# The Past, Present and Future of P4

Deb Chatterjee

START

**2014**

The paper "Programming "Programming Protocol Protocol Independent Packet Processors" is published in ACM Sigcomm

**2015**

P4 introduced its first first version P4$_{14}$

**2017**

P4RunTime, the control control plane component, component, is announced.
Around the same time, time, the evolved language spec P4$_{16}$ is published

**2018**

P4 joins ONF

**2019**

Intel acquires Barefoot

**2023**

Intel stops development of future generations of Tofino

**2024**

P4 moves from ONF to Linux Foundation

# Cascade Glacier

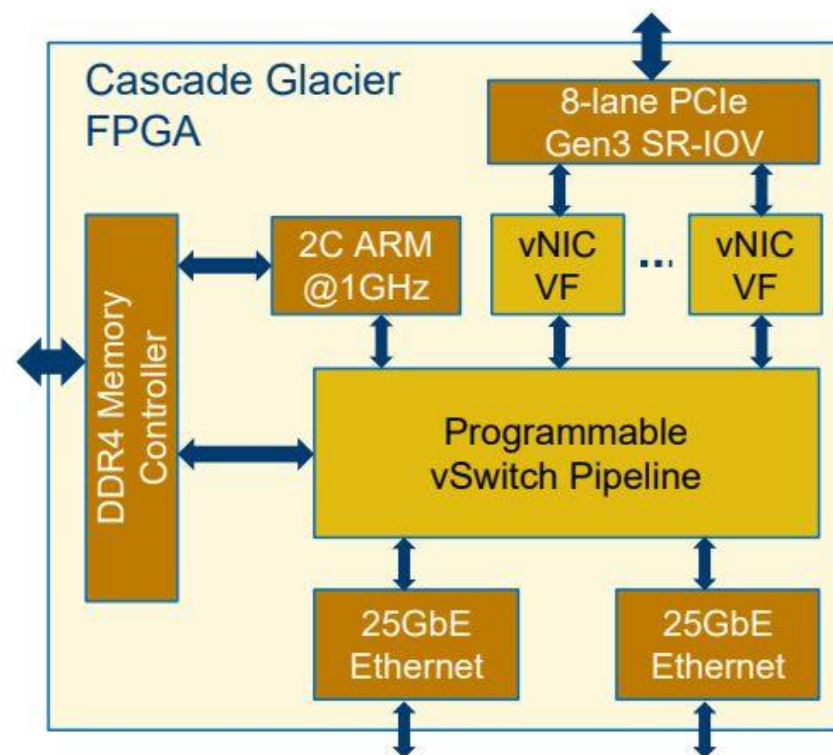## Programmable vSwitch Offload

- *Customer vSwitches & Open vSwitch*
- 128Gb Internal Switch Fabric
- Millions of flows @ 12.5 Mpps x2
- Programmable encap/decap & NAT
- Connection tracking & ACLs
- Pipeline generated from P4

## Embedded CPU Cores

- *vSwitch slow path & NIC management*

## Virtio-net Hardware Offload

- *Supports Existing Linux & DPDK VMs*
- Up to 40Gb/sec, 12.5 Mpps
- 128 VFs, Multiple queues/VF
- TSO, CSUM, transmit rate limiters
- Live migration between HW & SW



Cascade Glacier FPGA block diagram:
- 8-lane PCIe Gen3 SR-IOV
- 2C ARM @1GHz
- vNIC VF ... vNIC VF
- DDR4 Memory Controller
- Programmable vSwitch Pipeline
- 25GbE Ethernet / 25GbE Ethernet

**Arria-10 25GbE Board**
- Production by Q2'18
- Dual 25G SFP28
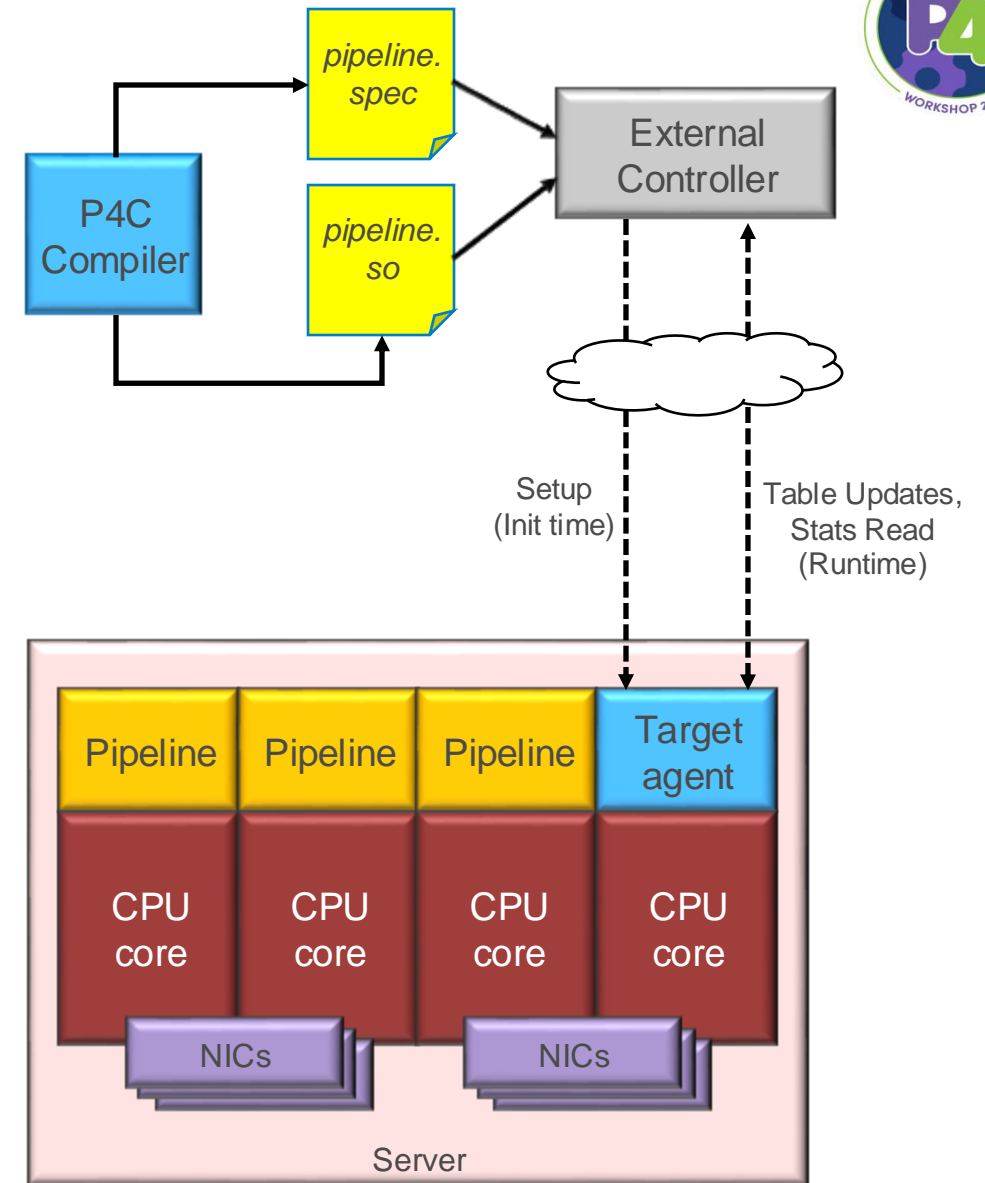- PCIe gen3 x8
- DDR4 Memory
- <45W, passive heatsink

# P4-DPDK



### Step 1: Offline

- The P4C compiler (i.e. p4c-dpdk) translates the input *pipeline.p4* file into an intermediate representation *pipeline.spec* (plain text file) or *pipeline.so* (binary file).

### Step 2: Run-time

- The external controller connects to the target agent to load the P4 blob: in the baseline "interpreted" mode, this is the *pipeline.spec* text file; in the optimized "compiled" mode, this is the *pipeline.so* binary file.

- The target agent creates the P4 objects for each pipeline based on the P4 blob and maps each pipeline to a CPU core.

- Each CPU core executes one or multiple pipelines by running the associated instructions for each input packet. Each pipeline is single threaded.

- The external controller performs table updates and reads the statistics.

# IPDK, with P4 as the pillar for everything

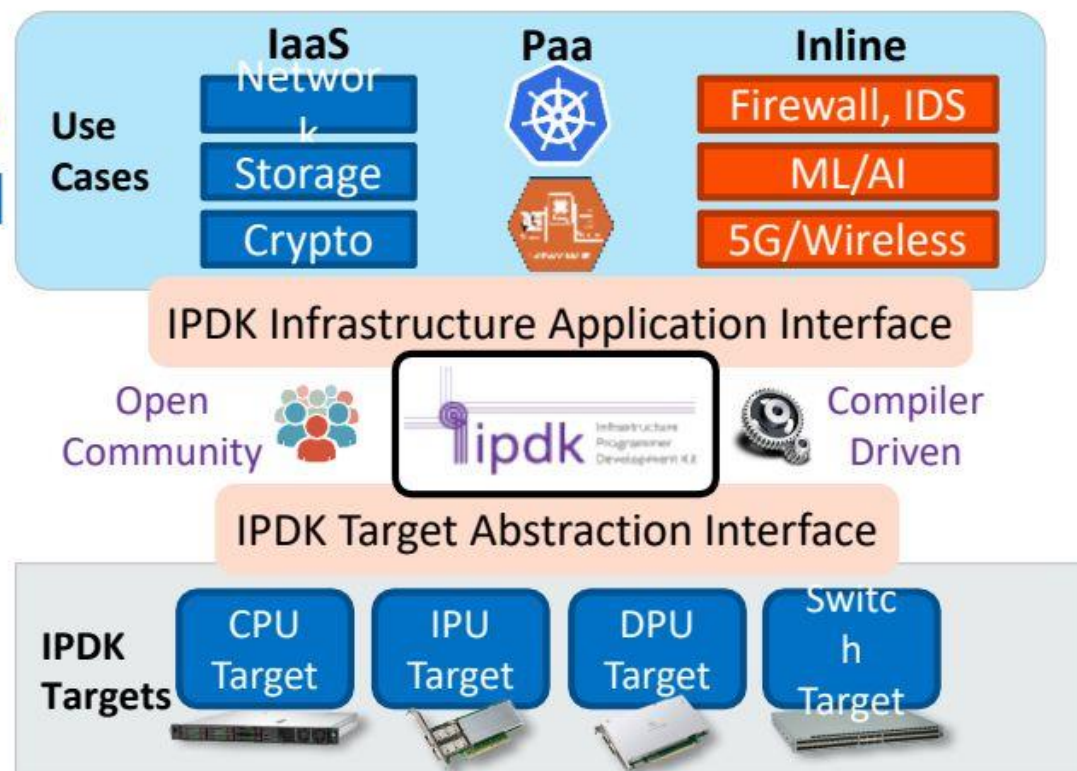1. **Infrastructure-as-a-Service**
   Virtual Networking, Storage & Crypto
   Across VMs, containers & bare metal

2. **Platform-as-a-Service**
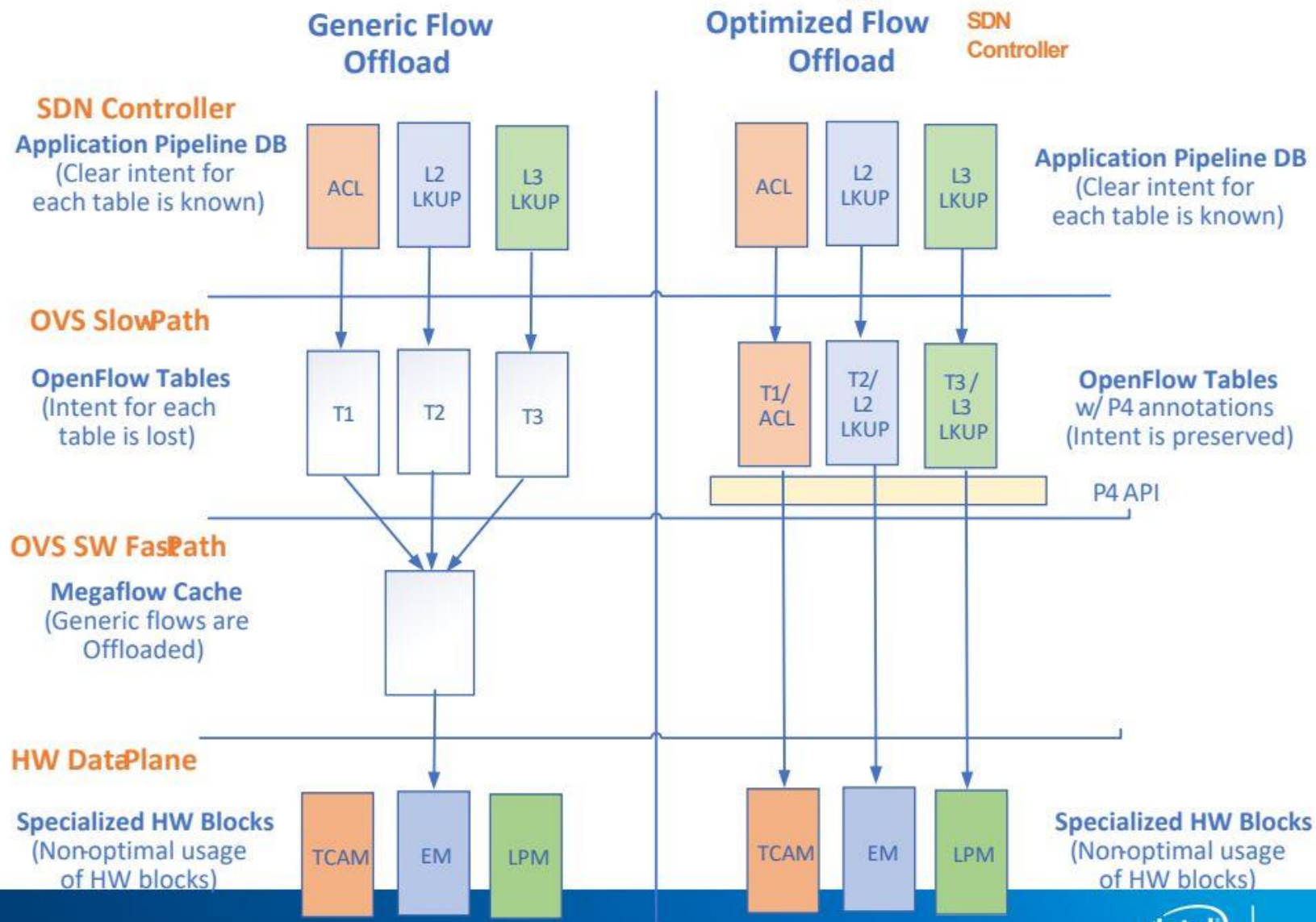   Container Networking (Kubernetes)
   Sidecars (Envoy, MongoDB)

3. **Inline Acceleration**
   Firewall, IDS, Network Telemetry
   5G/Wireless Infrastructure, AI/ML

Use Cases

IaaS: Network, Storage, Crypto
Paa
Inline: Firewall, IDS; ML/AI; 5G/Wireless

IPDK Infrastructure Application Interface

Open Community — ipdk — Compiler Driven

IPDK Target Abstraction Interface

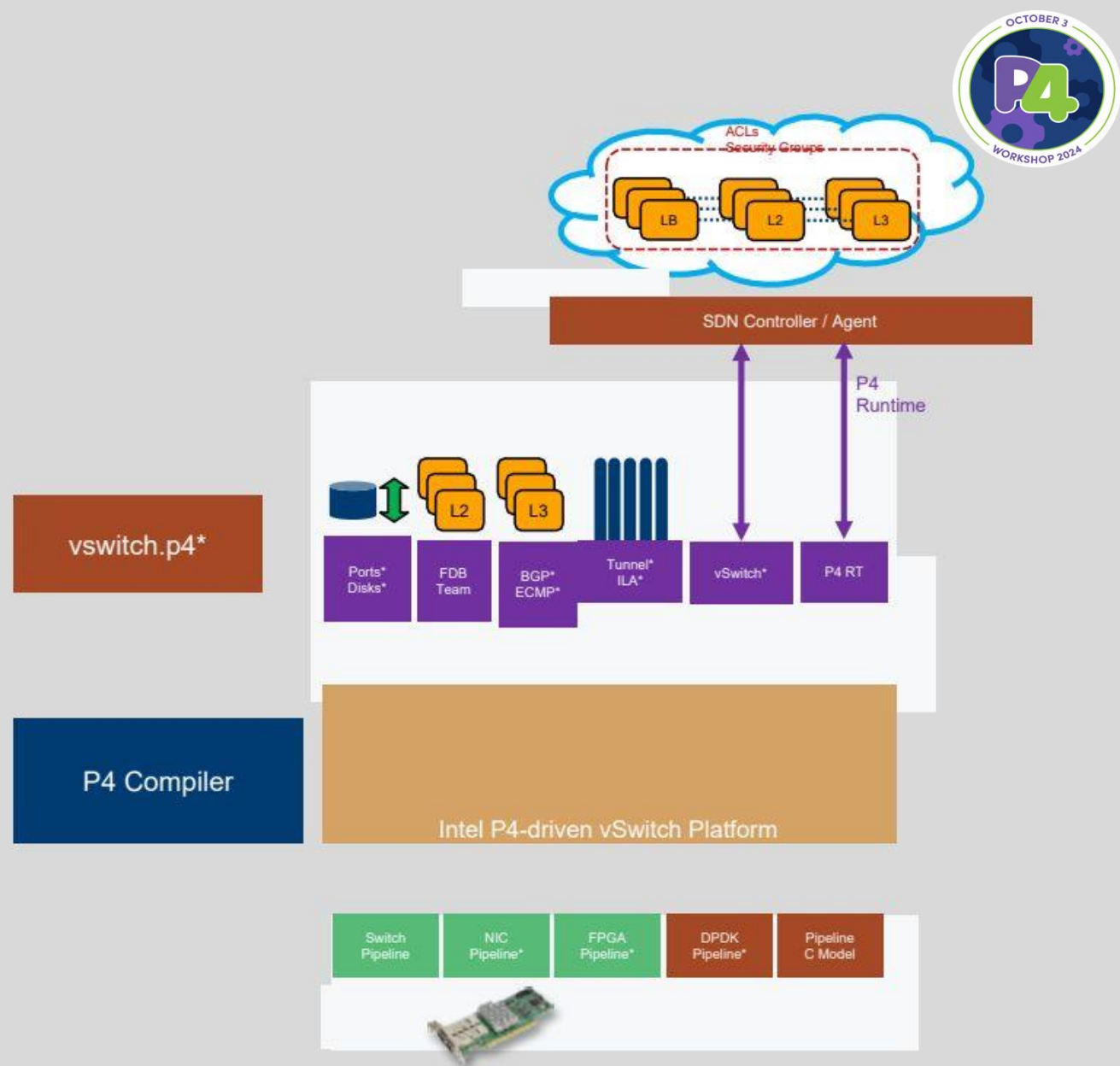IPDK Targets: CPU Target, IPU Target, DPU Target, Switch Target

# OVS HW offload challenges

- OVS Megaflow cache - great SW optimization idea, can be slow, potential of DDOS attack
- The reactive behavior requires high slow path performance; maintaining high flow set up rate is a challenge
- The aggregation followed by disaggregation -> the intent is lost
- Better -> skip the caches and offload directly to the HW tables
- Solves the other associated problems

**Generic Flow Offload**

**Optimized Flow Offload**    SDN Controller

**SDN Controller**

**Application Pipeline DB**
(Clear intent for each table is known)

| ACL | L2 LKUP | L3 LKUP |
| --- | --- | --- |

| ACL | L2 LKUP | L3 LKUP |
| --- | --- | --- |

**Application Pipeline DB**
(Clear intent for each table is known)

**OVS SlowPath**

**OpenFlow Tables**
(Intent for each table is lost)

| T1 | T2 | T3 |
| --- | --- | --- |

| T1/ ACL | T2/ L2 LKUP | T3/ L3 LKUP |
| --- | --- | --- |

**OpenFlow Tables**
w/ P4 annotations
(Intent is preserved)

P4 API

**OVS SW FastPath**

**Megaflow Cache**
(Generic flows are Offloaded)

**HW DataPlane**

**Specialized HW Blocks**
(Non-optimal usage of HW blocks)

| TCAM | EM | LPM |
| --- | --- | --- |

| TCAM | EM | LPM |
| --- | --- | --- |

**Specialized HW Blocks**
(Non-optimal usage of HW blocks)

OCTOBER 3
P4
WORKSHOP 2024

intel

# P4-driven vSwitch platform

- As HW vendor, we need to support many platforms
  - Switch, NIC, FPGA, SW

- As HW vendor, we need to support many vSwitches
  - Custom vSwitches not upstreamed
  - Windows GFT, Vmware NSX-T, VPP and many others

- P4-driven vSwitch platform
  - Is common, extensible, and programmable
  - Any specific vSwitch can be supported by simply adding a thin shim layer
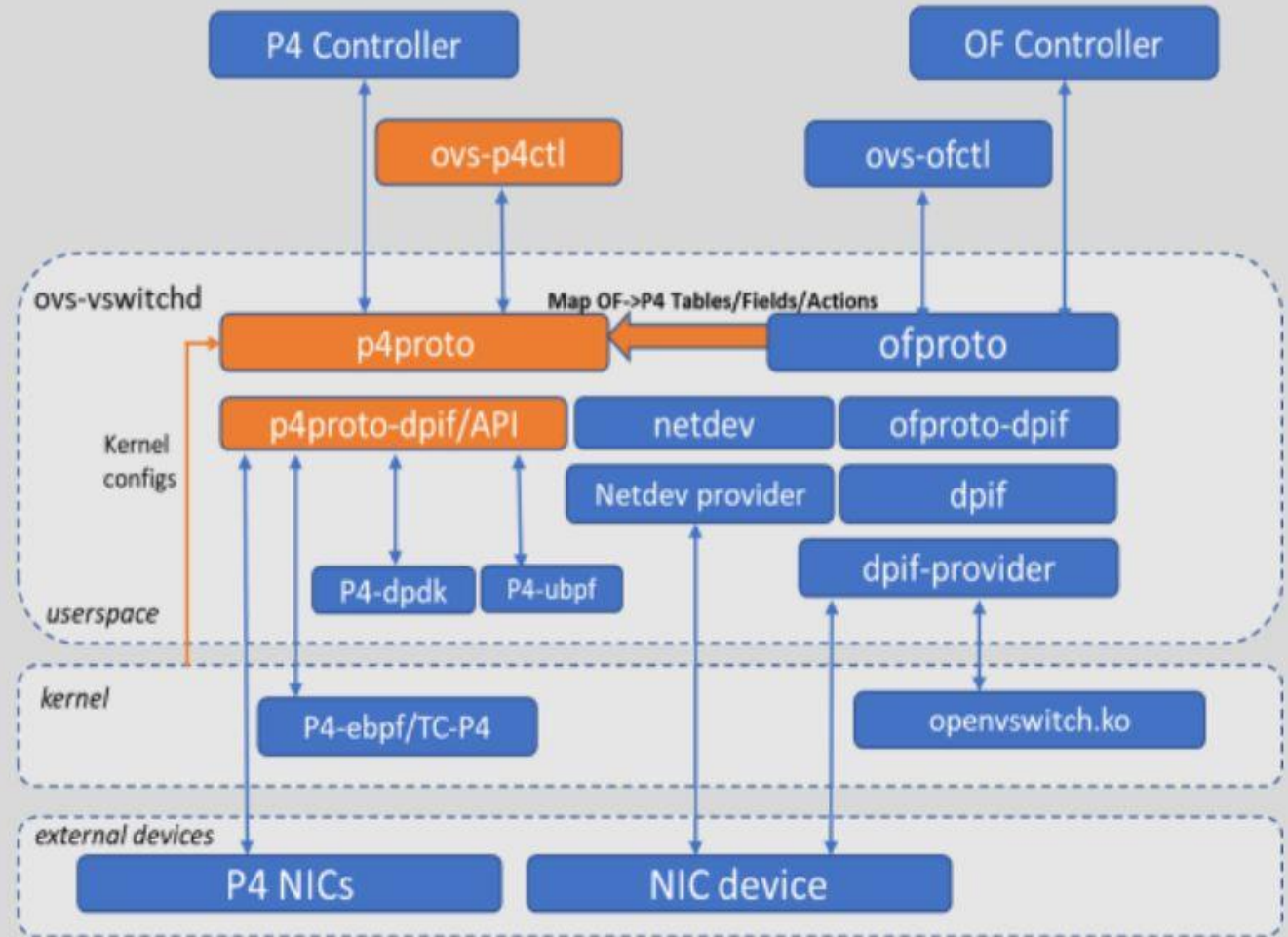
# P4-driven vSwitch platform

**Control planes**:

- OVS – Maps OVS configuration to P4 Tables ( E.g. Vxlan)
- P4Runtime + Openconfig – Configures P4 tables explicitly (E.g. Container load-balancing)
- Kernel – Maps Kernel configurations (via SAI) to P4 Tables (E.g. ECMP w/ FRR)
- All three control planes can used to program the same P4 target.
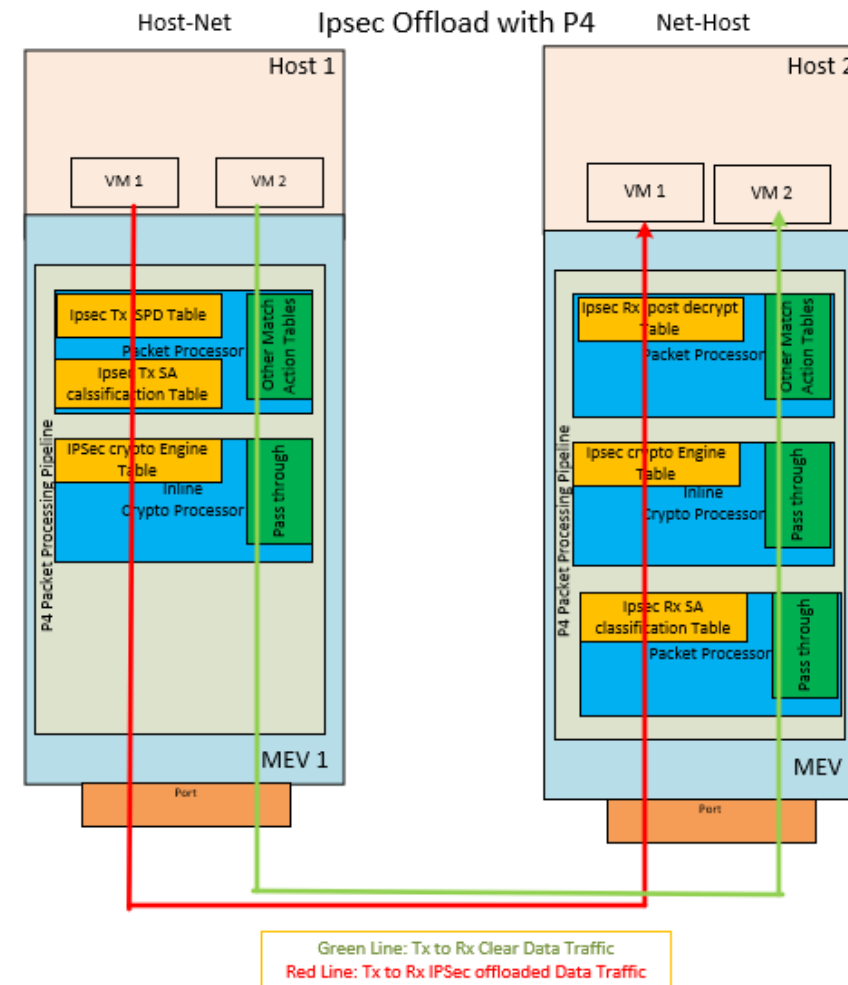- Multiple P4Runtime clients can connect and program different P4 pipelines

**Data Planes:**

- Physical NICs (Tofino, Intel P4 NICs etc)
- P4-DPDK (userspace)
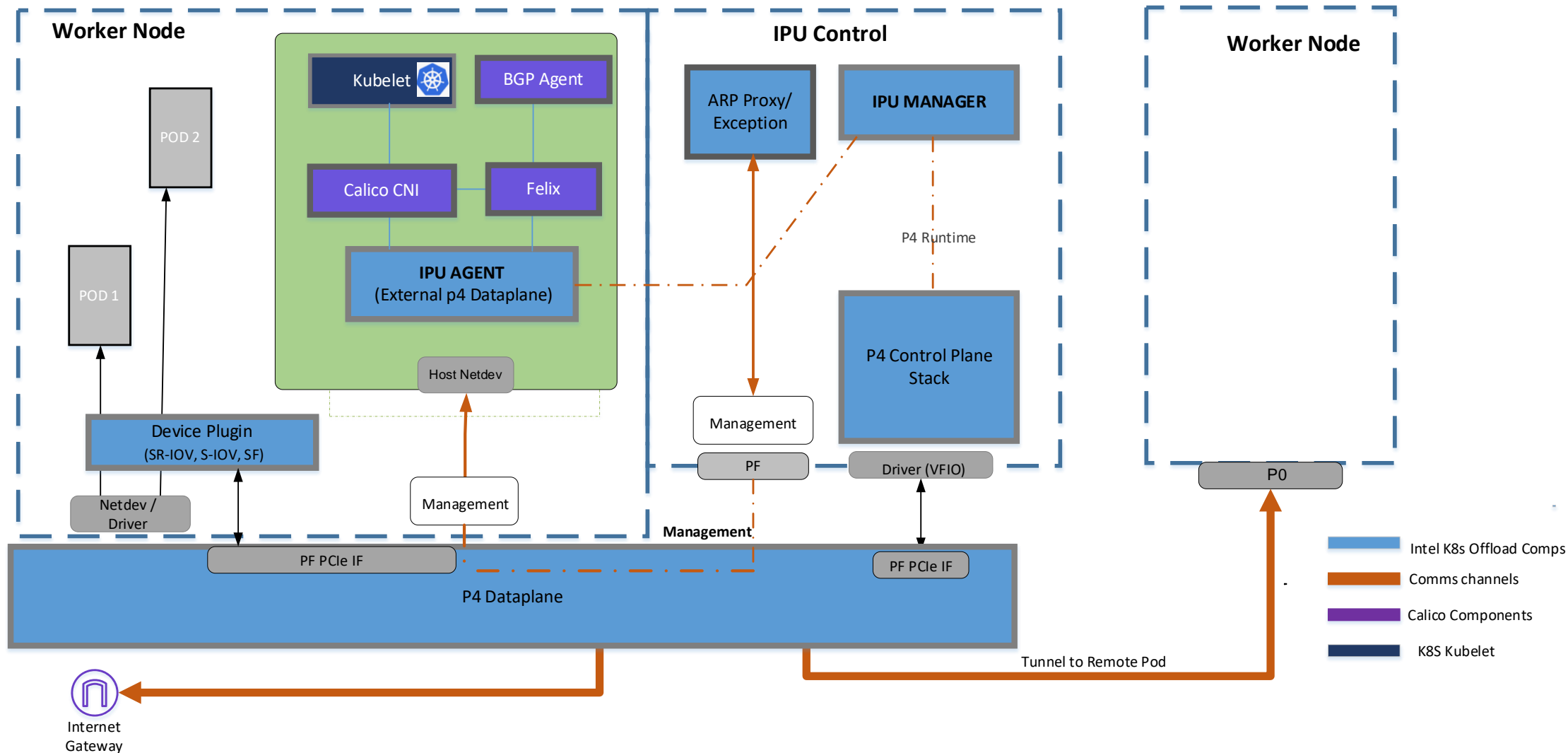- P4-ebpf (kernel)

# P4-programmable IPSec

- Both Tunnel and Transport mode IPSec traffic, along with non-IPSec traffic, can be offloaded simultaneously on the same NIC.
- All traffic passes through the inline crypto processor, and the P4 match-action tables determine which traffic to encrypt or decrypt, while passing the rest as clear traffic.
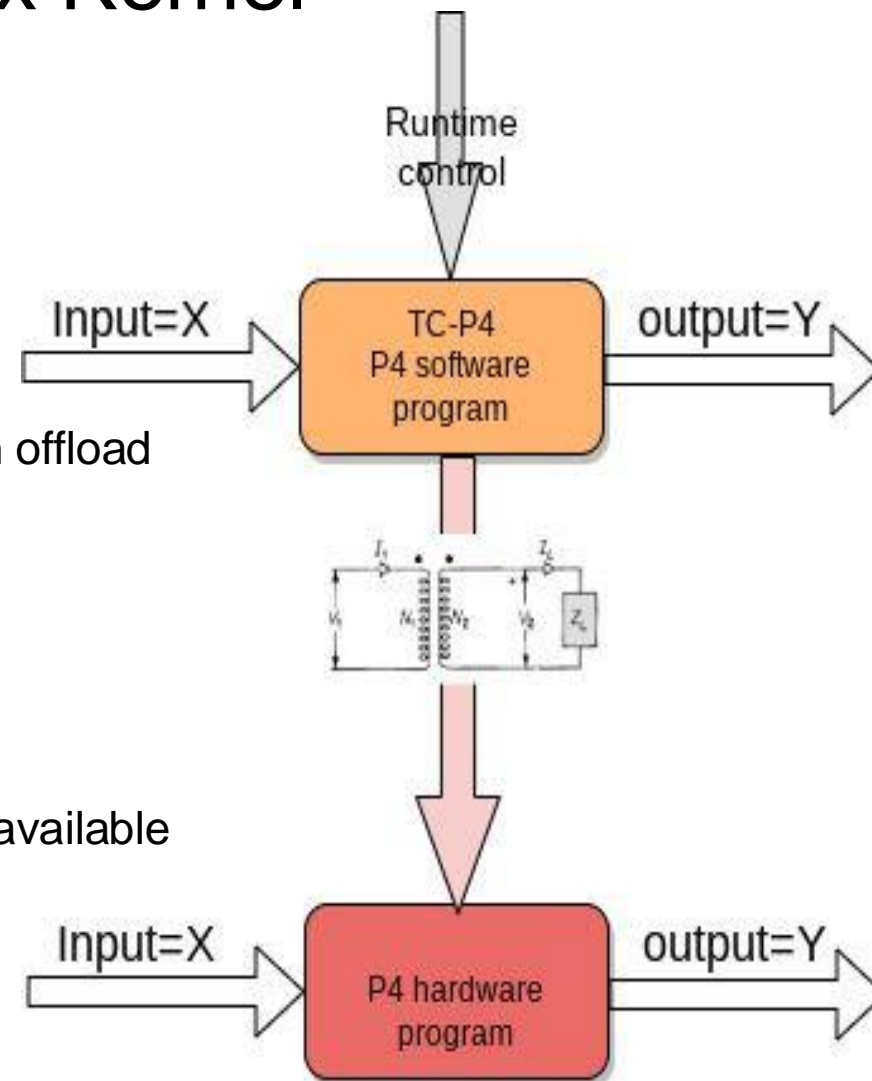- P4 was very useful in supporting features such as custom VXLAN tags



Ipsec Offload with P4

Green Line: Tx to Rx Clear Data Traffic
Red Line: Tx to Rx IPSec offloaded Data Traffic

# P4-enabled K8s offload

# P4-TC: P4 support in Linux Kernel

- ## Datapath definition using P4
  - Generate the datapath for both s/w and vendor h/w
    - Functional equivalence between sw and hw
- ## P4 Linux kernel-native implementation
  - Kernel TC-based software datapath and Kernel-based HW datapath offload
    - Understood Infra tooling which already has deployments
  - Seamless software and hardware symbiosis
  - Functional equivalence whether offloading or s/w datapaths
    - Bare Metal, VMs, or Containers
  - Ideal for datapath specification
    - test in s/w container, VM, etc) then offload when hardware is available

# Our experience with P4

- The learning curve of P4 is possibly a myth!
- More than the expressibility of P4, it's quite often the underlying limitations of the HW that becomes the limiting factor
- P4 leaves a lot of details to be figured out and changed at the compiler level, which is not a bad thing
- The language can benefit from having additional objects recognized, such as queues
- Perfect portability, defined as "taking P4 written for device A, and running it unmodified on device B" may be an elusive target
- New use cases, such as being able to support machine learning, is worth exploring. Some of these newer use cases will require support for not just abstraction of packets, but also abstraction of flows and beyond
- New P4 possibilities are endless. We should keep maintaining the momentum on P4!

# Thank You