# Cisco          P4

Flexibility

Data-Plane
Streaming
Telemetry

Future
Proof

**P4**

**NX-OS**

Scale

Speed &
Agility

ASIC

Switch OS

Run Time API

Driver

ASIC

NX-OS

Auto Generated API

Driver

P4

Tofino

VS.

:       ASIC

:       Match-Action

# Nexus 3000

| 3000/3100 Series | 3200 Series | 3400 Series | 3500 Series | 3600 Series |
|---|---|---|---|---|

| Trident II | Trident II+ | Trident III | Tomahawk | Tomahawk II | Barefoot Tofino | Cisco Monticello | Jericho+ |
|---|---|---|---|---|---|---|---|
| 3048TP | 31108PC/TC-V | 3132C-Z | 3264Q | 3264C-E | 34180YC | 3548/24 | 36180YC-R |
| 3172PQ/TQ | 3132Q-V | | 3232C | | | 3548/24-XL | 3636C-R |
| 3172PQ/TQ-XL | | | | | 3464C | | |
| 31128PQ | | | | | (64x100G) | | |
| 3132Q | | | | | | | |

# Nexus 3400 Series

## 34180YC

48p 10/25G SFP28 + 6p 40/100G QSFP28
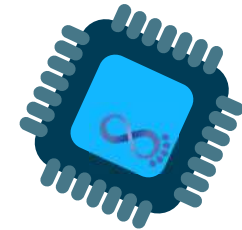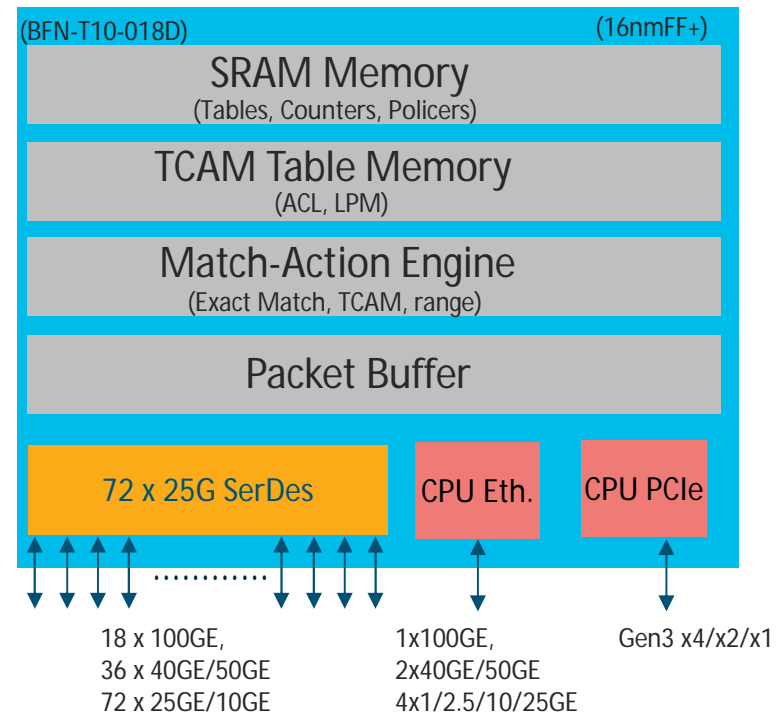Barefoot Tofino 1.8T ASIC

## 3464C

64p 40/100G QSFP28
Barefoot Tofino 6.4T ASIC
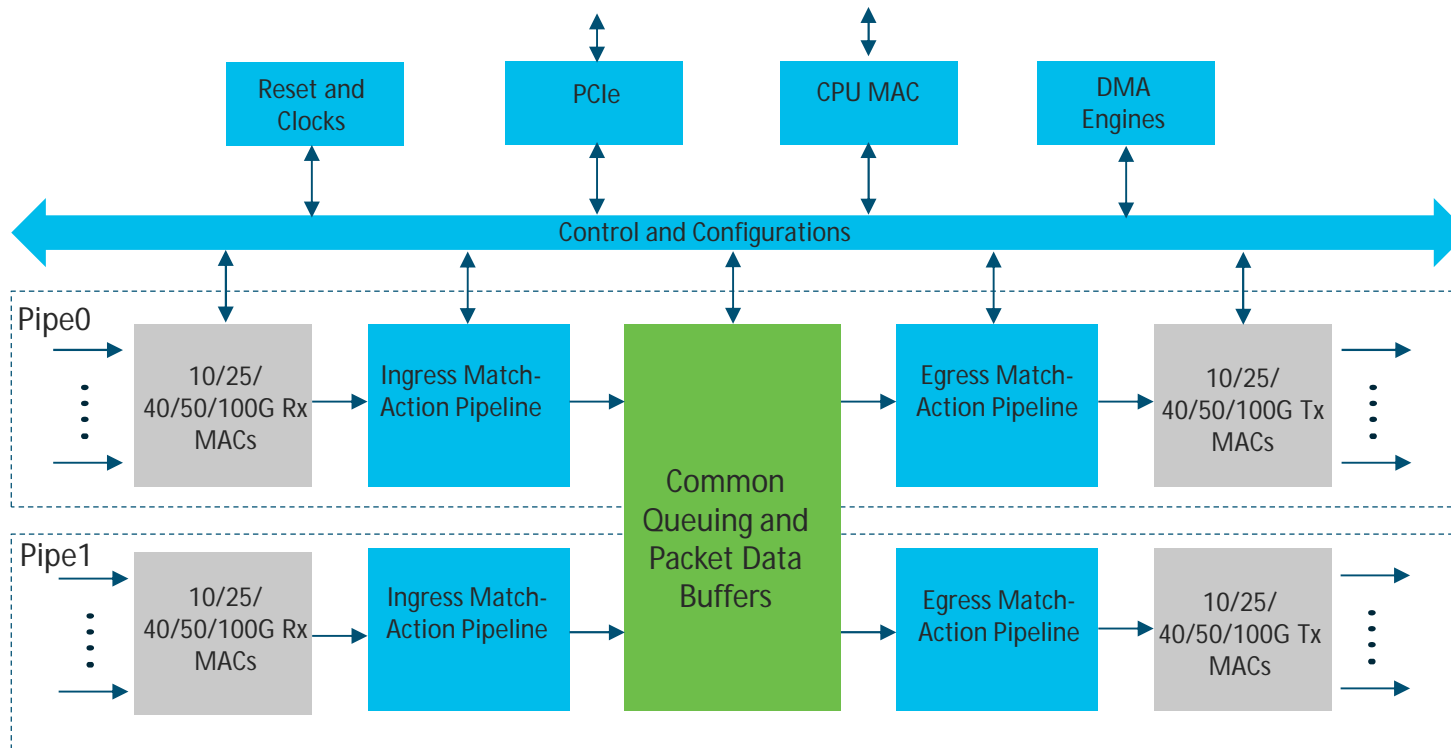
# Barefoot Tofino 1.8T ASIC Architecture

- BFN-T10-018D from Tofino family

- 1.8Tbps Single Chip Ethernet Switch

- 2 Pipes @0.9 Tbps

- P4-programmable pipeline

- Single 20 MB Unified Packet Buffer

- Customizable Low latency

(BFN-T10-018D)                                          (16nmFF+)

**SRAM Memory**
(Tables, Counters, Policers)

**TCAM Table Memory**
(ACL, LPM)

**Match-Action Engine**
(Exact Match, TCAM, range)

**Packet Buffer**

72 x 25G SerDes      CPU Eth.     CPU PCIe

18 x 100GE,            1x100GE,           Gen3 x4/x2/x1
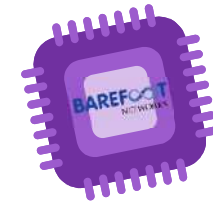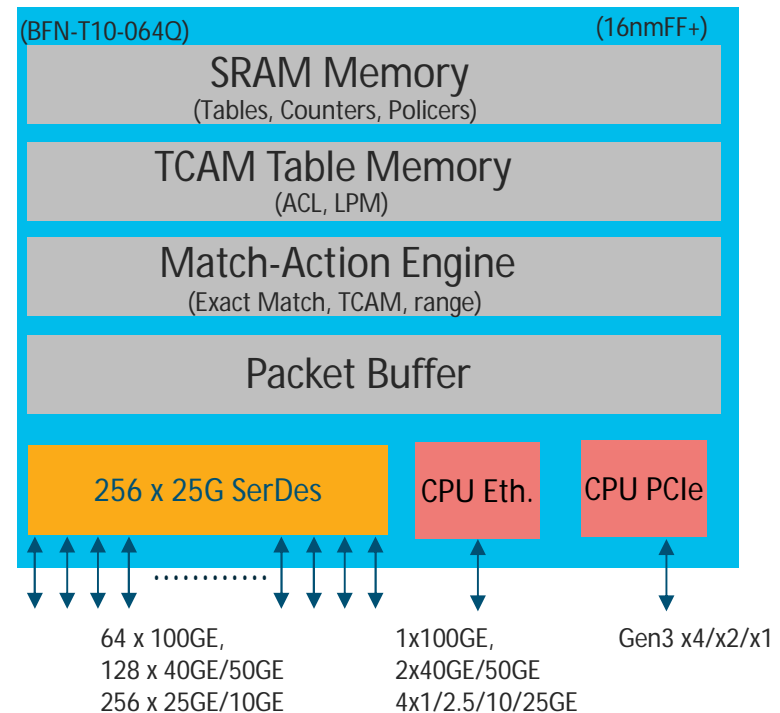36 x 40GE/50GE      2x40GE/50GE
72 x 25GE/10GE      4x1/2.5/10/25GE
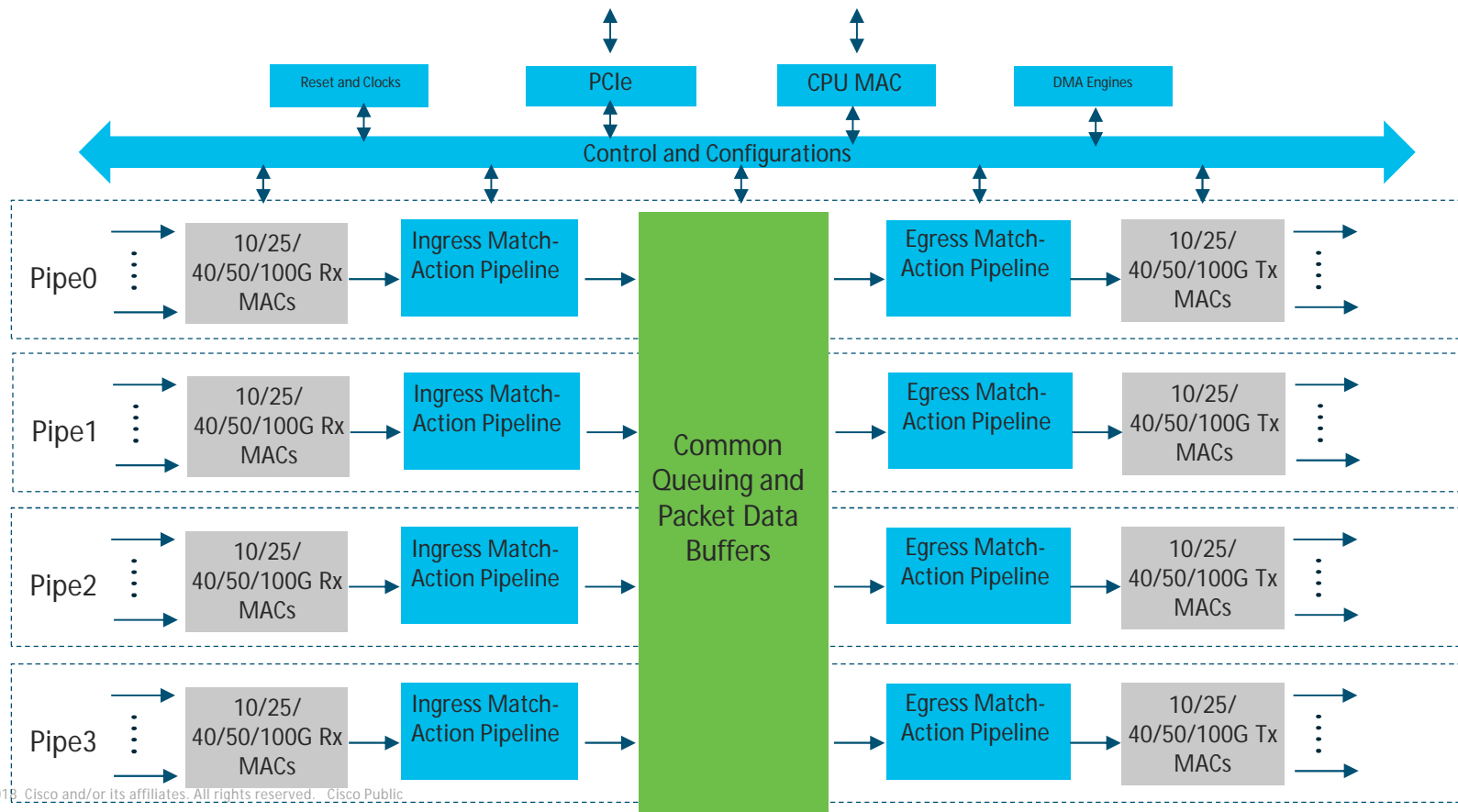
# Tofino 1.8T dual-pipeline Block Diagram

# Barefoot Tofino 6.4T ASIC Architecture

- BFN-T10-064Q from Tofino family

- 6.4 Tbps Single Chip Ethernet Switch

- 4 Pipes @1.6 Tbps

- P4-programmable pipeline

- Single 22 MB Unified Packet Buffer

- Customizable Low latency



(BFN-T10-064Q)                                    (16nmFF+)

**SRAM Memory**
(Tables, Counters, Policers)

**TCAM Table Memory**
(ACL, LPM)

**Match-Action Engine**
(Exact Match, TCAM, range)

**Packet Buffer**

| 256 x 25G SerDes | CPU Eth. | CPU PCIe |

64 x 100GE,
128 x 40GE/50GE
256 x 25GE/10GE

1x100GE,
2x40GE/50GE
4x1/2.5/10/25GE

Gen3 x4/x2/x1

# Tofino 6.4T quad-pipeline Block Diagram



Reset and Clocks

PCIe

CPU MAC

DMA Engines

Control and Configurations

**Pipe0**
10/25/40/50/100G Rx MACs → Ingress Match-Action Pipeline → Common Queuing and Packet Data Buffers → Egress Match-Action Pipeline → 10/25/40/50/100G Tx MACs

**Pipe1**
10/25/40/50/100G Rx MACs → Ingress Match-Action Pipeline → Egress Match-Action Pipeline → 10/25/40/50/100G Tx MACs

**Pipe2**
10/25/40/50/100G Rx MACs → Ingress Match-Action Pipeline → Egress Match-Action Pipeline → 10/25/40/50/100G Tx MACs

**Pipe3**
10/25/40/50/100G Rx MACs → Ingress Match-Action Pipeline → Egress Match-Action Pipeline → 10/25/40/50/100G Tx MACs
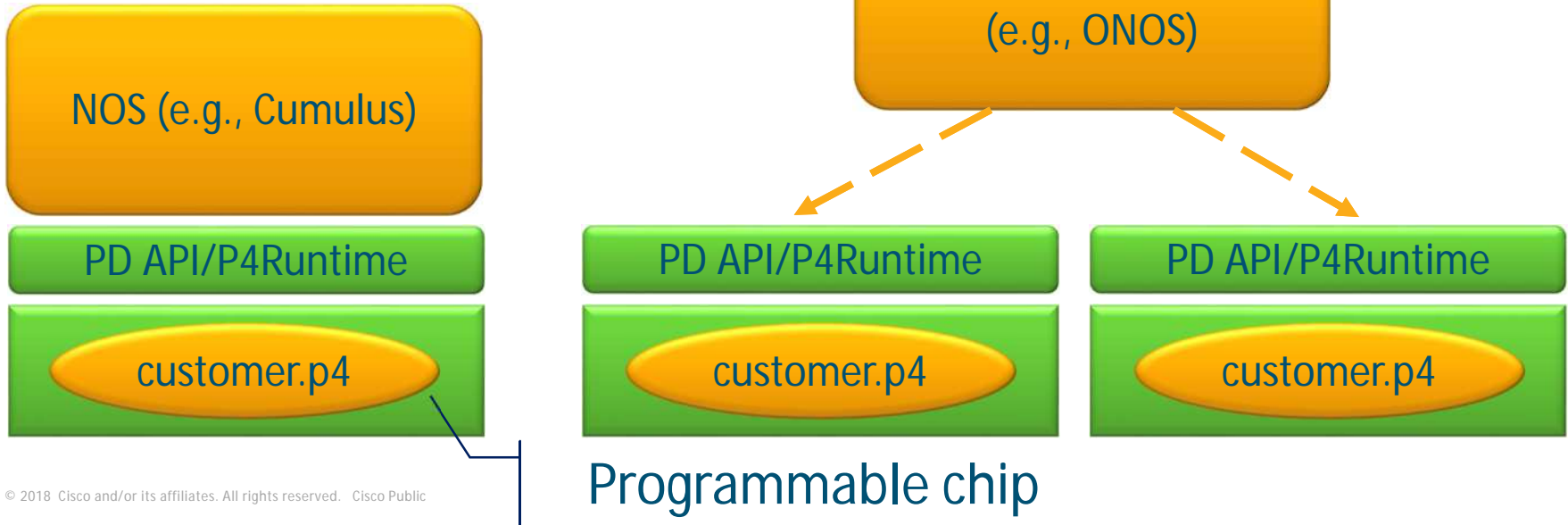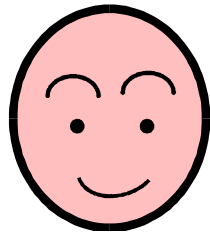
Platform vendor (Cisco)
Chip vendor (Barefoot)
Customer/open source

- Maximum flexibility
- Maximum disruption/risk/work

**NOS (e.g., Cumulus)**

PD API/P4Runtime

customer.p4

**Remote controller/NOS (e.g., ONOS)**

PD API/P4Runtime
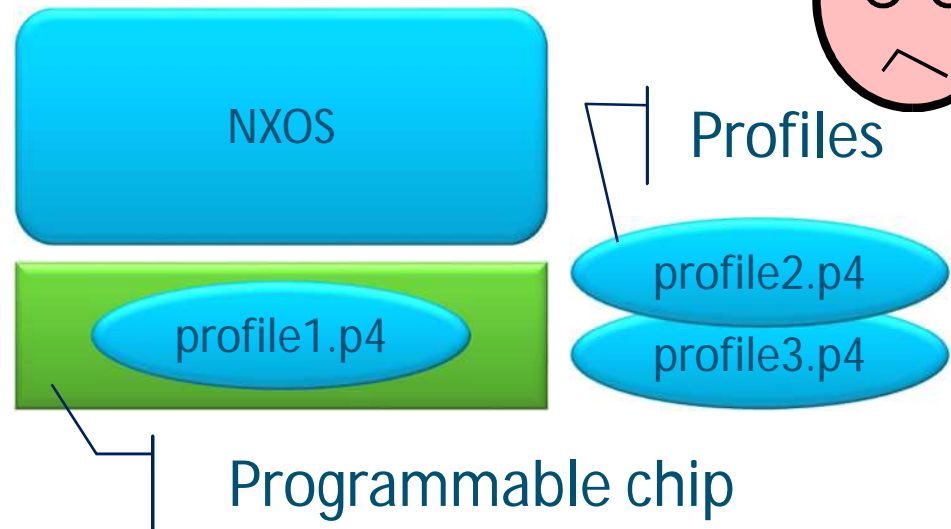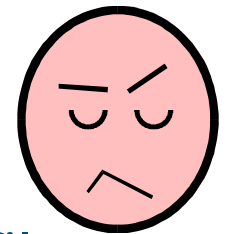
customer.p4

PD API/P4Runtime

customer.p4

Programmable chip

- Deployment as usual
  - Familiar features and interfaces
- Resource optimization
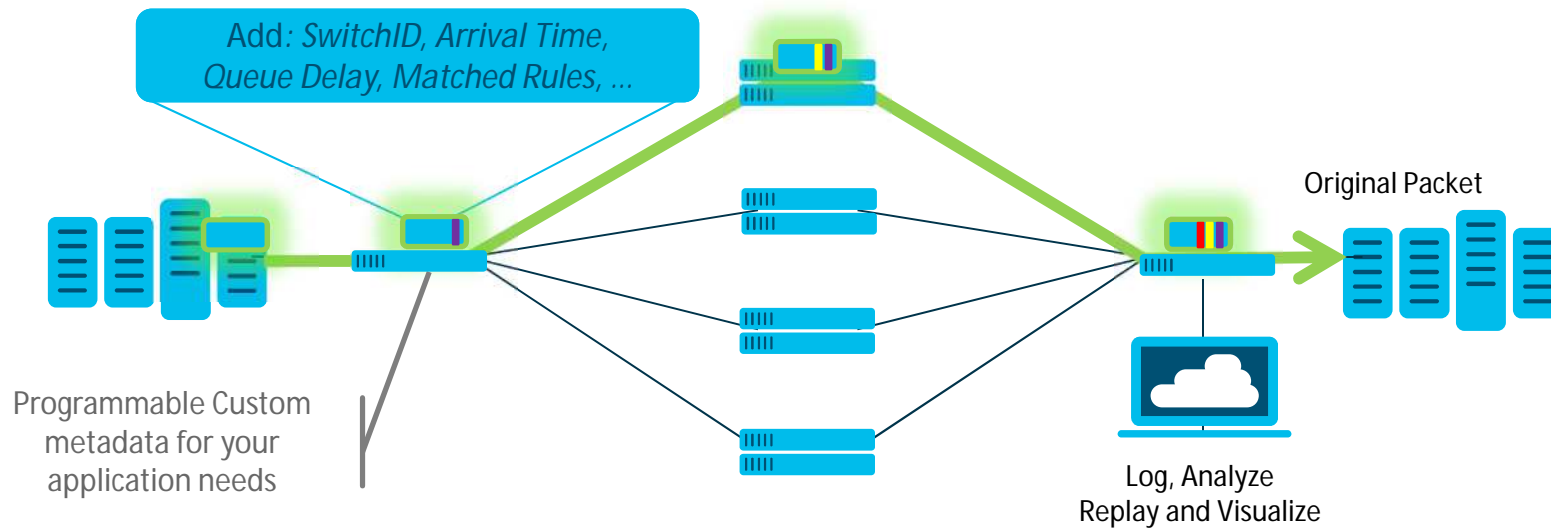- Future proof
- Feature agility
- Streaming telemetry

Platform vendor (Cisco)
Chip vendor (Barefoot)
Customer/open source

- No flexibility
  - No custom feature and protocol support

NXOS

Profiles

profile1.p4

profile2.p4

profile3.p4

Programmable chip

# Inband Network Telemetry (INT)

Add: *SwitchID, Arrival Time, Queue Delay, Matched Rules, ...*

Original Packet

Programmable Custom metadata for your application needs

Log, Analyze Replay and Visualize

# Inband Network Telemetry (INT)

- INT Source

      INT

- INT Transit         2

- INT Sink         3

- INT Sink     INT         INT

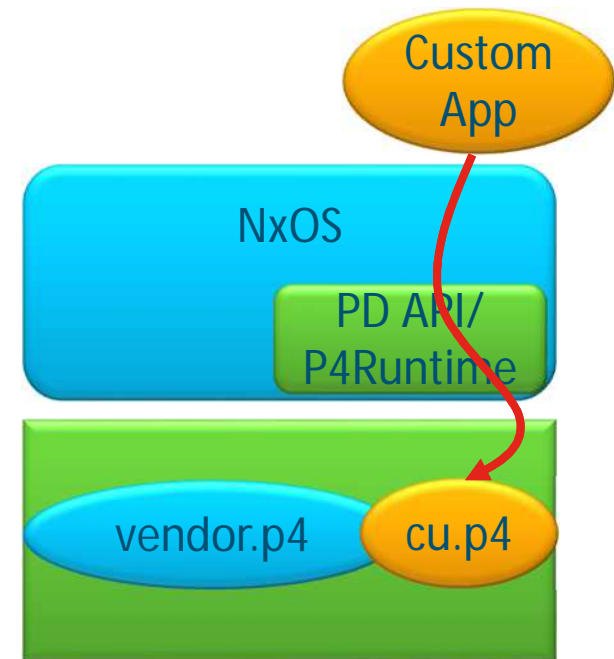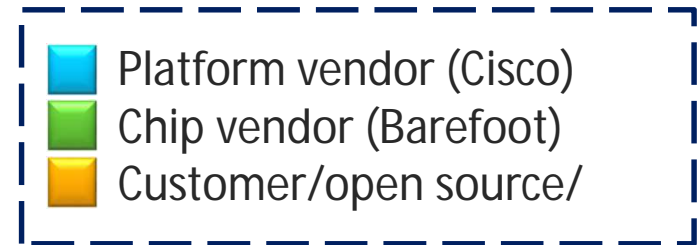| Flow Watch List (zoom-in view per 5-tuple of flow + DSCP bits) – 2K | Drop Watch List (Drop due to various drop reasons) - 512 | Queue Watch List (queue depth or latency exceeds configured threshold) |
|---|---|---|
| Switch ID | Switch ID | Switch ID |
| Hop latency | Ingress Port ID | Hop latency |
| Queue ID + Queue occupancy | Egress Port ID | Queue ID + Queue occupancy |
| Ingress timestamp | Queue ID | |
| Egress timestamp | Drop Reason | |

- Node-to-Node: Reserved DSCP bit will be inserted temporarily in data packets to indicate that packets also carry INT data
- Node-to-Collector: A UDP encapsulation is used to pack collected INT stack at INT Sink and send to collector. Flow-affinity is maintained to send same flow-record to same collector for easy processing
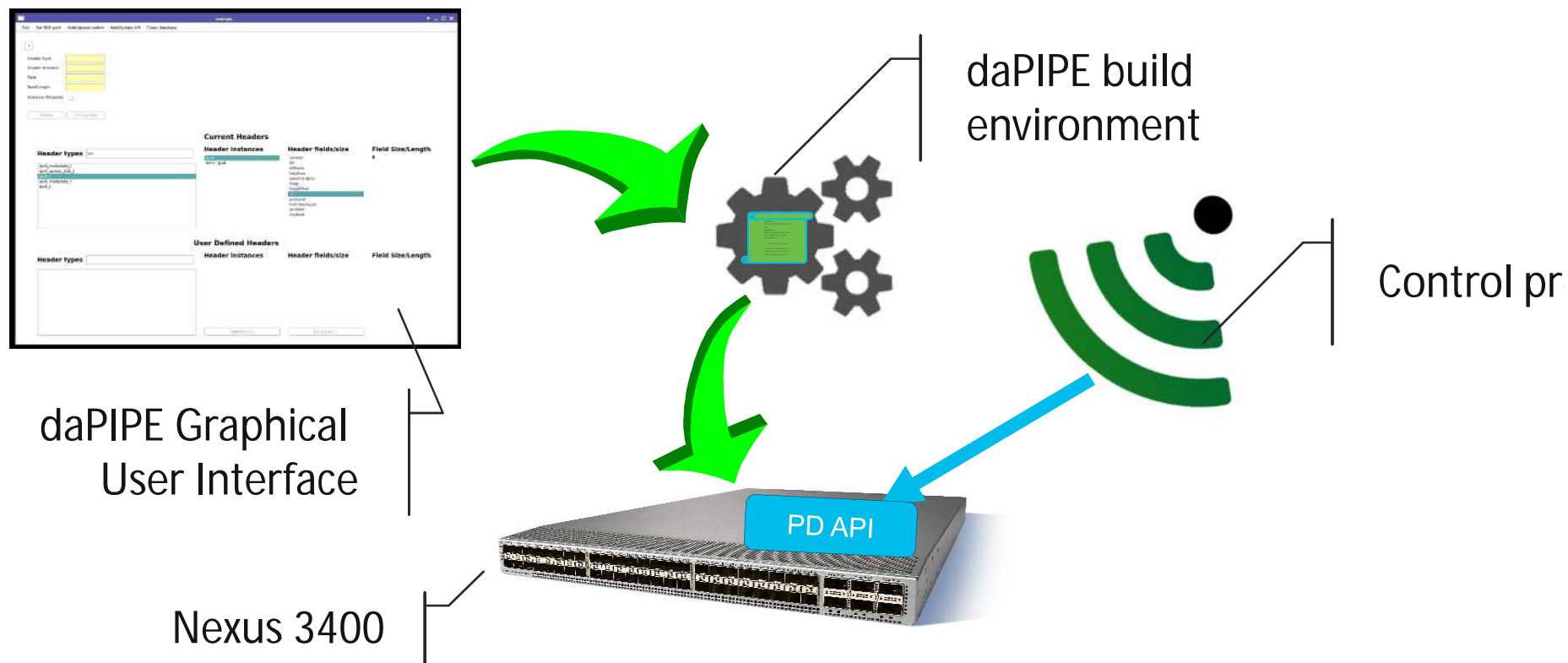
- Best of breed

- Deployment as usual
  - Familiar features and interfaces

- Minimum development effort
  - Leverage existing functions in building new features

**Minimize disruption and risk!**

Platform vendor (Cisco)
Chip vendor (Barefoot)
Customer/open source/

Custom App

NxOS

PD API/ P4Runtime

vendor.p4   cu.p4

# data Plane Incremental Programming Environment



daPIPE build environment

Control pr

PD API

daPIPE Graphical User Interface

Nexus 3400

# Customer Programming Workflow



Cu.c

Favorite SDE

Cu.exe

NetOS

NxAPI

PD-API.o

Development environment

vendor.p4

Cu.p4

Constraint Checker

P4 Compiler

Data_plane.bin

20

# Operating System Support



**Cisco Apps**
- BGP
- OSPF

**Customer Apps**
- Cfg
- Ctrl plane

Guest Shell (container)

SW (mostly) control plane

Infrastructure

HAL

NXOS

Controlled data plane API access

APIs generated by compiling P4

HW data plane

Cisco.p4

Cu.p4

Programmable ASIC