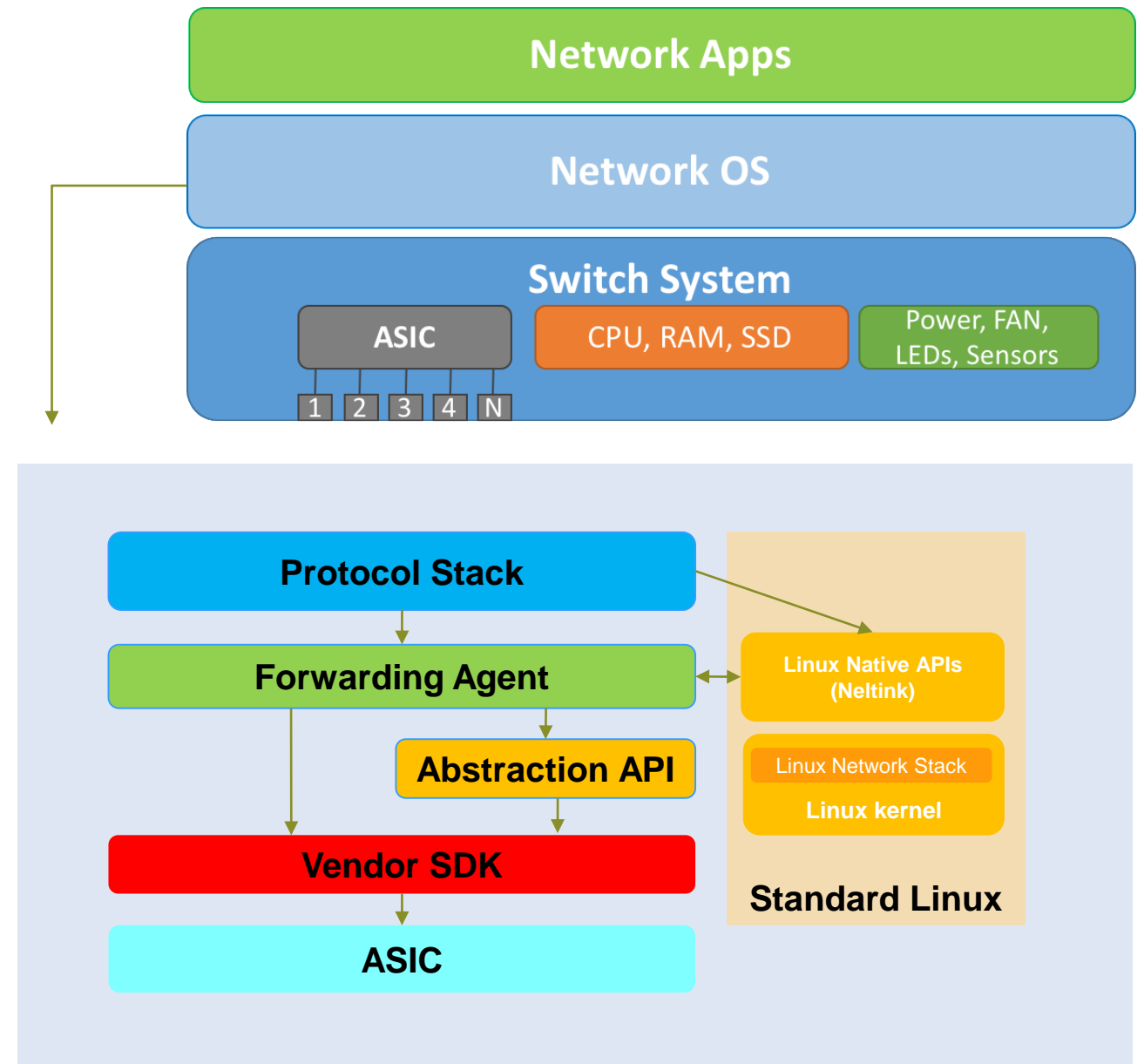25/50GbE ➜ 400GbE ToR
(Q2 2020)

- Best-in-Class Buffers
- Best-In-Class Virtualization
- Best-In-Class Telemetry
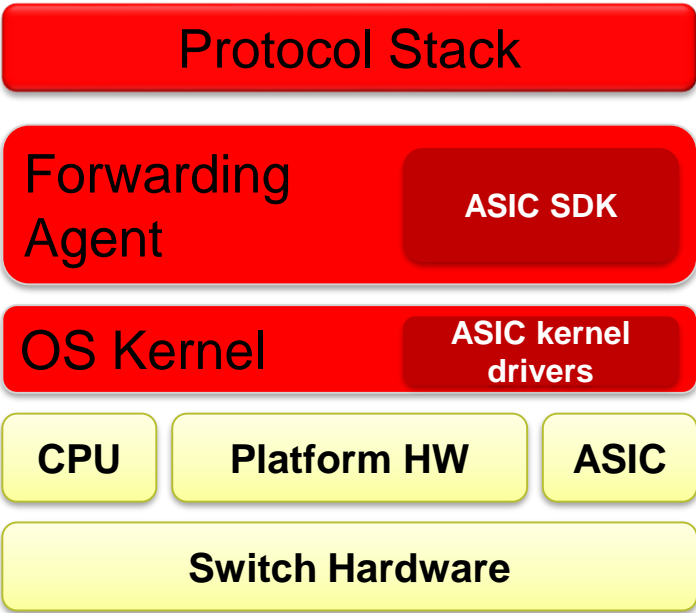
# Switch NOS Reference Architecture

- Protocol Stack
  - Network protocols (RIB)
    - Bridge, STP, OSPF, BGP.
  - Forming FIB our of RIB
- Forwarding Agent
  - Middleware between Protocol Stack & ASIC
  - Programming FIB into the ASIC (HW offload)
  - Uses special API to communicate with ASIC
- API
  - Proprietary ASIC vendor SDK
  - Broadcom, Mellanox, Cavium, Marvell, etc.
  - Trend to standardize and open SDK
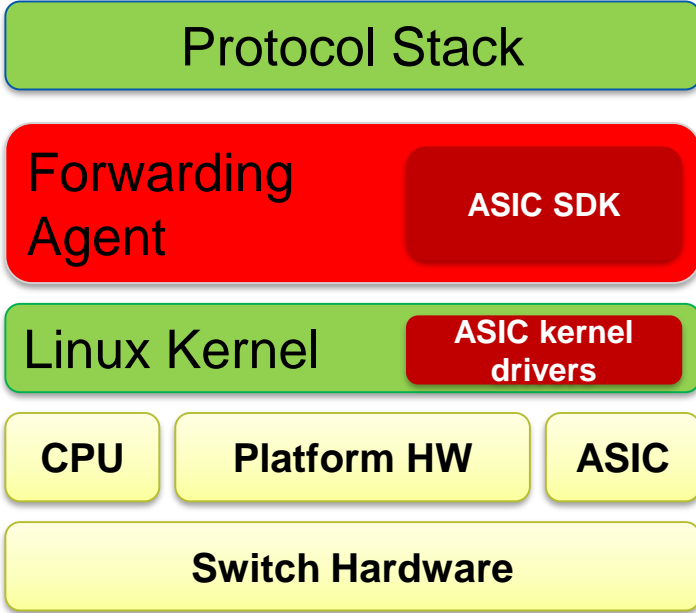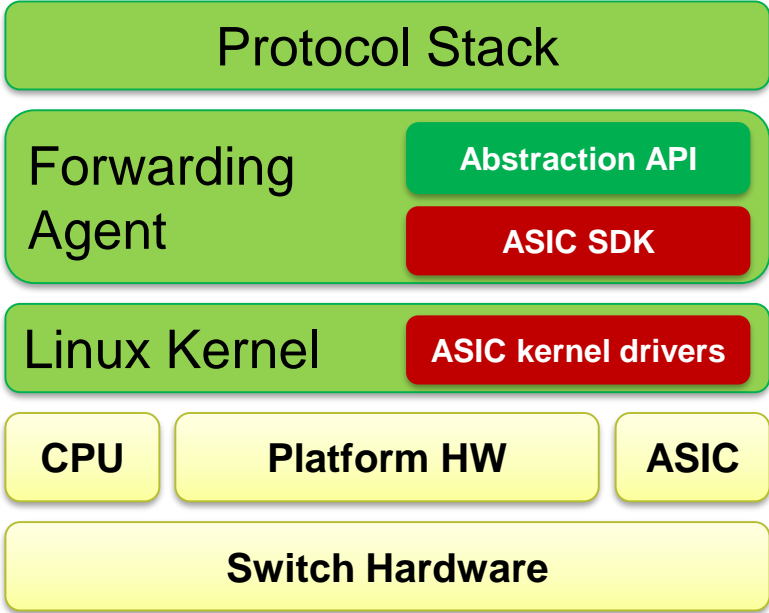  - SAI, OpenNSL, OF-DPA, P4, etc.



Network Apps

Network OS

Switch System — ASIC | 1 2 3 4 N | CPU, RAM, SSD | Power, FAN, LEDs, Sensors

Protocol Stack

Forwarding Agent

Abstraction API

Vendor SDK

ASIC

Linux Native APIs (Neltink)

Linux Network Stack

Linux kernel

Standard Linux

# Switch NOS architecture examples

**Traditional NOS**

**Commercial NOS**
based on Linux

**Open NOS**
based on Linux

| Protocol Stack | Protocol Stack | Protocol Stack |

**Traditional NOS:**

| Forwarding Agent | ASIC SDK |
| OS Kernel | ASIC kernel drivers |

| CPU | Platform HW | ASIC |
| Switch Hardware |

**Commercial NOS:**

| Forwarding Agent | ASIC SDK |
| Linux Kernel | ASIC kernel drivers |

| CPU | Platform HW | ASIC |
| Switch Hardware |

**Open NOS:**

| Forwarding Agent | Abstraction API / ASIC SDK |
| Linux Kernel | ASIC kernel drivers |

| CPU | Platform HW | ASIC |
| Switch Hardware |

Open/Free SW    Proprietary SW    Binary blob

# MLNX-OS traditional Industry like OS



- GUI
- CLI
- SNMP

# Cumulus – Mellanox Partnership

App App App

Network OS

Open Hardware

Cisco
Arista
Juniper

splunk>
elastic
Cumulus NetQ
influxdata
ANSIBLE
puppet

Cumulus Linux
Network Operating System

Best Hardware

① **Economical scalability**
With commodity hardware and a standardized Linux stack, achieving a lower TCO by up to 60%

② **Built for the automation age**
Making networking repeatable and consistent

③ **Standardized toolsets**
Easily enable Linux tools: automation, monitoring, analytics…

④ **Choice and flexibility**
50+ hardware platforms, from 11 vendors, and 2 silicon

# Cumulus

- Cumulus Linux is a Debian Jessie based , much lighter in the size

- Linux is a Linux, you can google it, you can use man command to learn about the different commands

- Except of the code that controls the silicon , all are open source and additions that are pushed to the upstream (in process)

- Cumulus are providing systems wrappers commands for a better user interface
- NCLU



```
cumulus@tor-11[~]# net ?
    abort    :   abandon changes since last commit
    add      :   add a configuration line
    clear    :   clear counters, BGP neighbors, etc
    commit   :   save pending changes
    del      :   delete a configuration line
    help     :   show this screen and exit
    pending  :   view pending changes
    show     :   show command output
```

# SONiC - Software for Open Networking in the Cloud

- SONiC is a collection of software packages installed on Linux running on a network hardware switch which make it a complete, functional router targeted at data center networks. Runs on Debian 8 'Jessie' distribution.

- SONiC is supported by the community and all code is shared in public github
  https://github.com/Azure/SONiC
  https://github.com/Azure/SONiC/wiki/Architecture

- SONiC deployment
  - in Microsoft production datacenters today and in Mellanox IT
  - Alibaba is planning soon as well
  - EMEA

# SONiC Architecture cont.

- The Switch State Service (SwSS) is a collection of software that provides a database interface for communication with and state representation of network applications and network switch hardware.

# SONiC participants (From Azure blog)

## Mellanox is the ONLY vendor to contribute at all levels

# Mellanox Contribution



MAC Aging

Performance on Redis DB

Port Split

Everflow enhancement

WHAT JUST HAPPENED

FRR as default routing stack

PMON Refactoring

DHCP Relay

PFC WD

BGP-EVPN support (type 5)

Upgrade docker engine to 18.09

COPP

LAG, LLDP, QoS, COS,

Fast Reload
[Microsoft & Mellanox]

Management
• AAA
• TACAS

Warm Reboot

Management
• SNMP
• NTP

Config DB and basic SONiC CLI

New Platform APIs
*For Mellanox platforms

Egress ACL

Mirroring

Tunnel decap

Basic L2, L3 Quagga

L3 RIF counter support

IPinIP v6 for IPv4 and IPv6

ACL

Asymmetric PFC

Virtual path for streaming telemetry

Upgrade each docker to stretch version

Debian 9 and Kernel 4.9 upgrade

Routing Stack Graceful Restart

Incremental Config (IP, LAG, Port admin)

Transceiver parameter tuning

Basic VRF
Microsoft & Mellanox]

L3 VXLAN
Microsoft & Mellanox]

IPinIP v6 including all v4 and v6 combinations

# HW offload without SDK = Linux Switch

**Open NOS based on Linux**

**Open Linux as a NOS**

| Protocol Stack |
|---|

| Forwarding Agents | Abstraction APIs |
|---|---|
| | Proprietary ASIC SDK |

| Linux Kernel | ASIC kernel drivers |
|---|---|

| CPU | Platform HW | ASIC |
|---|---|---|

| Switch Hardware |
|---|

**Open source**

| Protocol Stack |
|---|

*Netlink*

| Linux Kernel | Linux Network Stack |
|---|---|
| | Open ASIC kernel driver |

| CPU | Platform HW | ASIC |
|---|---|---|

| Switch Hardware |
|---|

# Linux Switch architecture



**User Space**

- Linux Management Applications
- Linux Routing Applications
- Linux Operating System

**Kernel Space**

- Linux Kernel w/ Switchdev Driver

**Hardware**

# Available today [Linux Kernel 5.0] – feature list

- **Visibility and Maintainability**
  - [ER]SPAN
  - Temperature
  - Fans
  - LED Control
  - ethtool (port counter, FW version, transceiver data)
  - Resource queries
  - RIF counters

- **Protocols**
  - Bridge - 802.1D
  - VLAN   - 802.1Q
  - LAG
  - LLDP
  - IGMP snooping
  - Unicast IPv4/IPv6 router
  - ECMP
  - DCB
  - QoS
  - IGMP flood control
  - sFlow
  - 256 VRFs
  - GRE tunnelling
  - Multicast IPv4/IPv6 router
  - IPv4/IPv6 weighted ECMP
  - VRRP
  - VxLAN

- **ACL**
  - tc-flower offload
  - Actions: Drop, Forward, Counters, Trap, TC_ACT_OK
  - TC chain template
  - Keys: Port, DMAC, SMAC, Ethertype, IP proto, SIP DIP (IPv4/6), TCP/UDP, L4 port, VLAN-ID, PCP, DCSP, VLAN valid, TCP flags

- **Misc**
  - 'devlink' tool
  - Port splitter
  - Shared buffer configuration
  - Internal secured FW upgrade
  - ECN: RED and PRIO

# Open NOS in production?
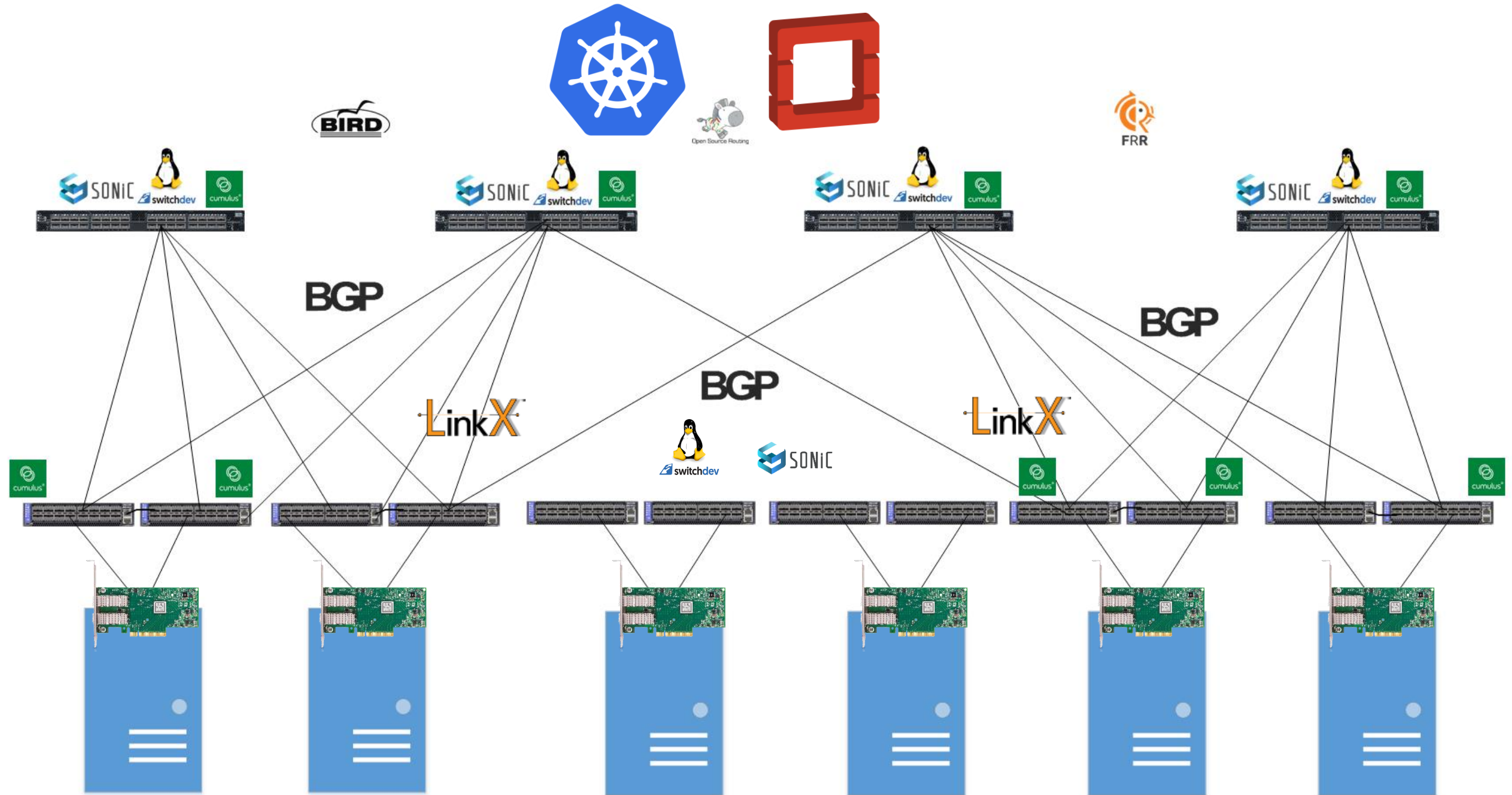
- **Why not?**
  - Microsoft, Alibaba - SONiC
  - Facebook - FBOSS
  - Ngenix (CDN) - Linux/switchdev
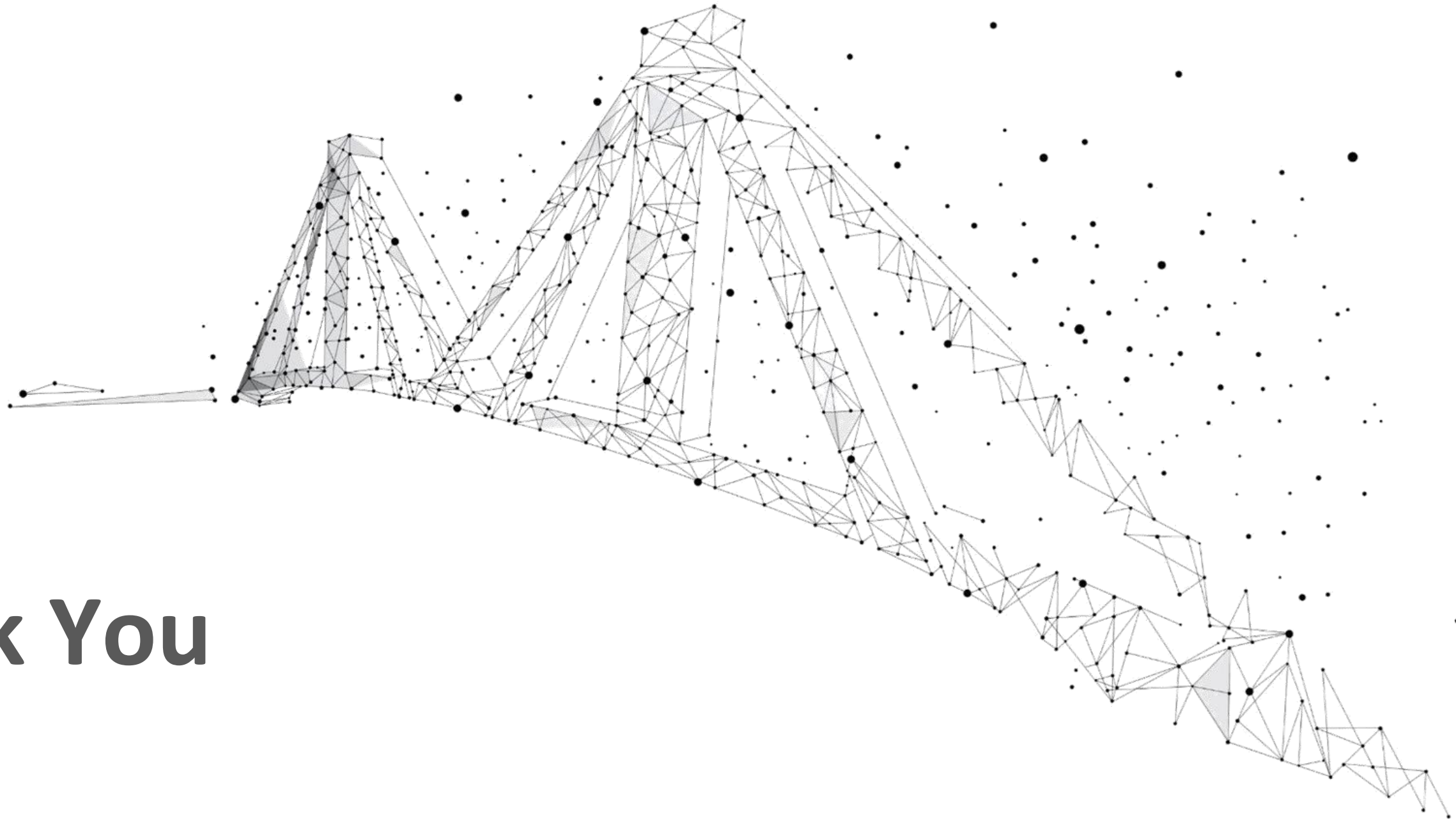  - Russia Biggest Bank (Cloud) - Linux/siwtchdev
  - …

- **Vendor Support**
  - SONiC - Mellanox, Dell, Edge-Core, Arista, *Cisco*, … 20+
  - Linux/switchdev - Linux, Mellanox, Cumulus, ALT …

- **Feature Richness**
  - L2 bridging
    - Linux bridge, MSTPd
  - L3 routing
    - Open OSPF/BGP implementations - 20+ yrs
    - Quagga/FRR/Bird/…
  - Tunneling - GRE/VXLAN
    - Linux kernel, OVS
    - EVPN - FRR/GoBGP
  - Security/Isolation - ACL, VRF
    - iptables, Linux TC, Linux NS/VRF
  - Management, monitoring
    - SNMPd, hsFlowd, Grafana

# Thank You