# COMMS SUPPORT IN INTEL® ETHERNET 800 SERIES

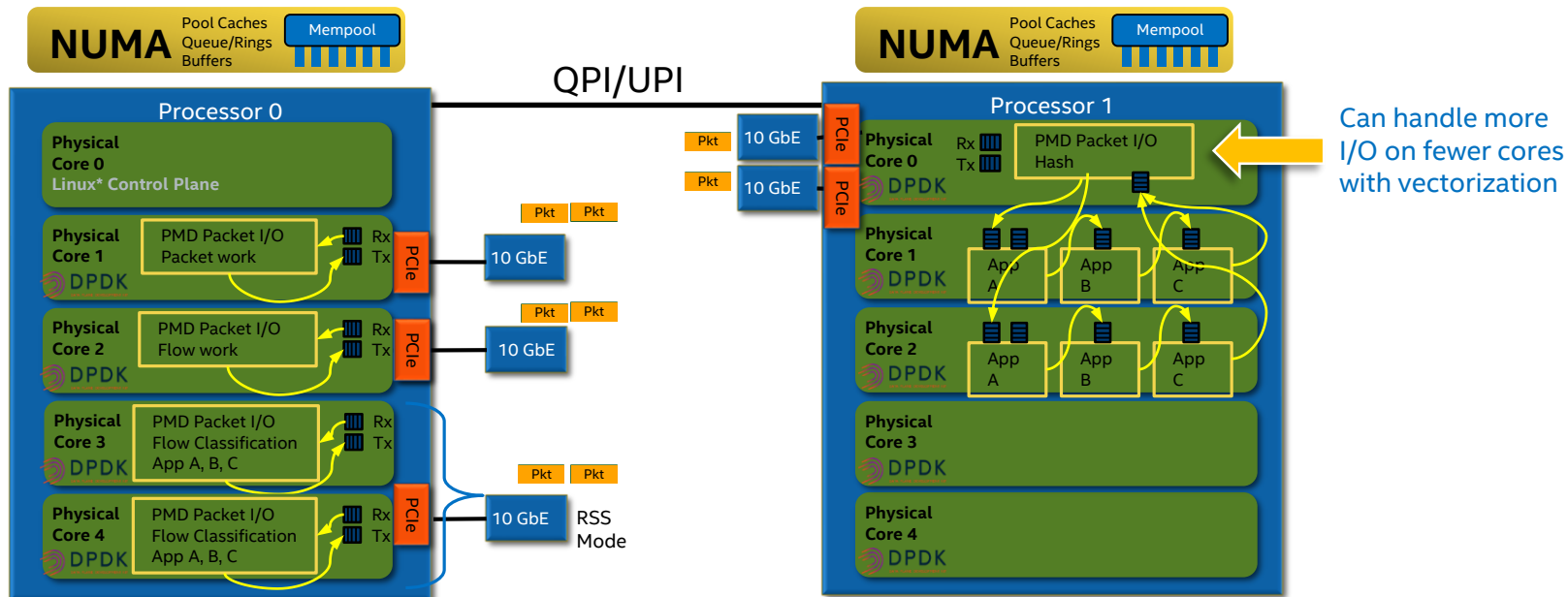NPG PRC PP SW Team
June 2020

# Agenda

- Background

- Programable Pipeline matters in comms

- VNF Use cases

# PCIe* Connectivity and Core Usage

## Using run-to-completion or pipeline software models



**Run to Completion Model**
- I/O and Application workload can be handled on a single core
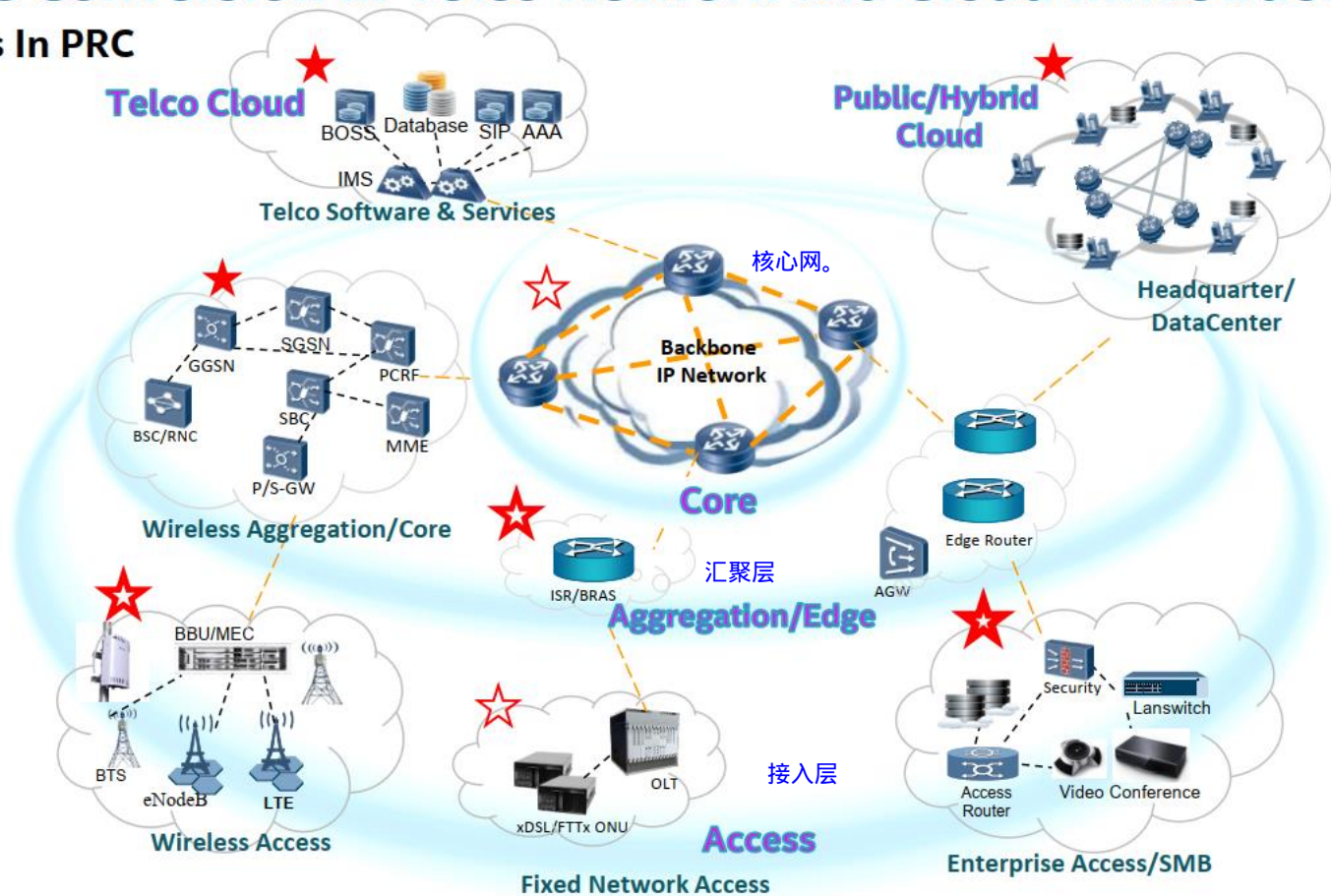- I/O can be scaled over multiple cores

**Pipeline Model**
- I/O application disperses packets to other cores
- Application work performed on other cores

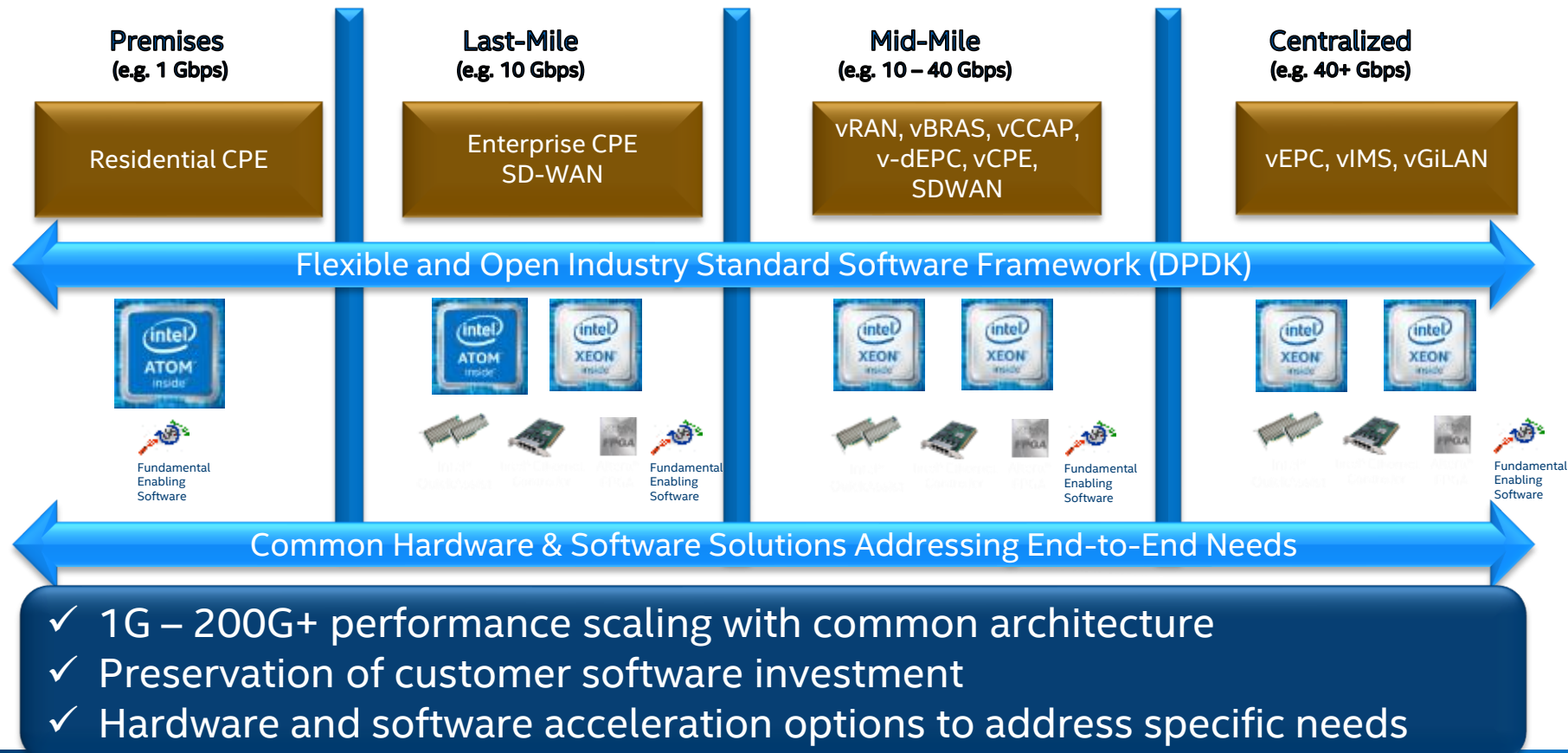Can handle more I/O on fewer cores with vectorization

```
default_hugepagesz=1G hugepagesz=1G hugepages=16 hugepagesz=2M hugepages=2048 isolcpus=1-11,22-33 nohz_full=1-11,22-33 rcu_nocbs=1-11,22-33
```
Note: nohz_full and rcu_nocbs is to disable Linux* kernel interrupts, and it's important for zero-packet loss test. Generally, 1G huge pages are used for performance test.

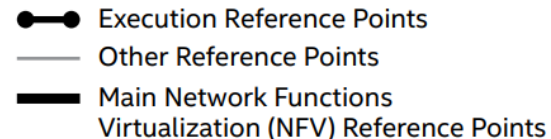# Architecture Conversion In Telco Network and Cloud Infrastructure

★ **IA Adoptions In PRC**

# Meeting Scalable Data Plane Needs

| **Premises**<br>(e.g. 1 Gbps) | **Last-Mile**<br>(e.g. 10 Gbps) | **Mid-Mile**<br>(e.g. 10 – 40 Gbps) | **Centralized**<br>(e.g. 40+ Gbps) |
|---|---|---|---|
| Residential CPE | Enterprise CPE<br>SD-WAN | vRAN, vBRAS, vCCAP,<br>v-dEPC, vCPE,<br>SDWAN | vEPC, vIMS, vGiLAN |

**Flexible and Open Industry Standard Software Framework (DPDK)**

intel ATOM inside — Fundamental Enabling Software

intel ATOM inside / intel XEON inside — Fundamental Enabling Software

intel XEON inside / intel XEON inside — Fundamental Enabling Software

intel XEON inside / intel XEON inside — Fundamental Enabling Software

**Common Hardware & Software Solutions Addressing End-to-End Needs**

- ✓ 1G – 200G+ performance scaling with common architecture
- ✓ Preservation of customer software investment
- ✓ Hardware and software acceleration options to address specific needs

# Workloads

| Location | Workloads |
|---|---|
| Access | Base Station, OLT, Router |
| Aggregation/Edge | BNG, Gateway, UPF |
| Core | UPF |
| Cloud | vSwitch, Load Balancer... |

# NFV Architecture Framework

# Map Network Platform Capability to SW Stack



**NFVi SW Stack**
- Northbound I/F (Vf-Ni)
- Service Agent
- HAL I/F (Kernel I/F, DPDK I/F)
- Driver

Refer to Orche. I/F & VNFM I/F (out of scope)

**VNF SW Stack**
- VeNf-Vnfm I/F
- UP/CP SW
- Vn-Nf API
- HAL I/F (Kernel I/F, DPDK I/F)
- Driver

**Device/Driver Interface (Device Specific)**

Network Platform (e.g. CVL)

map    NFV

... jingjing

**NFVi**
- {VM | Container}
- Northbound I/F
- Service Agent
- HAL I/F
- Kernel Provider
- DPDK Provider
- Custom. Provider

1 x Infra. Daemon

**VNFs**
- {VM | Container}
  - VNF UP SW Stack
  - HAL I/F
  - IAVF w/ ext.
- {VM | Container}
  - VNF CP SW Stack
  - HAL I/F
  - IAVF

M x User Plane VNFs

N x Control Plane VNFs

Driver Maps NP Capability to HAL

SW

NFVi    VF

Platform

Device I/F    CFG    PKT QPs    CFG    PKT QPs    CFG    PKT QPs

PF    S VF    S VF

RSS/FDIR    RSS/FDIR

VIM O&M    S

ACL

Switch

Parser

CVL    e.g. AAA, DHCP, PPPoED, etc.

LAN

# PROGRAMABLE PIPELINE MATTERS IN COMMS

DDP

# OS Protocol vs Comms Protocol

| OS Protocols | Comms Protocols | |
|---|---|---|
| VLAN | S-VLAN | OSPF |
| IPv4 | C-VLAN | BGP |
| IPv6 | MPLS | GTP |
| ARP | IPv4 | PFCP |
| ICMP | IPv6 | PPPoE |
| UDP | ARP | L2TP |
| TCP | ICMP | ESP |
| VXLAN | UDP | ...... |
| | TCP | Customized Protocols |

# Programmable Pipeline Matters



Worker 0 | Worker 1 | Worker 2 | Worker 3

Distributor Worker <-

worker

Customize Protocols

**Default** pipeline relies on host to parse undefined protocols

| Destination Address | Source Address | Undefined Protocol Start of Payload | ? | ? |
|---|---|---|---|---|

Parsed fields — Payload

Worker 0 | Worker 1 | Worker 2 | Worker 3

Load Distribution looking deeper

parser

Customize Protocols

The pipeline parser looks deeper in to the packets

| Destination Address | Source Address | Defined Protocol Outer Header | Inner Header | Payload |
|---|---|---|---|---|

Parsed fields — Payload

# Dynamic Device Personalization (DDP) Profile

2

**Normal Mode**

**OS Default Package**

NVM Default Package

OR

3

Comms Package

**Comms** ➕

**OS Default Package**

NVM Default Package

**Driver**

1

Safe Mode

**NVM Default Package**

20. 05

- **NVM Default Package**

Is a bare minimum package which allows reliable connectivity with basic protocols
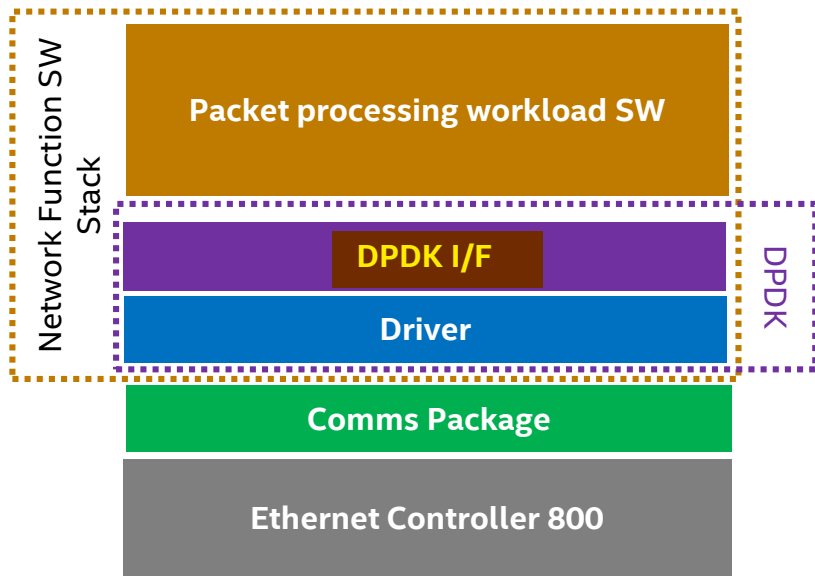
- **OS Default Package**

Baseline support for well known protocols and configurations

- **Comms Package**

Comms protocols and configuration built on top of (incremental) existing OS default package

# Bridge SW Stack to Programable Pipeline



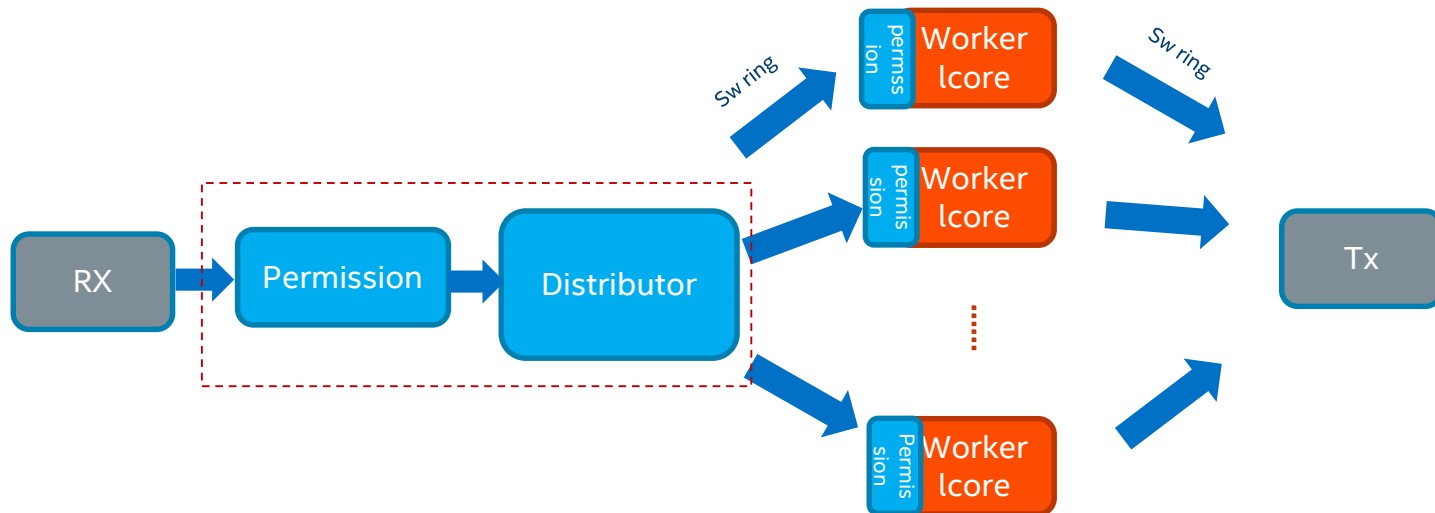DPDK is the bridge between packet processing workload and hardware programable pipeline capability.

- Download package
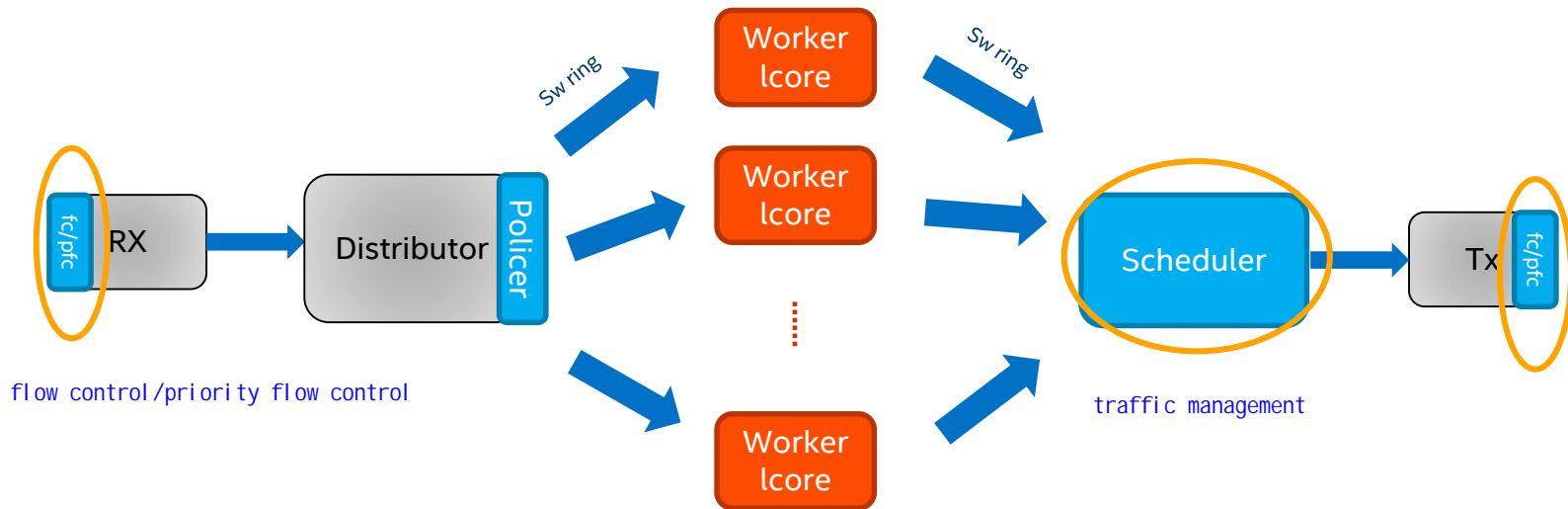- Run time table and rule programing

DPDK        api

# VNF USE CASES

# Simple Pipeline of packet process – Classify

# Simple Pipeline of packet process - QoS



RX

flow control/priority flow control

Distributor

Policer

Sw ring

Worker lcore

Worker lcore

Worker lcore

Sw ring

Scheduler

traffic management

fc/pfc

Tx

fc/pfc

# vBNG Data Plane Processing Stage Example

# vBNG requirement to CVL

# vOLT Data Plane Processing Stage Example

# vOLT Data Plane Processing Stage Example

# vOLT requirement to CVL

contrl plane  vm/container        data plane  vm/container.

**vOLT CP VF**
(multicast replicator)

**DP VF (PON)**

N x ONU
N x ONU

TC

**DP VF**

N x ONU
N x ONU

TC

...

One RX core in run to completion mode with TX QoS offloaded to HW using 8 TX queues/8 TC in ETS mode.

Default?

Multicast

sub-port

sub-port

- single VLAN: switch on the VLAN (S-VLAN)
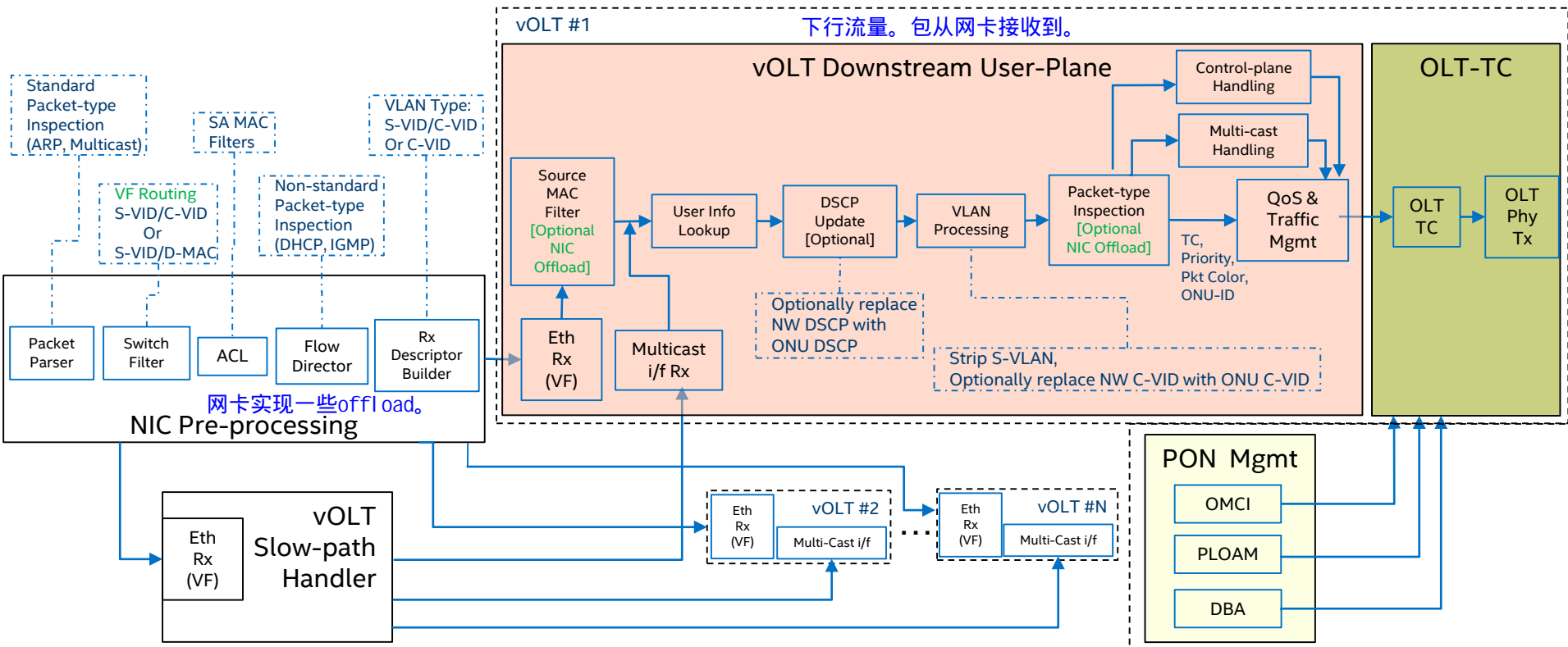- QinQ: switch on S-VLAN + C-VLAN
- QinQ: switch on S-VLAN only
- single VLAN: switch on S-VLAN + MAC DA
- QinQ: switch on S-VLAN + MAC DA
- switch on MAC DA

PF Switch directs control plane traffic to CP queues:
- Default traffic
- Multicast traffic

DCPF

**PF**

Switch rules

LAG

https://jira.devtools.intel.com/browse/DPDK-12400

# 5G UPF – software-based pipeline

workload



**5G Control Plane Domain with Service Based Architecture**

| Access & Mobility Mgmt Function (AMF) | Session Mgmt Function (SMF) | Other CP nodes (PCF, UDM, AUSF, etc) |

ACL lookup / action → GTP decap & Echo Req/Rsp processing → PFCP Session lookup → PDU Session lookup → Filter known Ptypes to CP → PFD search, PDR search → FAR lookup & Fwd/Drp/Buf enfrc

**5G User Plane Function (UPF)**

FAR policy enforcement (IP header, DSCP marking, etc.) ← TCP/HTTP Proxy, NAT ← Value Added Services HHE, DPI, Tethering detect, HTTP redirect, Sub-FW, URL filtering, Statistics, etc. ← URR look up / action (Charging) ← QER look up / action (QoS)

# vEPC/UPF

# COMMS PROTOCOL ENABLE IN DPDK

# Driver Enabling Software

| | Mode | Network Function | Host Driver | VNF HAL I/F | NFVi HAL I/F | Problems |
|---|---|---|---|---|---|---|
| ✅ | non-IOV *virtualization* | DPDK PF | N/A | DPDK I/F (e.g. rte_flow, rte_tm & etc.) | N/A | |
| ✅ | IOV | DPDK AVF *AVF-Intel VF* | Kernel LAN | DPDK I/F (e.g. rte_flow, rte_tm & etc.) | Kernel I/F (e.g. ethtools, iprouter, tc & etc.) | VNF Accl. requires extra device advanced features against kernel mainstream readiness. *vf* |
| ❌ | IOV | DPDK AVF | DPDK PF | DPDK I/F (e.g. rte_flow, rte_tm & etc.) | DPDK I/F (e.g. rte_flow, rte_tm & etc.) | Hard to sustain with incremental complexity of PF/VF co-existence by Host PF PMD kernel mainstream user space driver develop framework does not support SRIOV management. |

## Driver Enabling SW Strategy

- Minimize replicated effort on device specific function enabling

- Incremental building re-useful SW for multi-gen NIC coverage

- NPG/ND co-design, KMD/UMD co-existence

# **DCF** + iAVF advanced Feature is the way



- DCF
  - Device Config Function
  - Advanced Cap over trust VF
  - Functional at named VF
  - Single entity per port
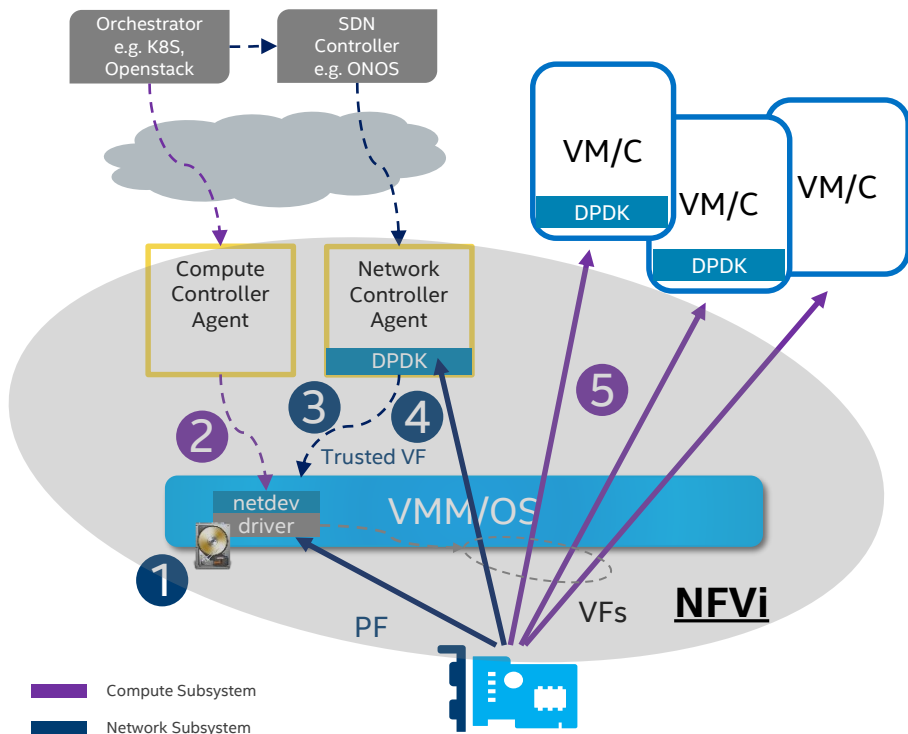  - AQ-CMD over virtnnl
- iAVF
  - Extend virtnnl to support advanced feature
  - More protocols
  - RSS input set change
  - FD supporting

# View of Deployment

hai yue                    DCF                        use case



1. Load DEF/COMM Package
2. PF driver generates SR-IOV VFs
   *echo 4 > /sys/class/net/…/sriov_numvfs*
3. Turn-on trust mode on a dedicated VF
   *ip link set dev $eth $vf_id trust on*
4. Assign trusted VF (VF0) to network controller agent
   *… options : dpdk-devargs=$BDF,cap=dcf,representor=[x]*
5. Assign other VFs to VMs/Containers

# Requirements for Comms NFV

| Feature | Details | Reference VNF | | |
| | | Wireline (vBNG) | Wireless (vEPC) | Cable (vCMTS) |
|---|---|---|---|---|
| Forward PPPoE Session signalling packets to VNF Control Plane VF | PPPoE Session packets with non-IP payload (APIs to | X | | |
| Forward packets with specific MAC DA to a User Plane VF | UP VF has no MAC address assigned, switch filter is needed to force packets to a VF | X | | |
| Forward IP multicast packets to VNF Control Plane VF | vOLT requirement, low priority | | | X |
| Forward L2 multicast packets to VNF Control Plane VF | vOLT requirement, low priority | | | X |
| Forward packets with specific MAC DA + VLAN combination to a User Plane VF | | | | X |
| Forward packets with specific VLAN to a User Plane VF | For single or outer VLAN incase of QinQ | | | X |
| Forward packets with specific QinQ combination to a User Plane VF | | | | X |
| Forward packets with IP protocol to VNF Control Plane VF | Match on IP Protocol only (IGMP, L2TP) | | | X |
| Forward packets with IP protocol tunnels to VNF Control Plane VF | Match on IP Protocol (ESP, L2TP) and IP src/dst addresses | | X | X |
| Forward all not-matched traffic to the default VF | | | | X |
| | | | | |
| **Hardware ACL requirements** | | | | |
| **Firewall rules for IP packets** | | | | |
| Deny IP SRC subnet | deny ip 192.0.2.0 0.0.0.255 any | X | | |
| Permit IP DST, TCP DST | permit tcp any host 192.168.201.103 eq smtp | X | | |
| Deny L4 (UDP/TCP) DST | deny tcp any any eq smtp | X | | |
| UPD SRC/DST range | permit uap any gt 1023 192.168.201.0 0.0.0.255 gt 1023 | | | |
| | | | | |
| **Drop L2 packets according to firewall rules** | | | | |
| Drop L2 packets if MAC SA is not in the whitelist | 2 to 4 valid MAC SA | | | X |
| | | | | |
| **Flow Director requirements** | | | | |
| Forward specific UDP DST port to a Queue/Qgroup | | | PFCP | |
| Tag packets with specific UDP DST port | | | | DHCP |
| Tag packets with specific IP protocol | | | | IGMP |
| Tag packets with specific L2 Ethertype | | | | ARP |
| Map IPv4 DSCP to RSS Qgroup | | | X | |
| Map IPv6 TC to RSS Qgroup | | | X | |
| Map GTP-U QFI to RSS Qgroup | | | X | |
| Forward specific IP protocol to a Queue/Qgroup | | | | L2TP |
| Forward specific L2TP session to a Queue/Qgroup | | | | L2TP |