

Deliberately Congesting a Switch for Better Network Functions Performance

Mariano Scazzariello ^{*}, Tommaso Caiazzi [‡], Marco Chiesa ^{*}

^{*} KTH Royal Institute of Technology, Sweden – [‡] Roma Tre University, Italy

Network Functions Are Pervasive

Firewall



NAT

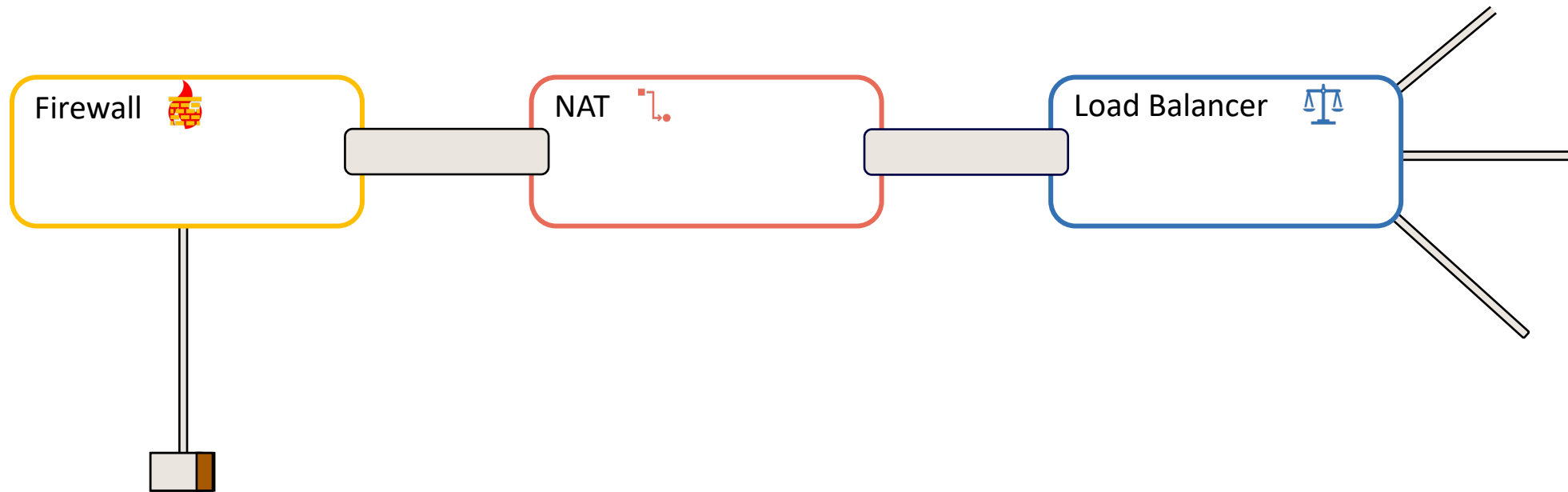


Load Balancer

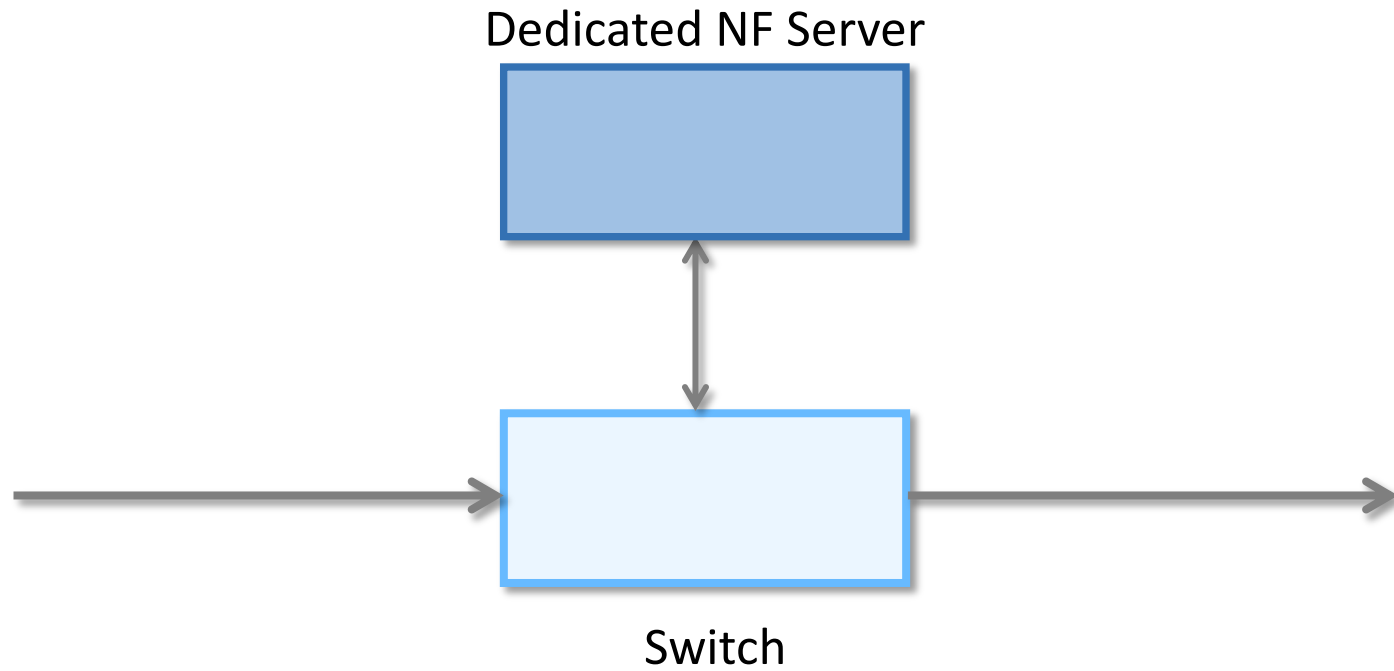


Network Functions Are Pervasive

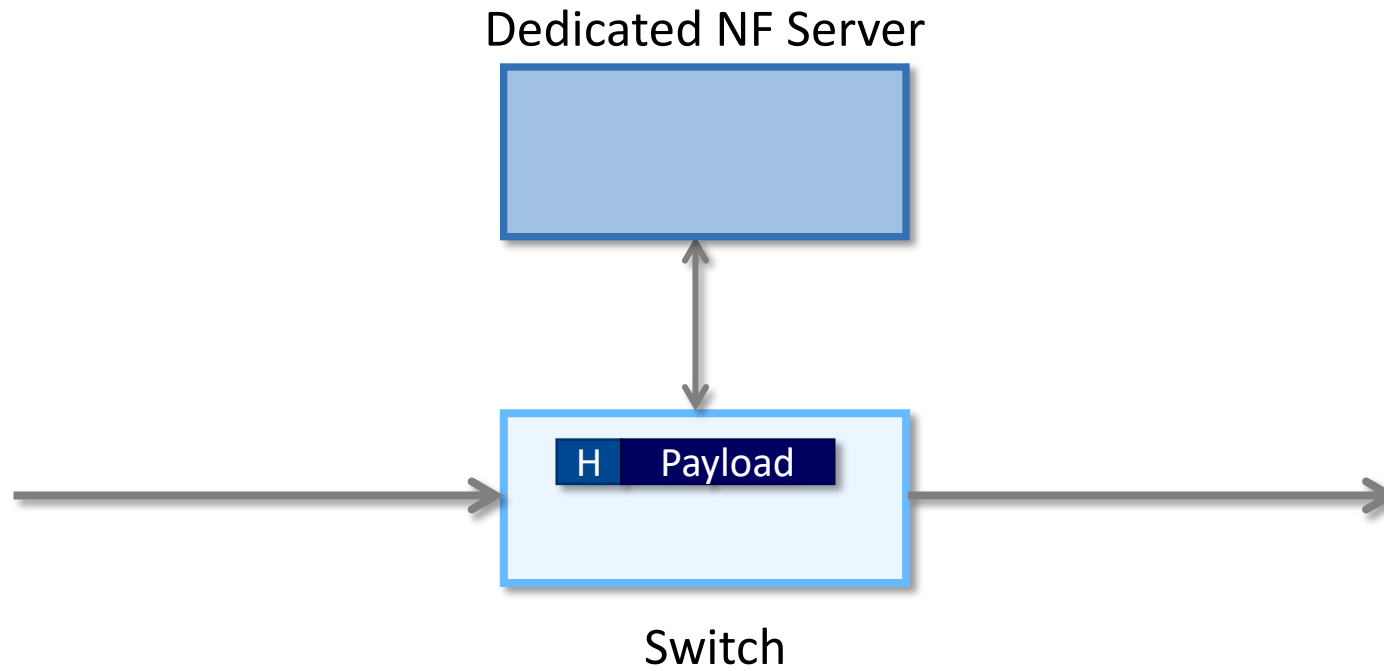
Network Functions Virtualization is an essential architectural paradigm of today's networks



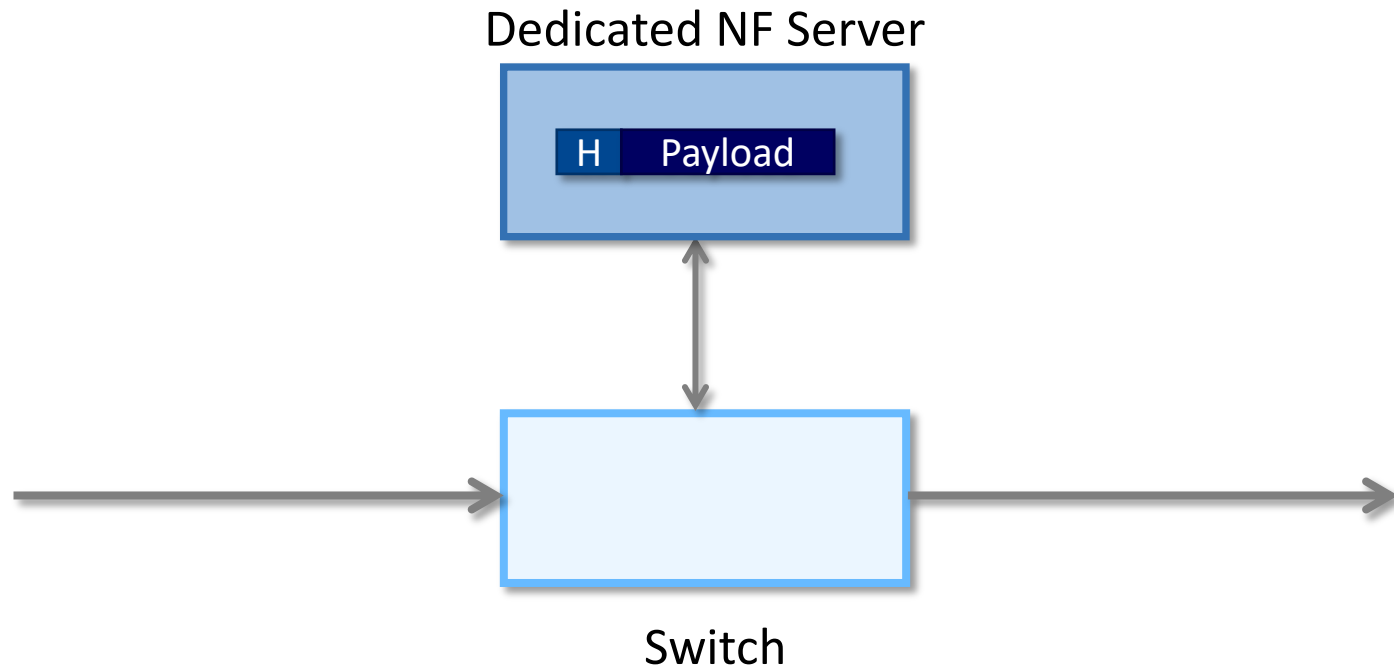
Advanced NF Deployments



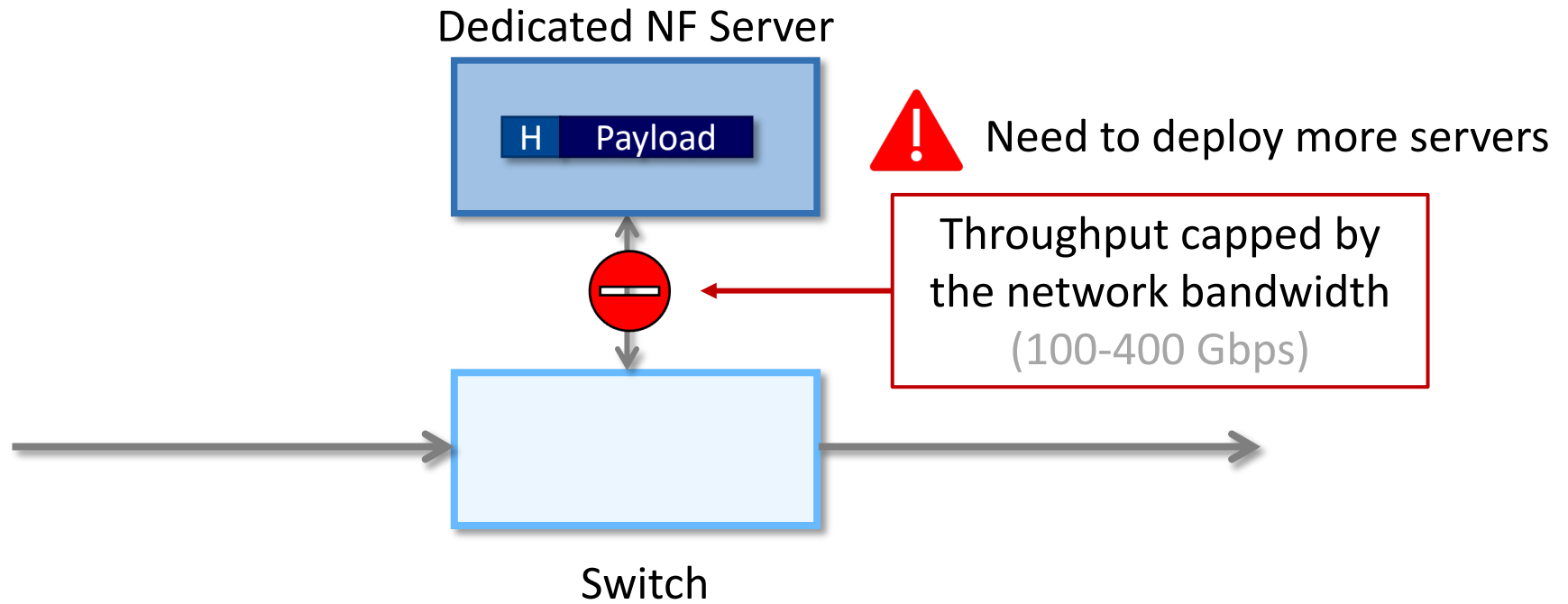
Advanced NF Deployments



Advanced NF Deployments



The Bandwidth Limit of Advanced NF Deployments

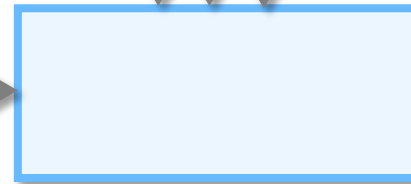
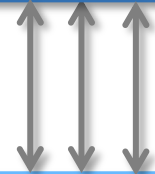
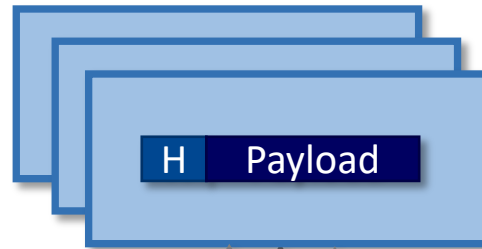


The Bandwidth Limit of Advanced NF Deployments



Not cost-effective!

Dedicated NF Servers



Switch



Need to deploy more servers

Throughput capped by the network bandwidth (100-400 Gbps)

Packet-Splitting Approaches

State-of-the-art packet-splitting solutions increase NFs throughput



Improve NF Performance

Better CPU cache exploitation



Free up Bandwidth

Reduce the number of dedicated NF servers

Lower energy consumption



Complex Deployment

Require end-host modifications

Require available shared resources

Packet-Splitting Approaches

State-of-the-art packet-splitting solutions increase NFs throughput

Storing Payloads in the On-NIC Memory (nicmem)



Storing Payloads on Shared Resources (Ribosome)










Storing a Few Bytes in ASIC SRAM (PayloadPark)



Packet-Splitting Approaches

State-of-the-art packet-splitting solutions increase NFs throughput

- Storing Payloads in the On-NIC Memory (nicmem)  
- Storing Payloads on Shared Resources (Ribosome)   
- Storing a Few Bytes in ASIC SRAM (PayloadPark)  

Packet-Splitting Approaches

State-of-the-art packet-splitting solutions increase NFs throughput

Storing Payloads in the On-NIC Memory (nicmem)



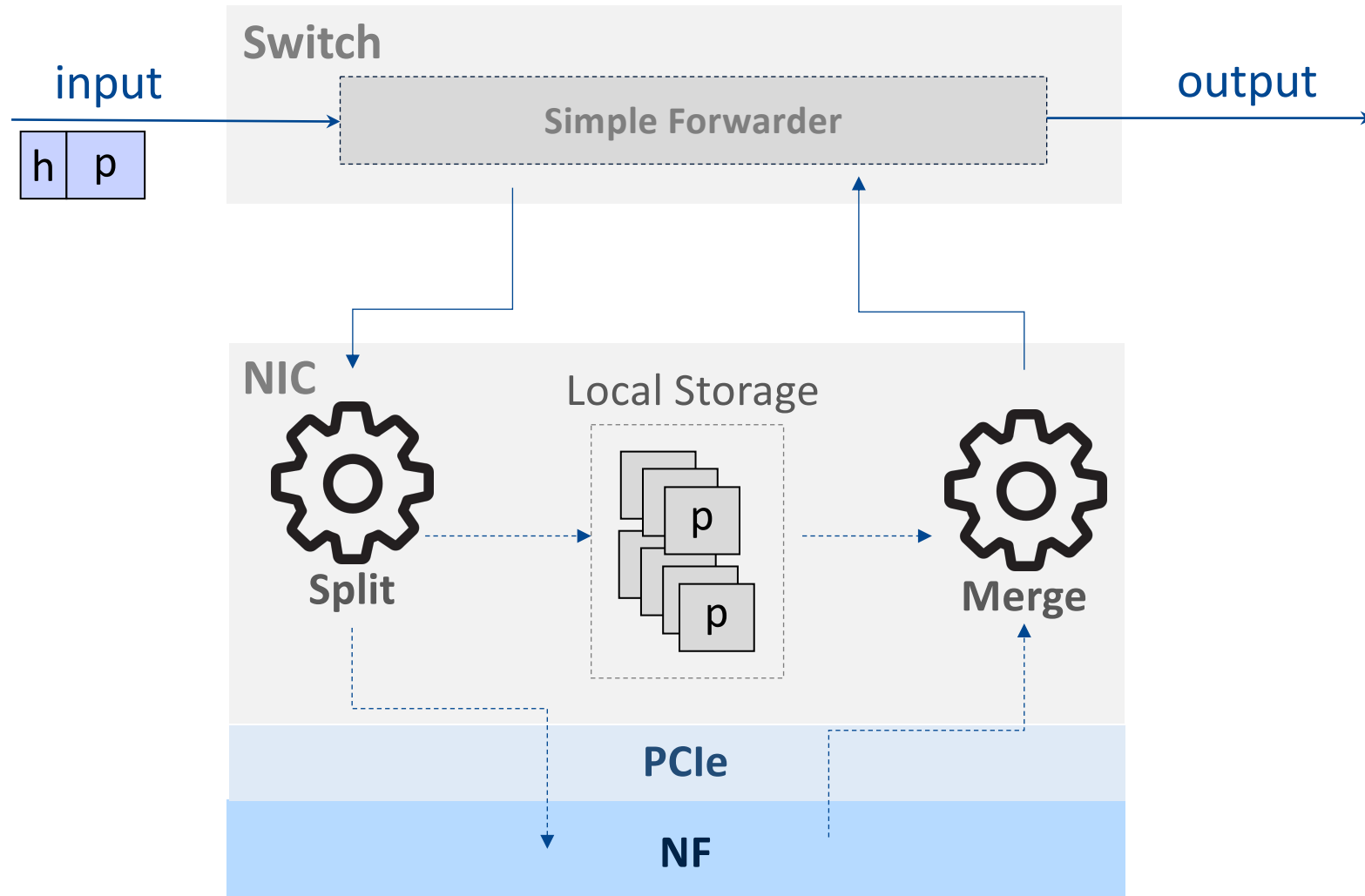
Storing Payloads on Shared Resources (Ribosome)



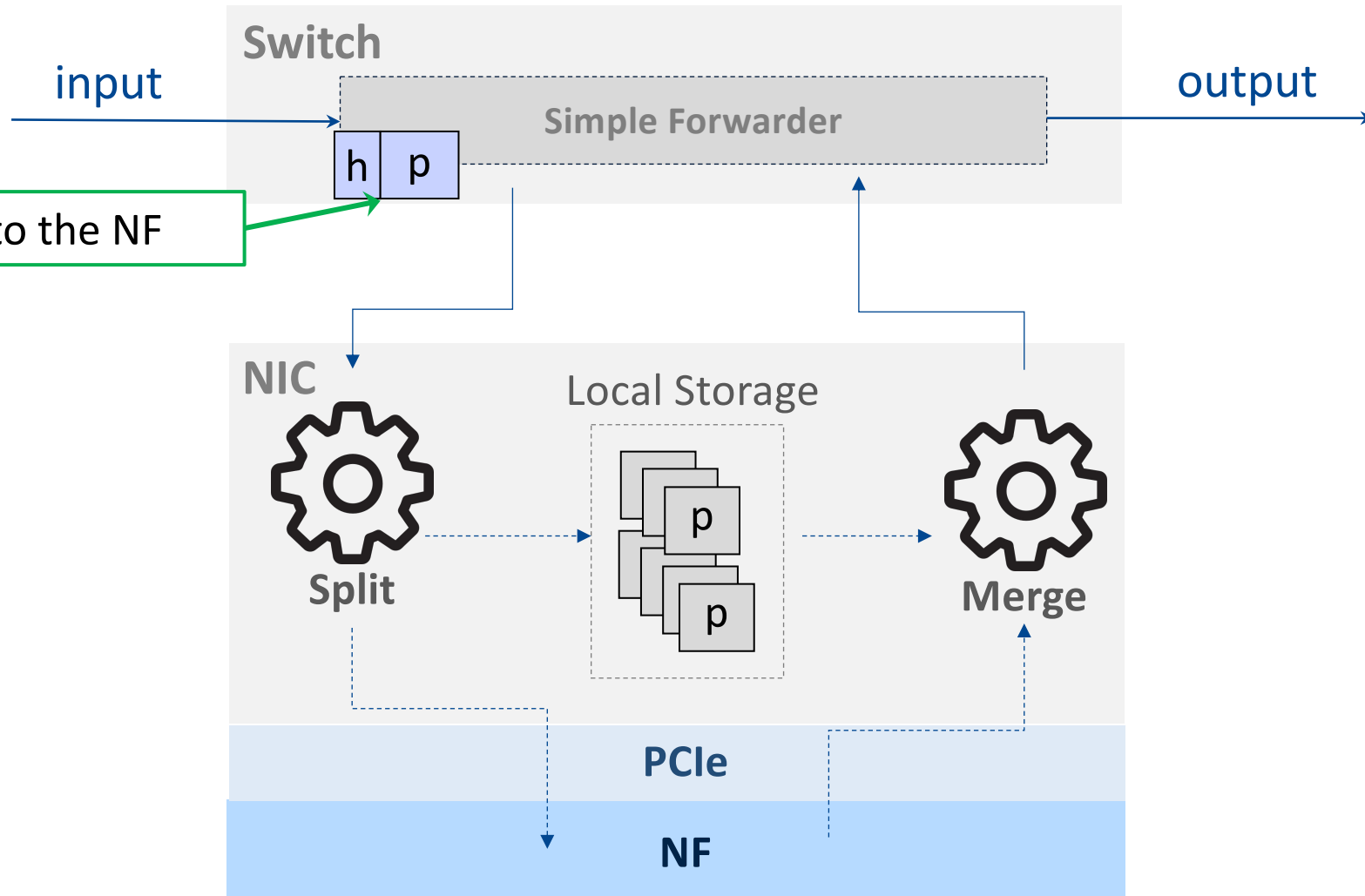
Storing a Few Bytes in ASIC SRAM (PayloadPark)



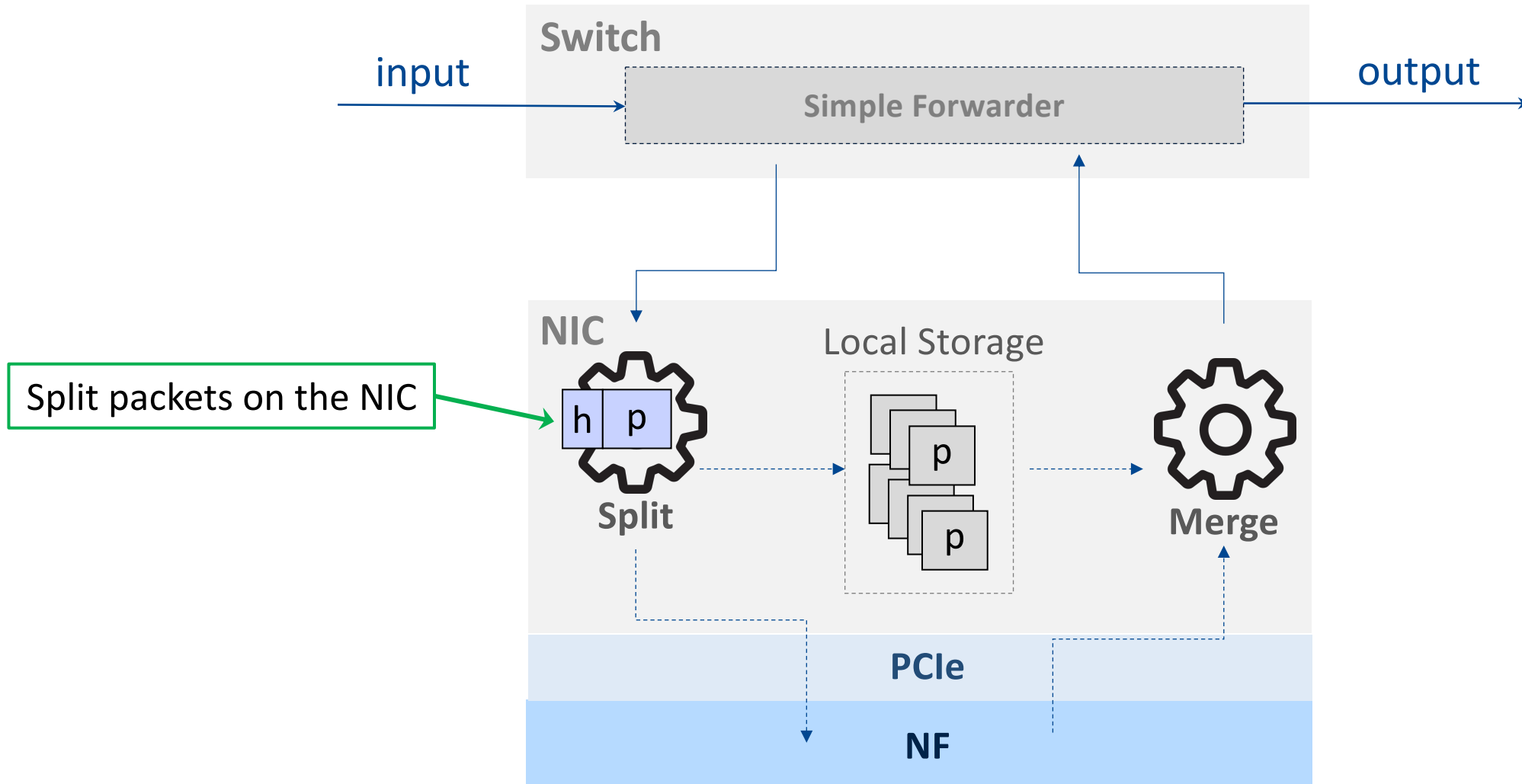
Storing Payloads in the On-NIC Memory (nicmem)



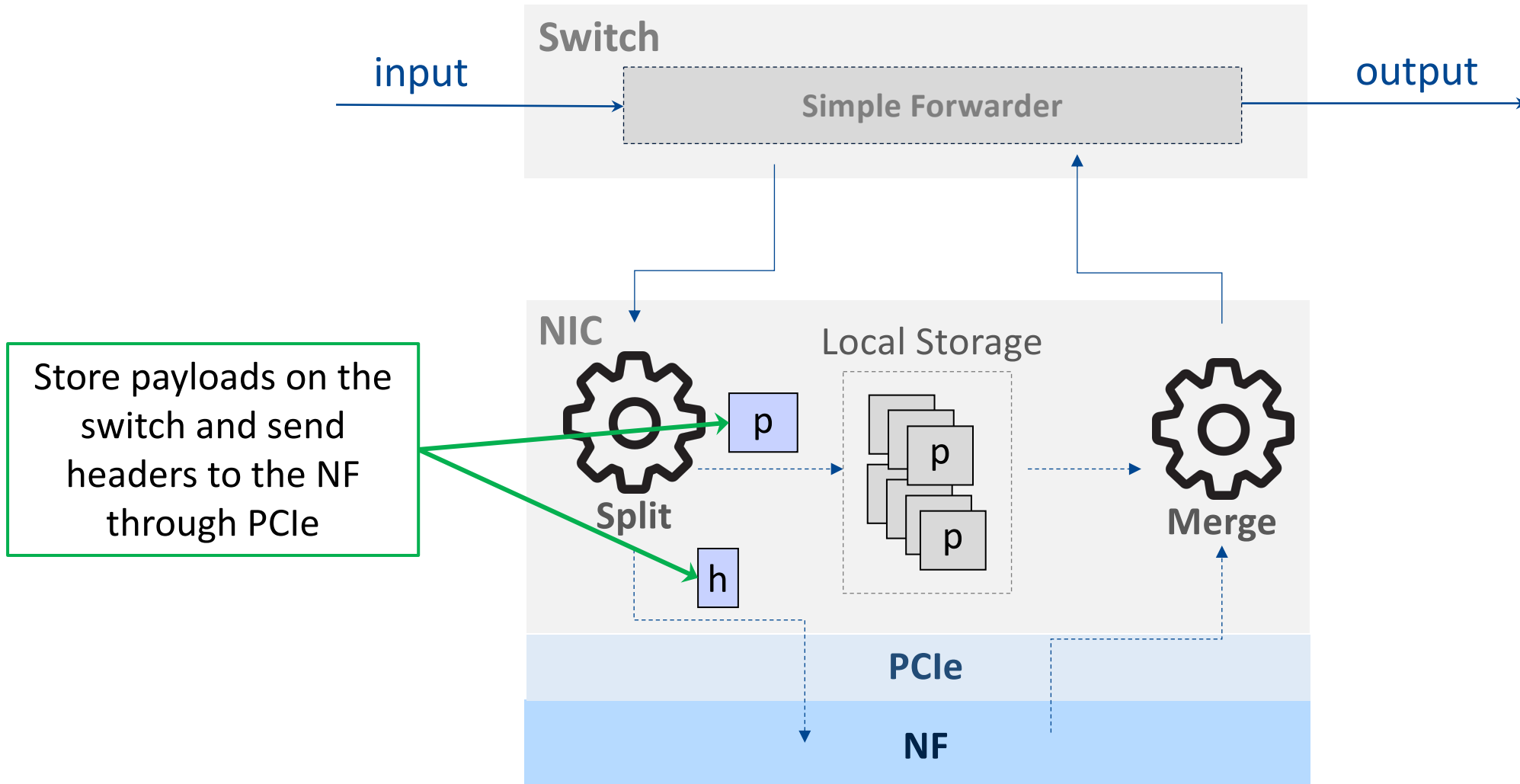
Storing Payloads in the On-NIC Memory (nicmem)



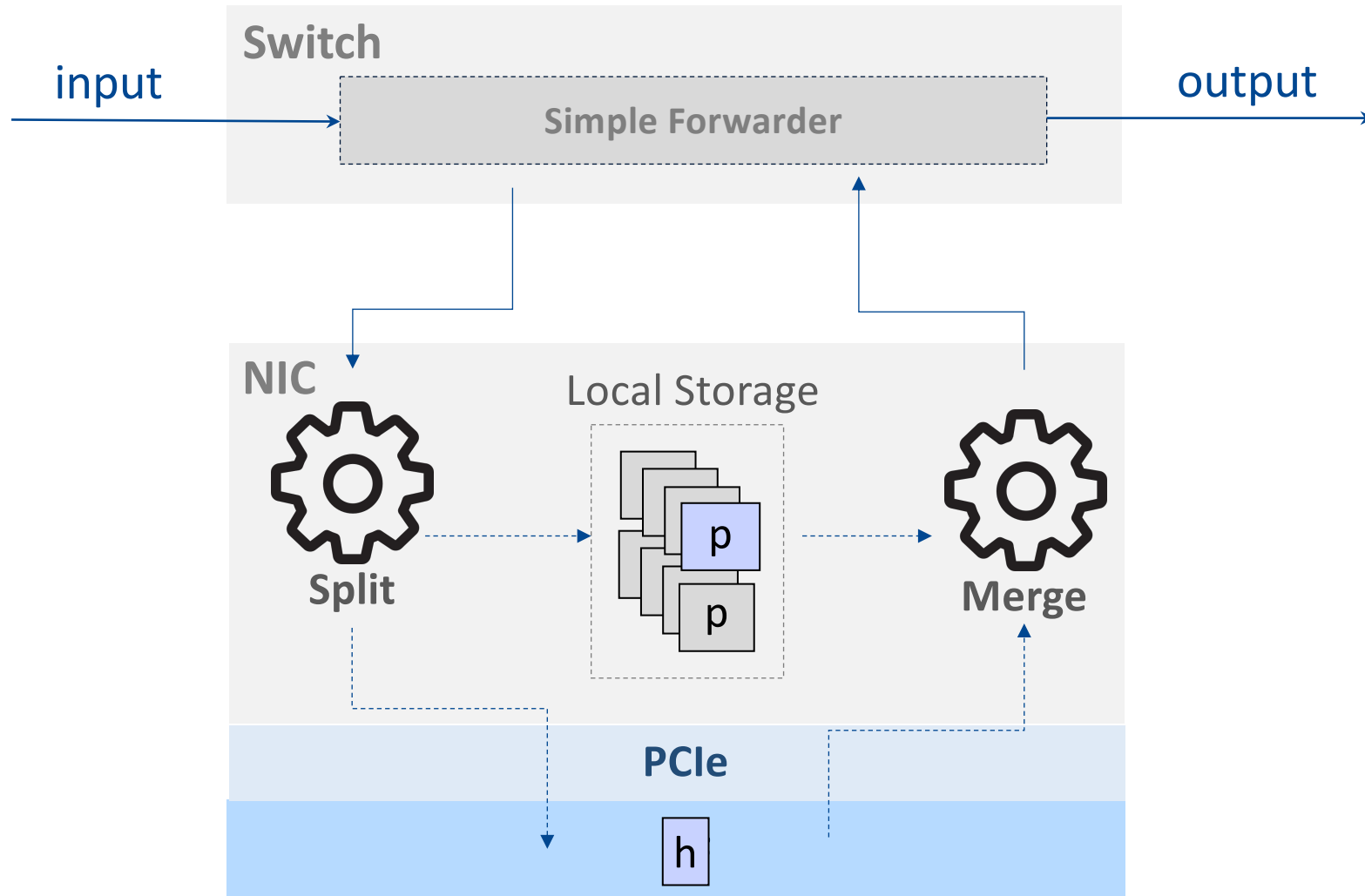
Storing Payloads in the On-NIC Memory (nicmem)



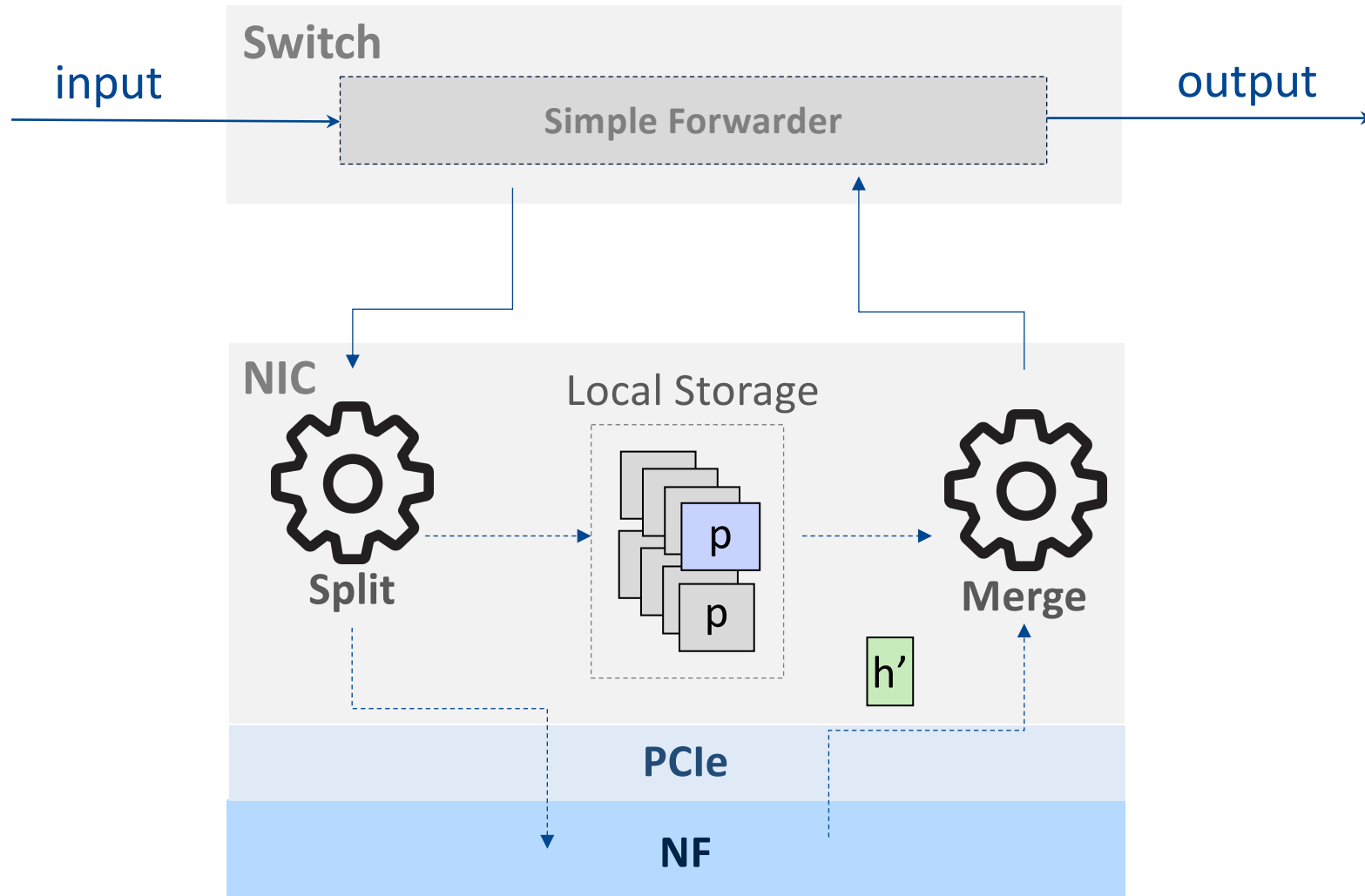
Storing Payloads in the On-NIC Memory (nicmem)



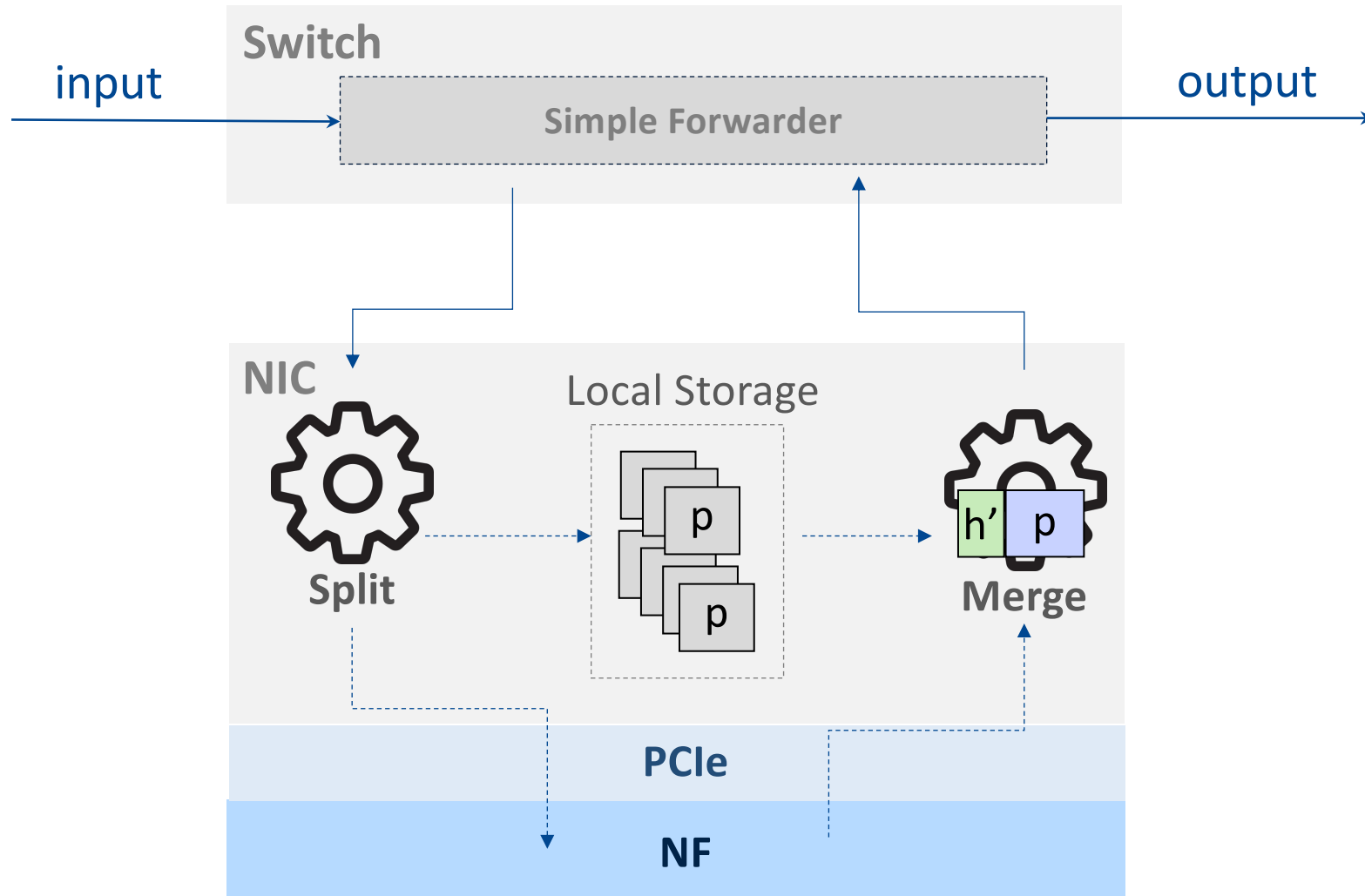
Storing Payloads in the On-NIC Memory (nicmem)



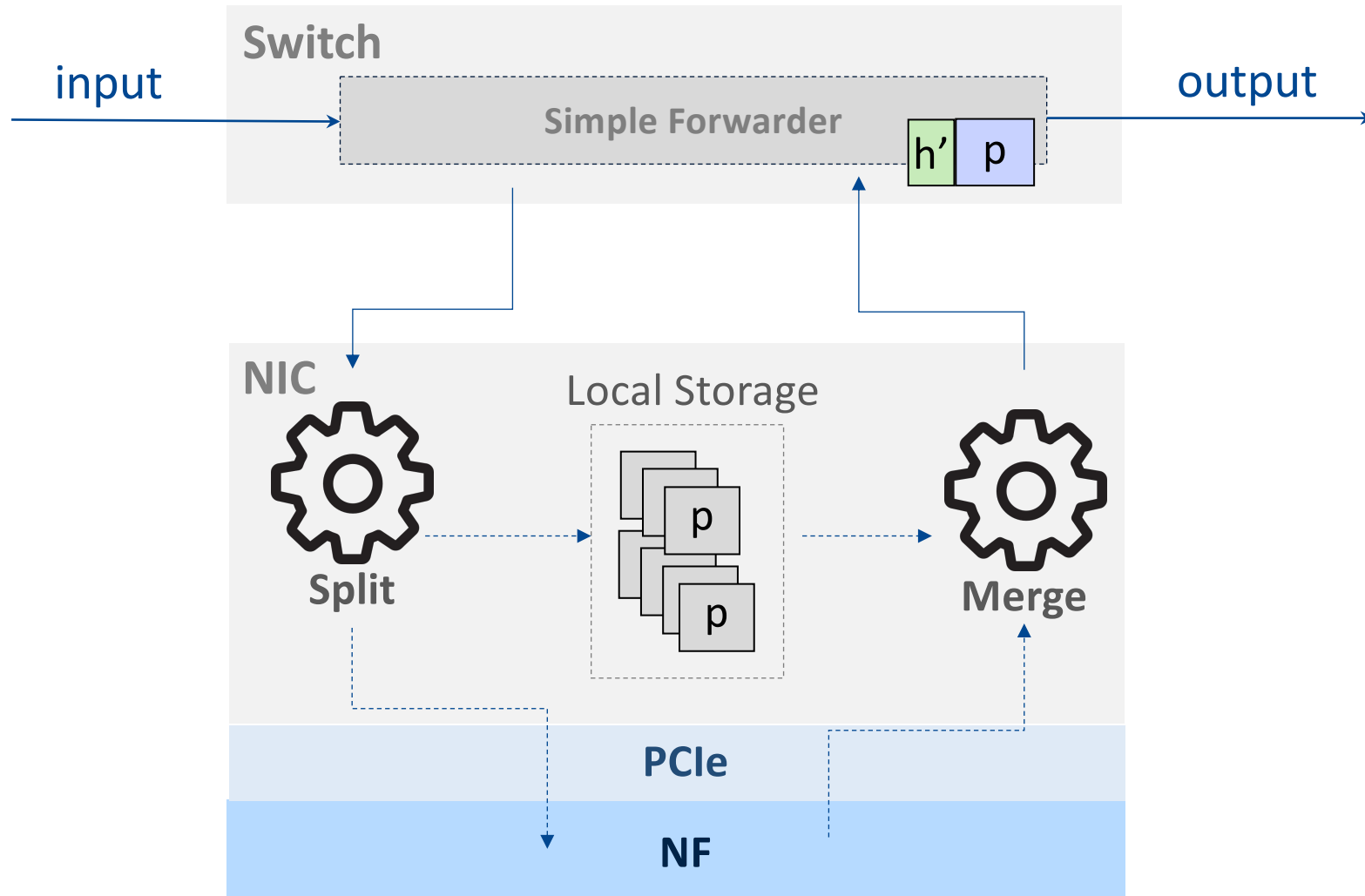
Storing Payloads in the On-NIC Memory (nicmem)



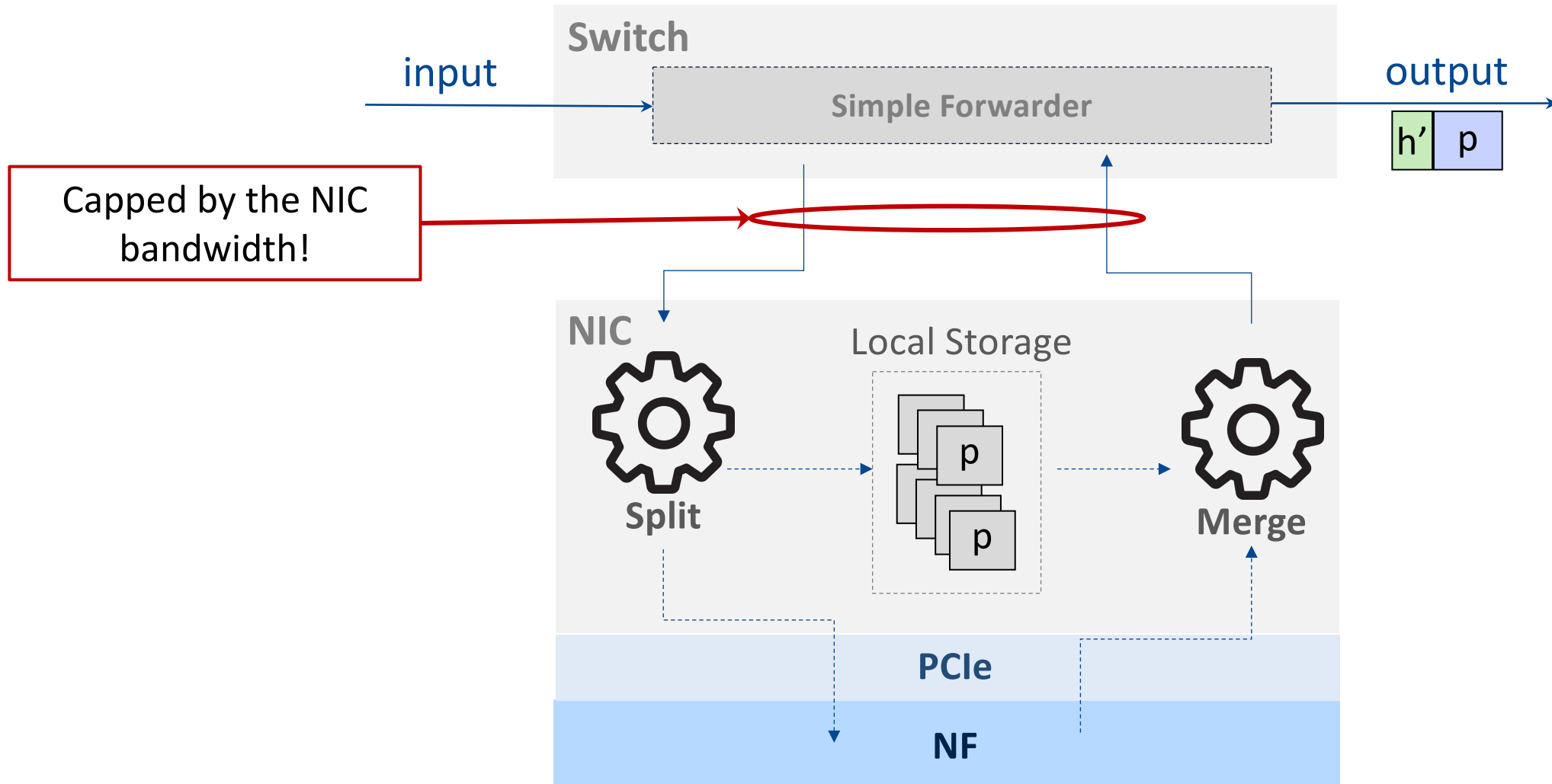
Storing Payloads in the On-NIC Memory (nicmem)



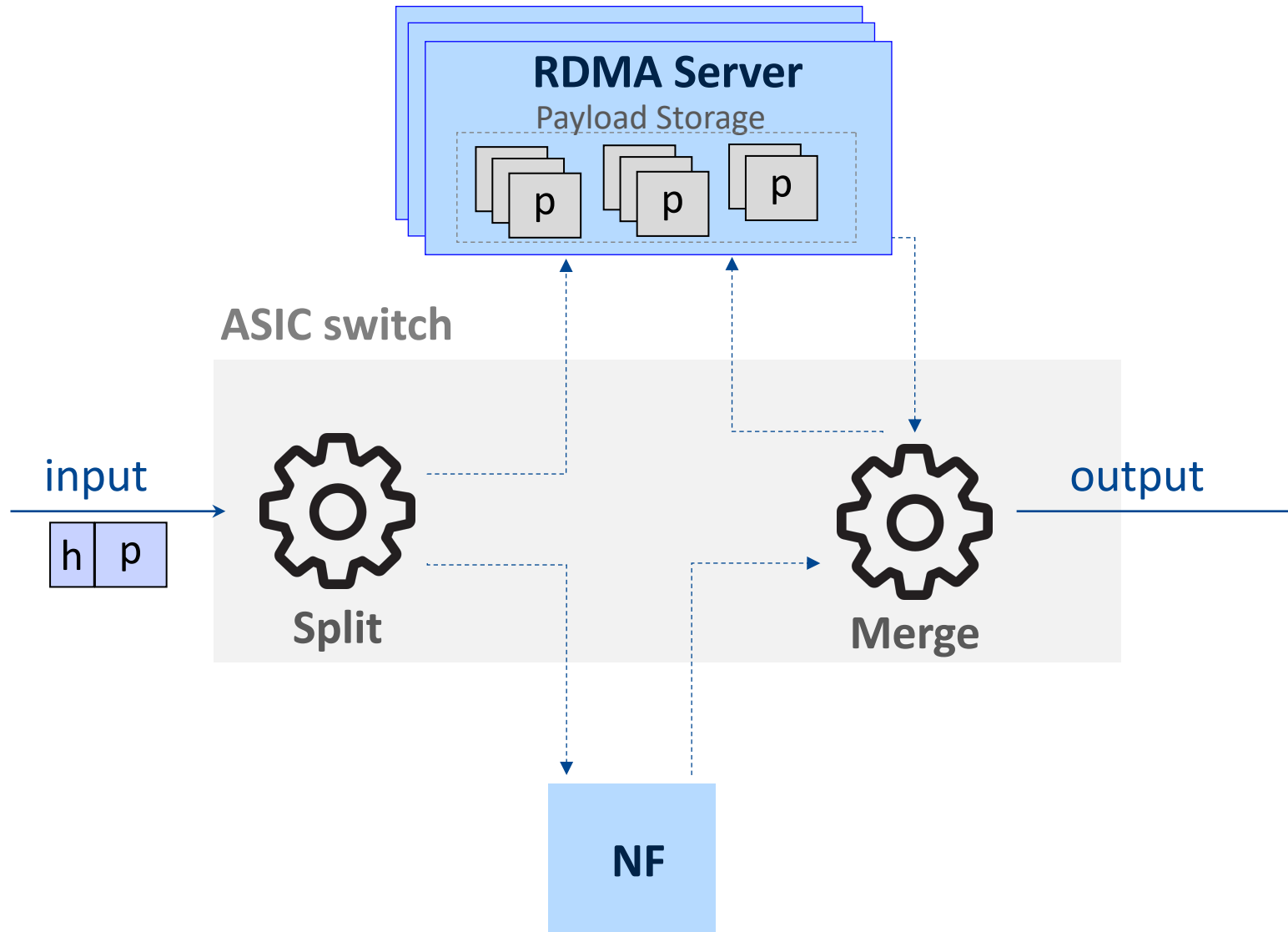
Storing Payloads in the On-NIC Memory (nicmem)



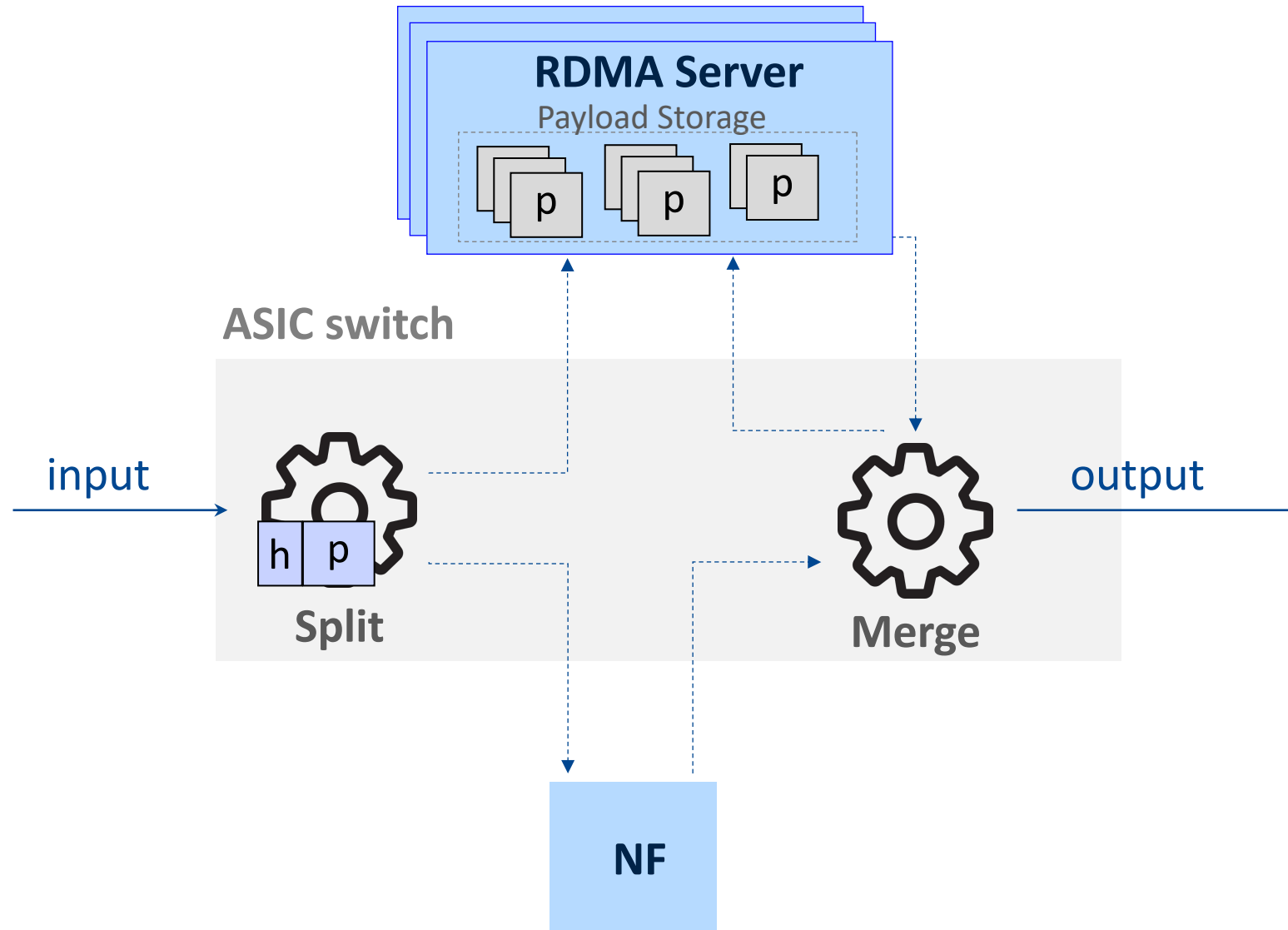
Storing Payloads in the On-NIC Memory (nicmem)



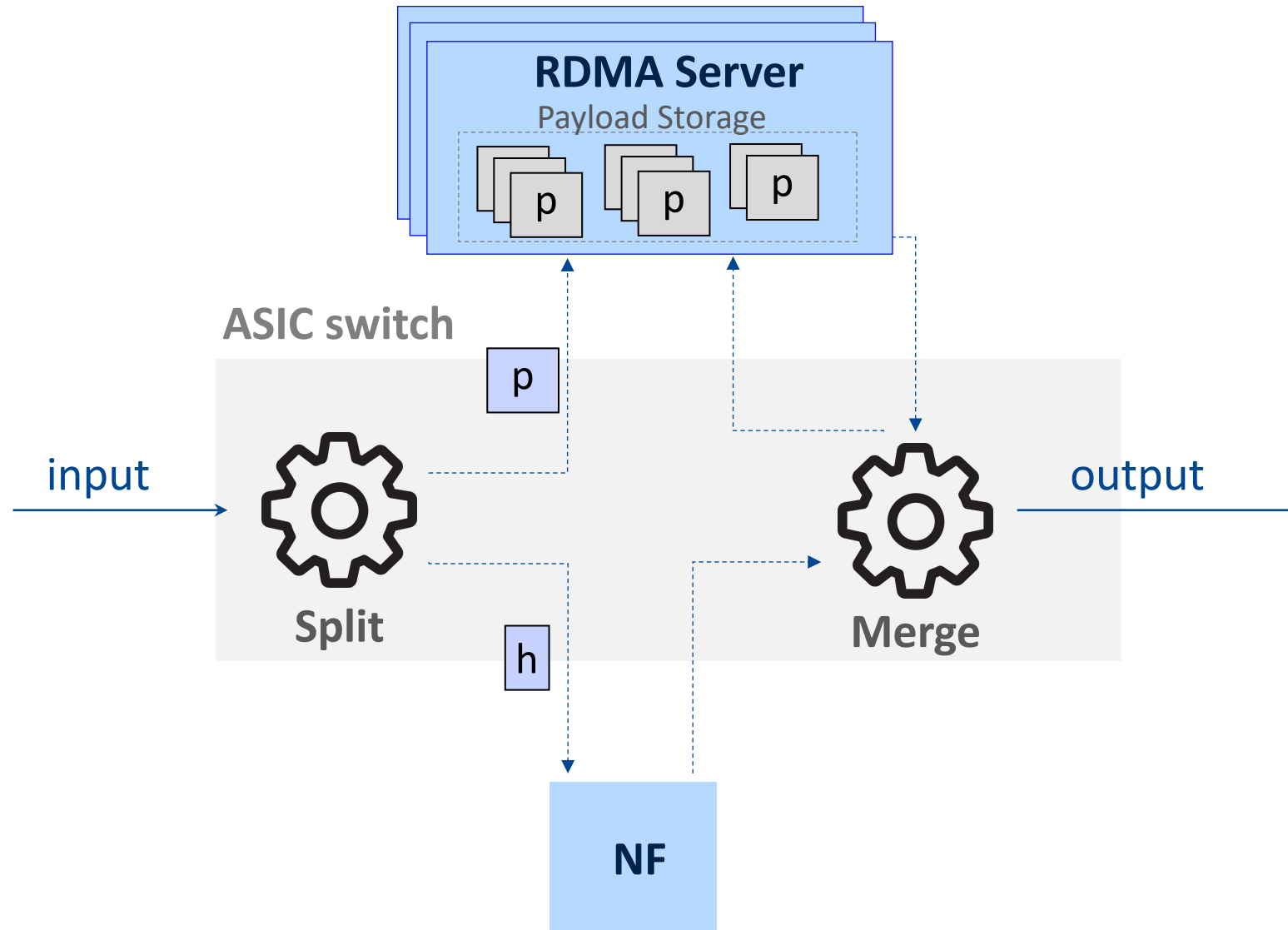
Storing Payloads on Shared Resources (Ribosome)



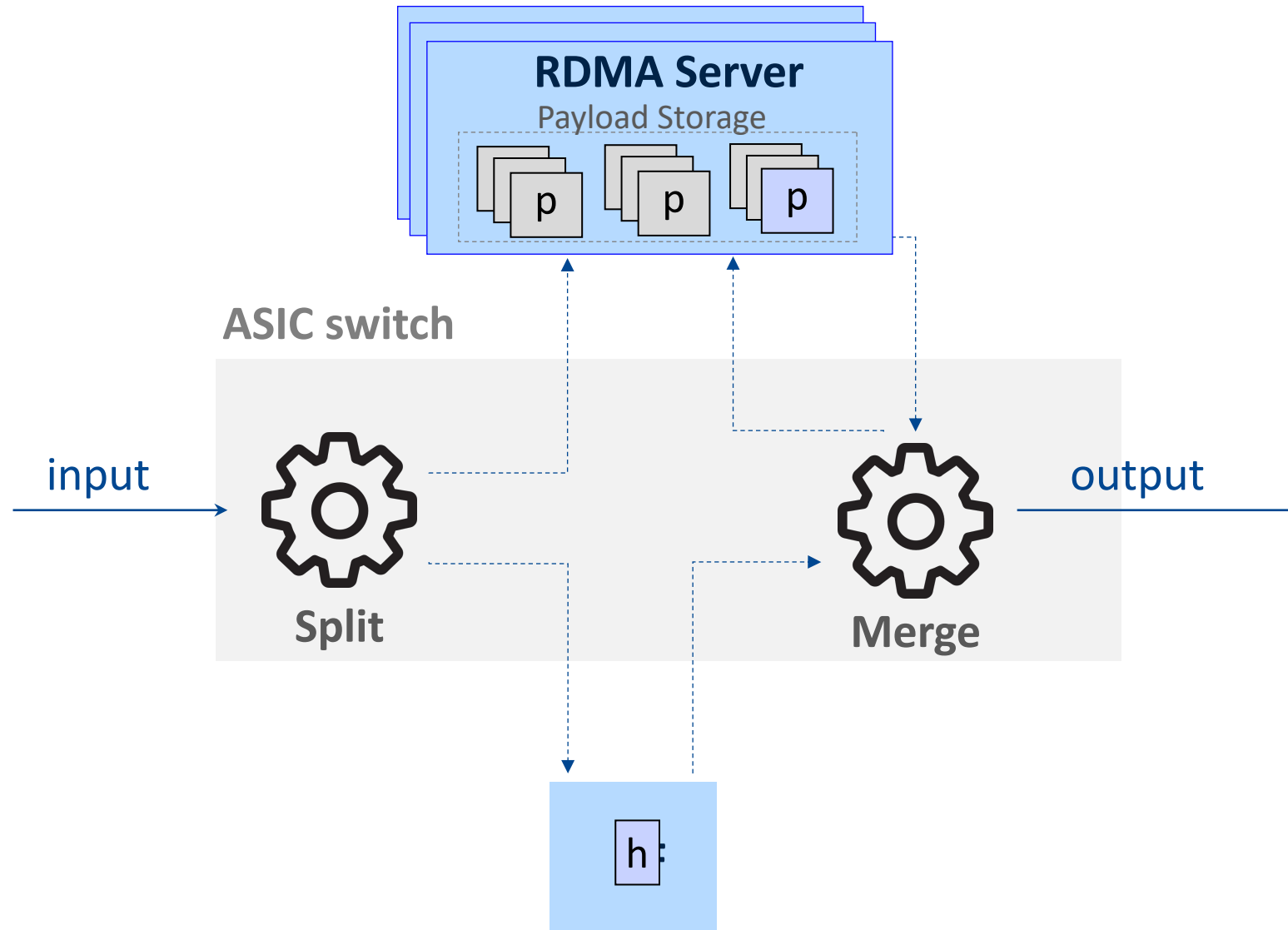
Storing Payloads on Shared Resources (Ribosome)



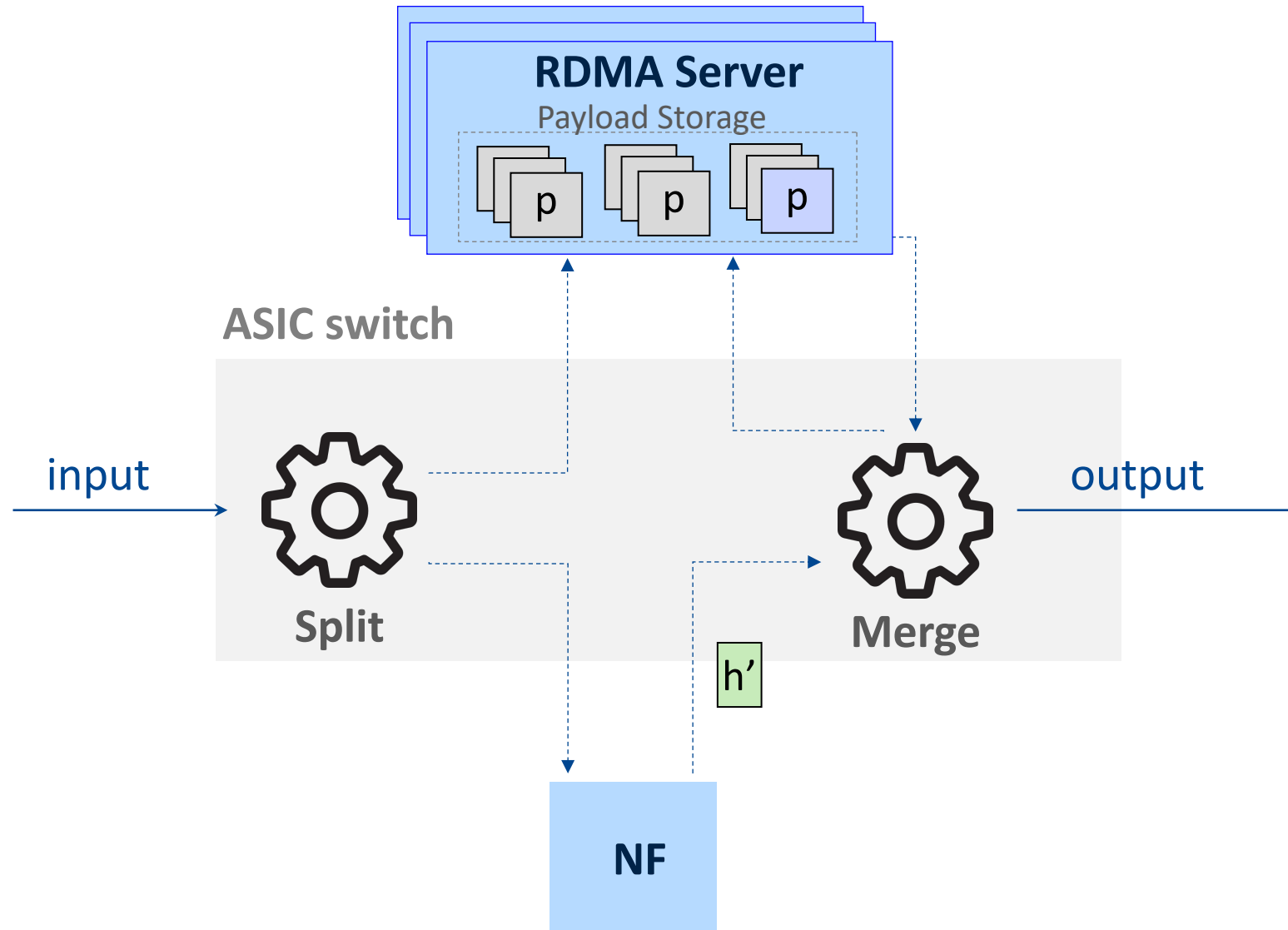
Storing Payloads on Shared Resources (Ribosome)



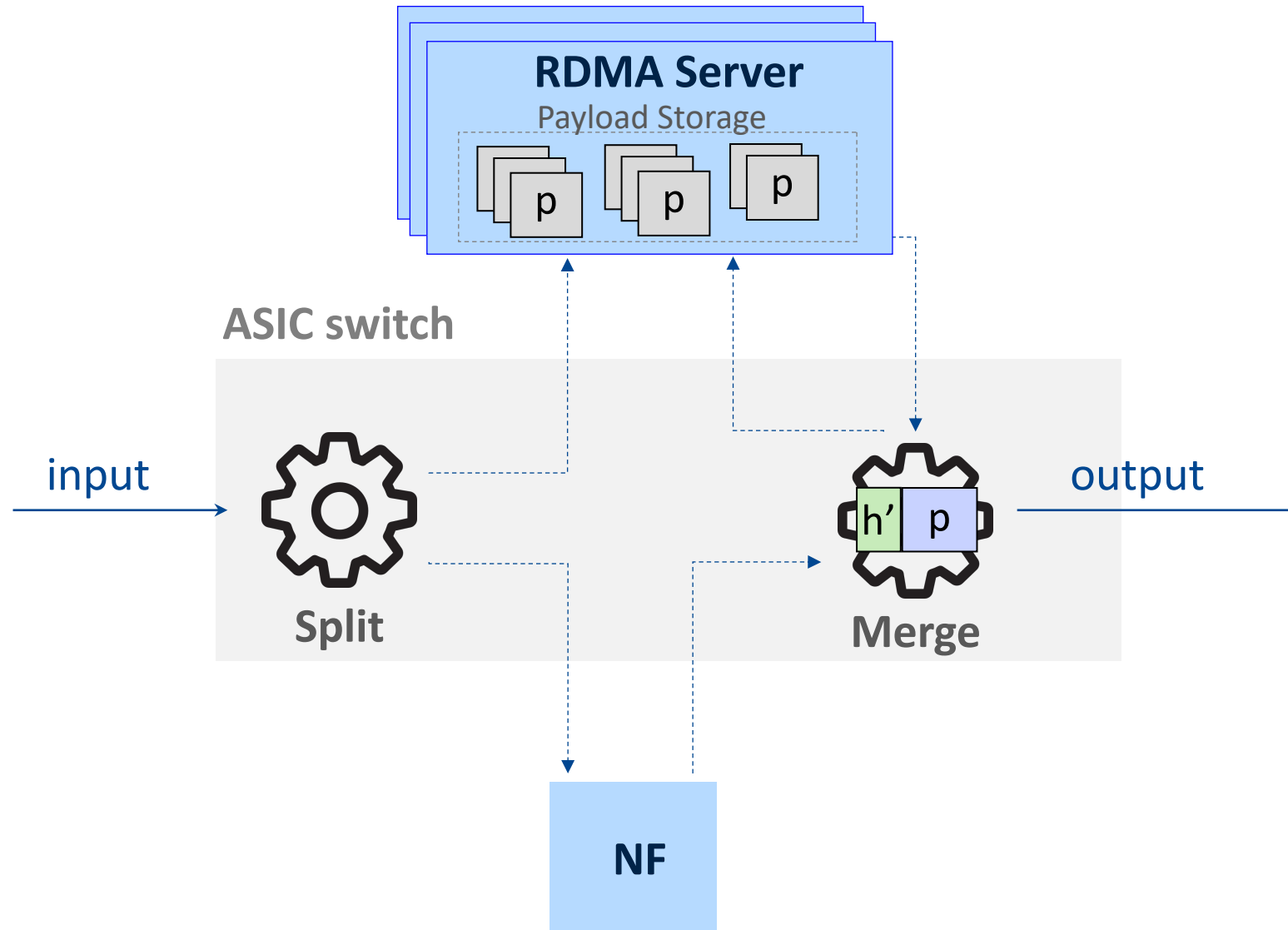
Storing Payloads on Shared Resources (Ribosome)



Storing Payloads on Shared Resources (Ribosome)



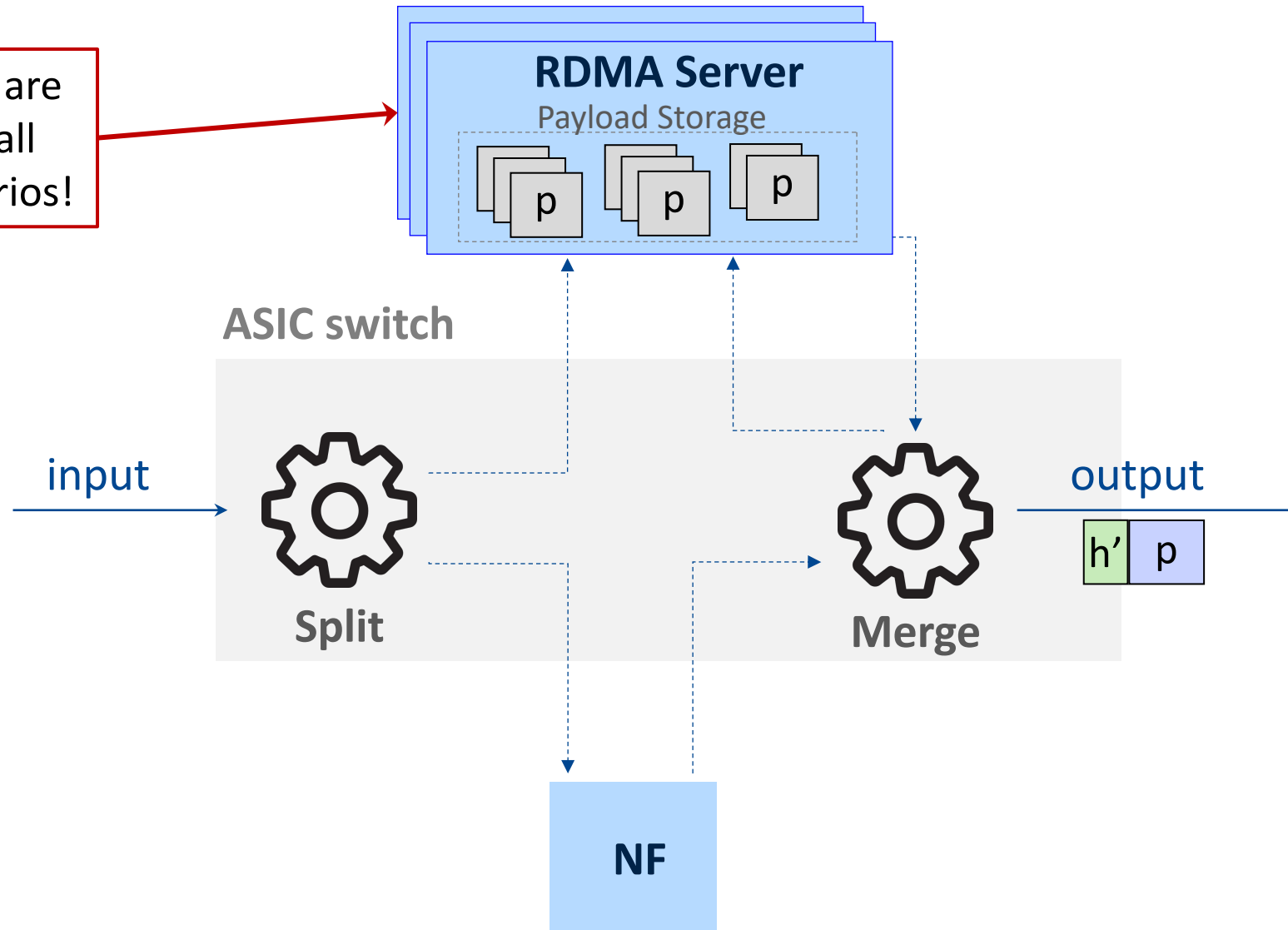
Storing Payloads on Shared Resources (Ribosome)



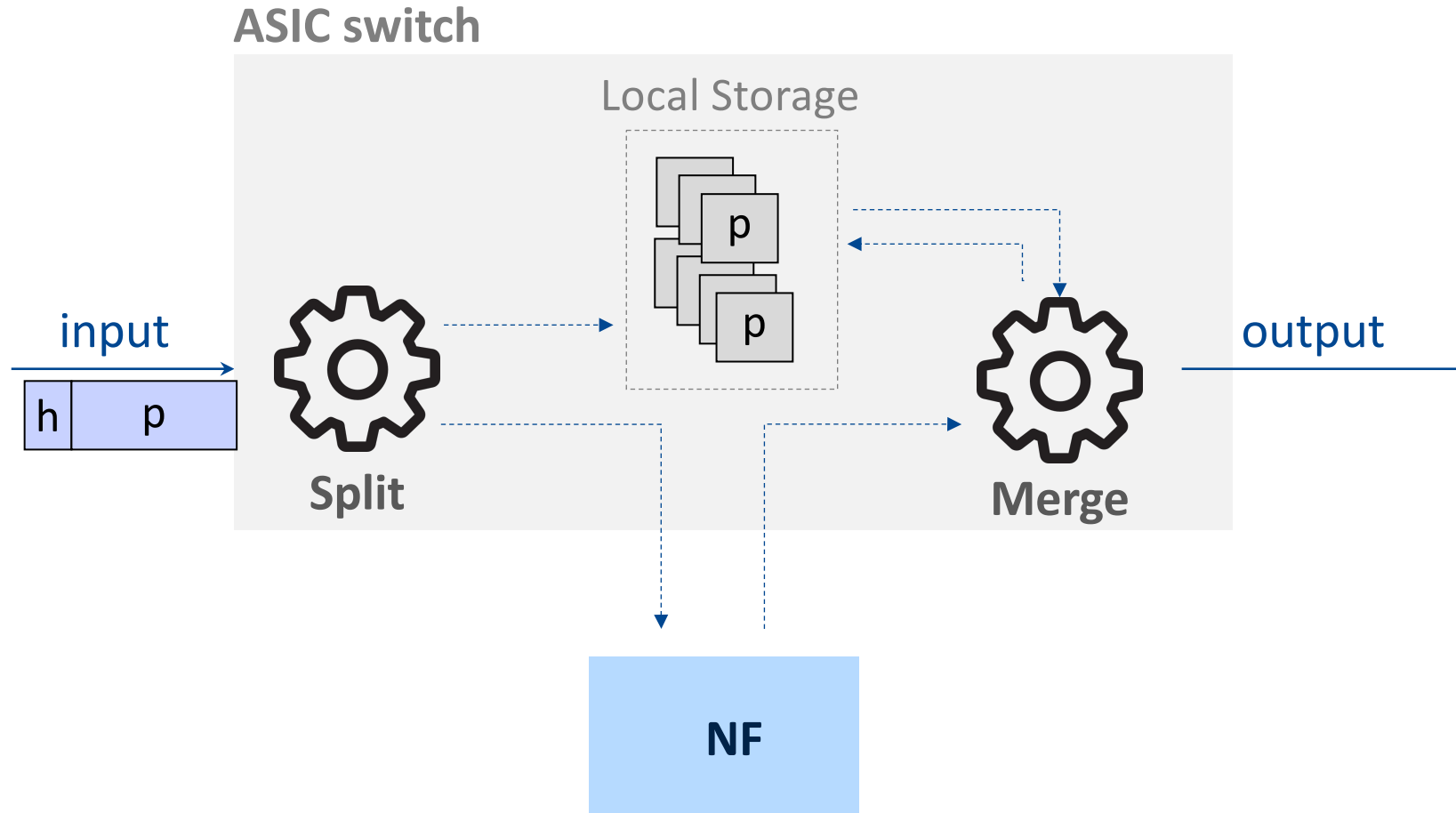
Storing Payloads on Shared Resources (Ribosome)



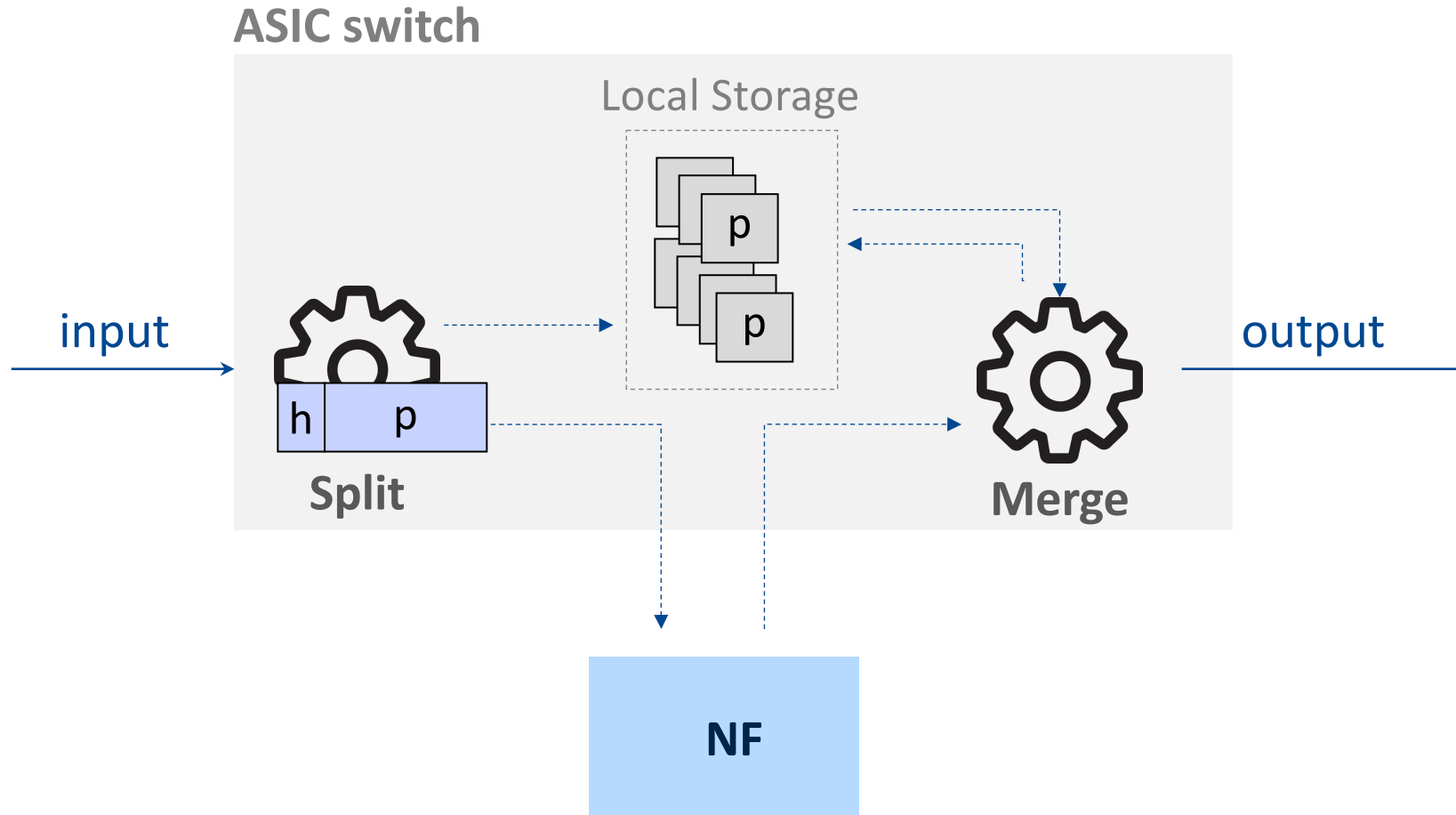
Shared resources are not available in all operational scenarios!



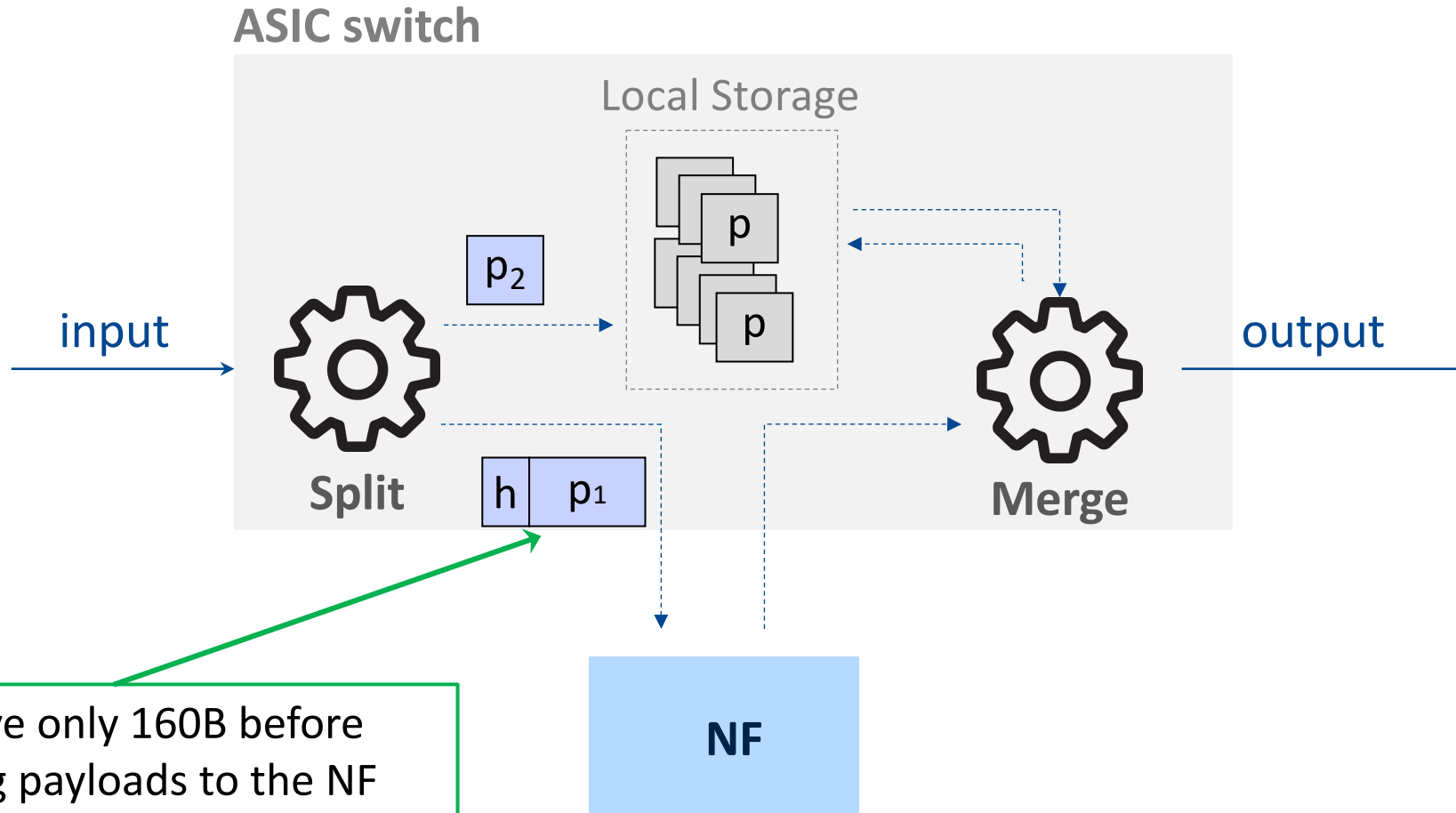
Storing a Few Bytes in ASIC SRAM (PayloadPark)



Storing a Few Bytes in ASIC SRAM (PayloadPark)

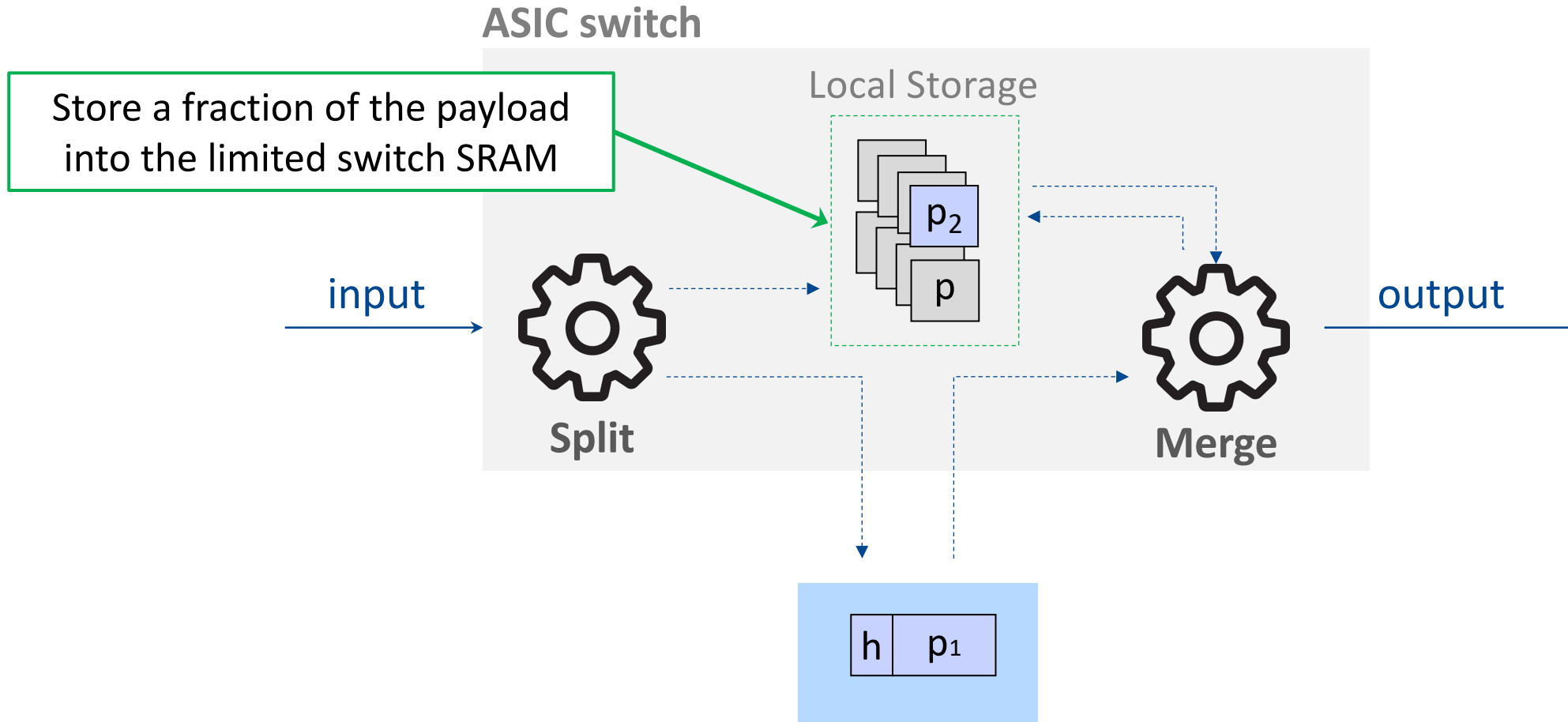


Storing a Few Bytes in ASIC SRAM (PayloadPark)

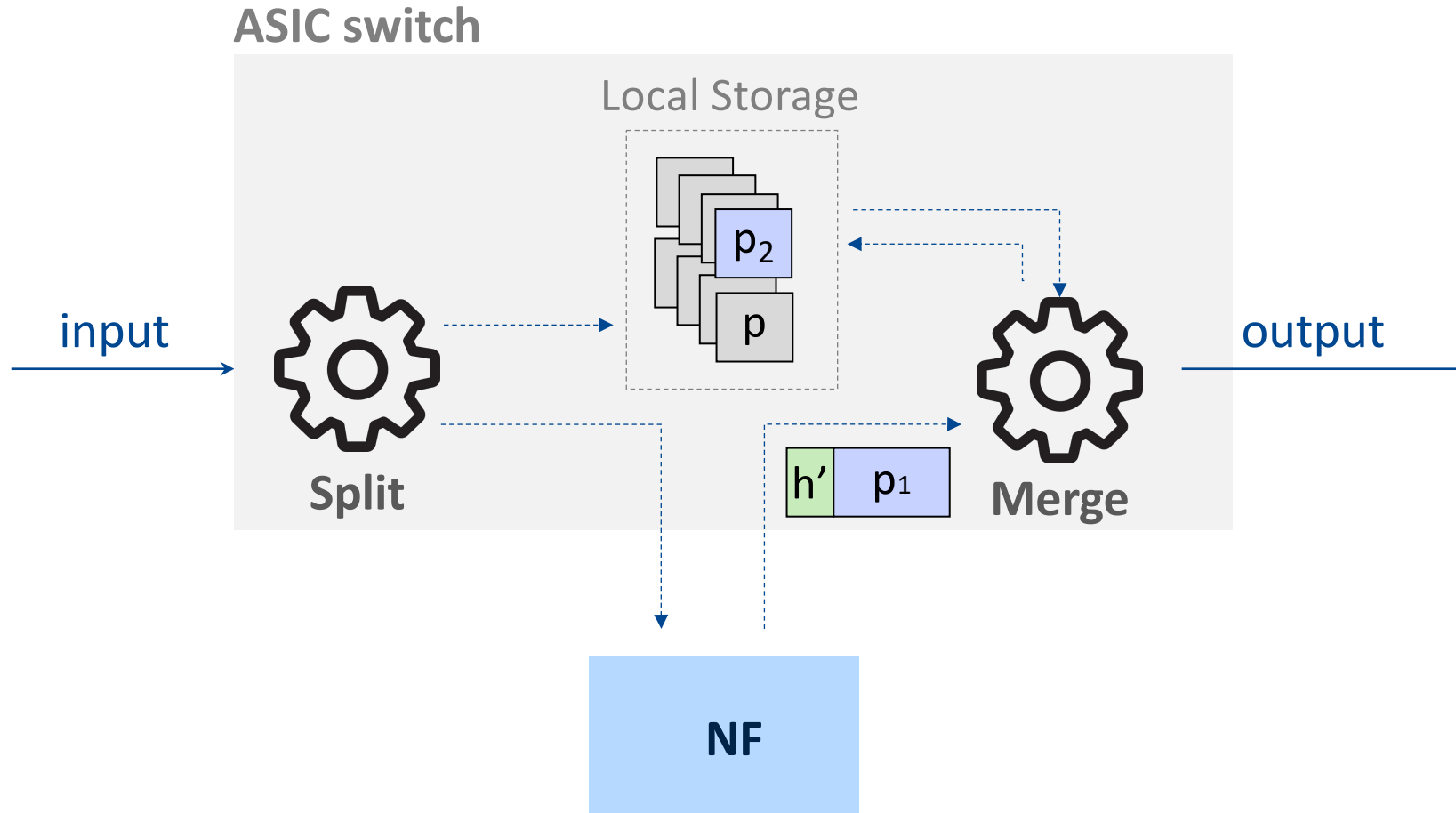


Remove only 160B before sending payloads to the NF

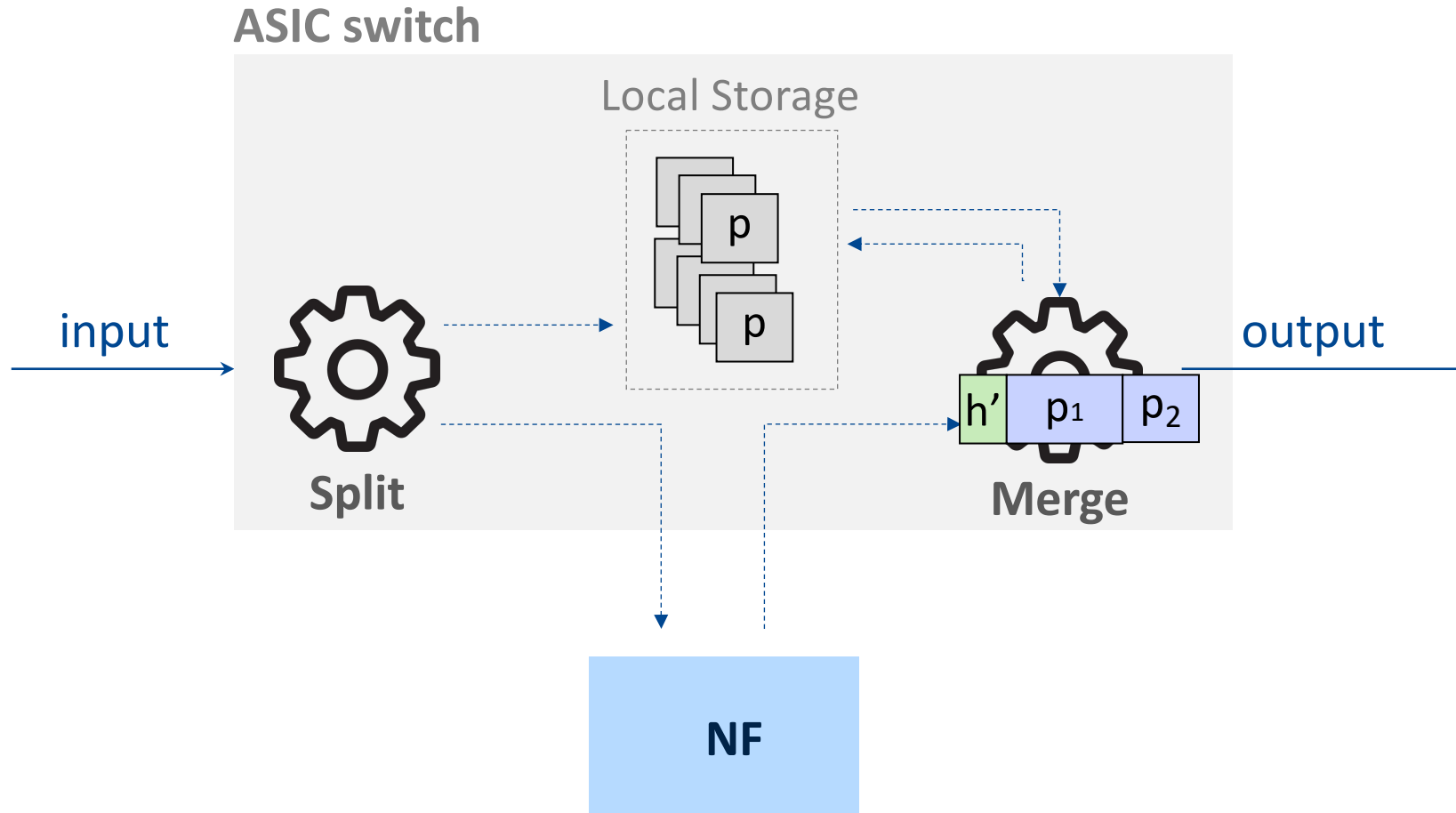
Storing a Few Bytes in ASIC SRAM (PayloadPark)



Storing a Few Bytes in ASIC SRAM (PayloadPark)



Storing a Few Bytes in ASIC SRAM (PayloadPark)

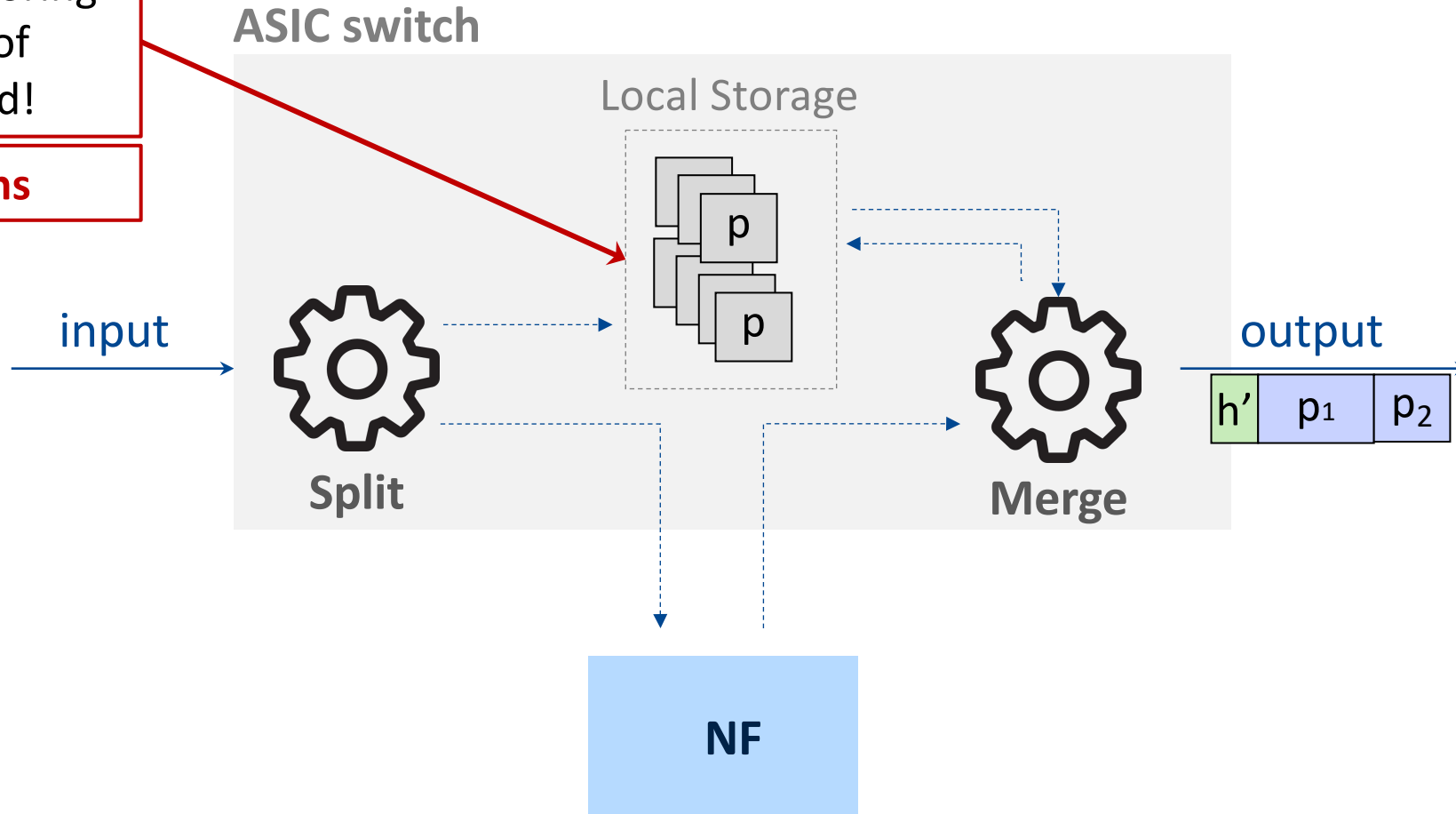


Storing a Few Bytes in ASIC SRAM (PayloadPark)



SRAM allows storing only a part of each payload!

Limited gains



Can we **store the entire payload** within the switch without requiring external devices?

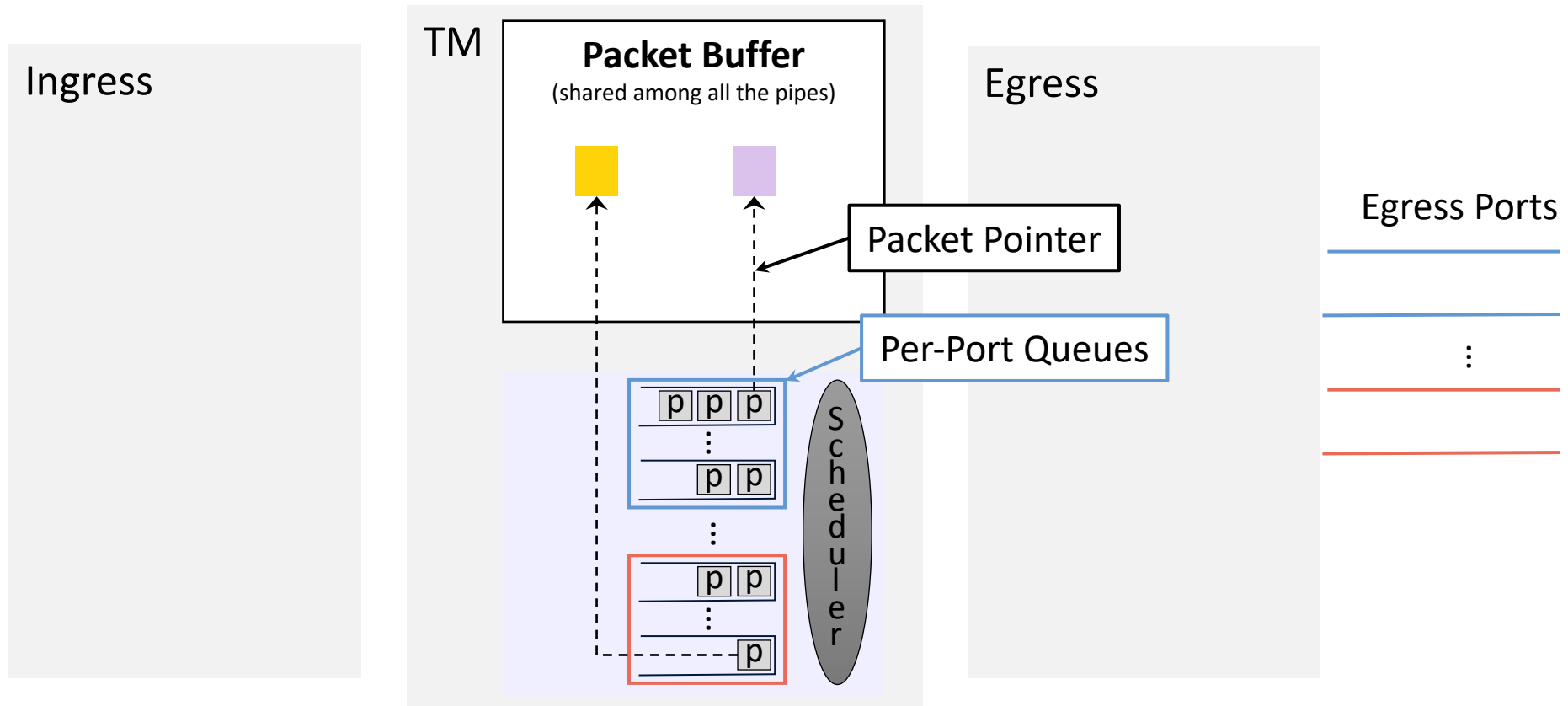


Core Idea: A Queue-Based Packet Storage

Exploit the switch shared buffer to store payloads

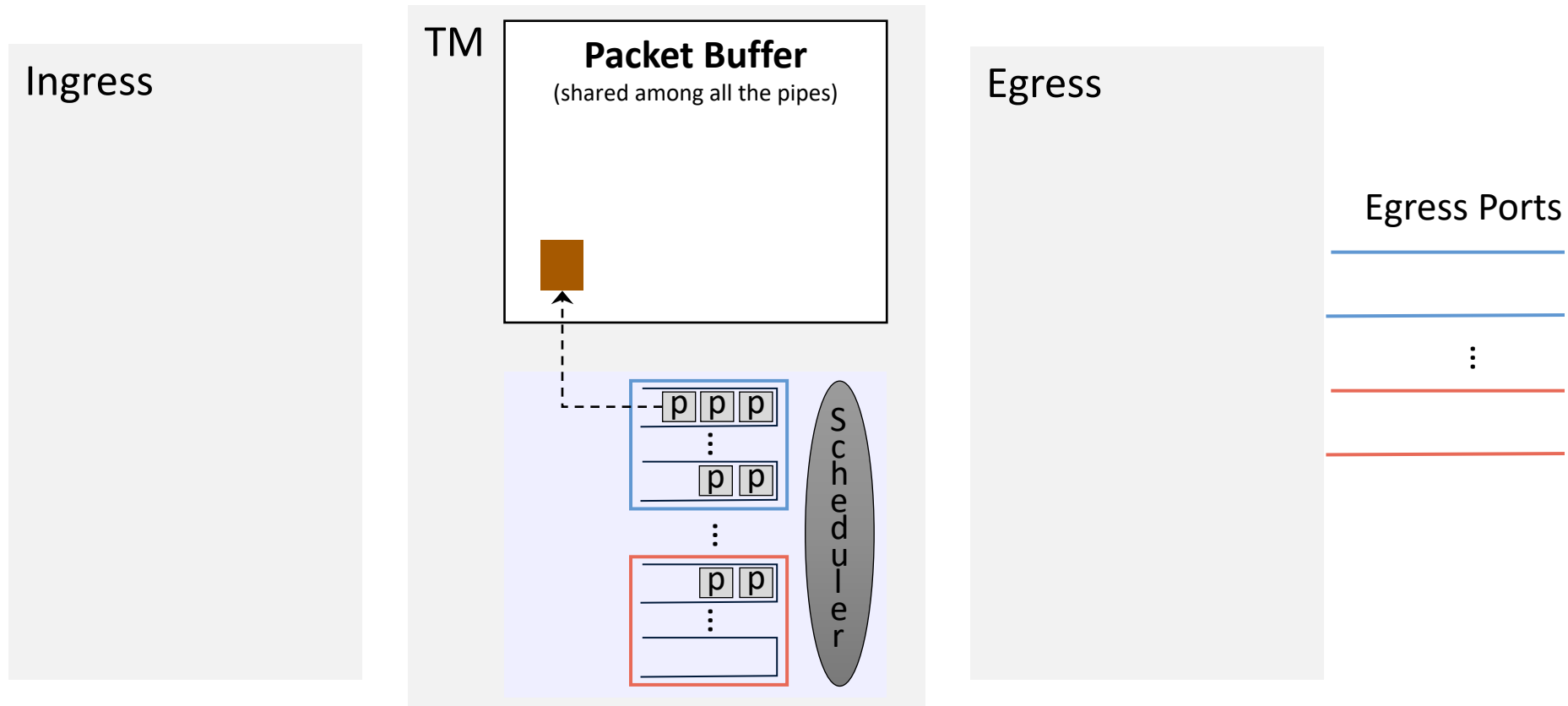
Core Idea: A Queue-Based Packet Storage

Exploit the switch shared buffer to store payloads



Core Idea: A Queue-Based Packet Storage

Exploit the switch shared buffer to store payloads



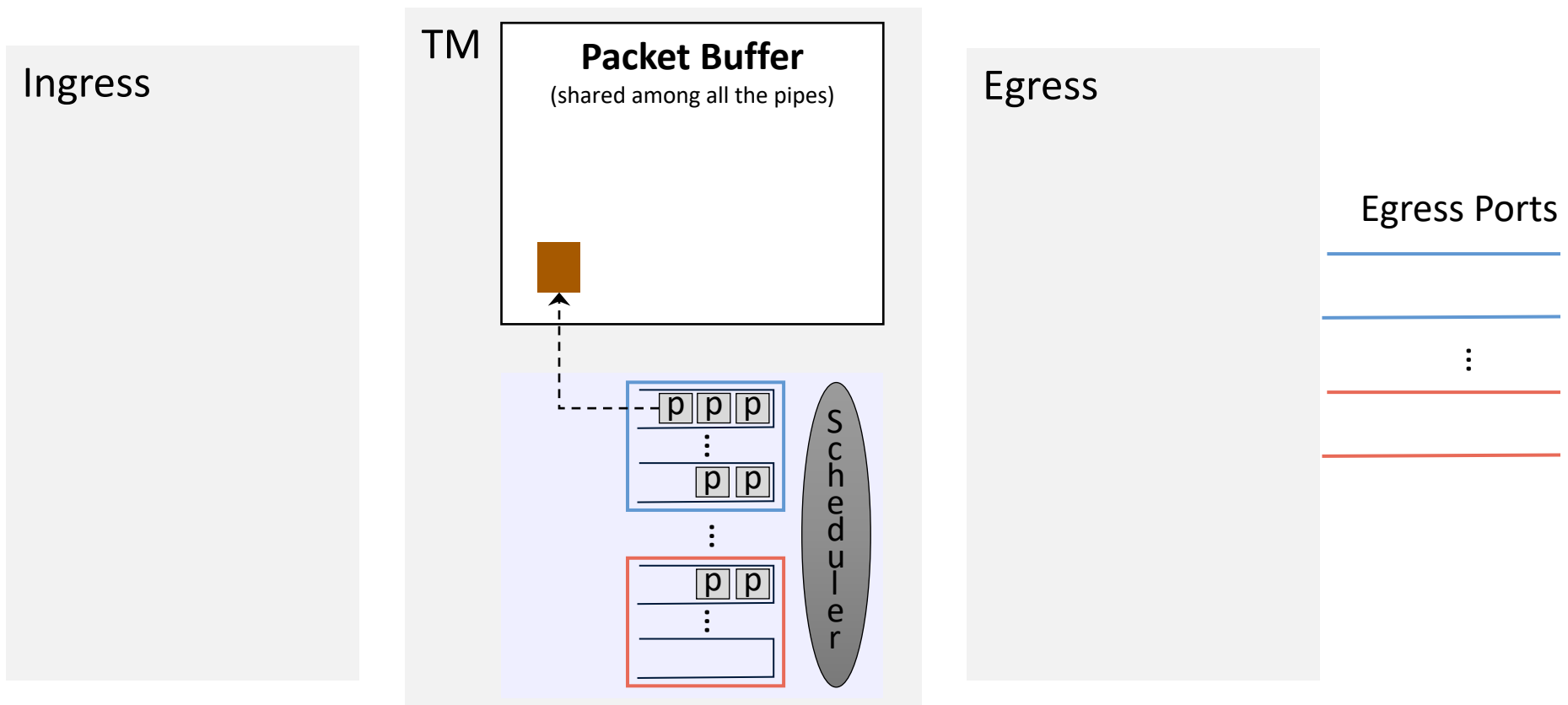
Core Idea: A Queue-Based Packet Storage

Exploit the switch shared buffer to store payloads

✓ Large shared buffer

✗ No programmatic access

⊙ How can we access the packet buffer?



How can we access the packet buffer?



Enqueued packets **implicitly** control the buffer memory

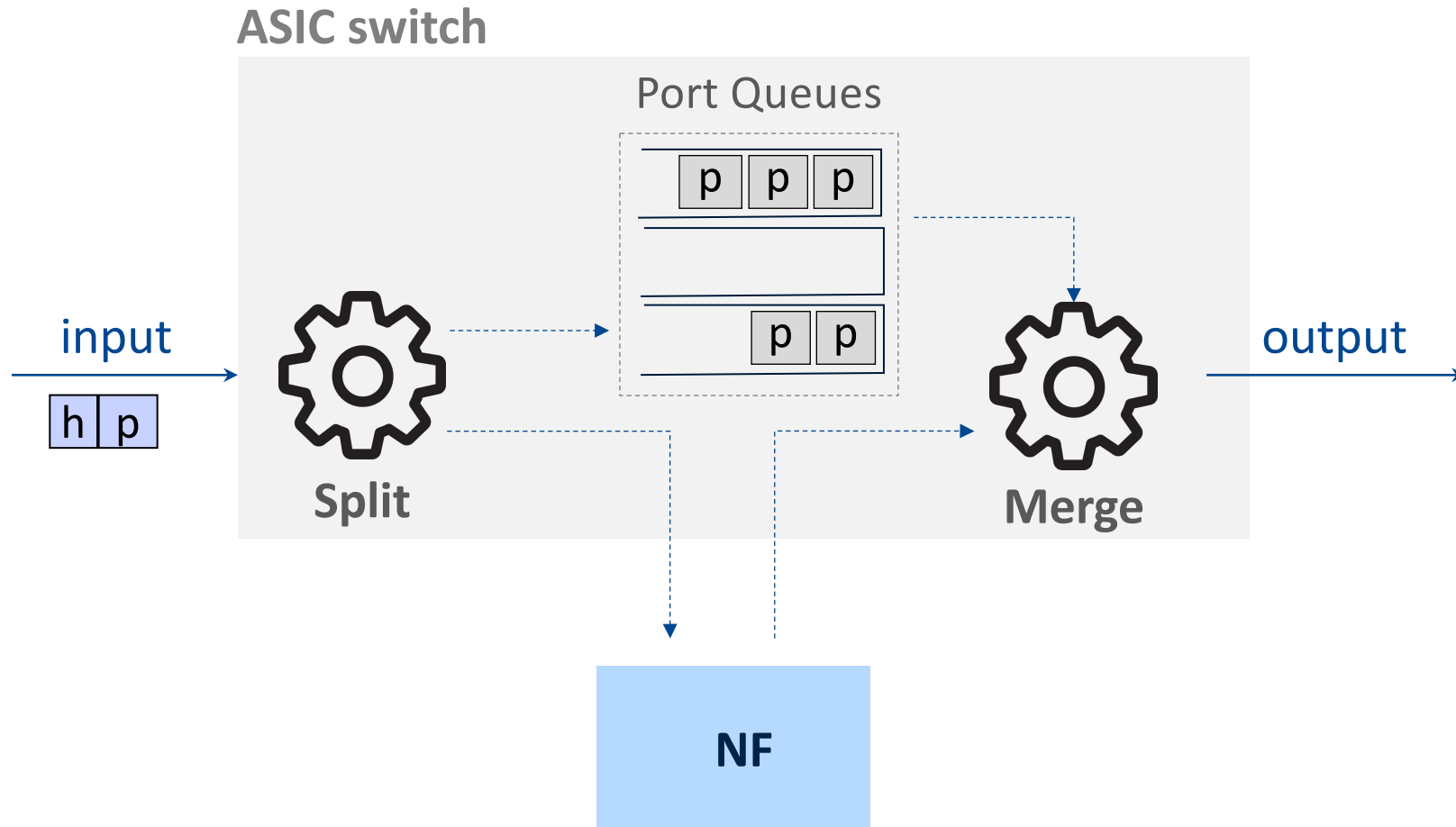


Hold payloads in the port queues until
the processed header from the NF is returned!

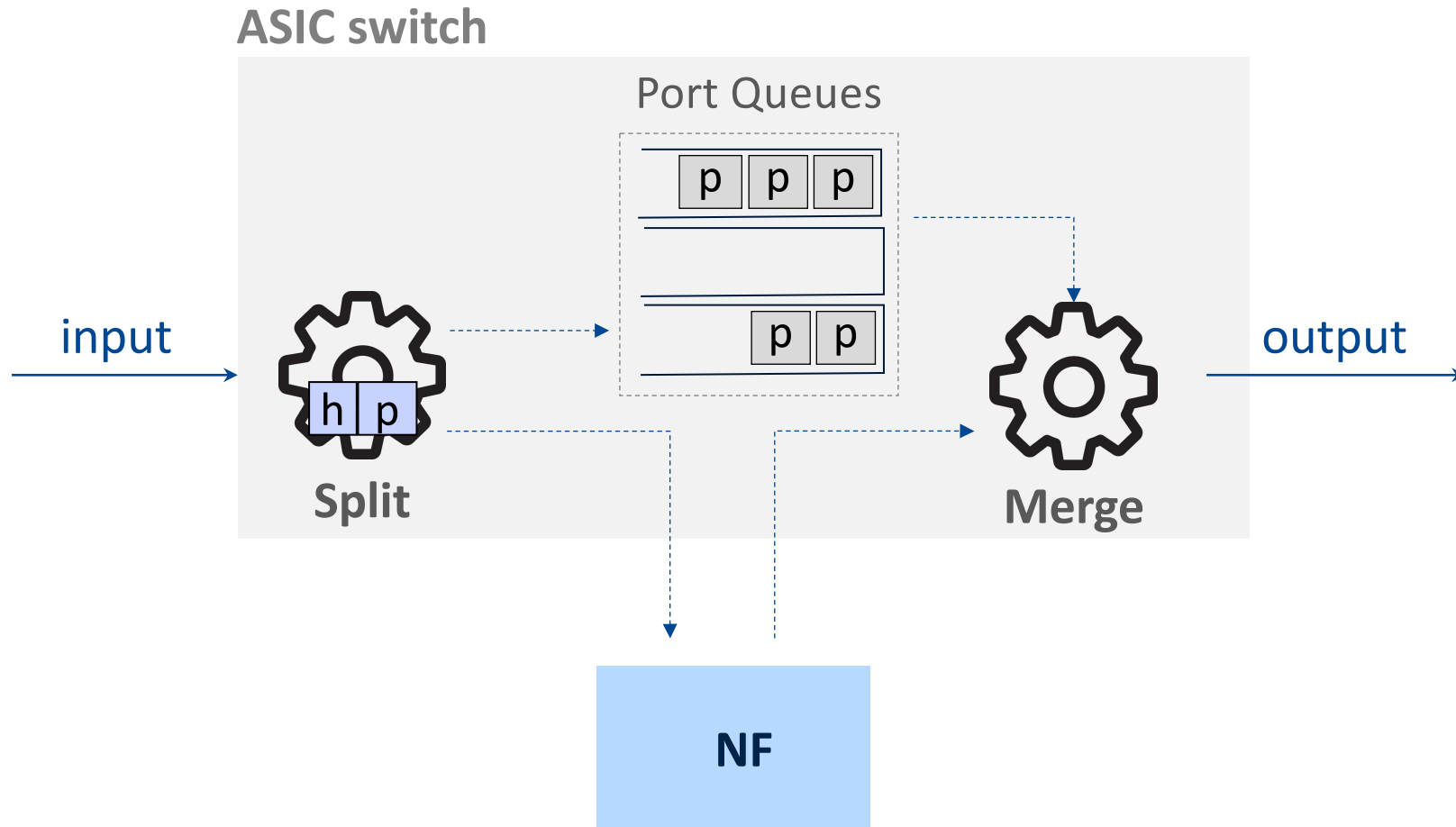


Would it work for real?

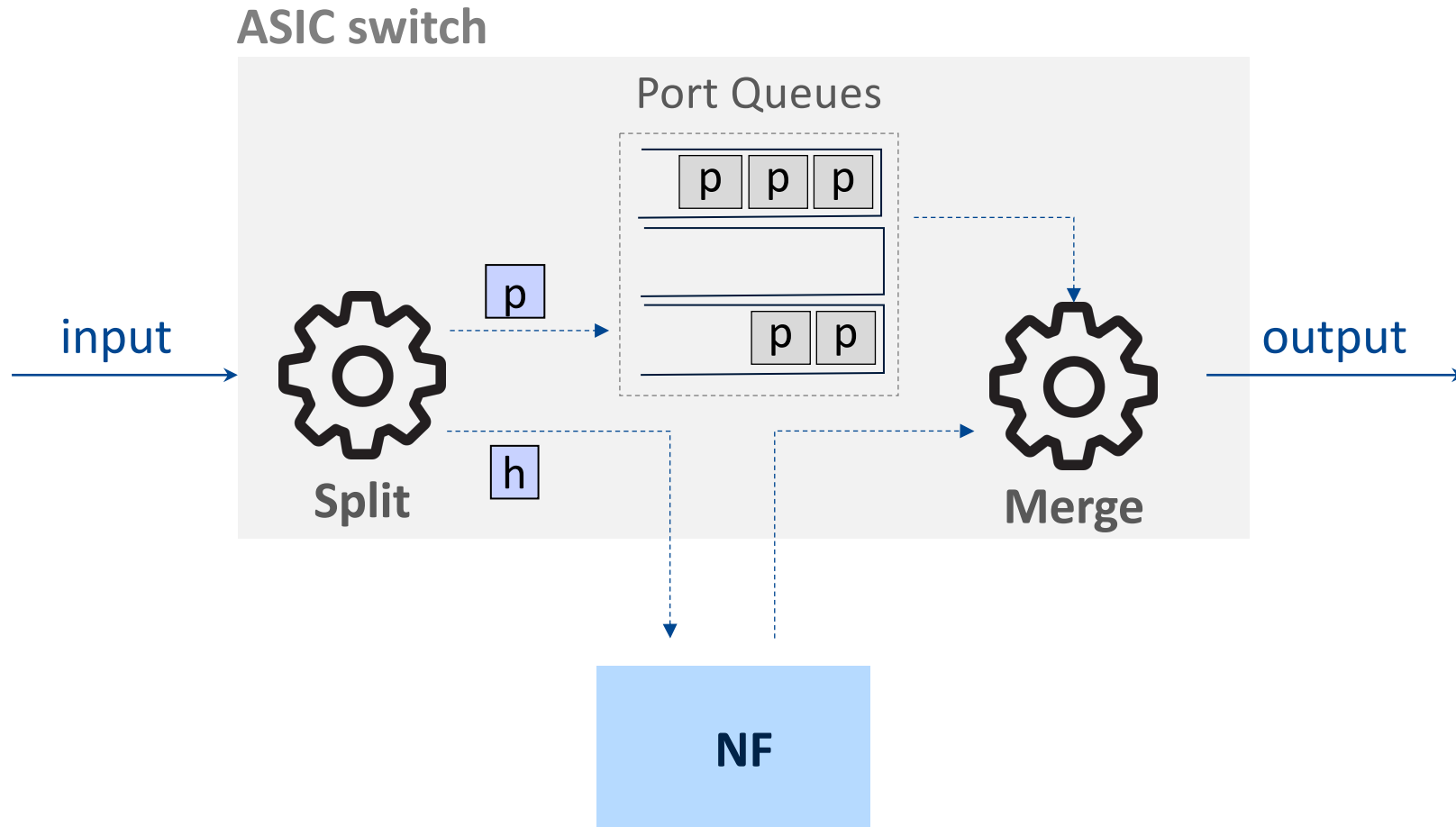
Challenges



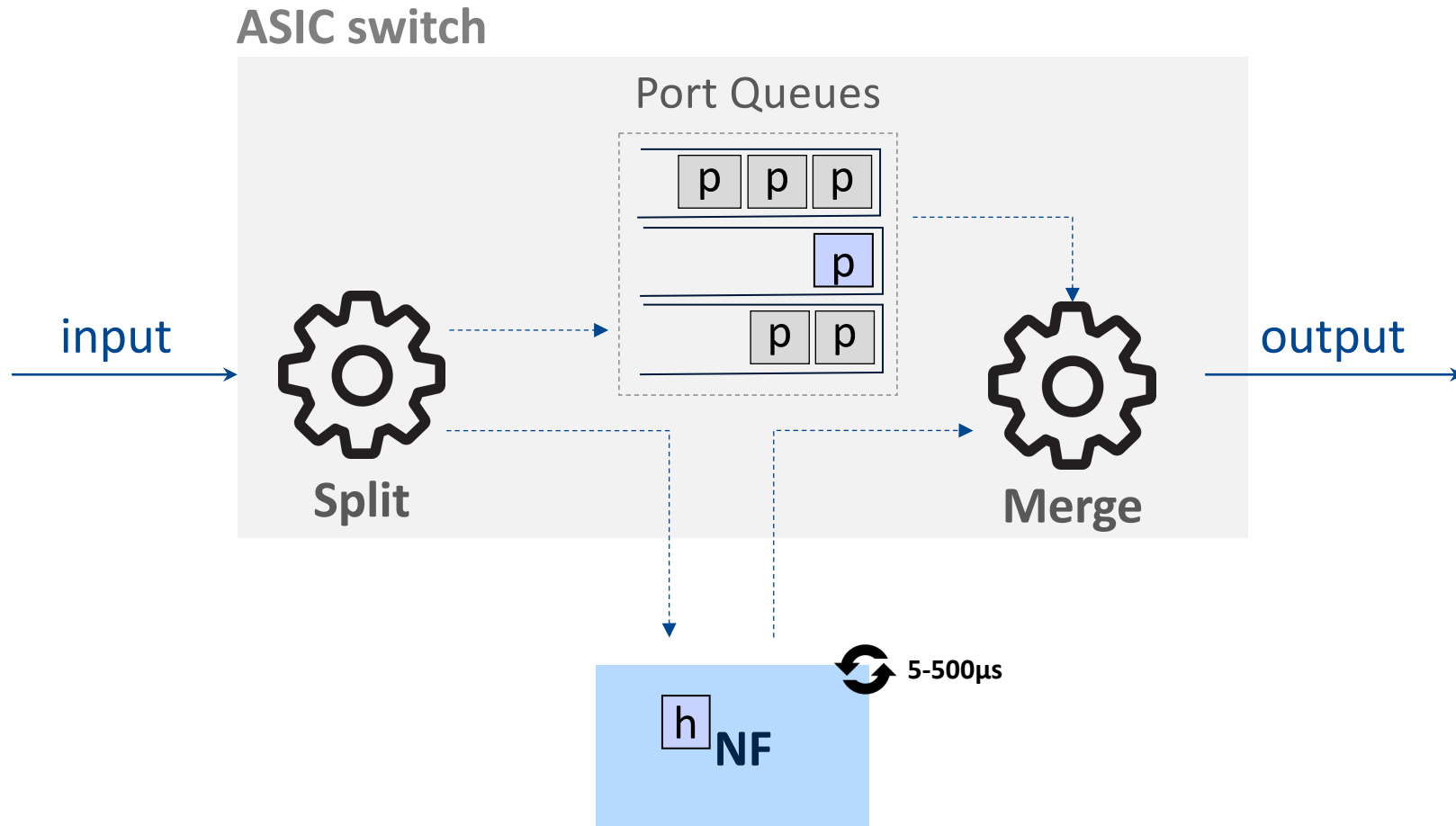
Challenges



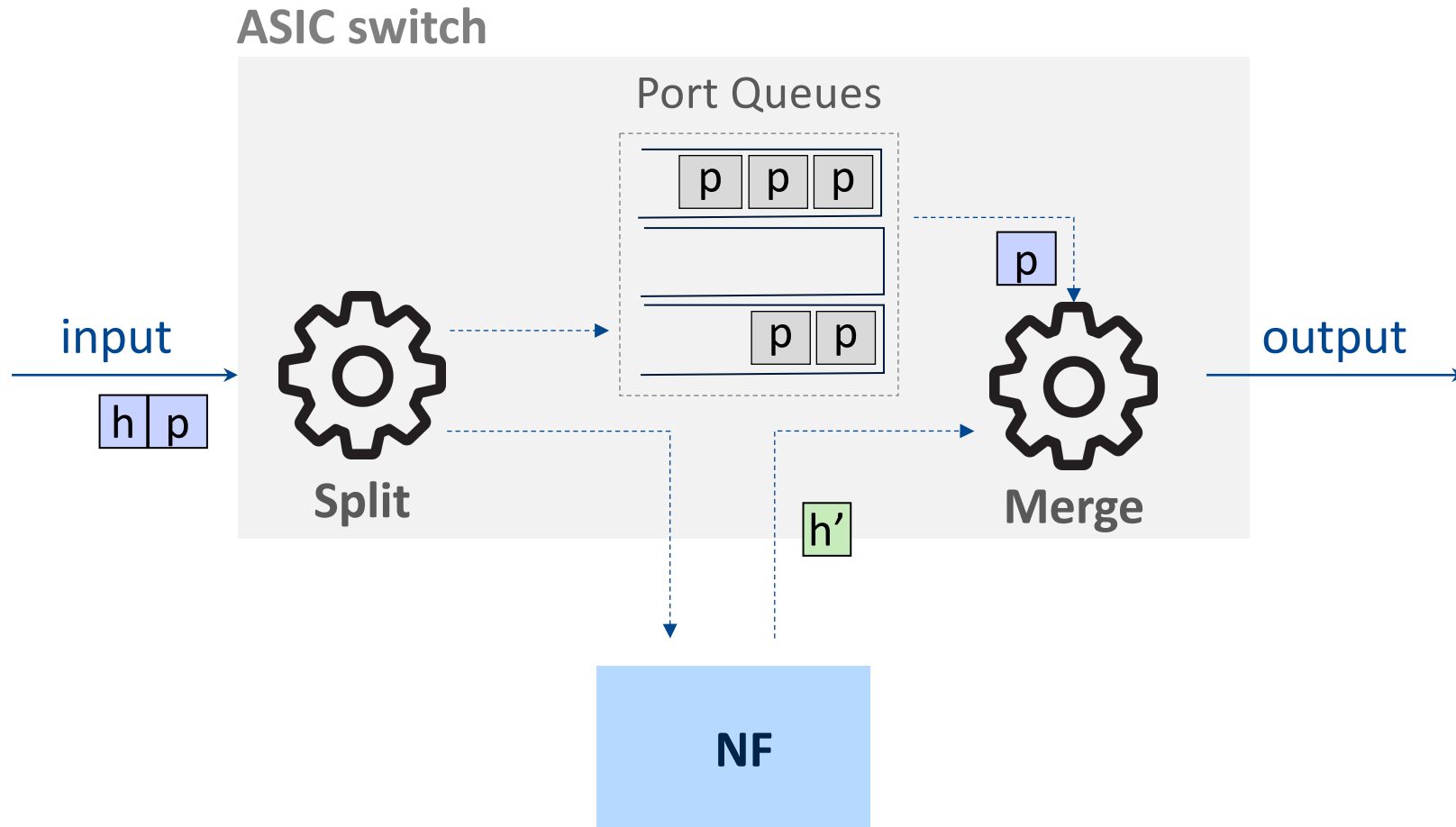
Challenges



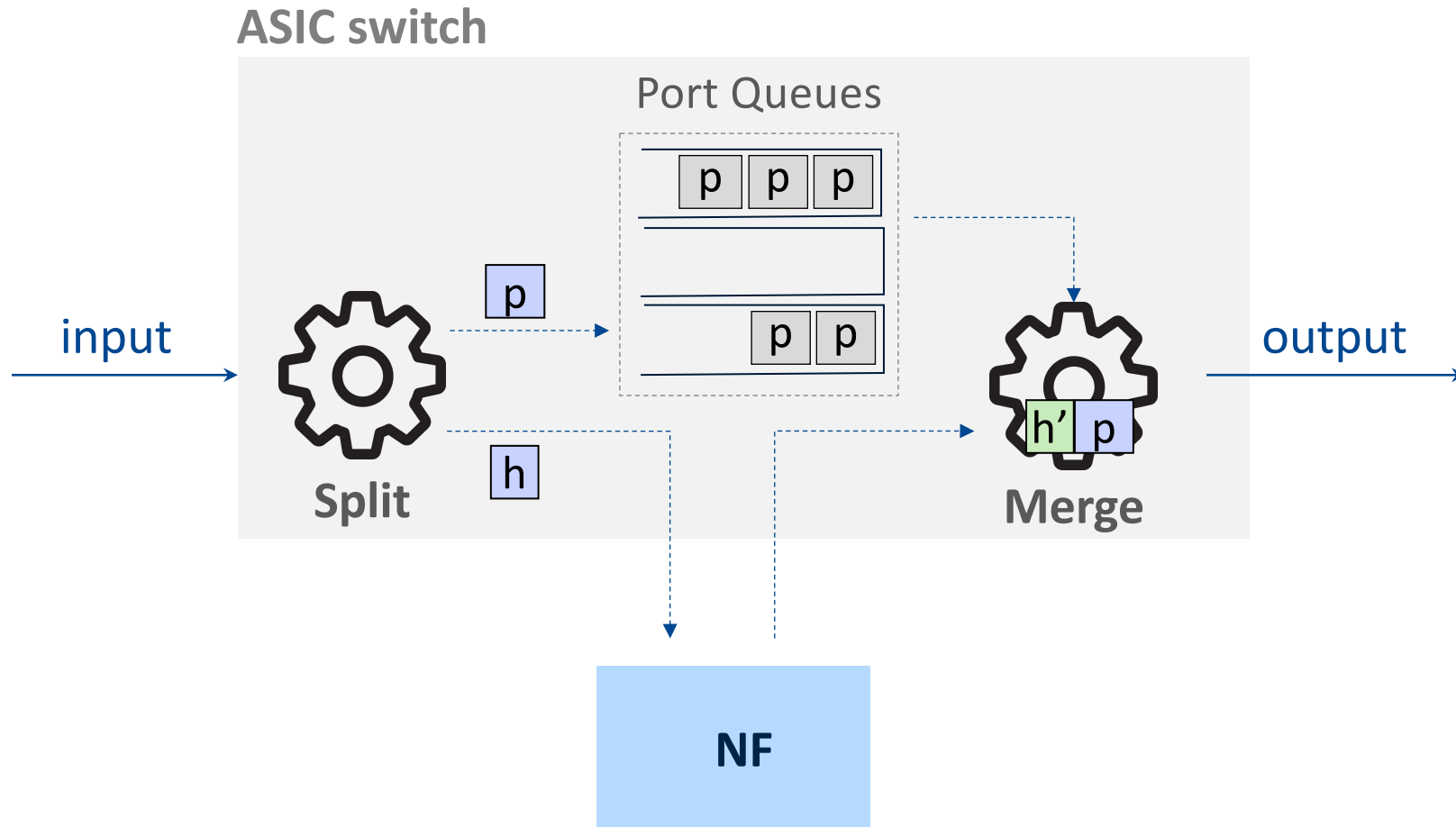
Challenges



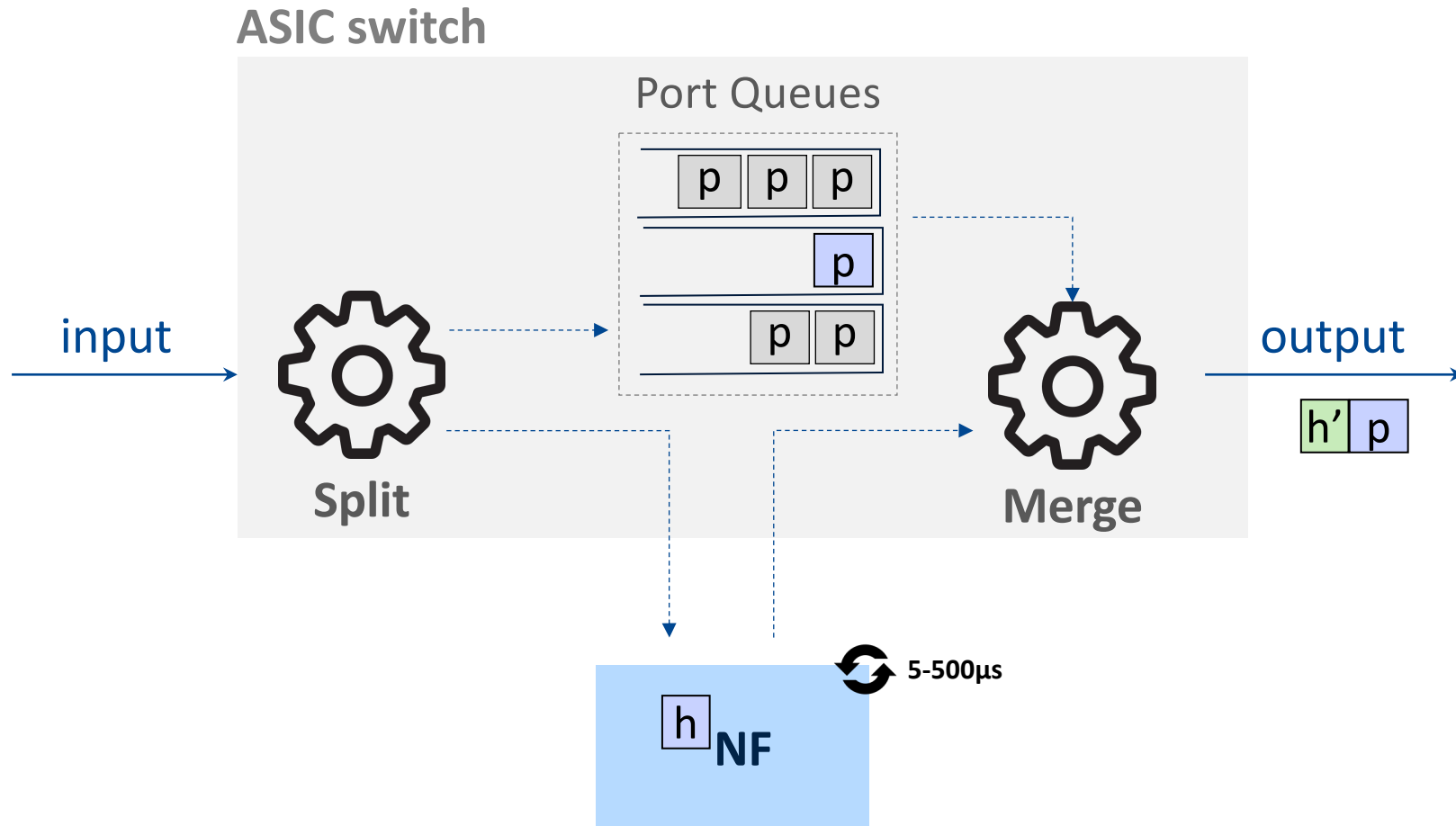
Challenges



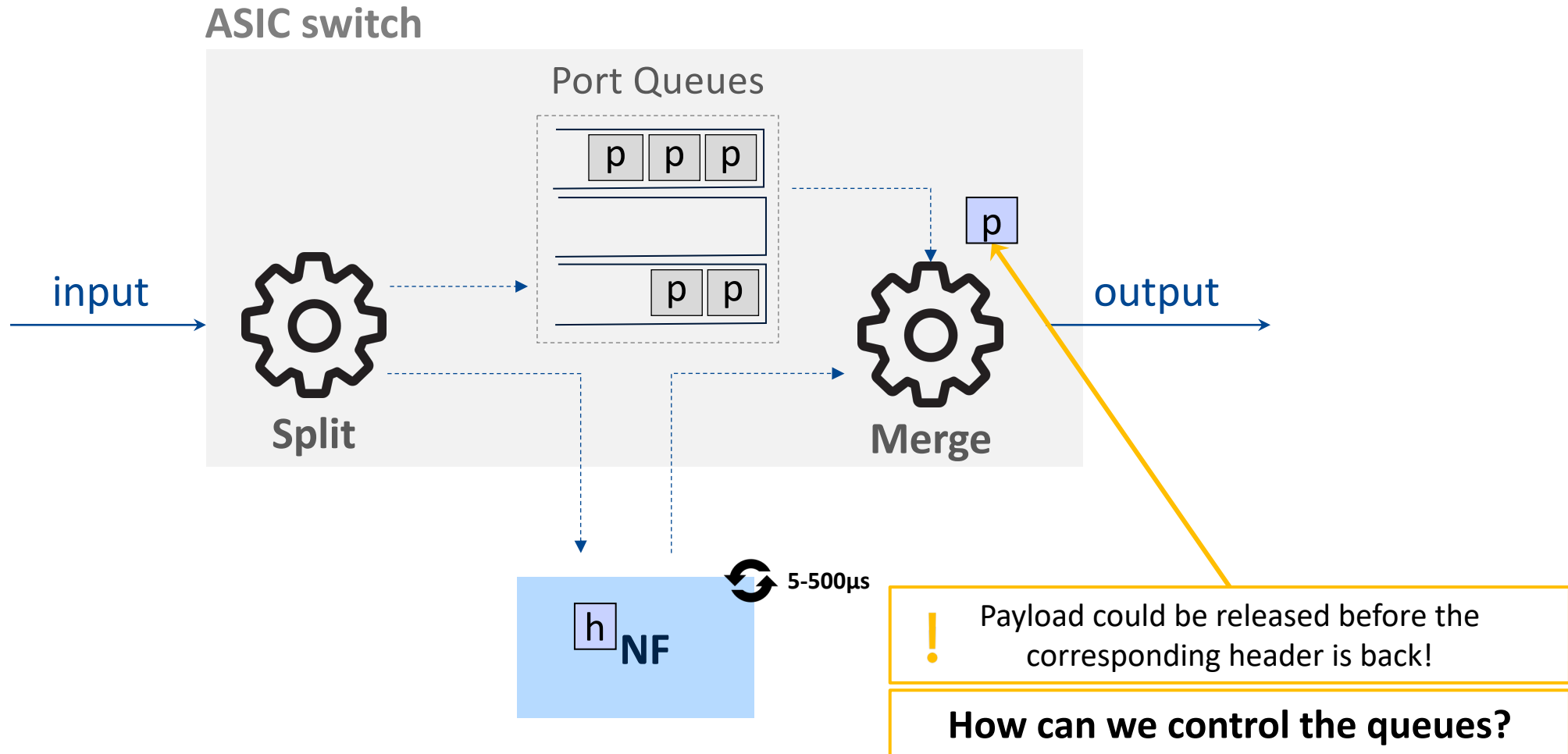
Challenges



Challenges



Challenges



How can we control the queues?



Congest the egress queues to delay payloads release by the time required for NF processing

How can we control the queues?



Congest the egress queues to delay payloads release by the time required for NF processing



Recirculate payloads until headers are ready

How can we control the queues?



Congest the egress queues to delay payloads release by the time required for NF processing



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues

How can we control the queues?



Congest the egress queues to delay payloads release by the time required for NF processing



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues

TurboSwitch



TurboSwitch



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues

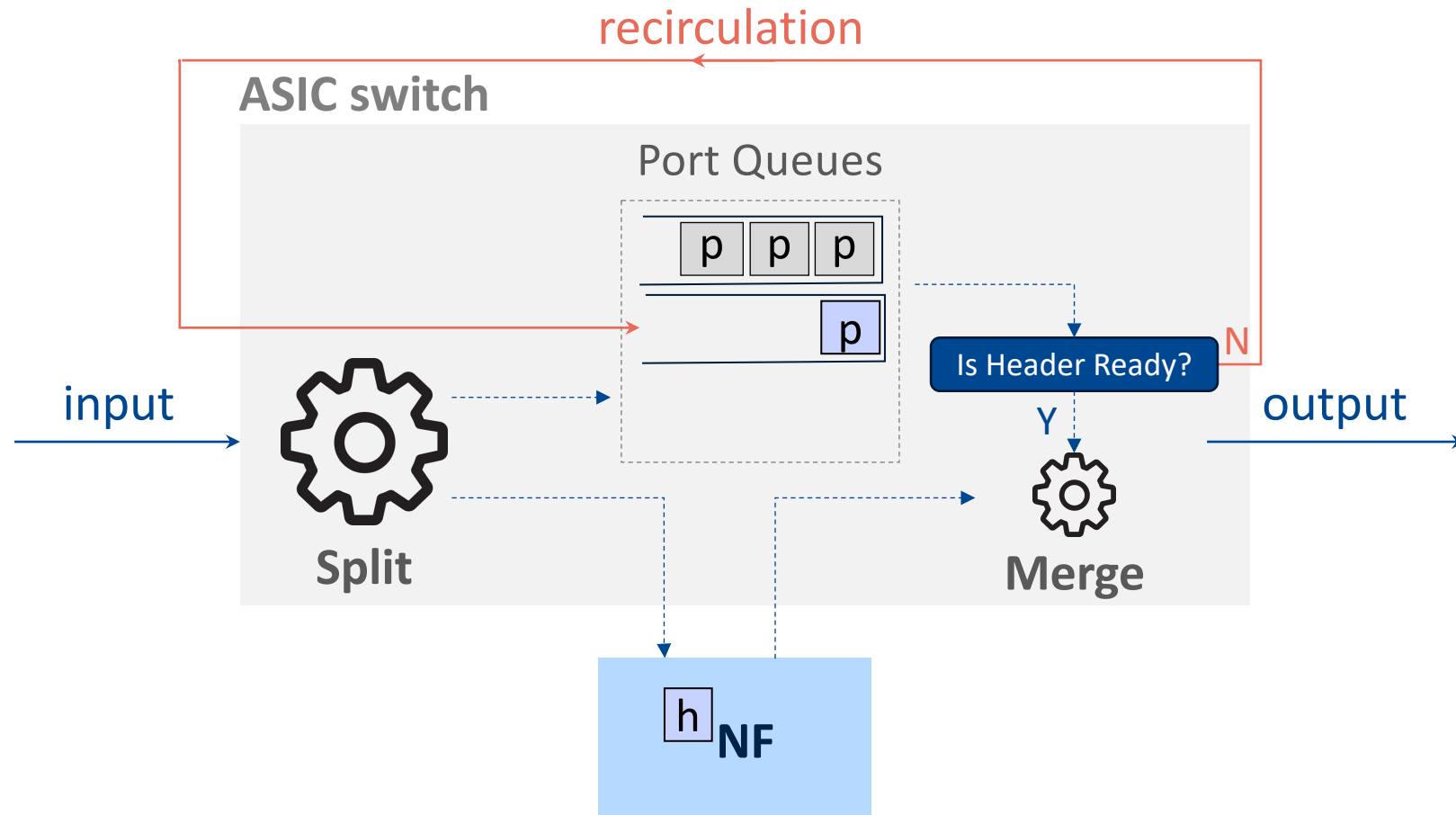
TurboSwitch



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



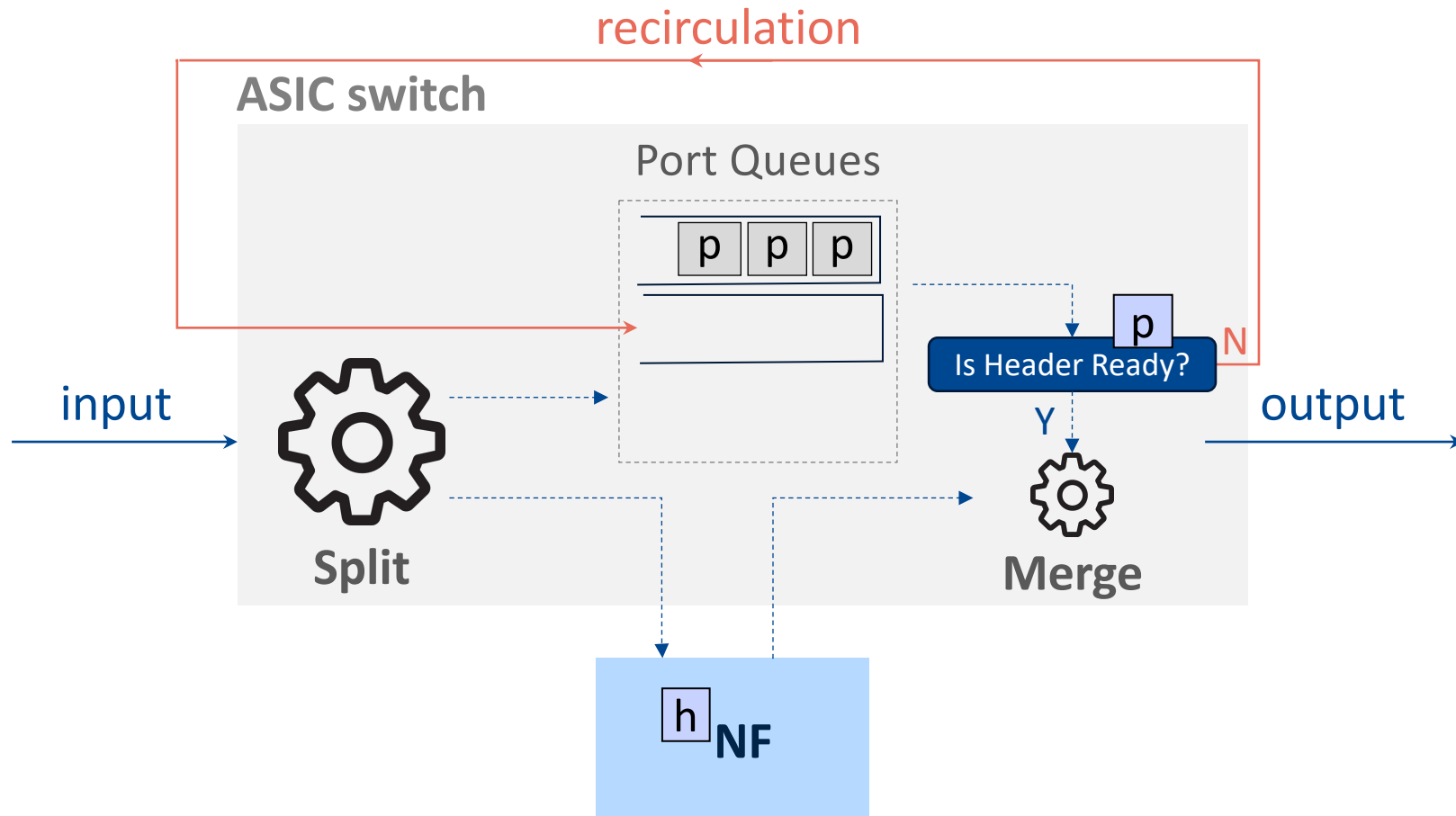
TurboSwitch



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



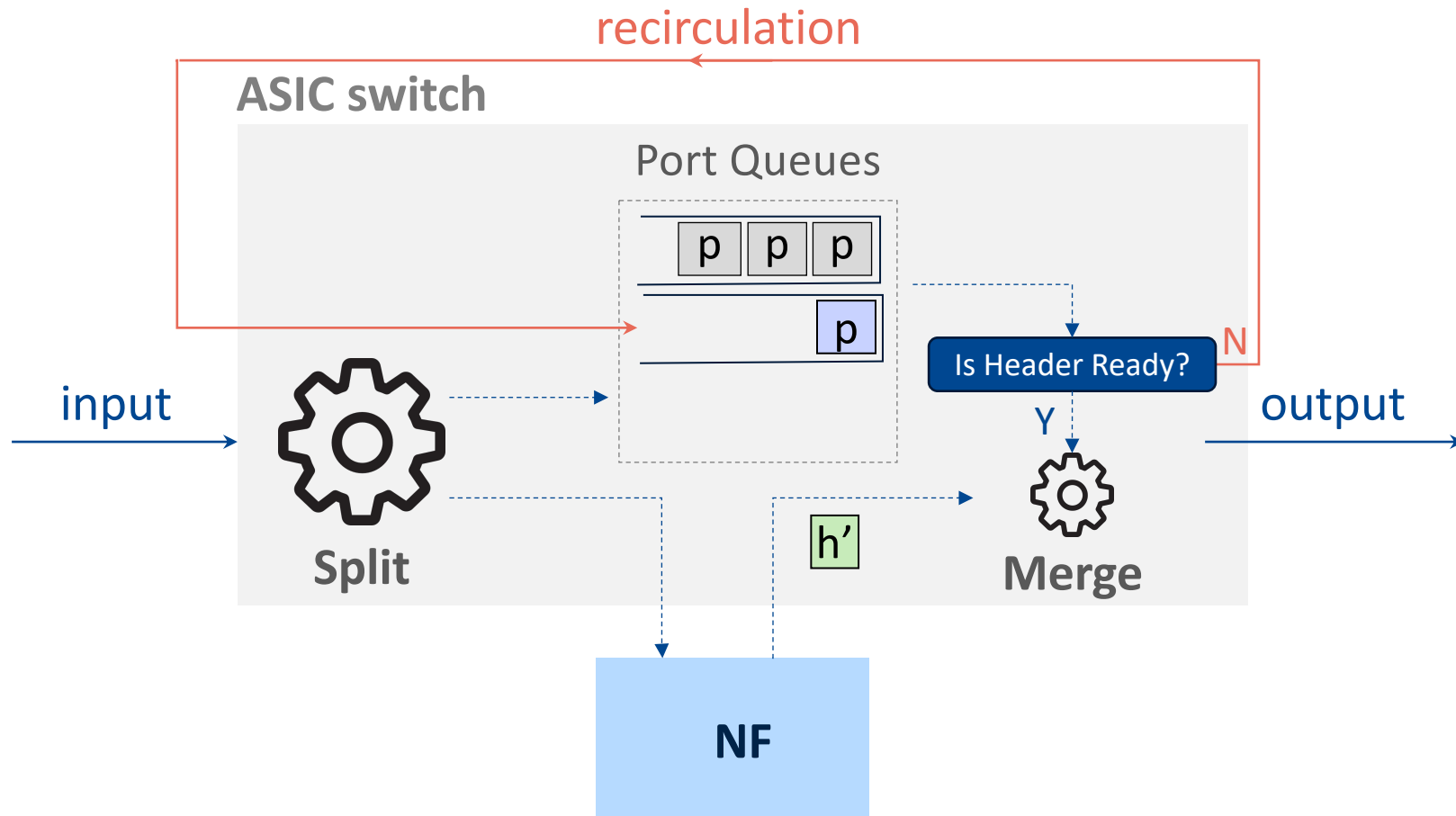
TurboSwitch



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



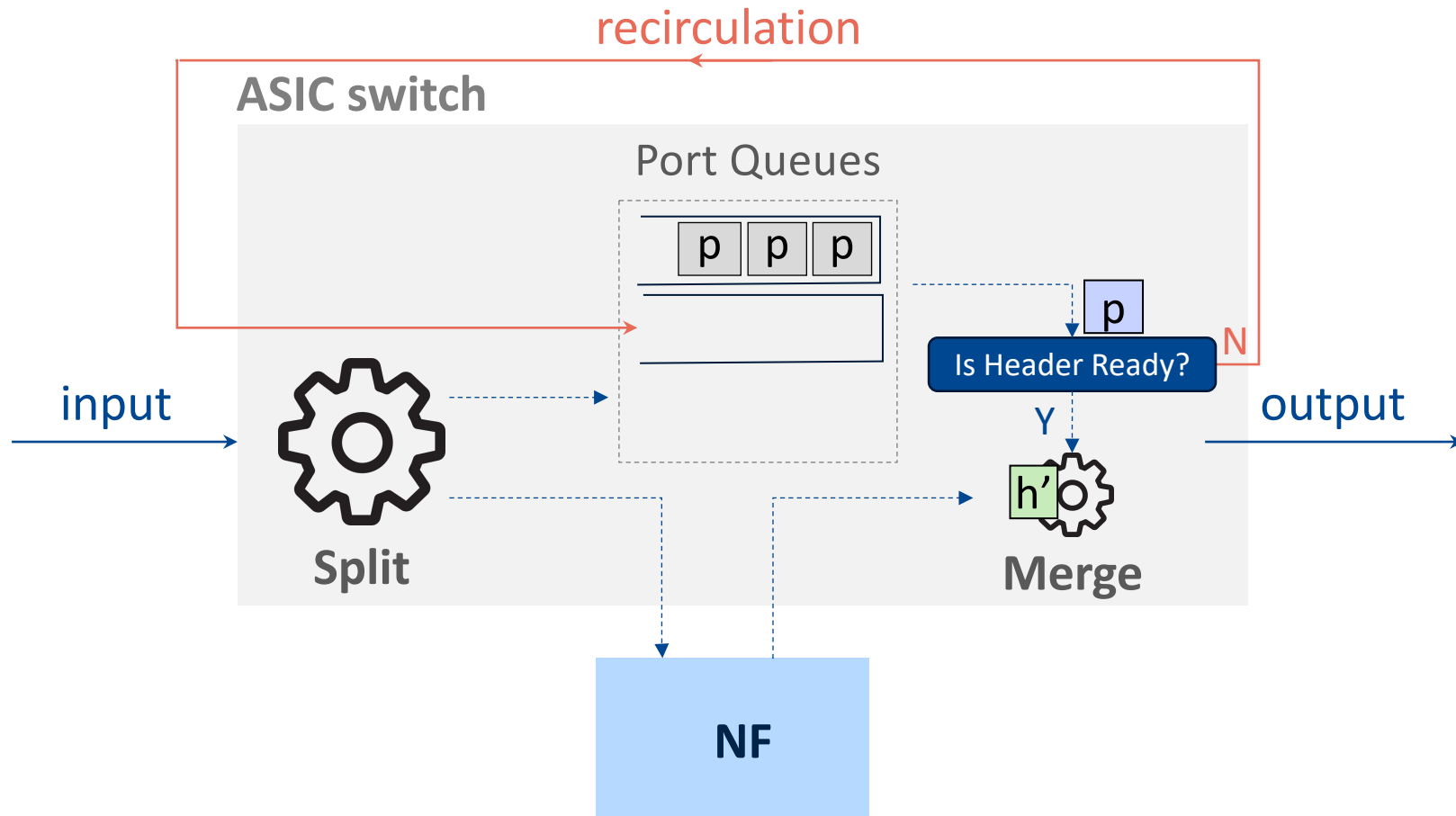
TurboSwitch



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



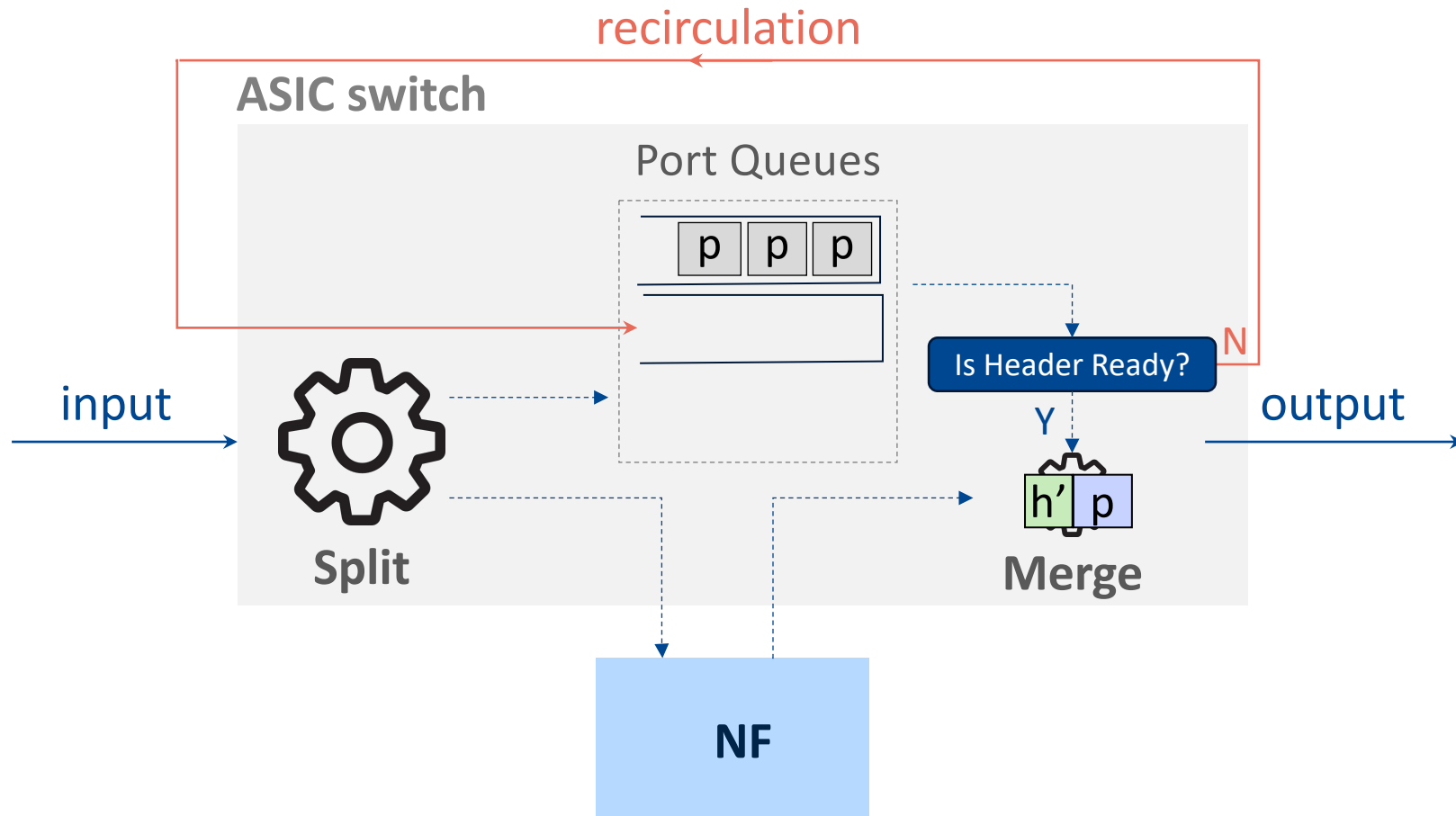
TurboSwitch



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



TurboSwitch



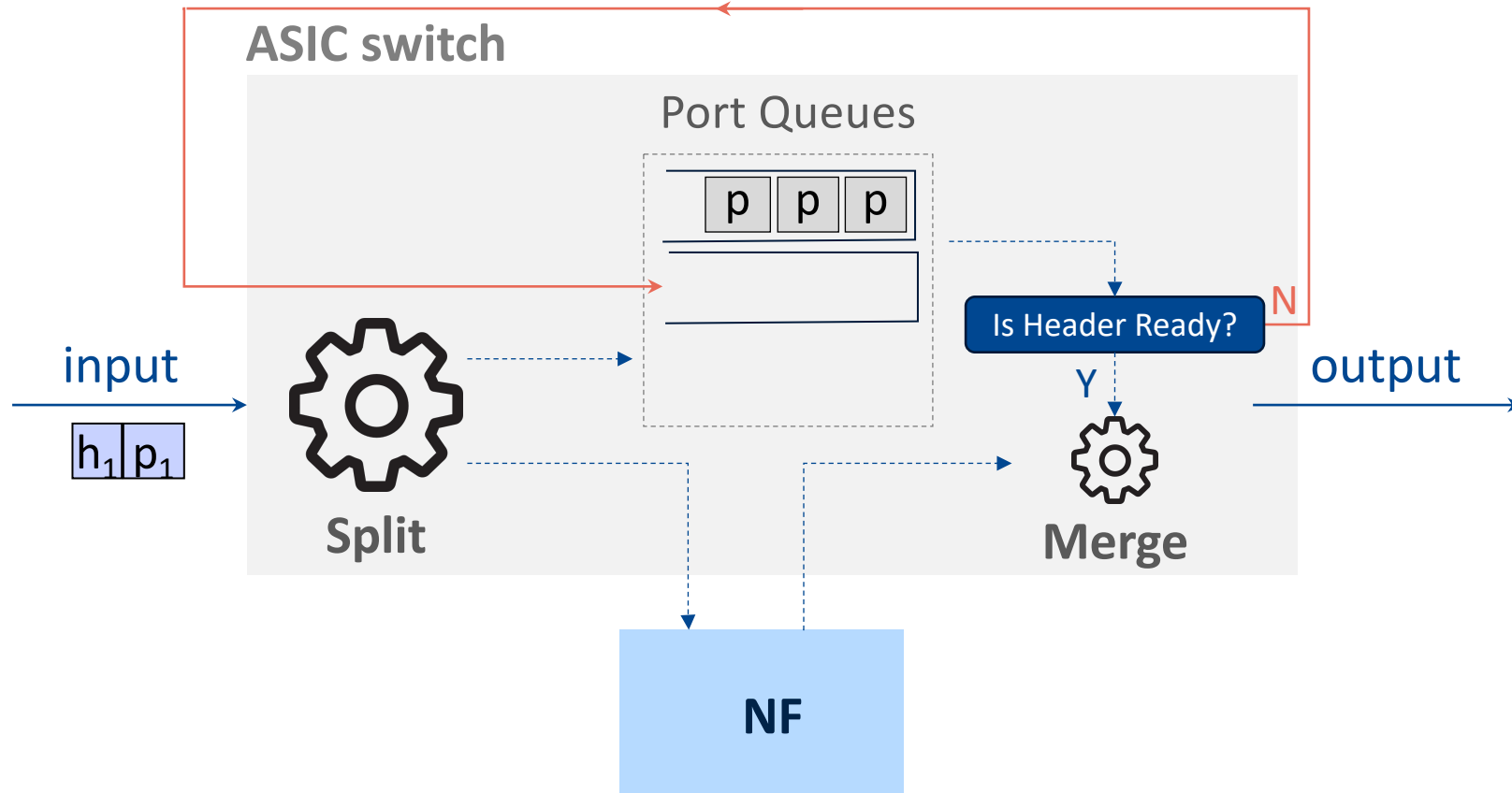
Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



Queue self-clocking



TurboSwitch



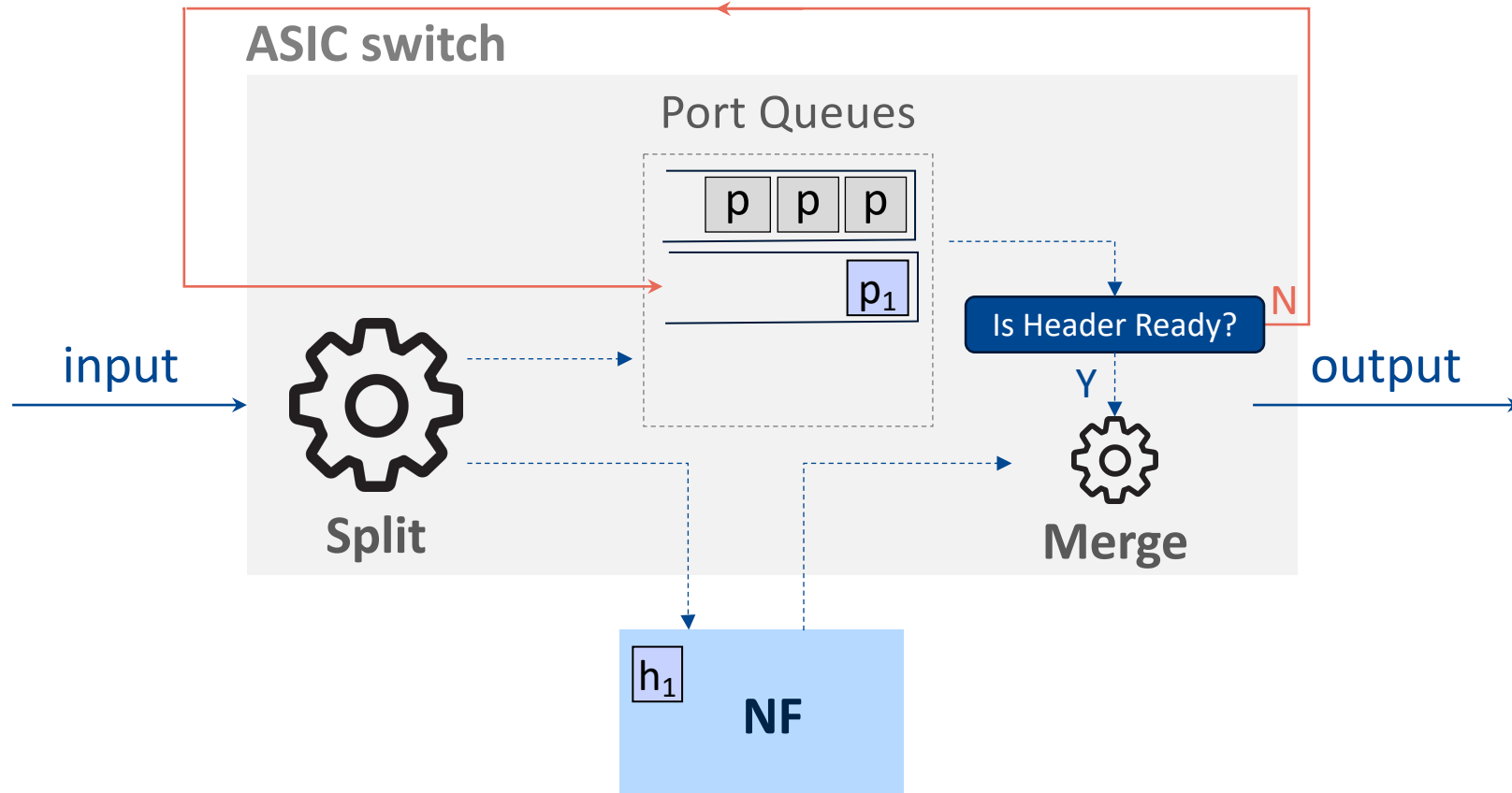
Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



Queue self-clocking



TurboSwitch



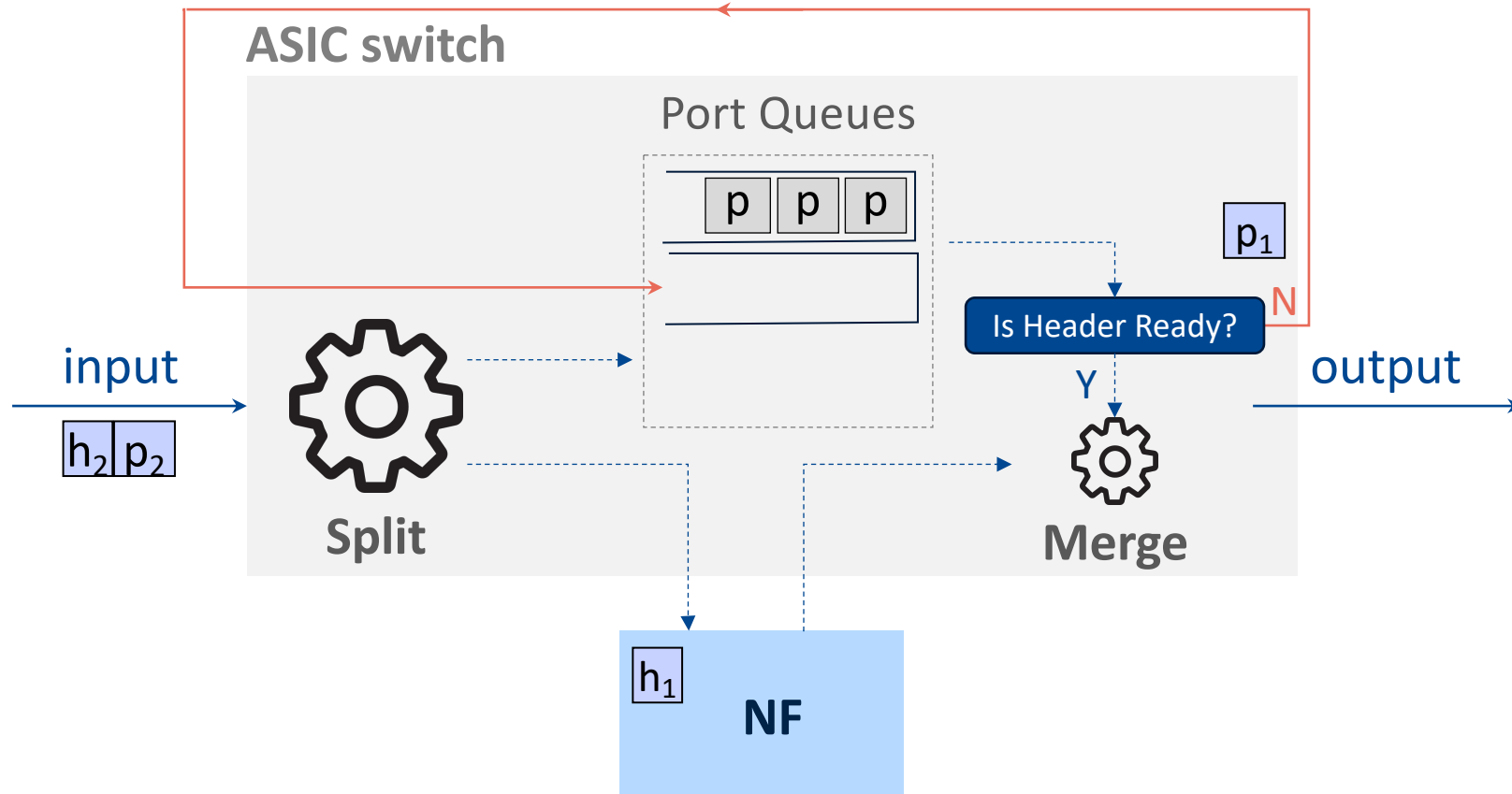
Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



Queue self-clocking



TurboSwitch



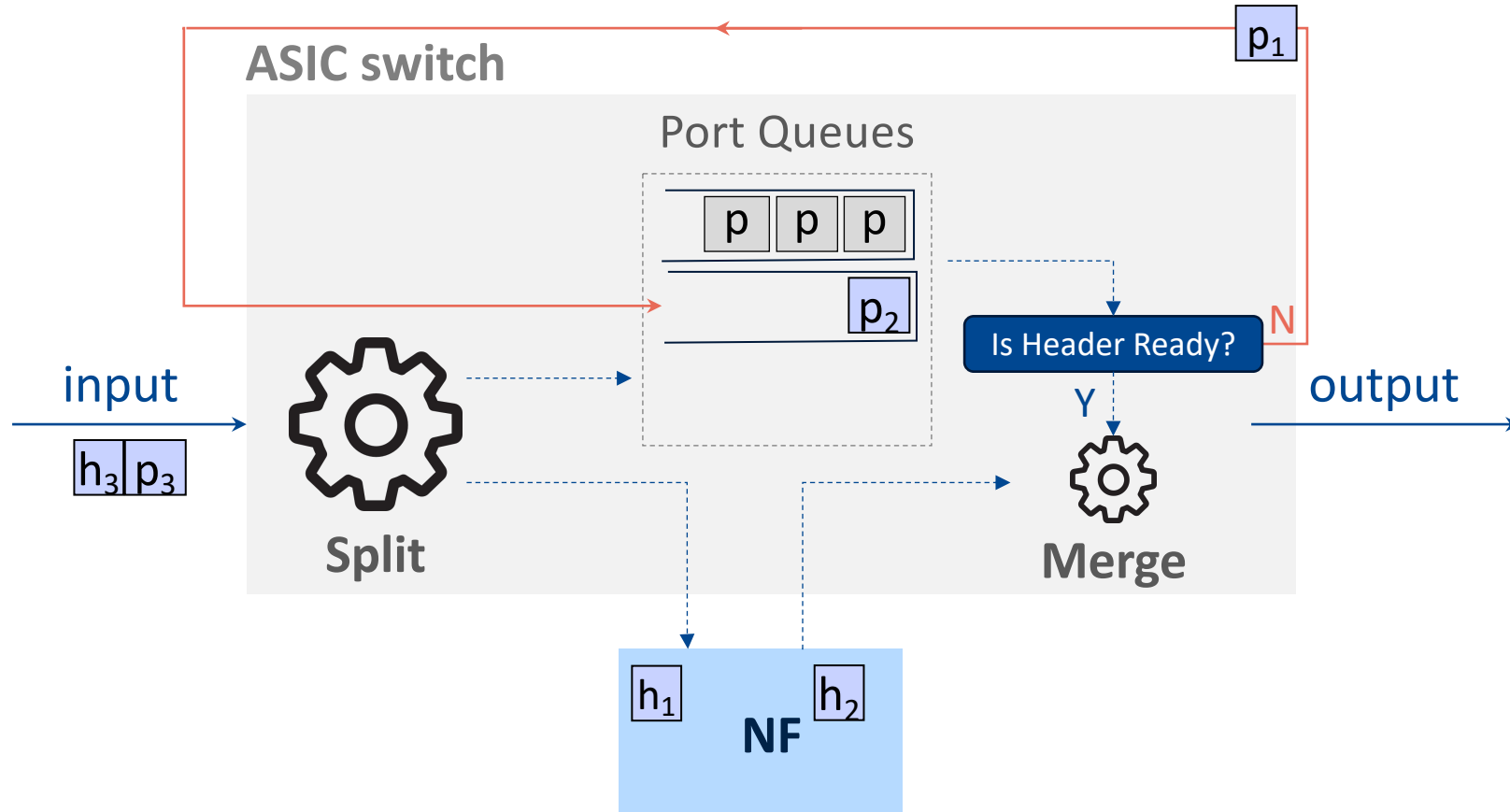
Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



Queue self-clocking



TurboSwitch



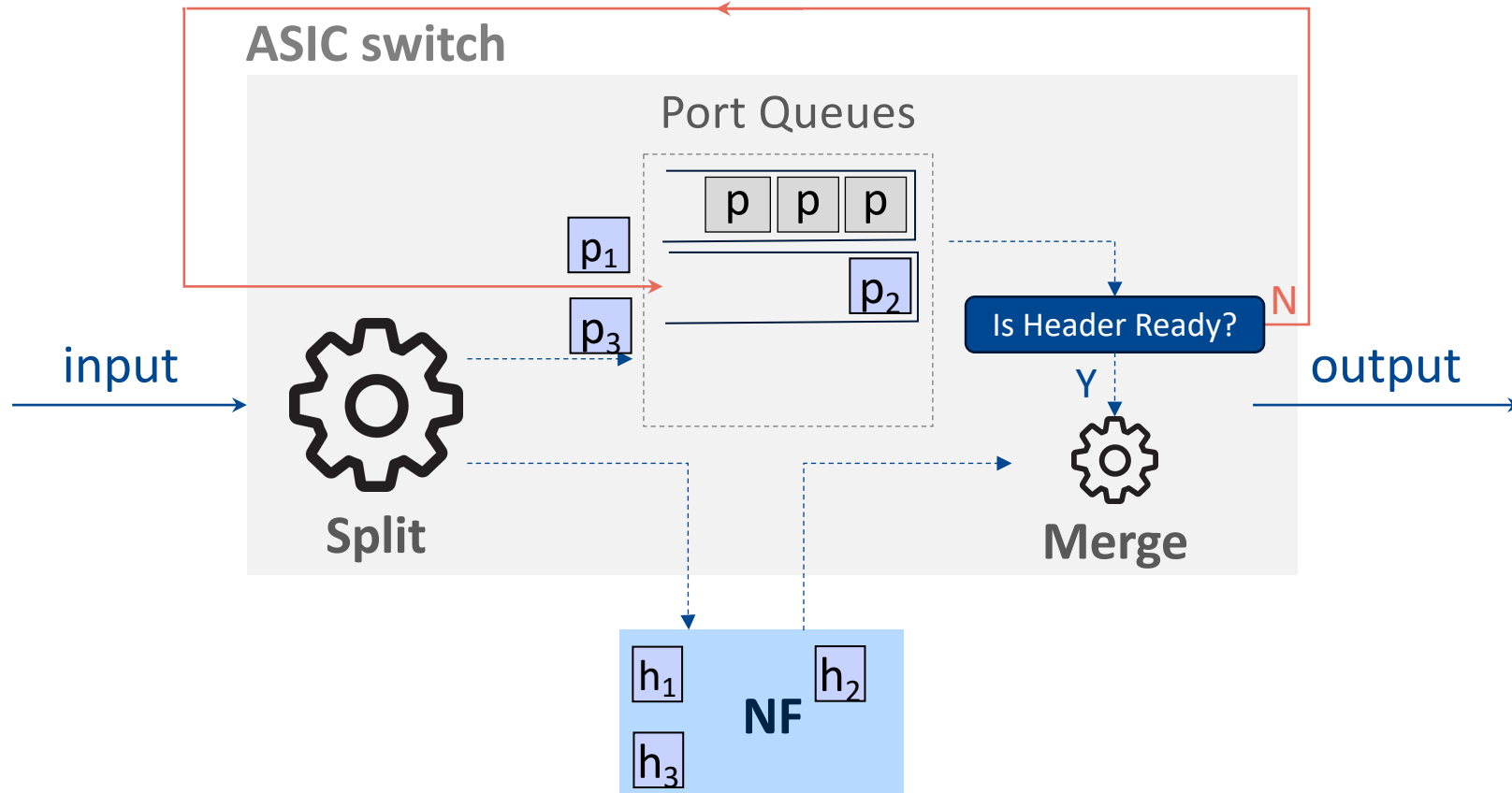
Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



Queue self-clocking



TurboSwitch



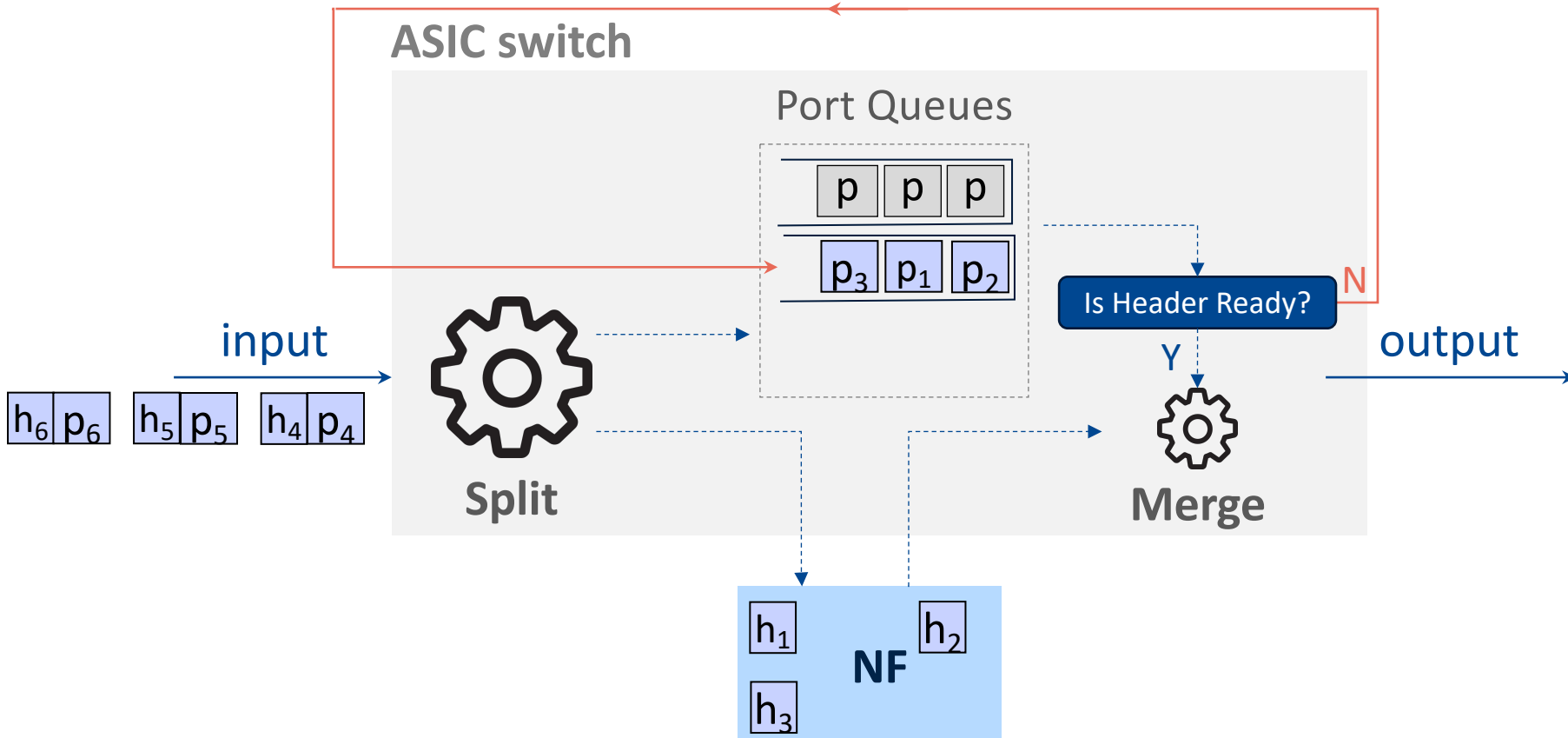
Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



Queue self-clocking



TurboSwitch



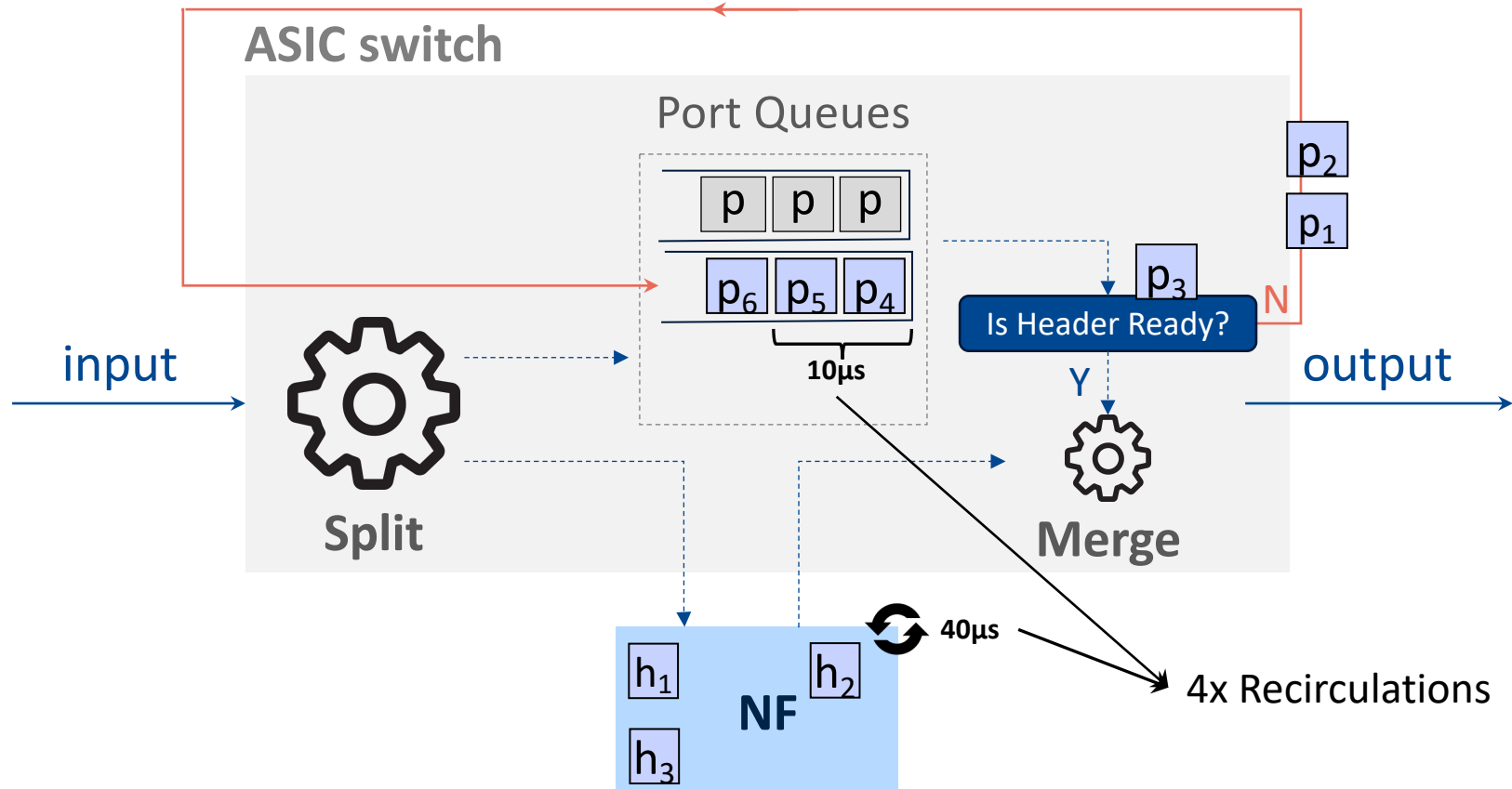
Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



Queue self-clocking



TurboSwitch



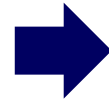
Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



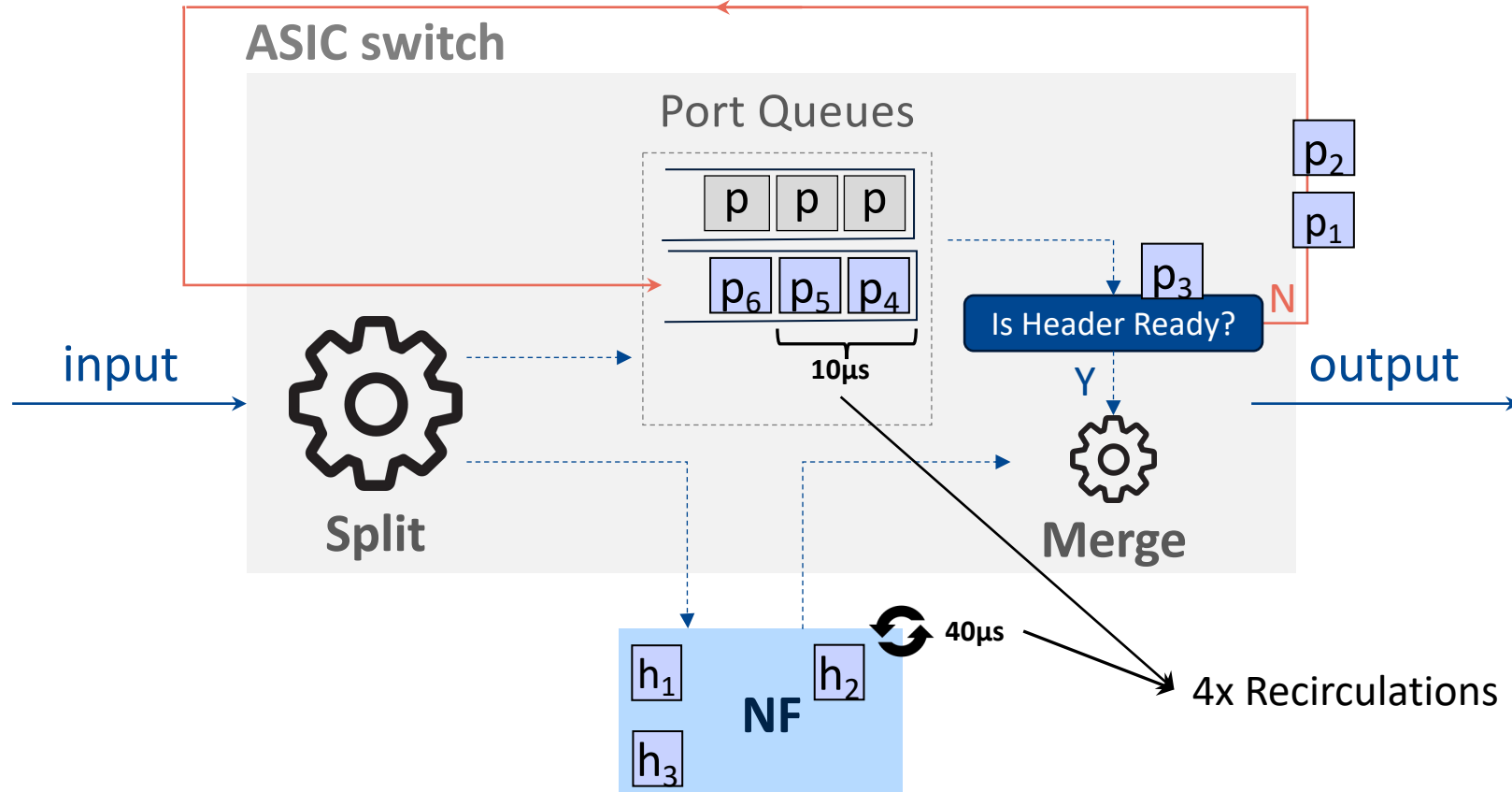
Queue self-clocking



At steady state, # recirculations is predictable



Queue latency converges to the correct value w/o setting any parameter!



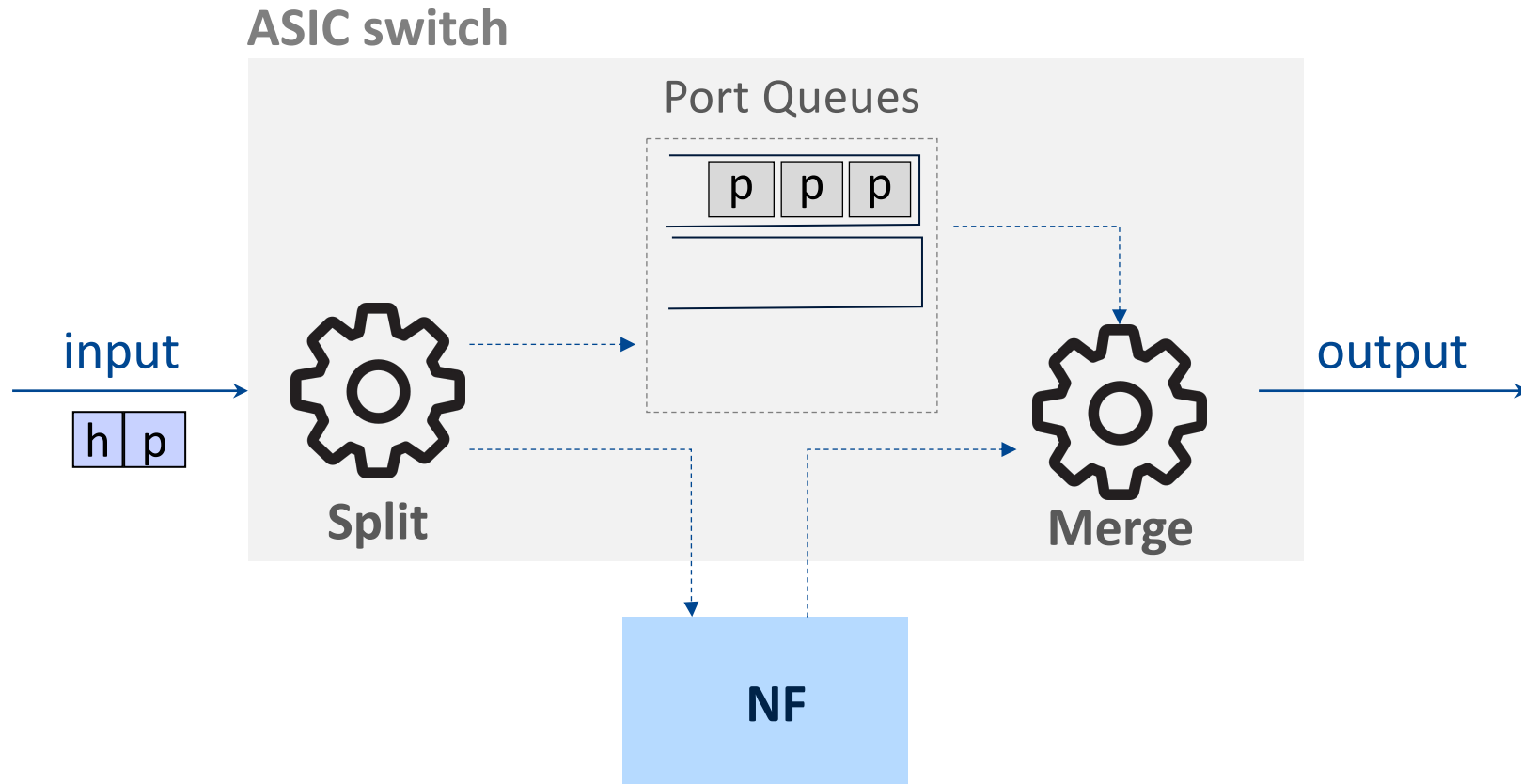
TurboSwitch



Create copies of payloads to fill up the queues



Recirculate payloads until headers are ready



TurboSwitch

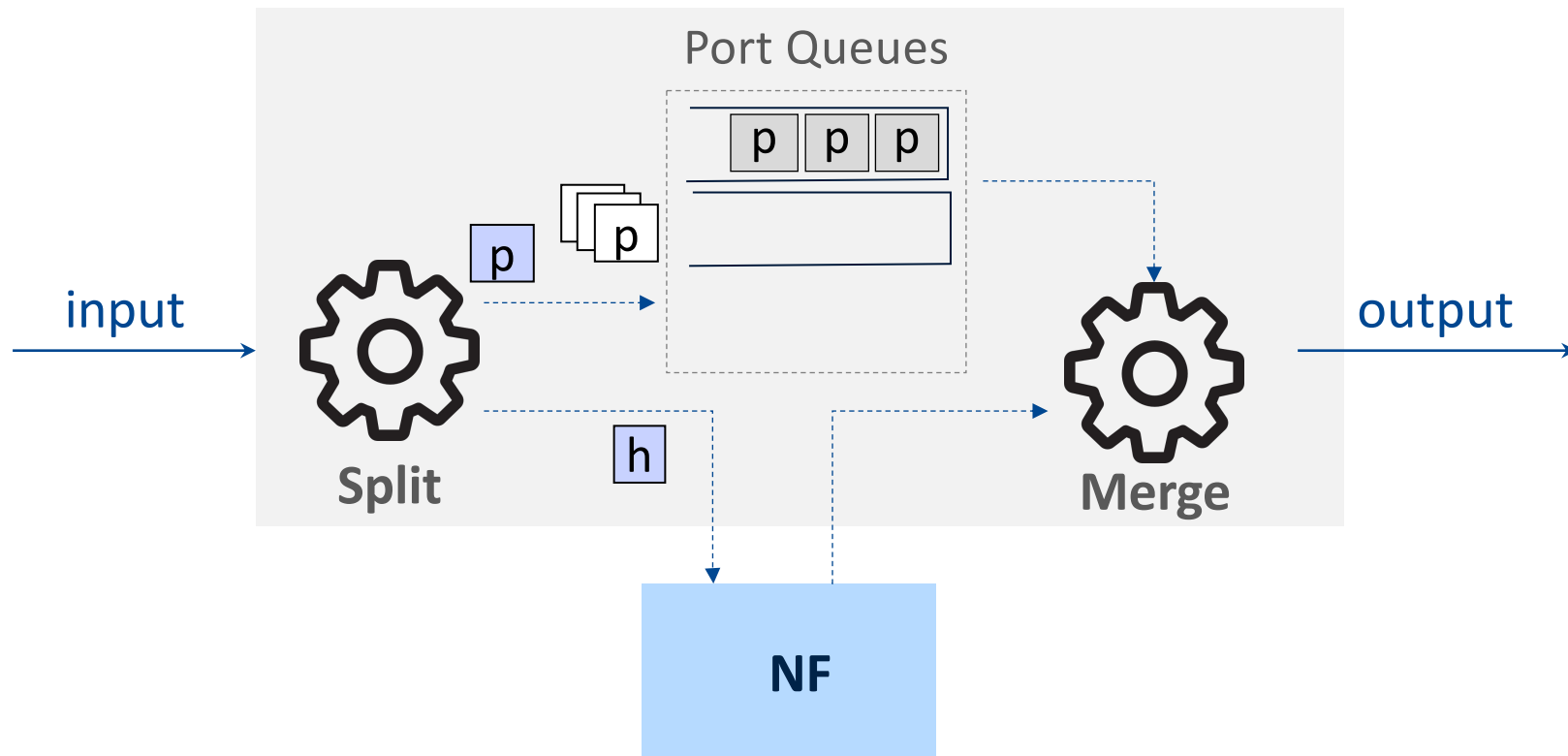


Create copies of payloads to fill up the queues



Recirculate payloads until headers are ready

ASIC switch



TurboSwitch

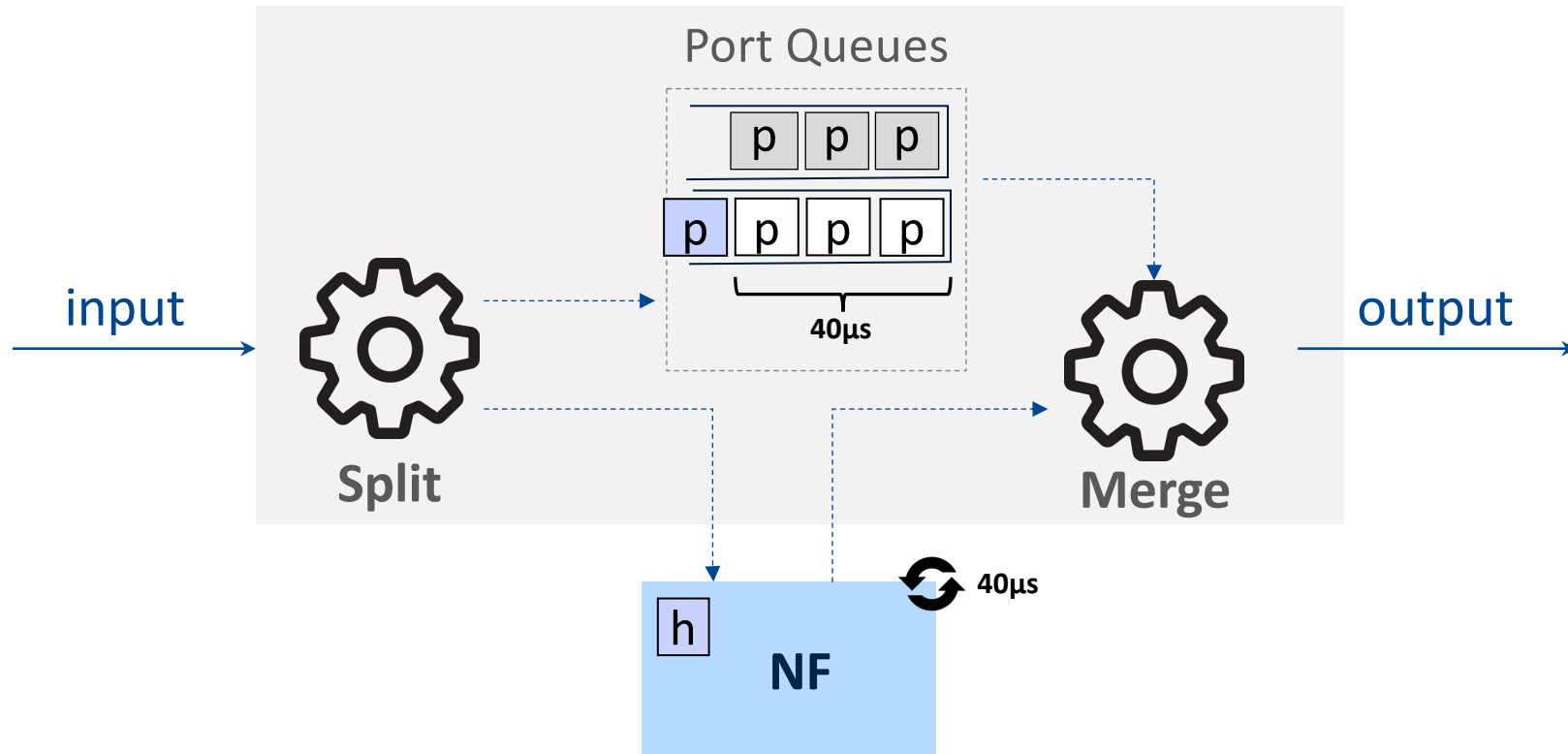


Create copies of payloads to fill up the queues



Recirculate payloads until headers are ready

ASIC switch



TurboSwitch

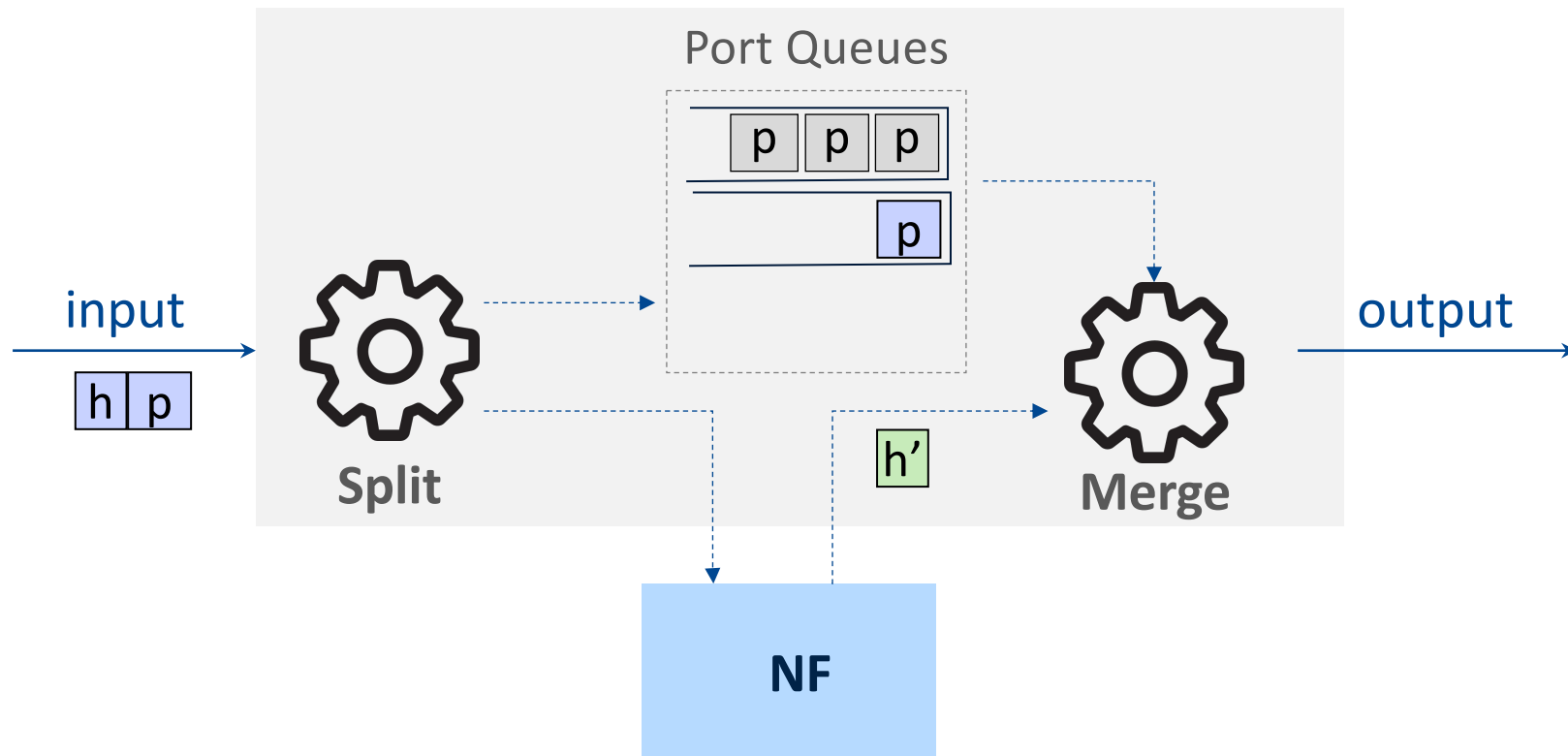


Create copies of payloads to fill up the queues



Recirculate payloads until headers are ready

ASIC switch



TurboSwitch

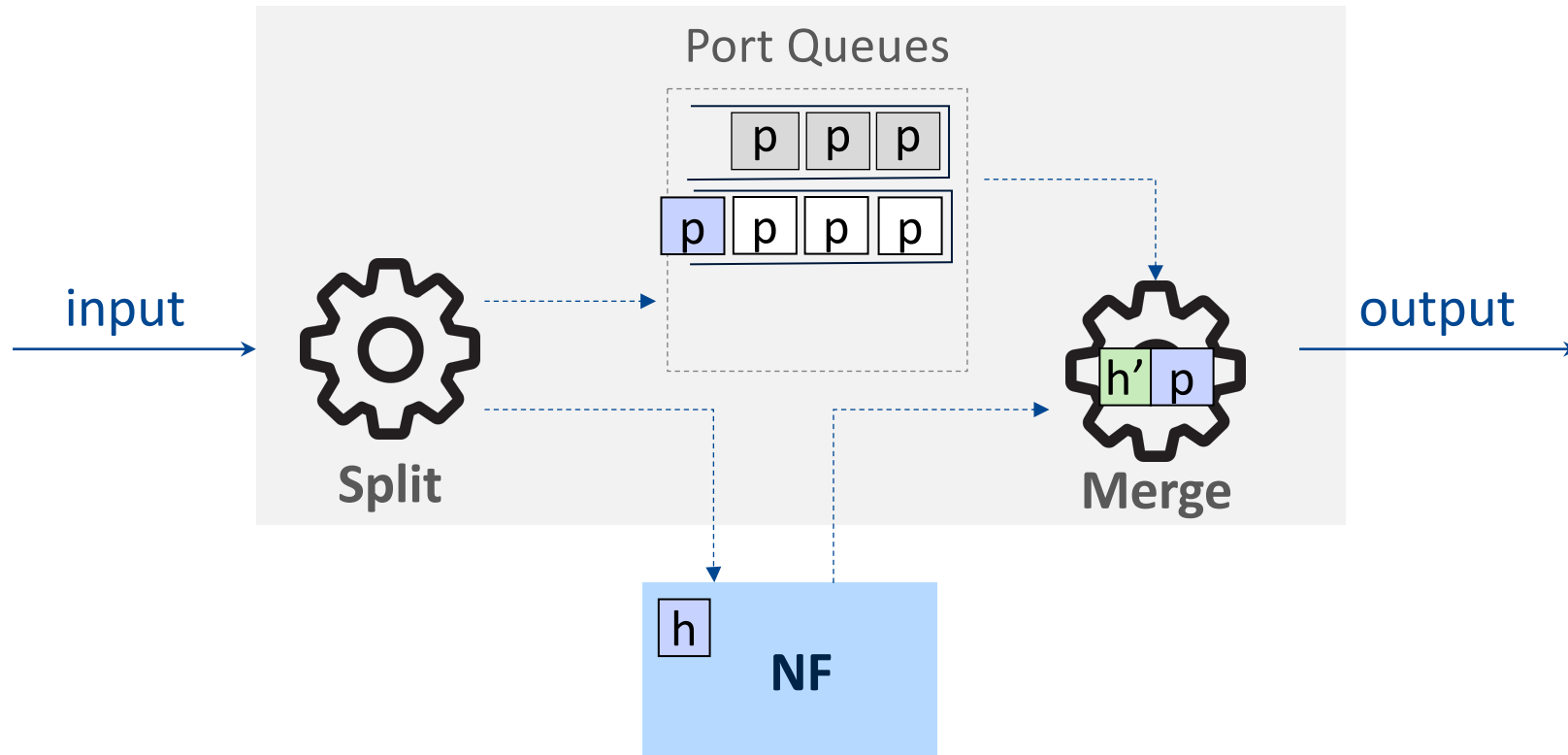


Create copies of payloads to fill up the queues



Recirculate payloads until headers are ready

ASIC switch



TurboSwitch

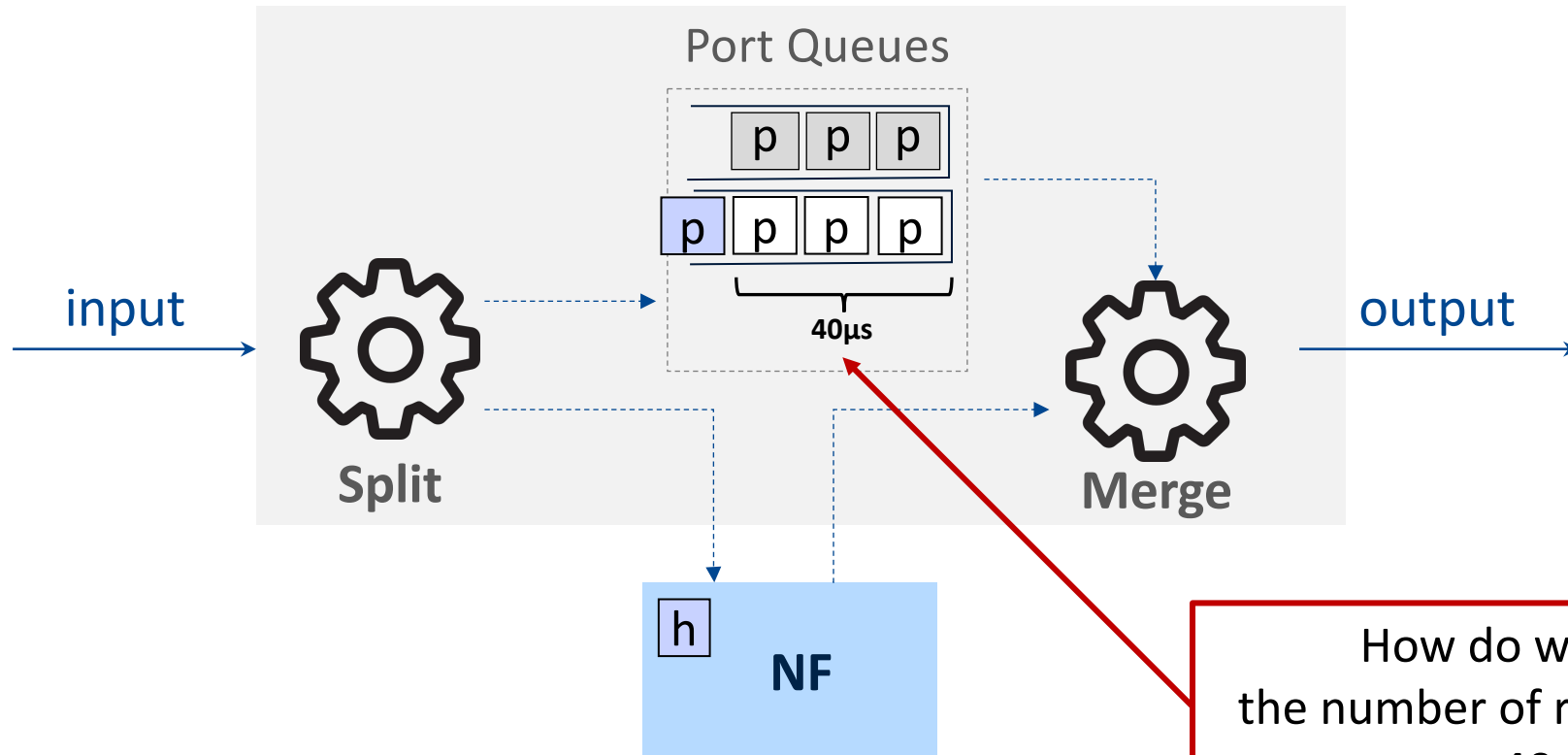


Create copies of payloads to fill up the queues



Recirculate payloads until headers are ready

ASIC switch



How do we control the number of required copies to get $40\mu\text{s}$ delay?

TurboSwitch

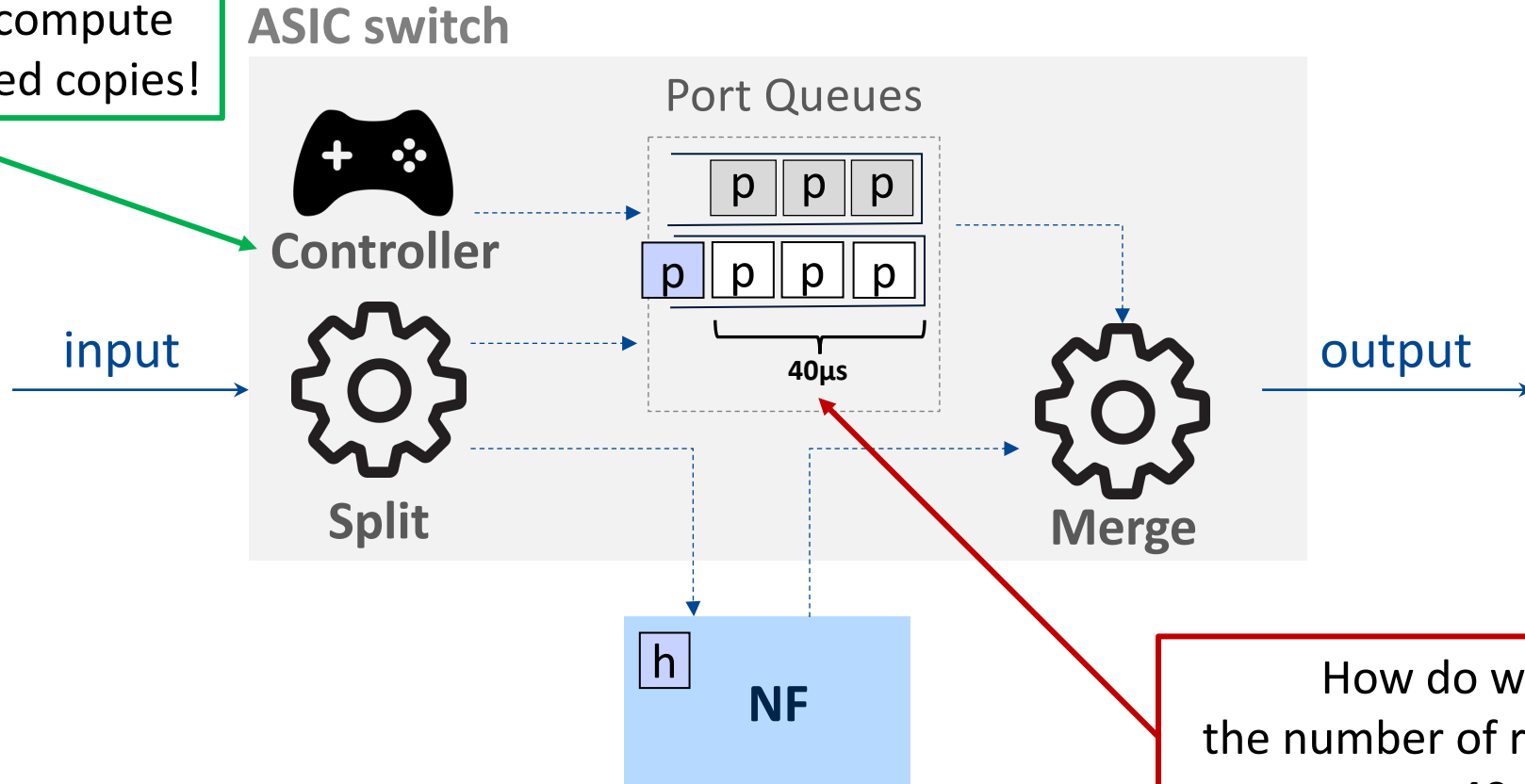


Create copies of payloads to fill up the queues



Recirculate payloads until headers are ready

Use Queue Stats to compute the number of required copies!



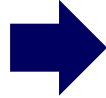
How do we control the number of required copies to get 40 μ s delay?

Evaluation

Evaluation

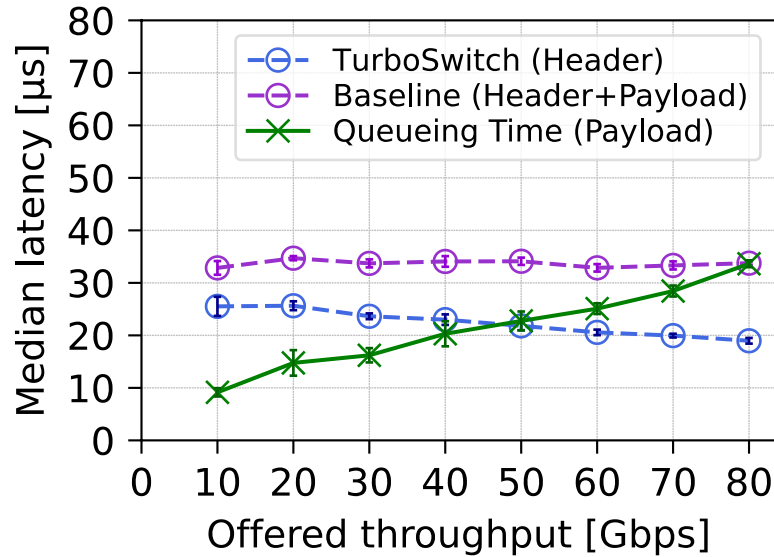


Recirculate payloads



- Intel Tofino ASIC
- FastClick Forwarder NF \Rightarrow 20 μ s Average Latency
- 100Gbps wiring

What are the NF latency gains of TurboSwitch?



Self-clocking to the correct delay!

33% reduction

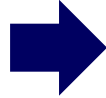


Is TurboSwitch *really* self-clocking?

Evaluation

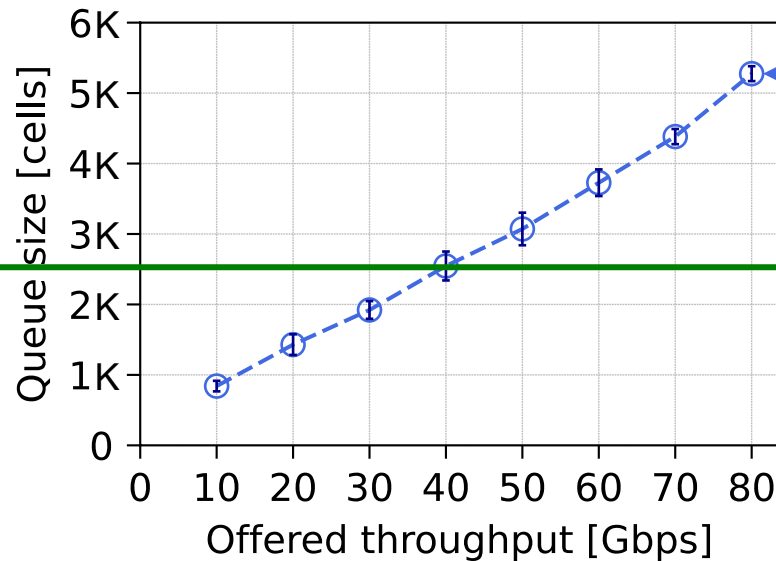
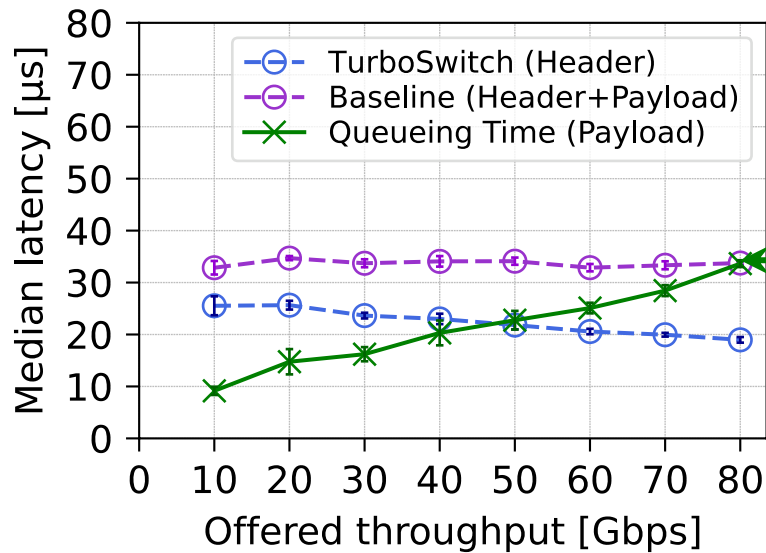


Recirculate payloads



- Intel Tofino ASIC
- FastClick Forwarder NF \Rightarrow 25 μ s Average Latency
- 100Gbps wiring

Is TurboSwitch *really* self-clocking?



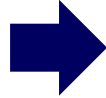
Queue size grows linearly w.r.t. input throughput

Linear relationship between payload queueing time and queue size

Evaluation

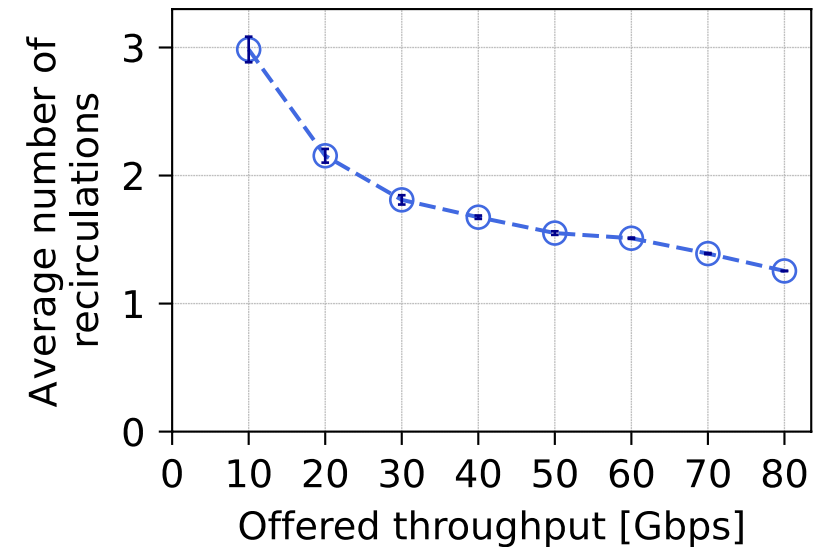
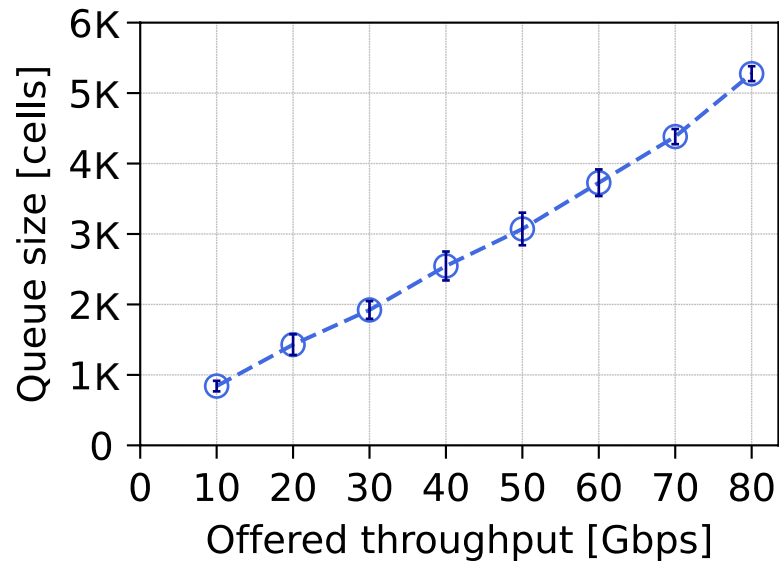
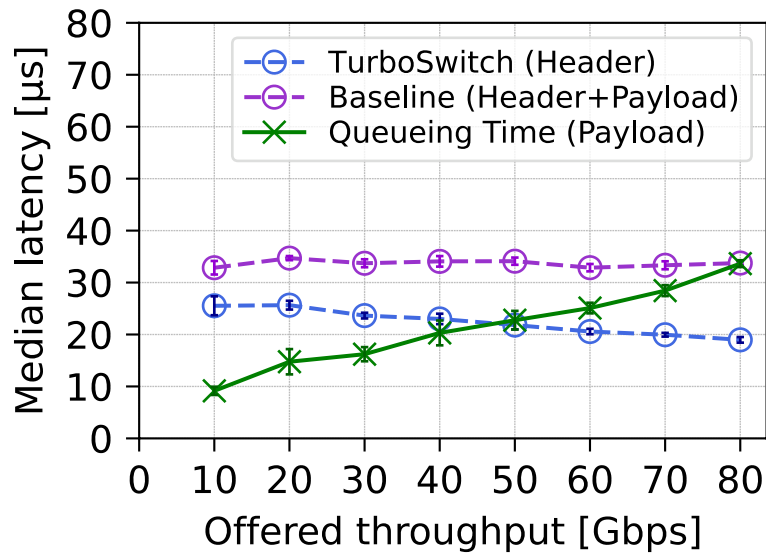


Recirculate payloads



- Intel Tofino ASIC
- FastClick Forwarder NF \Rightarrow 25 μ s Average Latency
- 100Gbps wiring

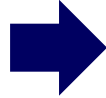
Is TurboSwitch *really* self-clocking?



Evaluation

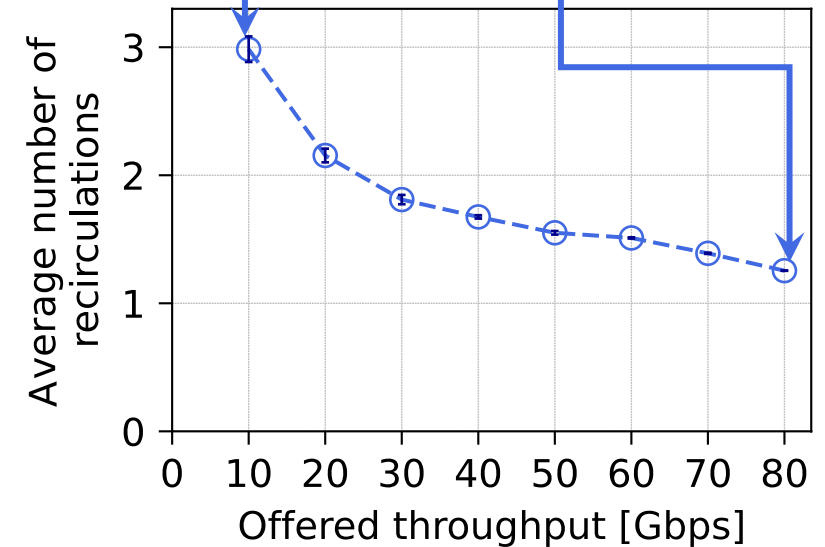
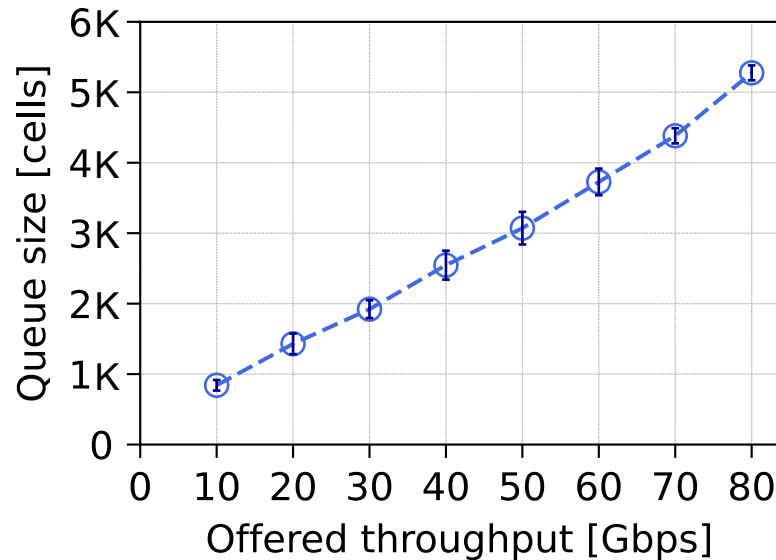
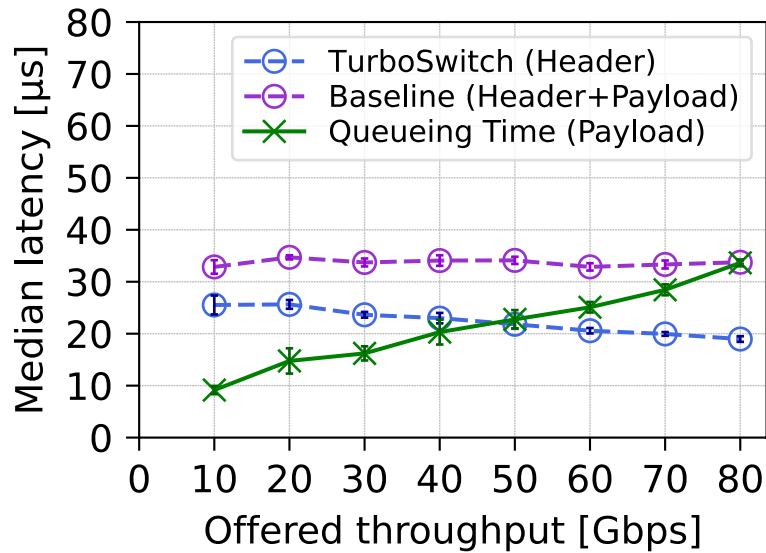


Recirculate payloads



- Intel Tofino ASIC
- FastClick Forwarder NF \Rightarrow 25 μ s Average Latency
- 100Gbps wiring

Is TurboSwitch *really* self-clocking?



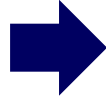
More recirculations to keep a steady state

~1 recirculation per packet

Evaluation



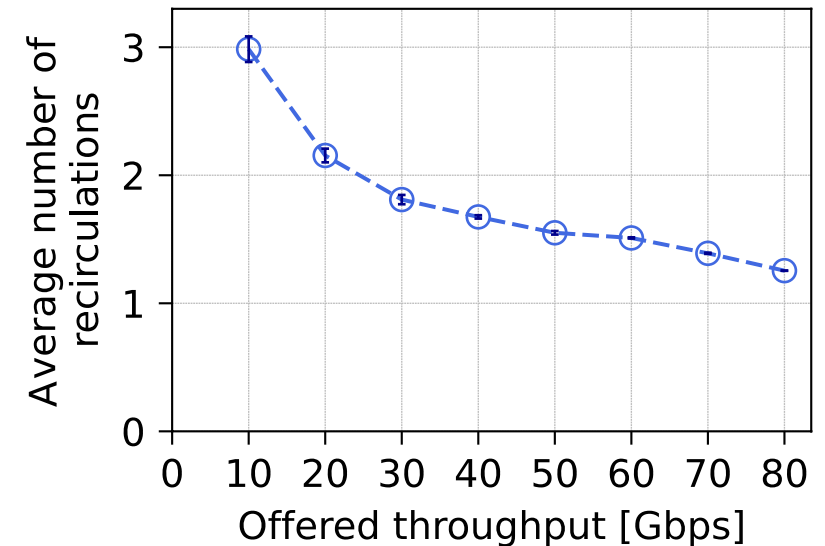
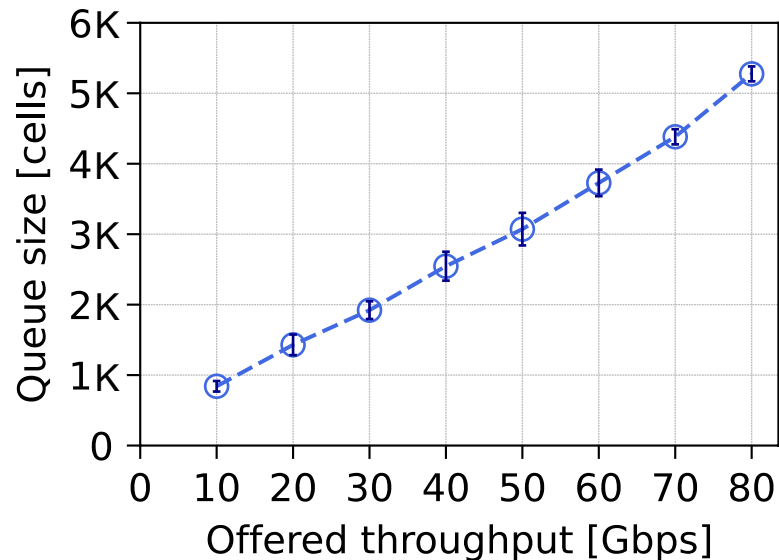
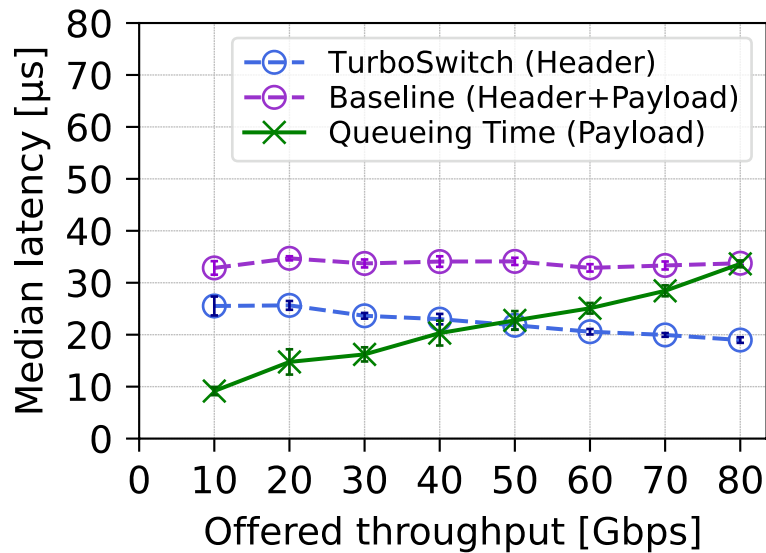
Recirculate payloads



- Intel Tofino ASIC
- FastClick Forwarder NF \Rightarrow 25 μ s Average Latency
- 100Gbps wiring

Is TurboSwitch *really* self-clocking?

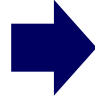
Turboswitch **improves** NF performance by inducing **self-clocked congestion!**



Evaluation



Create copies of payloads



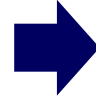
- Python Simulator – 20 μ s NF Processing Latency
- 100Gbps w/ Poisson arrival (64B-1.5KB pkt size)
- Controller \Rightarrow 1 μ s Queue Stats Update Delay

Can TurboSwitch's controller enforce the correct delay to payloads?

Evaluation

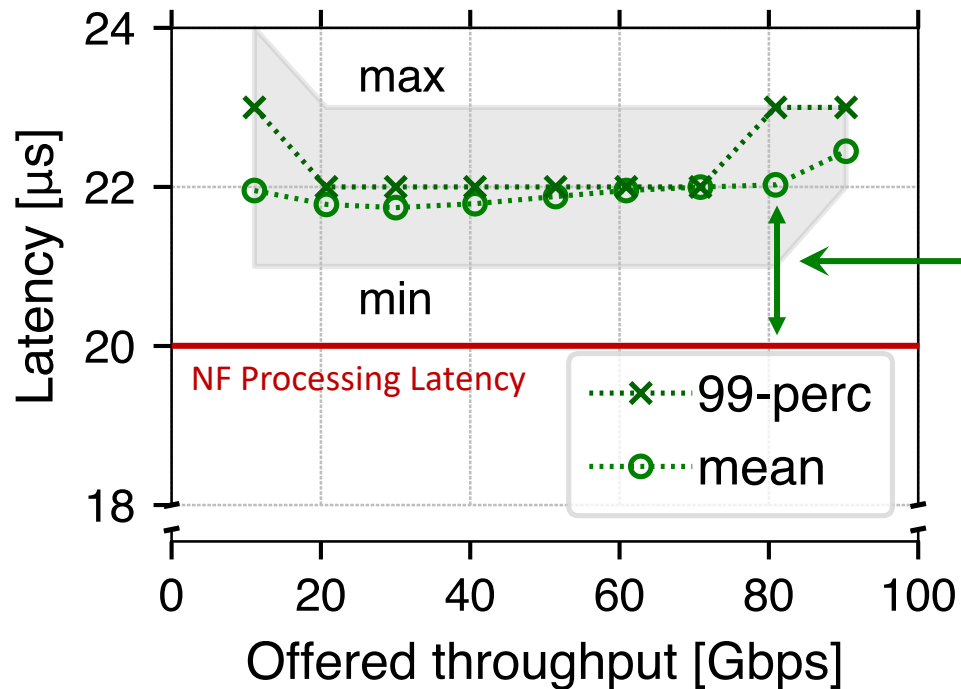


Create copies of payloads



- Python Simulator – 20 μ s NF Processing Latency
- 100Gbps w/ Poisson arrival (64B-1.5KB pkt size)
- Controller \Rightarrow 1 μ s Queue Stats Update Delay

Can TurboSwitch's controller enforce the correct delay to payloads?

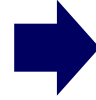


Stable payload latency at $\sim 22\mu$ s

Evaluation

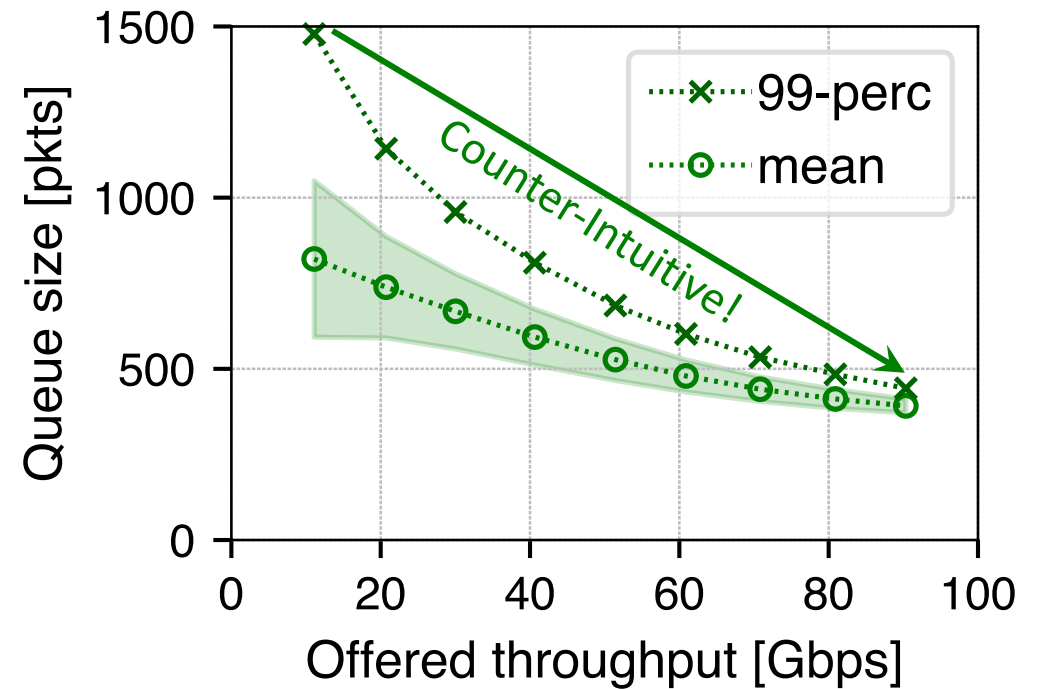
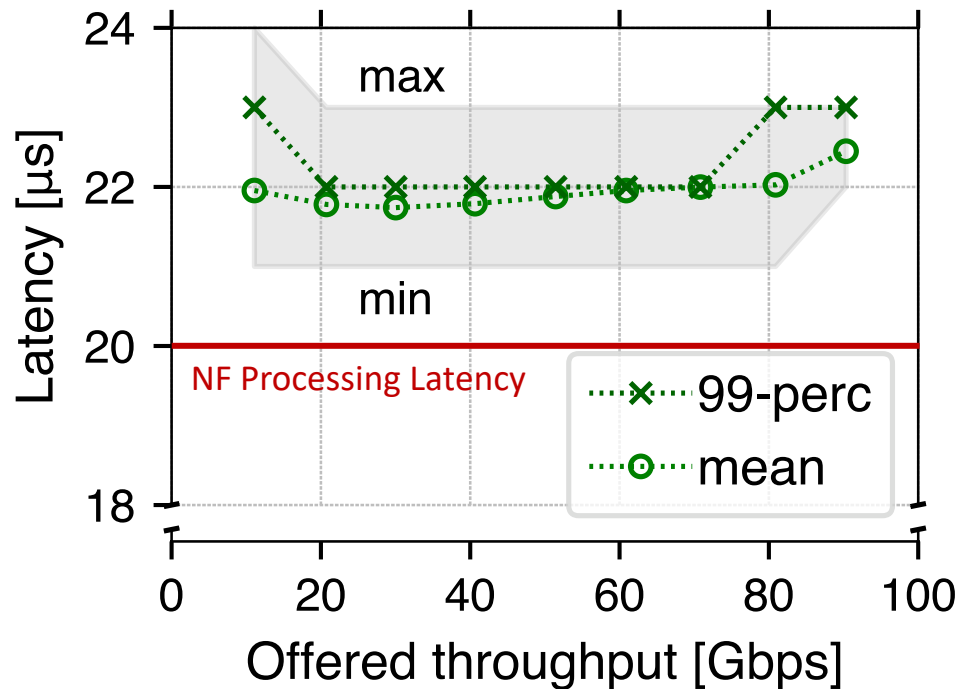


Create copies of payloads



- Python Simulator – 20 μ s NF Processing Latency
- 100Gbps w/ Poisson arrival (64B-1.5KB pkt size)
- Controller \Rightarrow 1 μ s Queue Stats Update Delay

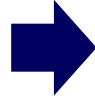
Can TurboSwitch's controller enforce the correct delay to payloads?



Evaluation



Create copies of payloads

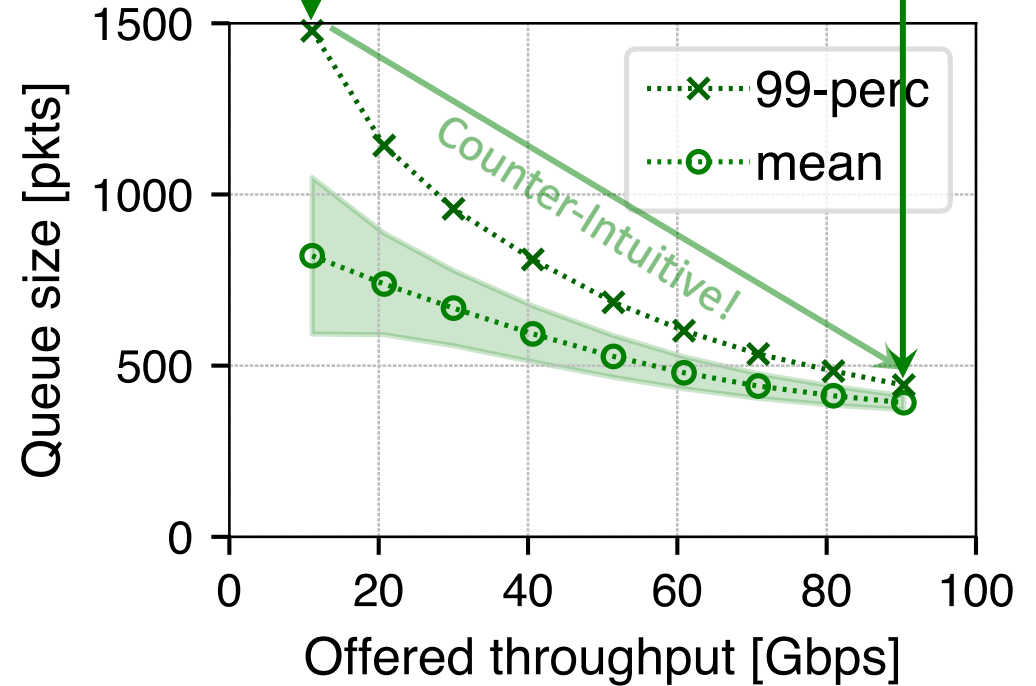
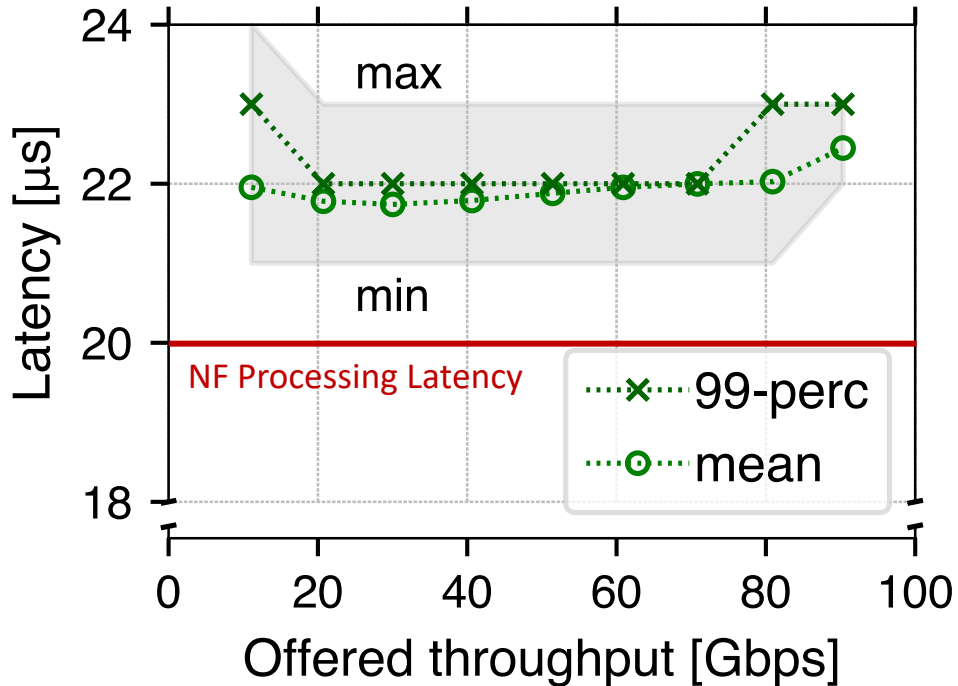


- Python Simulator – 20 μ s NF Processing Latency
- 100Gbps w/ Poisson arrival (64B-1.5KB pkt size)
- Controller \Rightarrow 1 μ s Queue Stats Update Delay

Can TurboSwitch's controller enforce the correct delay to payloads?

↑ Less incoming packets, controller creates **more copies**

↓ More incoming packets, controller creates **less copies**



Evaluation



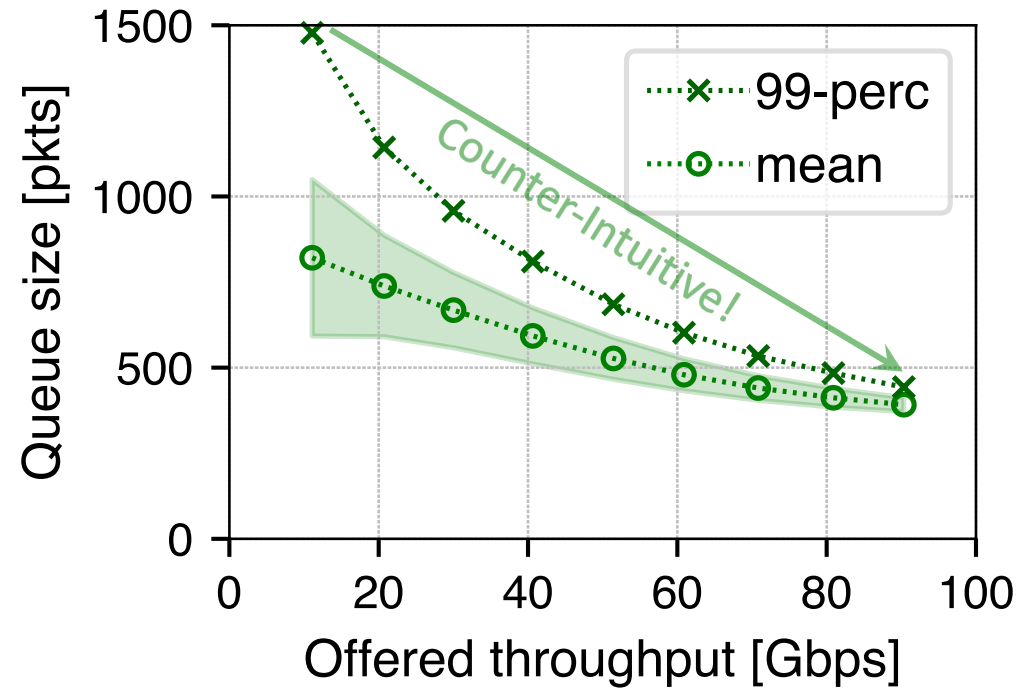
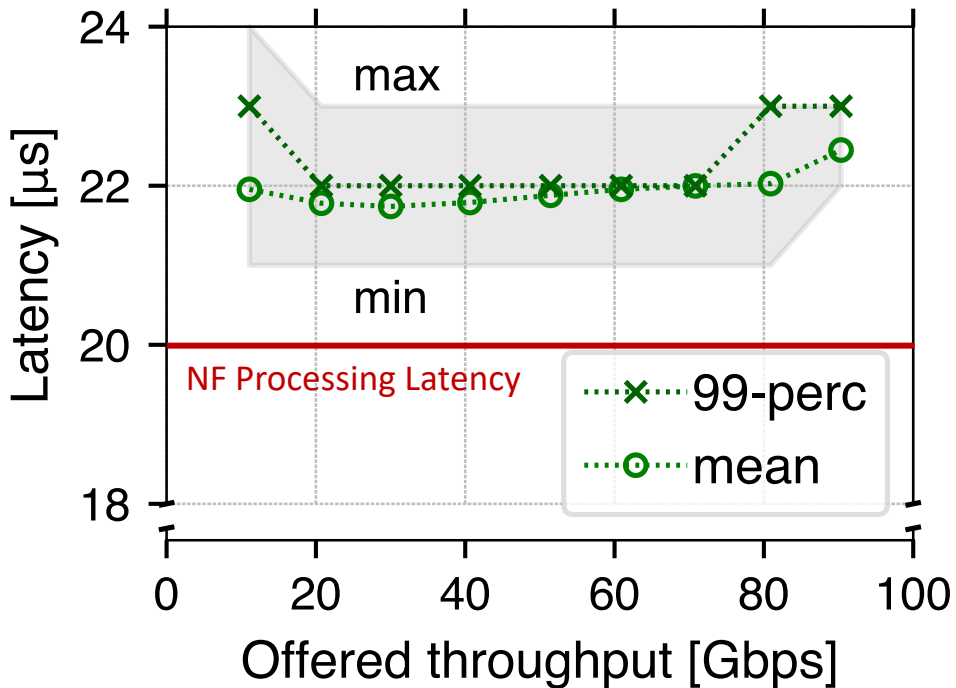
Create copies of payloads



- Python Simulator – 20 μ s NF Processing Latency
- 100Gbps w/ Poisson arrival (64B-1.5KB pkt size)
- Controller \Rightarrow 1 μ s Queue Stats Update Delay

Can TurboSwitch's controller enforce the correct delay to payloads?

TurboSwitch **correctly sizes queues to enforce** the NF processing latency!



Conclusions

The **TurboSwitch** system:

- Presents a packet-splitting approach that **does not require external resources**
- Proposes a **queue-based buffer** to store payloads on existing ASIC switches
- **Reduces** NF processing time with **zero configuration**
- Spurs a question on whether **expose APIs** for controlling the forwarding buffer