BISDN GmbH

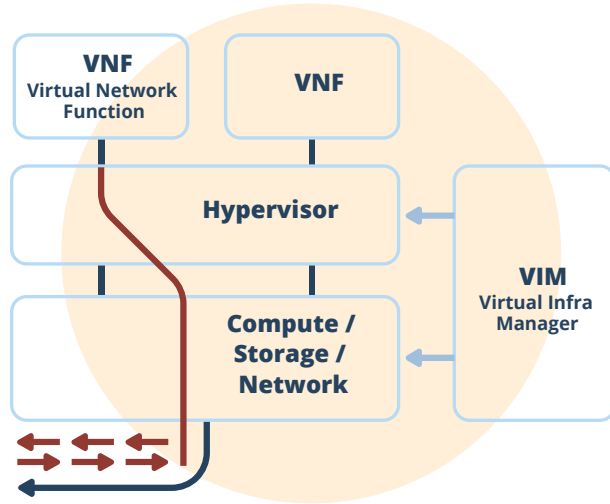# Open Source vBNG Architecture and Performance Evaluation

Hagen Woesner

ONF Broadband Spotlight Presentation, June 2020

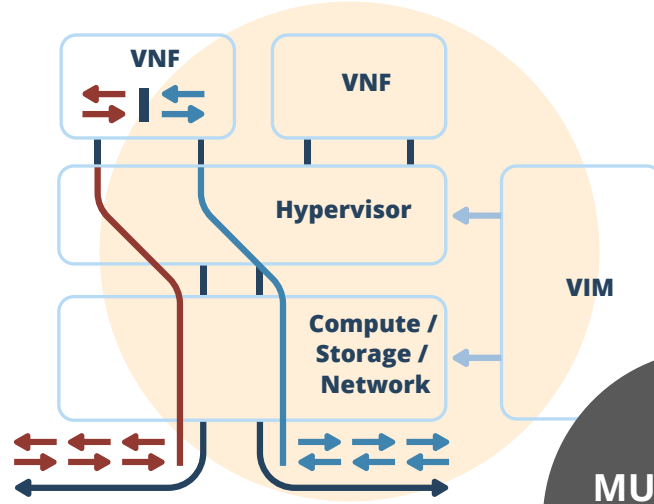# Challenge: Traditional clouds are unsuited for telco

At least twice the bandwidth required for telco services



Traditional cloud: **End of the wire**

VNF
Virtual Network Function

VNF

Hypervisor

VIM
Virtual Infra Manager

Compute / Storage / Network

Most cloud applications are web services

Telco cloud: **Bump in the wire**

VNF

VNF
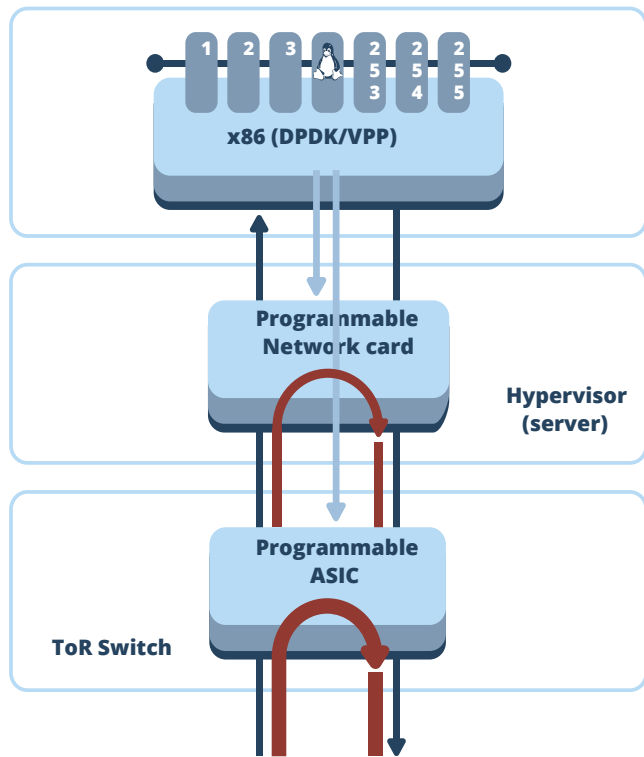
Hypervisor

VIM

Compute / Storage / Network

Most telco applications are traffic filters or gateways

**MUCH higher throughput required!**

# Solution: Offload packet handling to infrastructure

SDN required to 'remote control' ASICs and FPGAs

**x86 (DPDK/VPP)**

**Programmable Network card**

**Hypervisor (server)**

**Programmable ASIC**

**ToR Switch**

## Control/User Plane Separation (CUPS)

- Decompose network functions to become ‚cloud native'

## Common user plane abstraction: Linux pipeline

- baseboxd SDN controller translates Linux pipeline to whitebox
- No separate north-bound interface to controller – no SDN "apps"
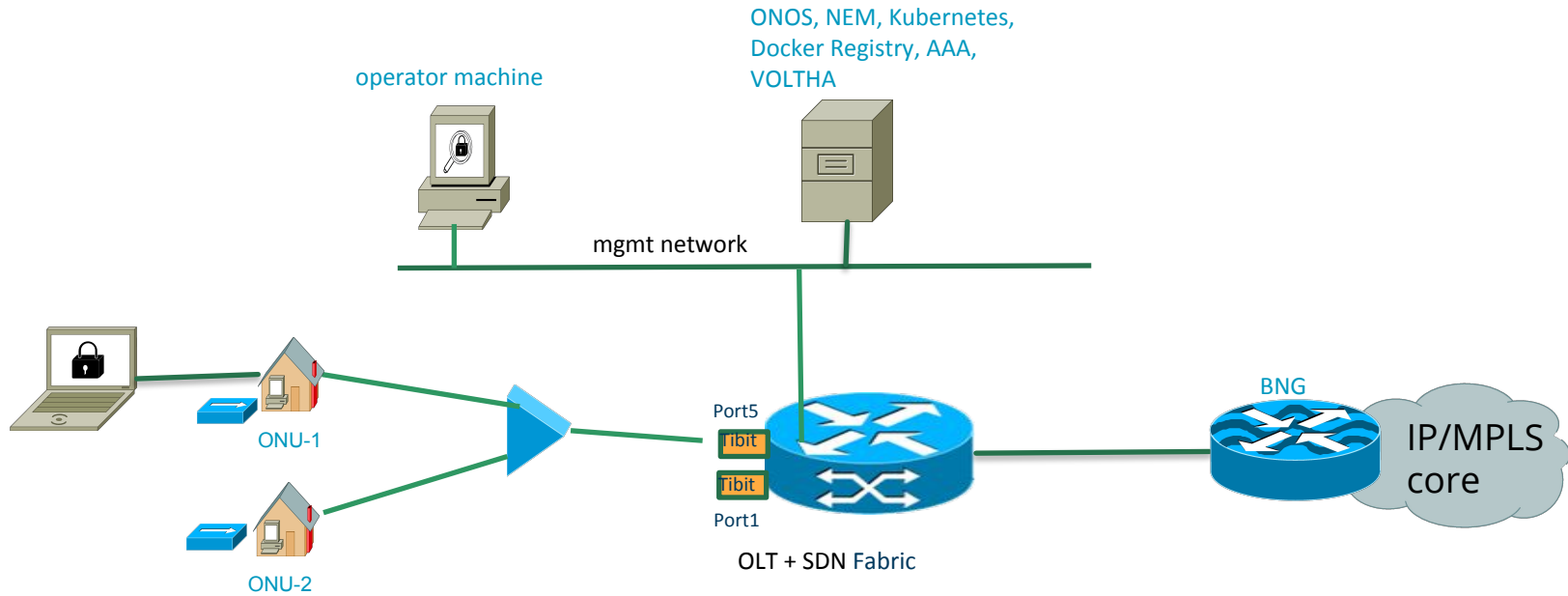- Re-use of millions of lines of Linux code

## Hardware acceleration (up to factor 1.000)

- Step 1: DPDK* inside VNF – factor 10
- Step 2: Programmable NIC on server – factor 10
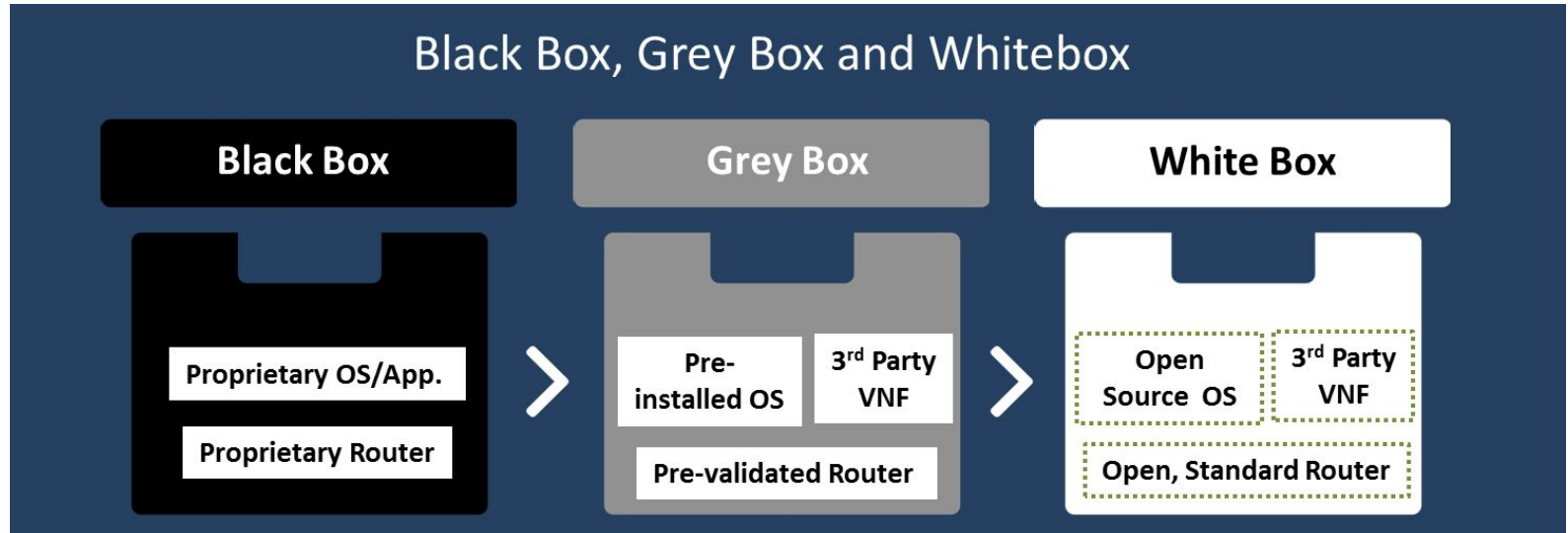- Step 3: Whitebox switch – factor 10

* DPDK = Data Plane Developer Kit – Intel's SDK for forwarding acceleration

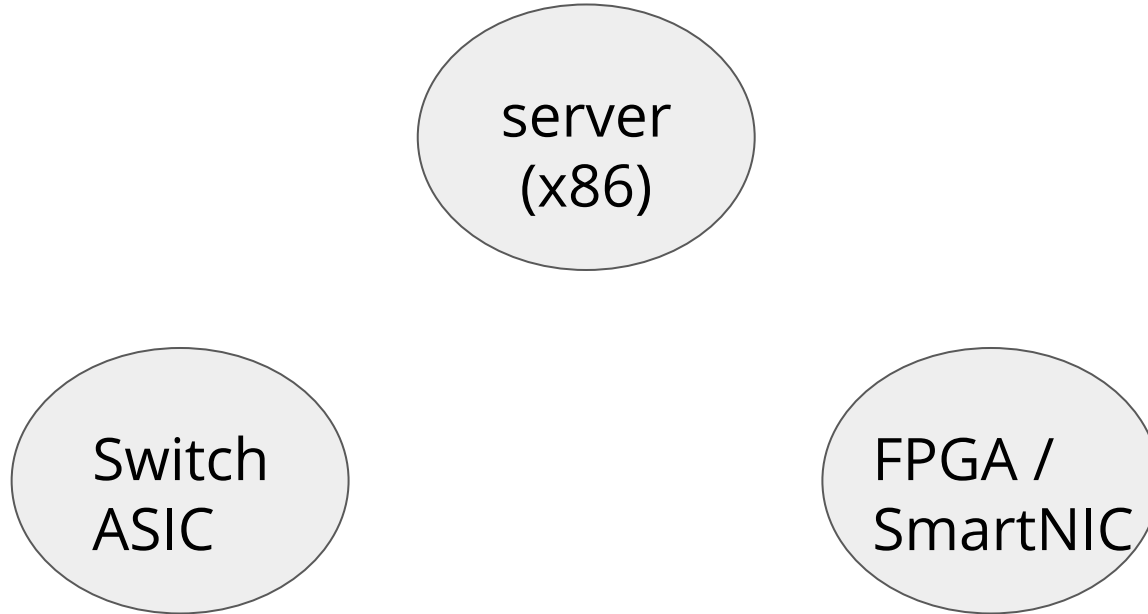**benefit from NFV *and* run high-speed**

# SEBA POD Topology



operator machine

ONOS, NEM, Kubernetes, Docker Registry, AAA, VOLTHA

mgmt network

ONU-1

ONU-2

Port5

Tibit

Tibit

Port1

OLT + SDN Fabric

BNG

IP/MPLS core

BISDN

# COTS brings cost down



source: https://www.lanner-america.com/latest-news/network-disaggregation/

# architectural options for BNG

server
(x86)

Switch
ASIC

FPGA /
SmartNIC

BISDN

# architectural options for BNG
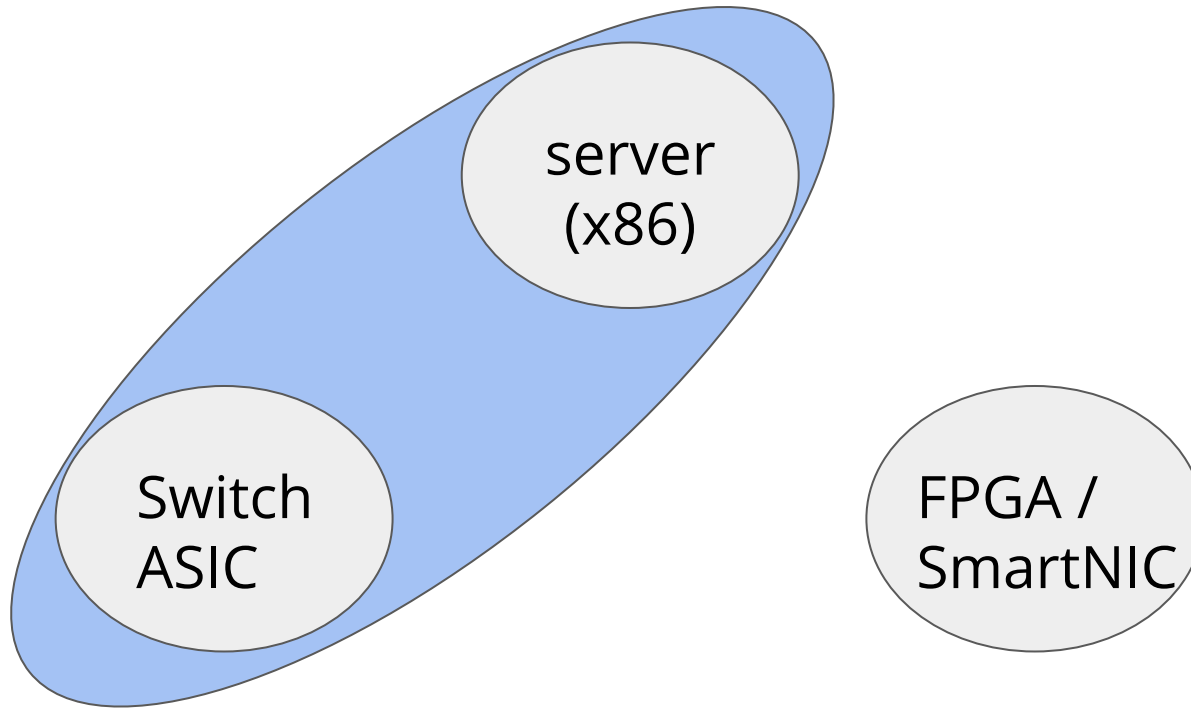
# architectural options for BNG



Switch ASIC

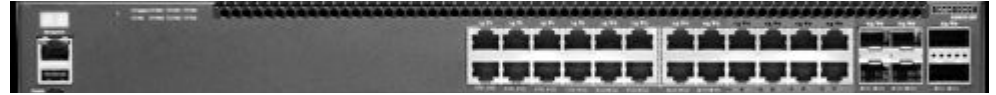server (x86)

FPGA / SmartNIC

# architectural options for BNG

# BISDN supports low-cost, high-volume platforms

- 1G/10G
  - AS 4610
  - Broadcom Helix4

- 10G/40G
  - AG7648
  - Broadcom Trident2

- 25G/100G
  - AG5648
  - Broadcom Tomahawk+

BISDN

# not every 'white box' is COTS

**Deutsche Telekom and Delta announced Open BNG platform at OCP Summit 2019**

http://www.delta-emea.com/news/pressDetail.aspx?secID=3&pID=1&typeID=1;2;8&itemID=9565&tid=0&hl=en-GB

**AGCX422S**
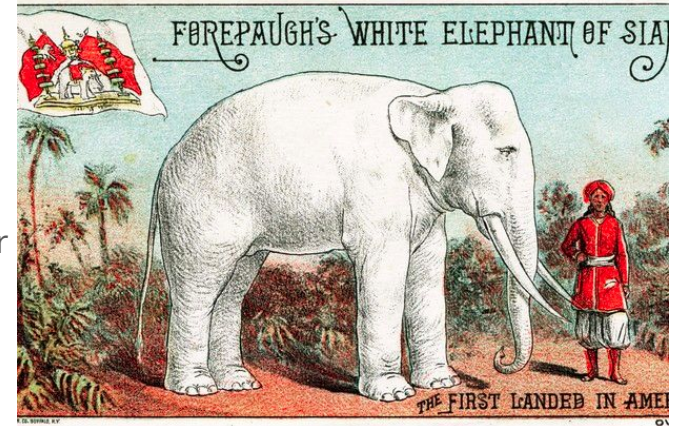**22*QSFP28 and 4*QSFP-DD**

**AGCVA48S**
**4*SFP+, 48*SFP28 and 10*QSFP28**

➢ Dual Q2C,
➢ additional counter processor
➢ Xeon-D inside

FOREPAUGH'S WHITE ELEPHANT OF SIA

THE FIRST LANDED IN AMER

custom designs won't fly **unless usable for mass market**.
(personal opinion of Dr.-Ing. Hagen Woesner, Berlin)

**BISDN**

# SEBA Overview - need to integrate into OSS/BSS

# vOLTHA Architecture

# Schematic of grpc/OF coexistence in k8s

- currently working on vOLTHA 1.6
- control plane extracted into VXLAN port
- transport via data port (not the internal CPU)

**kubernetes-cluster**

**k8s calico networking**

**vOLTHA-core**
**create OLT abstraction**

**TiBiT-Adapter**
**discover OLT devices**

**gRPC switch**
**configurator**

**vTEP**

**create**
**vTEP**

**switch**

**OLT**

**OLT**

**mgmt**
**traffic**

**vTEP**

**DATA-PORT**

**MGMT-PORT**

**gRPC Server**

**ofagent**

**OF-DPA**

**ASIC**

BISDN

# COTS BNG architecture - switch+server

Extend switch pipeline into server - balance the tasks to where they can be handled best



x86 server

Intel XXV710 | Intel XXV710 | Intel XXV710 | Intel XXV710

SuperMicro X11DPI-N, 2 x CPU Intel Xeon Gold 5120 (Skylake)
w/ 4 Intel XXV710DA2 cards

Tibit XGSPON module could go here,
or any other OLT uplink

25G
25G
25G

100G

Delta AG5648v1
(BRCM Tomahawk+)
BISDN Linux, baseboxd controller for L2 access and L3 core routing to spine

Access

Core

BISDN

# Packets do not leave hardware pipeline



control plane

downstream pipeline

delegate complex tasks like traffic shaping, NAT, firewall into DPDK worker processes

upstream pipeline

DDP: Split control from data traffic in hardware

DDP NIC

table0
table1

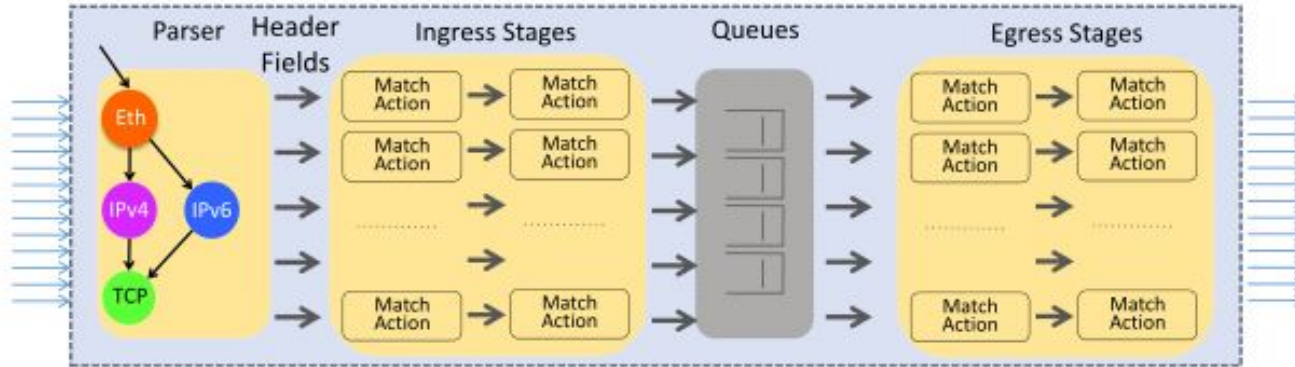ASIC pipeline

ASIC pipeline

keep simple tasks like fan-out, routing, MPLS, in whitebox ASIC
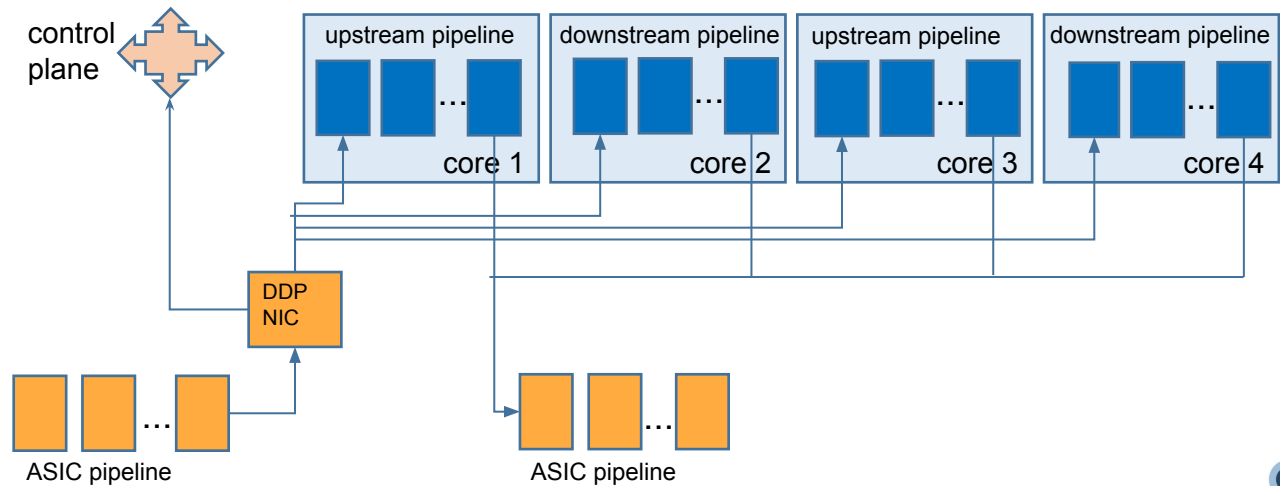
BISDN

# P4 pipeline / folded

- similarity in the pipeline structure
- both ingress and egress stage are handled in the same switch
- control interface is PFCP (NOT p4runtime), but similar paradigm
- use local OpenFlow controller for switch (github.com/bisdn/basebox)
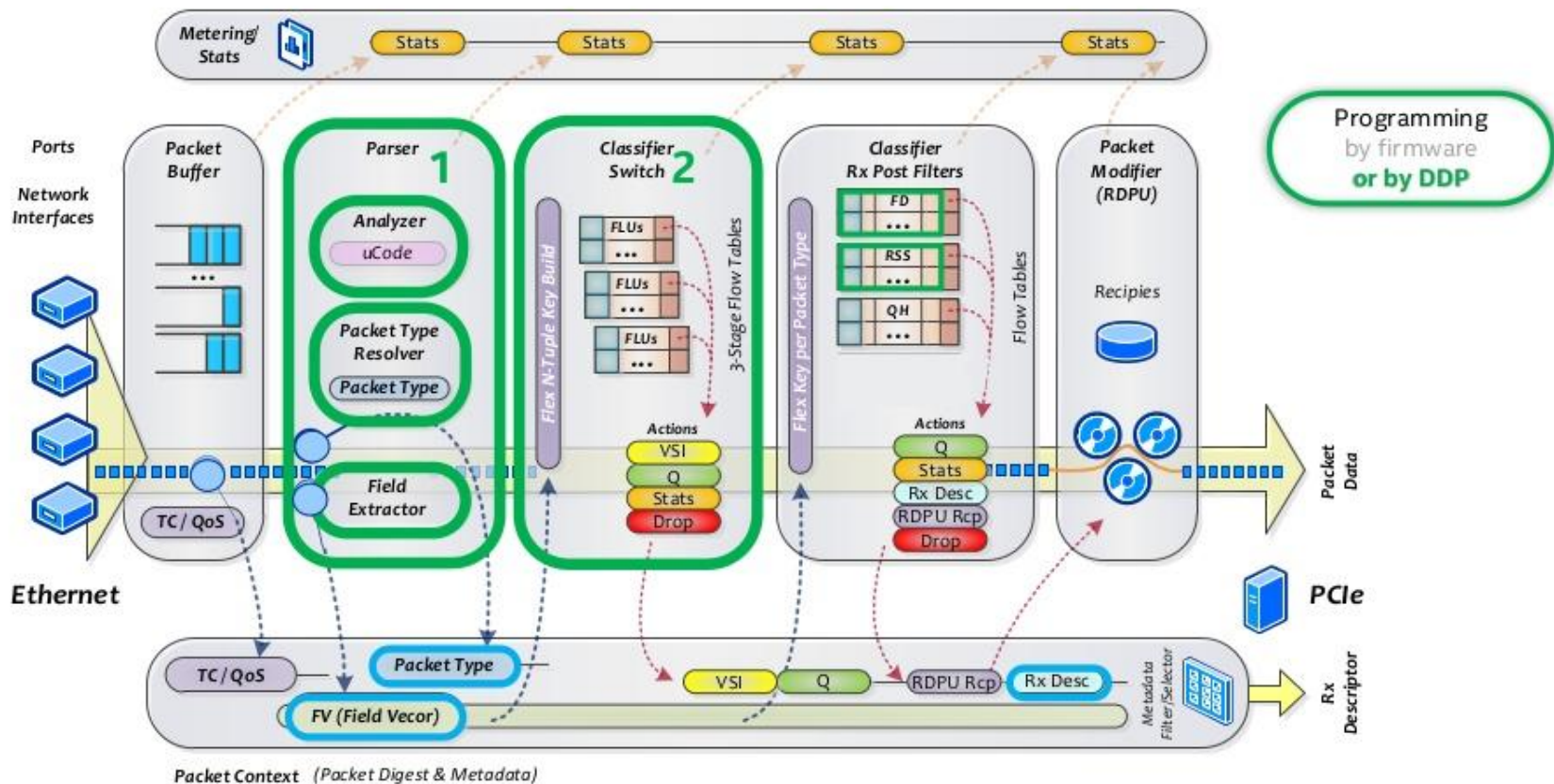


BISDN

# horizontal scaling

- scale to 16 parallel BNG instances on a single server
  - 2 instances per 25G port

# Dynamic Device Personalization - DDP

- program filters into NIC pipeline, profiles available for PPPoE, IPSec(?), GTP
  - cloud filters can be set/programmed with ethtool
- Here we use some specific filter config (made by Intel)
  - create 5 VFs
  - VF_0 is 'control VF', bind to kernel driver, on to VXLAN (remote control server)
  - VF_1, VF_3 downlink
  - VF_2, VF_4 uplink
- initialize NIC such that PPPoE control traffic is sent to CP VFs

BISDN

# Intel® Ethernet 700 Series Rx Programmable Pipeline

# vBNG Base Config - Uplink Packet Flow stages

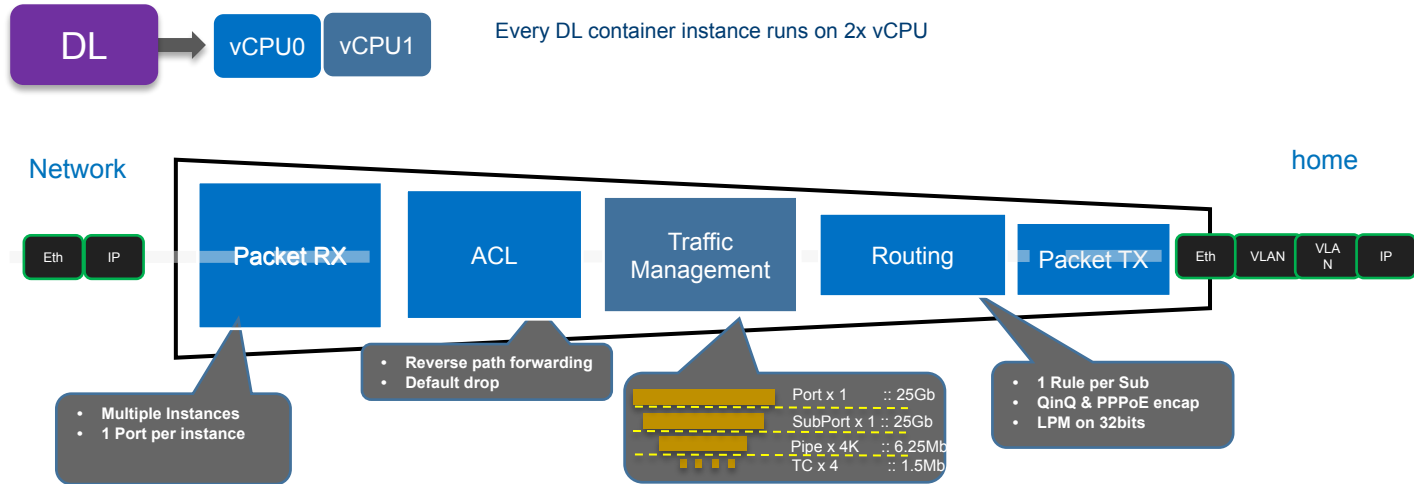**UL** → **vCPU0**  Every UL container instance runs full packet flow on single vCPU

home
Network

| Eth | VLAN | VLAN | IP |

Packet RX | ACL | Flow Classification | Metering/ Policing | Routing | Packet TX

| Eth | IP |

- **Multiple Instances**
- **1 Port per instance**

- **5 Tuple Lookup**
- **100 Infra REJ Rules**
- **15 Martian REJ rules**
- **15 SMTP ACC Rules (P0)**
- **Default SMTP REJ rule (P1)**
- **Default ACC**

- **1 FC rule per Sub**
- **Classify on QinQ**
- **Total rules per instance = Total subs supported per instance (e.g. 4K)**

- **1 Policer per Sub**

- **1 total route per instance**
- **1 Rule per Port**

ACL lookup tuple needs to be extended to block unwanted ICMP traffic

**BISDN**

# vBNG Base Config- Downlink packet flow stages

DL → vCPU0 vCPU1

Every DL container instance runs on 2x vCPU

Network

home

| Eth | IP |

**Packet RX** | **ACL** | **Traffic Management** | **Routing** | **Packet TX**

| Eth | VLAN | VLAN | IP |

- **Multiple Instances**
- **1 Port per instance**

- **Reverse path forwarding**
- **Default drop**

| Port x 1 | :: 25Gb |
| SubPort x 1 | :: 25Gb |
| Pipe x 4K | :: 6.25Mb |
| TC x 4 | :: 1.5Mb |

- **1 Rule per Sub**
- **QinQ & PPPoE encap**
- **LPM on 32bits**

BISDN

# vBNG Full Config - Uplink Packet Flow stages

UL → vCPU0    Every UL container instance runs full packet flow on single vCPU

home
Eth | VLAN | VLAN | IP

Packet RX | ACL | Flow Classification | Metering/ Policing | DSCP | NAT | Routing | Packet TX

Network
Eth | IP

- **Multiple Instances**
- **1 Port per instance**

- **5 Tuple Lookup**
- **100 Infra REJ Rules**
- **15 Martian REJ rules**
- **15 SMTP ACC Rules (P0)**
- **Default SMTP REJ rule (P1)**
- **Default ACC**

- **1 FC rule per Sub**
- **Classify on QinQ**
- **Total rules per instance = Total subs supported per instance (e.g. 4K)**

ACL lookup tuple needs to be extended to block unwanted ICMP traffic

- **1 Policer per Sub**

- **4K entries mapped to DSCP 46**

- **Static NAPT translation**
- **4K unique private IP address -> 1 public IP address with 4K UDP ports**

- **1 total route per instance**
- **1 Rule per Port**

**BISDN**

# vBNG Full Config- Downlink packet flow stages



DL → vCPU0 vCPU1

Every DL container instance runs on 2x vCPU

Network

home

Eth | IP → Packet RX | ACL | Traffic Management | NAT | Routing | Packet TX → Eth | VLAN | VLAN | IP

- Multiple Instances
- 1 Port per instance

- Reverse path forwarding
- Default drop

Port x 1      :: 25Gb
SubPort x 1 :: 25Gb
Pipe x 4K    :: 6.25Mb
TC x 4        :: 1.5Mb

- Static NAPT reverse translation
- 1 public IP address with 4K UDP ports ->
- 4K unique private IP address

- 1 Rule per Sub
- QinQ & PPPoE encap
- LPM on 32bits

BISDN

# Example IP pipeline code

```
;***************DECLARE ACTION PROFILE********
table action profile AP5 ipv4 offset 270 fwd encap qinq_pppoe
;***************PB'S CREATION*****************
pipeline downstream|firewall period 10 offset_port_id 0 cpu 0
;***************PB LINKING********************
pipeline downstream|routing port in bsz 32 swq SWQU02
;***************TABLE CREATION****************
pipeline downstream|routing table match lpm ipv4 offset 286 size 4K action AP5
;************TABLE ASSOCIATION****************
pipeline downstream|routing port in 0 table 0
;****************************PIPELINE BLOCKS****************************;
pipeline downstream|routing table 0 rule add match default action fwd drop
pipeline downstream|routing table 0 rule add bulk
./bng_configs/bulk_rules/dl_bng/route_bulk_pppoe_ins_%S5.txt
```

```
******************************** vBNG-DOWNLINK ********************************
LINK0 RXQ0 --->|           |--> SWQU00 -SWQU00-->|           |SWQU01--->
               |           |                     |           |
               |  FIREWALL |                     |    NAT    |
               |           |                     |           |
               -+-                               -+-
                |                                 |
SWQU01 --->|           |--> SWQU02 SWQU02 -->|           |--> LINK0
           |           |                     |           |
           |   HQOS    |                     |  ROUTING  |
           |           |                     |           |
           -+-                               -+-
            |                                 |
******************************** vBNG-DOWNLINK ********************************
```

**BISDN**

# vBNG Control Plane attachment



**BNG Control Plane Node**

Flow events logged to Redis DB

redis

DP Update Client

Radius

accel-pppd

Session information is programmed into the vBNG Instance on add/remove event

Mgmt Network

VXLAN Tunnel

VXLAN Tunnel

Connection between DP node and CP node is encapsulated

OVS connects multiple ports to accel-pppd. SDN controller on OVS can take over the load balancing

Accel-ppp has as many input ports (and MAC addresses) as there are uplink VFs

**BNG Data Plane Node**

BNG POD

UL Instance
Container
VF
CPU

DL Instance
Container
CPU
VF

Control VF

Dataplane VFs

Dataplane PF

NIC forwards all control traffic (PPPoE Link Ctrl & DHCP) to the Control VF

NIC forwards session data traffic is forwarded to the relevant VF

Assume that Dataplane VF MAC addresses are replicated at accel-pppd's input ports
Static mapping, input port determines the DP instances to be called from python script

BISDN

# Packet Forwarding Control Protocol (PFCP)

# BNG with PFCP:

# Measurement setup
## using a Spirent TestCenter 5.01 at DT

- 2*100Gbit/s Spirent ports attached to single server
  - SuperMicro X11DPI-N, 2xCPU Intel Xeon Gold 5120
    - 14 core @ 2.2GHz
  - Delta Agema AG5648 (BRCM Tomahawk) does routing and switching
    - split to 8*25G Intel XXV710DA2 cards
    - baseboxd controller on switch

- bidirectional traffic access+core in one port

- asymmetric upstream to downstream traffic ratio
  - upstream between 10 and 25% of downstream
  - session setup speed limited by DPDK pipeline implementation
    - currently ~10 sessions/s per instance, will be fixed

BISDN

# Queue dimensioning

identified three bottlenecks:

- DPDK IP Pipeline
- Fortville ASIC on NIC
- and then "something somewhere" in PCI/CPU, practical limitations of PCI bandwidth to CPU.

main focus on downstream pipeline.

upstream traffic in general much lower than downstream, and no HQoS (we did measure upstream, too.)

**BISDN**

# Characterization of IP pipeline

single vBNG instance consists of two containers

- **upstream**: one hyperthreaded core
- **downstream**: two hyperthreaded cores
  - siblings on one CPU core, one for HQoS only, the other for all the rest

DDP allows splitting the traffic into more than one instance

- single instance cannot quite serve a full 25 GbE port
  - this depends on the CPU, we used a 2.2 GHz CPU from 2017

... just take two.

BISDN

# DPDK pipeline measurements
## throughput test (RFC 2544) - achievable throughput with zero packet loss

for 4K sessions, single instance can process 5Mpps. ⇒ 11.1 Gbit/s @256 byte

two instances in the same port process 10Mpps. ⇒ 22.2 Gbit/s @256 byte

dominating factor for pipeline processing is number of sessions (lookup time for routing and HQoS)



max Throughput in Frames/second (RFC 2544) — single instance



Max. Throughput in frames/sec (RFC2544) — 2 instances per port

# DPDK pipeline measurements
**throughput test (RFC 2544), average delay at peak load**

average delay stays below 1ms

max delay is typically twice the average, min delay around 30 microseconds



Average delay at max throughput (RFC2544)

single instance per 25G port



Average delay at max throughput (RFC2544)

2 instances per port

average delay

single instance, 512 PPP sessions

512 sessions

average delay

single instance, 2K PPP sessions

2K sessions

dashed line ···
is 1 ms
end-to-end
delay (incl.
switch)

delays above 1ms
typically mean
packet loss

1K sessions

4K sessions

# For 2 instances
## ... typical delay around 70 $\mu$s

- delay curves are flat until very high load is applied
- switch traversal is 3.8 $\mu$s, i.e. a total of 7.6 $\mu$s is included that is caused by the switch
  - similar number for Mellanox and Broadcom



average delay

two instances, 2K PPP sessions

2K sessions



average delay

two instances, 1K PPP sessions

1K sessions



average delay

two instances, 4K PPP sessions

4K sessions

# port vs. card limitations

single port can serve full line rate of a dual port XXV710 card

dual-port card is limited to ~20 Mpps

both ports can be filled up to 45-46 Gbit/s in total

above that, packet losses "before" the pipeline

36

**BISDN**

# single card to multiple cards

multiple cards on the same socket are not completely independent **in our server (!)**

loss-free operation until 80% offered load.

multiple sockets are indeed completely independent

total achievable loss-free throughput is 160Gbit/s
        combined up+downstream traffic

37

# 160 Gbit/s total throughput, lossless

8 ports, 2 instances per port, 16 G down/ 4G upstream traffic

# vBNG x86+whitebox implementation
**Lessons learned**

DDP allows forwarding of traffic into separate vBNG instances

- no hashing, no cache pollution, 4K entries fit into cache

delay is generally low, but moreover: **manageable**

- if delay is too high ⇒ add a pipeline

processing of 160 Gbit/s of traffic is possible on a commodity (low end) server

there are bottlenecks outside of the pipeline

- increasing the Rx/Tx queue size to 4K and
mempool to 320K helped a lot

39

BISDN

# Conclusions

- COTS is bringing cost of access equipment down

- Linux as common API over switches and servers

- stateful functions like OLTs can be separated from forwarding

- stateful functions like PPPoE session termination, HQoS, NAT should be separated from forwarding

- whiteboxes + pluggable OLTs + programmable NICs allow this

BISDN