

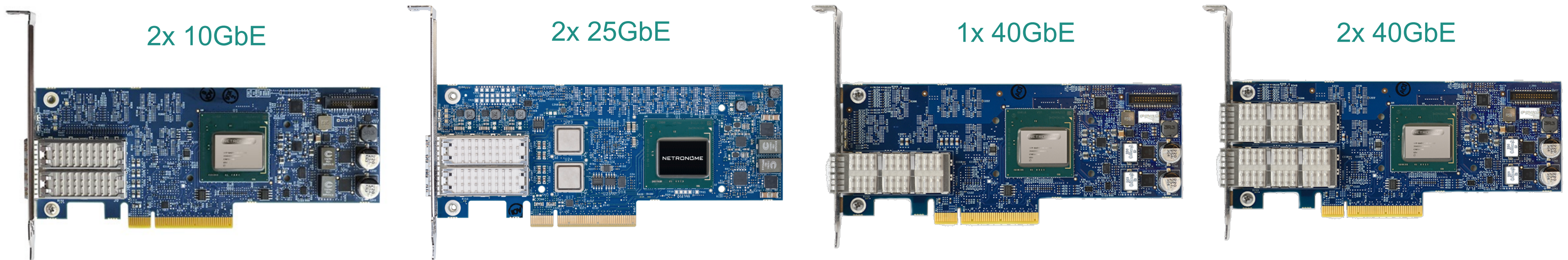
# SmartNIC Programming Models

Johann Tönsing

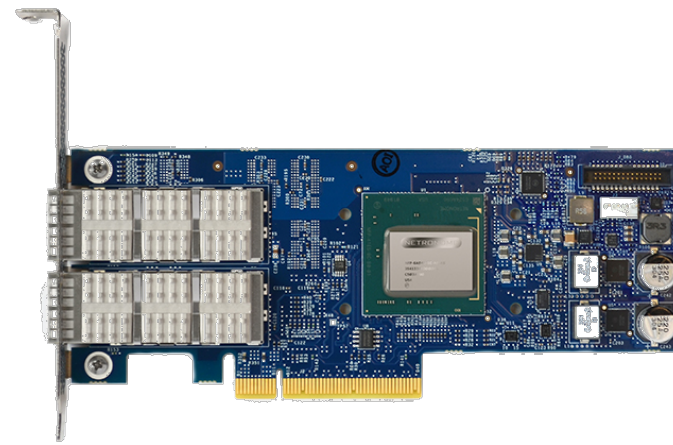
2016-11-09

- SmartNIC hardware
- Pre-programmed vs. custom (C and/or P4) firmware
- Programming models / offload models
- Switching on NIC, with SR-IOV / virtio data delivery
- SmartNIC performance + TCO
- Silicon and datapath software architectures
- Example code

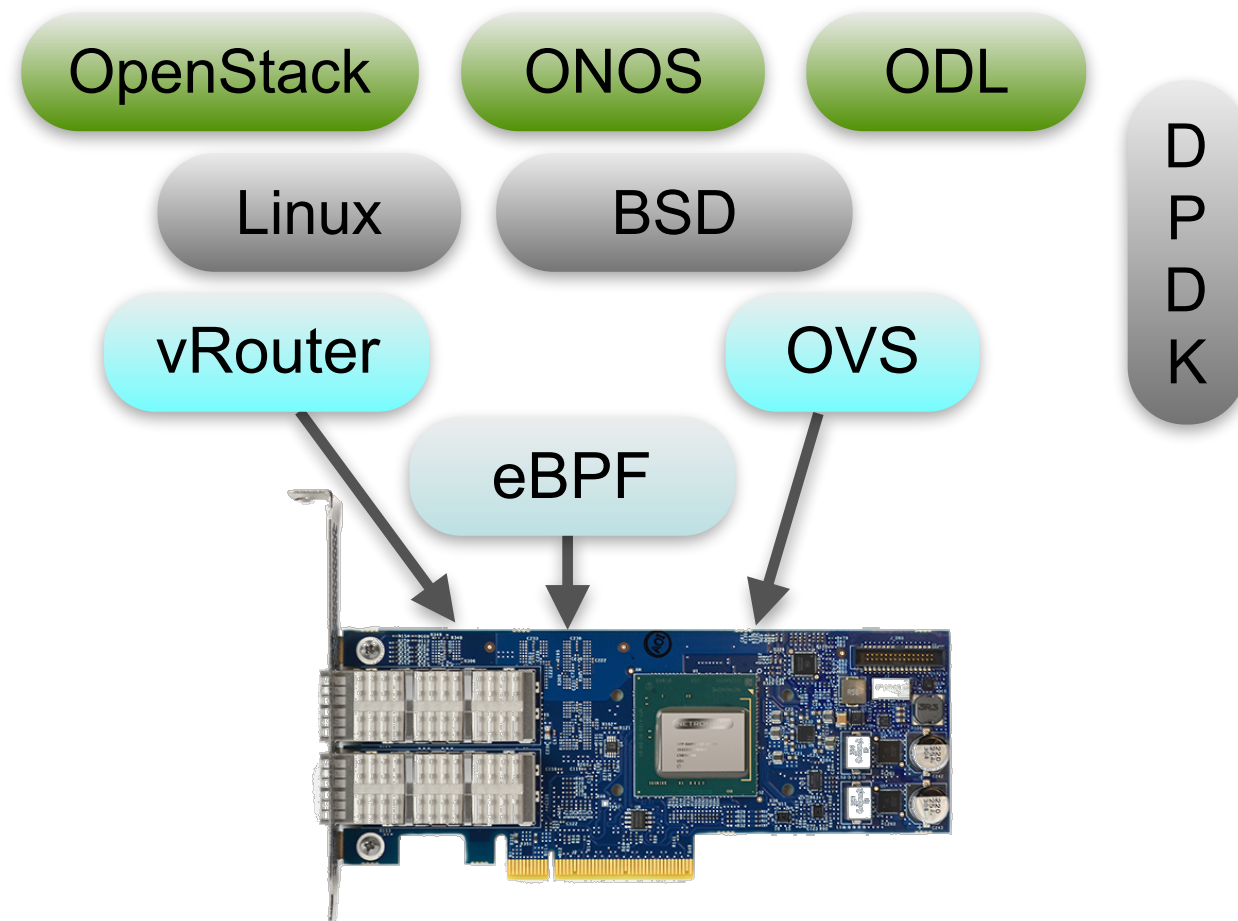
- Optimized for standard server based cloud data centers
- Low Profile Half Length PCIe form factor, power < 25W
- Based on Netronome's NFP-4xxx silicon (72 C programmable cores, 8 threads each)
- 2GB DRAM for lookup tables / state tables (millions of entries)
- *Dataplane fully implemented in software*



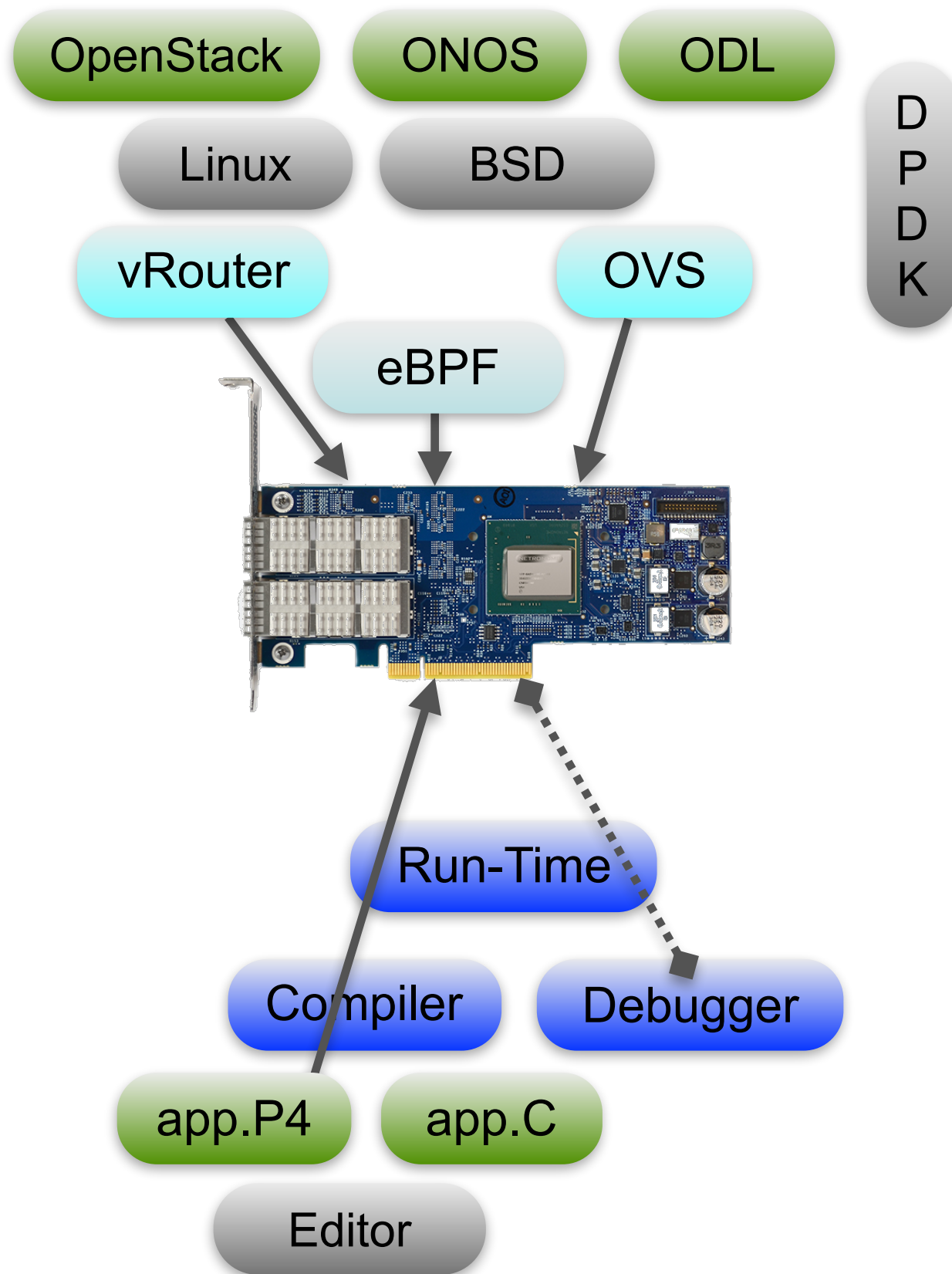
Also available: Agilio™ LX 2x40G / 1x100G with dual PCIe interfaces, 120 cores, 8GB DRAM...



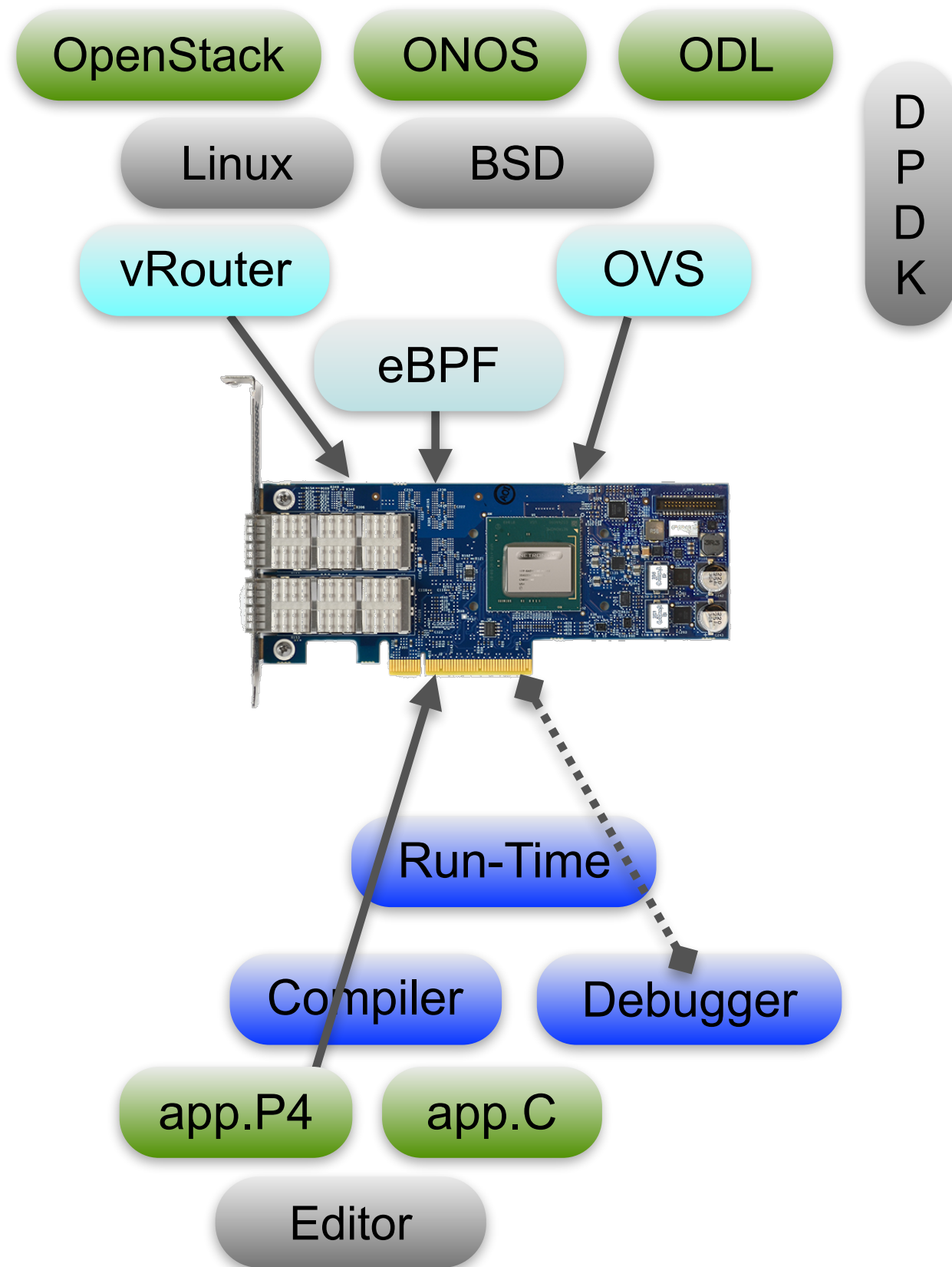
- **SmartNIC with dynamically downloadable firmware**



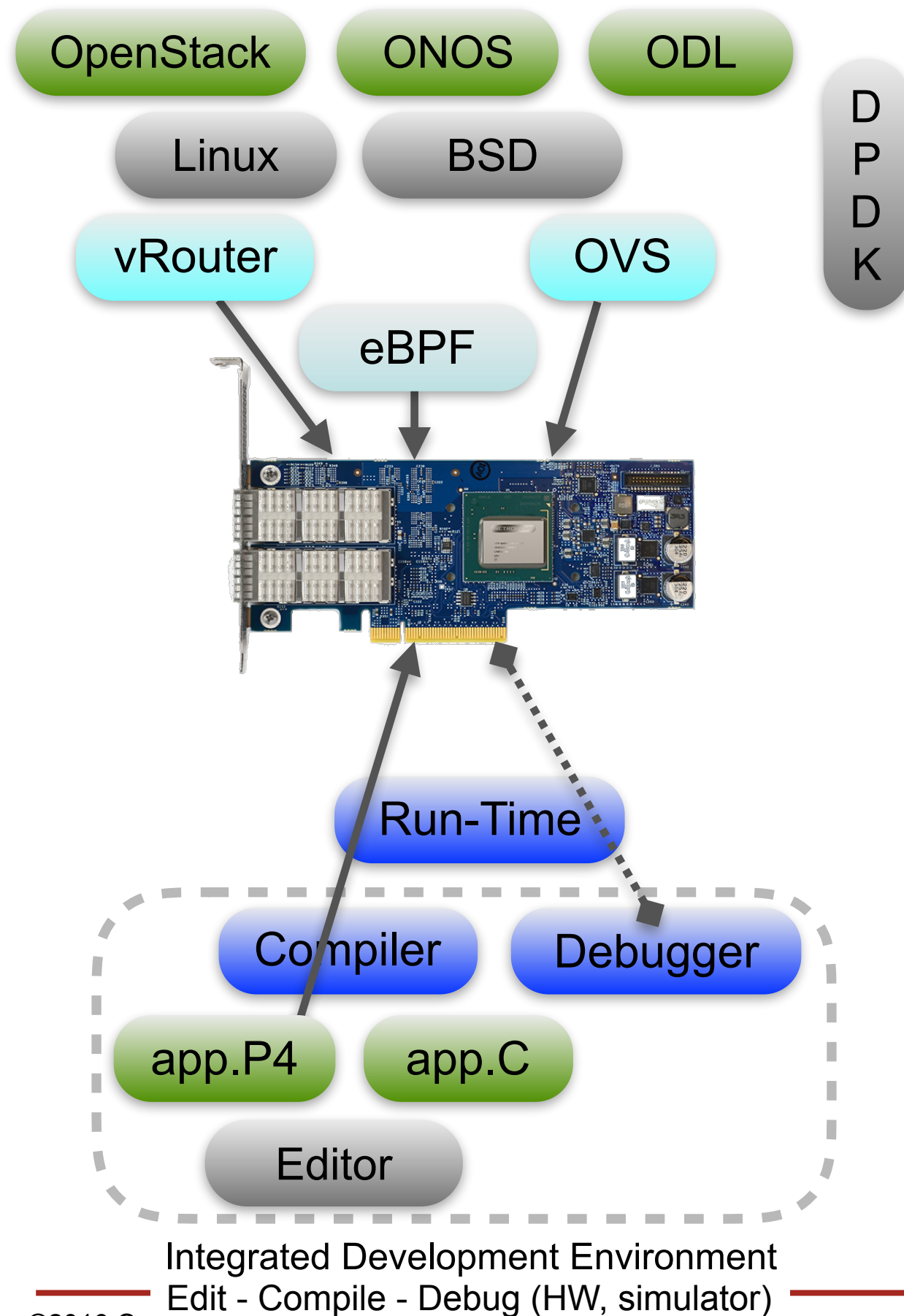
- **OVS / vRouter / eBPF+XDP datapath on host can be accelerated by SmartNIC**
  - *Firmware / drivers supplied by Netronome*
- **SmartNIC with dynamically downloadable firmware**



- **OVS / vRouter / eBPF+XDP datapath on host can be accelerated by SmartNIC**
  - *Firmware / drivers supplied by Netronome*
- **SmartNIC with dynamically downloadable firmware**
- **Firmware can be developed in P4 and/or C**
  - *Custom dataplane*



- **OVS / vRouter / eBPF+XDP datapath on host can be accelerated by SmartNIC**
  - *Firmware / drivers supplied by Netronome*
- **SmartNIC with dynamically downloadable firmware**
- **Firmware can be developed in P4 and/or C**
  - *Custom dataplane*
- **Hybrid - “sandbox / plugin” concept**
  - *Example: C plugin embedded in P4*

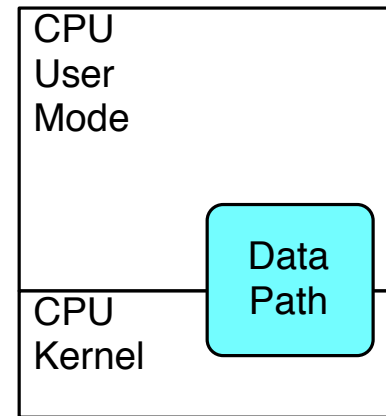


- **OVS / vRouter / eBPF+XDP datapath on host can be accelerated by SmartNIC**
  - *Firmware / drivers supplied by Netronome*
- **SmartNIC with dynamically downloadable firmware**
- **Firmware can be developed in P4 and/or C**
  - *Custom dataplane*
- **Hybrid - “sandbox / plugin” concept**
  - *Example: C plugin embedded in P4*



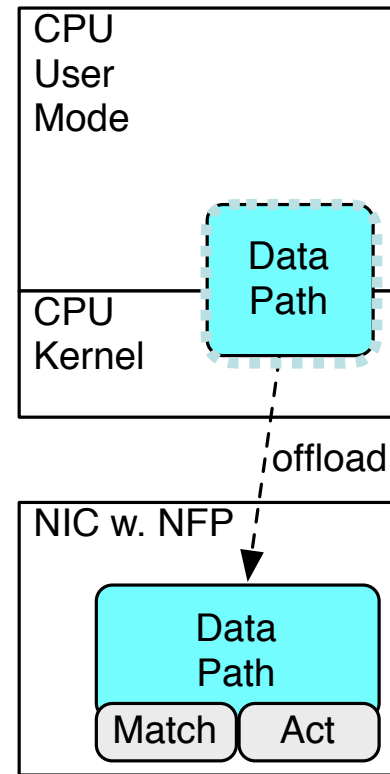
# Programming Models

**Offload / Acceleration**  
(transparent -  
drop in replacement)



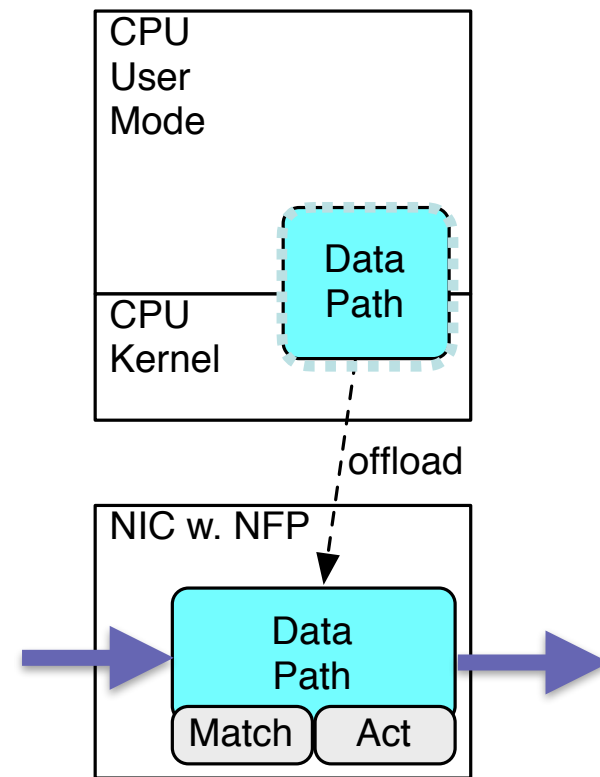
# Programming Models

## Offload / Acceleration (transparent - drop in replacement)



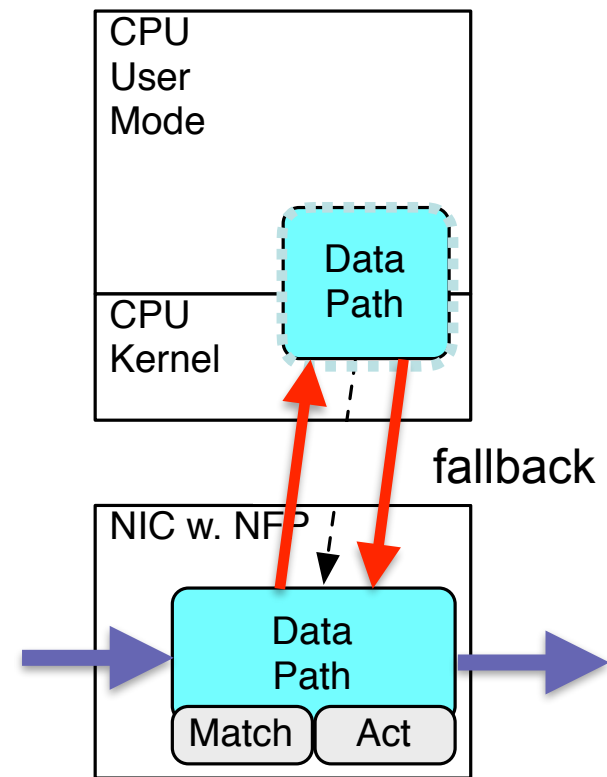
# Programming Models

**Offload / Acceleration**  
(transparent -  
drop in replacement)



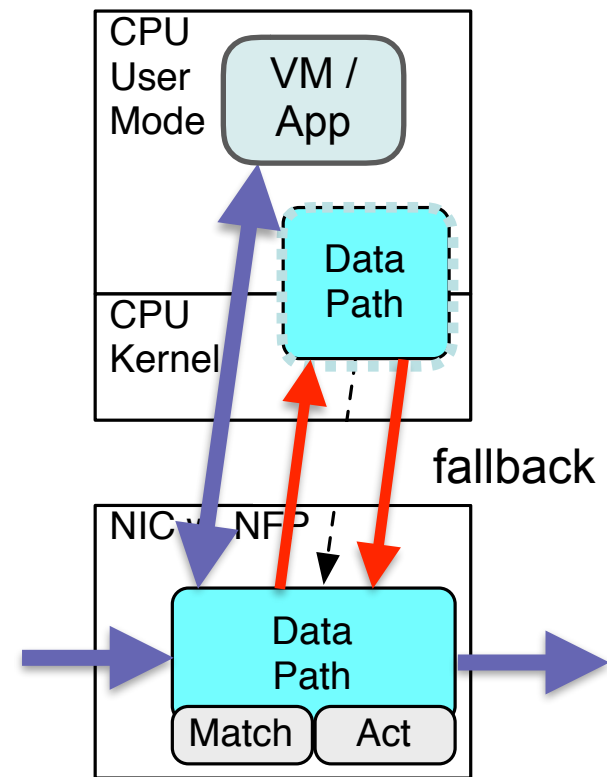
# Programming Models

**Offload / Acceleration**  
(transparent -  
drop in replacement)



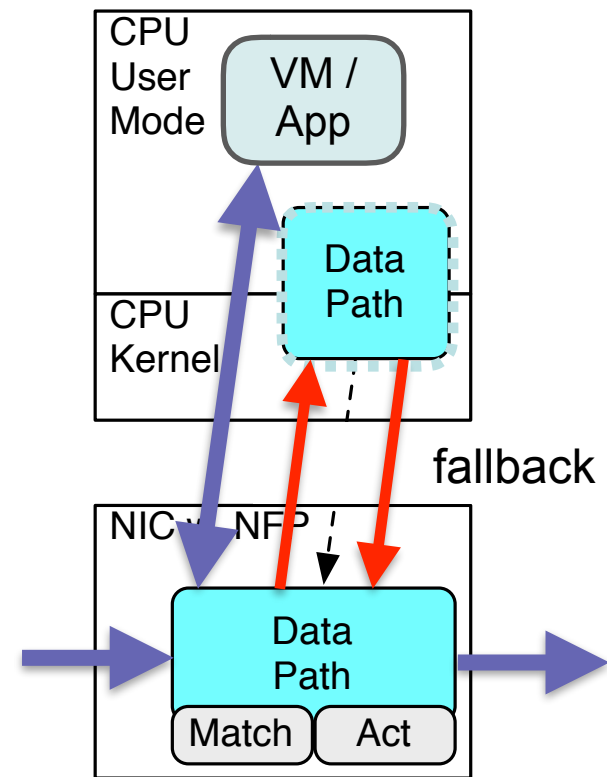
# Programming Models

**Offload / Acceleration**  
(transparent -  
drop in replacement)



# Programming Models

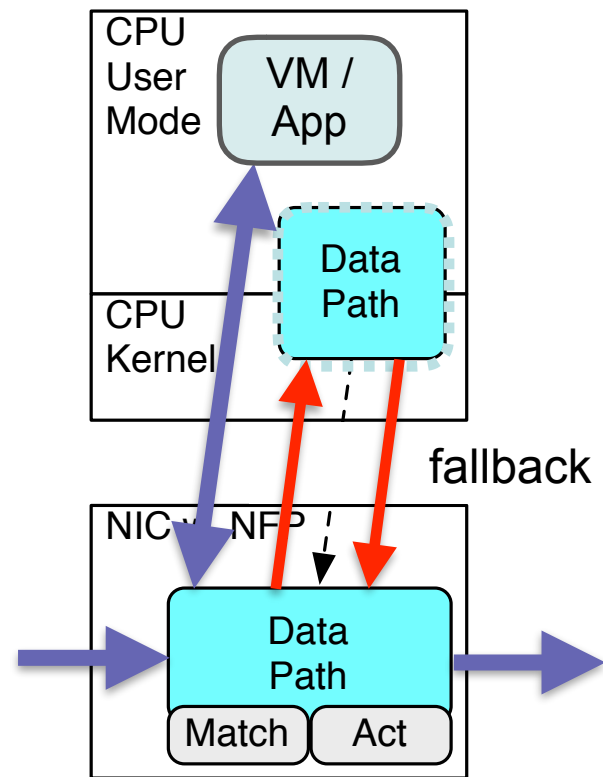
**Offload / Acceleration**  
(transparent -  
drop in replacement)



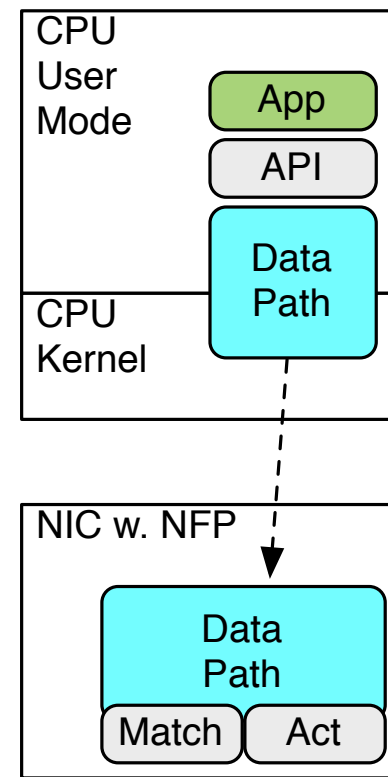
- t OVS/OpenFlow  
Match, Action,  
Tunnels
- Contrail vRouter  
Match, Action,  
Tunnels
- Contrack  
(Firewall)
- eBPF / XDP
- ...

# Programming Models

## Offload / Acceleration (transparent - drop in replacement)



## CPU Apps Calling APIs Compatible / Open Sourced vs. Vendor Extension

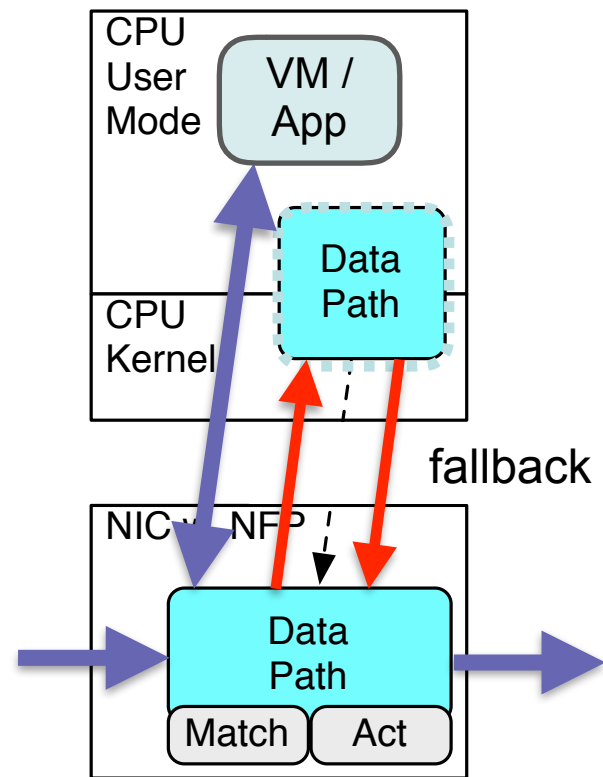


- t OVS/OpenFlow  
Match, Action,  
Tunnels
- Conrail vRouter  
Match, Action,  
Tunnels
- Contrack  
(Firewall)
- eBPF / XDP
- ...



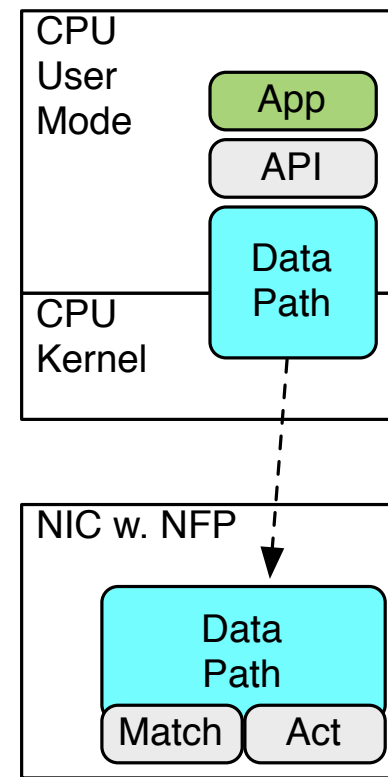
# Programming Models

## Offload / Acceleration (transparent - drop in replacement)



- t OVS/OpenFlow  
Match, Action,  
Tunnels
- Contrail vRouter  
Match, Action,  
Tunnels
- Contrack  
(Firewall)
- eBPF / XDP
- ...

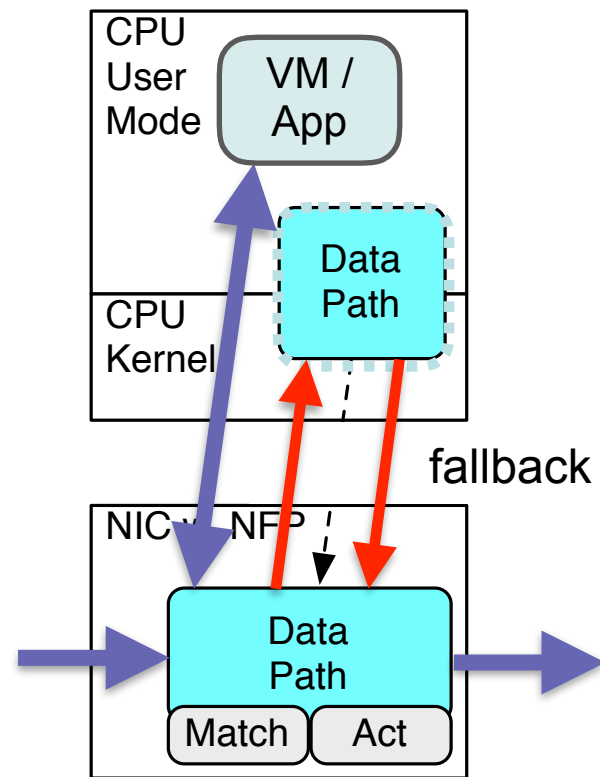
## CPU Apps Calling APIs Compatible / Open Sourced vs. Vendor Extension



- t DPDK  
Poll Mode Driver
- eBPF / XDP APIs
- Flow APIs -  
Match / Act / Tunnel
- Load Balancing APIs
- Crypto APIs
- ...

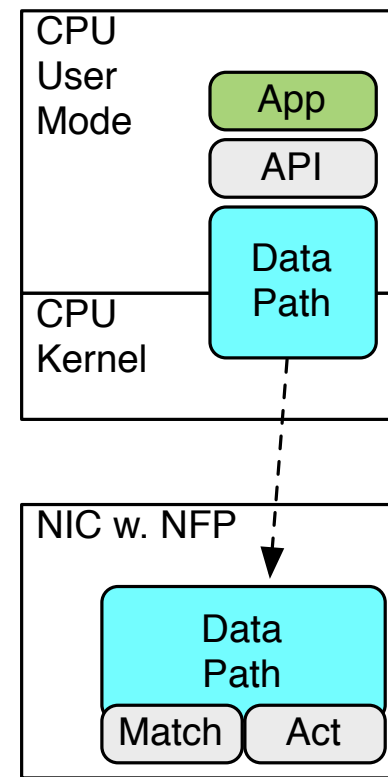
# Programming Models

## Offload / Acceleration (transparent - drop in replacement)



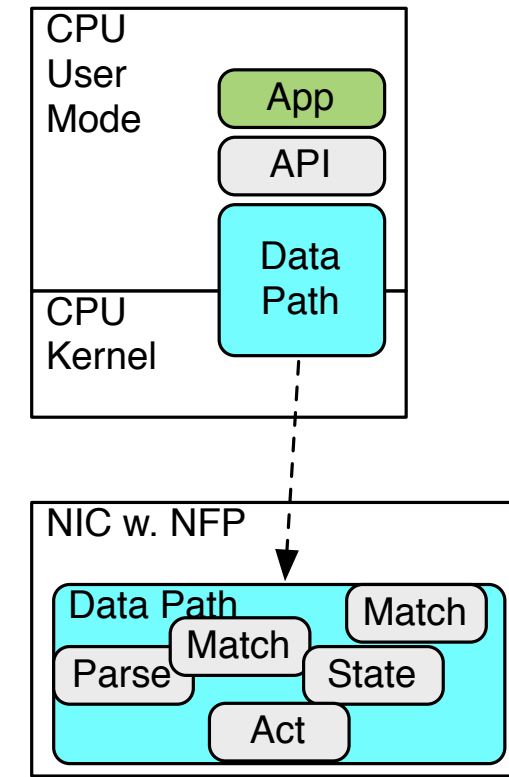
- t OVS/OpenFlow  
Match, Action, Tunnels
- Contrail vRouter  
Match, Action, Tunnels
- Contrack  
(Firewall)
- eBPF / XDP
- ...

## CPU Apps Calling APIs Compatible / Open Sourced vs. Vendor Extension



- t DPDK  
Poll Mode Driver
- eBPF / XDP APIs
- Flow APIs -  
Match / Act / Tunnel
- Load Balancing APIs
- Crypto APIs
- ...

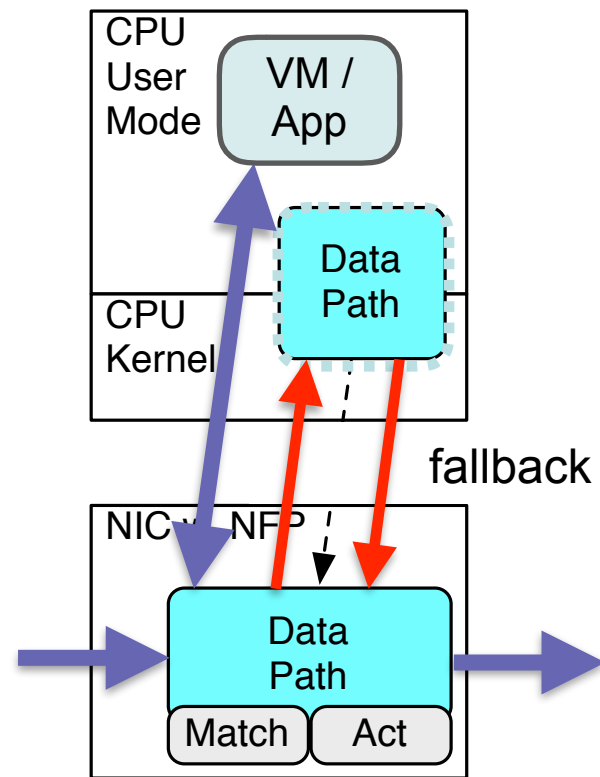
## Flexible Datapath Abstraction OpenFlow 2.x, P4, PIF, eBPF...



- t Protocol agnostic  
flexible parsing
- Arbitrary arrangement  
of matching tables
- Matching without  
tables
- State storage / retrieval
- Complex actions
- Event handling

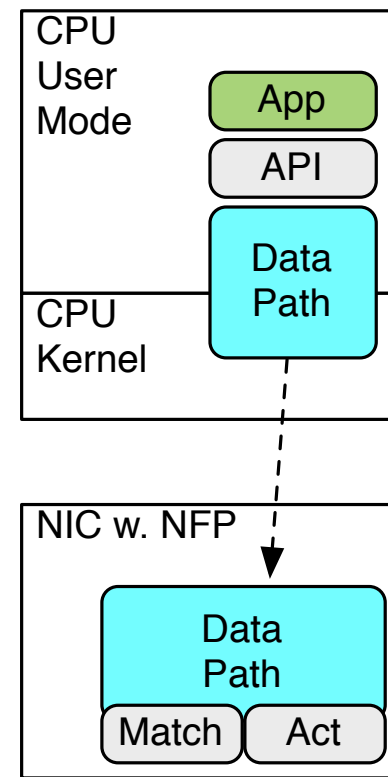
# Programming Models

**Offload / Acceleration**  
(transparent - drop in replacement)



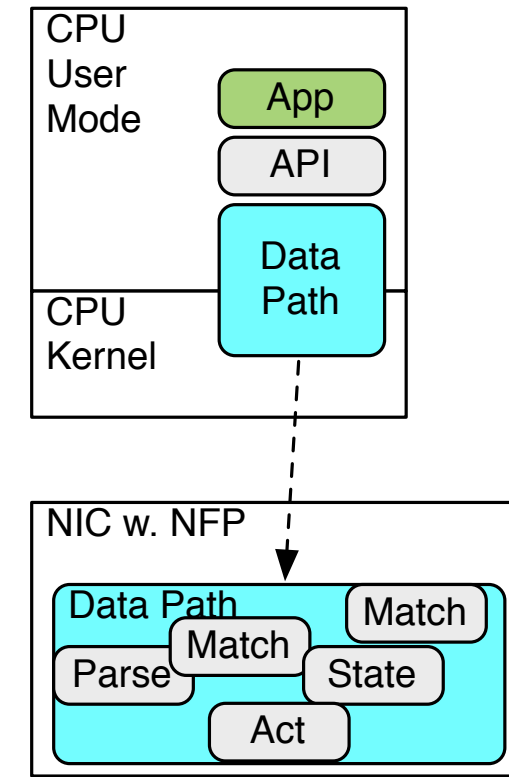
- t OVS/OpenFlow  
Match, Action, Tunnels
- Contrail vRouter  
Match, Action, Tunnels
- Contrack (Firewall)
- eBPF / XDP
- ...

**CPU Apps Calling APIs**  
Compatible / Open Sourced vs. Vendor Extension



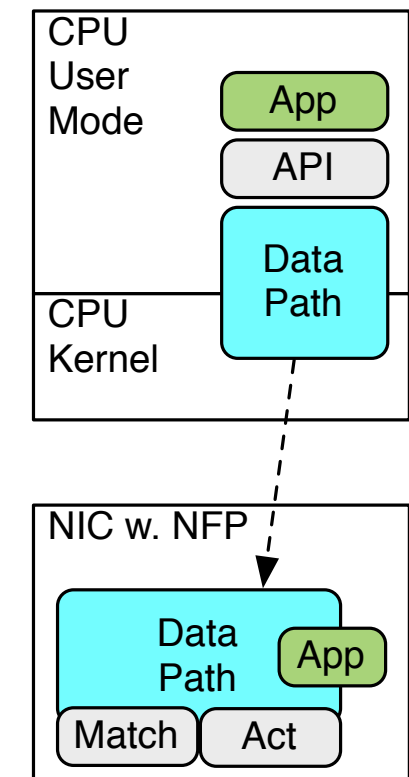
- t DPDK  
Poll Mode Driver
- eBPF / XDP APIs
- Flow APIs -  
Match / Act / Tunnel
- Load Balancing APIs
- Crypto APIs
- ...

**Flexible Datapath Abstraction**  
OpenFlow 2.x, P4, PIF, eBPF...



- t Protocol agnostic  
flexible parsing
- Arbitrary arrangement  
of matching tables
- Matching without  
tables
- State storage / retrieval
- Complex actions
- Event handling

**Hybrid: Datapath Extensions in CPU / NFP**  
In C, P4 / PIF, ...



- t Custom tunnel
- Custom action
- Custom matching
- ...

# Traditional Model: SR-IOV



(Nova, Neutron)

x86 Userspace

1

SR-IOV Configuration

Virtual Machine

Apps

netdev or DPDK

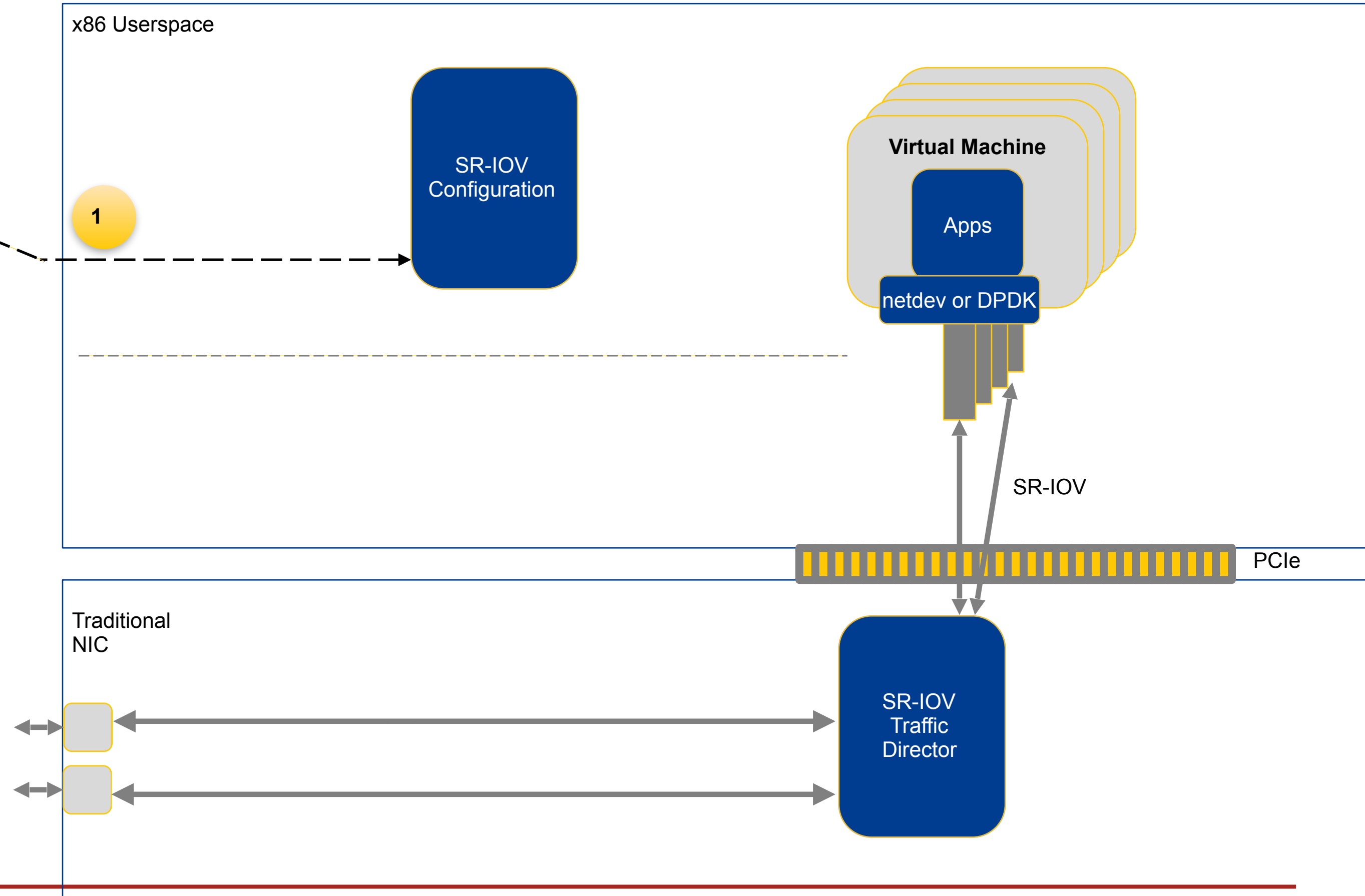
SR-IOV

PCIe

Traditional NIC

SR-IOV Traffic Director

1 Configuration from cloud management system



# Traditional Model: SR-IOV



(Nova, Neutron)

x86 Userspace

1

SR-IOV Configuration

2

Virtual Machine

Apps

netdev or DPDK

SR-IOV

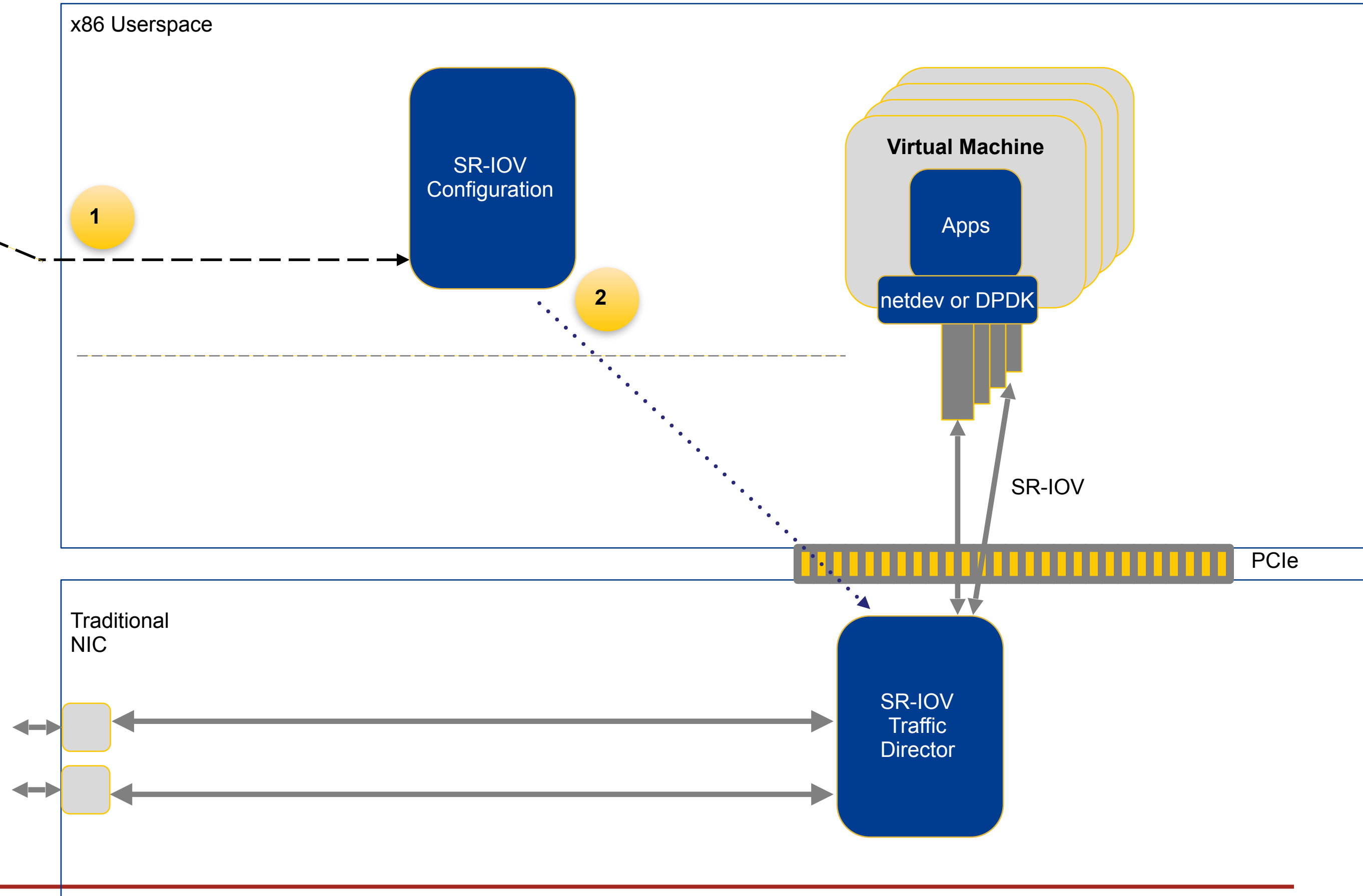
PCIe

Traditional NIC

SR-IOV Traffic Director

1 Configuration from cloud management system

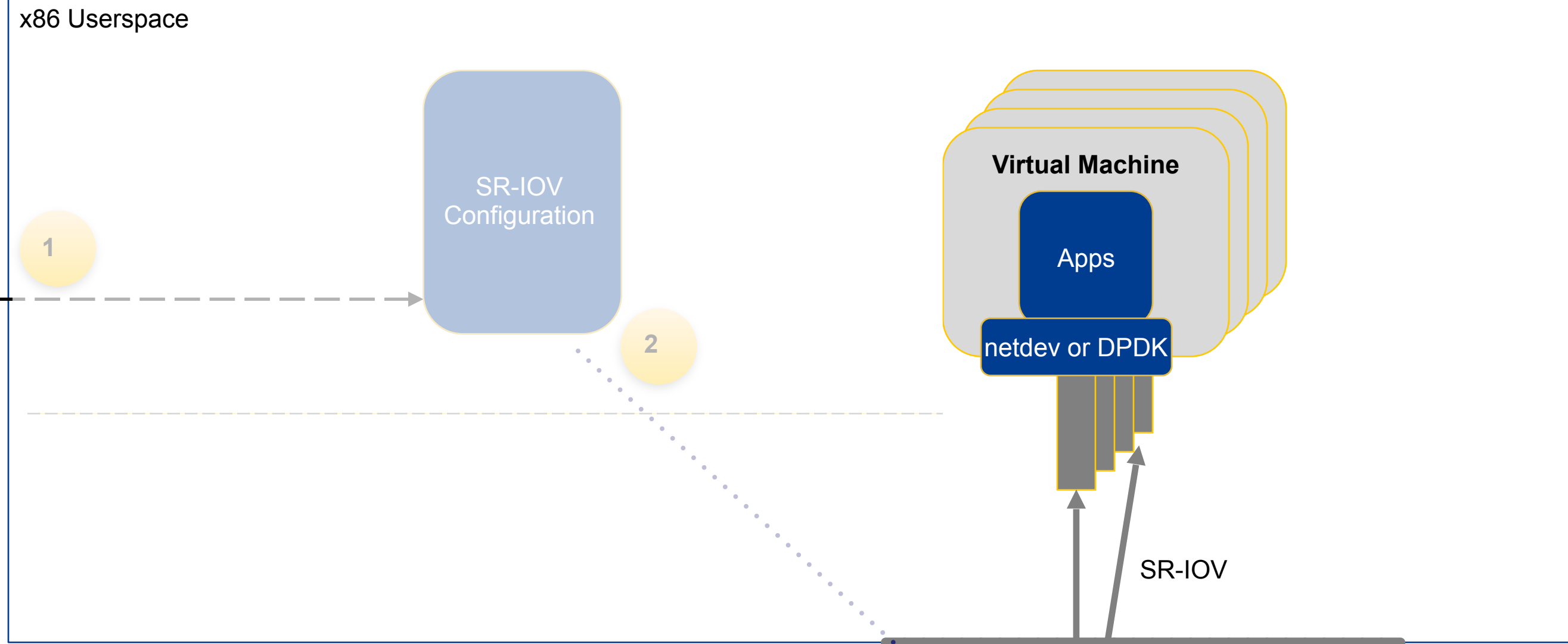
2 Hardware configuration



# Traditional Model: SR-IOV

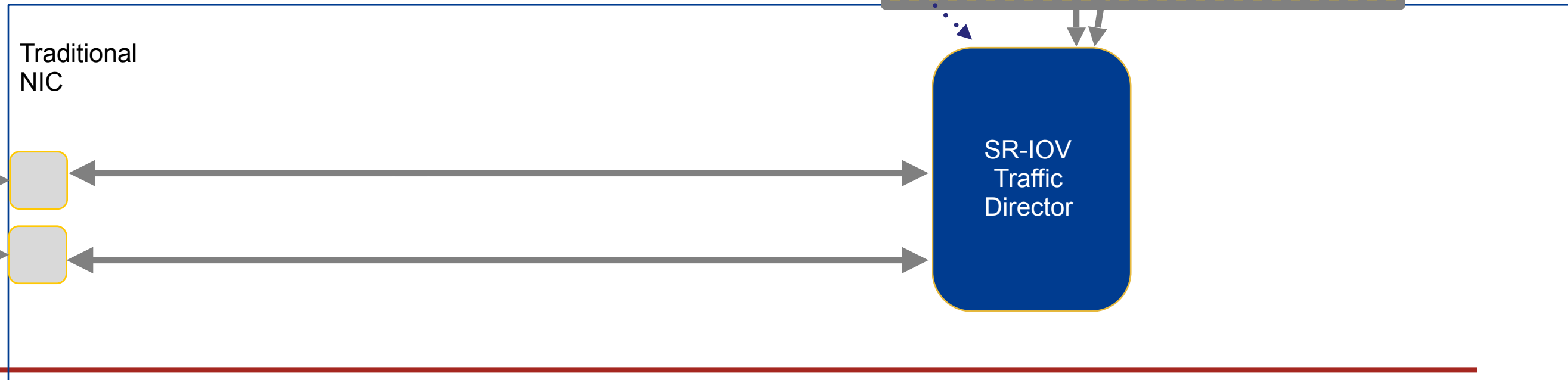


(Nova, Neutron)



1 Configuration from cloud management system

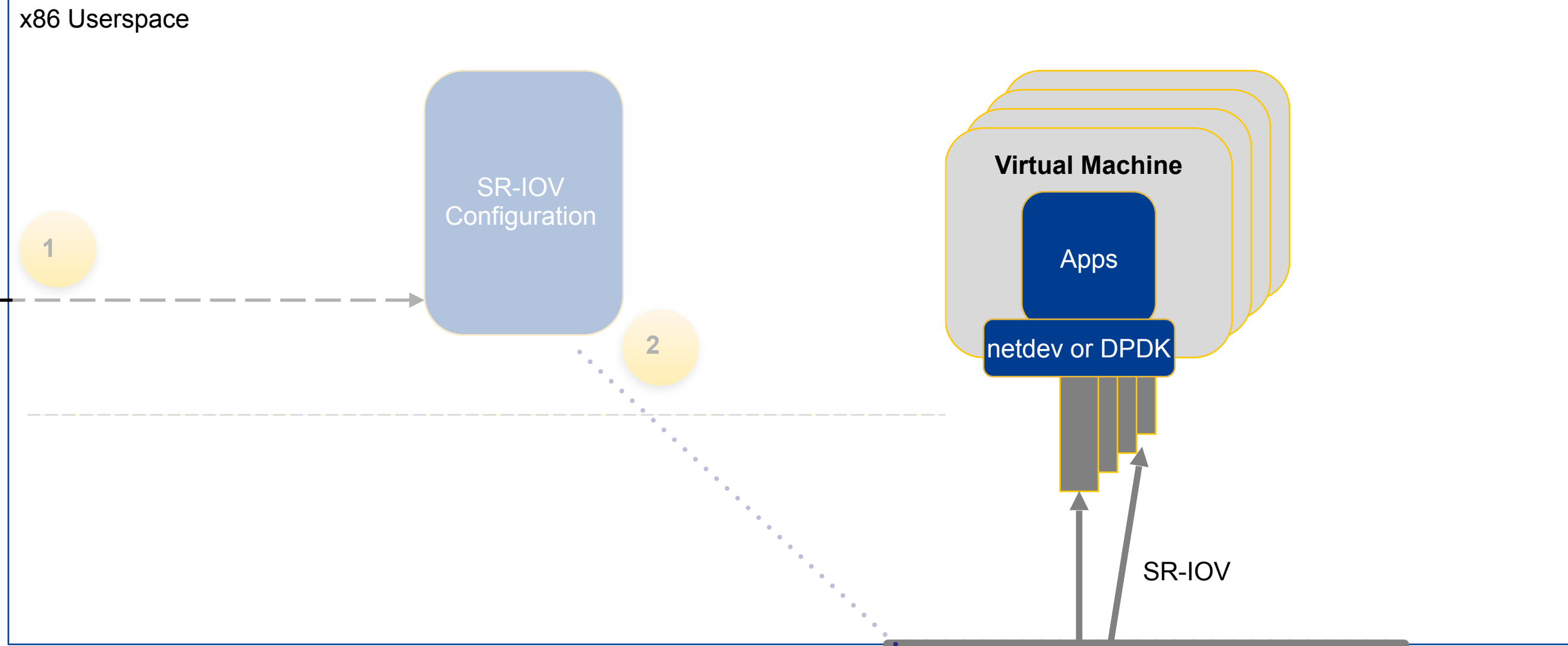
2 Hardware configuration



# Traditional Model: SR-IOV



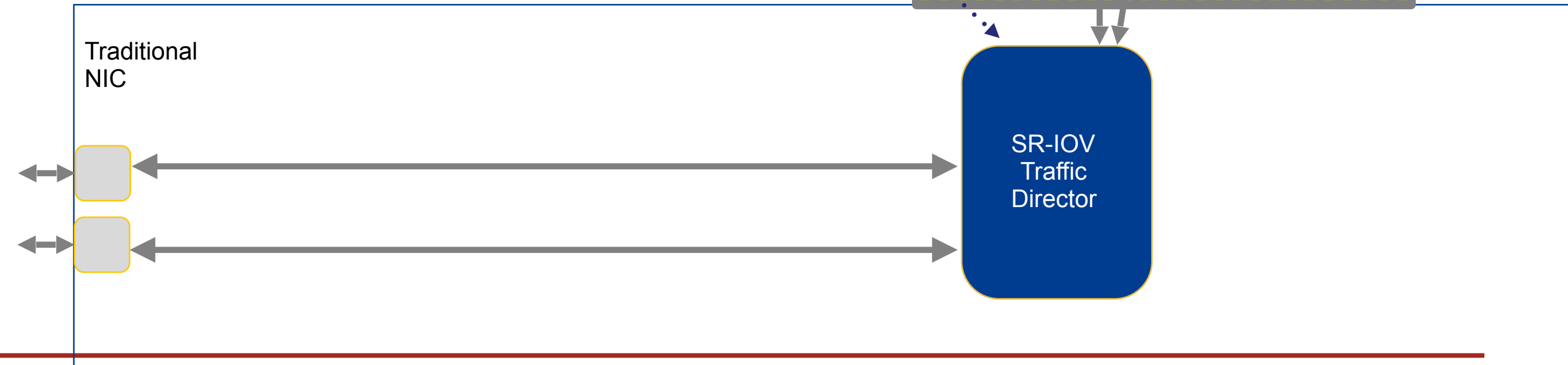
(Nova, Neutron)



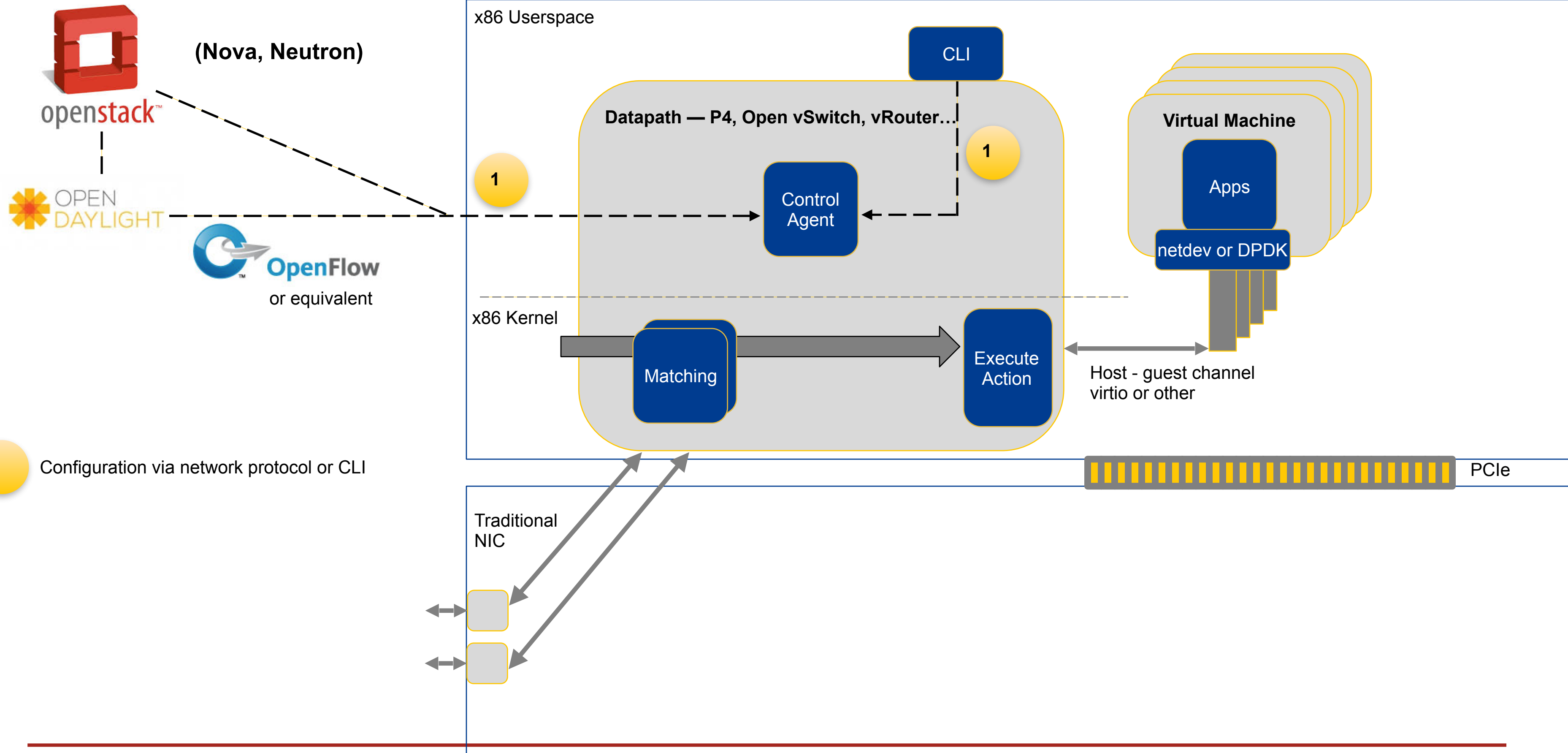
1 Configuration from cloud management system

2 Hardware configuration

- Low expressiveness**  
(MAC/VLAN based traffic directing)
- High performance**
- Poor manageability**  
(no VM migration)

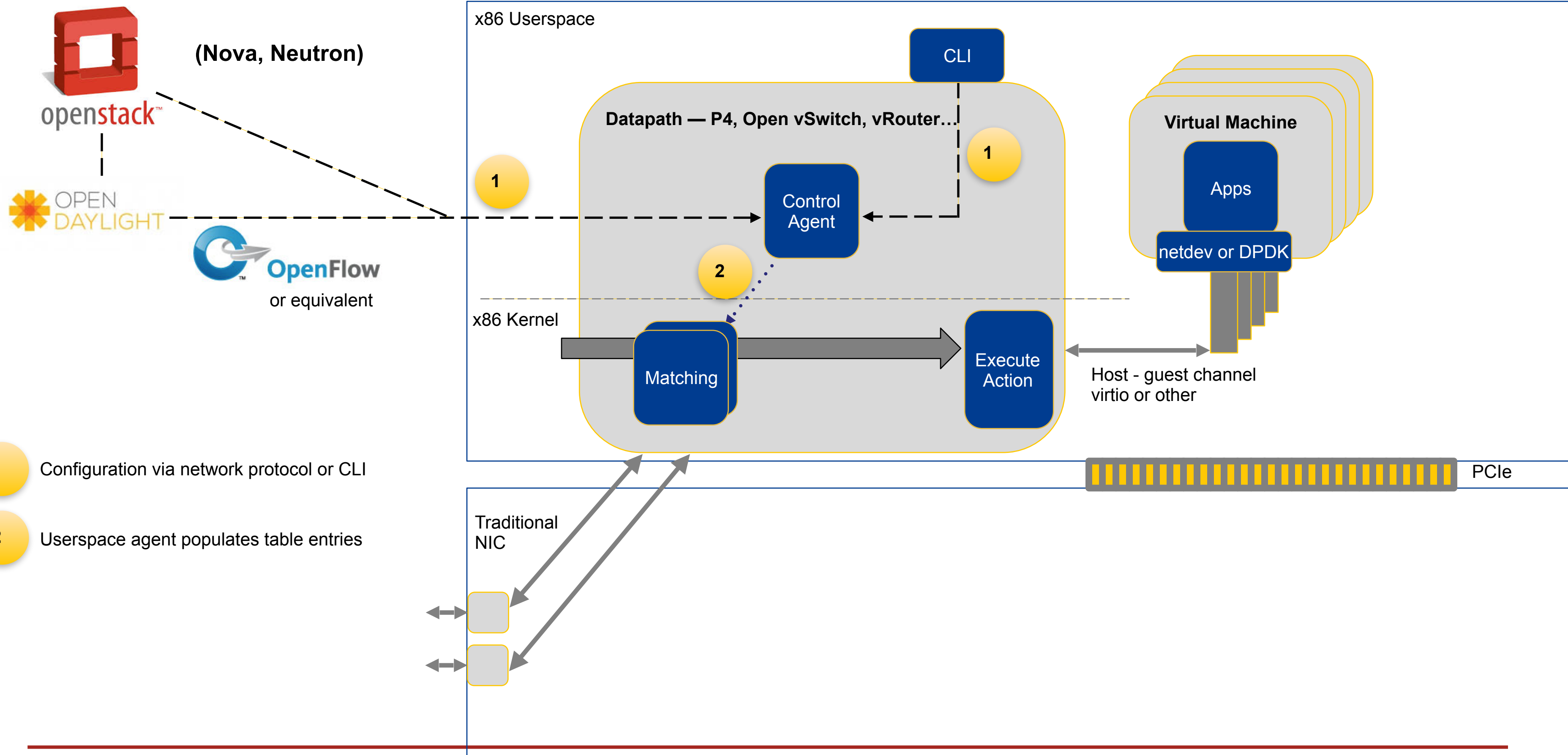


# Traditional Model: Host Runs Datapath

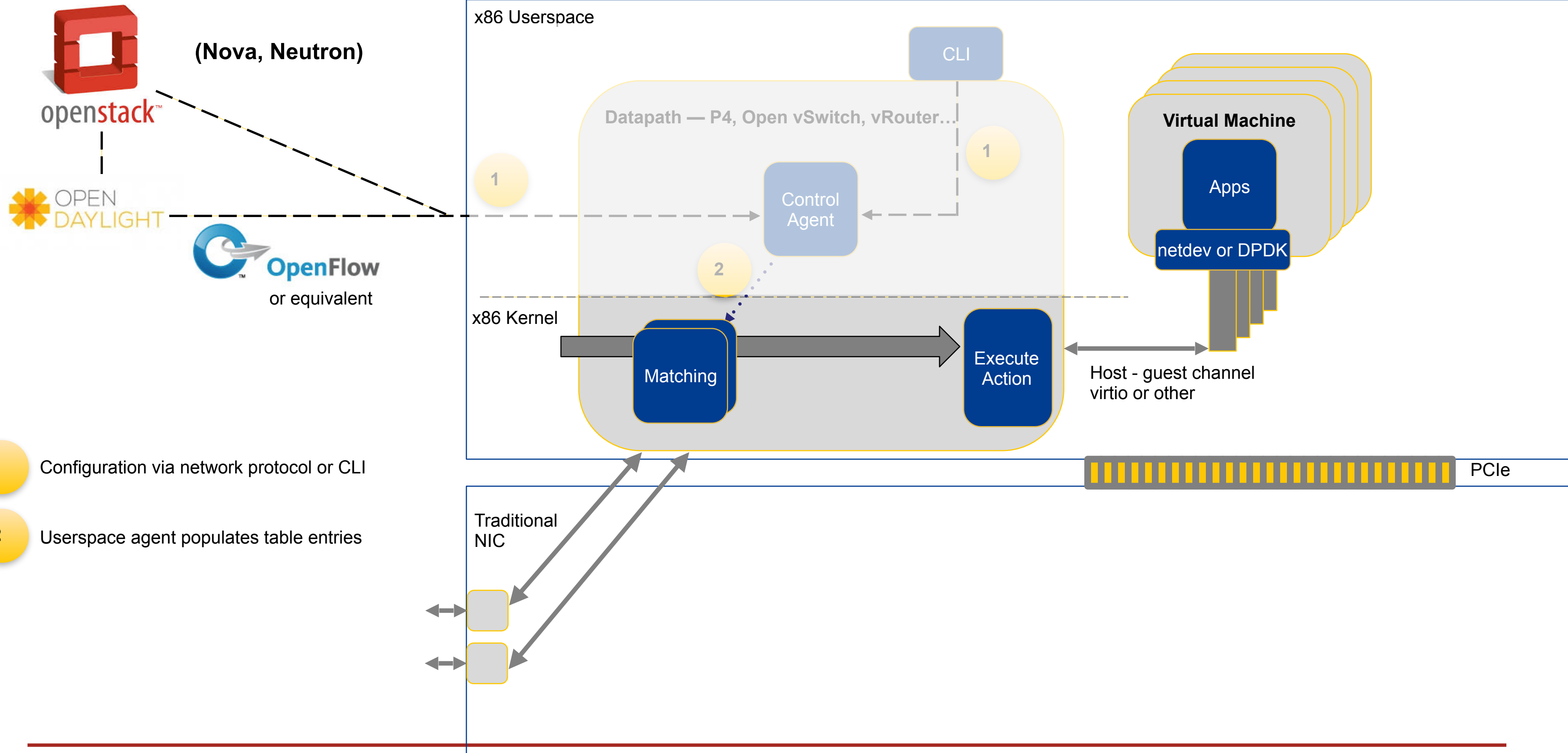




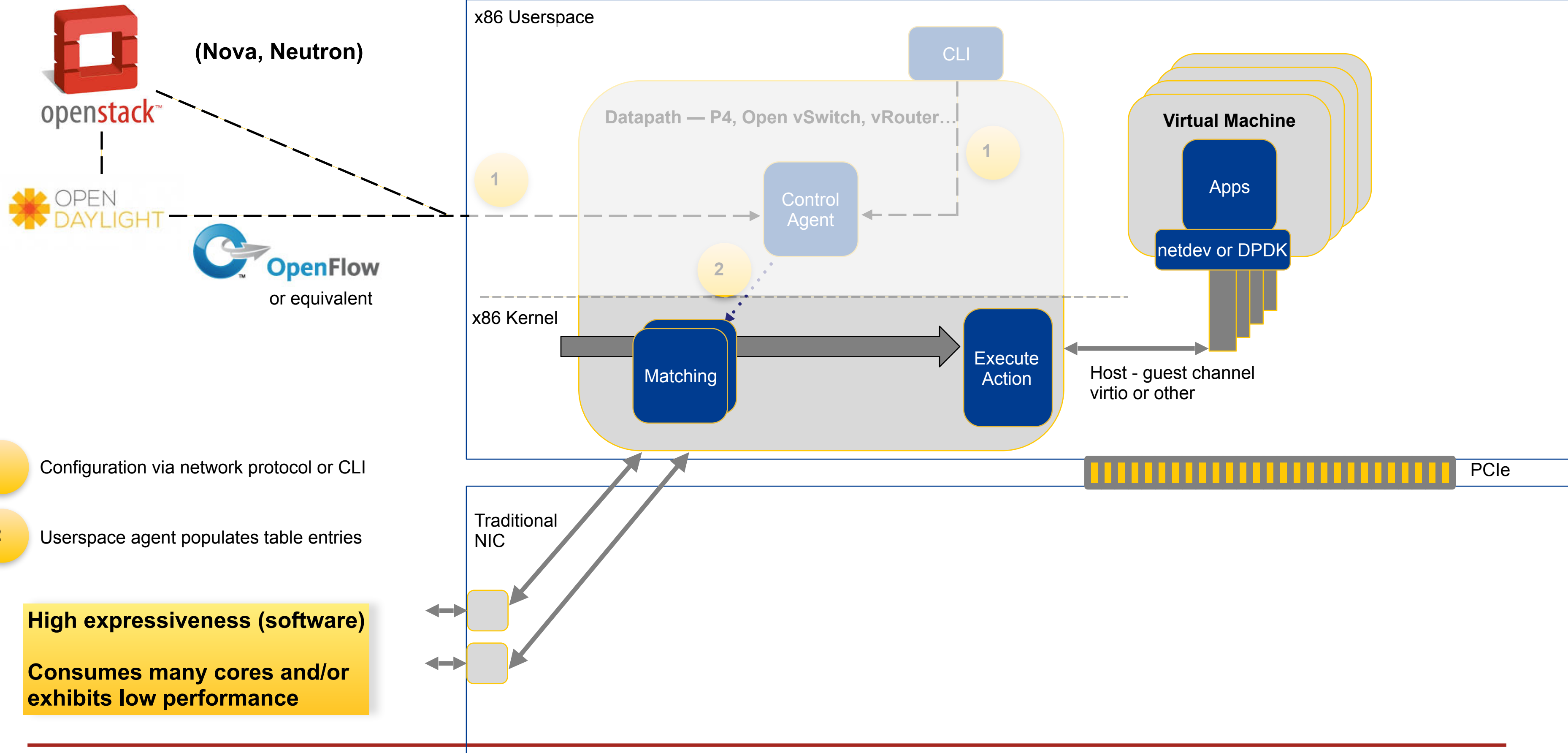
# Traditional Model: Host Runs Datapath



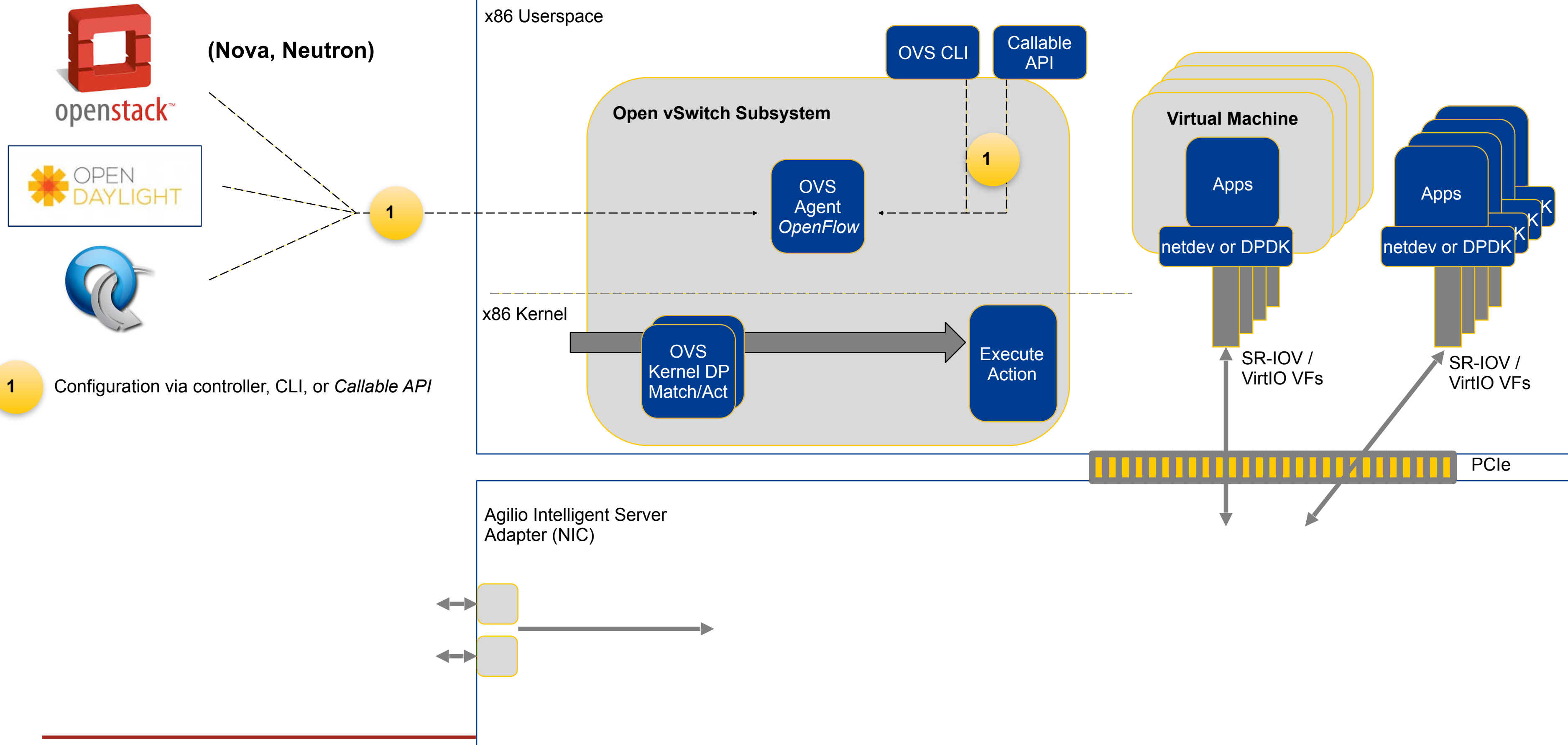
# Traditional Model: Host Runs Datapath



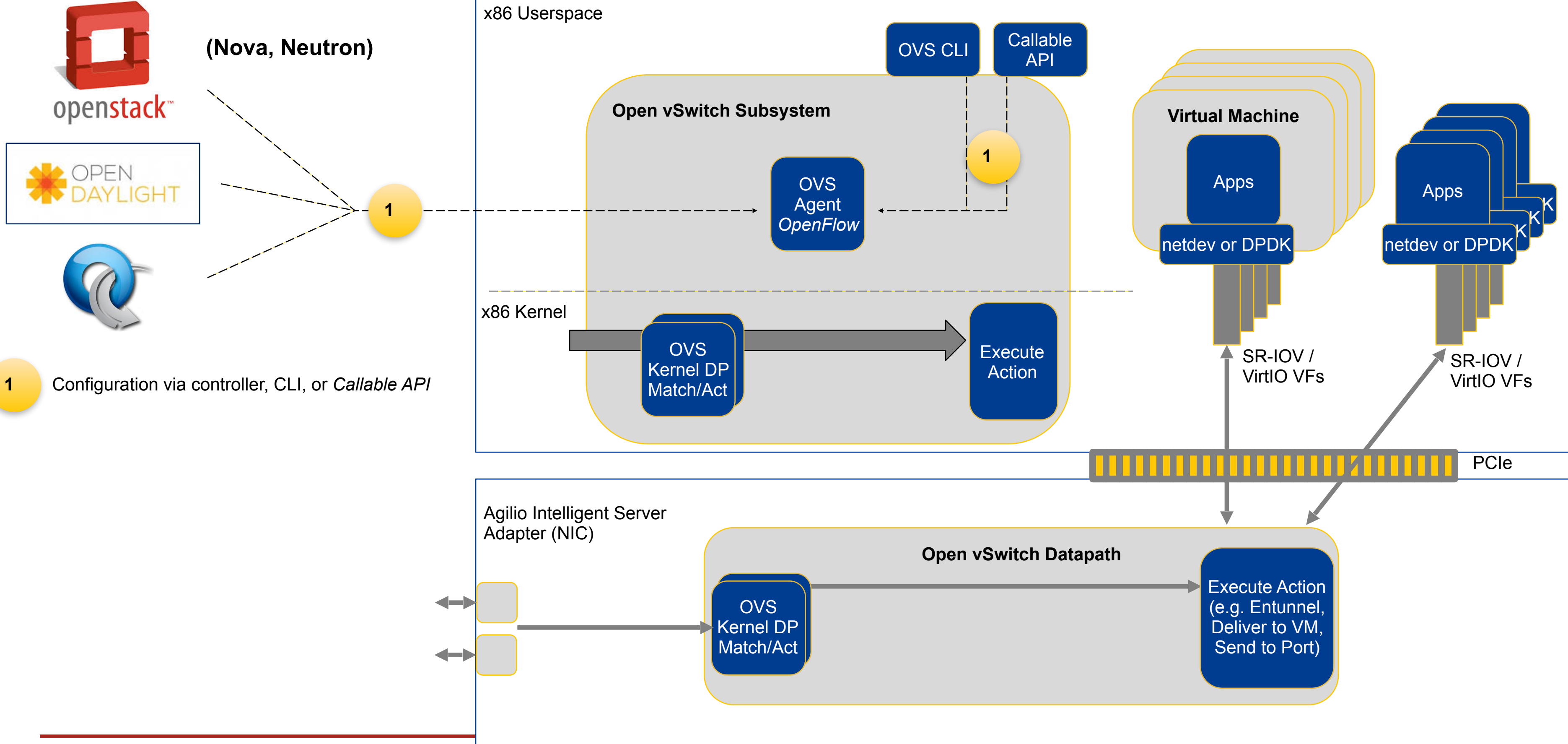
# Traditional Model: Host Runs Datapath



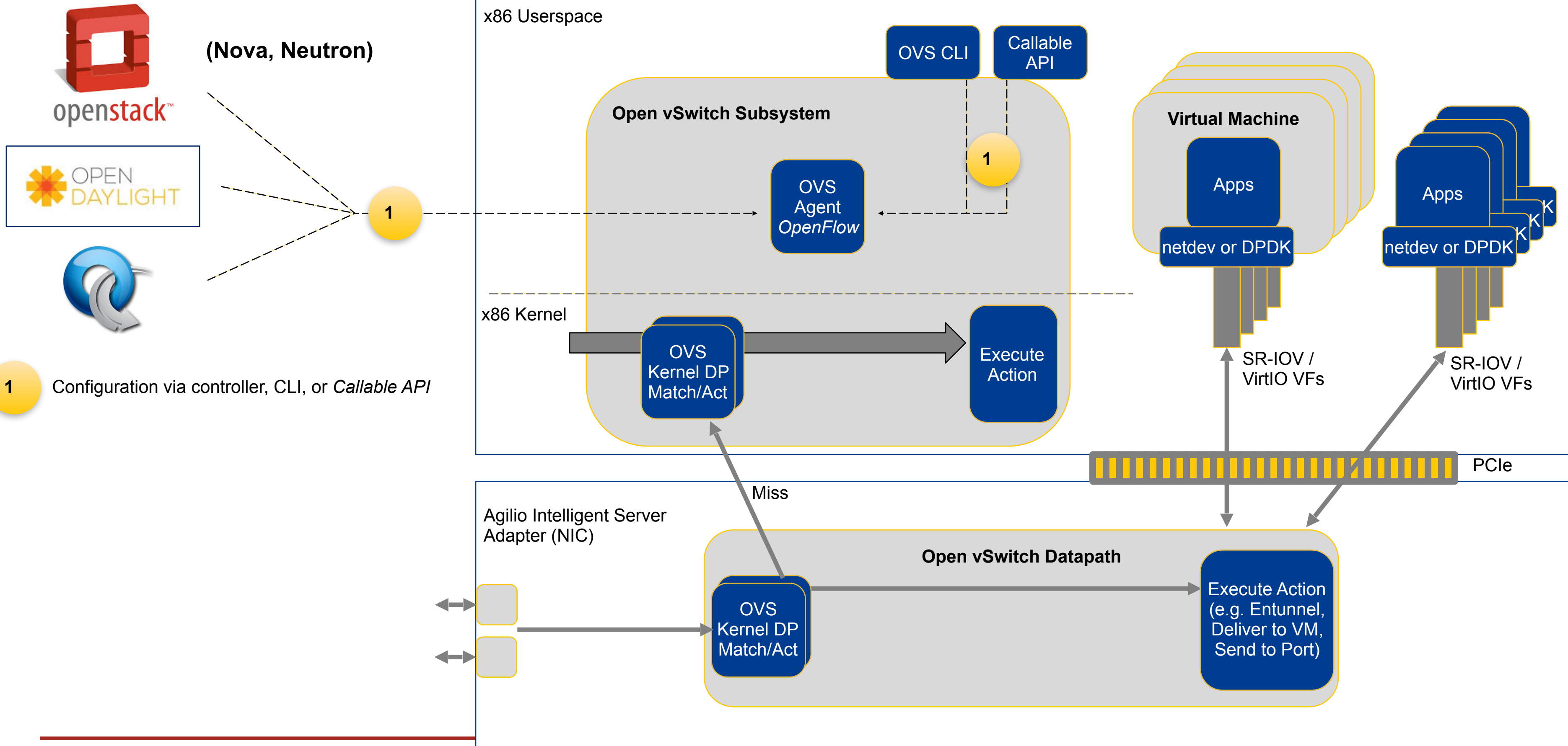
# Offload Model: Agilio™ OVS Acceleration



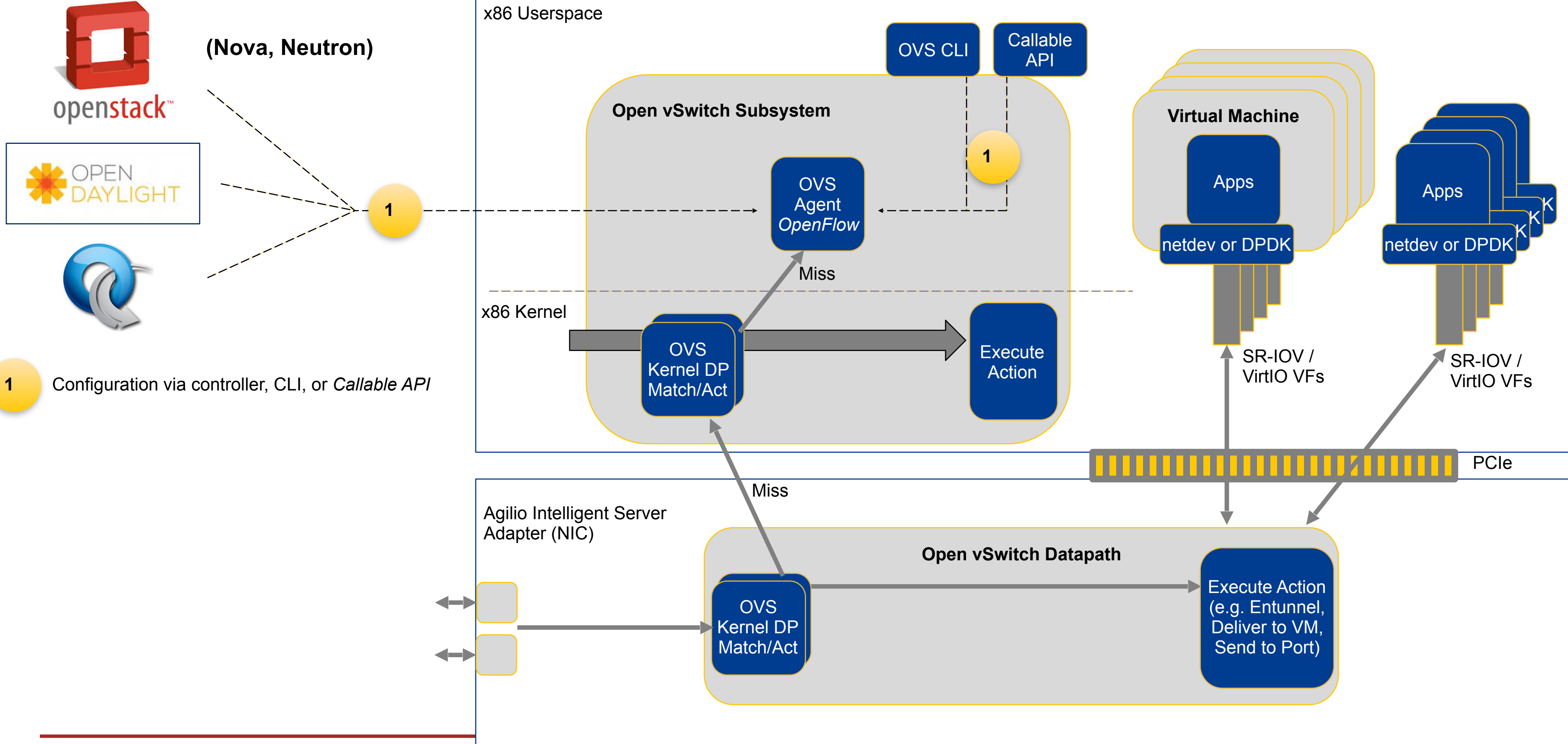
# Offload Model: Agilio™ OVS Acceleration



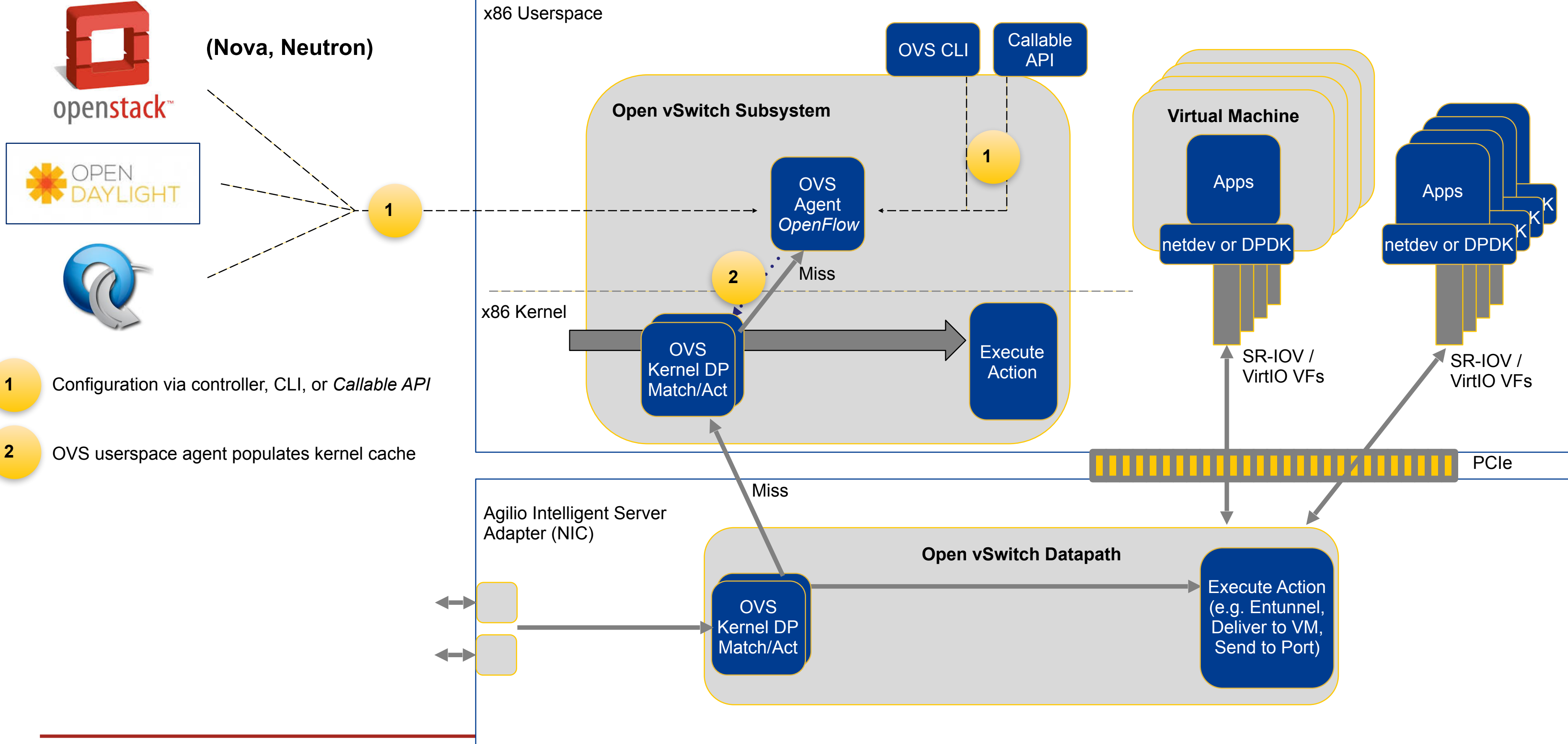
# Offload Model: Agilio™ OVS Acceleration



# Offload Model: Agilio™ OVS Acceleration

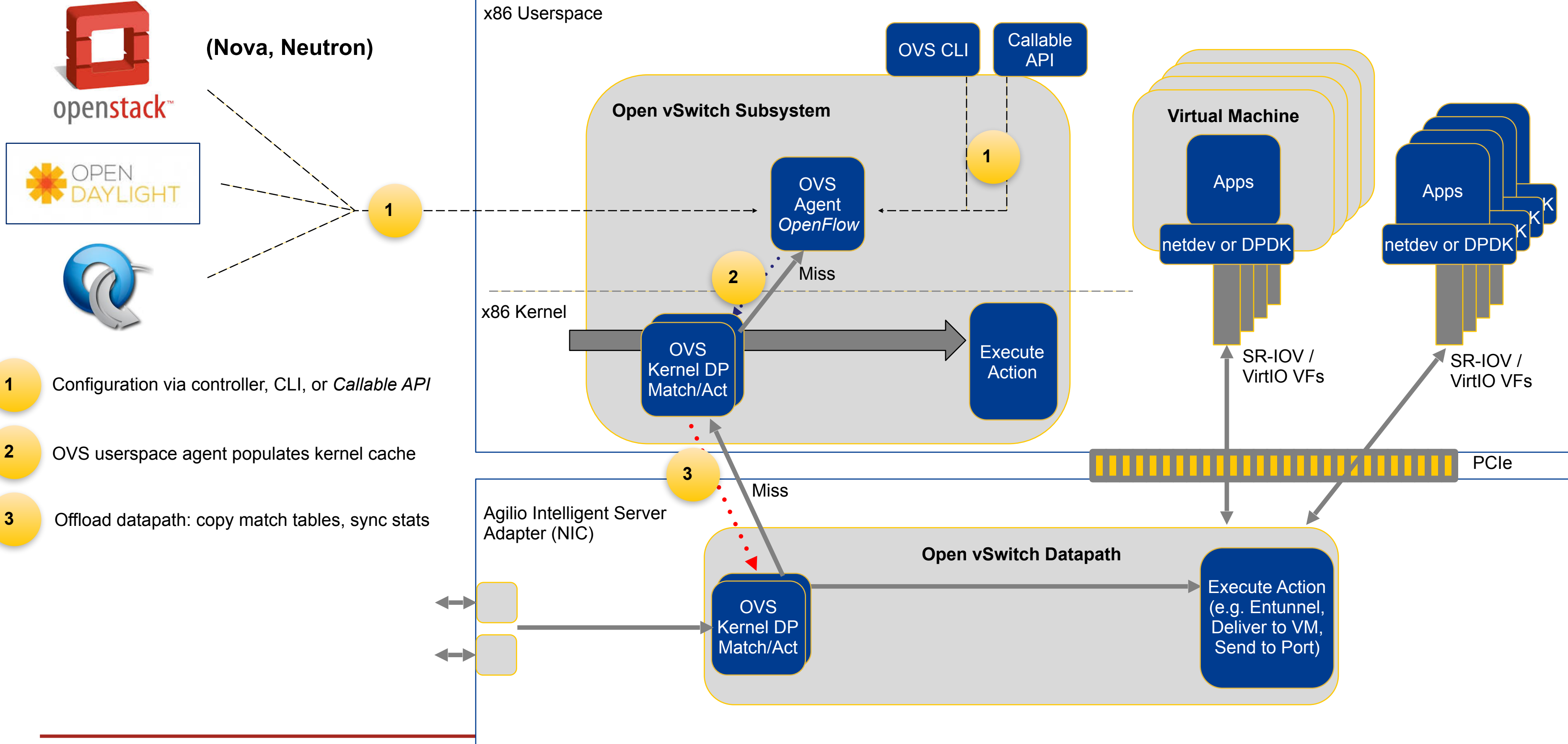


# Offload Model: Agilio™ OVS Acceleration

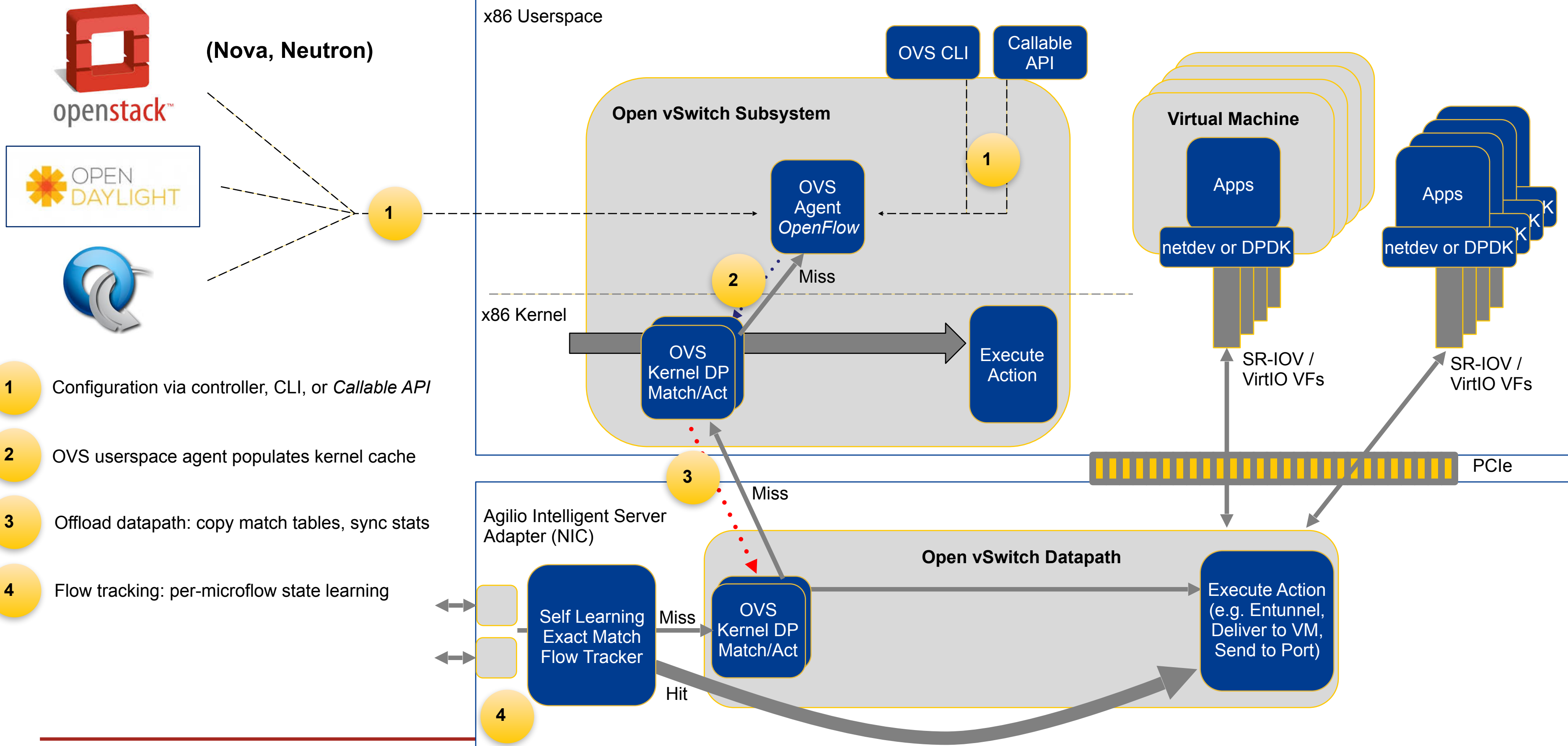




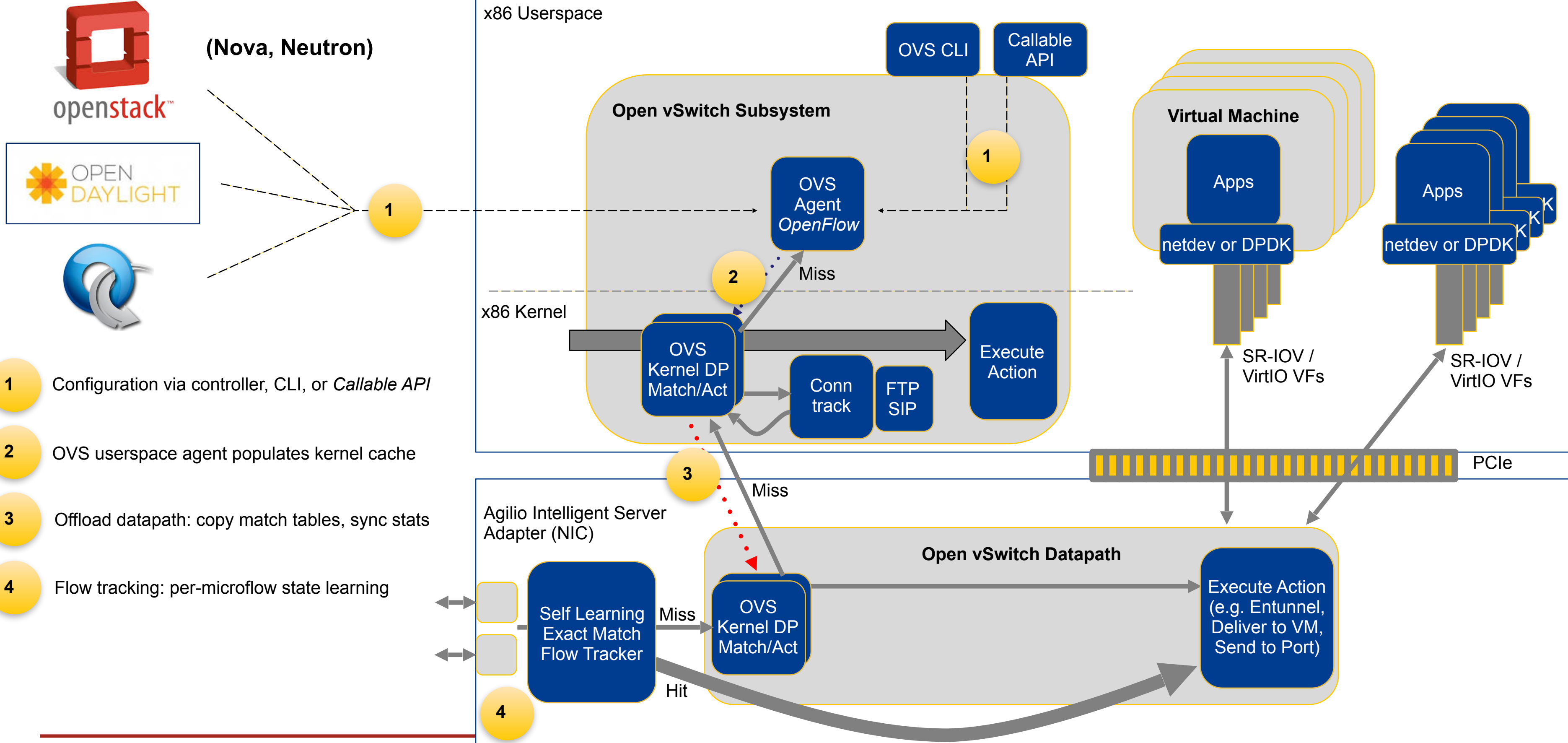
# Offload Model: Agilio™ OVS Acceleration



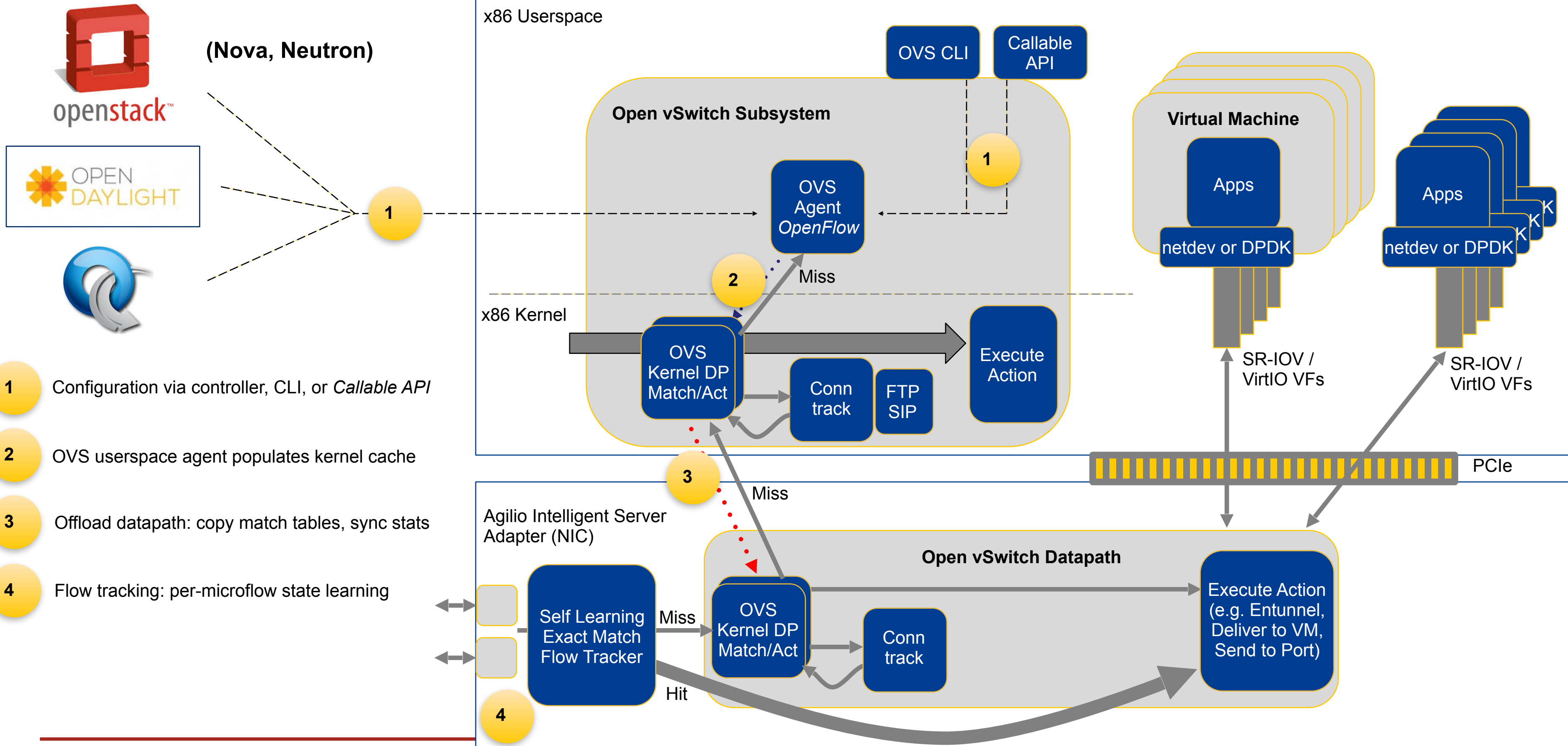
# Offload Model: Agilio™ OVS Acceleration



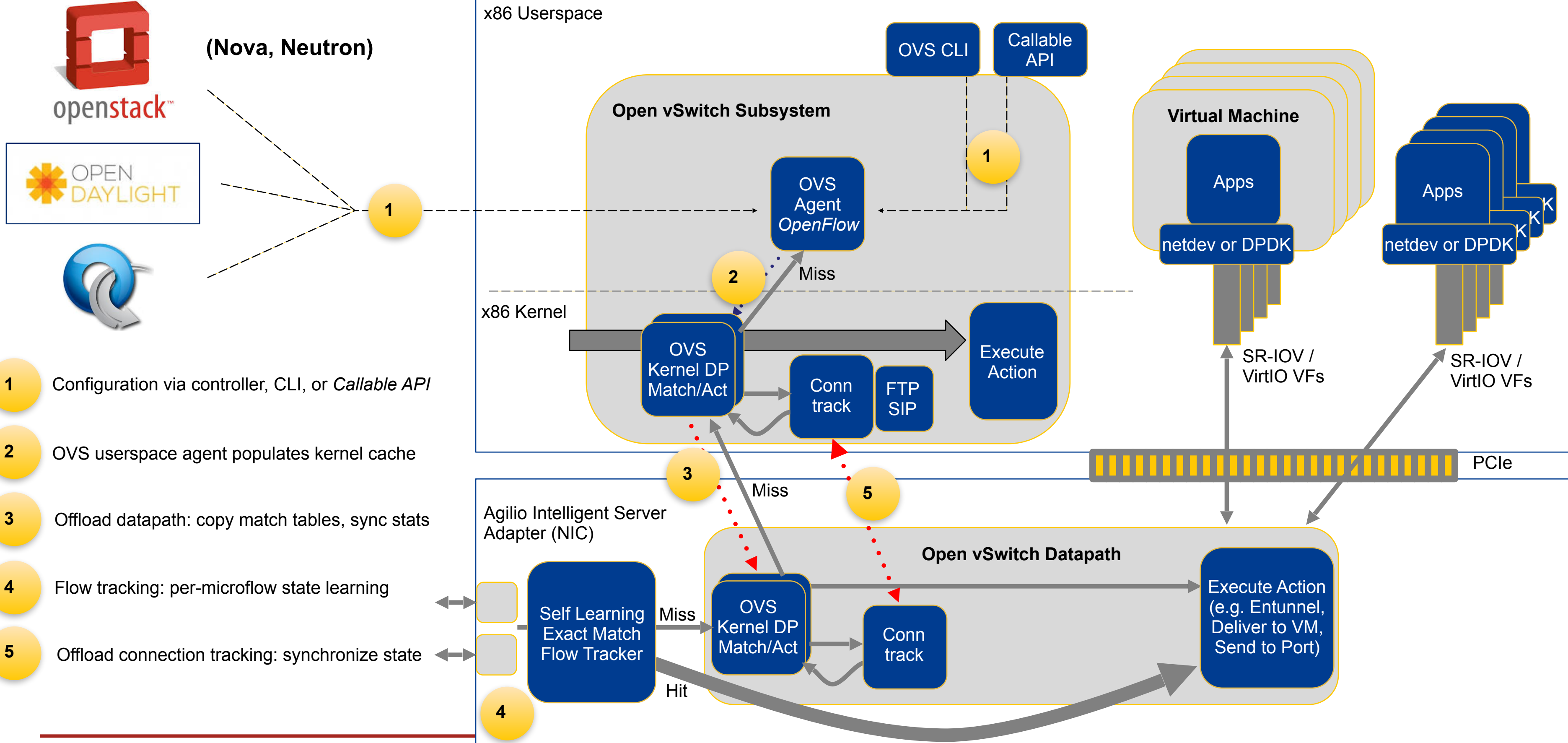
# Offload Model: Agilio™ OVS Acceleration



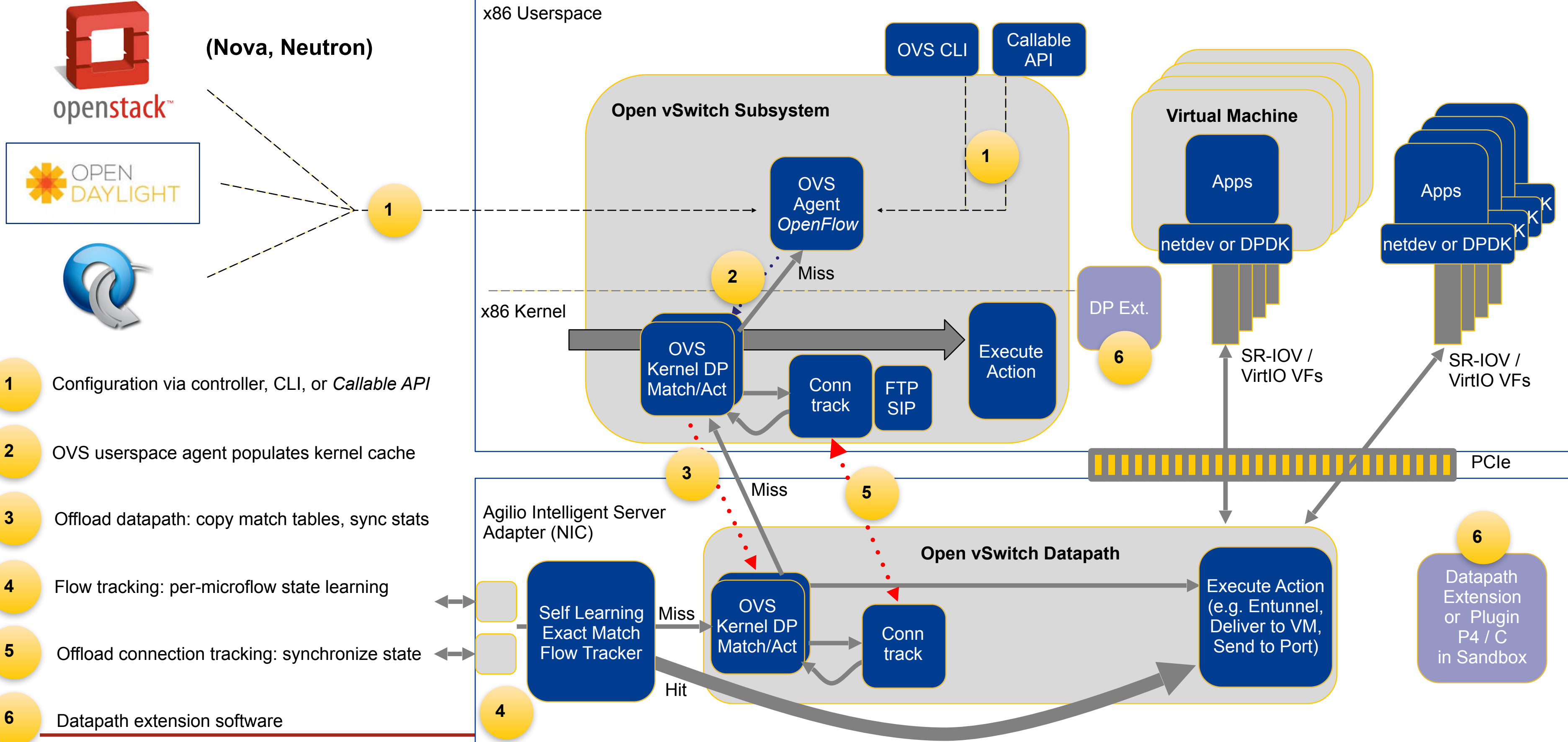
# Offload Model: Agilio™ OVS Acceleration



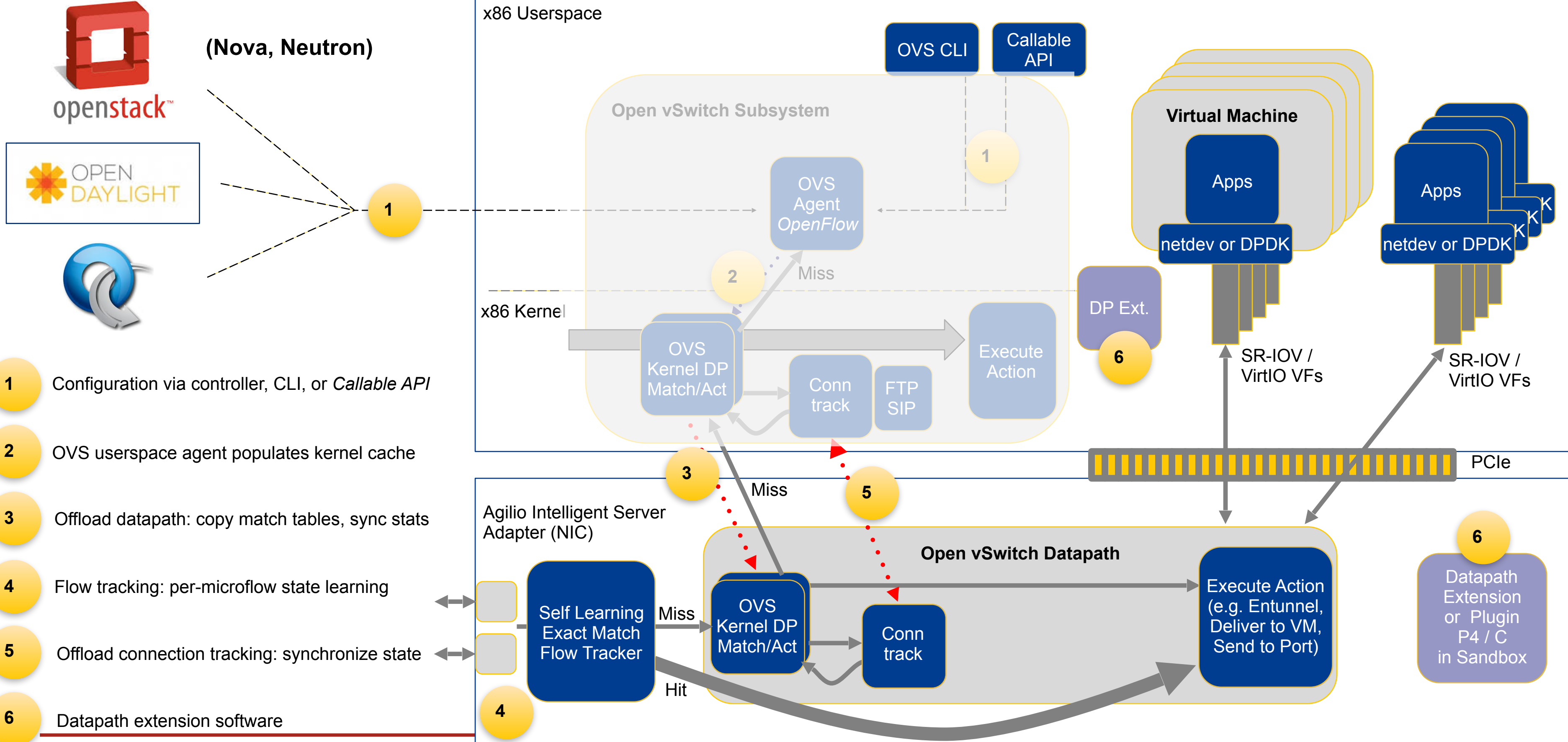
# Offload Model: Agilio™ OVS Acceleration



# Offload Model: Agilio™ OVS Acceleration



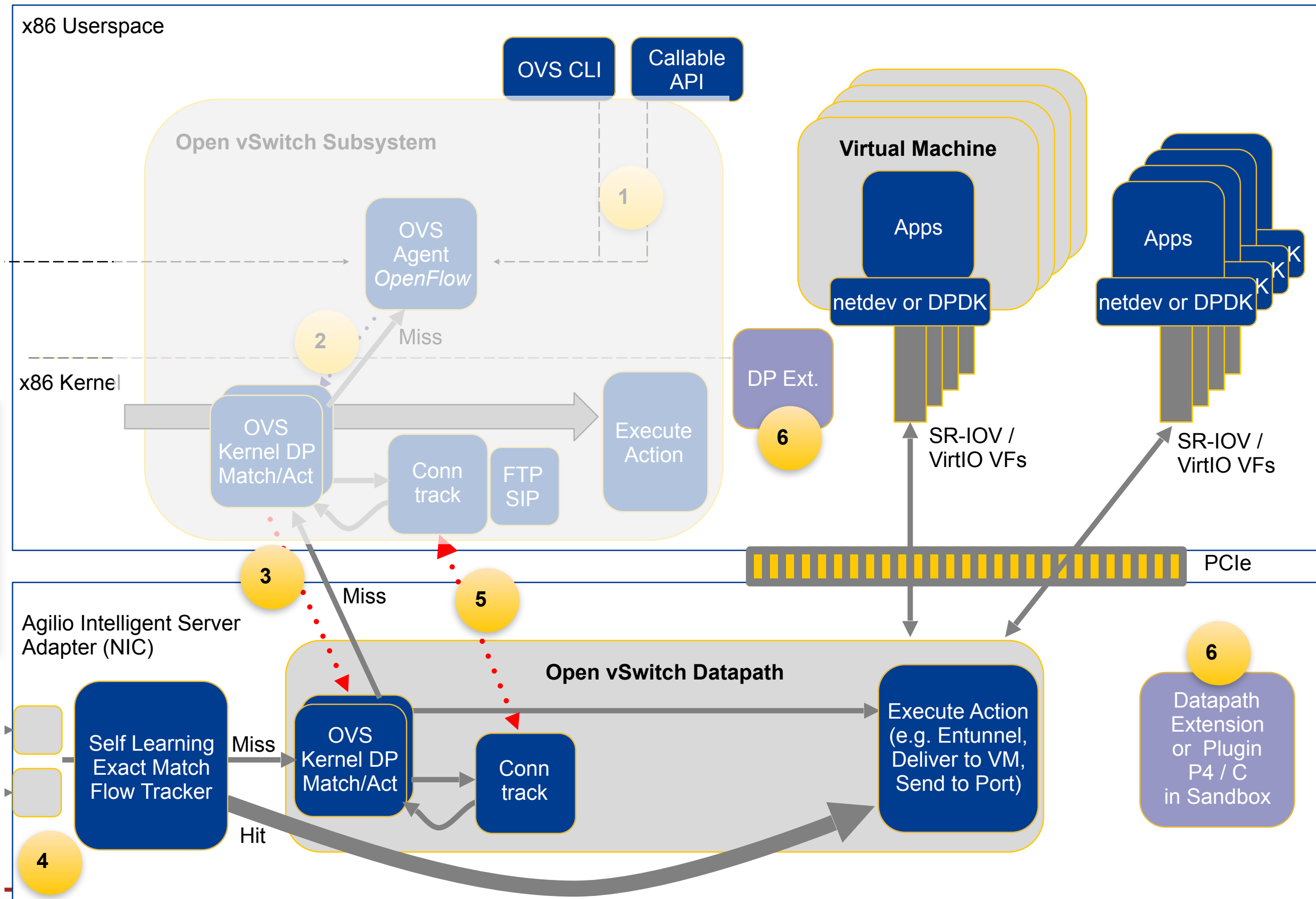
# Offload Model: Agilio™ OVS Acceleration



# Offload Model: Agilio™ OVS Acceleration

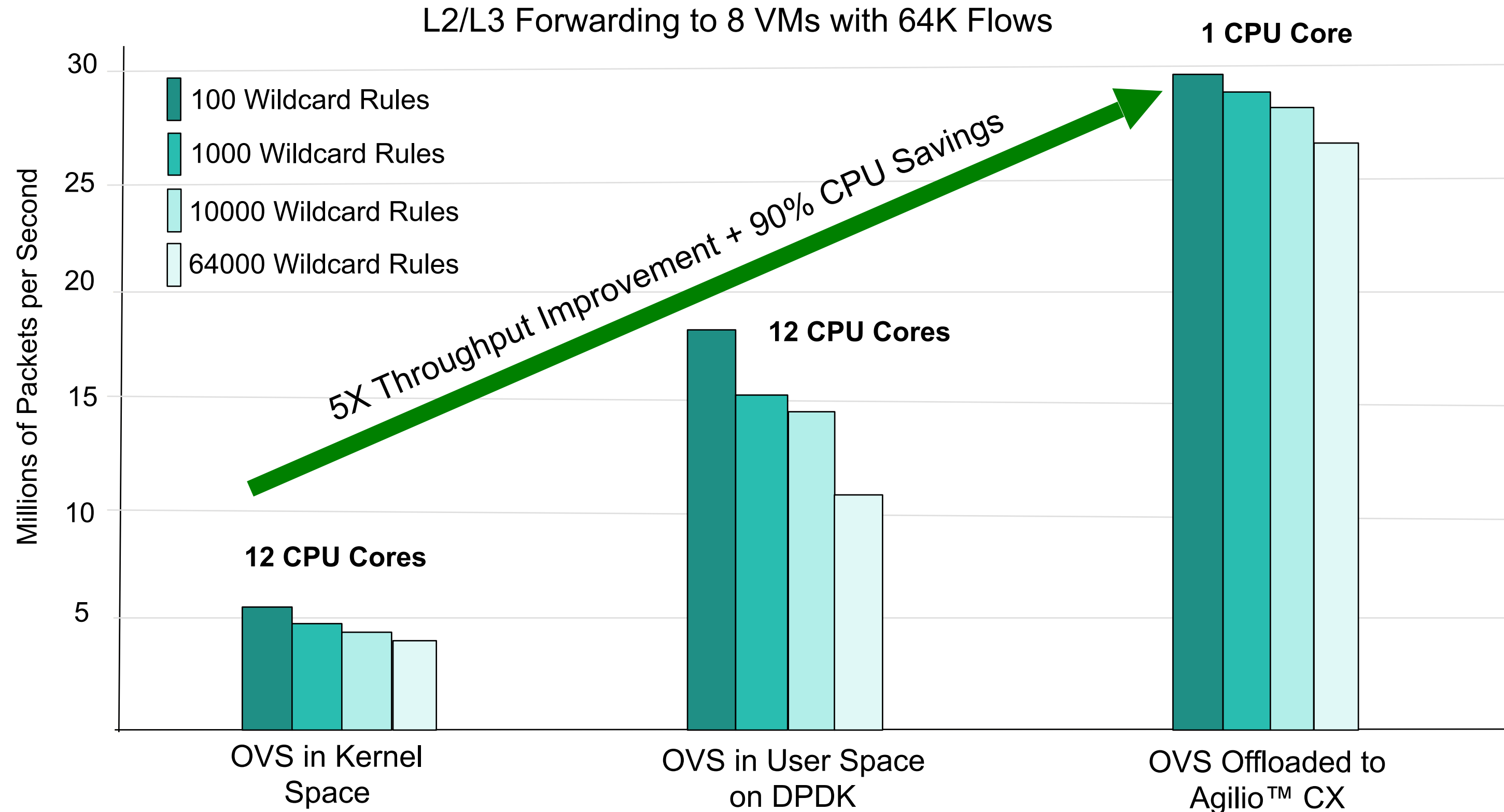
**Best of all worlds**

- Performance of SR-IOV
- Flexibility of virtio (VM migration)
- Performance and CPU core saving of switching on SmartNIC

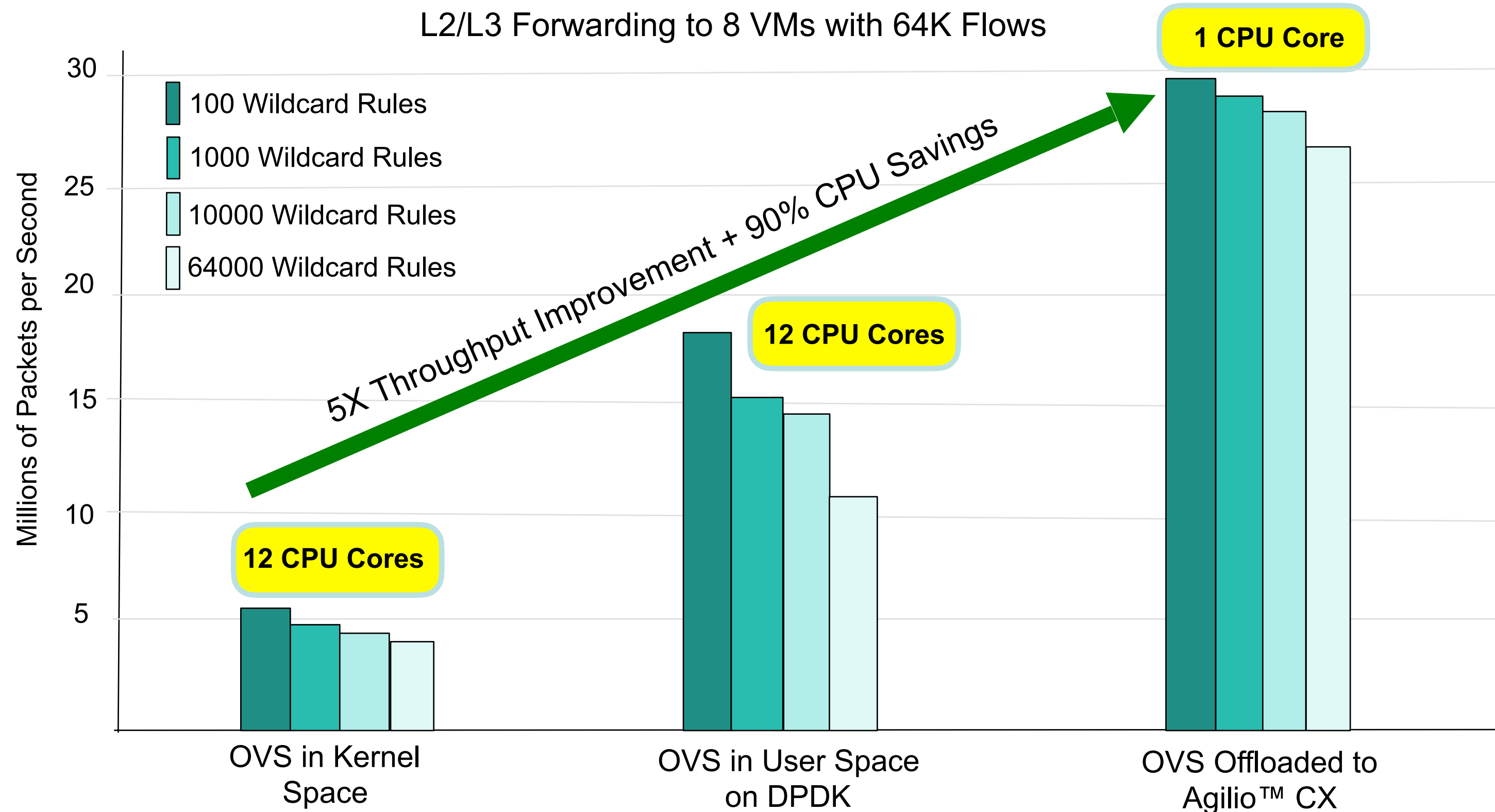




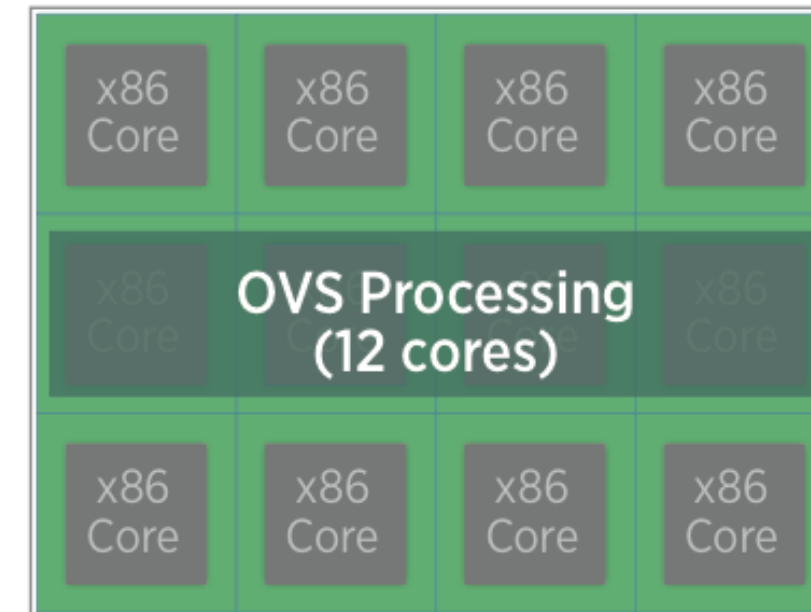
# Example: Throughput vs. Number of Rules



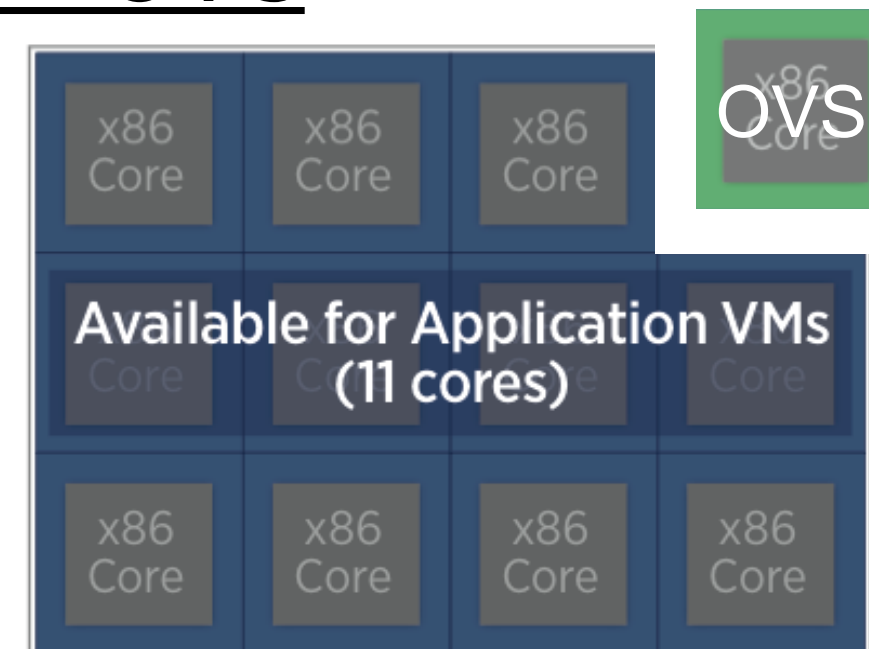
# Example: Throughput vs. Number of Rules



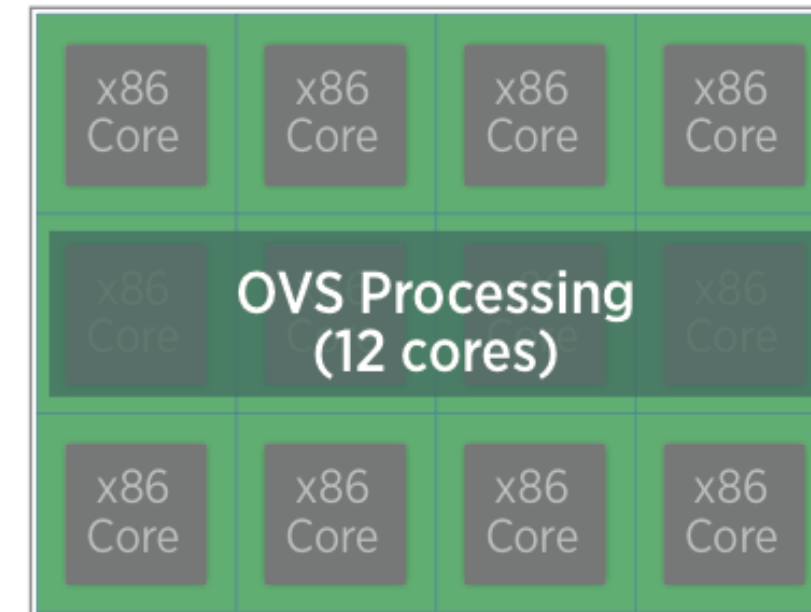
## Unaccelerated OVS (Kernel / User Mode)



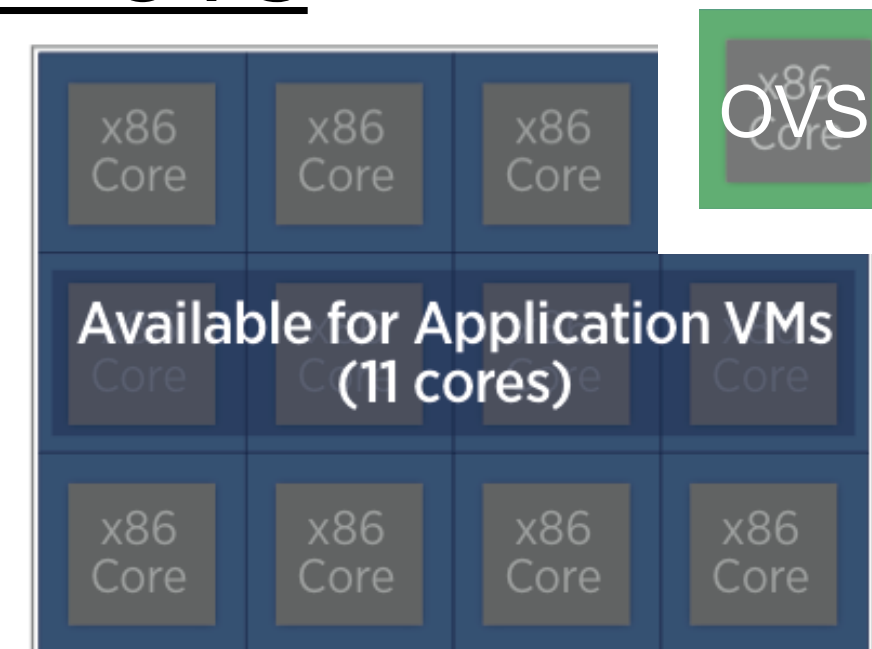
## Agilio™ OVS



## Unaccelerated OVS (Kernel / User Mode)

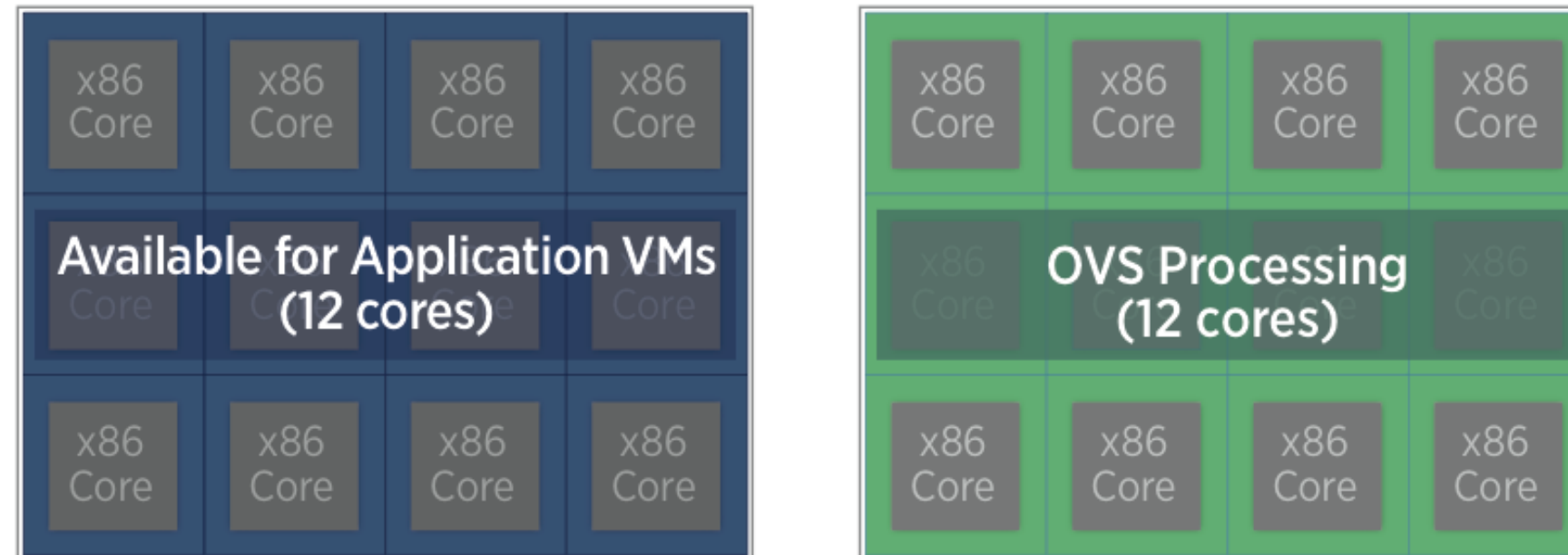


## Agilio™ OVS

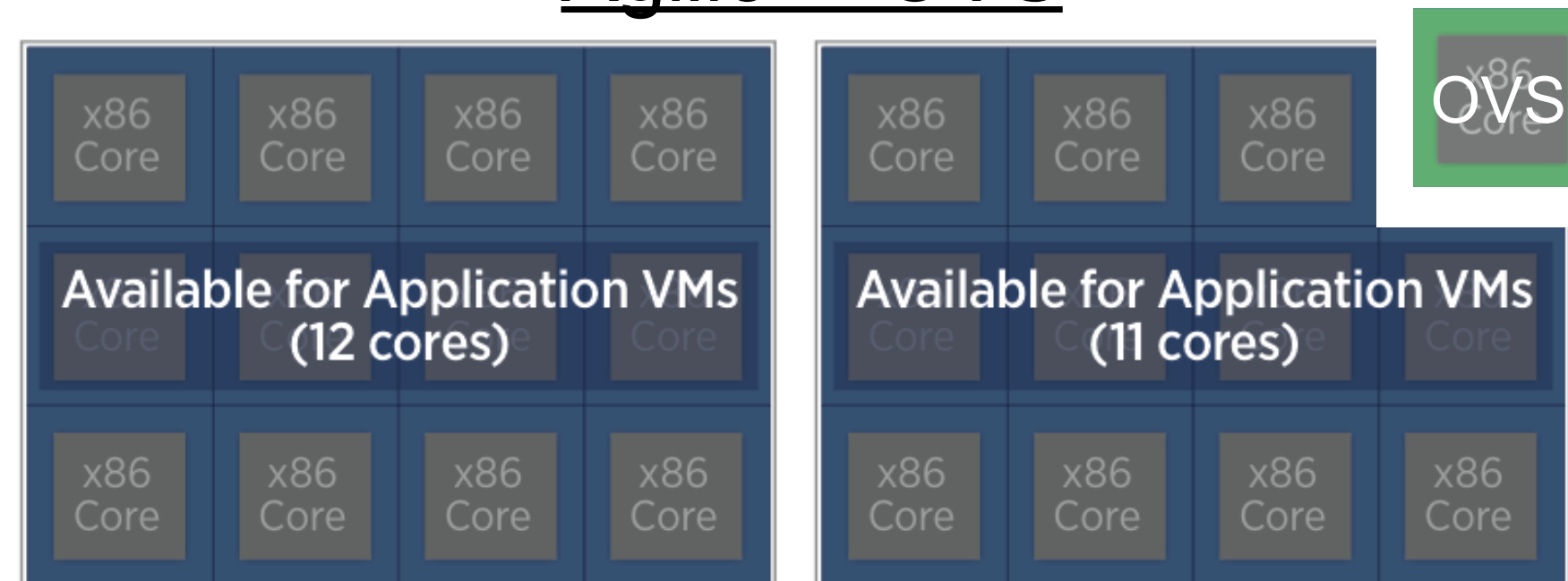


Typical result:  
replace 3-6 racks  
with 1 rack!

## Unaccelerated OVS (Kernel / User Mode)



## Agilio™ OVS

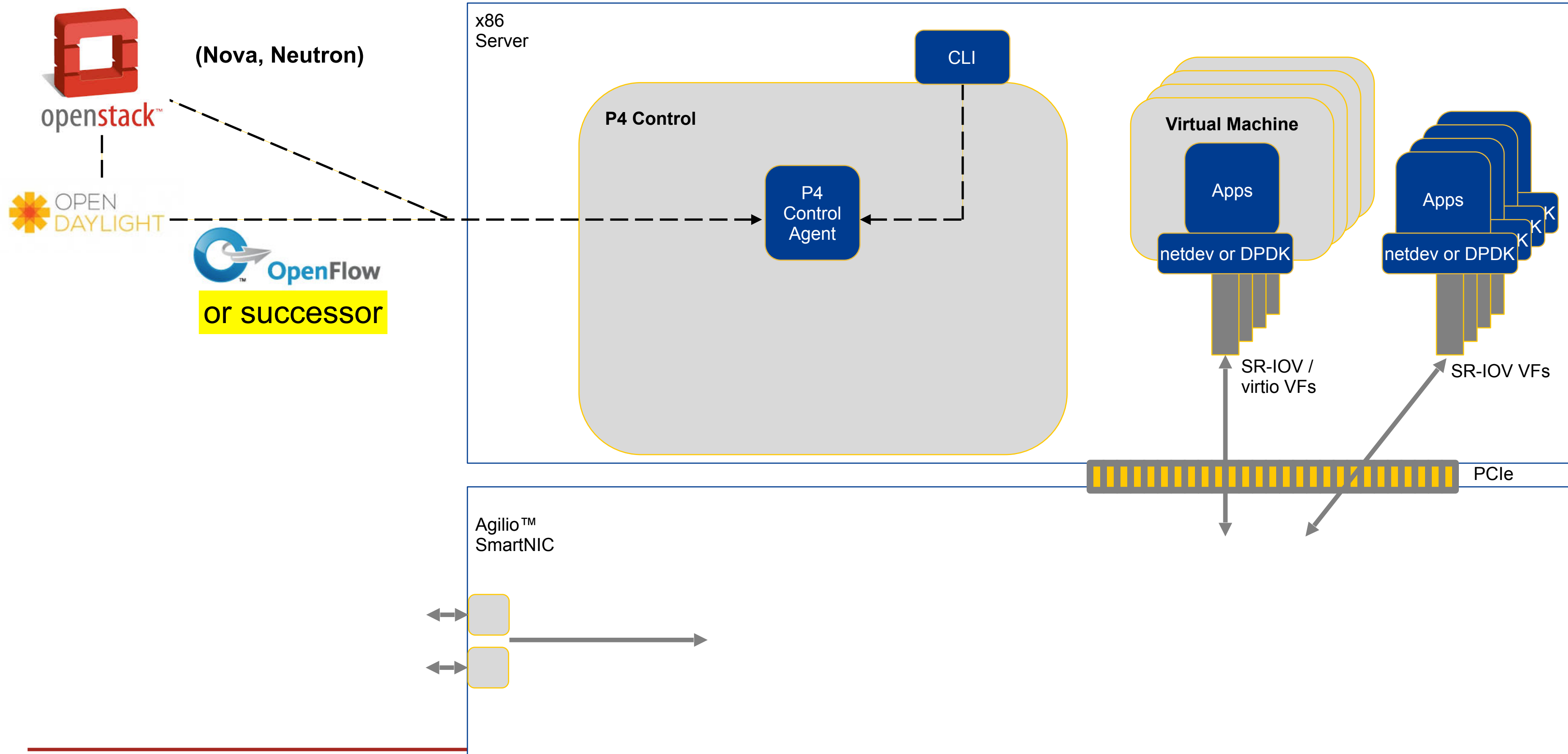


*Benefits for your use case:*

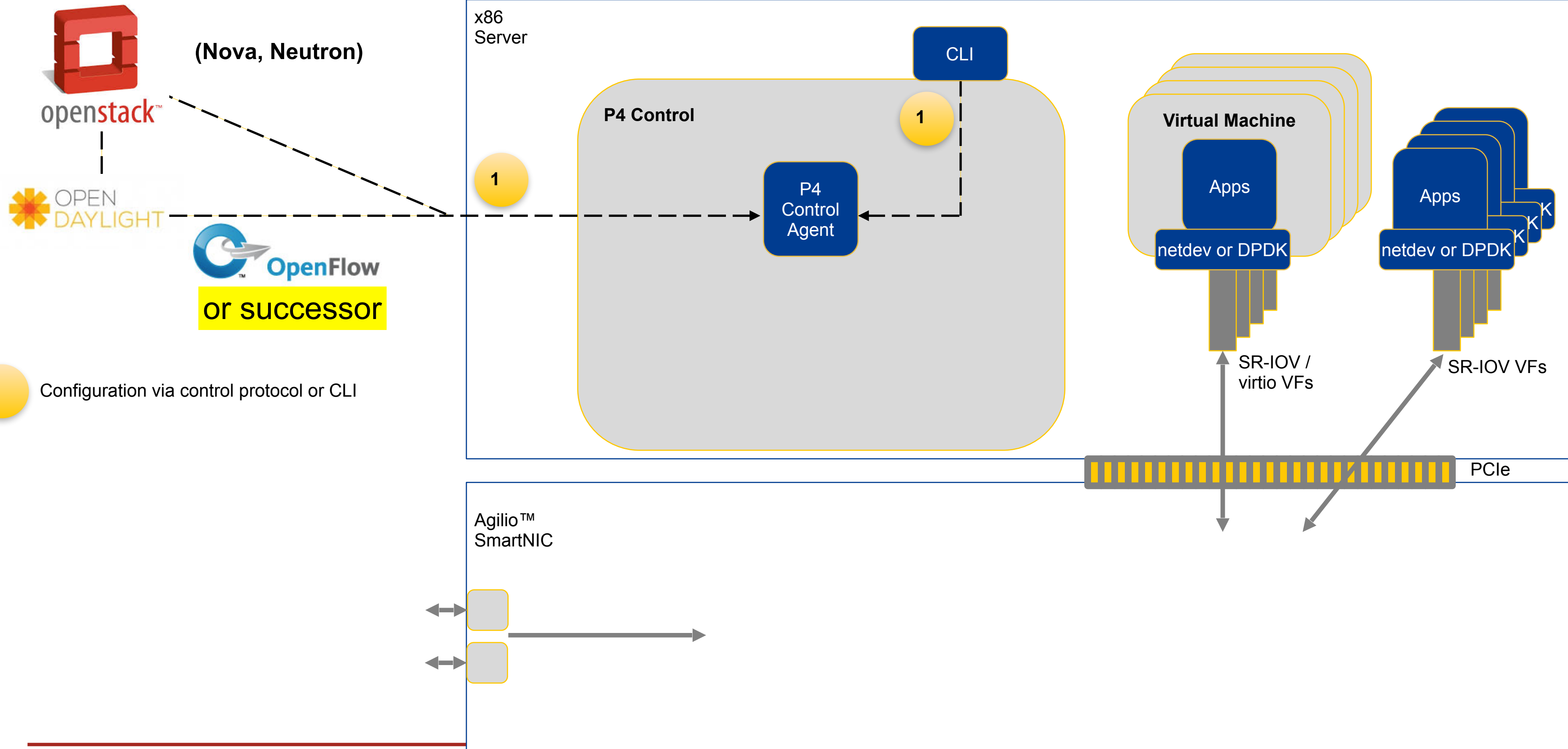
<https://www.netronome.com/products/ovs/roi-calculator/>

Typical result:  
replace 3-6 racks  
with 1 rack!

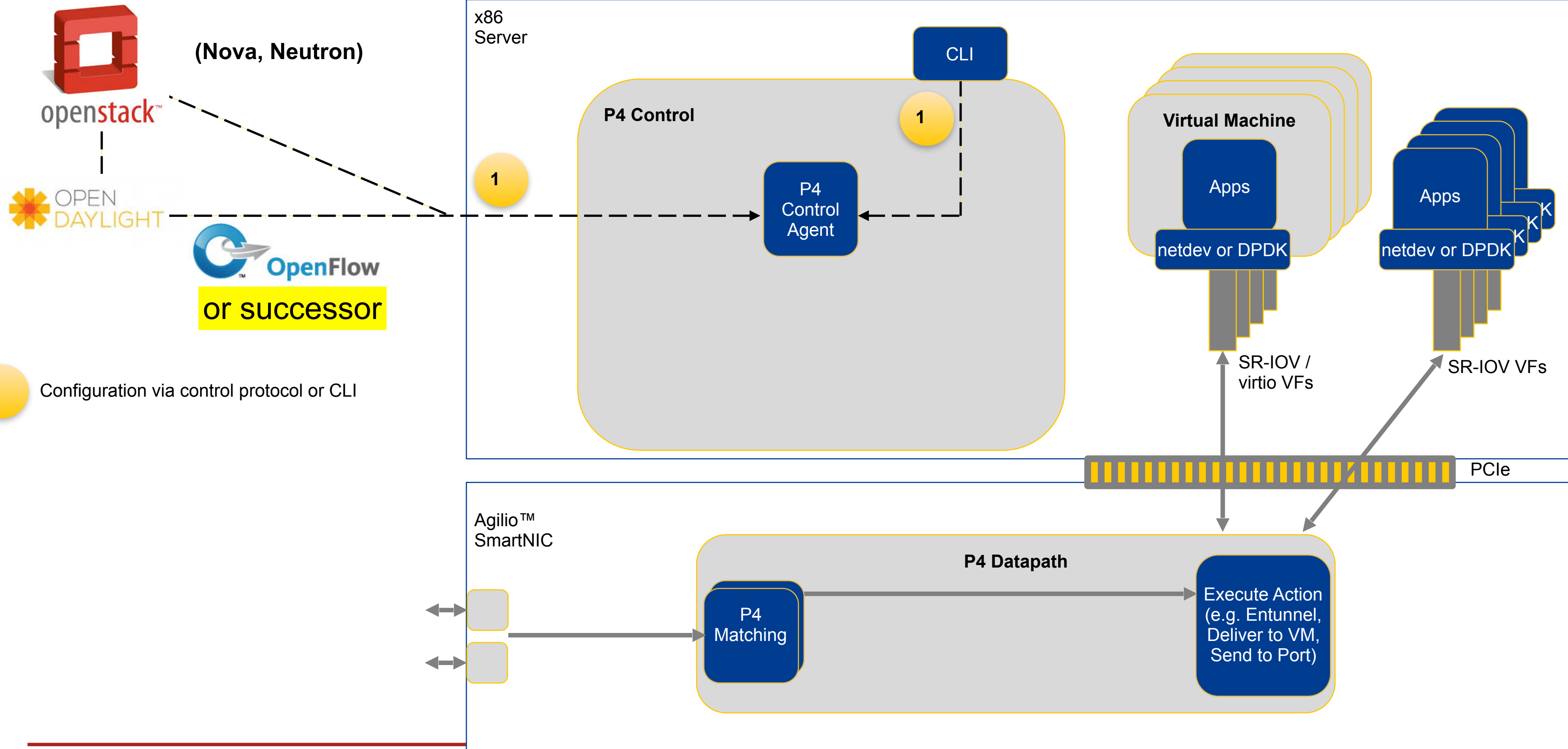
# P4 Datapath on SmartNIC



# P4 Datapath on SmartNIC

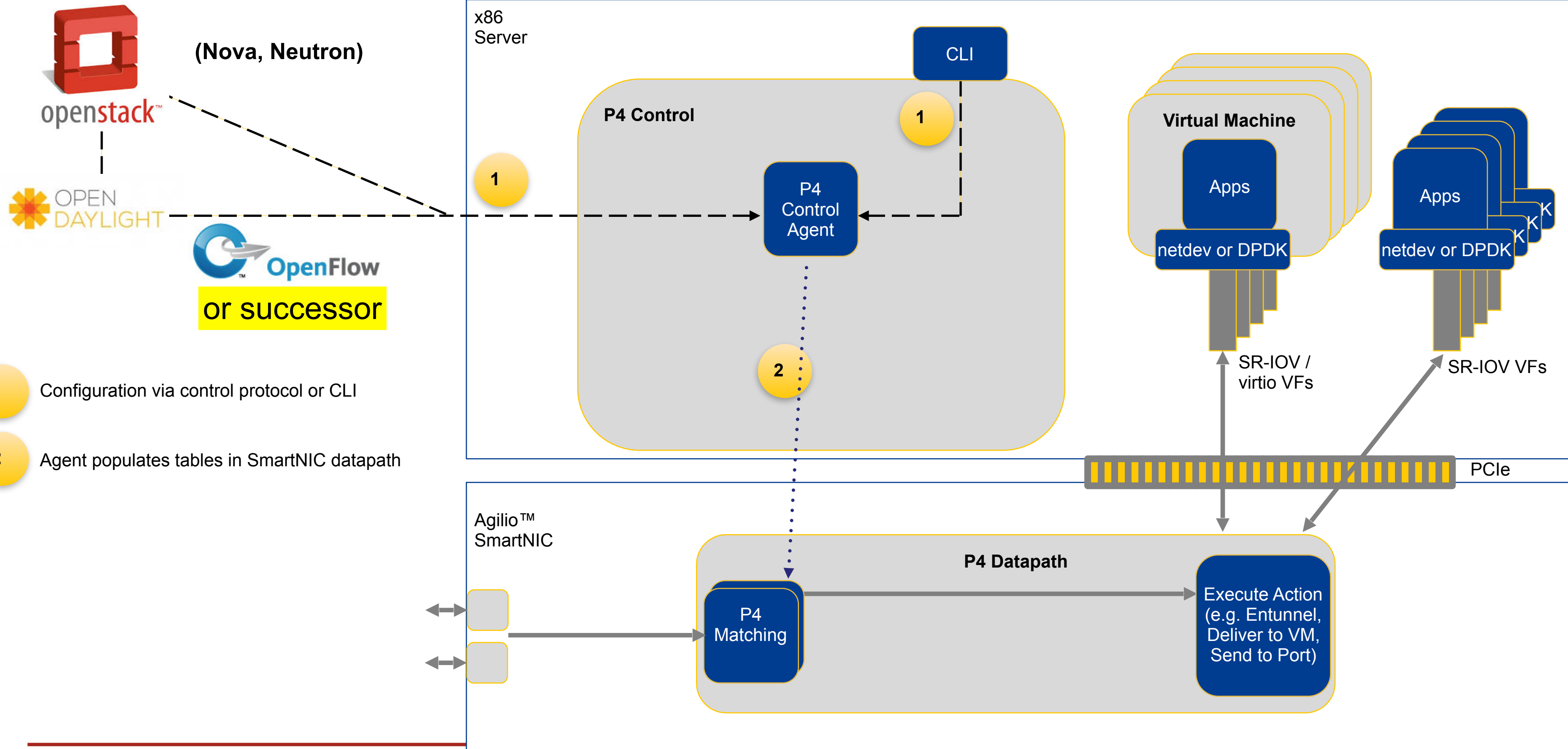


# P4 Datapath on SmartNIC





# P4 Datapath on SmartNIC



Forwarding /  
Virtual Switching  
Technology

SR-IOV

Intelligent  
Datapath

*P4 / vRouter /  
Open vSwitch...*

Forwarding / Virtual Switching Technology	Traditional Approach (Unaccelerated)
SR-IOV	<p><b>Limited expressiveness</b> to direct traffic to VMs (no support for general match/action rules, tunnel termination, stateful firewalling)</p> <p><b>High throughput</b></p> <p><b>No VM migration support</b> (difficult to manage)</p>
Intelligent Datapath  <i>P4 / vRouter / Open vSwitch...</i>	<p><b>High expressiveness</b> - match/action, tunnels, stateless/stateful firewalling etc.</p> <p><b>Limited throughput</b></p> <p><b>High CPU utilization</b> (e.g. 50% of cores)</p>

Forwarding / Virtual Switching Technology	Traditional Approach (Unaccelerated)	Fully Programmable SmartNIC Accelerated Approach
SR-IOV	<p><b>Limited expressiveness</b> to direct traffic to VMs (no support for general match/action rules, tunnel termination, stateful firewalling)</p> <p><b>High throughput</b></p> <p><b>No VM migration support</b> (difficult to manage)</p>	<p><b>High expressiveness</b> - match/action, tunnels, stateless/stateful firewalling etc. <b>and SR-IOV based data delivery to VMs</b></p> <p><b>High throughput</b></p> <p><b>Virtio</b> supporting VM migration (facilitating cloud optimization and upgrading)</p>
Intelligent Datapath  <i>P4 / vRouter / Open vSwitch...</i>	<p><b>High expressiveness</b> - match/action, tunnels, stateless/stateful firewalling etc.</p> <p><b>Limited throughput</b></p> <p><b>High CPU utilization</b> (e.g. 50% of cores)</p>	<p><b>Higher throughput</b> (~5x higher)</p> <p><b>Lower CPU utilization</b> (~10x lower)</p>

## Network Flow Processor 4xxx (used on Agilio-CX SmartNICs)

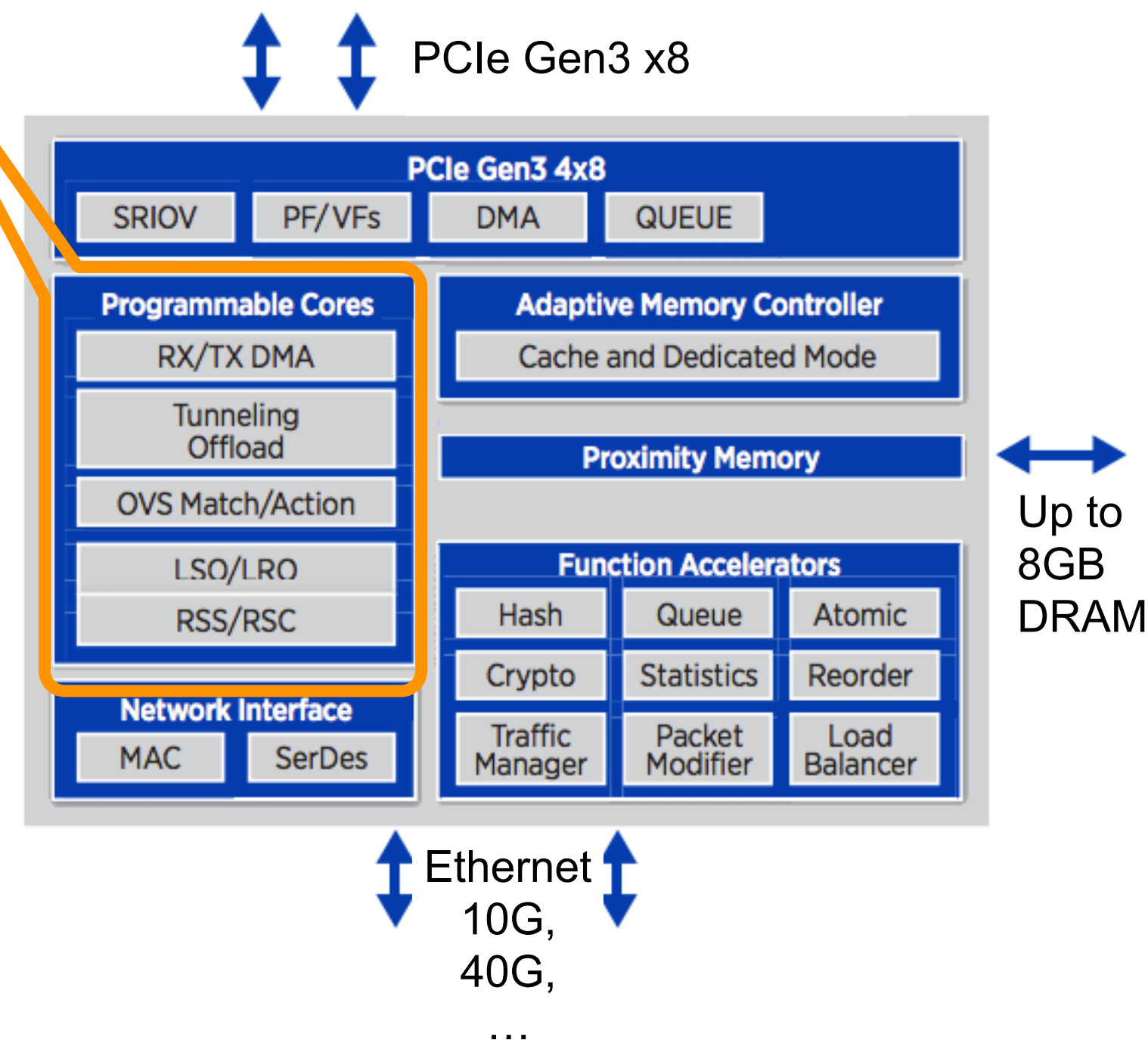
- Highly parallel multithreaded architecture (8 threads / core) for high throughput
- Purpose built Flow Processing Cores (72) maximize flexibility
- H/W accelerators further maximize efficiency (throughput/watt)

## Fully software defined feature set - examples:

- Network and PCIe SR-IOV / VirtIO RX/TX with stateless offloads
- Flexible tunneling support (e.g. VXLAN, GRE, VLAN, MPLS, NSH)
- Flexible Match/Action processing - many packet fields / protocols
- Highly scalable and fine grained security policies

External DDR3 accommodates millions of flows / rules

Convenient programmability using P4 and C



## Network Flow Processor 4xxx (used on Agilio-CX SmartNICs)

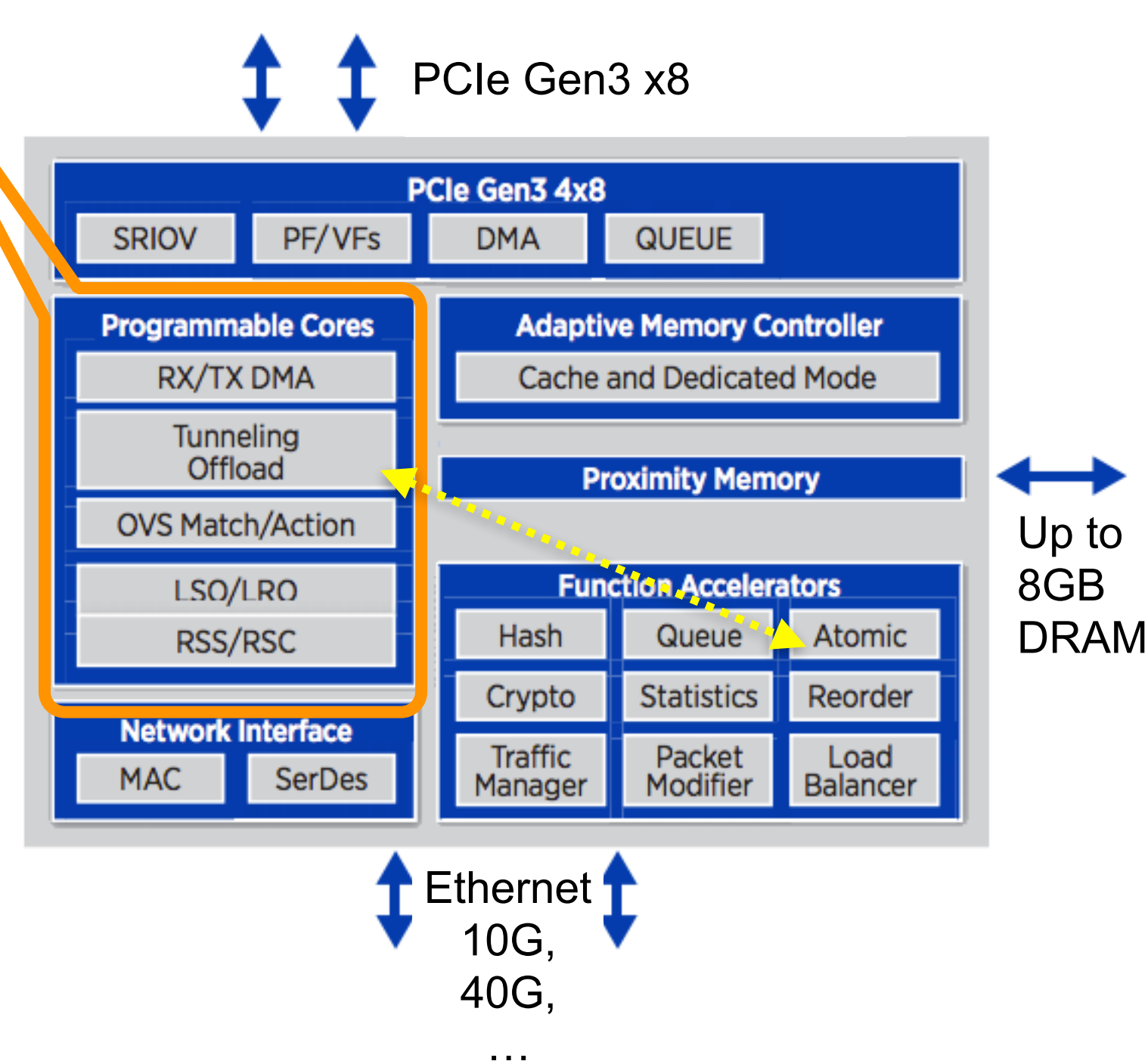
- Highly parallel multithreaded architecture (8 threads / core) for high throughput
- Purpose built Flow Processing Cores (72) maximize flexibility
- H/W accelerators further maximize efficiency (throughput/watt)

## Fully software defined feature set - examples:

- Network and PCIe SR-IOV / VirtIO RX/TX with stateless offloads
- Flexible tunneling support (e.g. VXLAN, GRE, VLAN, MPLS, NSH)
- Flexible Match/Action processing - many packet fields / protocols
- Highly scalable and fine grained security policies

External DDR3 accommodates millions of flows / rules

Convenient programmability using P4 and C



## Network Flow Processor 4xxx (used on Agilio-CX SmartNICs)

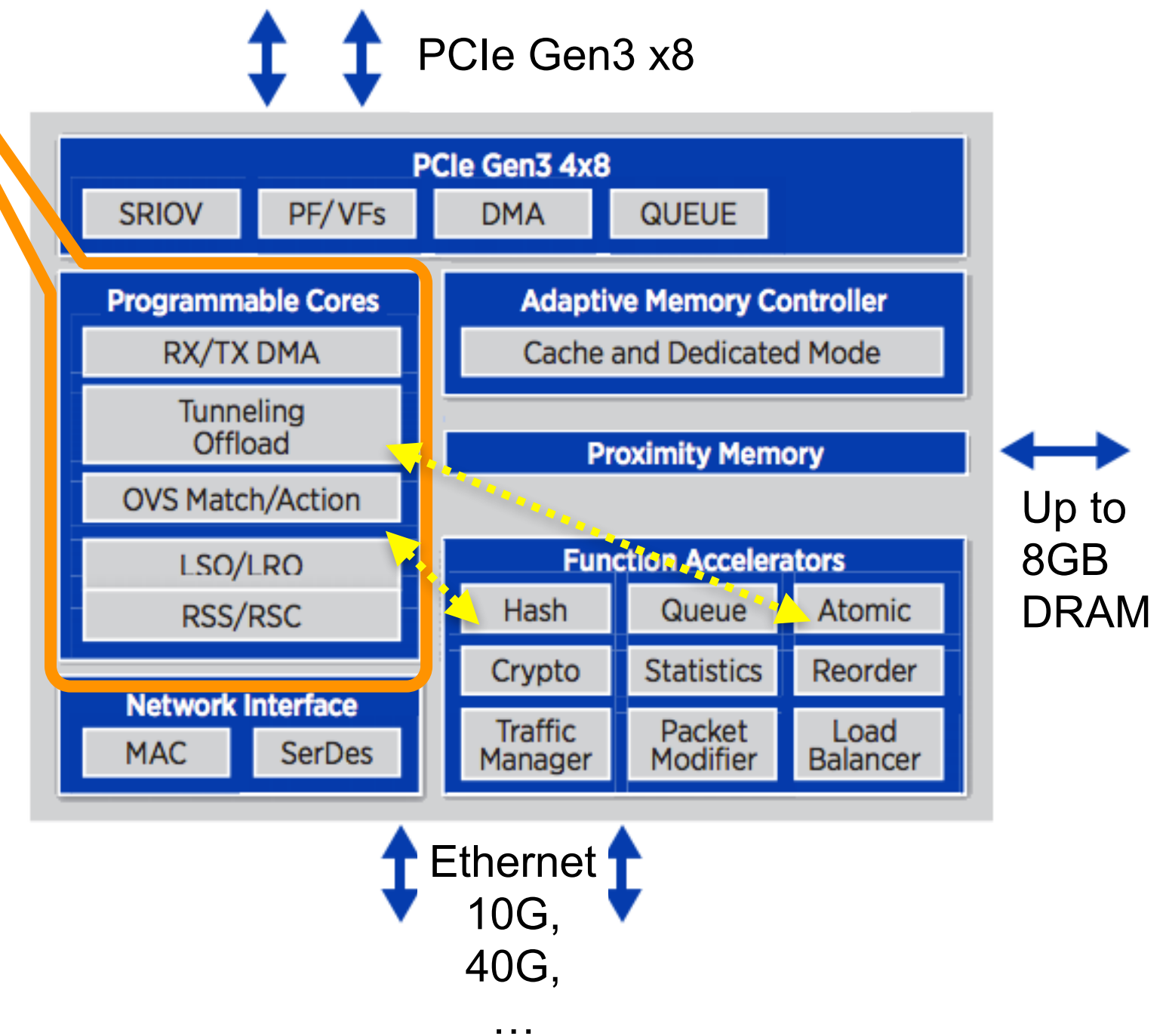
- Highly parallel multithreaded architecture (8 threads / core) for high throughput
- Purpose built Flow Processing Cores (72) maximize flexibility
- H/W accelerators further maximize efficiency (throughput/watt)

## Fully software defined feature set - examples:

- Network and PCIe SR-IOV / VirtIO RX/TX with stateless offloads
- Flexible tunneling support (e.g. VXLAN, GRE, VLAN, MPLS, NSH)
- Flexible Match/Action processing - many packet fields / protocols
- Highly scalable and fine grained security policies

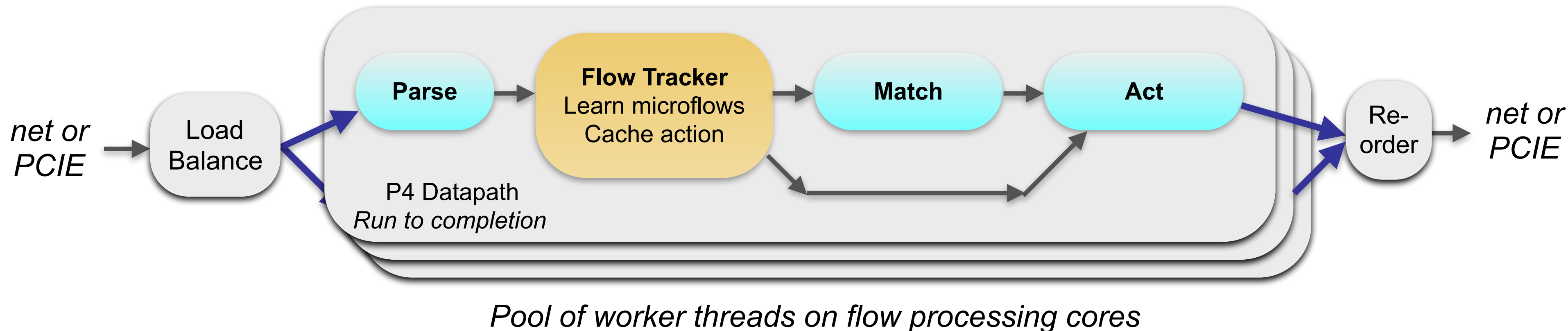
External DDR3 accommodates millions of flows / rules

Convenient programmability using P4 and C



- Load balancer distributes each packet to next available thread for optimum throughput
- Hardware assisted reordering ensures packet order is maintained
- Flow tracker statefully learns / tracks millions of sessions
- Matching performed using DRAM-backed tables - capacity millions of entries
- Actions efficiently performed in on-chip memory

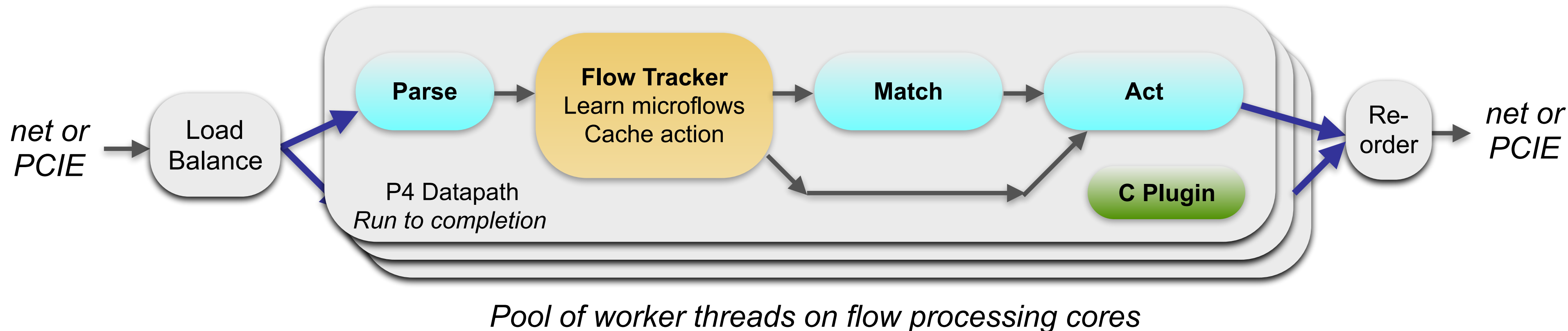
 = Ring / Work Queue (multi producer / consumer)





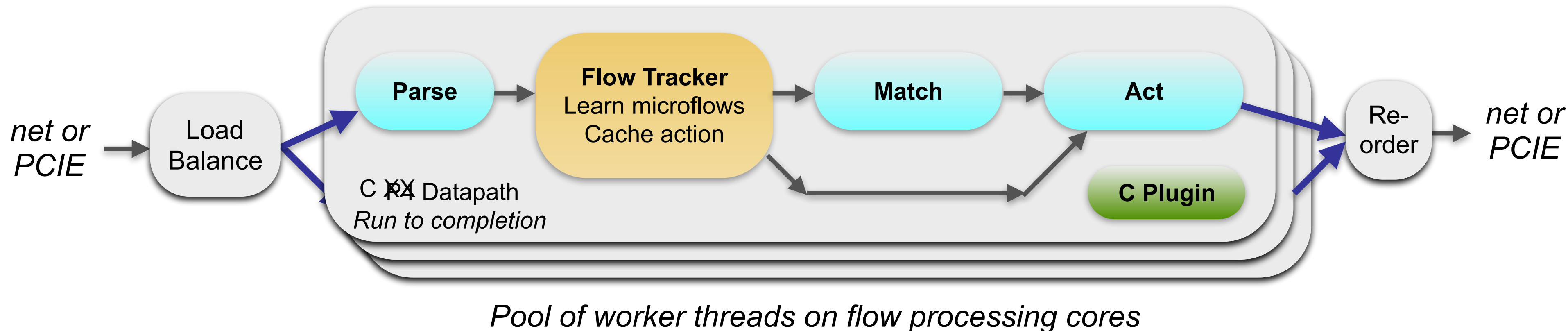
- Load balancer distributes each packet to next available thread for optimum throughput
- Hardware assisted reordering ensures packet order is maintained
- Flow tracker statefully learns / tracks millions of sessions
- Matching performed using DRAM-backed tables - capacity millions of entries
- Actions efficiently performed in on-chip memory

 = Ring / Work Queue (multi producer / consumer)



- Load balancer distributes each packet to next available thread for optimum throughput
- Hardware assisted reordering ensures packet order is maintained
- Flow tracker statefully learns / tracks millions of sessions
- Matching performed using DRAM-backed tables - capacity millions of entries
- Actions efficiently performed in on-chip memory

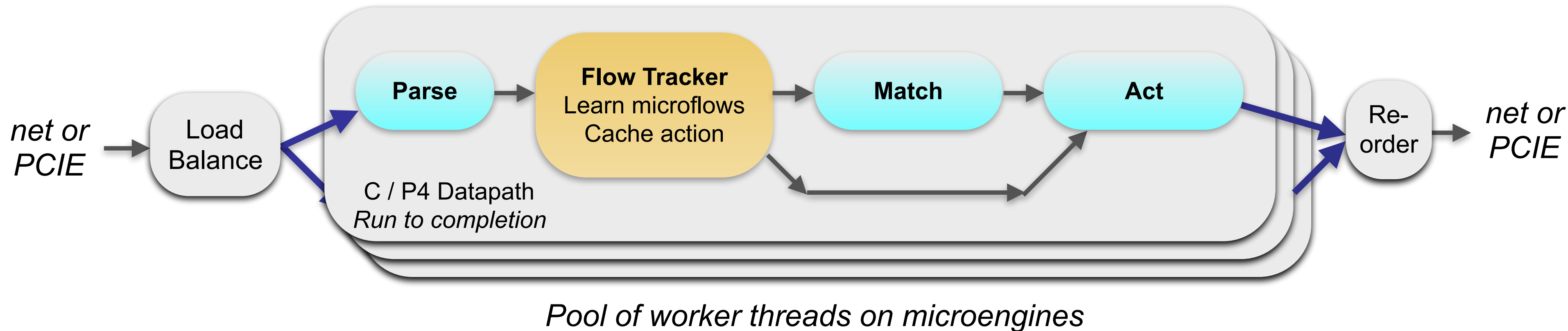
 = Ring / Work Queue (multi producer / consumer)



# Datapath Distributed over Microengines

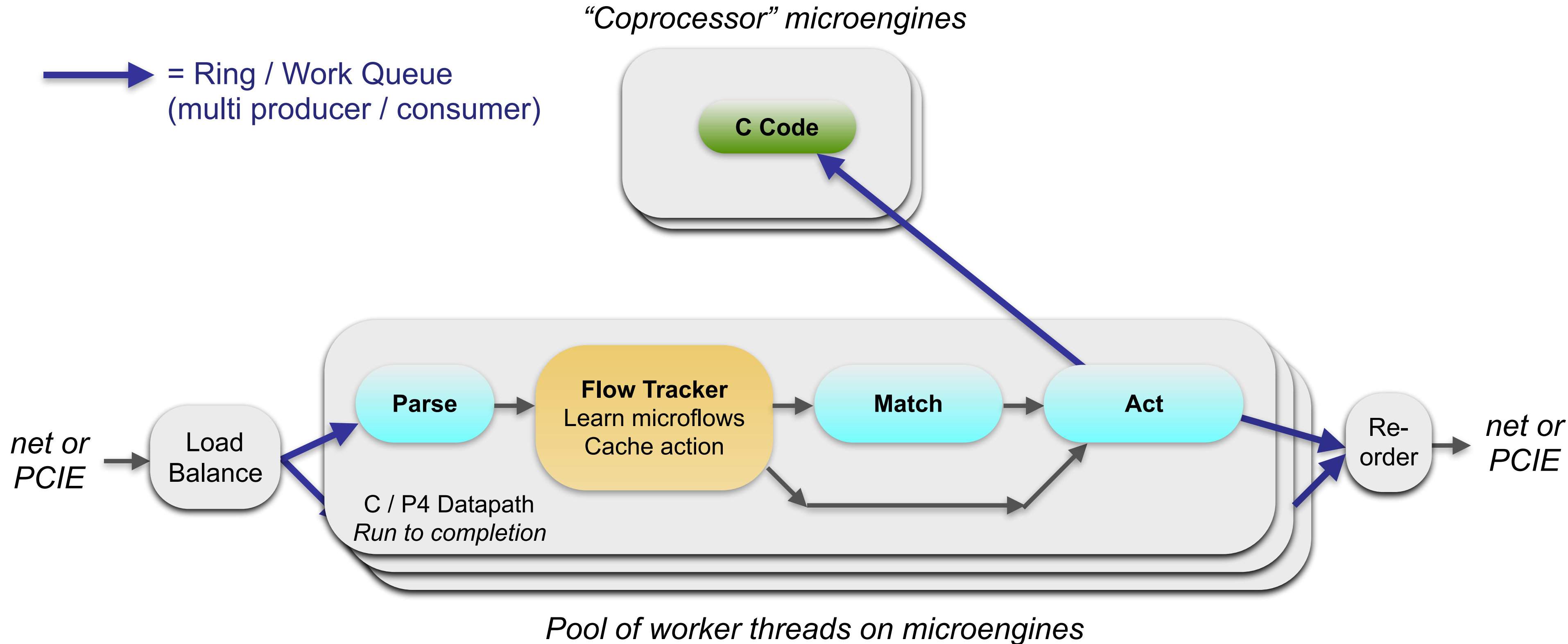
- Worker uses ring to forward packets to other microengine(s)

➔ = Ring / Work Queue  
(multi producer / consumer)



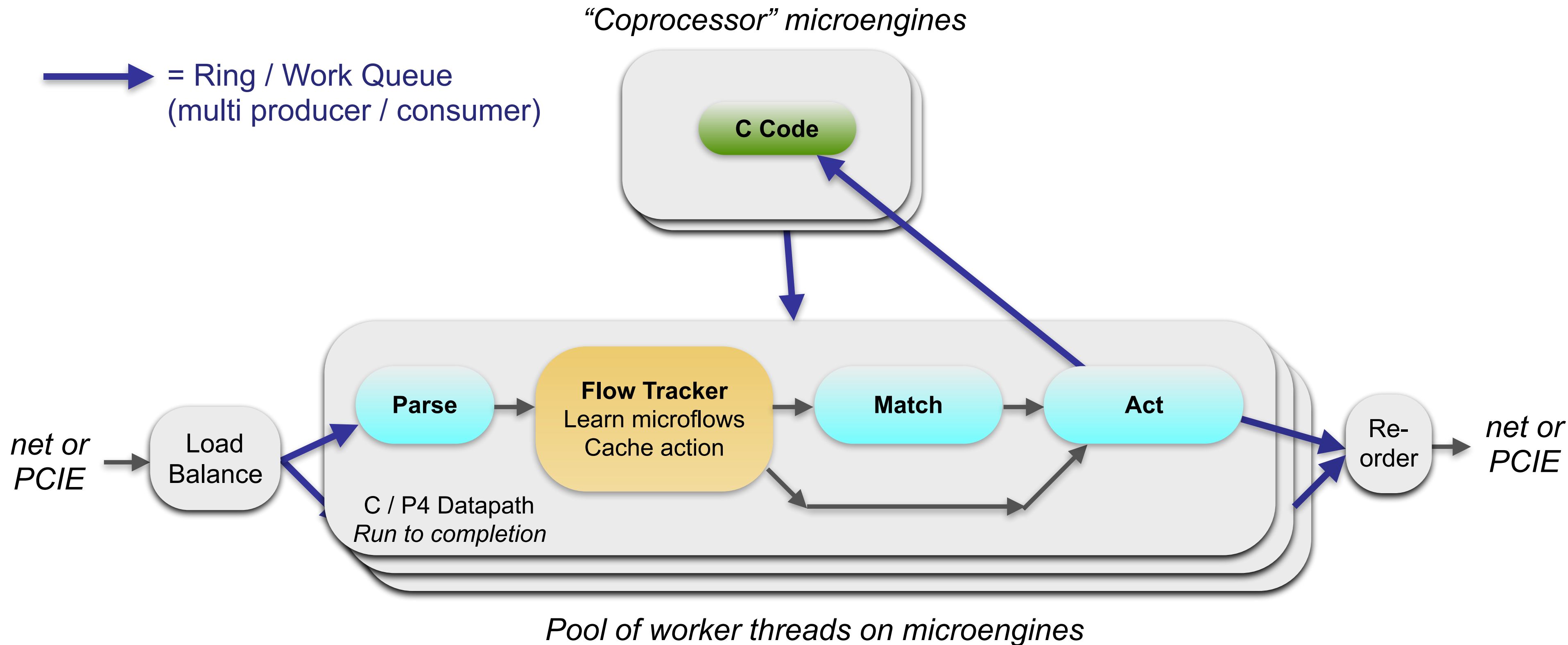
# Datapath Distributed over Microengines

- Worker uses ring to forward packets to other microengine(s)



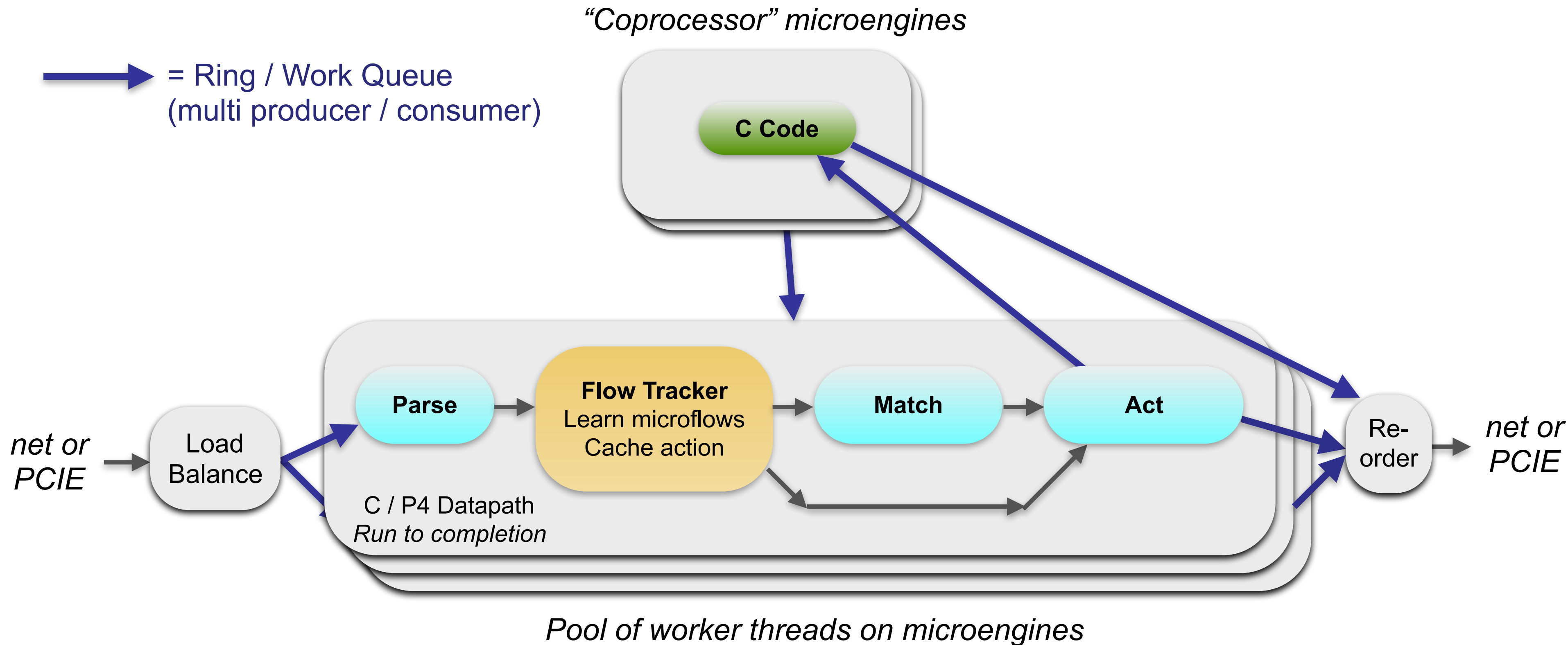
# Datapath Distributed over Microengines

- Worker uses ring to forward packets to other microengine(s)



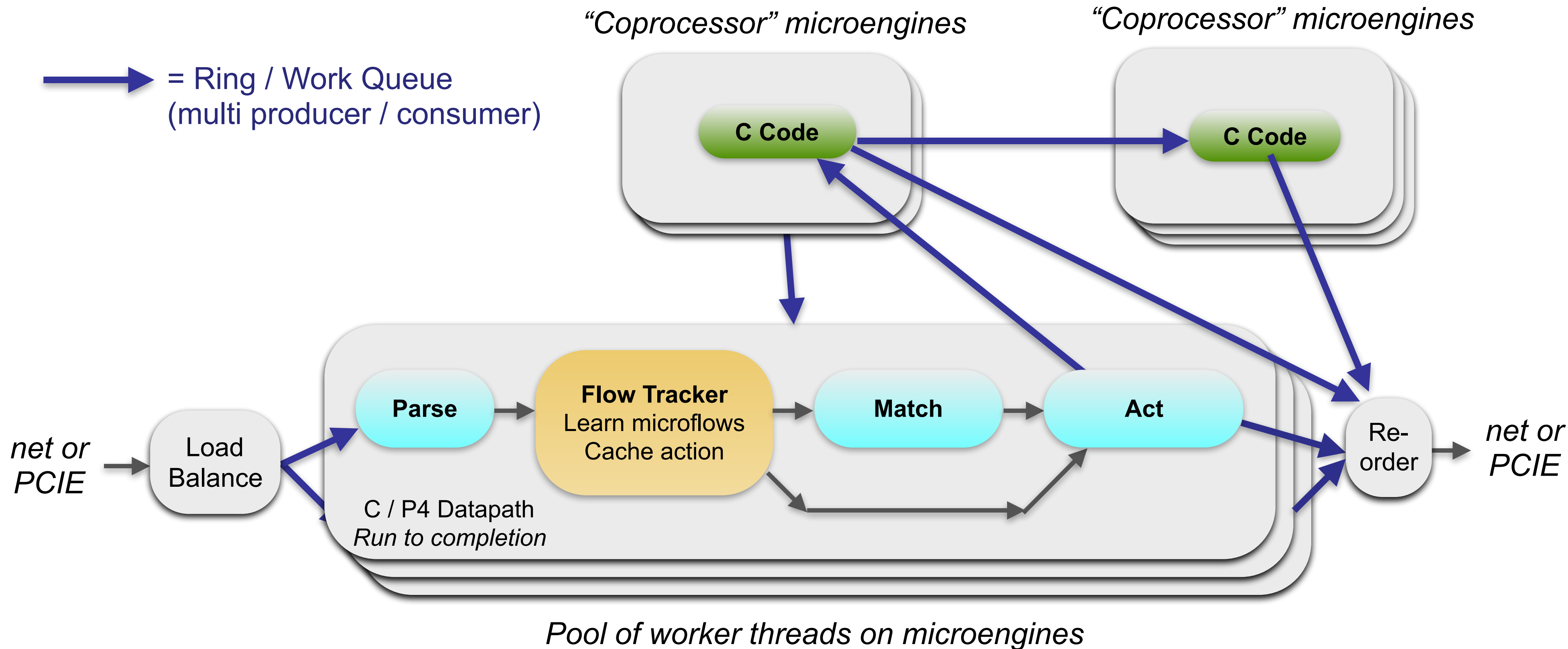
# Datapath Distributed over Microengines

- Worker uses ring to forward packets to other microengine(s)



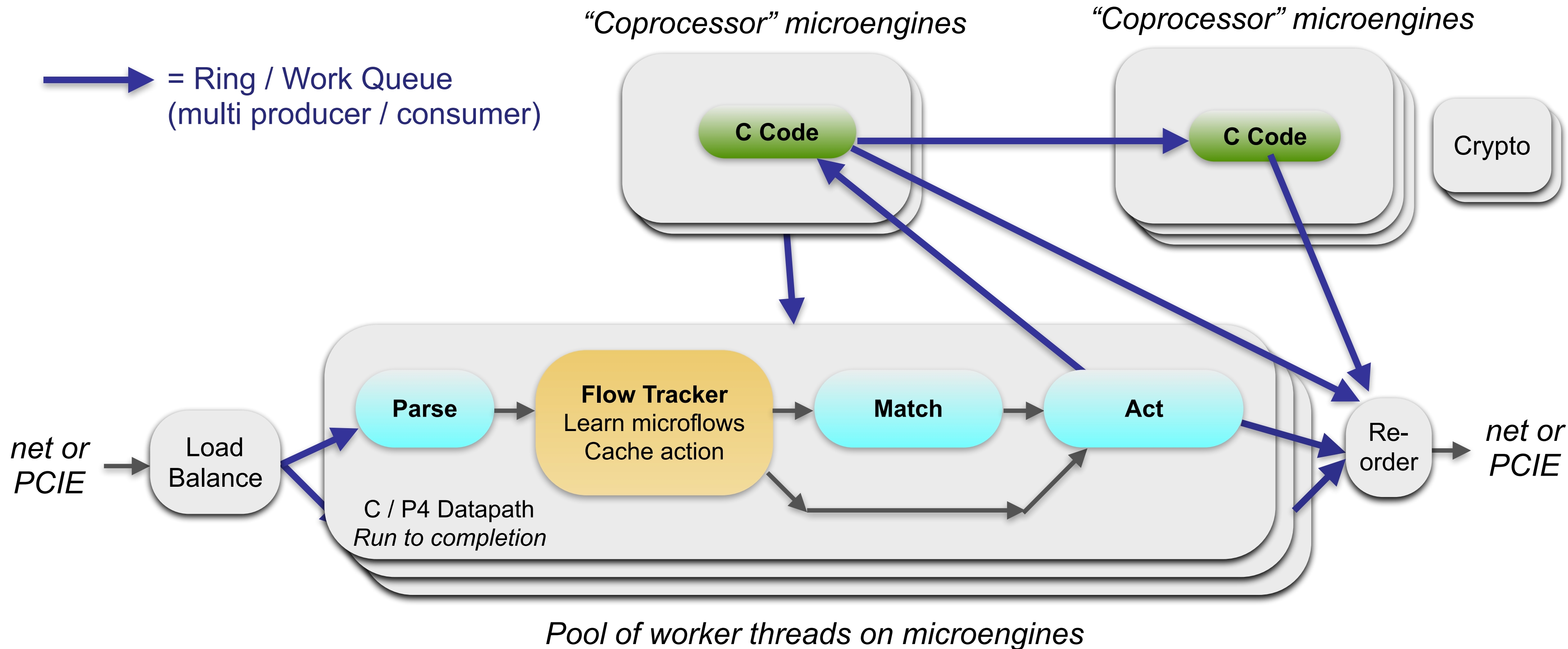
# Datapath Distributed over Microengines

- Worker uses ring to forward packets to other microengine(s)



# Datapath Distributed over Microengines

- Worker uses ring to forward packets to other microengine(s)





# Example: P4 “main” implementing a simple NIC

```
header_type eth_hdr {
  fields {
    dst : 48;
    src : 48;
    etype : 16;
  }
}
```

```
header eth_hdr eth;
```

```
parser start {
  return eth_parse;
}
```

```
parser eth_parse {
  extract(eth);
  return ingress;
}
```

```
action drop_act() {
  drop();
}
```

```
action fwd_act(port) {
  modify_field(standard_metadata.egress_spec,
  port);
}
```

```
table in_tbl {
  reads {
    standard_metadata.ingress_port : exact;
  }
  actions {
    fwd_act;
    drop_act;
  }
}
```

```
control ingress {
  apply(in_tbl);
}
```

# Example: C Code

```
static void
init_me()
{
    uint32_t cmask;
    cmask = disable_ctxs();
    pkt_init_rx();
    pkt_init_tx();
    enable_ctxs(cmask);
}

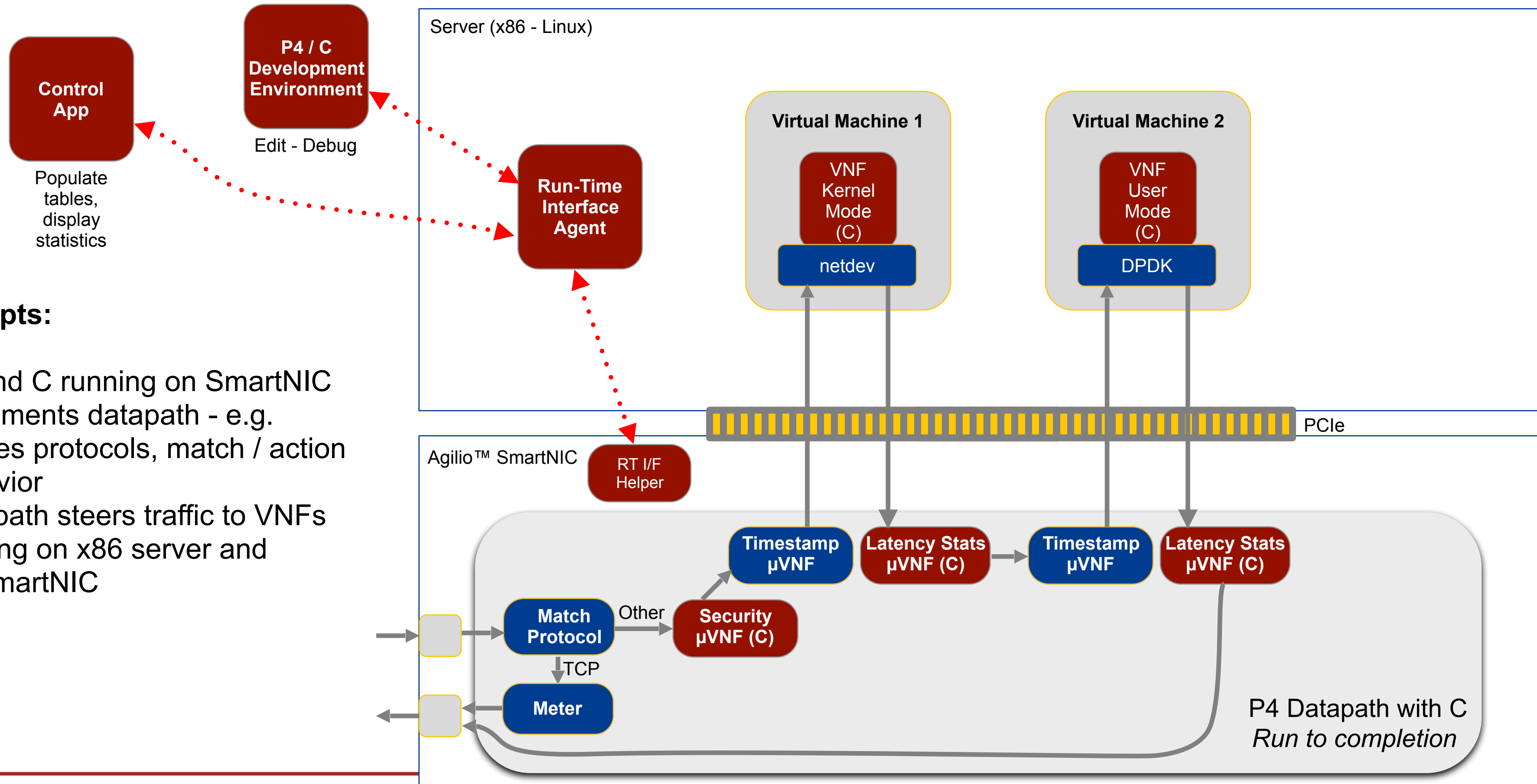
void
main()
{
    /* Configuration and initialization */
    if (ctx() == 0)
        init_me();

    for (;;) {
        /* half the threads receive from the wire and half from the host */
        if ((ctx() & 1) == 0) {
            if (pkt_rx_wire() == -1)
                goto transmit;
        } else {
            if (pkt_rx_host() == -1)
                goto transmit;
        }

        /* INSERT YOUR CODE HERE */
        /* Bump packet to the adjacent port */
        Pkt.p_dst = Pkt.p_src ^ 1;

        /* Attempt to send and drop if we encountered an error */
    transmit:
        pkt_tx();
    }
}
```

# Example of Fully Customized Datapath (P4 / C)



## Concepts:

- P4 and C running on SmartNIC implements datapath - e.g. defines protocols, match / action behavior
- Datapath steers traffic to VNFs running on x86 server and on SmartNIC

- Use Agilio™ SmartNICs with existing dataplanes
  - Use Agilio™ OVS (with / without Conntrack)
  - Use Agilio™ Contrail vRouter
  - Use Agilio™ eBPF/XDP
- Program Agilio™ SmartNICs (*following sessions - SDK, P4, OVS, eBPF/XDP*)
  - Program using P4, C, eBPF/XDP...
- Participate in open source and standards evolution:  
[openstack.org](http://openstack.org), [openvswitch.org](http://openvswitch.org), [opencontrail.org](http://opencontrail.org), [p4.org](http://p4.org), [iovisor.org](http://iovisor.org), [open-nfp.org](http://open-nfp.org),  
[opennetworking.org](http://opennetworking.org), [opensourcesdn.org](http://opensourcesdn.org), [opnfv.org](http://opnfv.org), [linuxfoundation.org](http://linuxfoundation.org)
  - Examples: P4 / OpenFlow callable run-time API (cross-body effort starting), acceleration APIs

***Increase flexibility, improve performance, free up server resources!***



NETRONOME

Thank You!

More information: [netronome.com](http://netronome.com)  
and: [open-nfp.org](http://open-nfp.org)