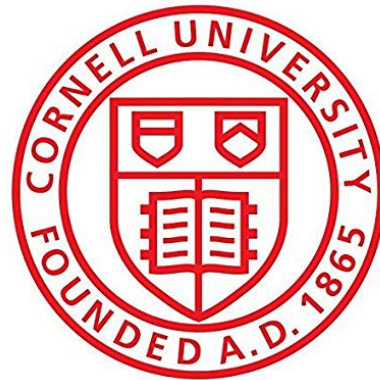


Towards High-speed Networking in the Post-Moore Era

Vishal Shrivastav

Ph.D. Thesis Defense
Cornell University 2020



Datacenters 101



Datacenters 101



Datacenters 101



Datacenters 101

 datacenter serves 2.6 billion active users daily



Datacenters 101

 datacenter serves 2.6 billion active users daily

In 2018, 40% of world total data was stored and processed inside datacenters

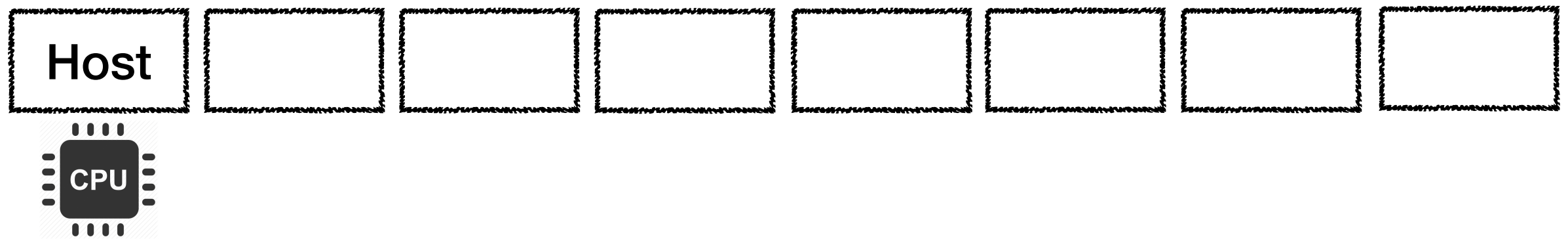


Inside a Typical Datacenter

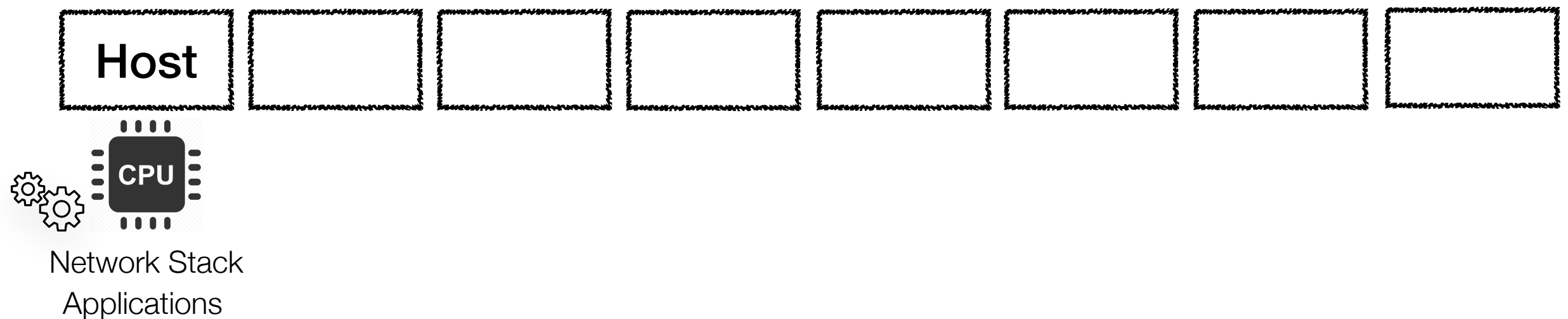
Inside a Typical Datacenter



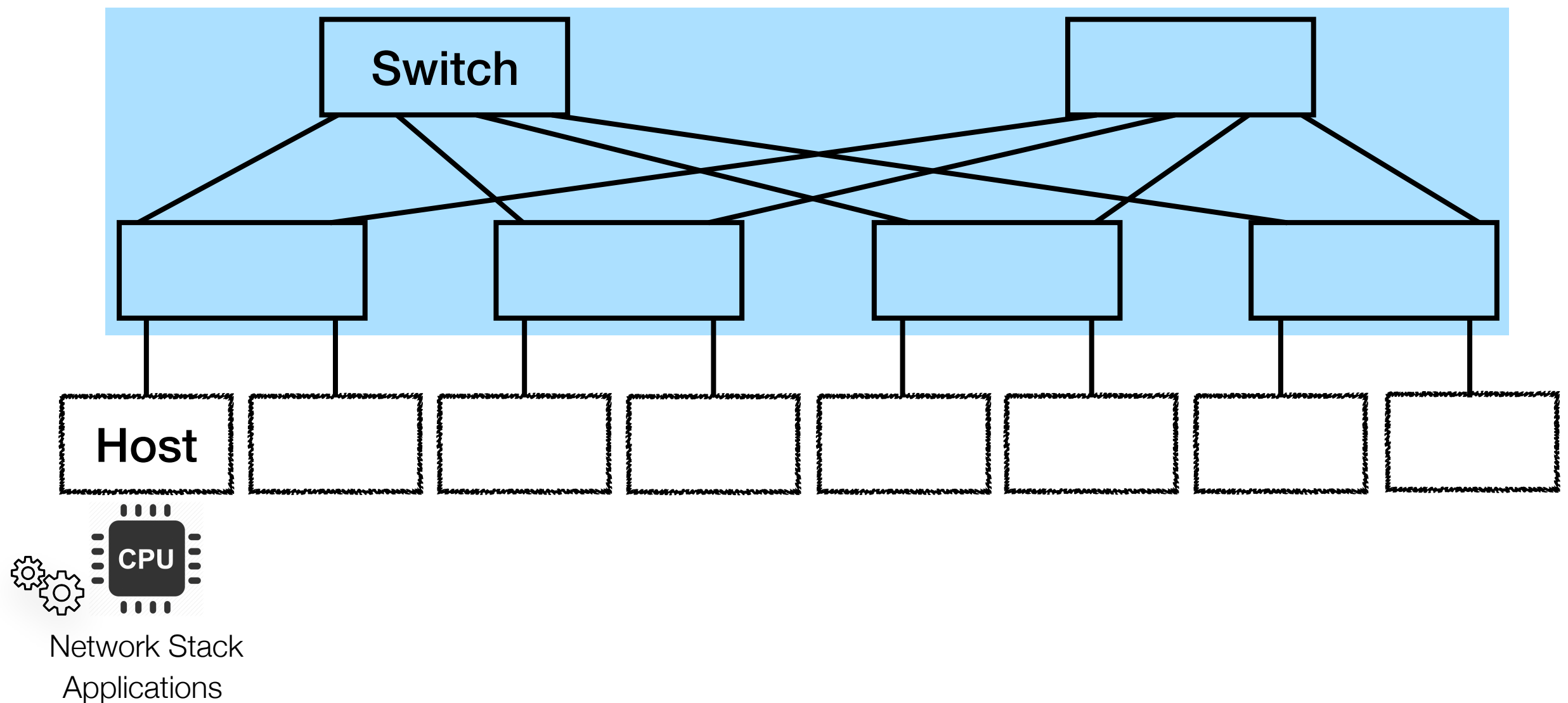
Inside a Typical Datacenter



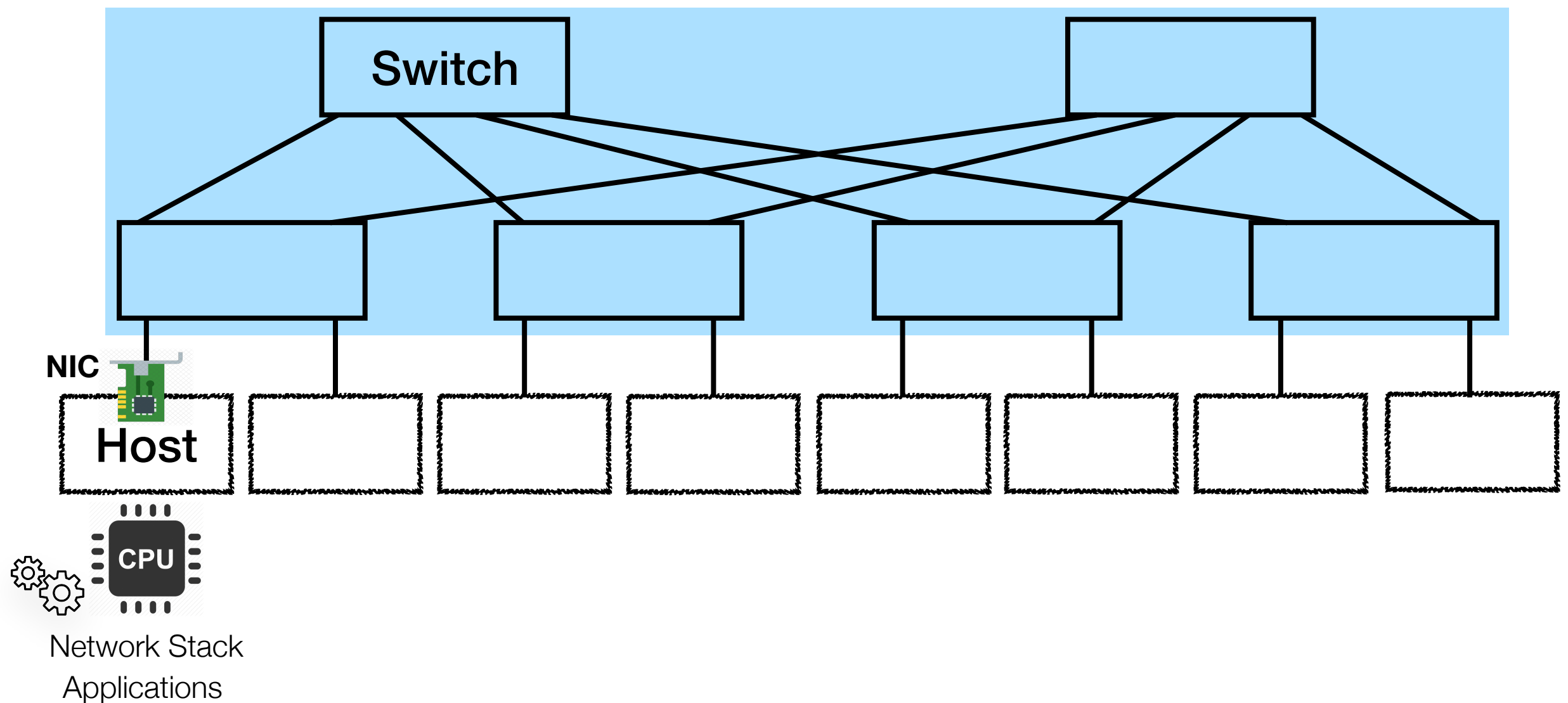
Inside a Typical Datacenter



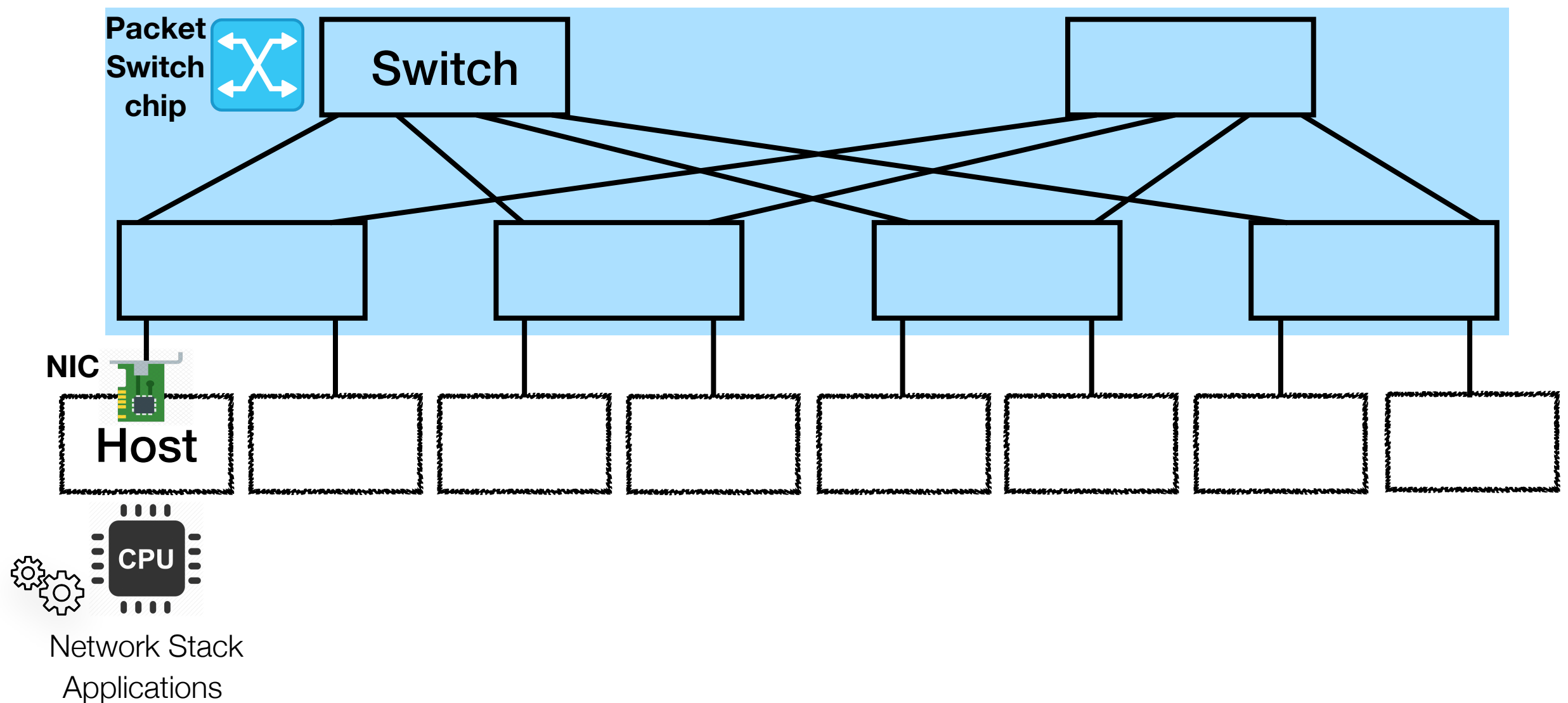
Inside a Typical Datacenter



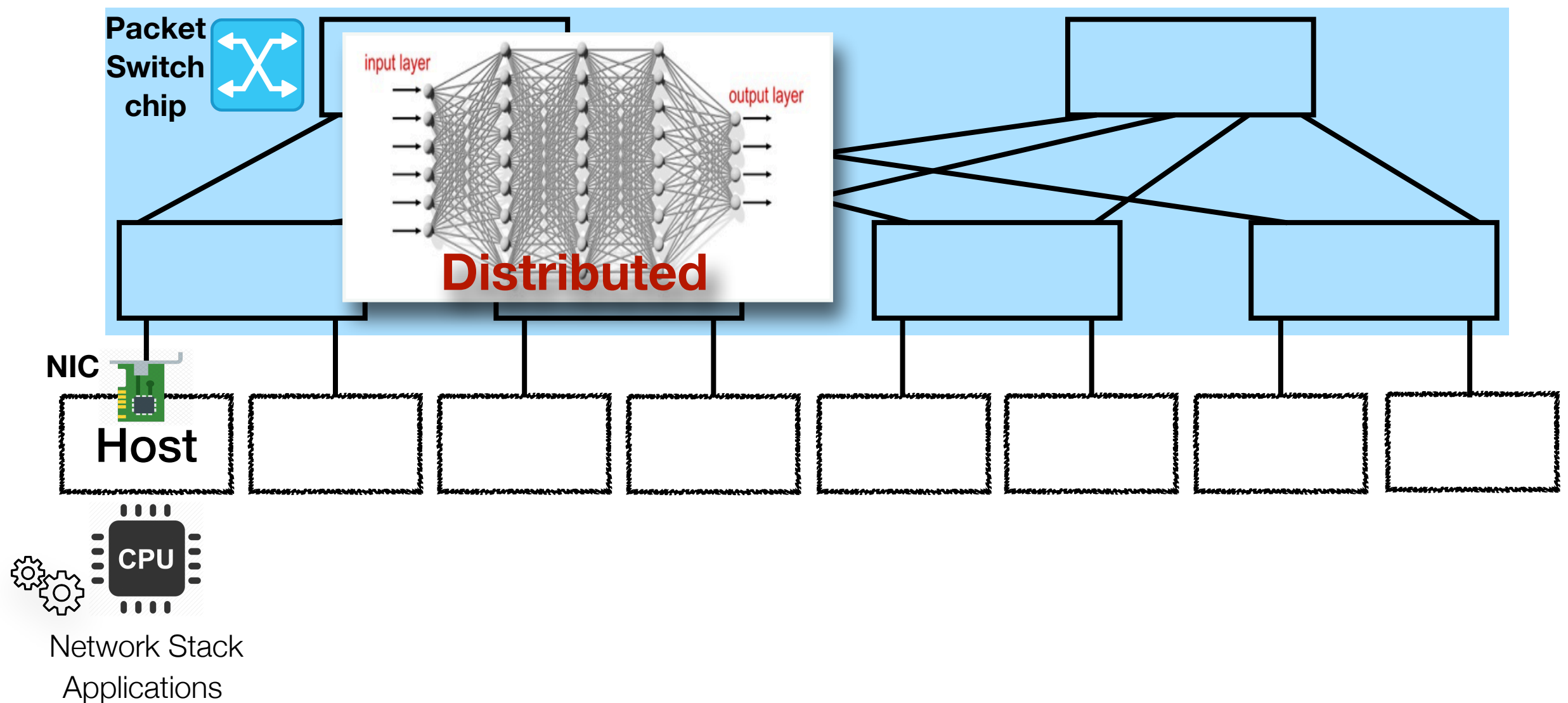
Inside a Typical Datacenter



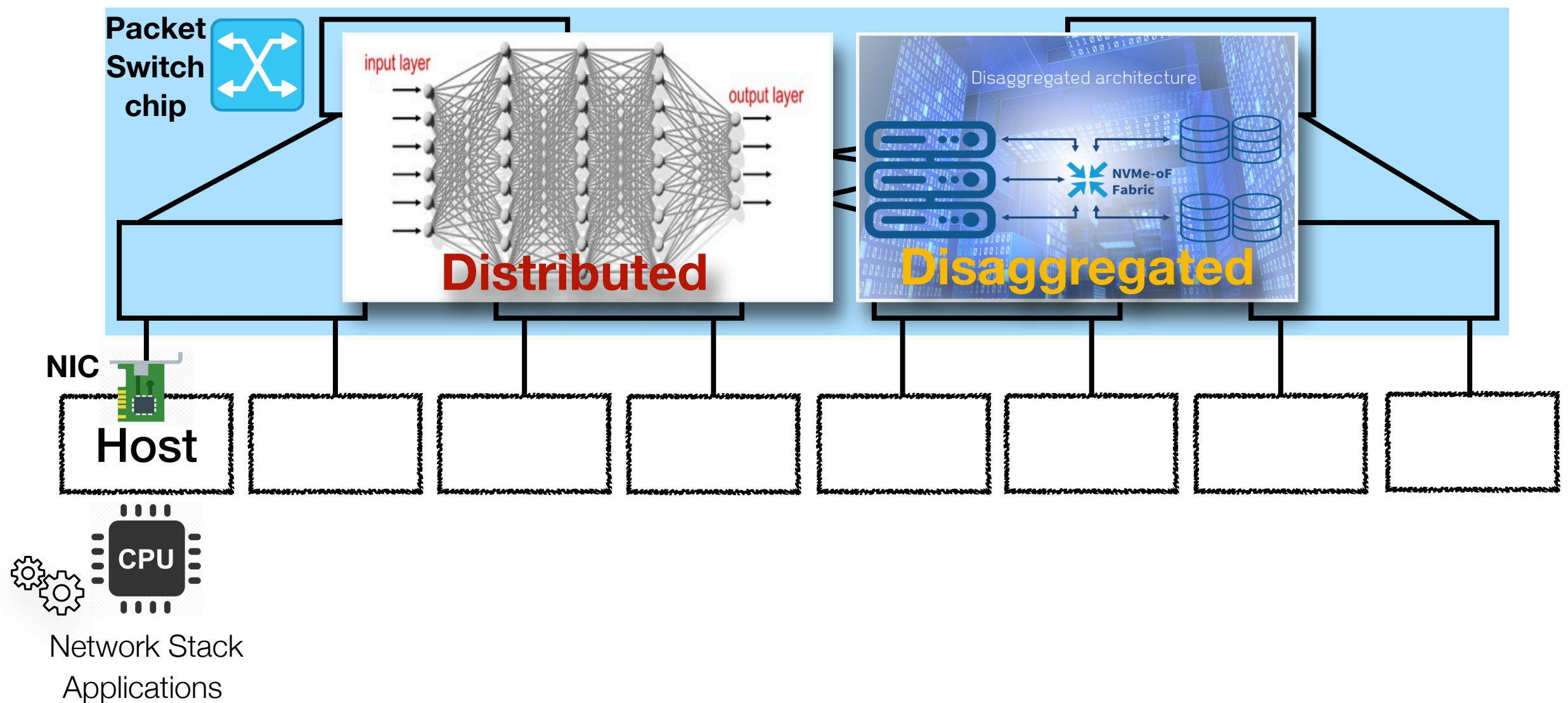
Inside a Typical Datacenter



Inside a Typical Datacenter



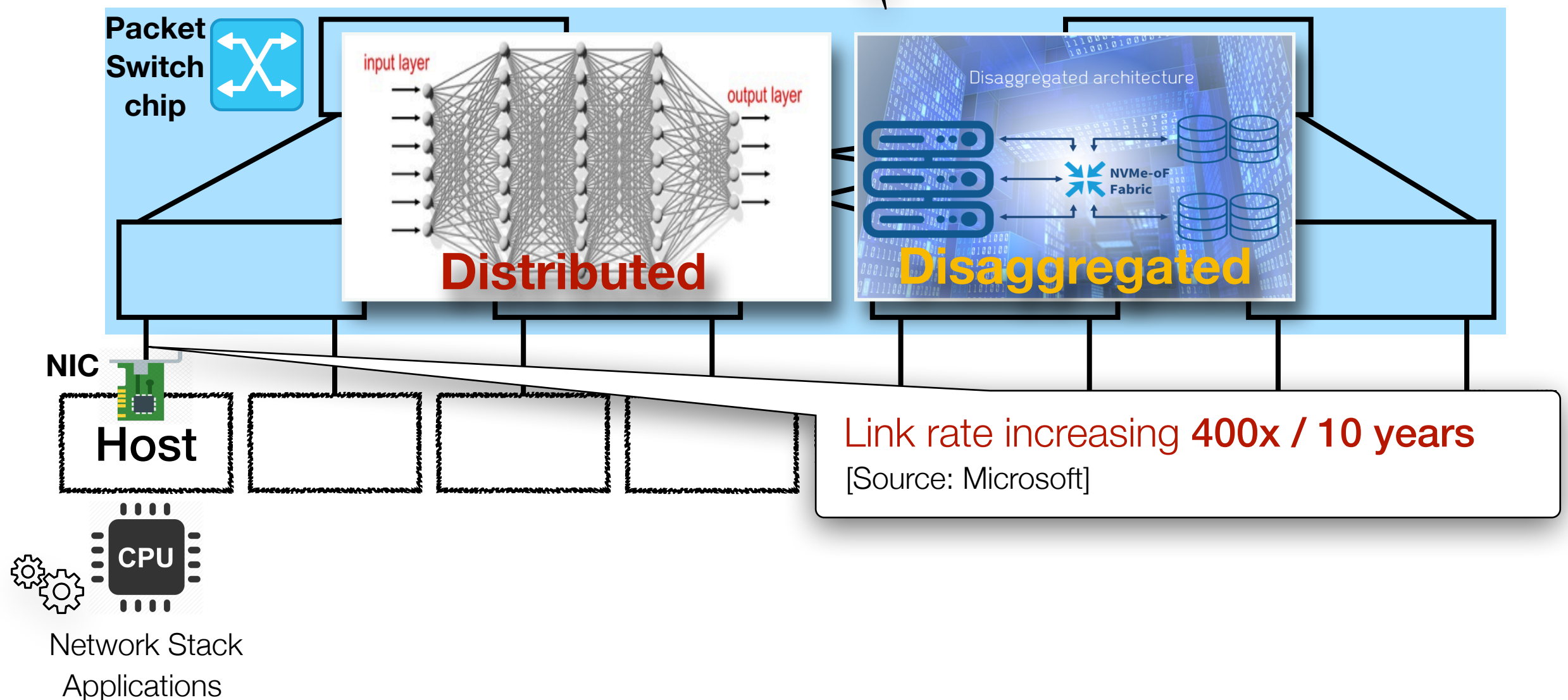
Inside a Typical Datacenter



Inside a Typical Datacenter

Bandwidth demand increasing **2x / year**

[Source: Google]



Link rate increasing **400x / 10 years**

[Source: Microsoft]

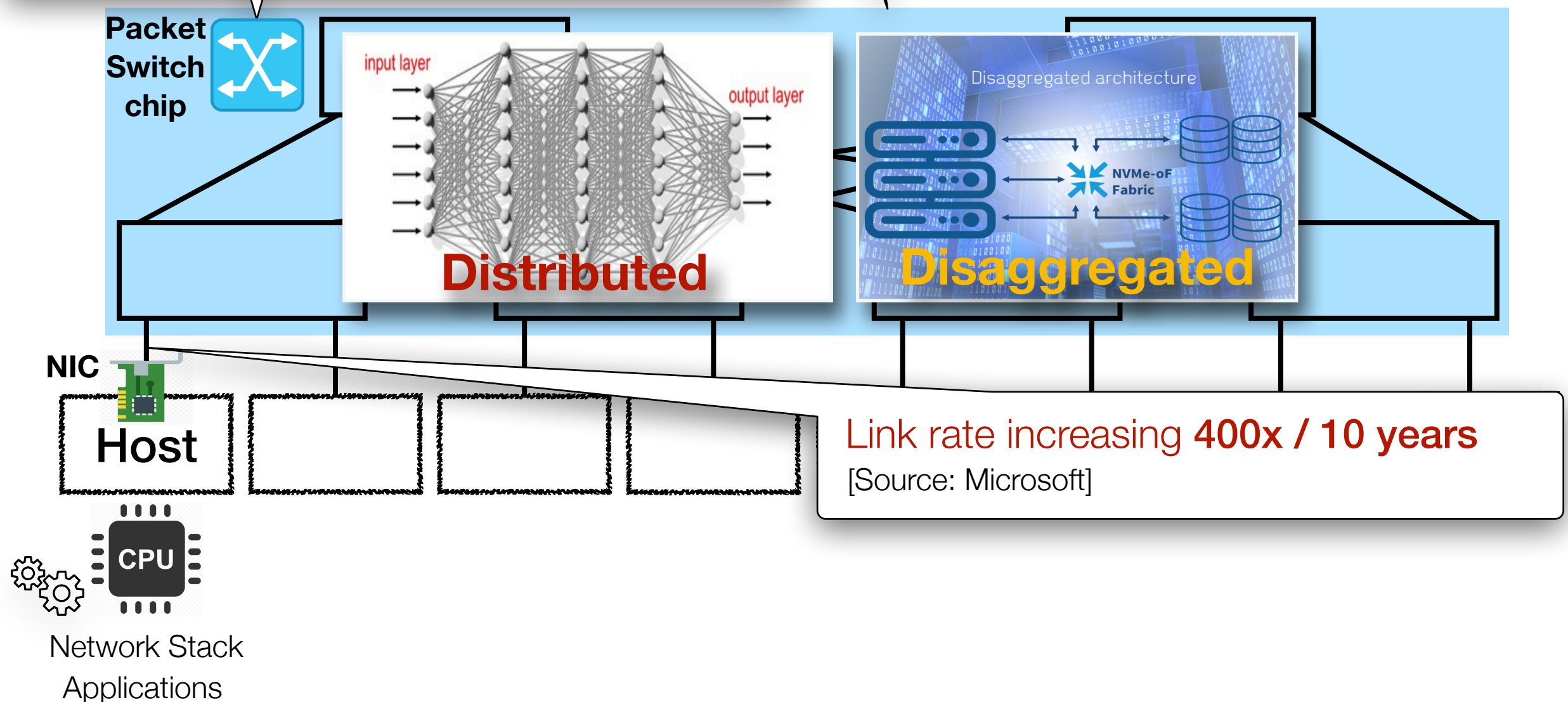
Inside a Typical Datacenter

Bandwidth demand increasing **2x / year**

[Source: Google]

Switch chip capacity increasing **2x / 2 years**

[Source: Broadcom]



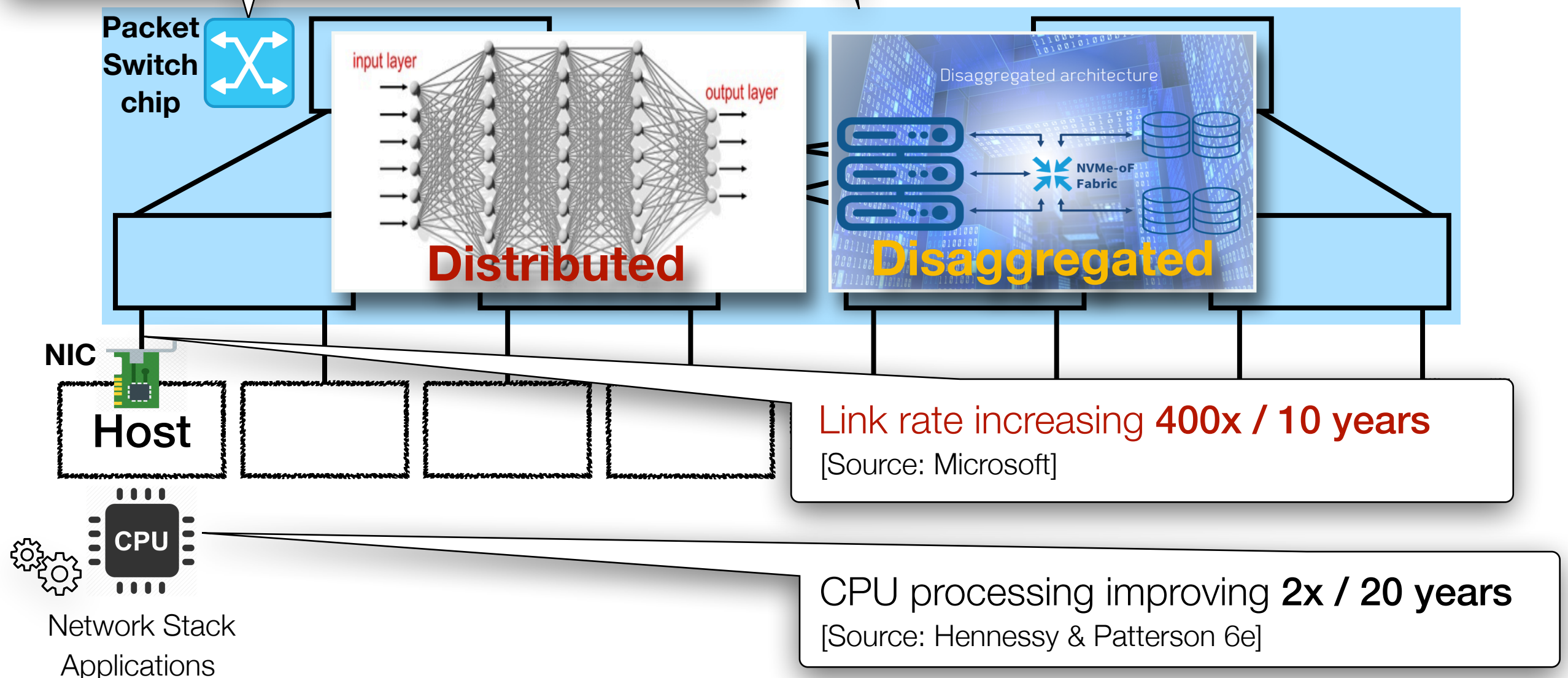
Inside a Typical Datacenter

Bandwidth demand increasing **2x / year**

[Source: Google]

Switch chip capacity increasing **2x / 2 years**

[Source: Broadcom]



Inside a Typical Datacenter

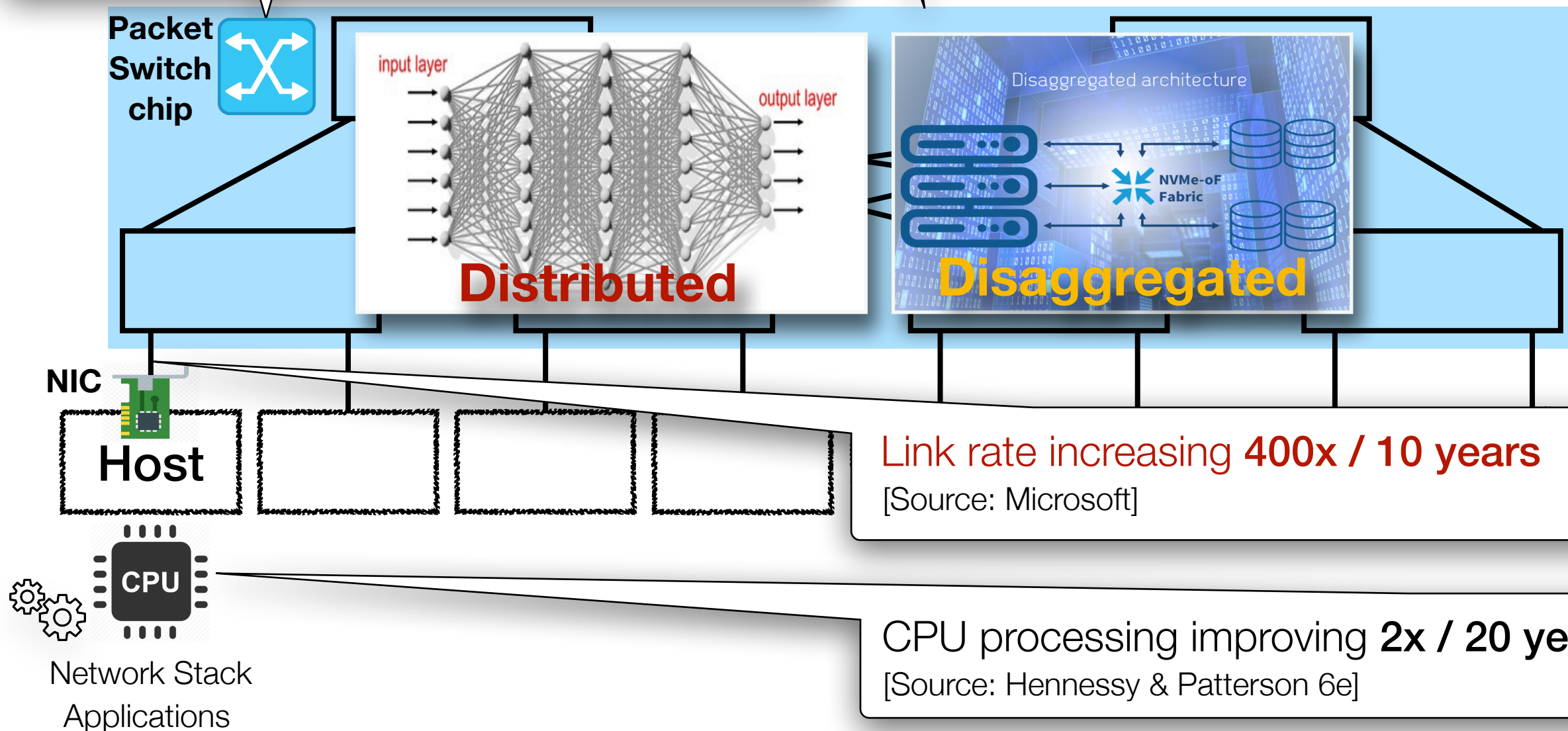
Bandwidth demand increasing **2x / year**

[Source: Google]

Switch chip capacity increasing **2x / 2 years**

[Source: Broadcom]

Slowdown in Moore's Law
End of Dennard scaling



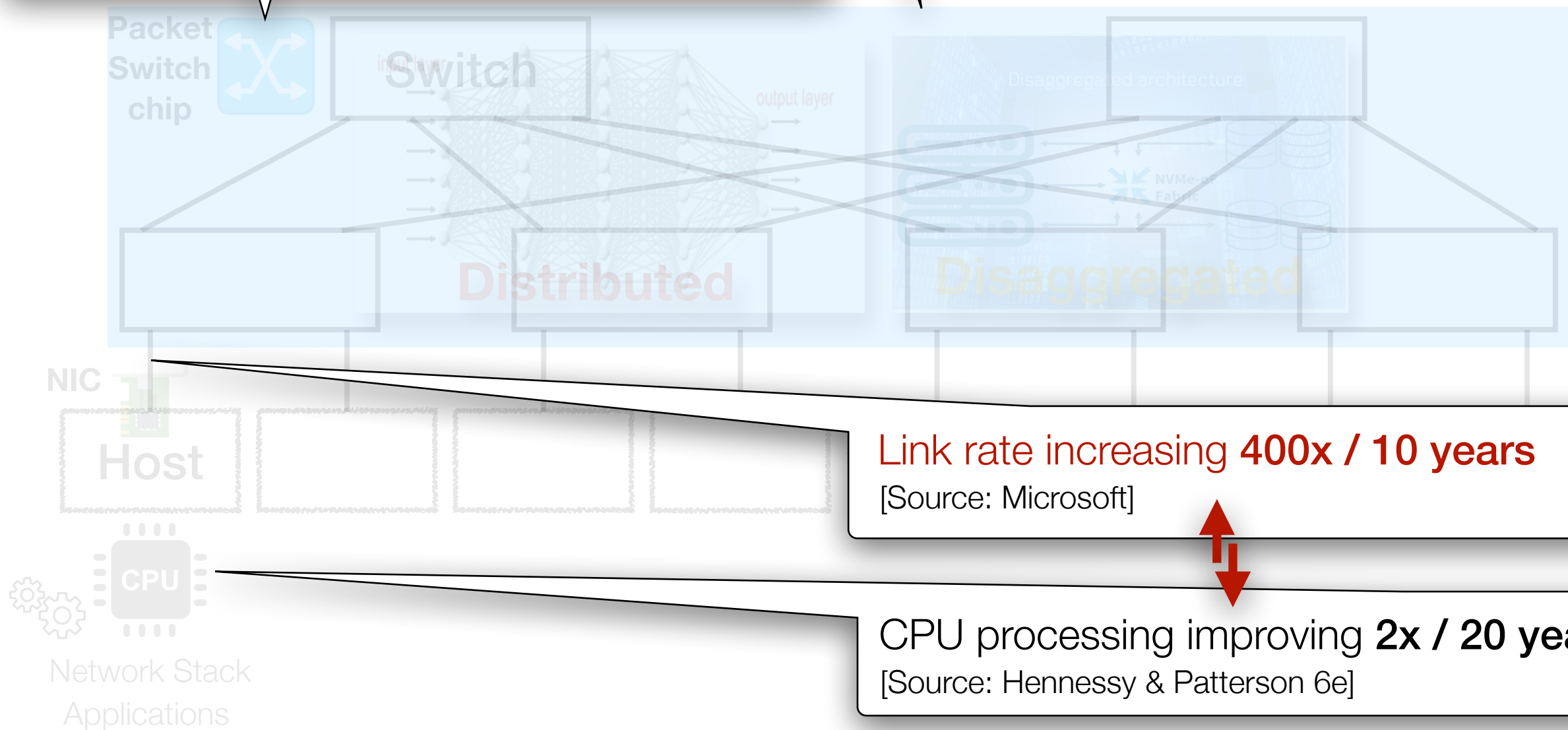
Inside a Typical Datacenter

Bandwidth demand increasing 2x / year

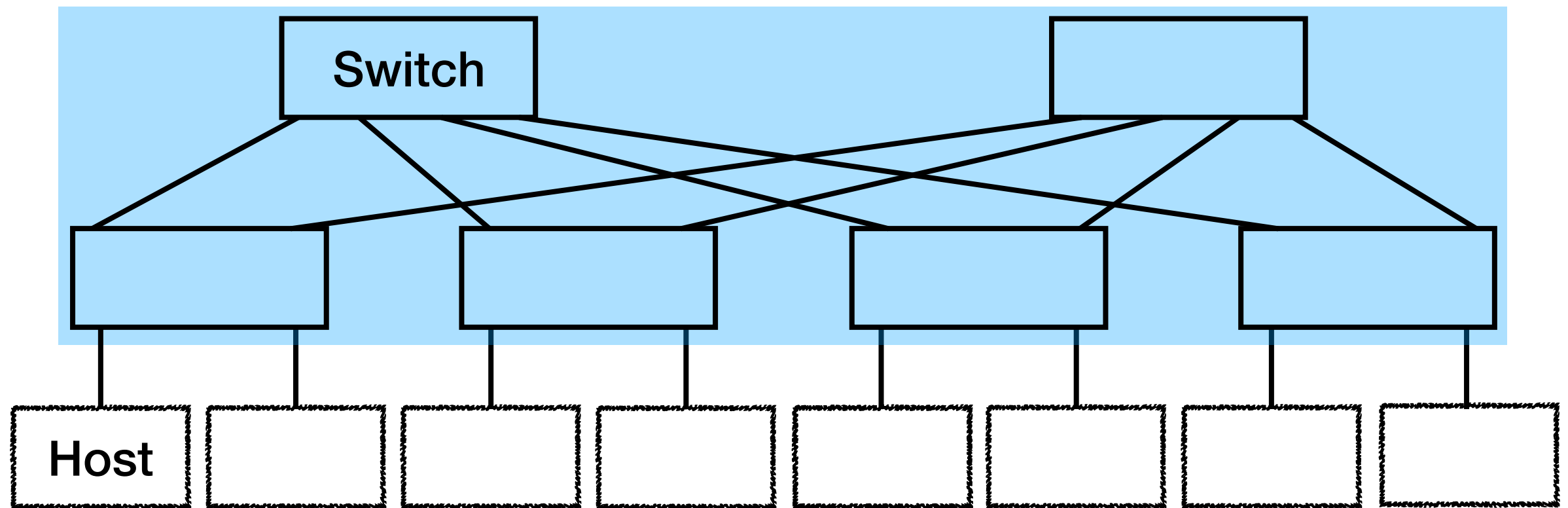
[Source: Google]

Switch chip capacity increasing 2x / 2 years

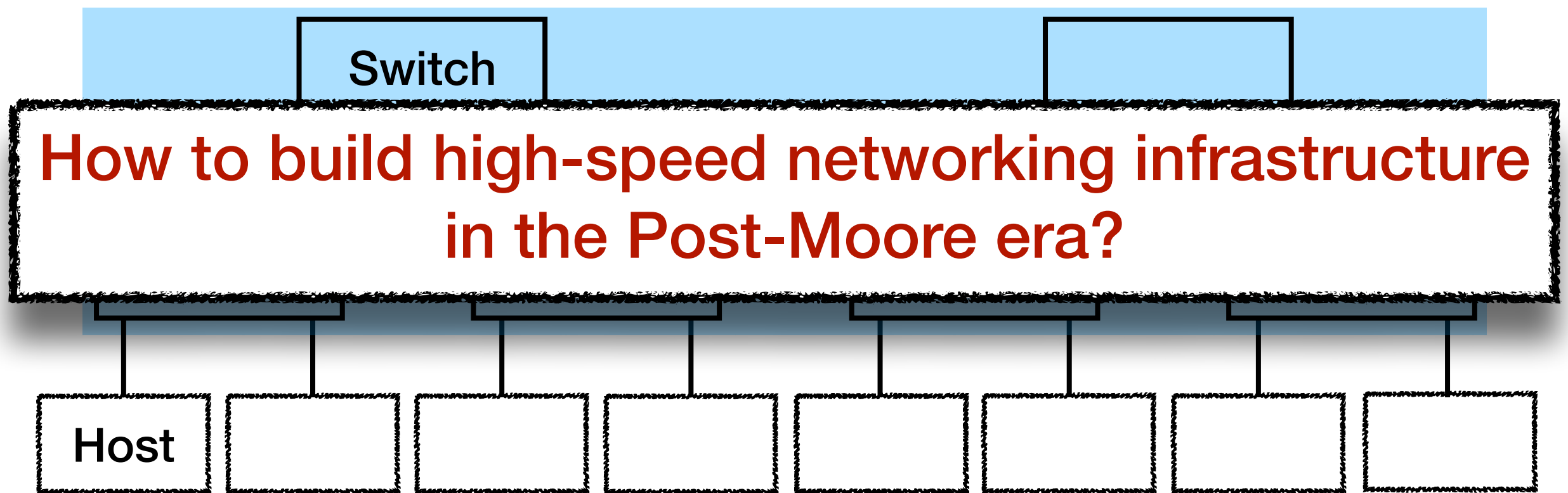
[Source: Broadcom]



Big Research Question



Big Research Question



Big Research Question

Switch

The diagram illustrates a network topology. At the top, a light blue horizontal bar contains a box labeled 'Switch'. Below this bar, a row of eight boxes is shown, each connected to the switch by a vertical line. The first box on the left is labeled 'Host', while the other seven are empty. A large, hand-drawn black box with a thick border is superimposed over the middle of the diagram, containing the text 'How to build high-speed networking infrastructure in the Post-Moore era?'. A curved arrow points from the right side of this box down to a callout box at the bottom left, which contains the text 'How to build high-speed end-host network stack?'.

How to build high-speed networking infrastructure
in the Post-Moore era?

Host

How to build high-speed end-host network stack?

Big Research Question

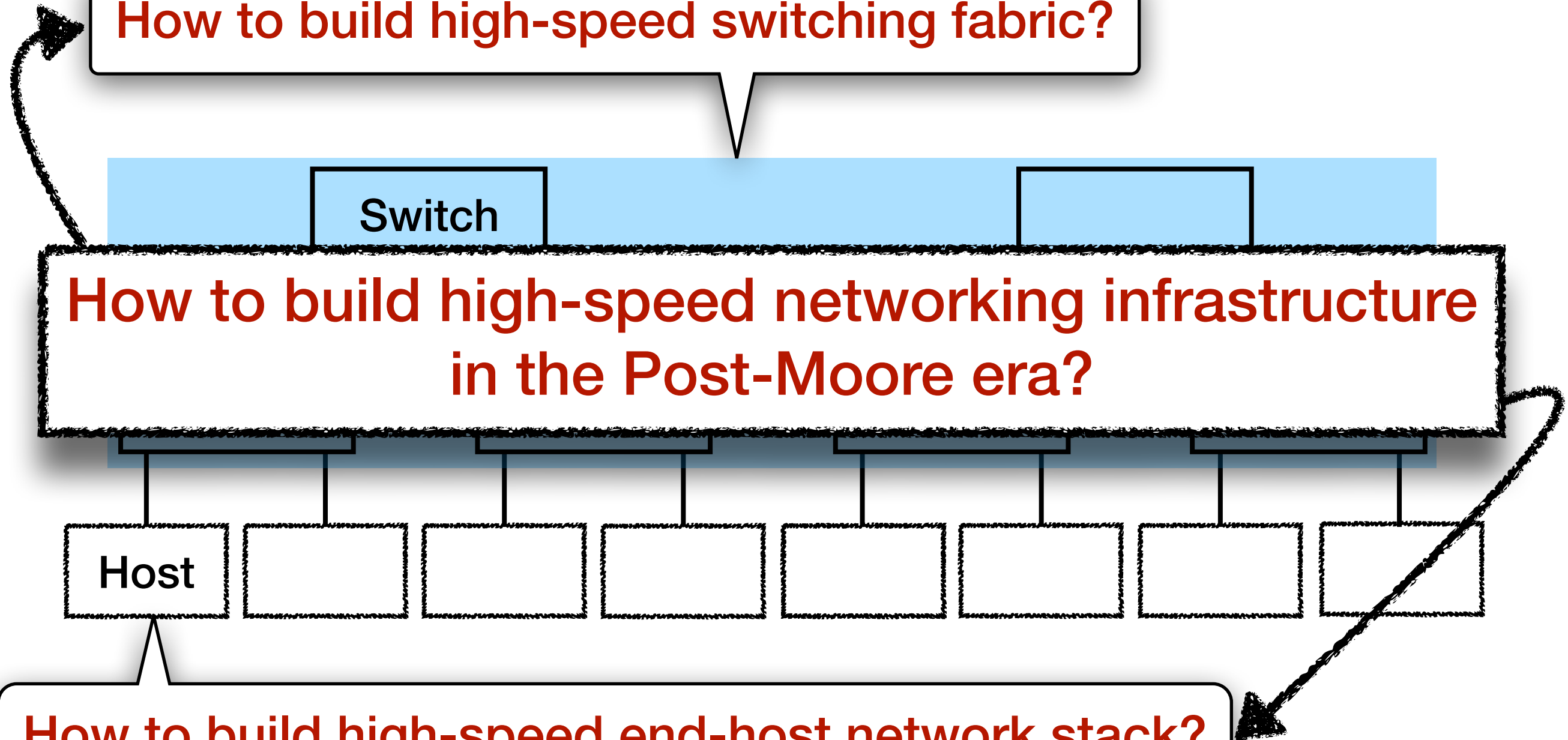
How to build high-speed switching fabric?

Switch

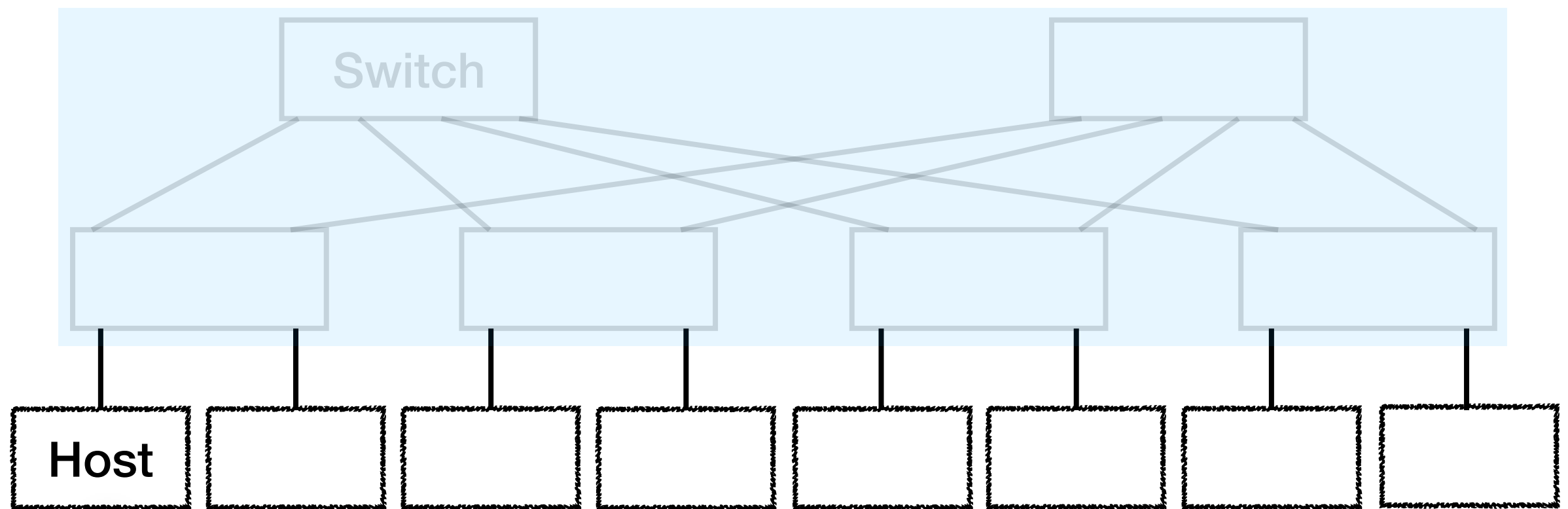
How to build high-speed networking infrastructure
in the Post-Moore era?

Host

How to build high-speed end-host network stack?

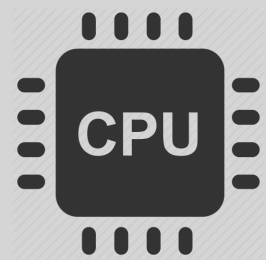


Part I : End-host Network Stack



How to build high-speed end-host network stack?

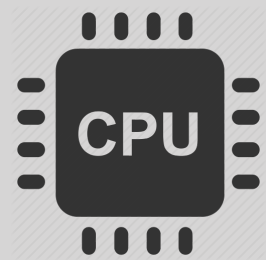
Part I : End-host Network Stack



Host

How to build high-speed end-host network stack?

Part I : End-host Network Stack



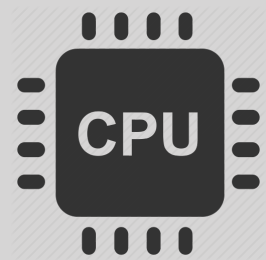
+

Domain-specific Processors



How to build high-speed end-host network stack?

Part I : End-host Network Stack



+

Domain-specific Processors

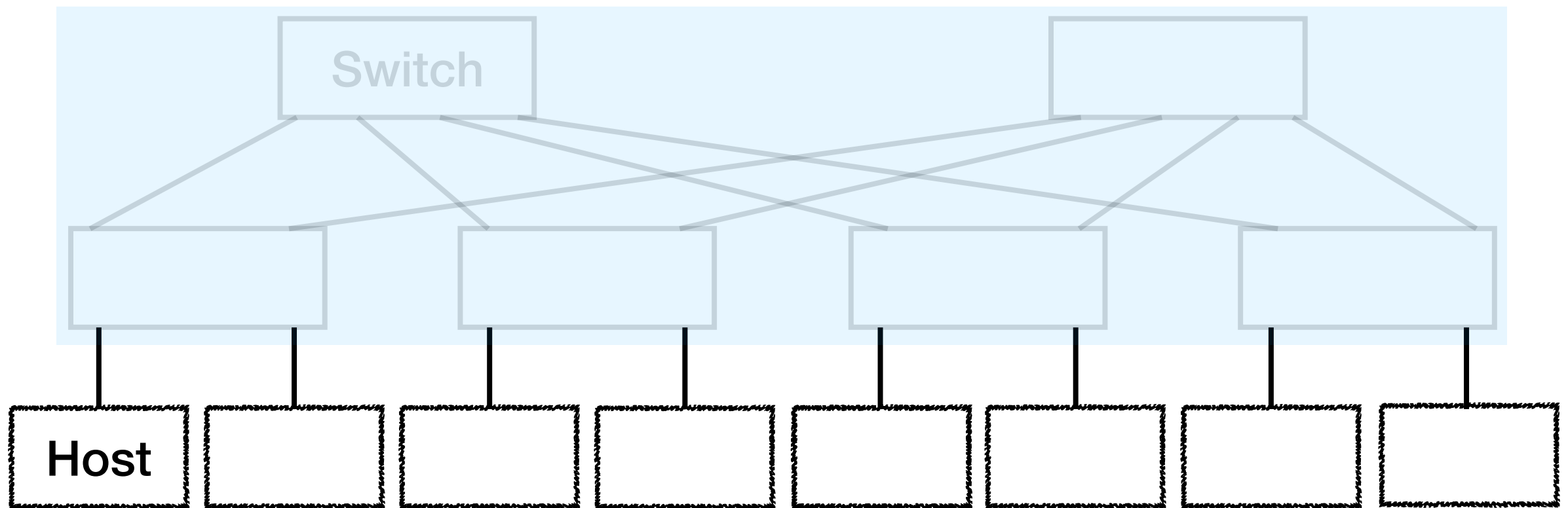
Programmability  **Performance**



How to build high-speed end-host network stack?

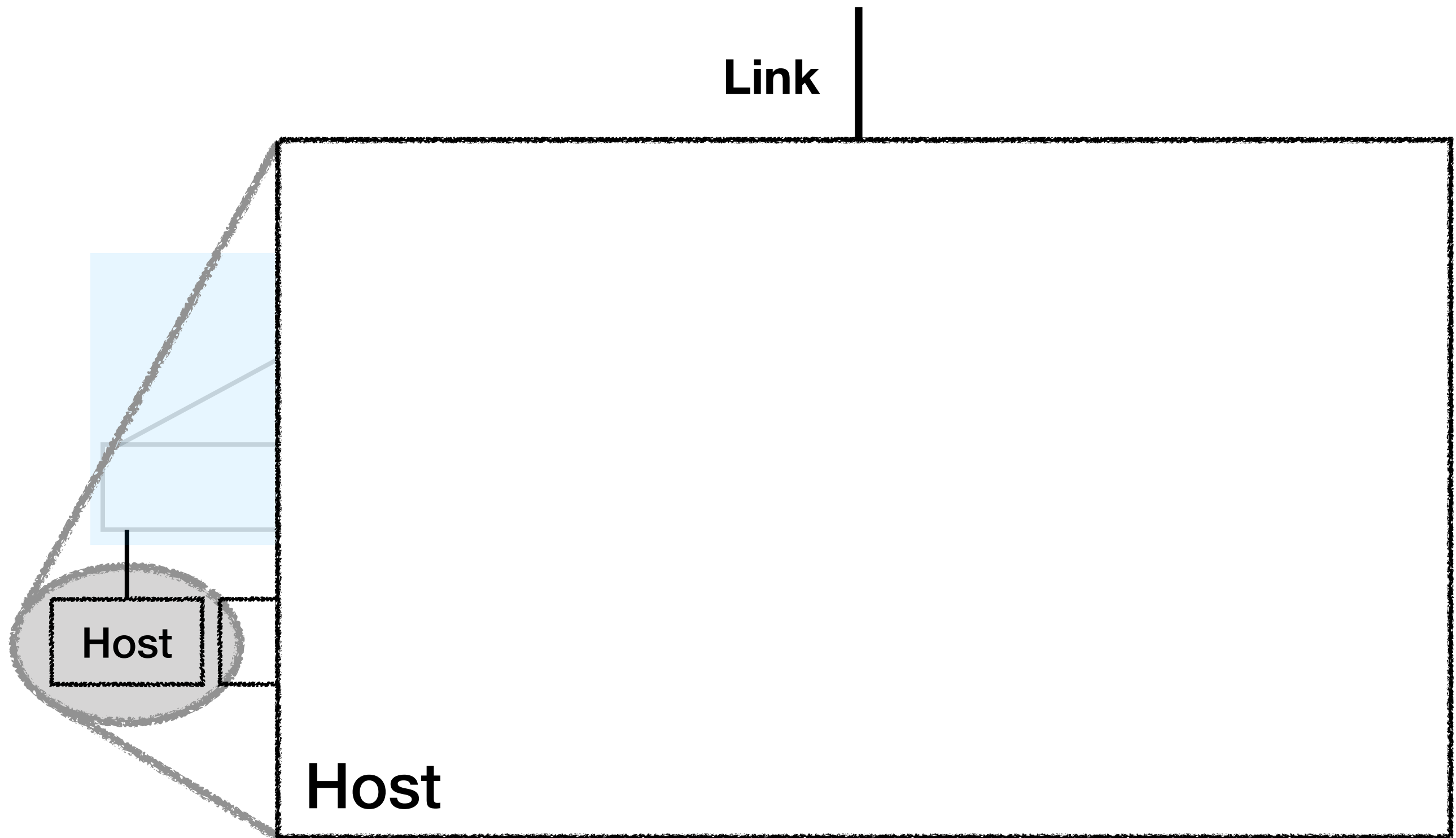
End-host Network Stack Operation :

Packet Scheduling



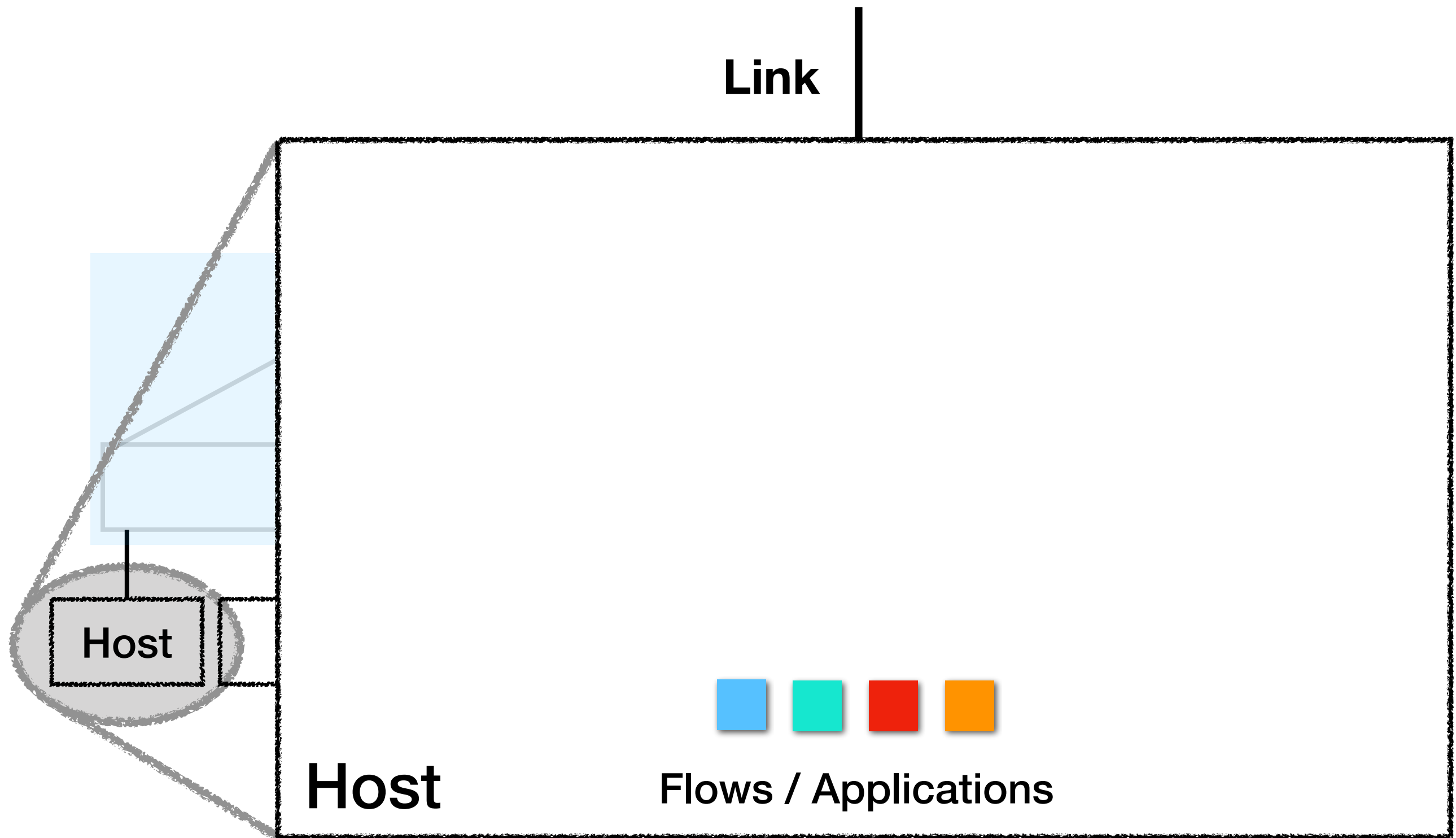
End-host Network Stack Operation :

Packet Scheduling



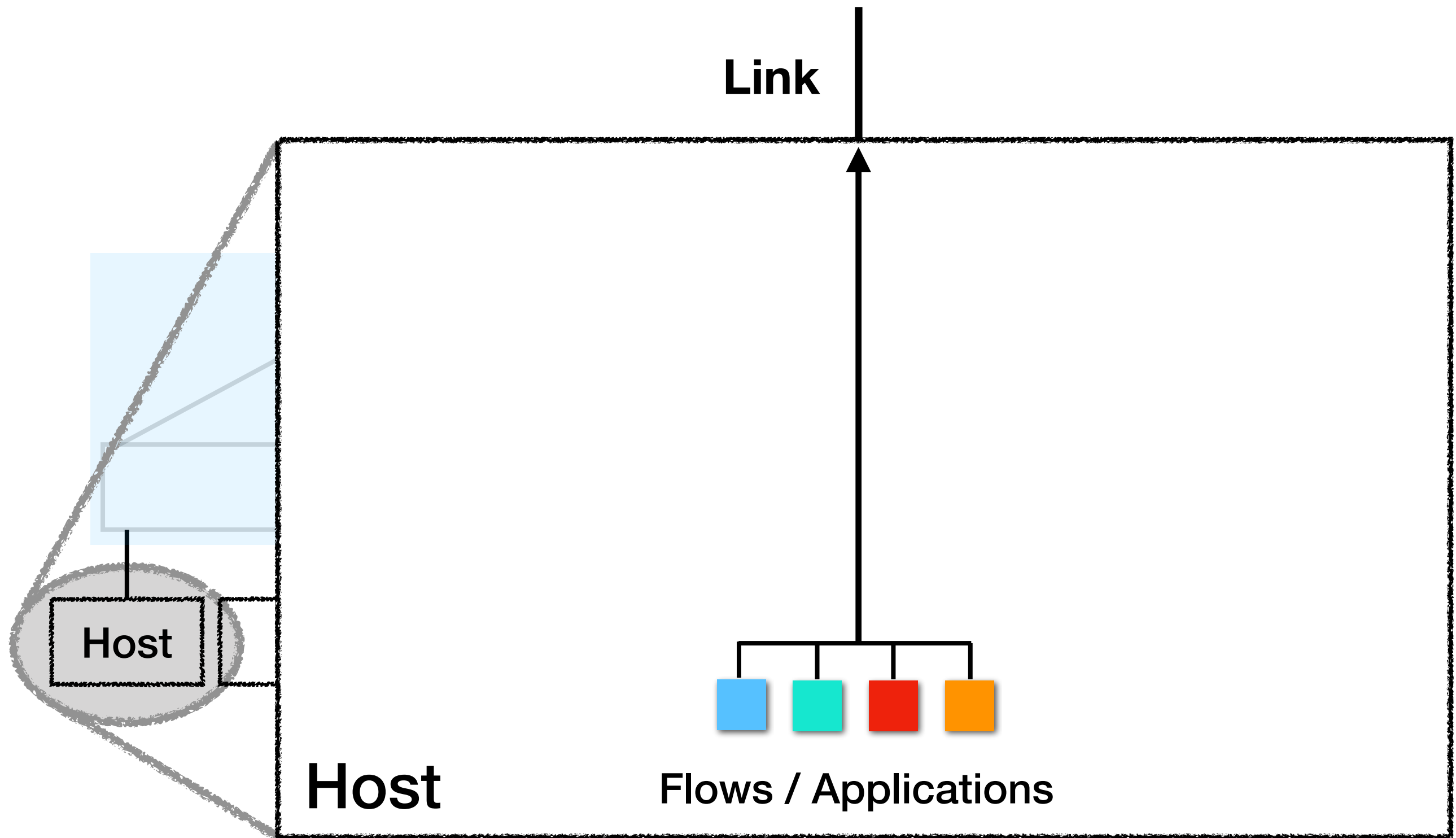
End-host Network Stack Operation :

Packet Scheduling



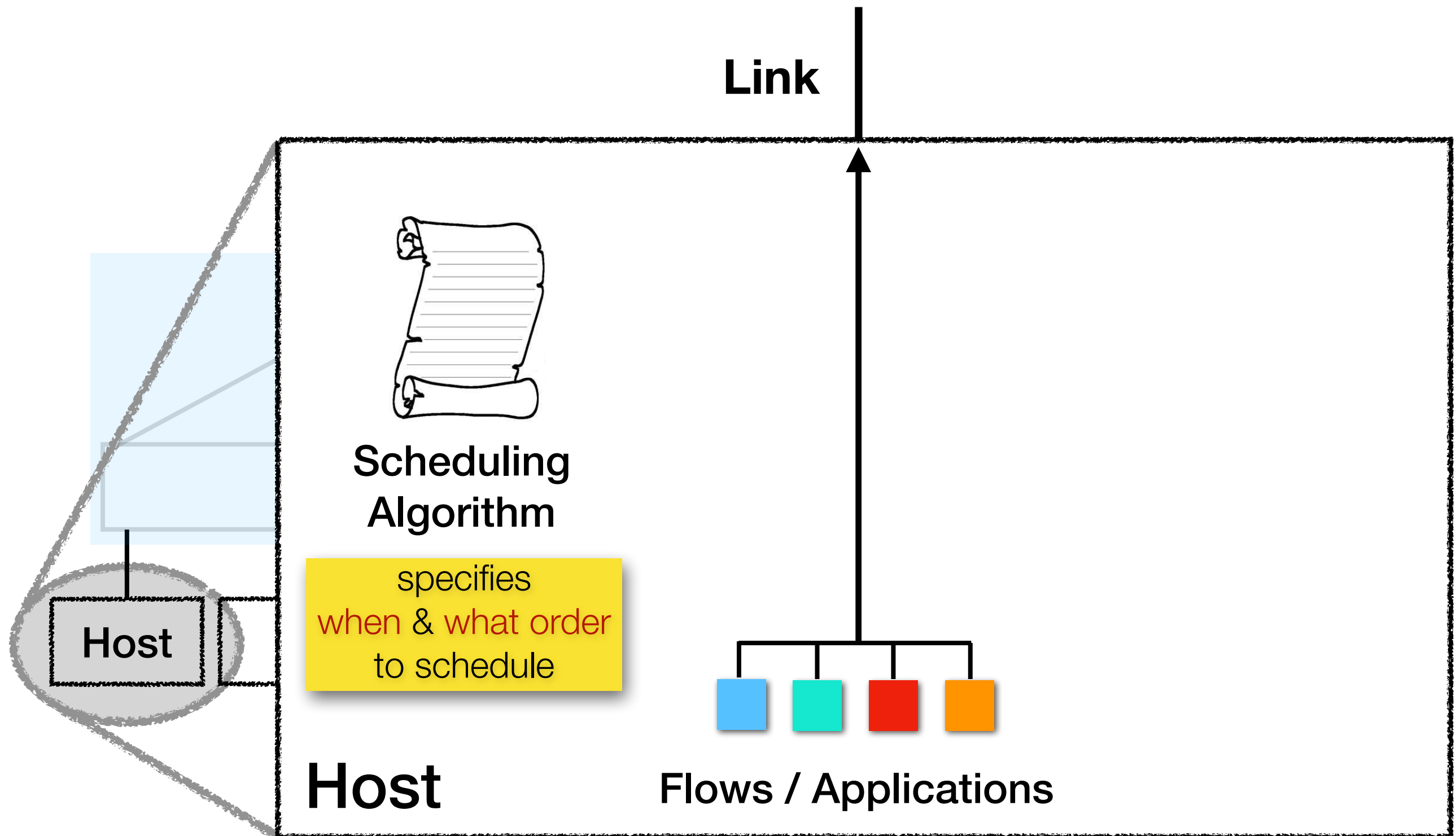
End-host Network Stack Operation :

Packet Scheduling

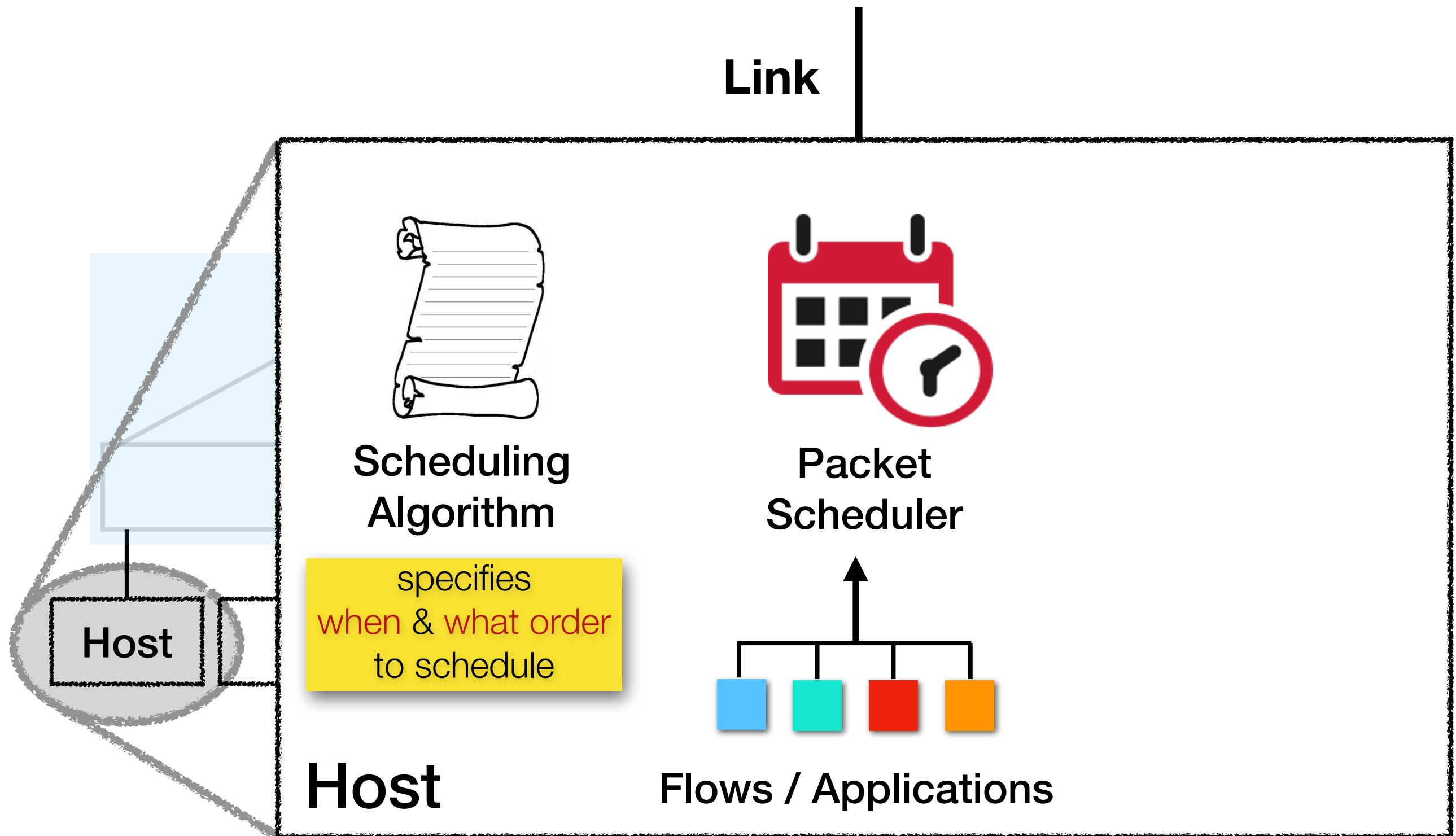


End-host Network Stack Operation :

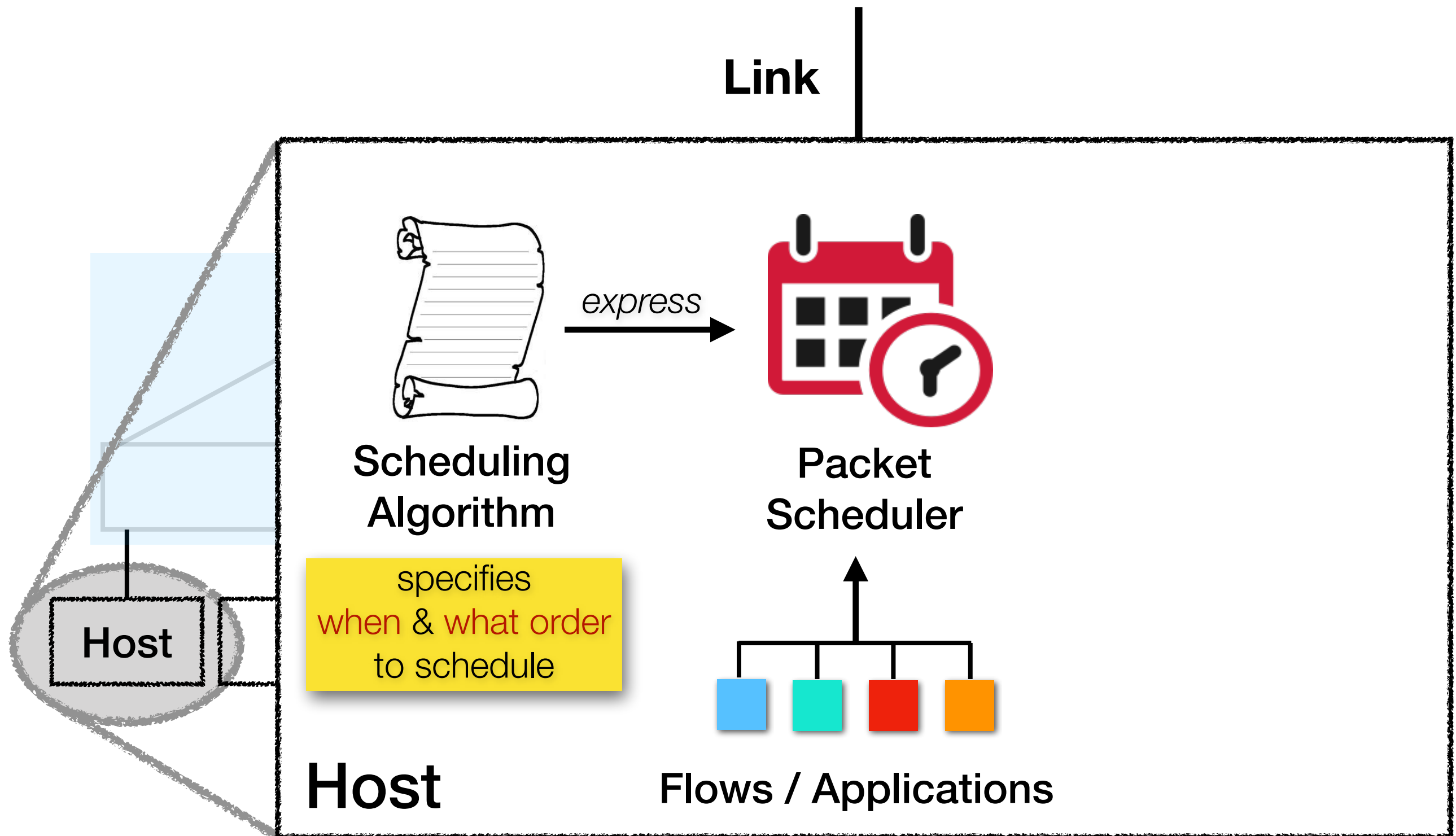
Packet Scheduling



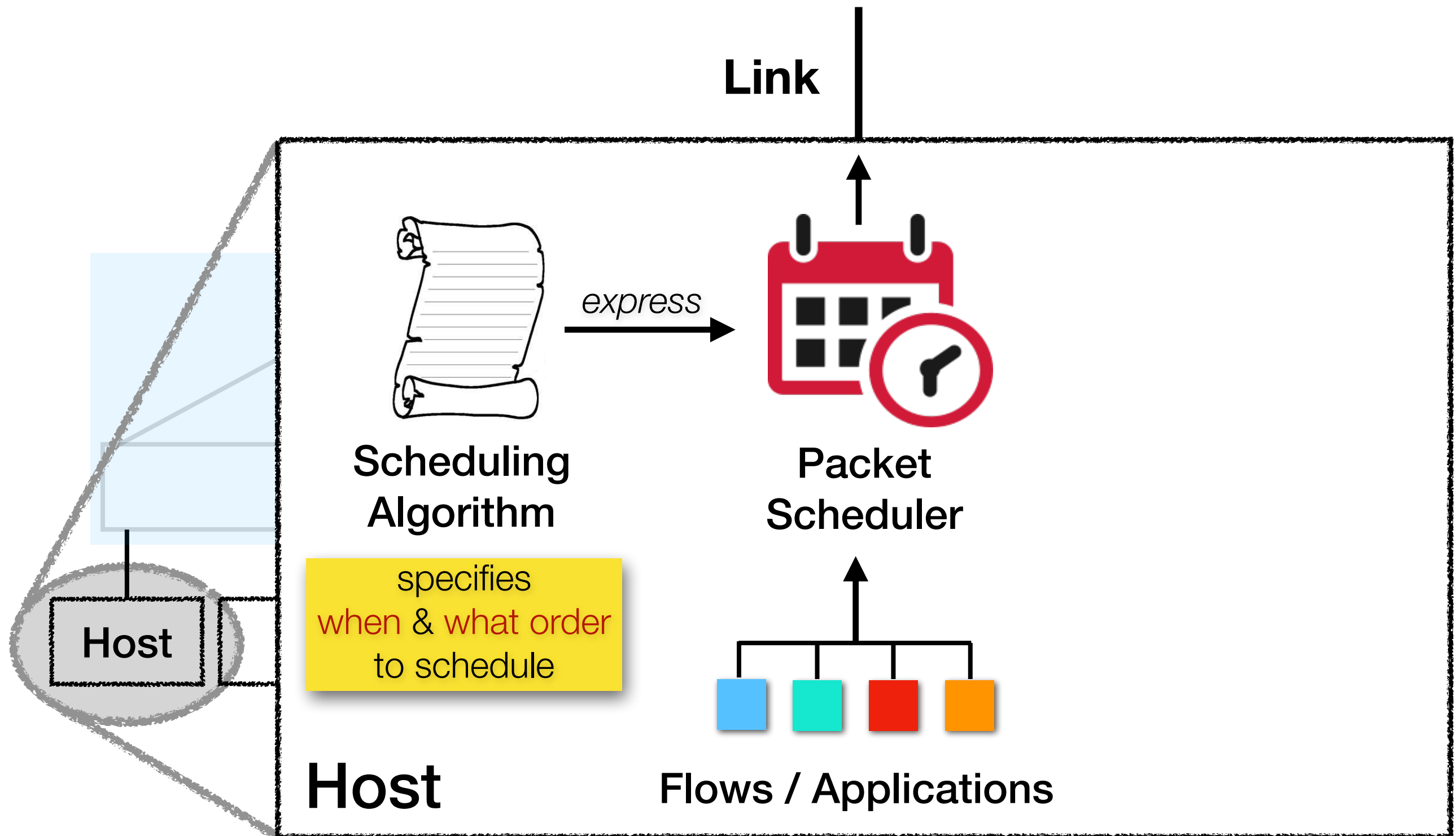
End-host Network Stack Operation : Packet Scheduling



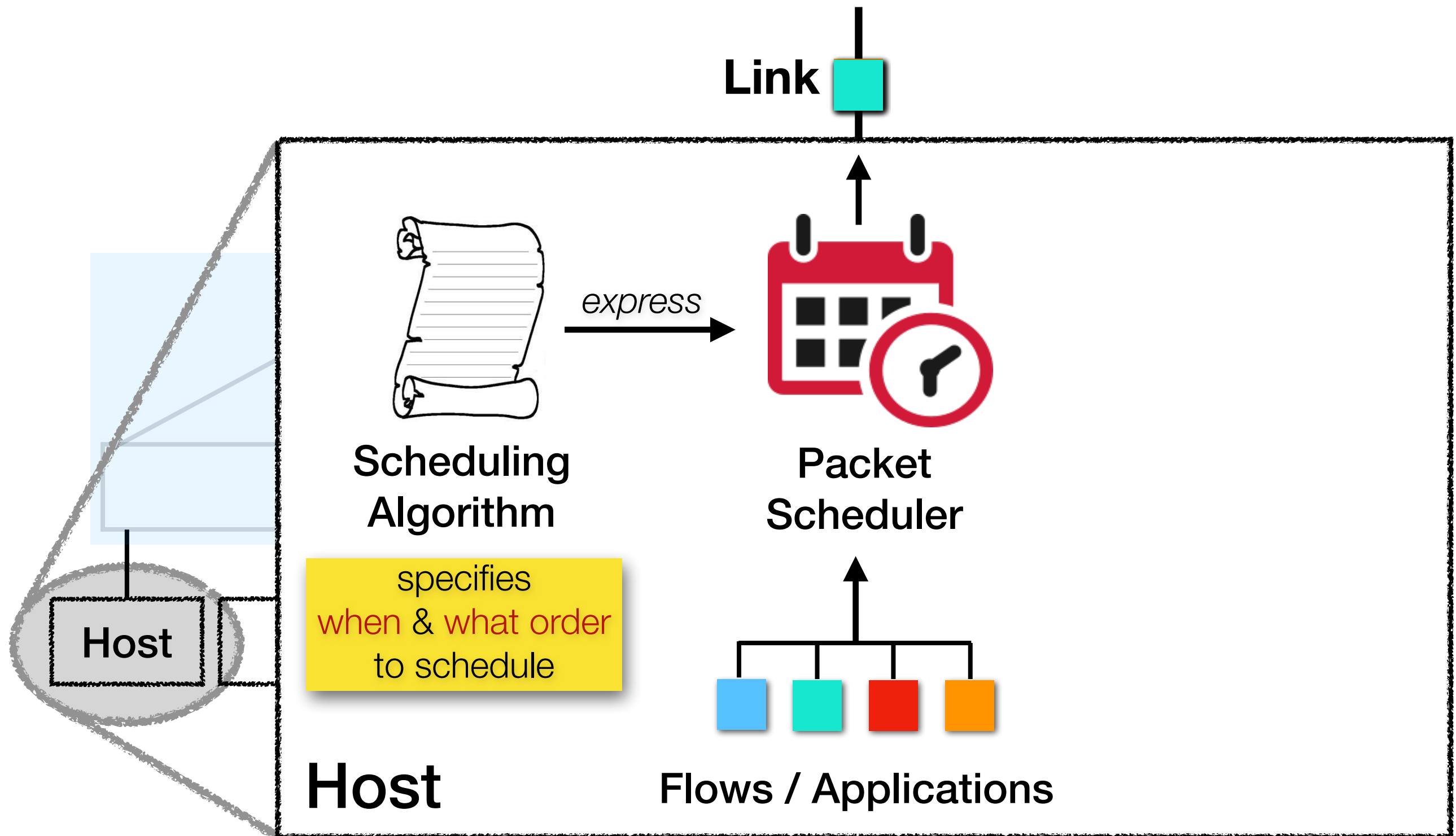
End-host Network Stack Operation : Packet Scheduling



End-host Network Stack Operation : Packet Scheduling

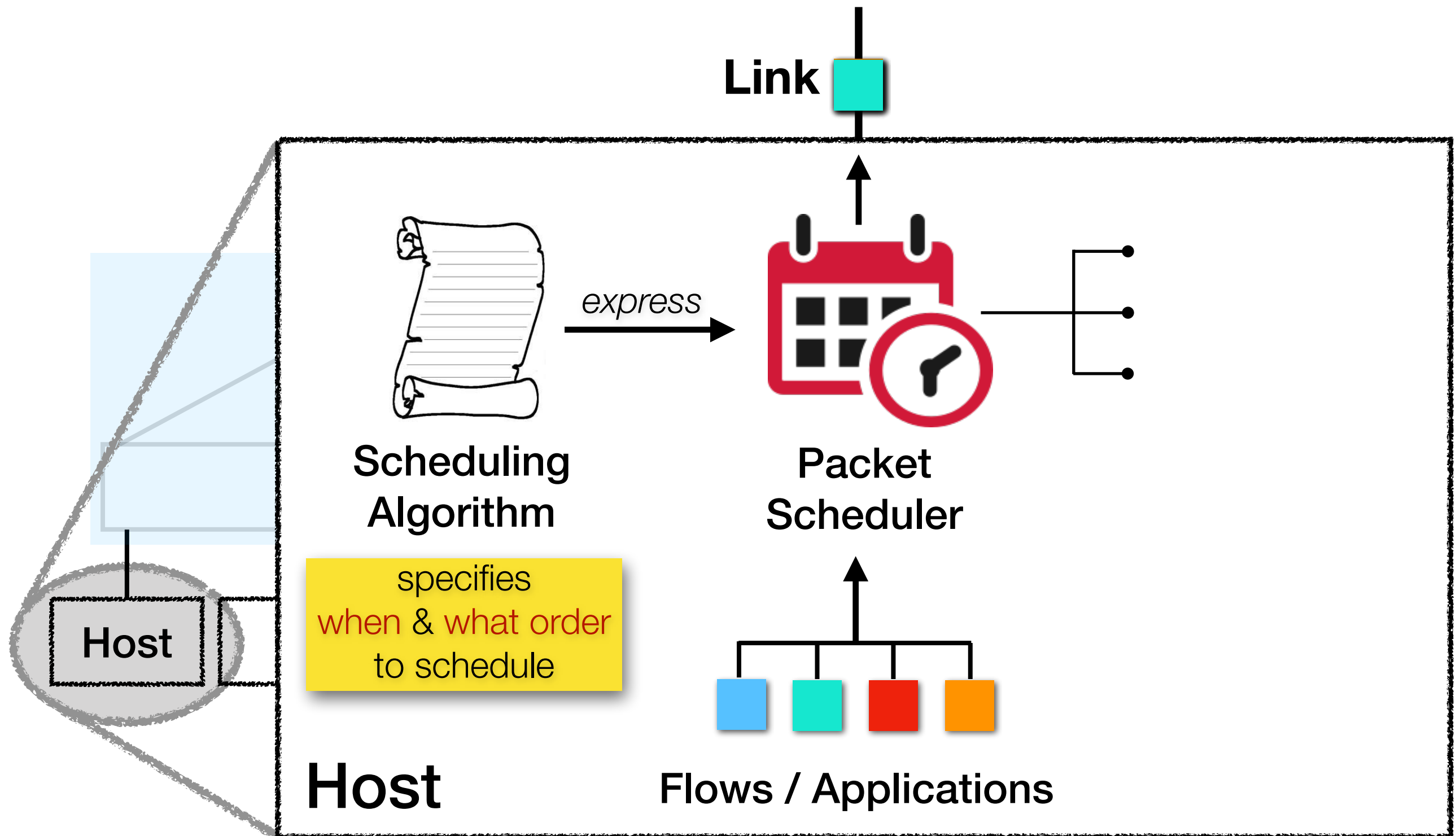


End-host Network Stack Operation : Packet Scheduling

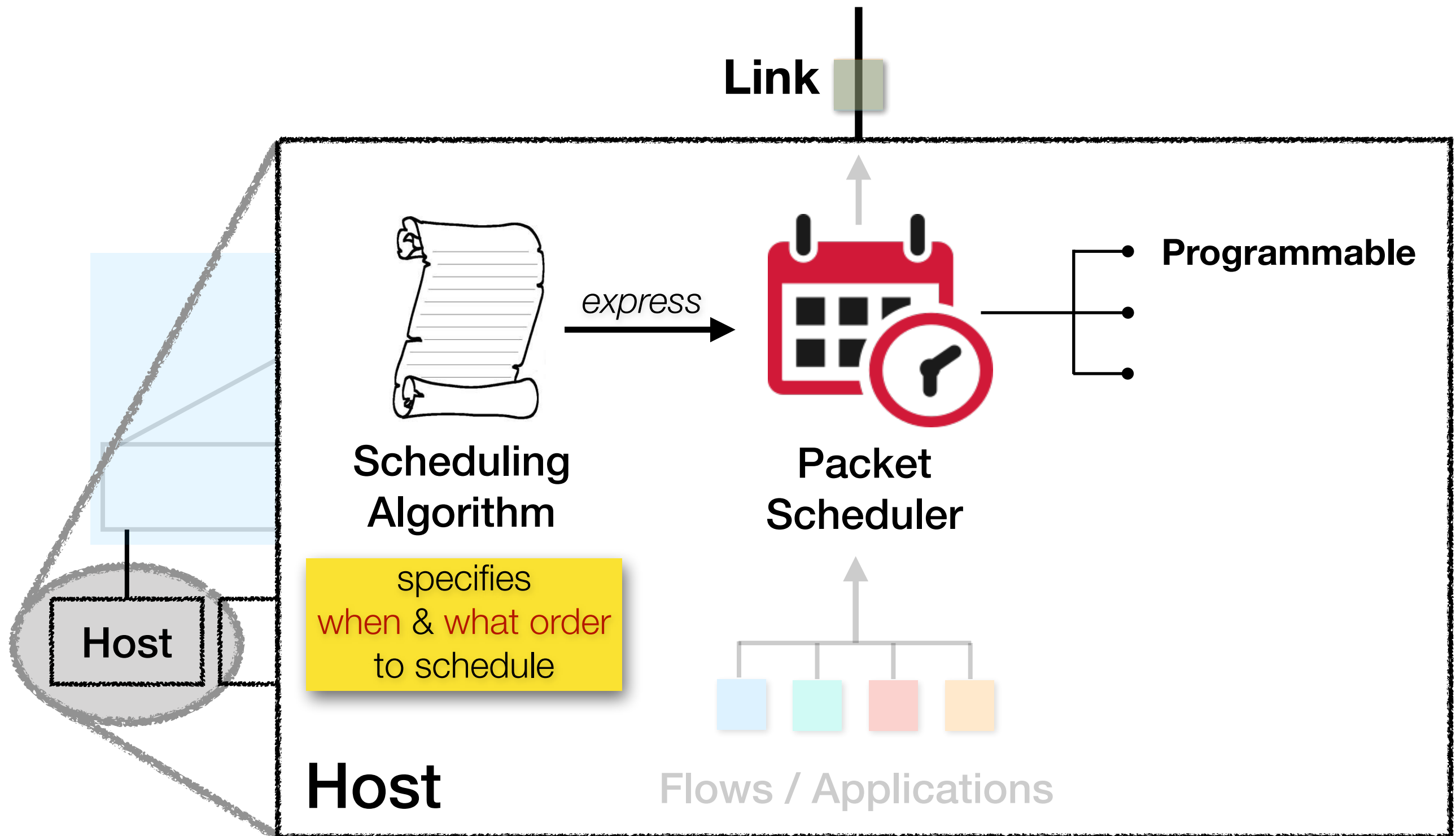


End-host Network Stack Operation :

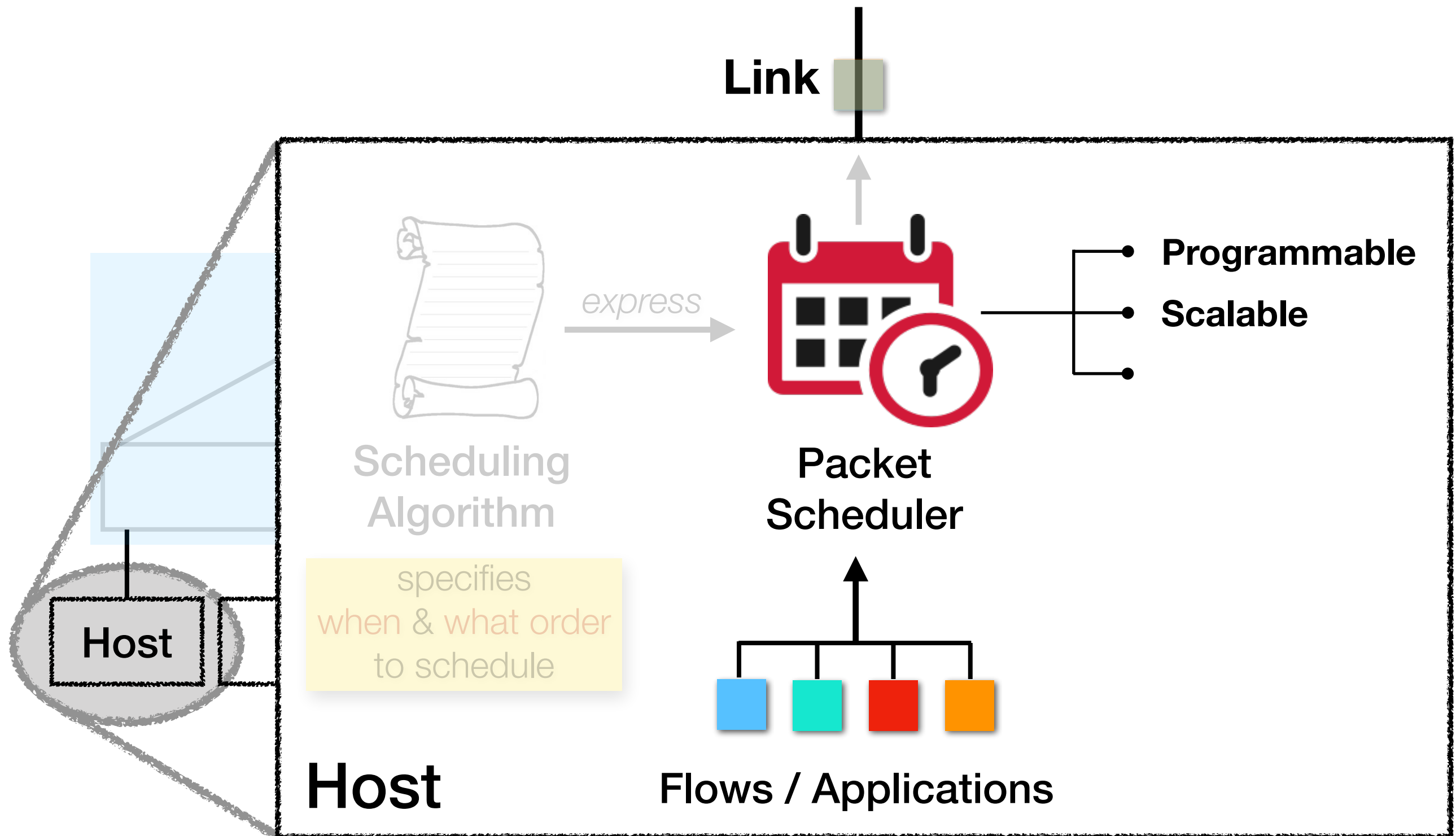
Packet Scheduling



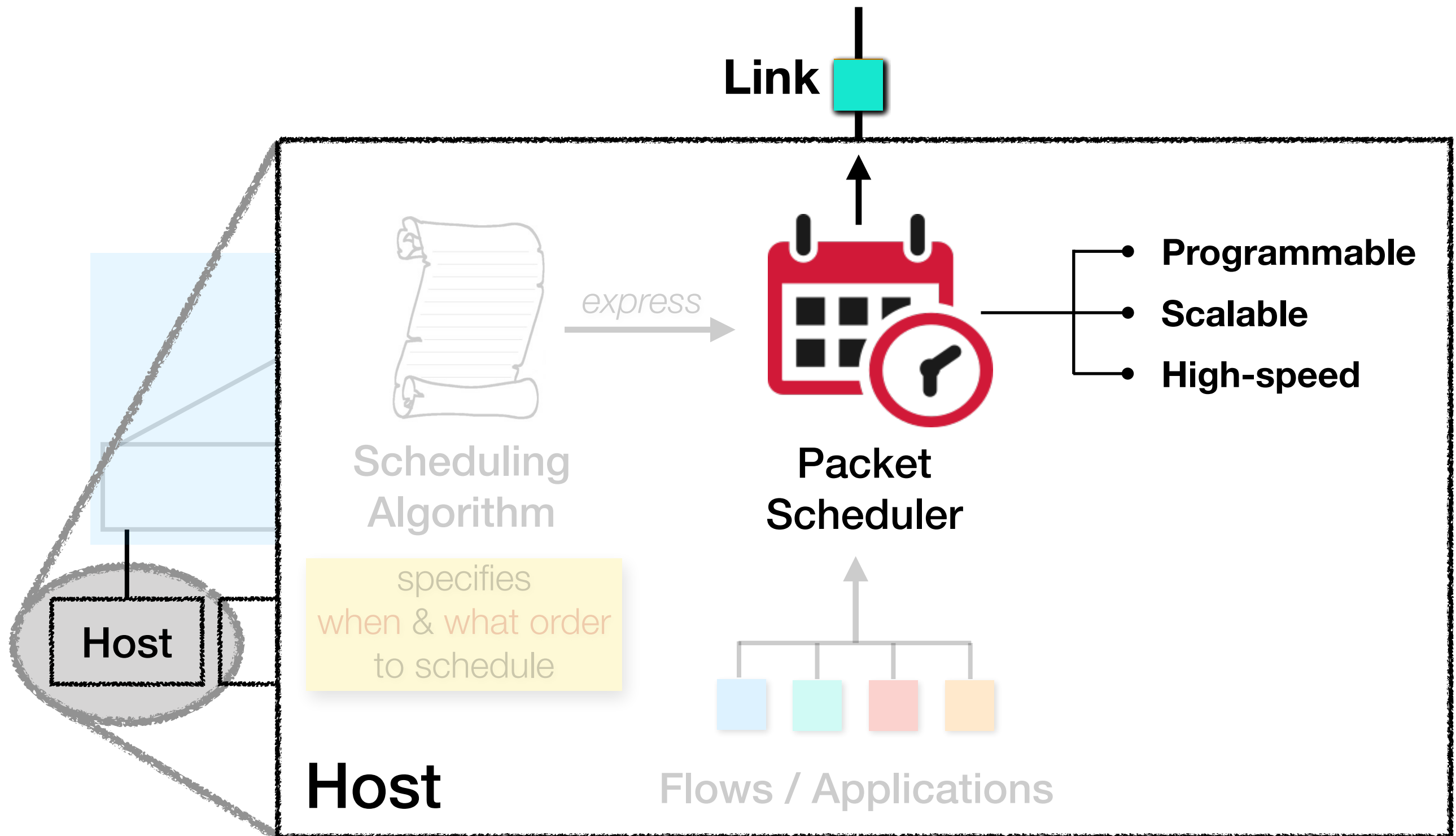
End-host Network Stack Operation : Packet Scheduling



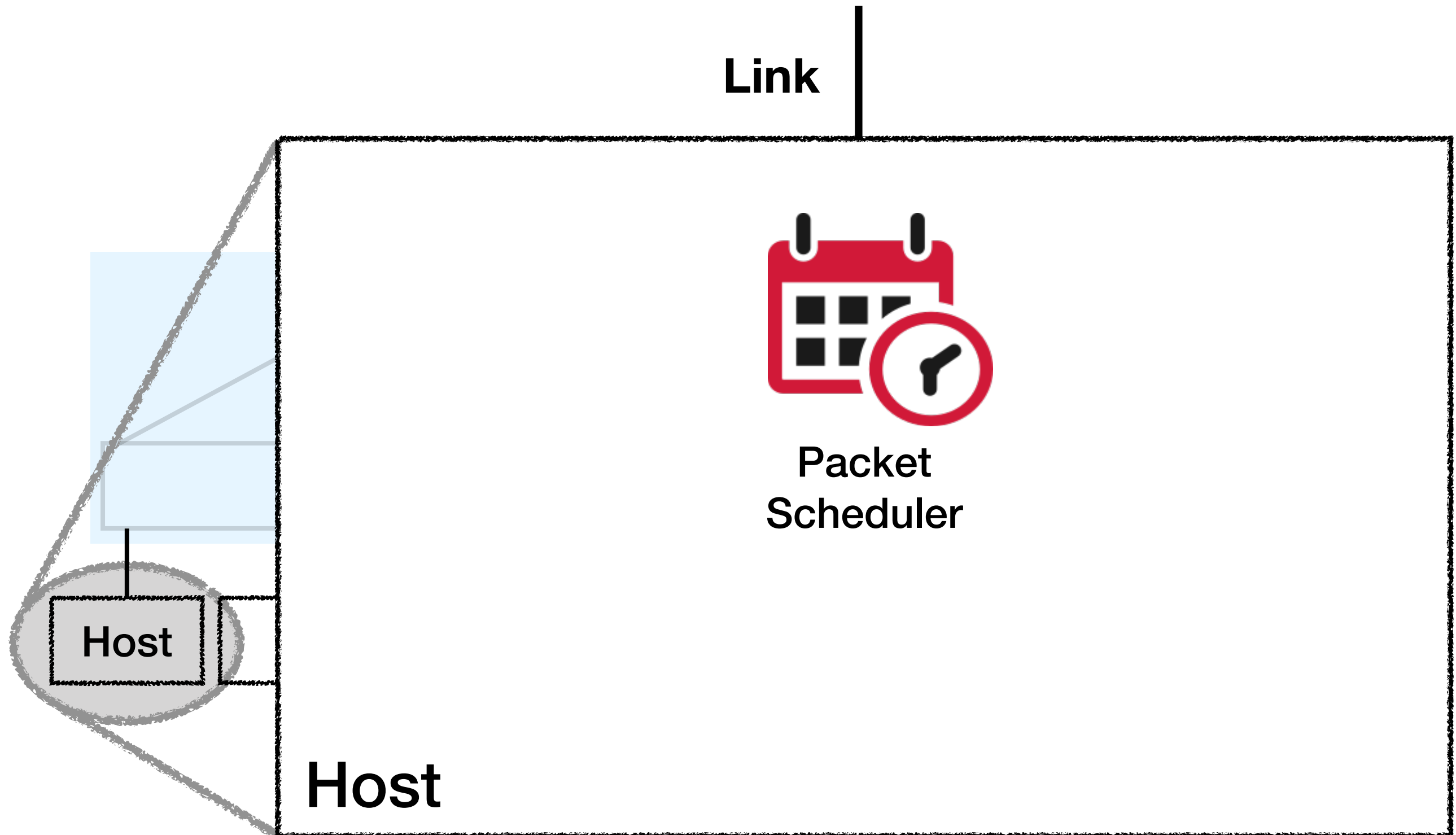
End-host Network Stack Operation : Packet Scheduling



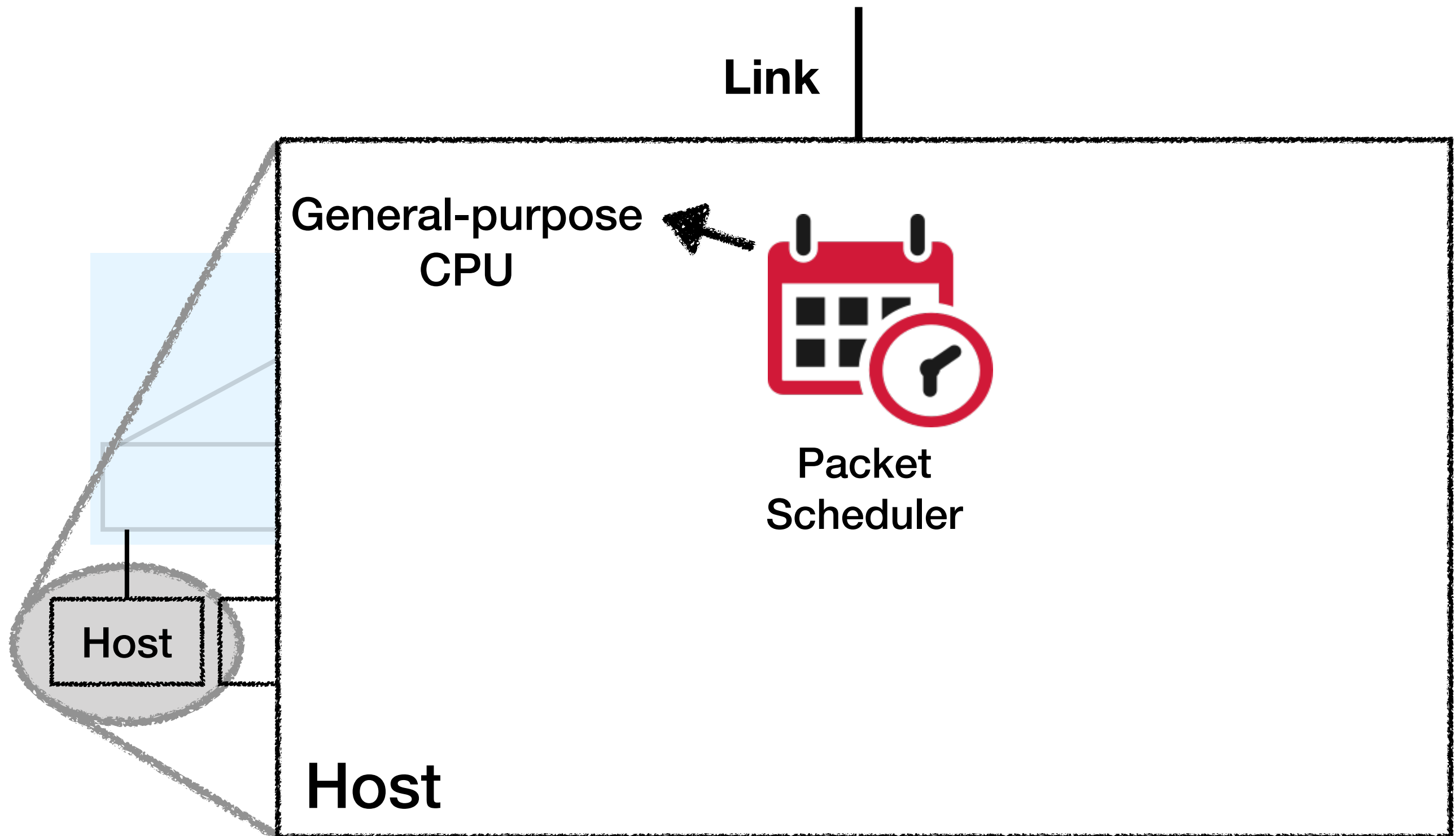
End-host Network Stack Operation : Packet Scheduling



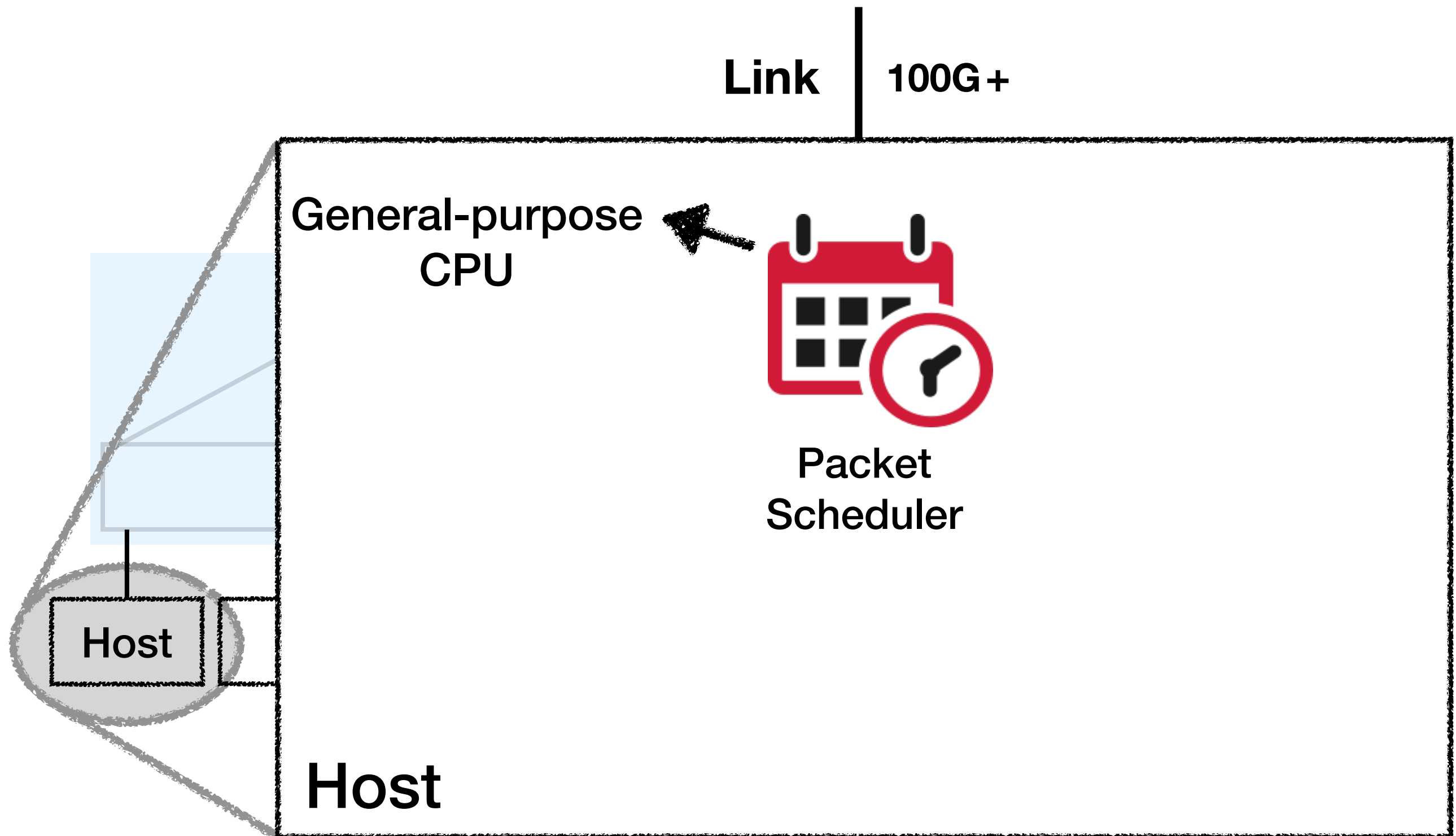
State-of-the-Art Packet Schedulers



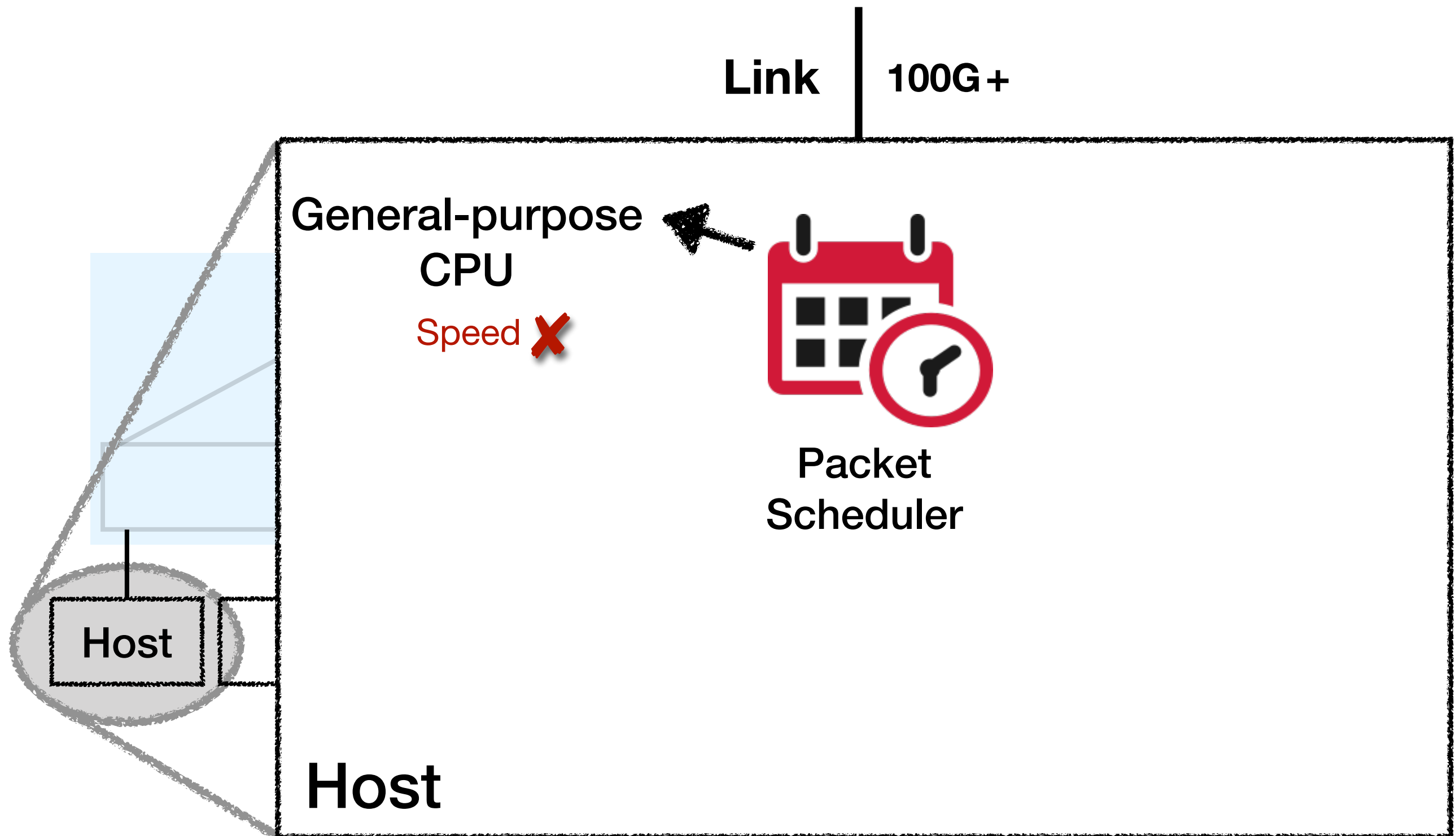
State-of-the-Art Packet Schedulers



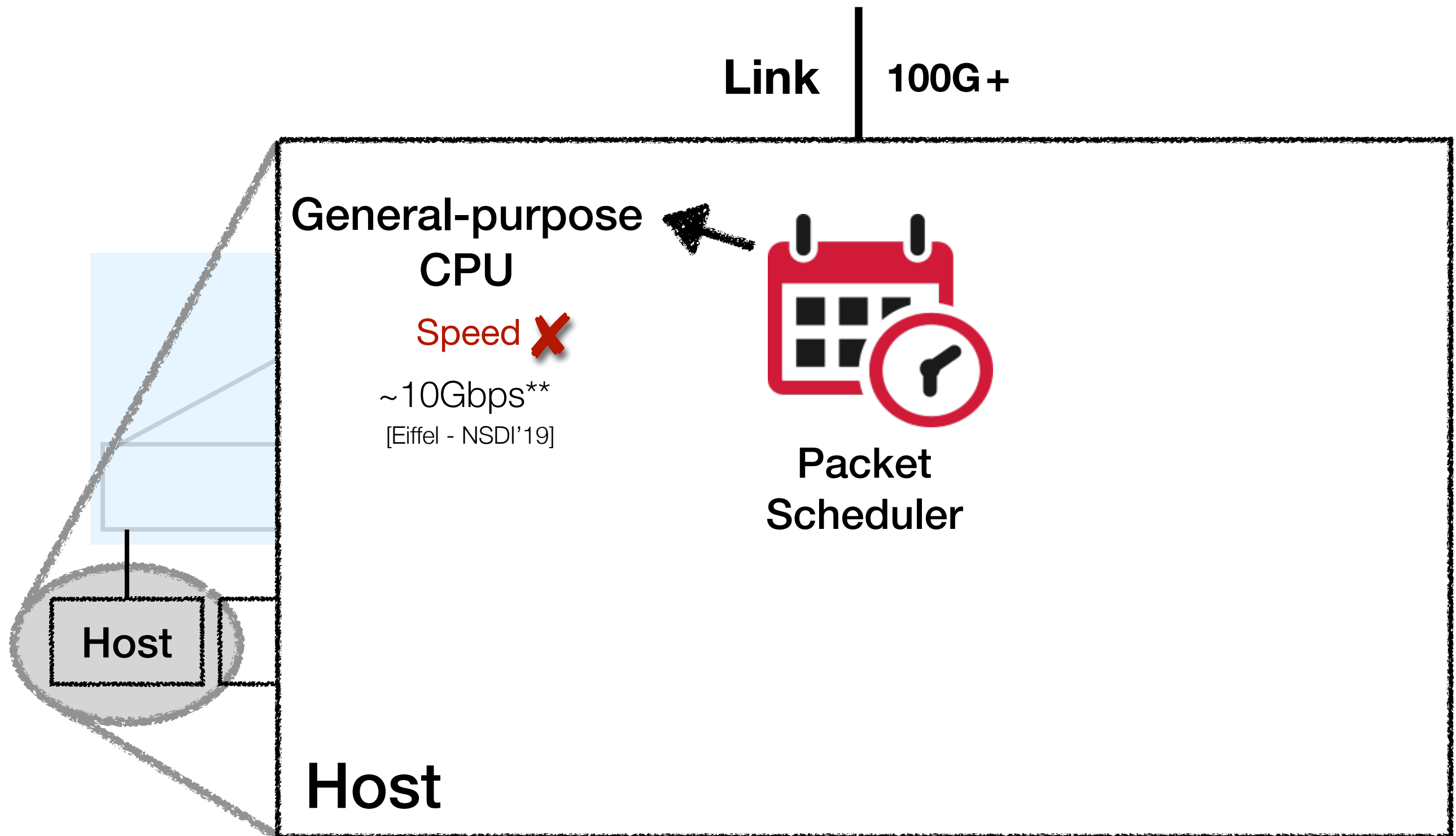
State-of-the-Art Packet Schedulers



State-of-the-Art Packet Schedulers

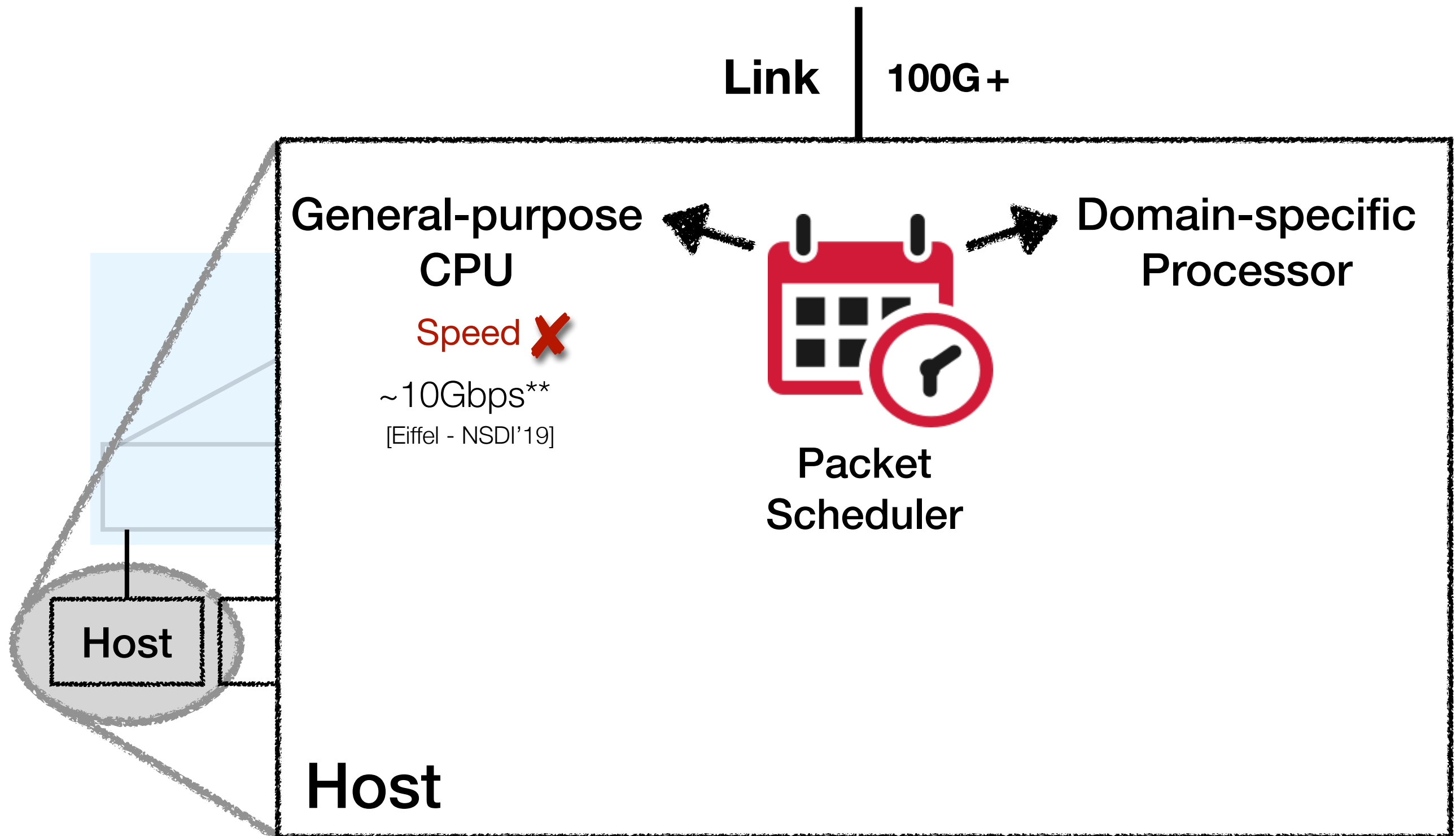


State-of-the-Art Packet Schedulers



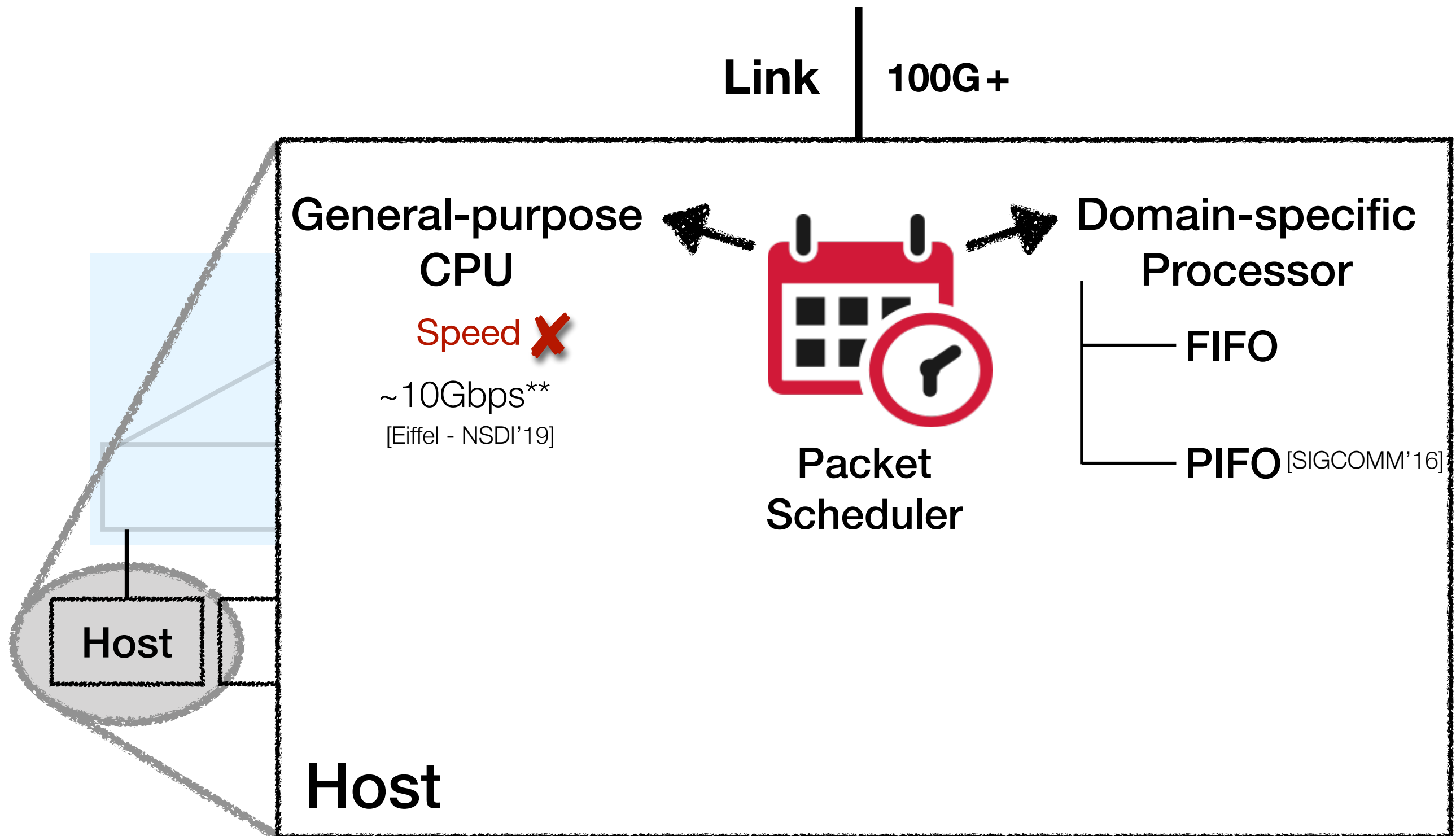
**1 core, 1500B packets, running variant of SJF algorithm

State-of-the-Art Packet Schedulers



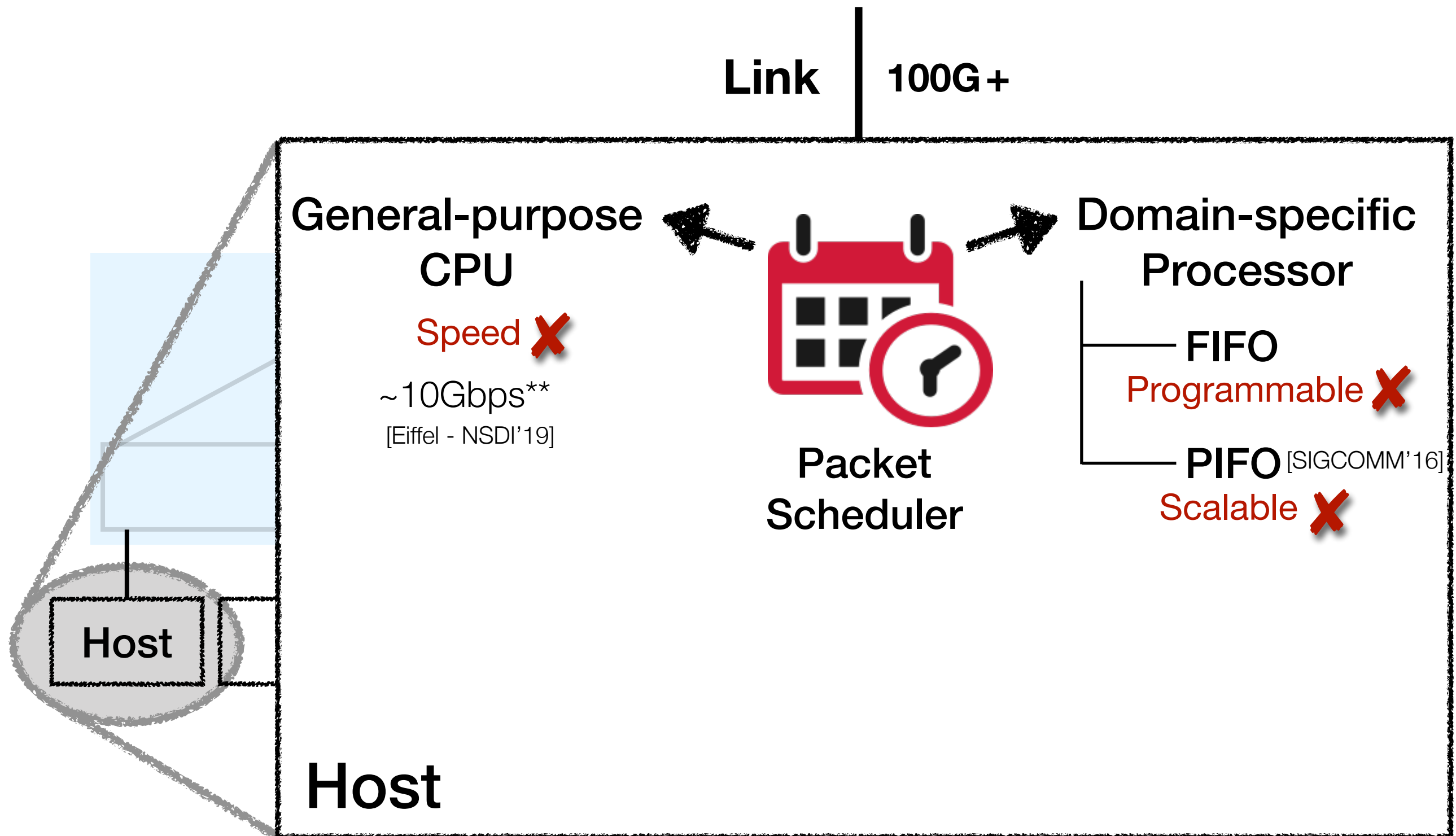
**1 core, 1500B packets, running variant of SJF algorithm

State-of-the-Art Packet Schedulers



**1 core, 1500B packets, running variant of SJF algorithm

State-of-the-Art Packet Schedulers



**1 core, 1500B packets, running variant of SJF algorithm

How to build a packet scheduler that is
simultaneously

Programmable, Scalable, High-speed?

|
express wide range
of packet scheduling
algorithms

|
schedule
10s of thousands
of flows

ref. [SENIC - NSDI'14]
[Carousel - SIGCOMM'17]

|
schedule
at line rate
100G+

How to build a packet scheduler that is simultaneously Programmable, Scalable, High-speed?

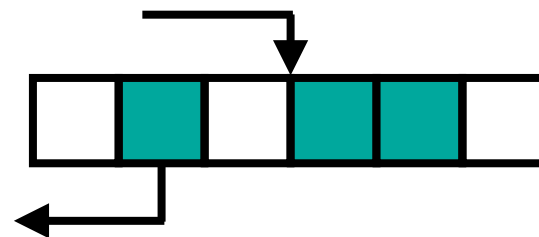
express wide range
of packet scheduling
algorithms

schedule
10s of thousands
of flows

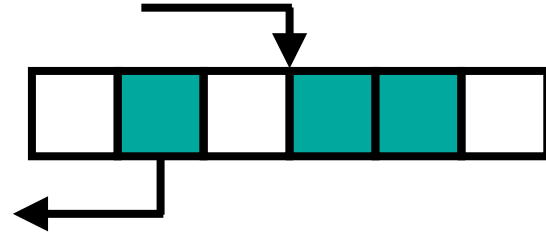
ref. [SENIC - NSDI'14]
[Carousel - SIGCOMM'17]

schedule
at line rate
100G+

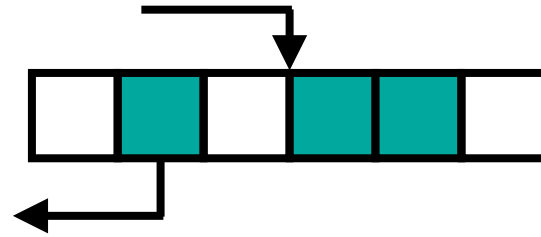
PIEO



PIEO



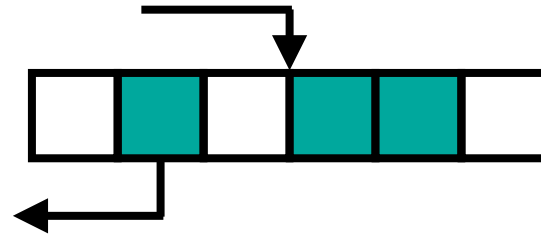
PIEO



Programmable

more expressive than any state-of-the-art packet scheduling primitive

PIEO



Programmable

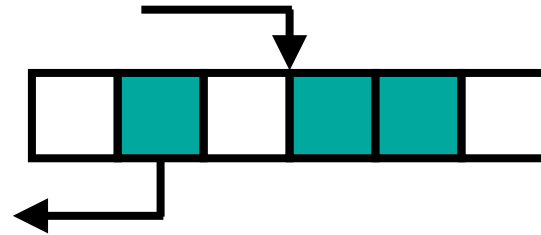
more expressive than any state-of-the-art packet scheduling primitive



Scalable

easily scales to 10s of thousands of flows

PIEO



Programmable

more expressive than any state-of-the-art packet scheduling primitive



Scalable

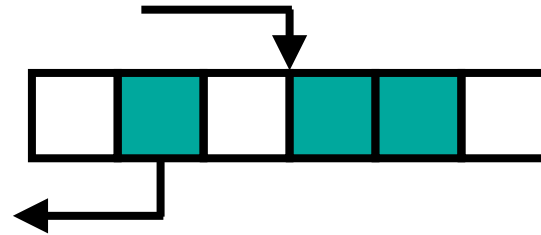
easily scales to 10s of thousands of flows



High-speed

makes scheduling decisions in $O(1)$ time [4 clock cycles]

PIEO



Programmable

more expressive than any state-of-the-art packet scheduling primitive

Abstraction



Scalable

easily scales to 10s of thousands of flows

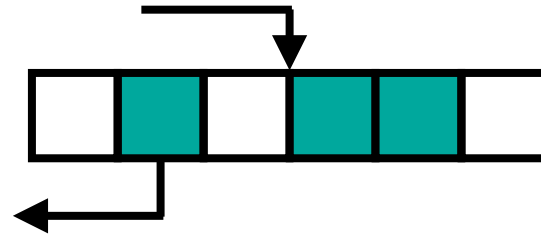
Hardware Design



High-speed

makes scheduling decisions in $O(1)$ time [4 clock cycles]

PIEO



Programmable

more expressive than any state-of-the-art packet scheduling primitive

Abstraction



Scalable

easily scales to 10s of thousands of flows

Hardware Design



High-speed

makes scheduling decisions in $O(1)$ time [4 clock cycles]

PIEO Scheduler Abstraction

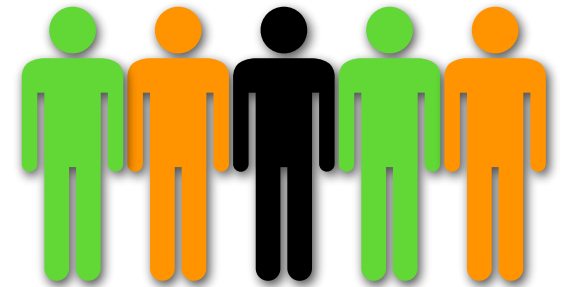
Scheduling Algorithms

PIEO Scheduler Abstraction



when an element becomes
eligible for scheduling?

Scheduling Algorithms

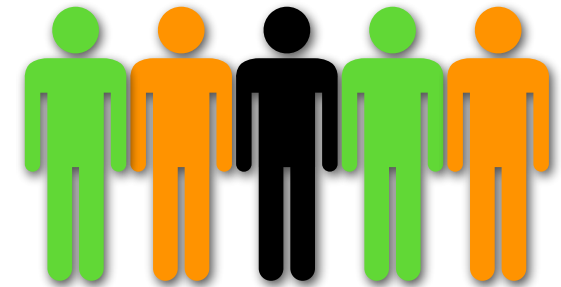


what order to schedule
amongst eligible elements?

PIEO Scheduler Abstraction



Scheduling Algorithms



when an element becomes
eligible for scheduling?
encode using a $t_{eligible}$ value

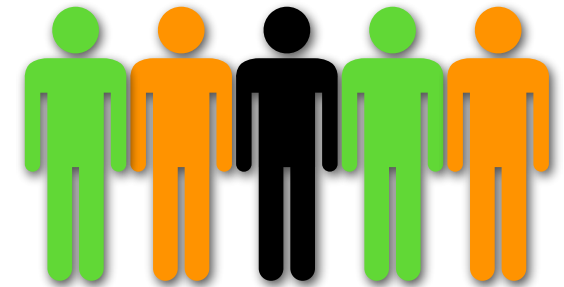


what order to schedule
amongst eligible elements?
encode using a *rank* value

PIEO Scheduler Abstraction



Scheduling Algorithms



when an element becomes
eligible for scheduling?

encode using a $t_{eligible}$ value

what order to schedule
amongst eligible elements?
encode using a *rank* value

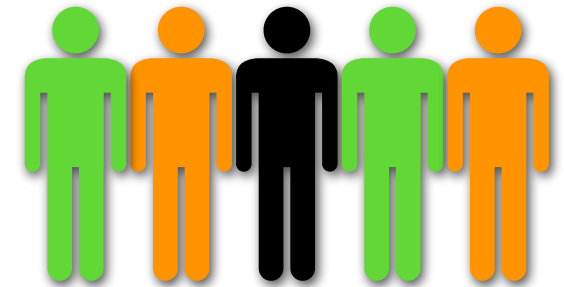
whenever the link is idle:

among all elements satisfying the eligibility predicate $t_{current} \geq t_{eligible}$:
schedule the smallest ranked element

PIEO Scheduler Abstraction



Scheduling Algorithms



when an element becomes eligible for scheduling?

encode using a $t_{eligible}$ value

what order to schedule amongst eligible elements?
encode using a *rank* value

whenever the link is idle:

among all elements satisfying the eligibility predicate $t_{current} \geq t_{eligible}$:
schedule the smallest ranked element

PIEO scheduler simply schedules the **smallest ranked eligible** element
at any given time

Push-In-Extract-Out Primitive

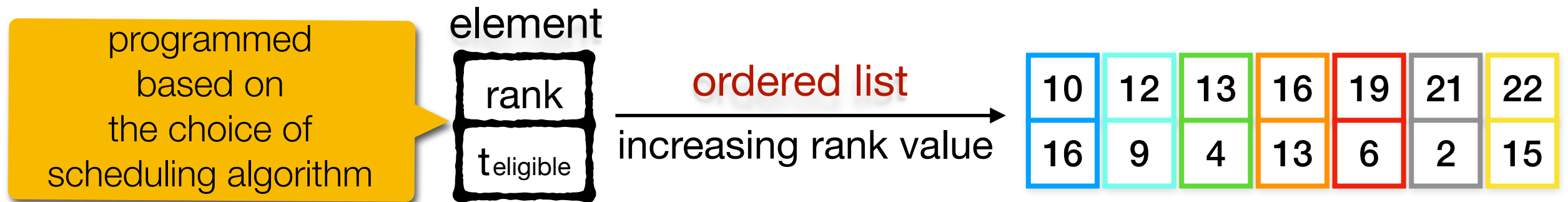
Push-In-Extract-Out Primitive

programmed
based on
the choice of
scheduling algorithm

element



Push-In-Extract-Out Primitive



Push-In-Extract-Out Primitive

programmed
based on
the choice of
scheduling algorithm

element



ordered list

increasing rank value

10	12	13	16	19	21	22
16	9	4	13	6	2	15

enqueue(

18
1

)

10	12	13	16	19	21	22
16	9	4	13	6	2	15

dequeue()

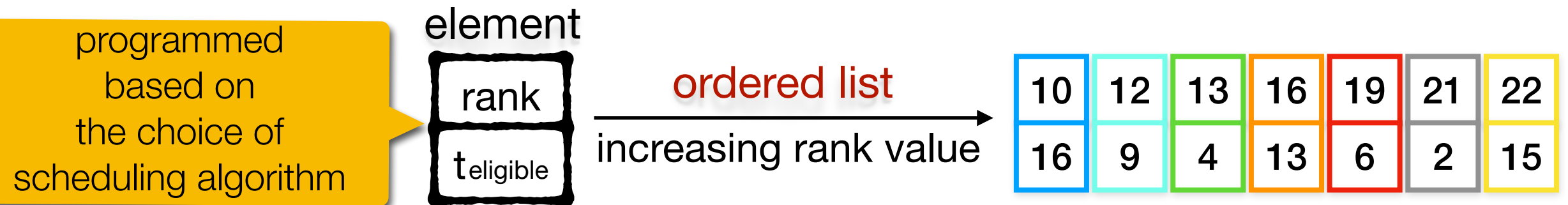
10	12	13	16	19	21	22
16	9	4	13	6	2	15

dequeue(

)

10	12	13	16	19	21	22
16	9	4	13	6	2	15

Push-In-Extract-Out Primitive



enqueue(

18
1

) "Push-In"

inserts element at position dictated by its rank value

10	12	13	16	19	21	22
16	9	4	13	6	2	15

dequeue()

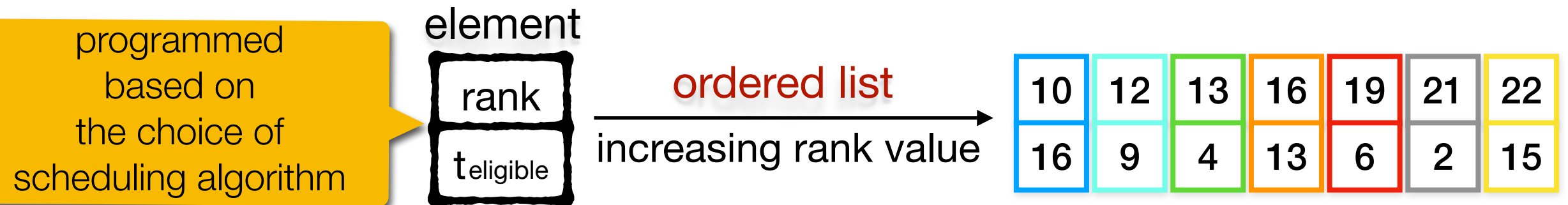
10	12	13	16	19	21	22
16	9	4	13	6	2	15

dequeue(

)

10	12	13	16	19	21	22
16	9	4	13	6	2	15

Push-In-Extract-Out Primitive



enqueue(

18
1

) "Push-In"

inserts element at position dictated by its rank value

10	12	13	16	18	19	21	22
16	9	4	13	1	6	2	15

dequeue()

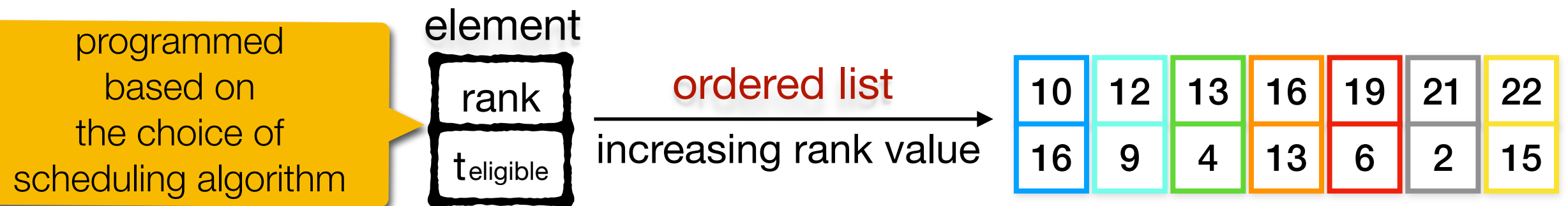
10	12	13	16	19	21	22
16	9	4	13	6	2	15

dequeue(

)

10	12	13	16	19	21	22
16	9	4	13	6	2	15

Push-In-Extract-Out Primitive



enqueue(

18
1

) “Push-In”

inserts element at position dictated by its rank value

10	12	13	16	18	19	21	22
16	9	4	13	1	6	2	15

dequeue() “Extract-Out”

returns “smallest ranked eligible” element

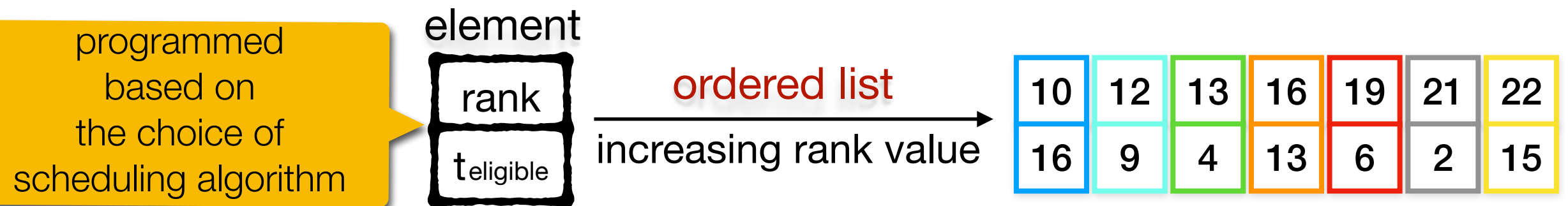
10	12	13	16	19	21	22
16	9	4	13	6	2	15

dequeue(

)

10	12	13	16	19	21	22
16	9	4	13	6	2	15

Push-In-Extract-Out Primitive



enqueue(

18
1

) "Push-In"

inserts element at position dictated by its rank value

10	12	13	16	18	19	21	22
16	9	4	13	1	6	2	15

$t_{\text{current}} = 7$

dequeue() "Extract-Out"

returns "smallest ranked eligible" element

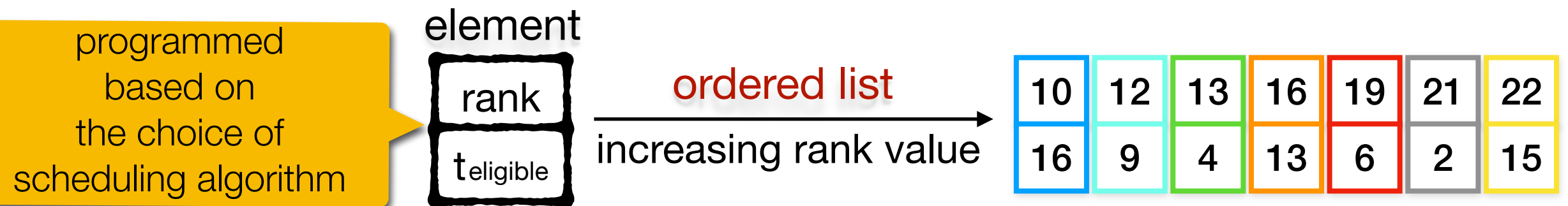
10	12	13	16	19	21	22
16	9	4	13	6	2	15

dequeue(

)

10	12	13	16	19	21	22
16	9	4	13	6	2	15

Push-In-Extract-Out Primitive



enqueue(

18
1

) "Push-In"

inserts element at position dictated by its rank value

10	12	13	16	18	19	21	22
16	9	4	13	1	6	2	15

$$t_{\text{current}} = 7$$

filter: $t_{\text{current}} \geq t_{\text{eligible}}$

dequeue() "Extract-Out"

returns "smallest ranked eligible" element

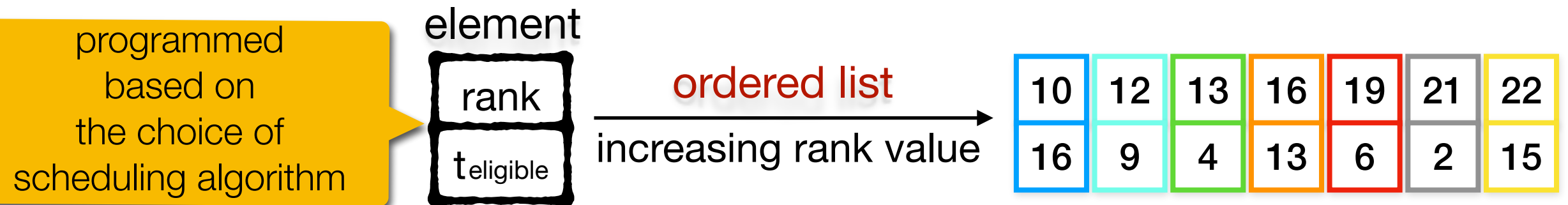
10	12	13	16	19	21	22
16	9	4	13	6	2	15

dequeue(

)

10	12	13	16	19	21	22
16	9	4	13	6	2	15

Push-In-Extract-Out Primitive



enqueue(

18
1

) "Push-In"

inserts element at position dictated by its rank value

10	12	13	16	18	19	21	22
16	9	4	13	1	6	2	15

$t_{\text{current}} = 7$

filter: $t_{\text{current}} \geq t_{\text{eligible}}$

dequeue() "Extract-Out"

returns "smallest ranked eligible" element

13
4

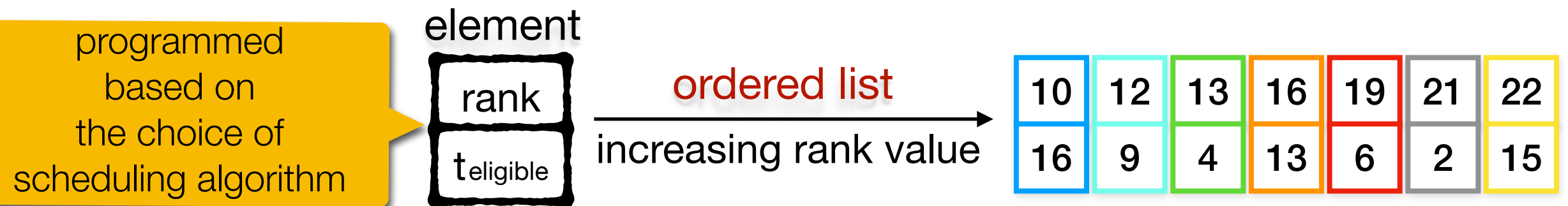
10	12	16	19	21	22
16	9	13	6	2	15

dequeue(

)

10	12	13	16	19	21	22
16	9	4	13	6	2	15

Push-In-Extract-Out Primitive



enqueue(

18
1

) "Push-In"

inserts element at position dictated by its rank value

10	12	13	16	18	19	21	22
16	9	4	13	1	6	2	15

$t_{\text{current}} = 7$

filter: $t_{\text{current}} \geq t_{\text{eligible}}$

dequeue() "Extract-Out"

returns "smallest ranked eligible" element

13
4

10	12	16	19	21	22
16	9	13	6	2	15

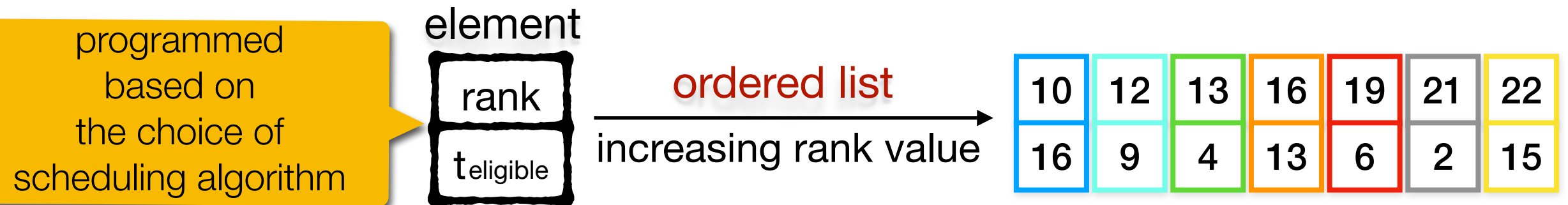
dequeue(

)

returns a specific element

10	12	13	16	19	21	22
16	9	4	13	6	2	15

Push-In-Extract-Out Primitive



enqueue(

18
1

) "Push-In"

inserts element at position dictated by its rank value

10	12	13	16	18	19	21	22
16	9	4	13	1	6	2	15

$t_{\text{current}} = 7$

filter: $t_{\text{current}} \geq t_{\text{eligible}}$

dequeue() "Extract-Out"

returns "smallest ranked eligible" element

13
4

10	12	16	19	21	22
16	9	13	6	2	15

dequeue(

)

returns a specific element

19
6

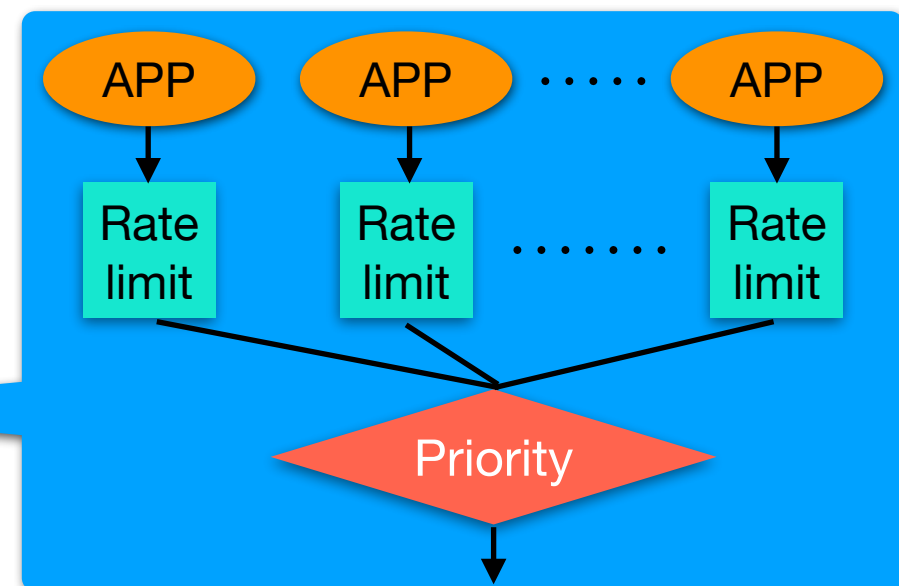
10	12	13	16	21	22
16	9	4	13	2	15

Expressiveness of PIEO

Expressiveness of PIEO

- Work conserving
 - e.g., DRR, WFQ, WF²Q
- Non-work conserving
 - e.g., Token Bucket, RCSP
- Hierarchical scheduling
 - e.g., HPFQ
- Asynchronous scheduling
 - e.g., Starvation avoidance, D³
- Priority scheduling
 - e.g., SJF, SRTF, LSTF, EDF
- Complex scheduling policies
 - mixture of shaping and ordering

e.g.



Expressiveness of PIEO

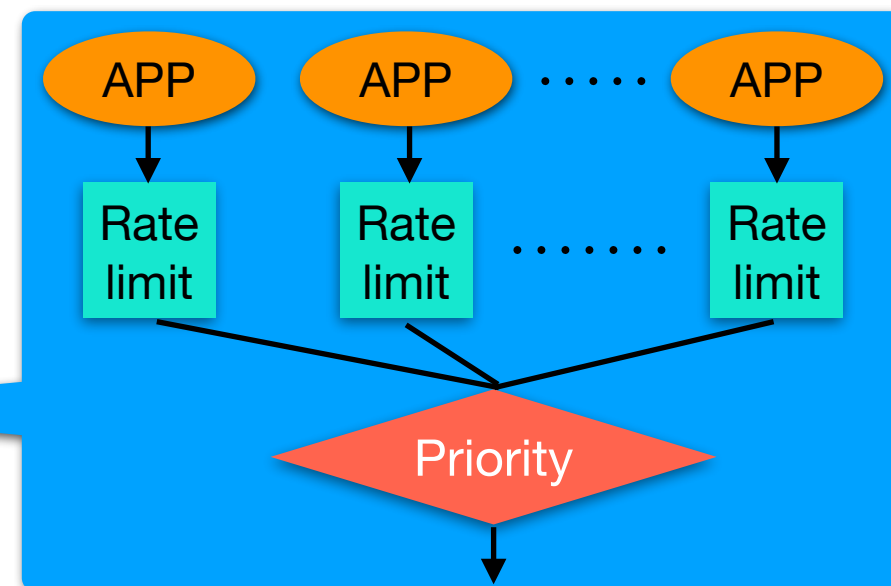
- Work conserving
 - e.g., DRR, WFQ, WF²Q
- Non-work conserving
 - e.g., Token Bucket, RCSP
- Hierarchical scheduling
 - e.g., HPFQ
- Asynchronous scheduling
 - e.g., Starvation avoidance, D³
- Priority scheduling
 - e.g., SJF, SRTF, LSTF, EDF
- Complex scheduling policies
 - mixture of shaping and ordering

for each element:

calculate **start_time** and **finish_time**

at time x , all elements s.t. **virtual_time(x)** \geq **start_time**:
schedule element with **smallest finish_time**

e.g.



Expressiveness of PIEO

- Work conserving
 - e.g., DRR, WFQ, WF²Q
- Non-work conserving
 - e.g., Token Bucket, RCSP
- Hierarchical scheduling
 - e.g., HPFQ
- Asynchronous scheduling
 - e.g., Starvation avoidance, D³
- Priority scheduling
 - e.g., SJF, SRTF, LSTF, EDF
- Complex scheduling policies
 - mixture of shaping and ordering

for each element:

calculate **start_time** and **finish_time**

at time x , all elements s.t. **virtual_time(x)** \geq **start_time**:
schedule element with **smallest finish_time**

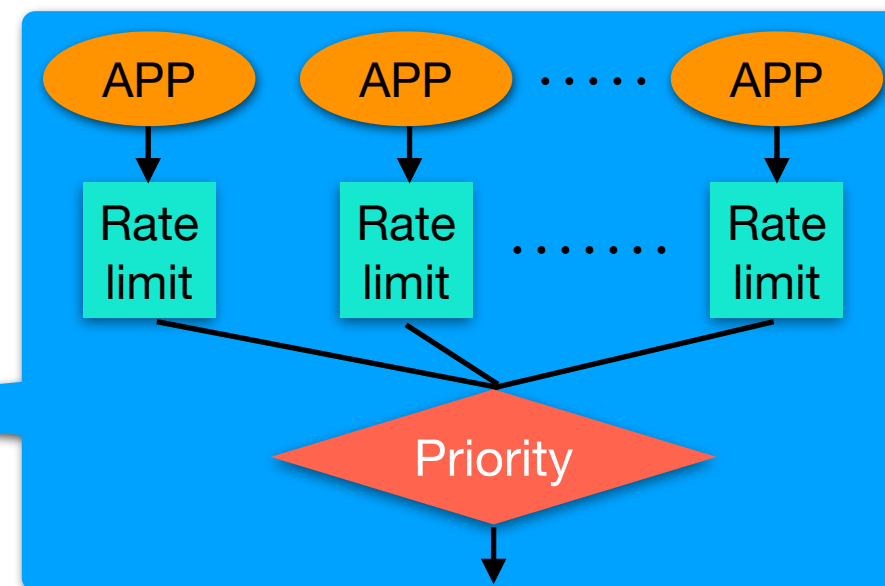
programming PIEO

$rank = finish_time$

$t_{eligible} = start_time$

Predicate for filtering at dequeue at time x :
 $(virtual_time(x) \geq t_{eligible})$

e.g.



Expressiveness of PIEO

- Work conserving
 - e.g., DRR, WFQ, WF²Q
- Non-work conserving
 - e.g., Token Bucket, RCSP
- Hierarchical scheduling
 - e.g., HPFQ
- Asynchronous scheduling
 - e.g., Starvation avoidance, D³
- Priority scheduling
 - e.g., SJF, SRTF, LSTF, EDF
- Complex scheduling policies
 - mixture of shaping and ordering

for each element:

calculate **start_time** and **finish_time**

at time x , all elements s.t. **virtual_time(x) \geq start_time:**
schedule element with **smallest finish_time**

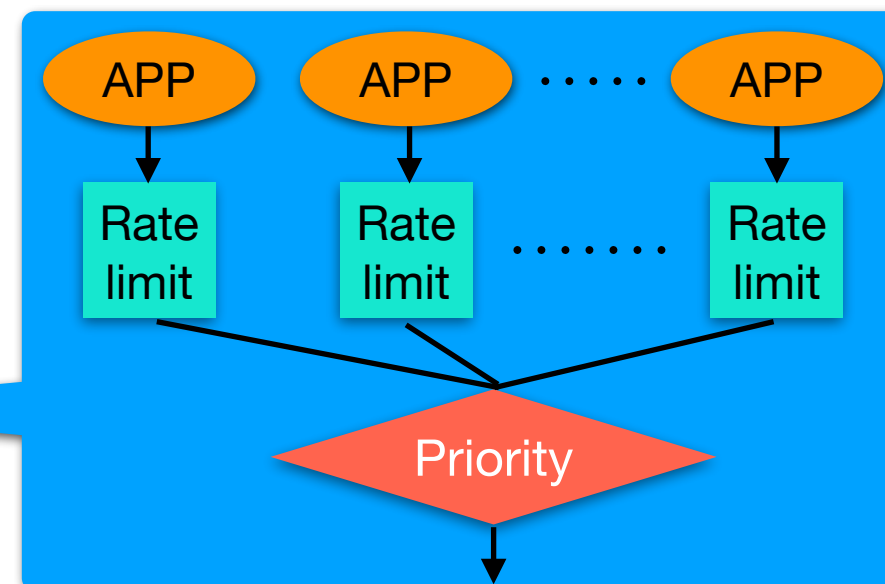
programming PIEO

$rank = finish_time$

$t_{eligible} = start_time$

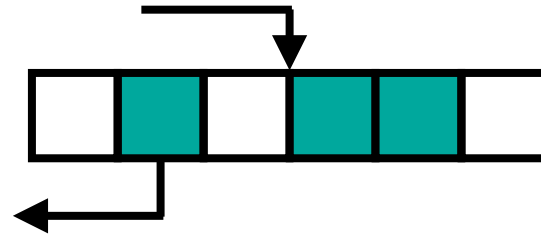
Predicate for filtering at dequeue at time x :
 $(virtual_time(x) \geq t_{eligible})$

e.g.



cannot express accurately using state-of-the-art (PIFO)

PIEO



Abstraction



Programmable

more expressive than any state-of-the-art packet scheduling primitive



Scalable

easily scales to 10s of thousands of flows



High-speed

makes scheduling decisions in $O(1)$ time [4 clock cycles]

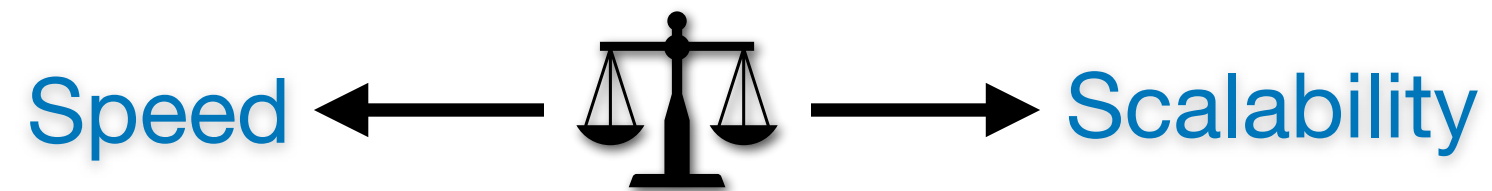
Hardware Design

Hardware Design

PIEO primitive relies on an ordered list datastructure

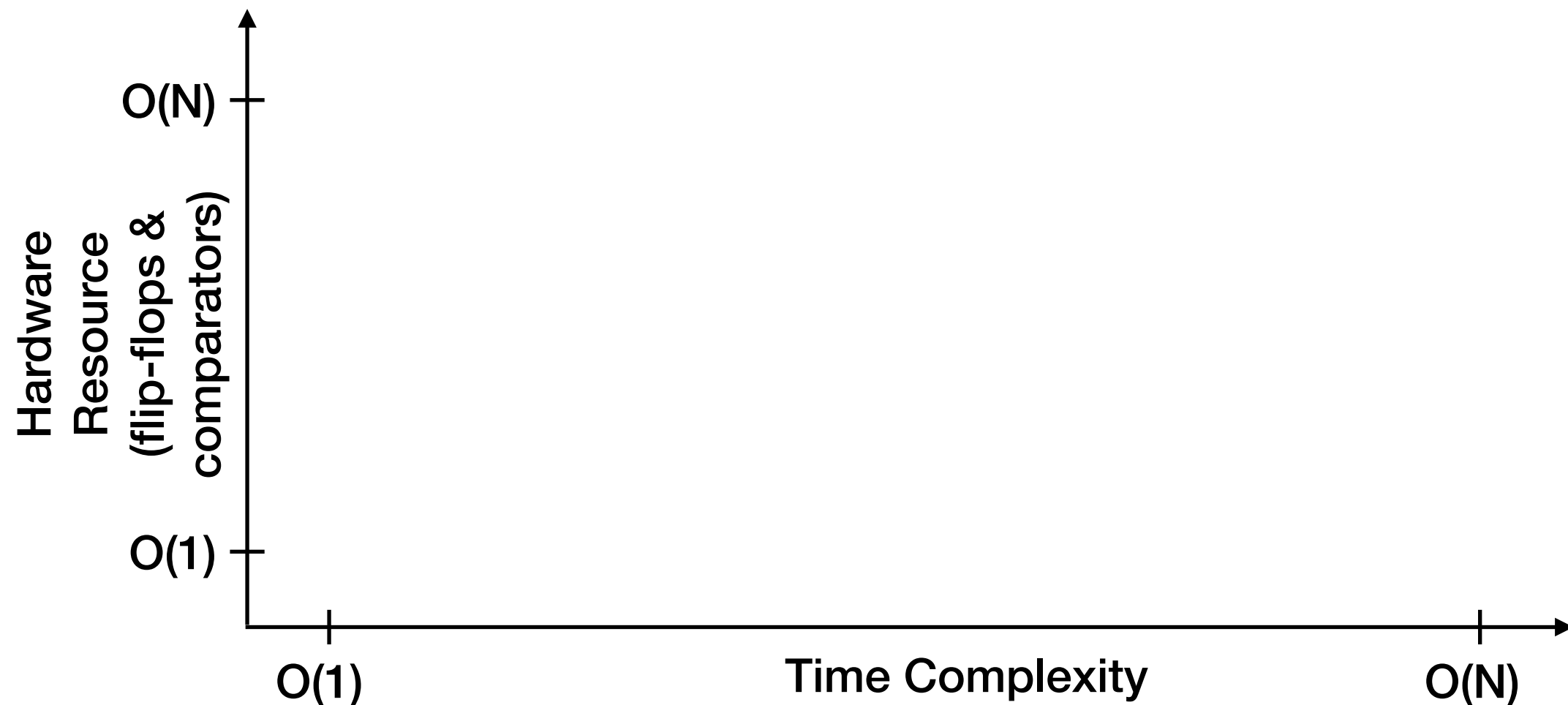
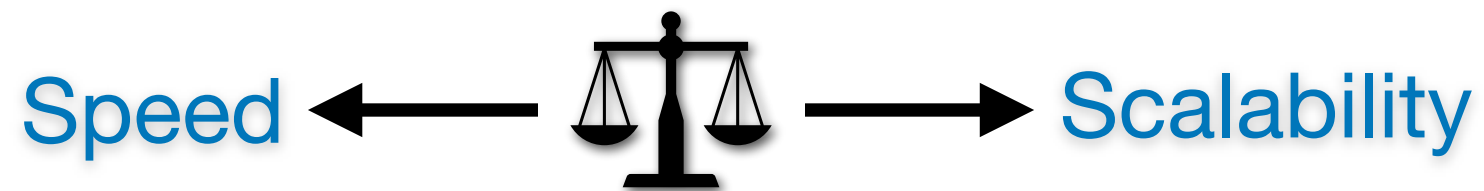
Hardware Design

PIEO primitive relies on an ordered list datastructure



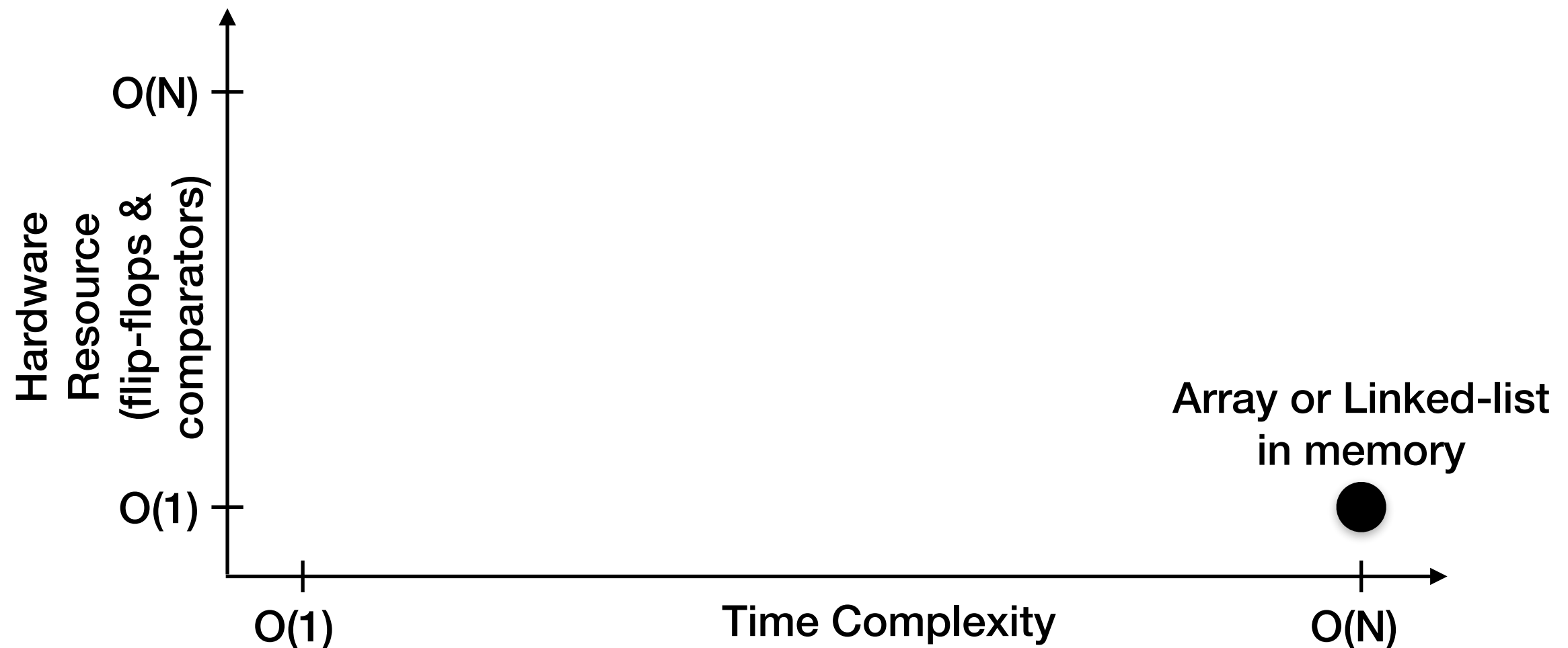
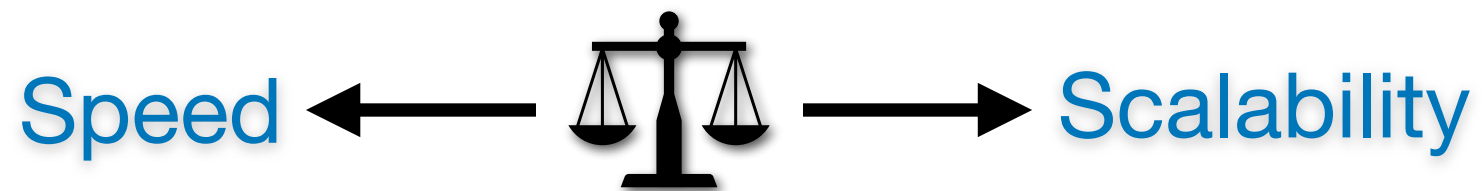
Hardware Design

PIEO primitive relies on an ordered list datastructure



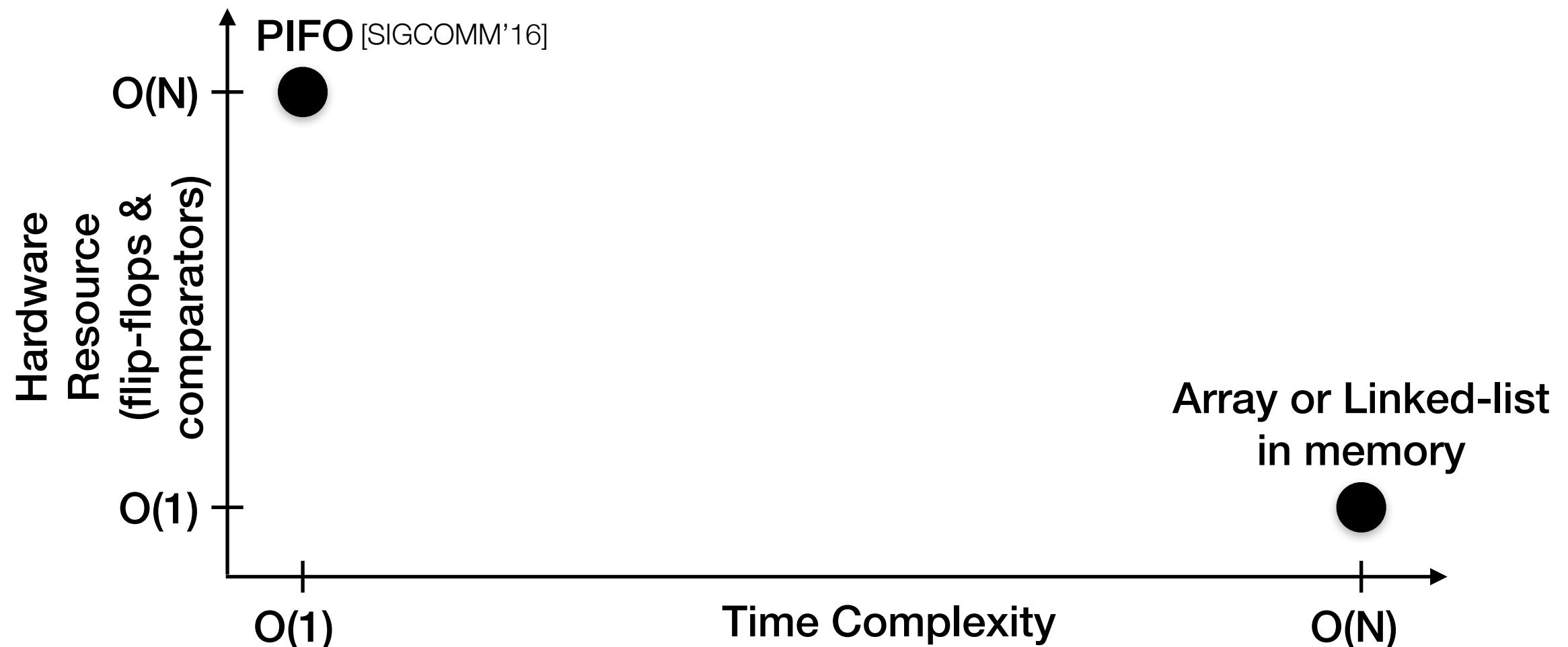
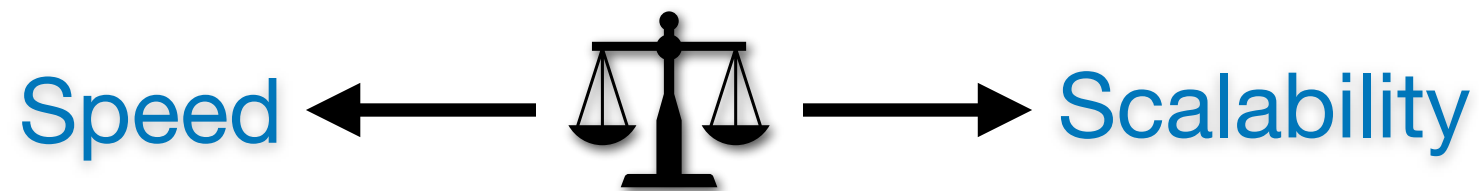
Hardware Design

PIEO primitive relies on an ordered list datastructure



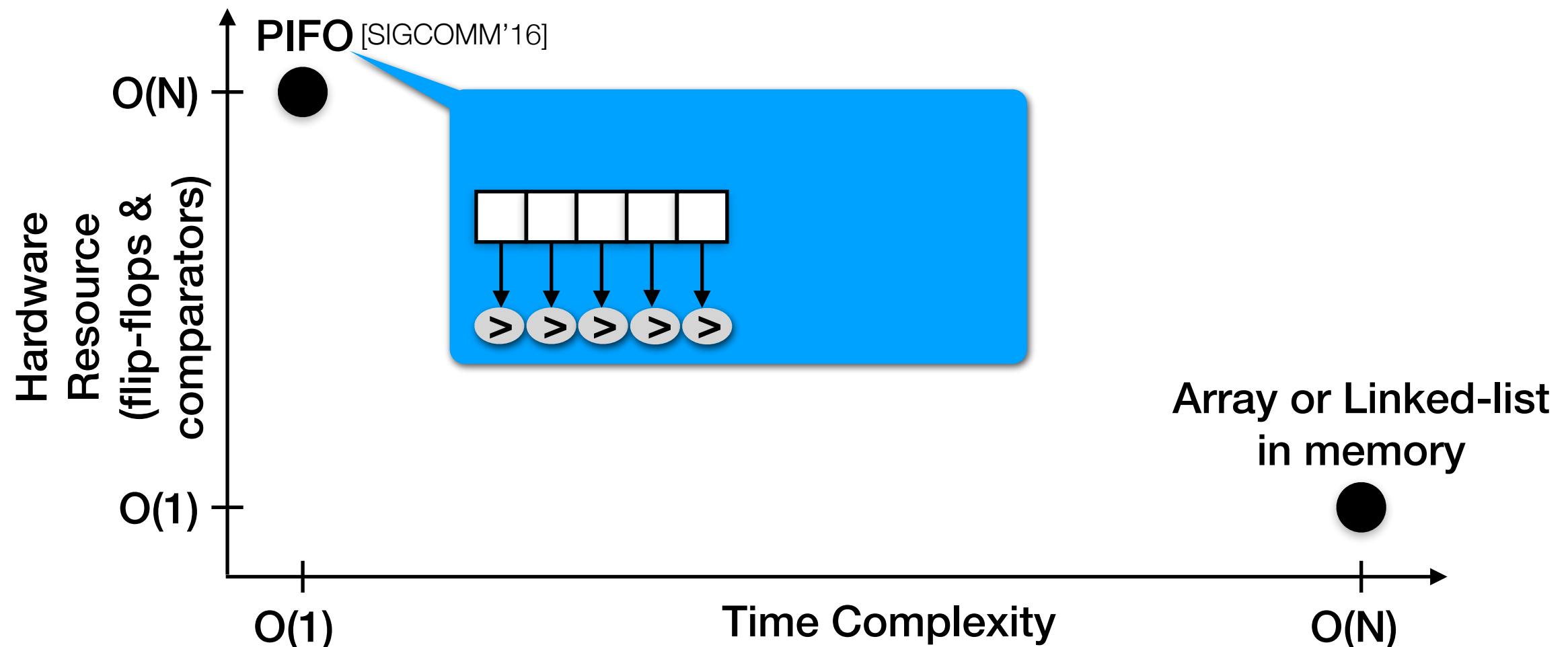
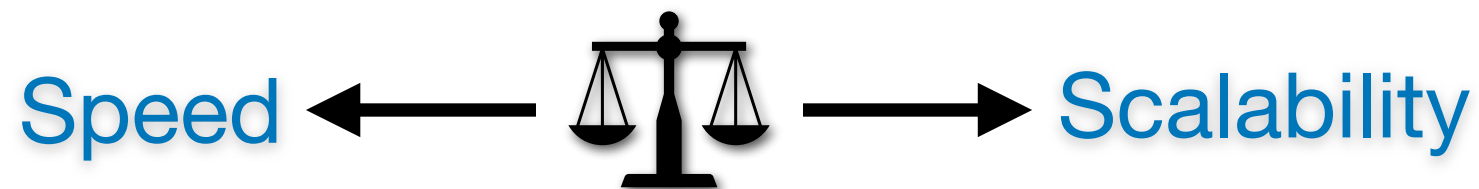
Hardware Design

PIEO primitive relies on an ordered list datastructure



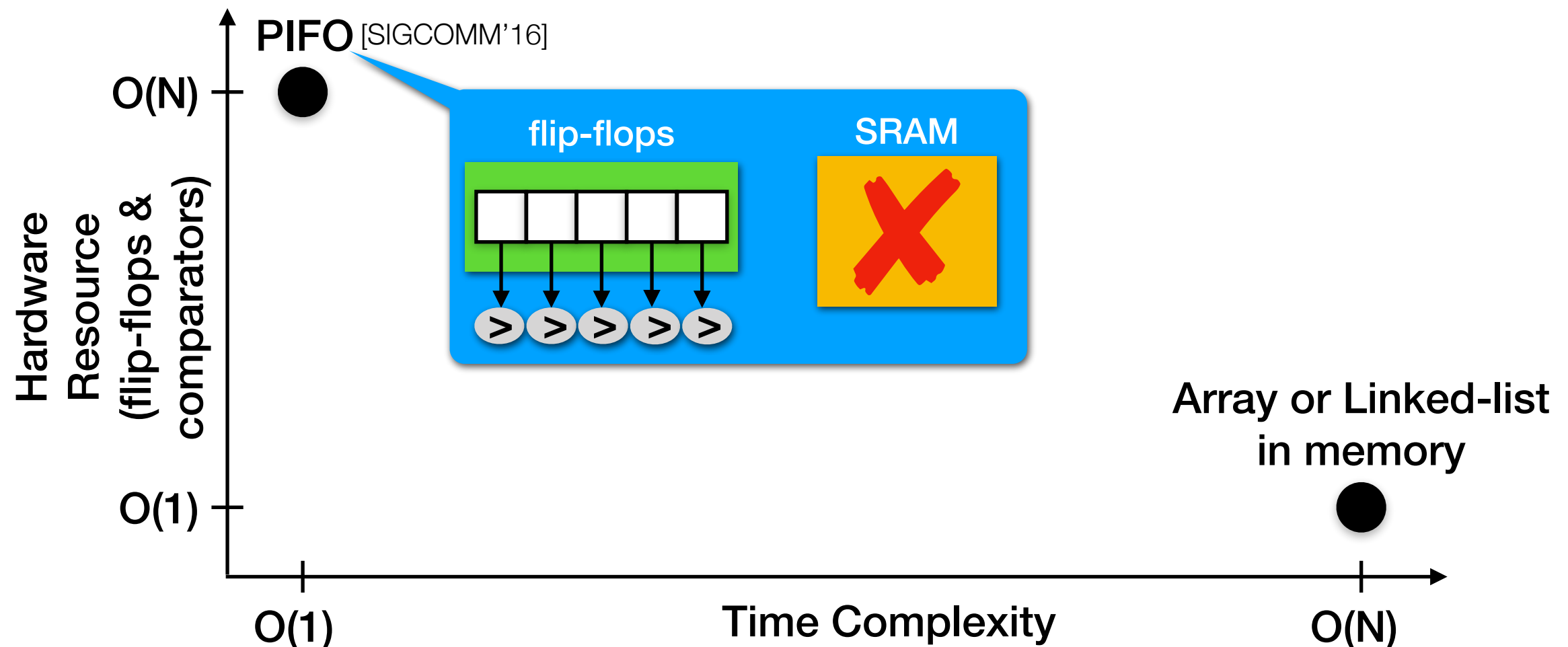
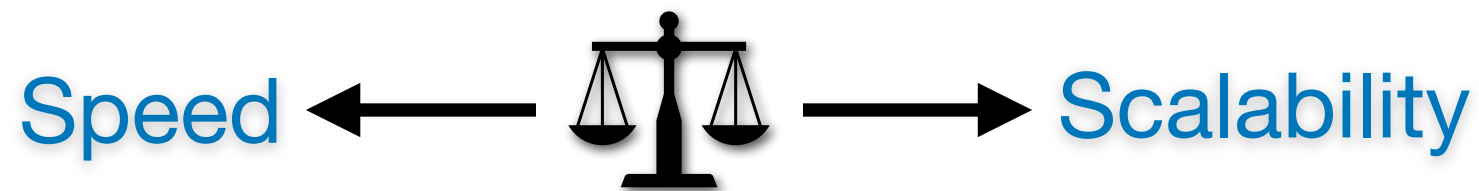
Hardware Design

PIEO primitive relies on an ordered list datastructure



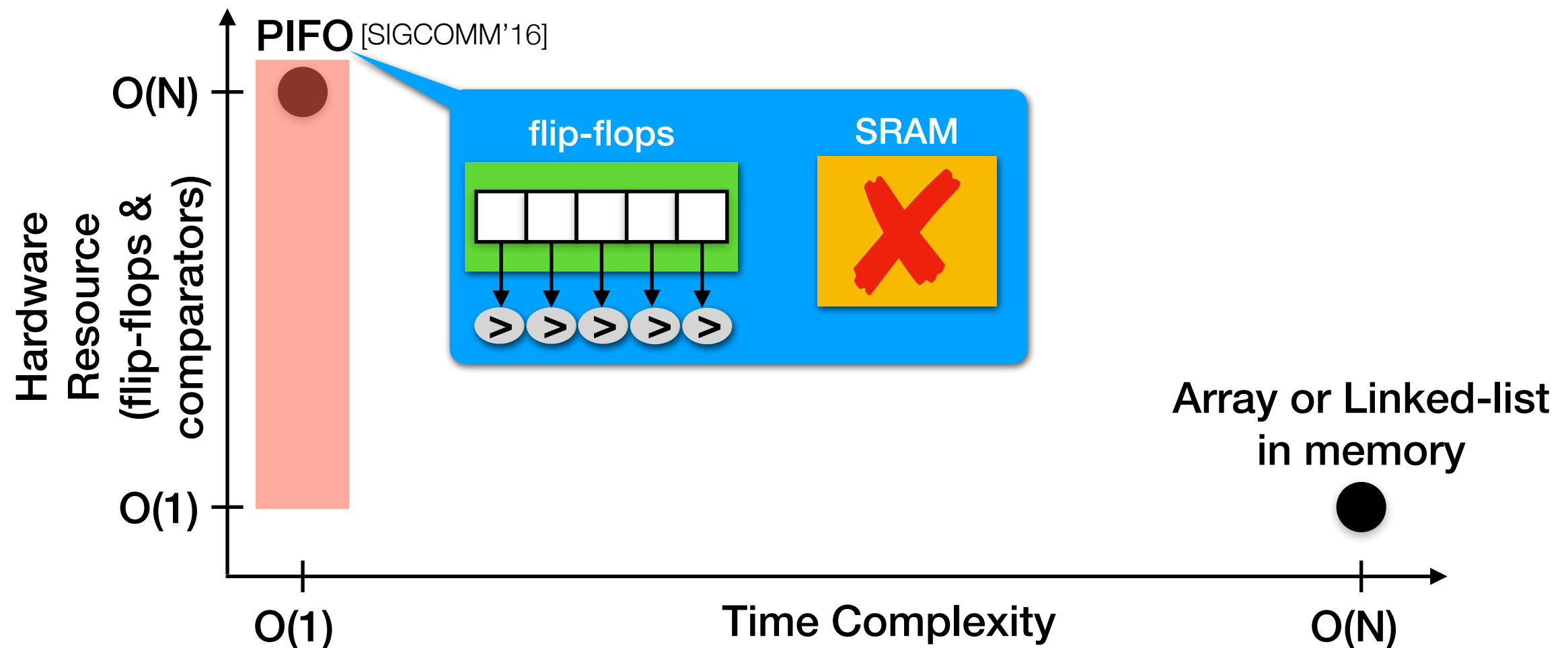
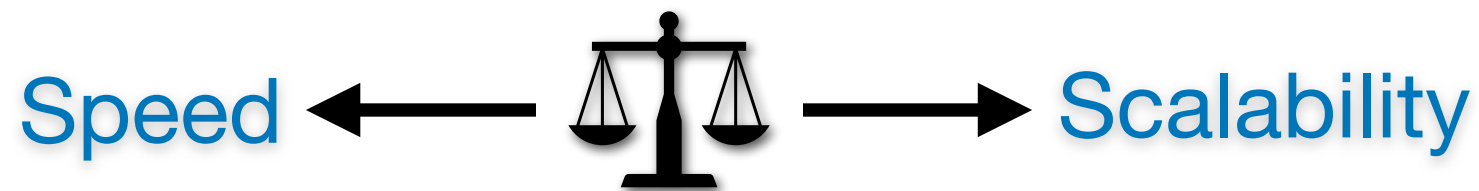
Hardware Design

PIEO primitive relies on an ordered list datastructure



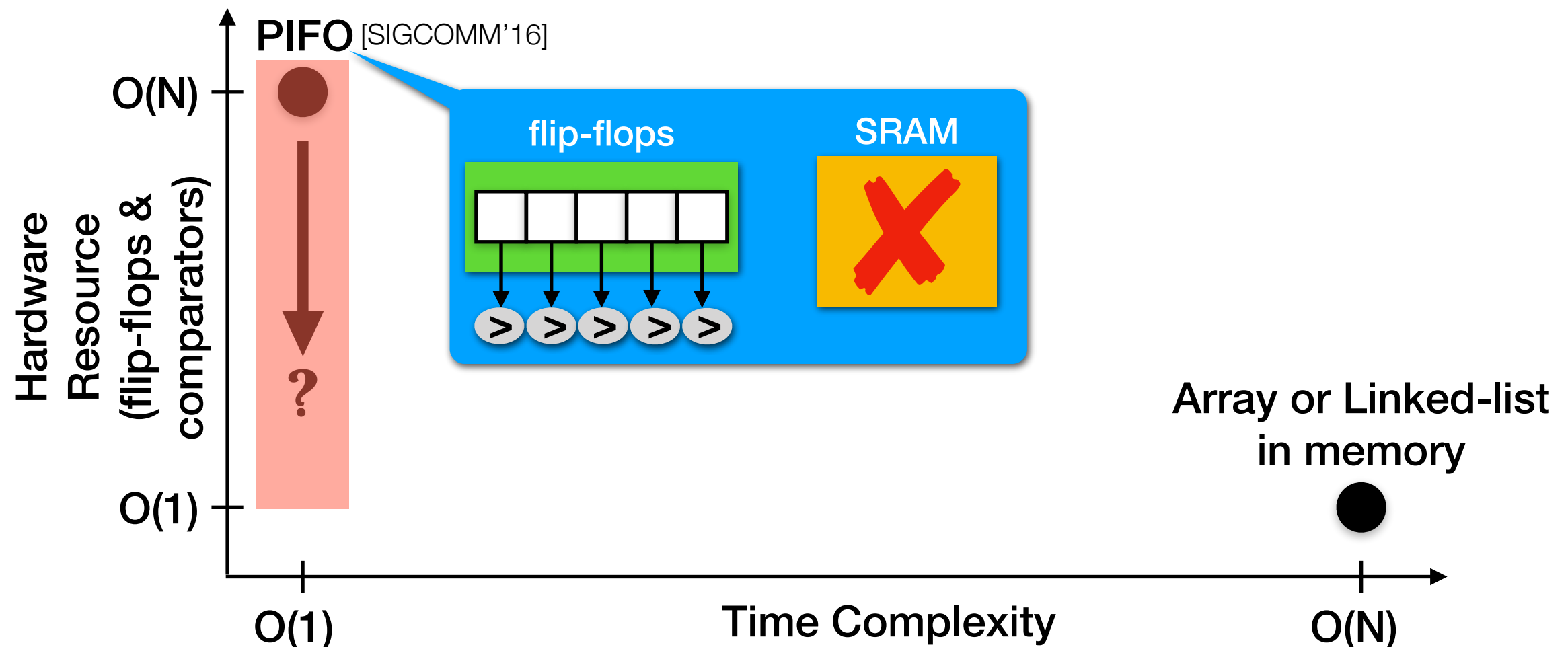
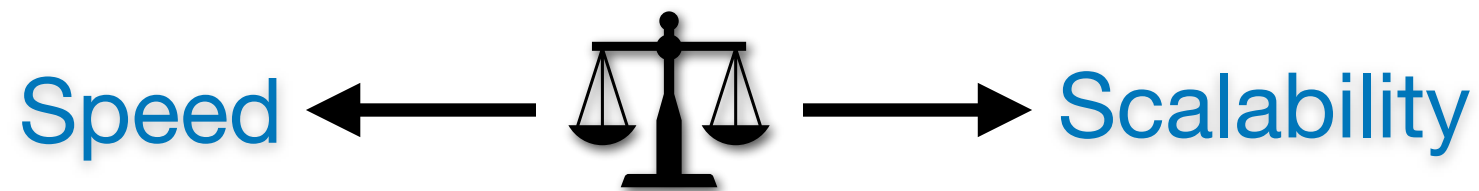
Hardware Design

PIEO primitive relies on an ordered list datastructure



Hardware Design

PIEO primitive relies on an ordered list datastructure



Is it fundamentally necessary to access and compare $O(N)$ elements in parallel to maintain an (exact) ordered list (of size N) in $O(1)$ time?

Is it fundamentally necessary to access and compare $O(N)$ elements in parallel to maintain an (exact) ordered list (of size N) in $O(1)$ time?

We present a design that can maintain an (exact) ordered list in $O(1)$ time, but only needs to access and compare $O(\sqrt{N})$ elements in parallel.

Is it fundamentally necessary to access and compare $O(N)$ elements in parallel to maintain an (exact) ordered list (of size N) in $O(1)$ time?

We present a design that can maintain an (exact) ordered list in $O(1)$ time, but only needs to access and compare $O(\sqrt{N})$ elements in parallel.

Key Insight

“All problems in computer science can be solved by another level of indirection”

—David Wheeler

Hardware Architecture

Flip-Flops



SRAM



Hardware Architecture

Flip-Flops

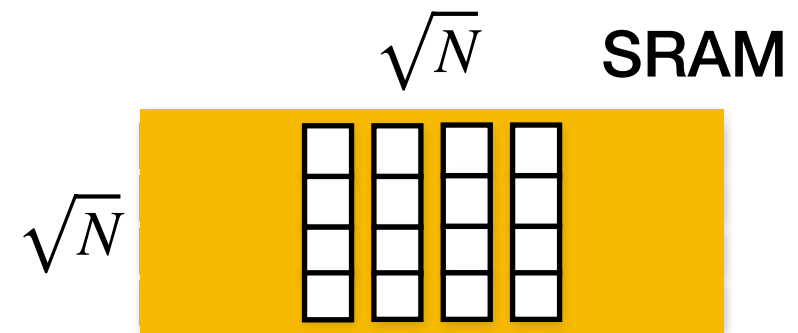


SRAM



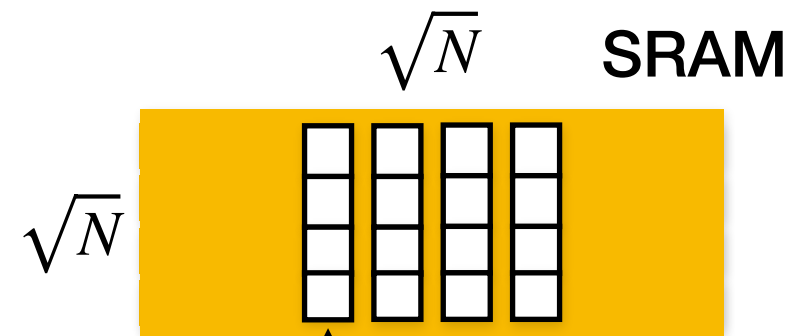
Hardware Architecture

Flip-Flops



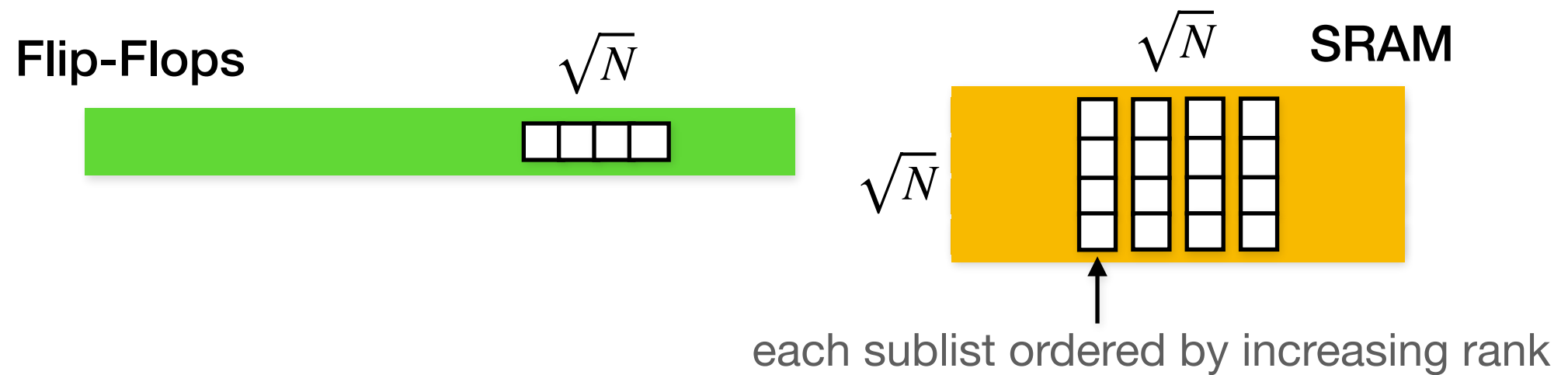
Hardware Architecture

Flip-Flops

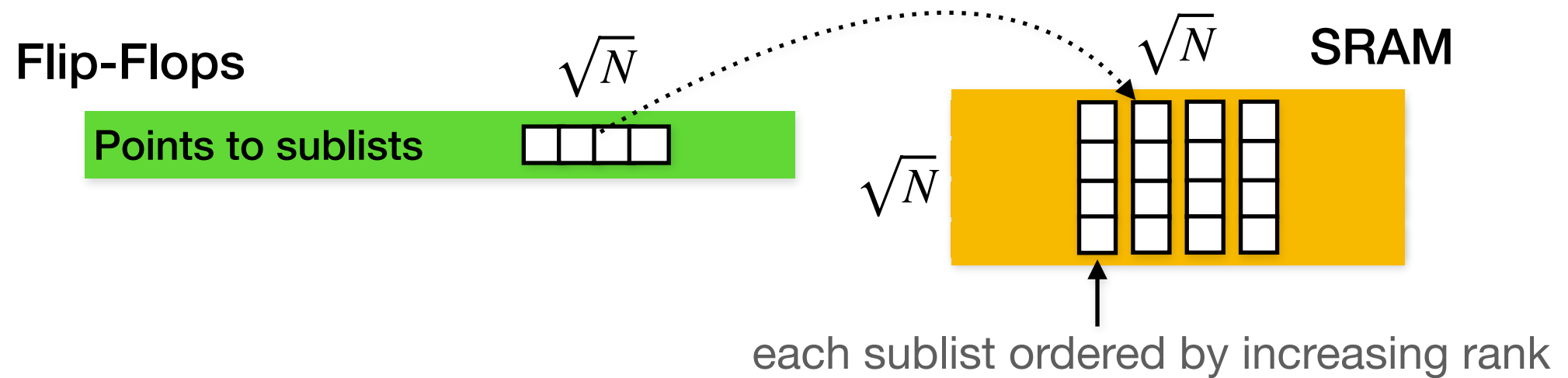


each sublist ordered by increasing rank

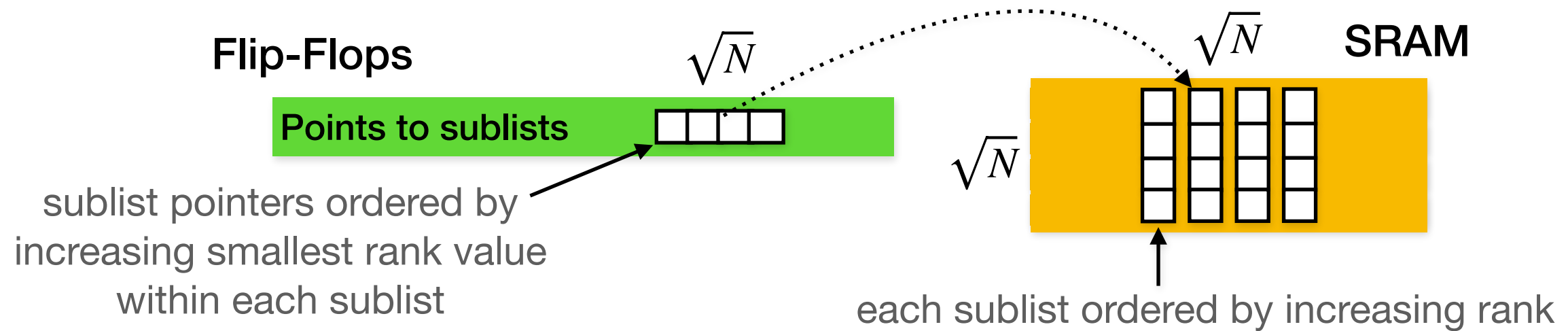
Hardware Architecture



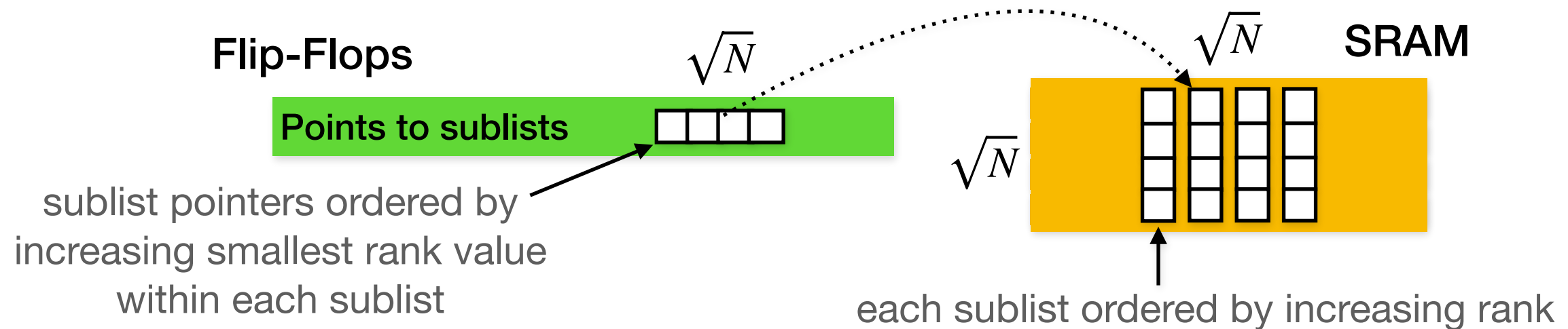
Hardware Architecture



Hardware Architecture

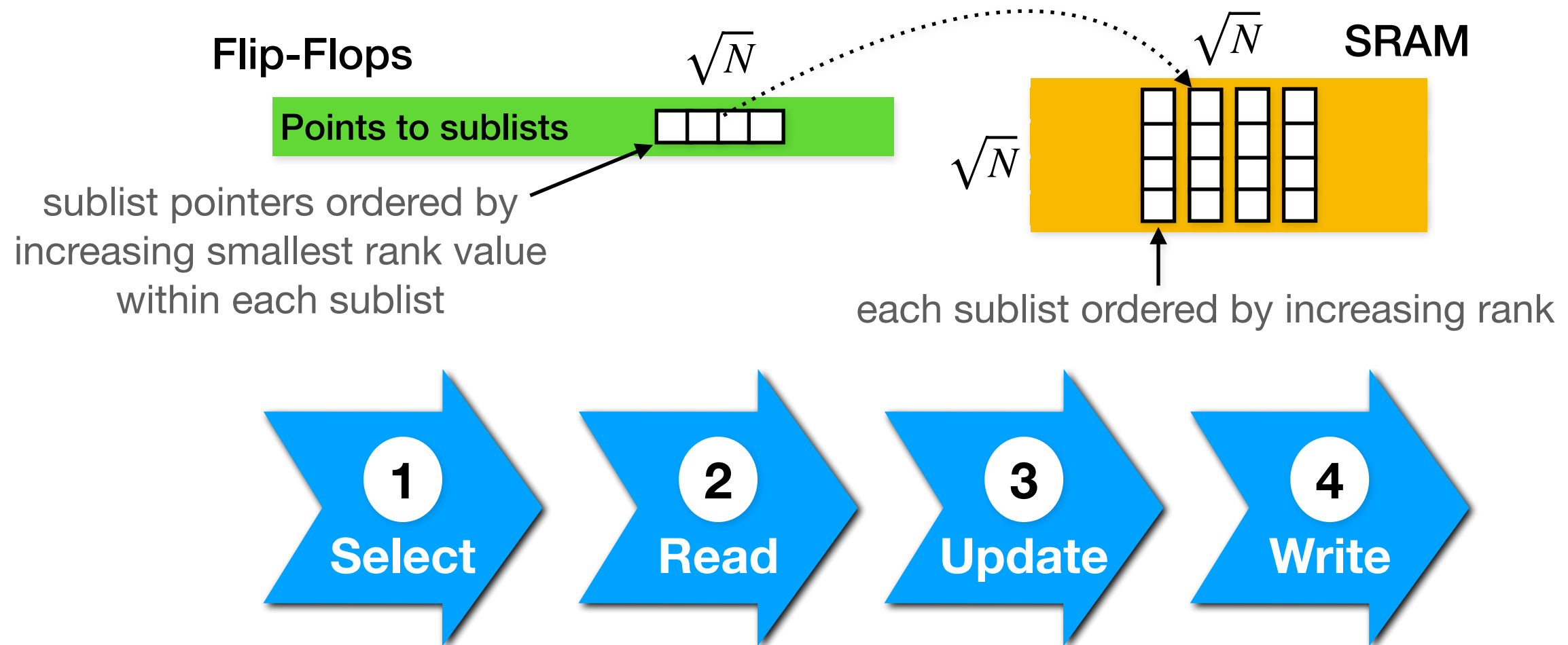


Hardware Architecture



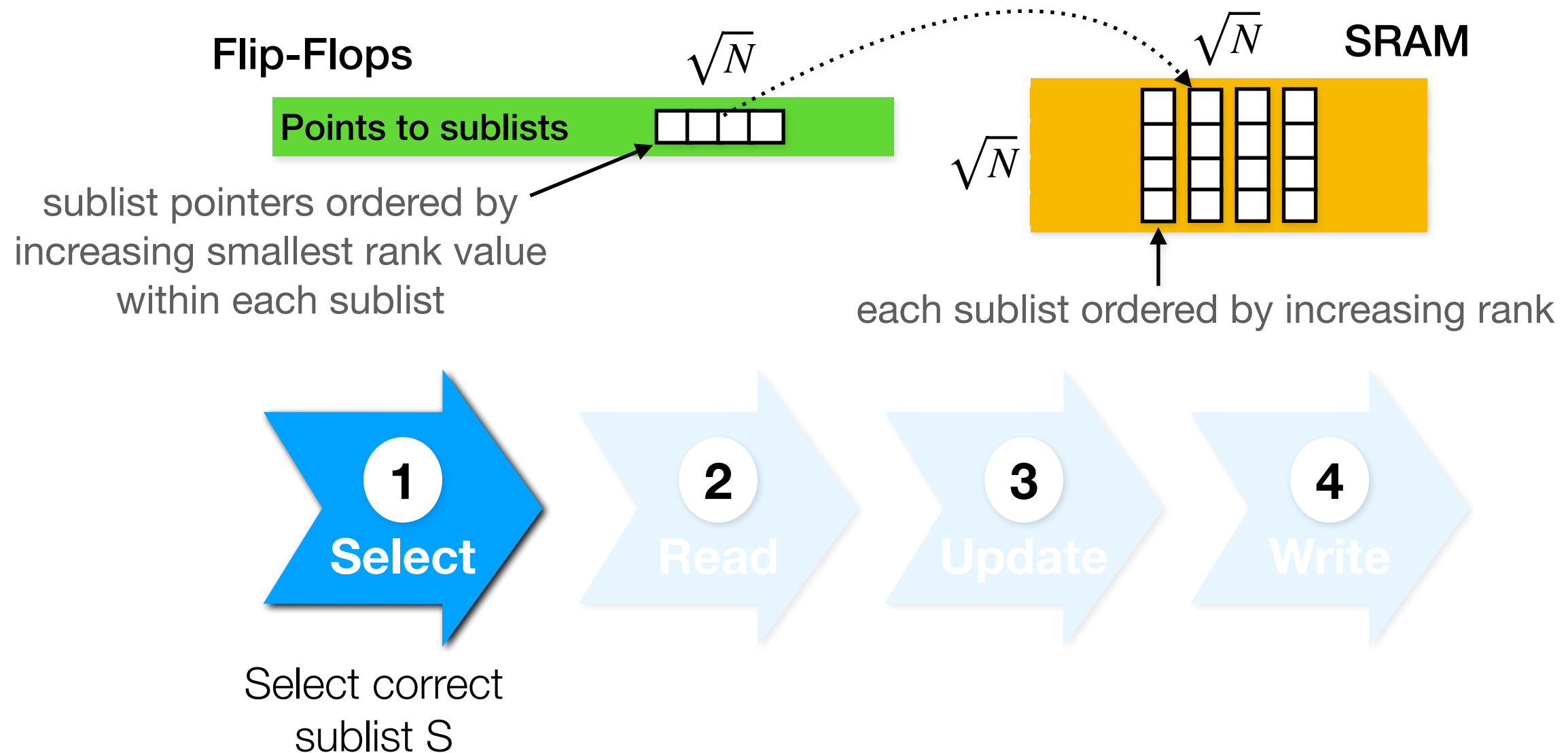
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



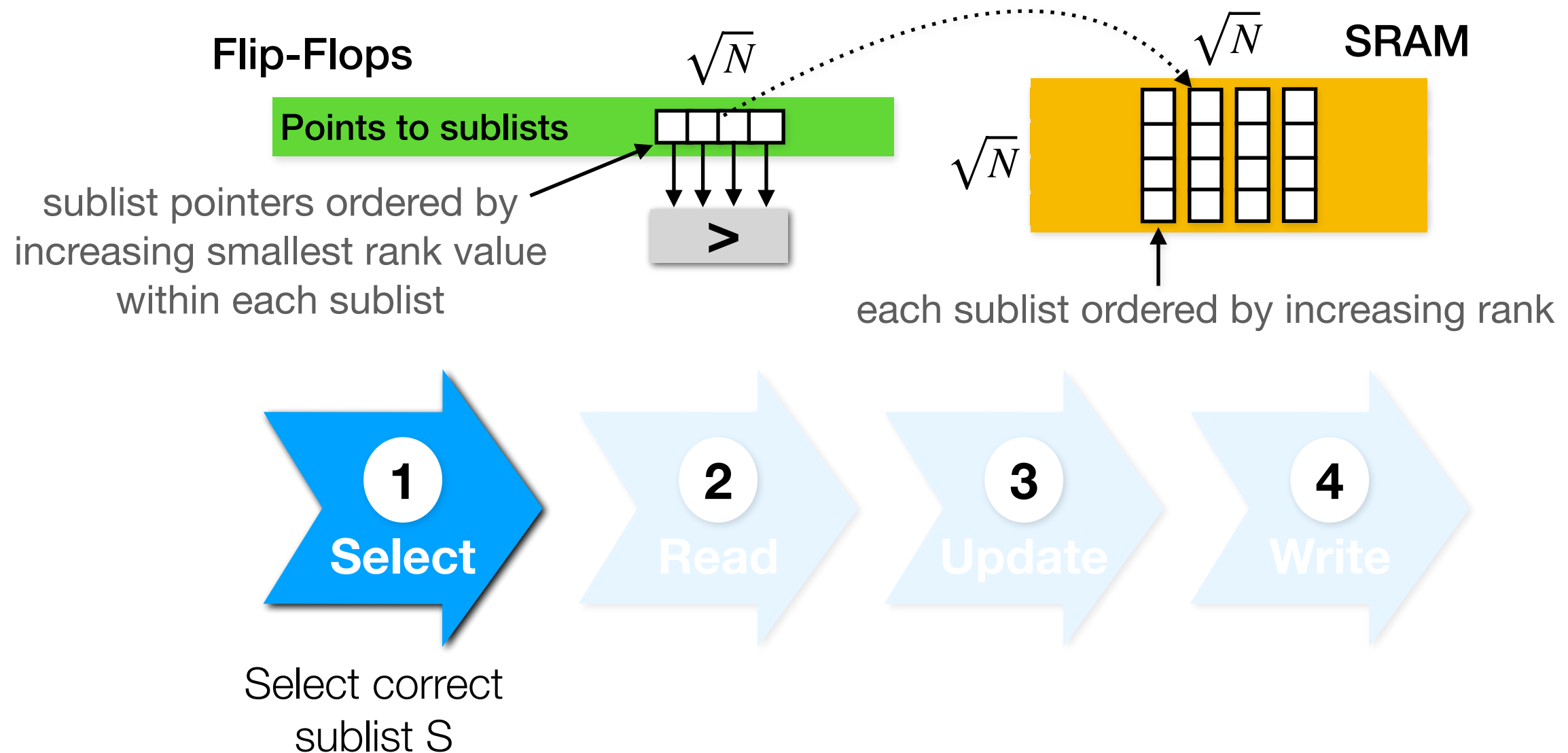
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



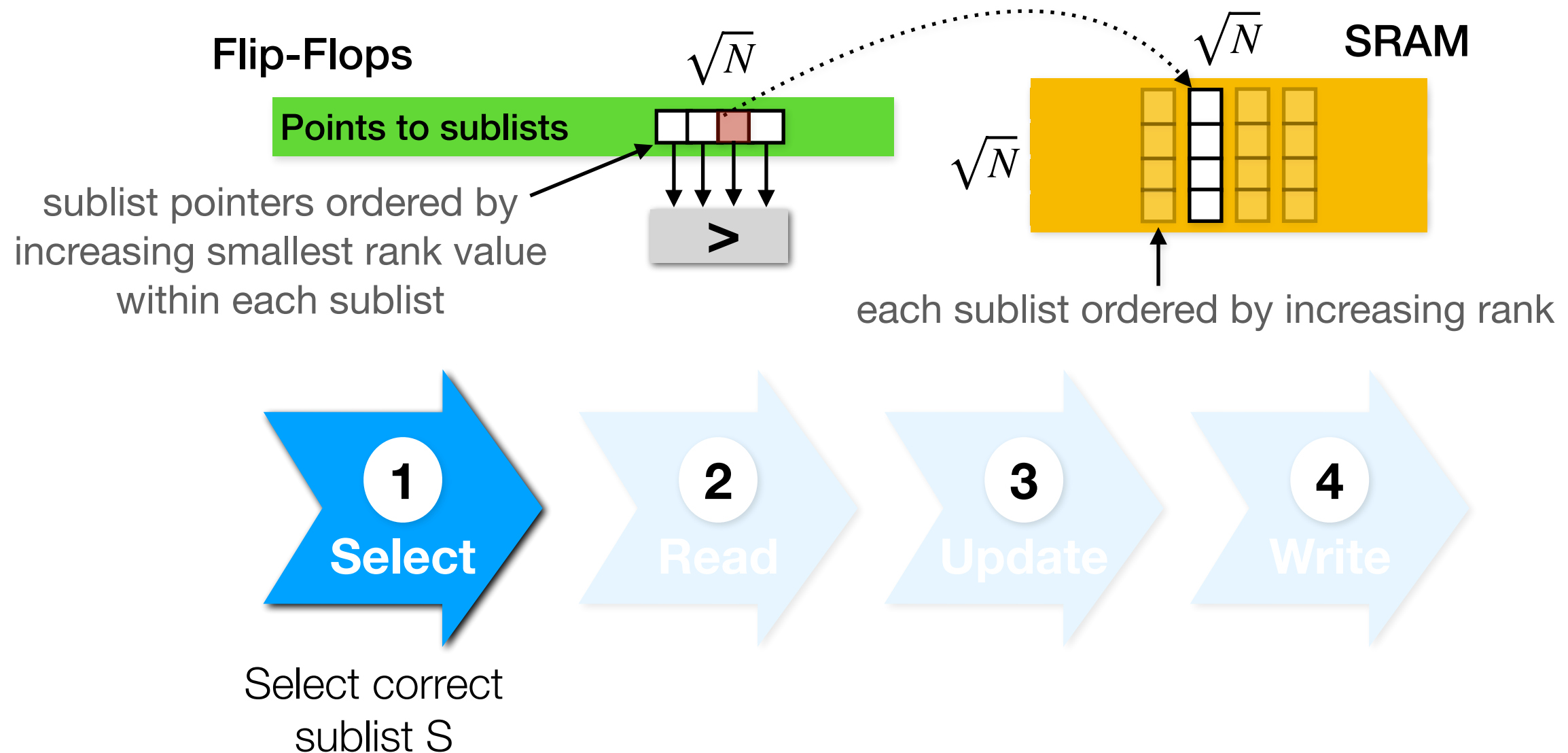
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



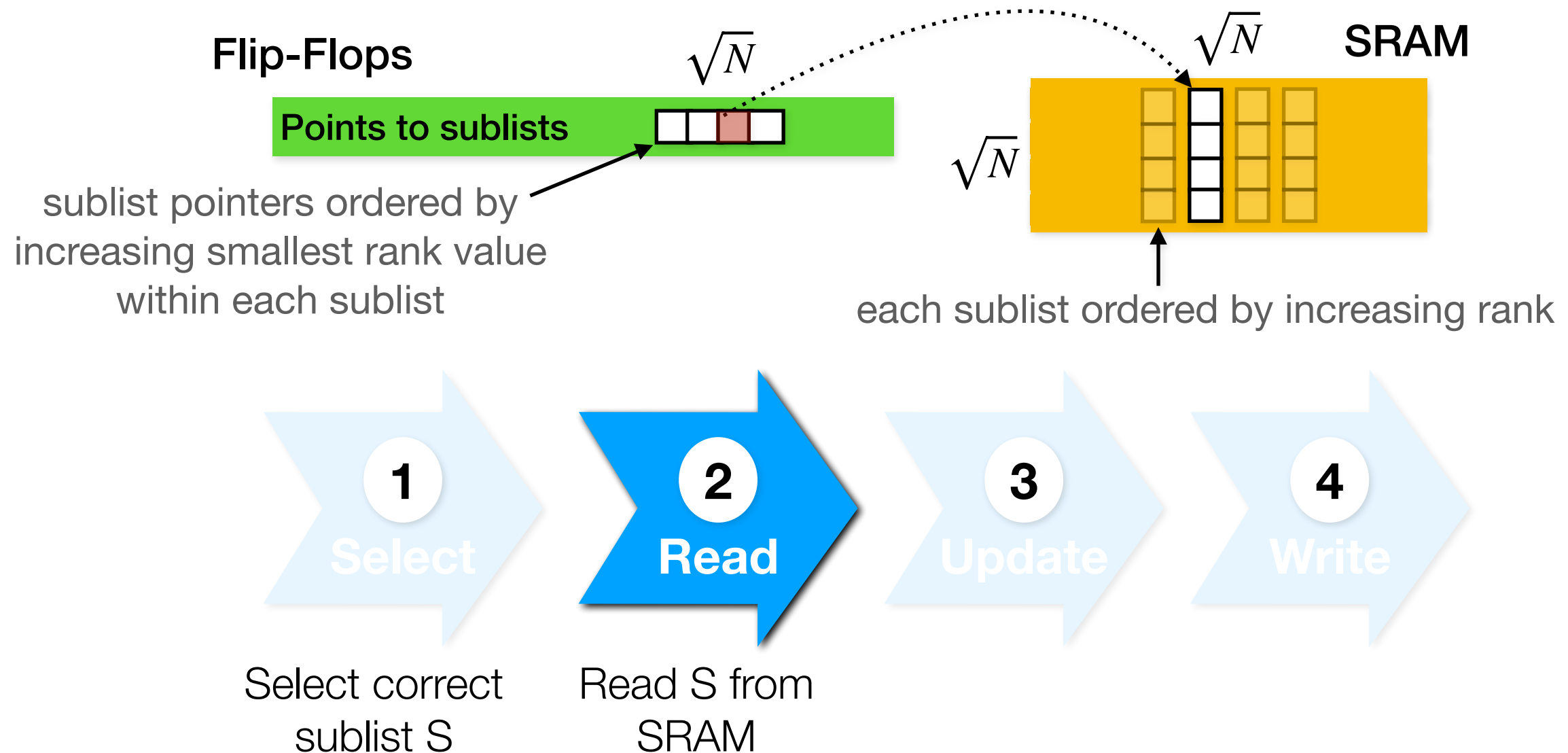
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



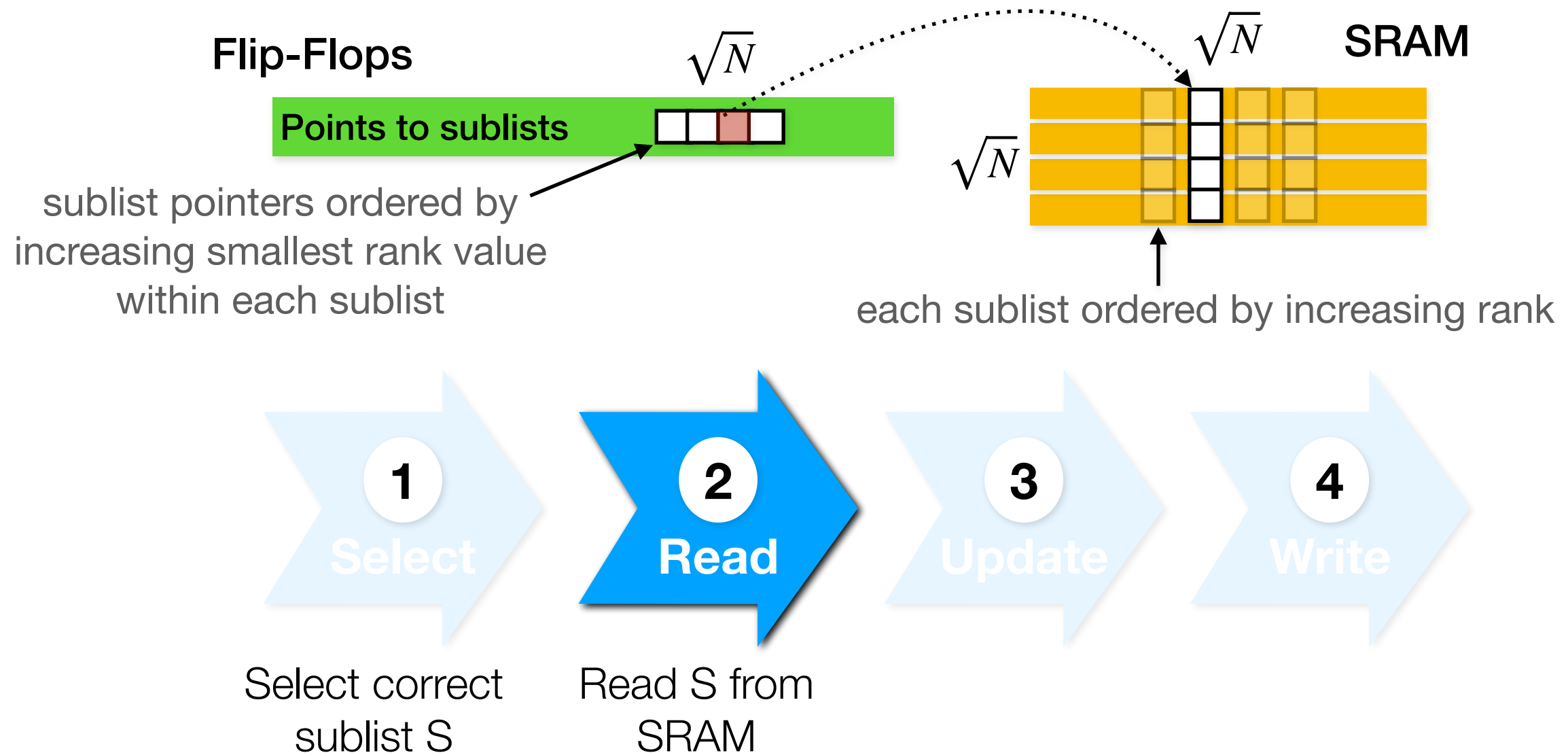
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



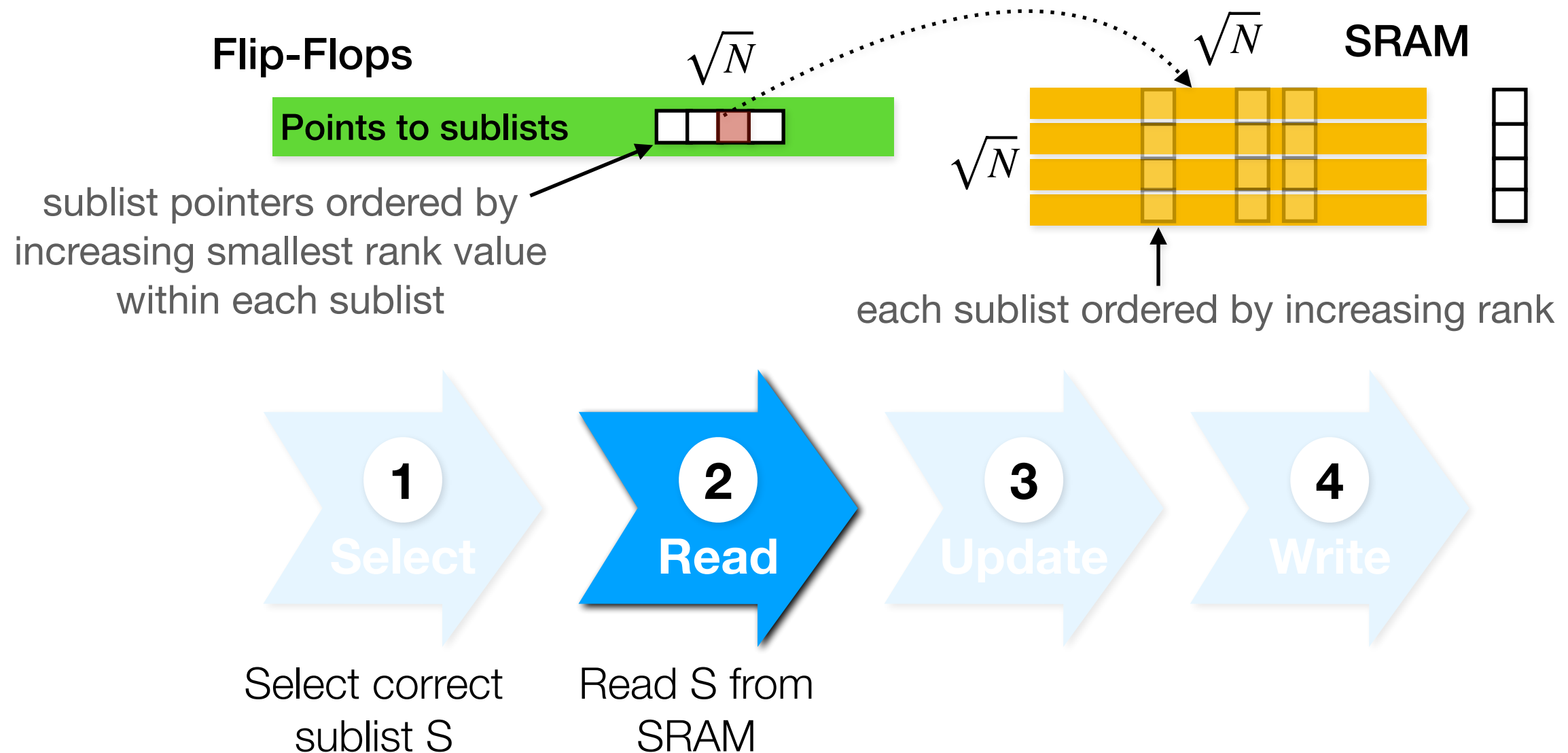
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



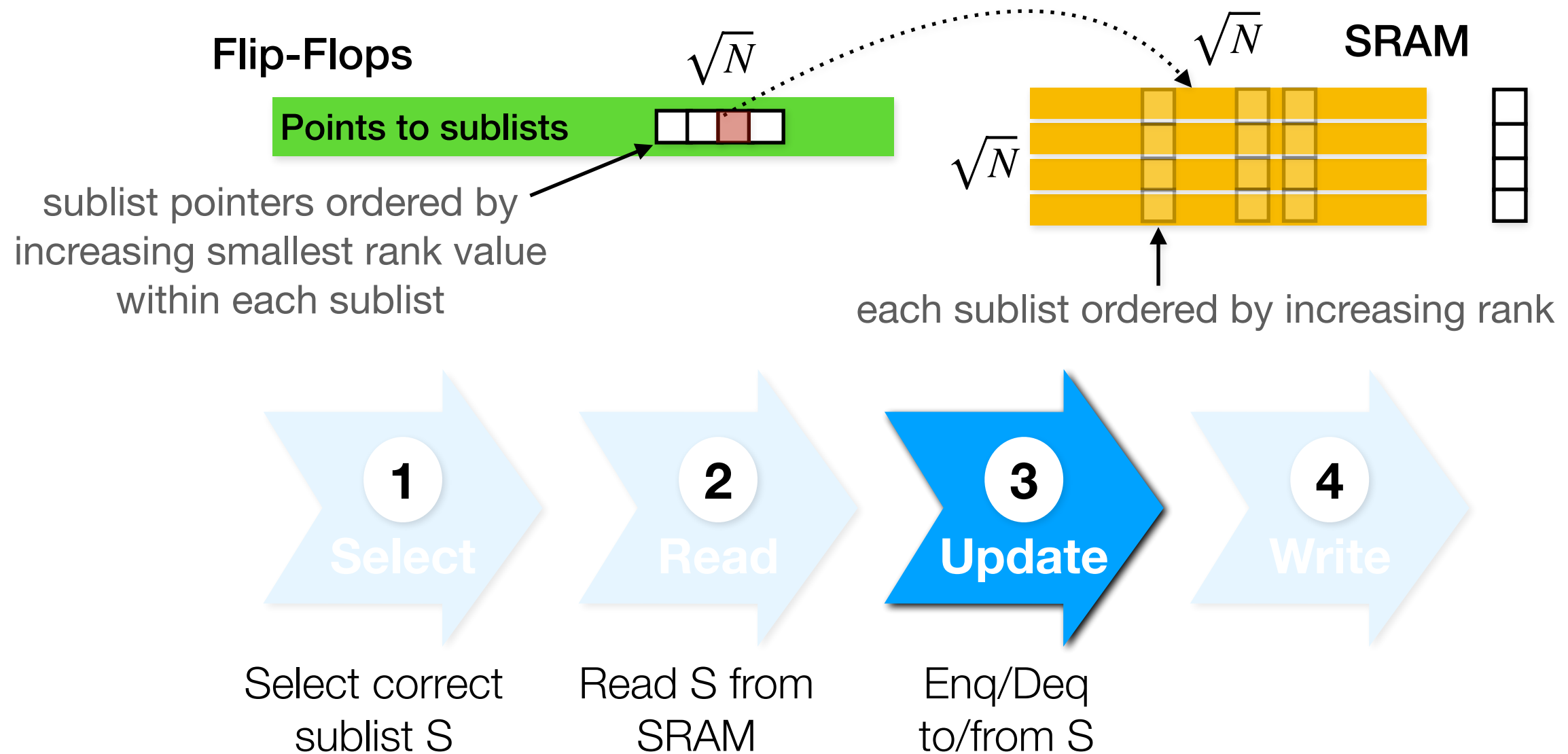
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



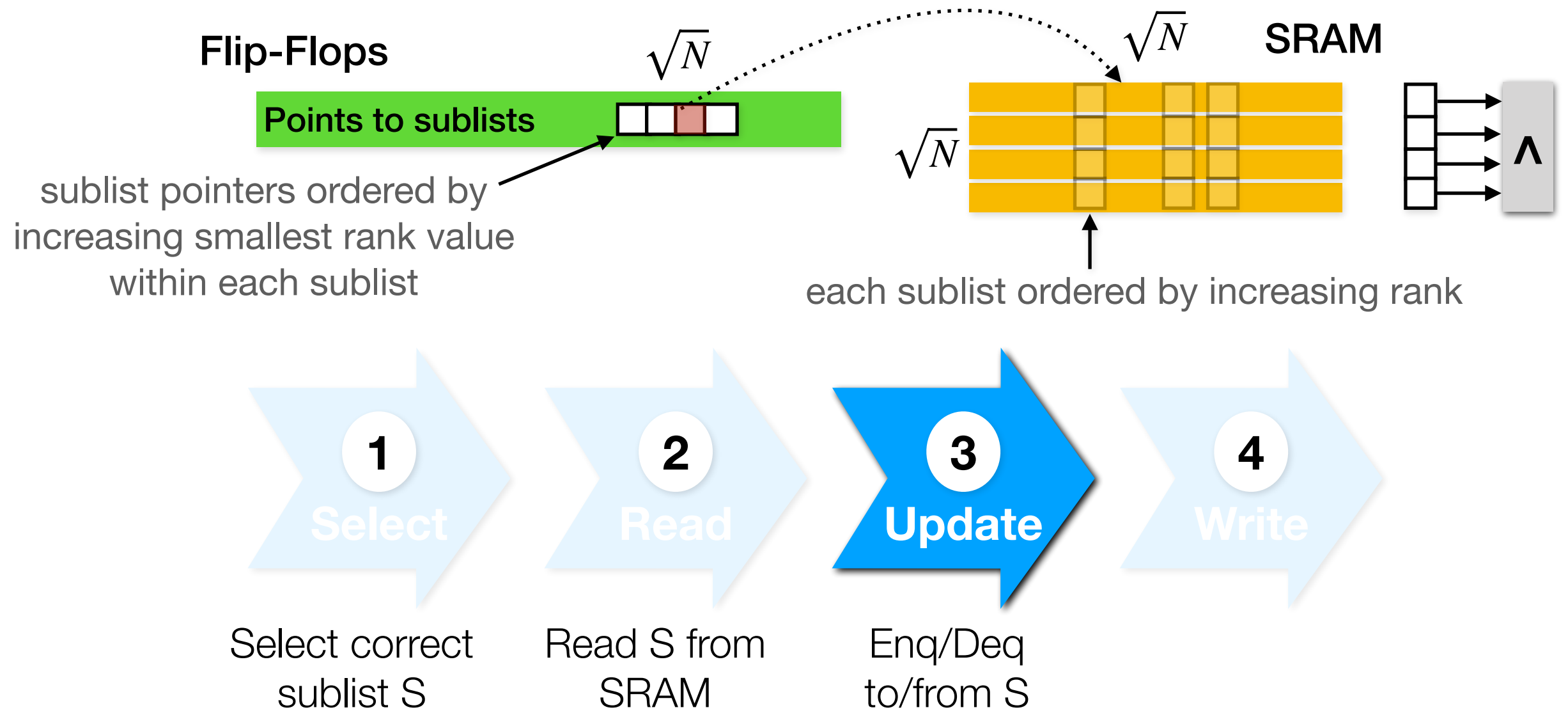
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



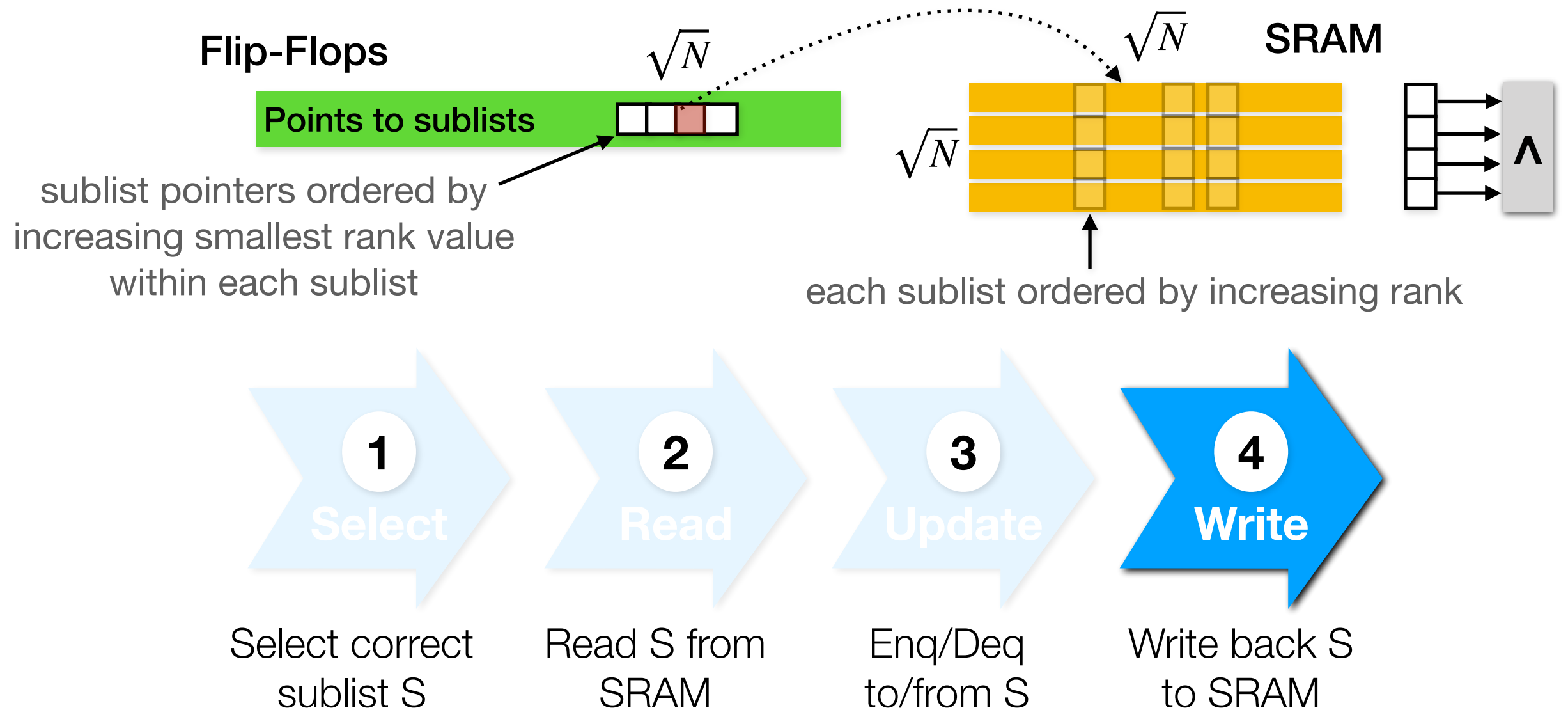
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



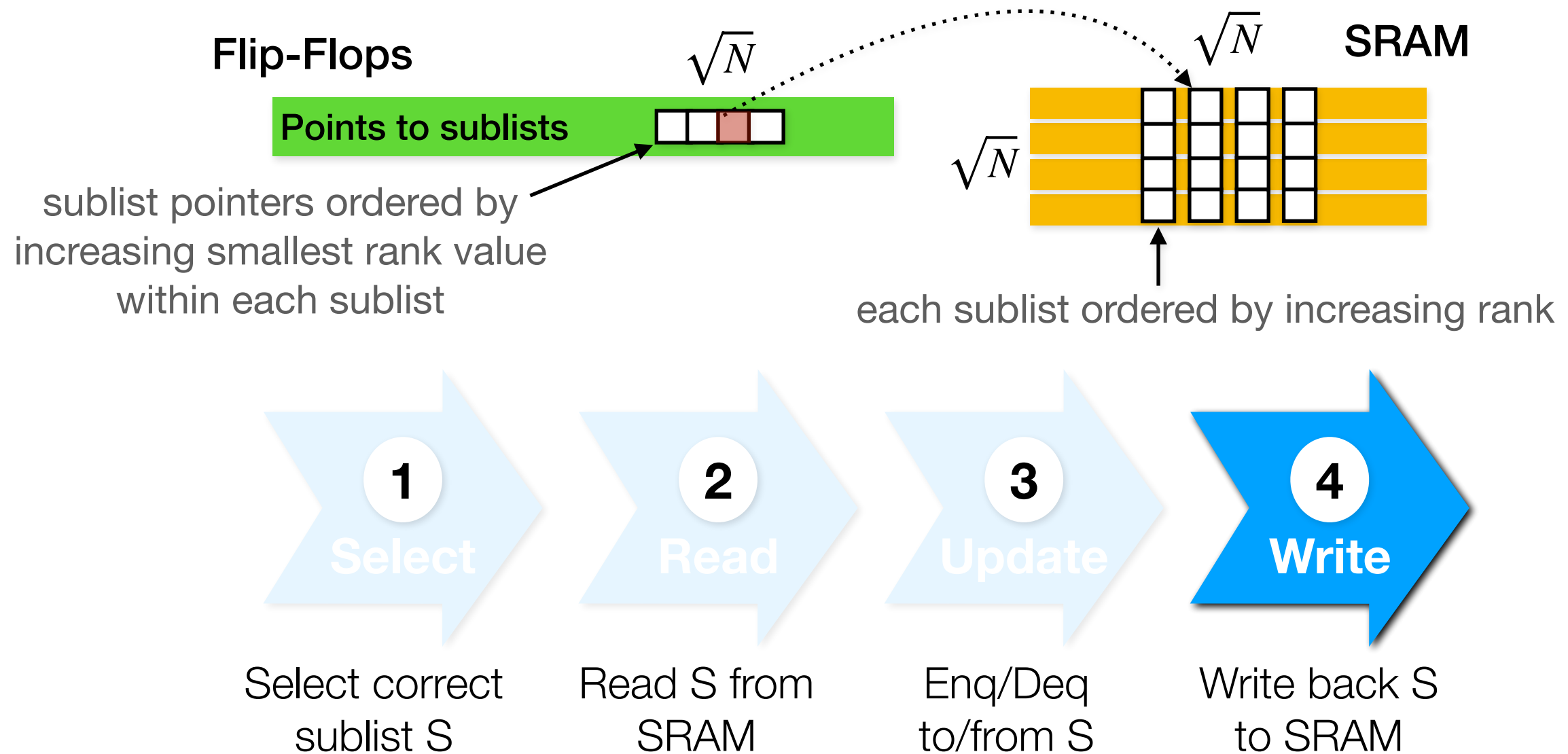
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



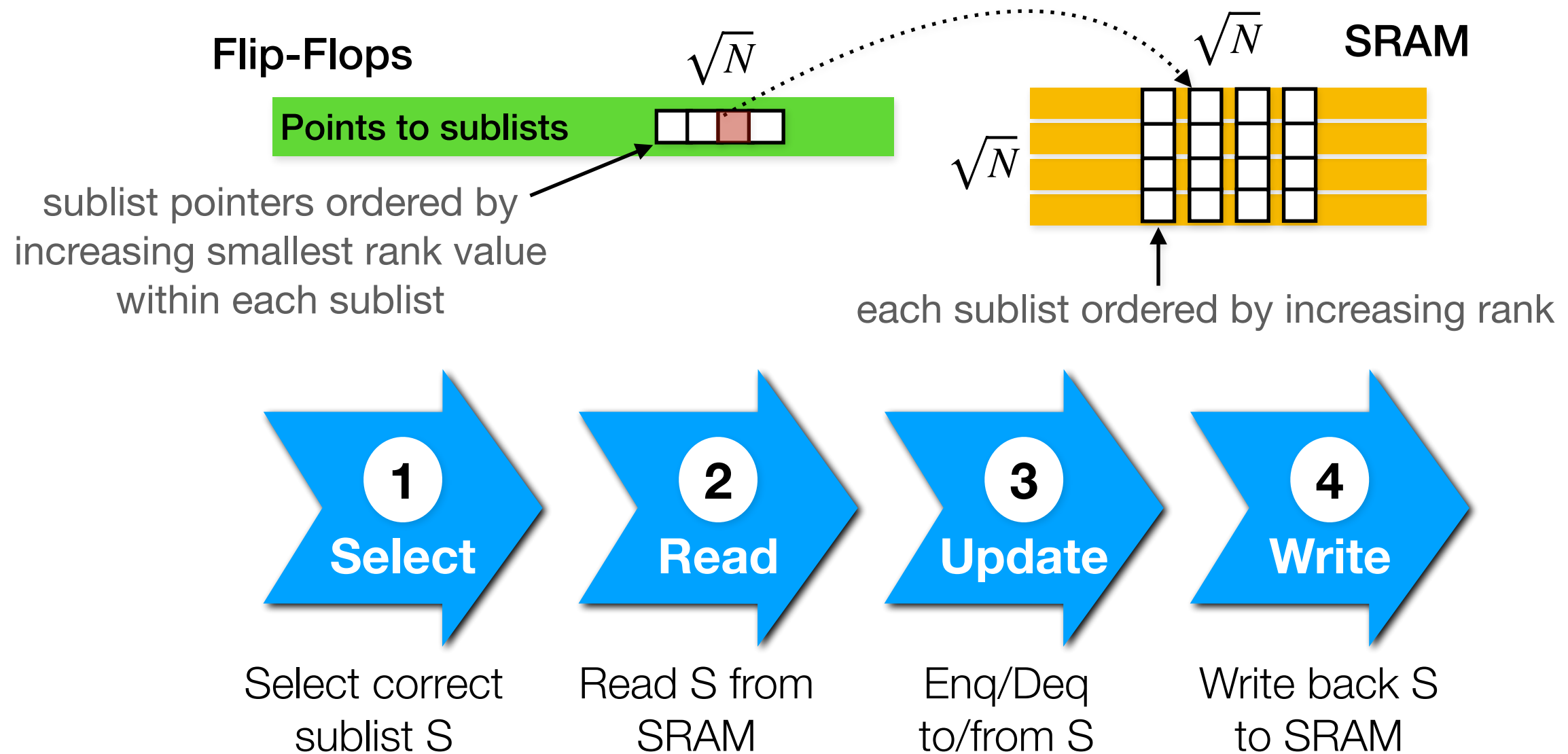
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



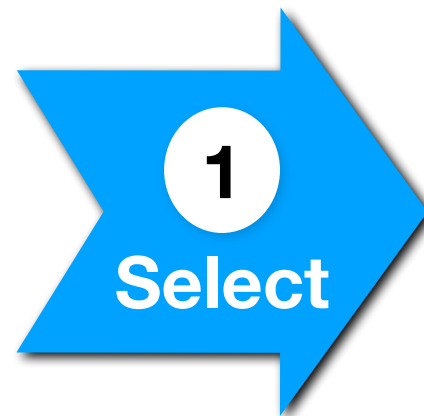
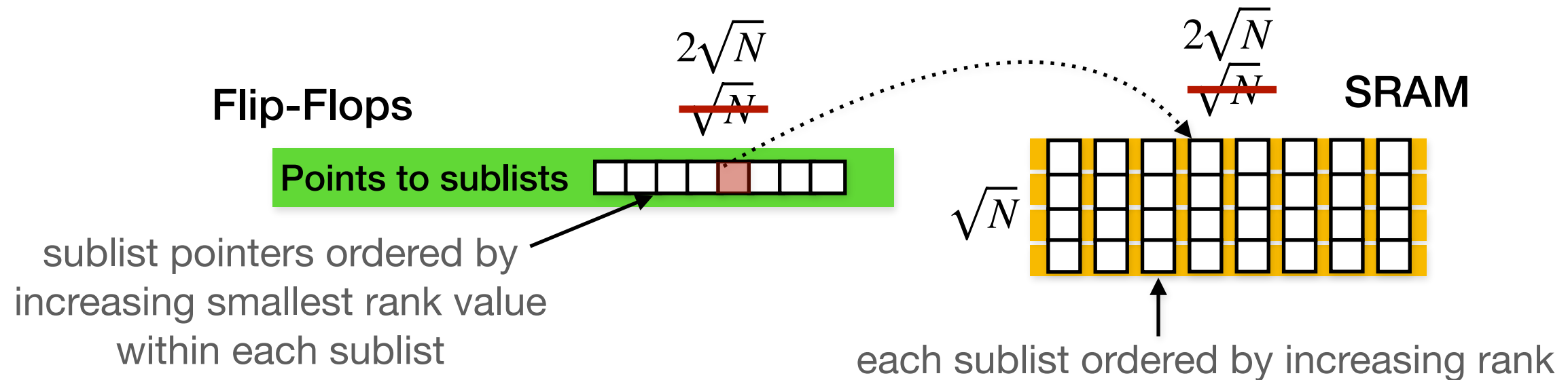
enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

Hardware Architecture



Select correct
sublist S



Read S from
SRAM



Enq/Deq
to/from S

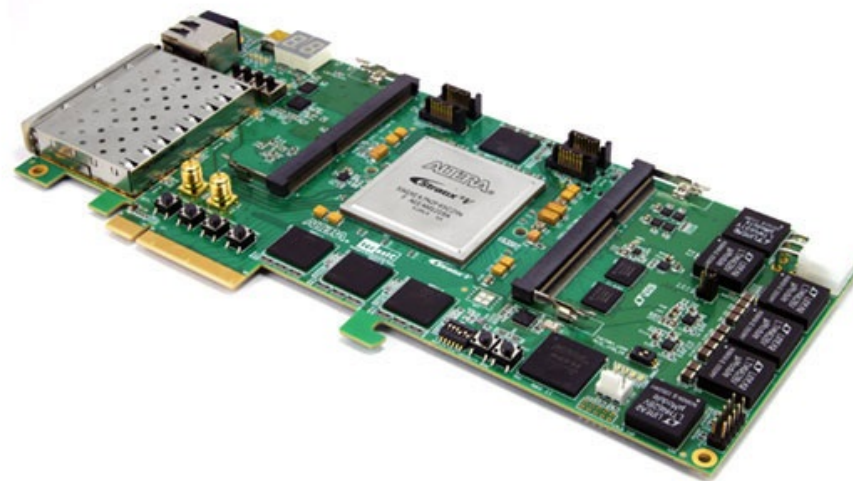


Write back S
to SRAM

enqueue(f), dequeue(), dequeue(f) each execute in exactly 4 clock cycles

... at the cost of 2x memory overhead

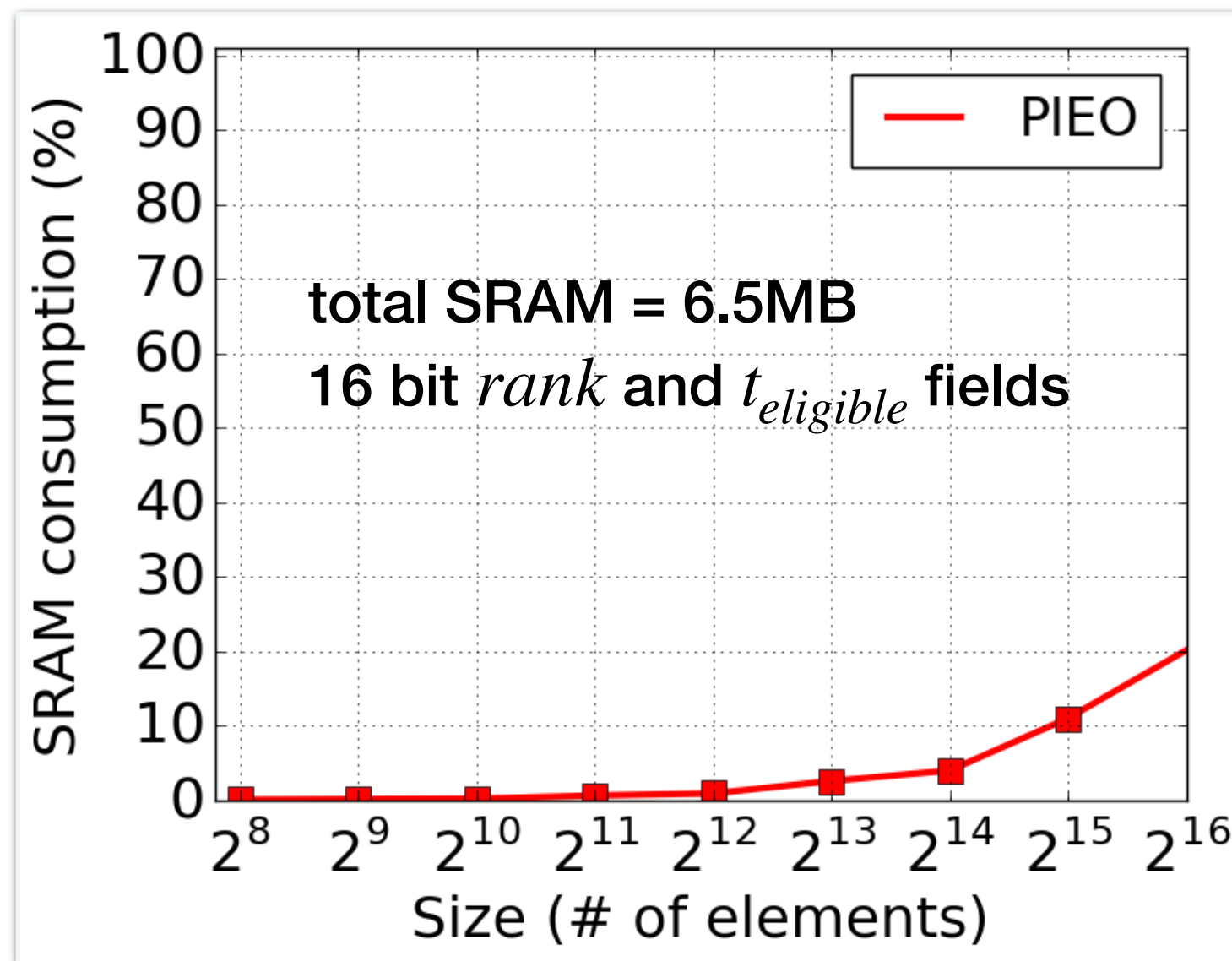
Implementation



- Implemented PIEO scheduler on a **Stratix V FPGA**
 - 234K logic modules (ALMs)
 - 6.5MB SRAM
 - 40Gbps interface bandwidth
- ~1300 LOCs in System Verilog

FPGA code for the implementation of PIEO scheduler is available at:
<https://github.com/vishal1303/PIEO-Scheduler>

Memory Consumption



Speed

PIEO on FPGA : 80 MHz (for 32K flows)

Speed

PIEO on FPGA : 80 MHz (for 32K flows) → 50ns per primitive operation

Speed

PIEO on FPGA : 80 MHz (for 32K flows) → 50ns per primitive operation

~200 Gbps scheduling rate for 1500B (MTU) packets

Speed

PIEO on FPGA : 80 MHz (for 32K flows) → 50ns per primitive operation

~200 Gbps scheduling rate for 1500B (MTU) packets

PIEO on ASIC : >1 GHz (expected)

Speed

PIEO on FPGA : 80 MHz (for 32K flows) \longrightarrow 50ns per primitive operation

~200 Gbps scheduling rate for 1500B (MTU) packets

PIEO on ASIC : >1 GHz (expected) \longrightarrow <4ns per primitive operation

Speed

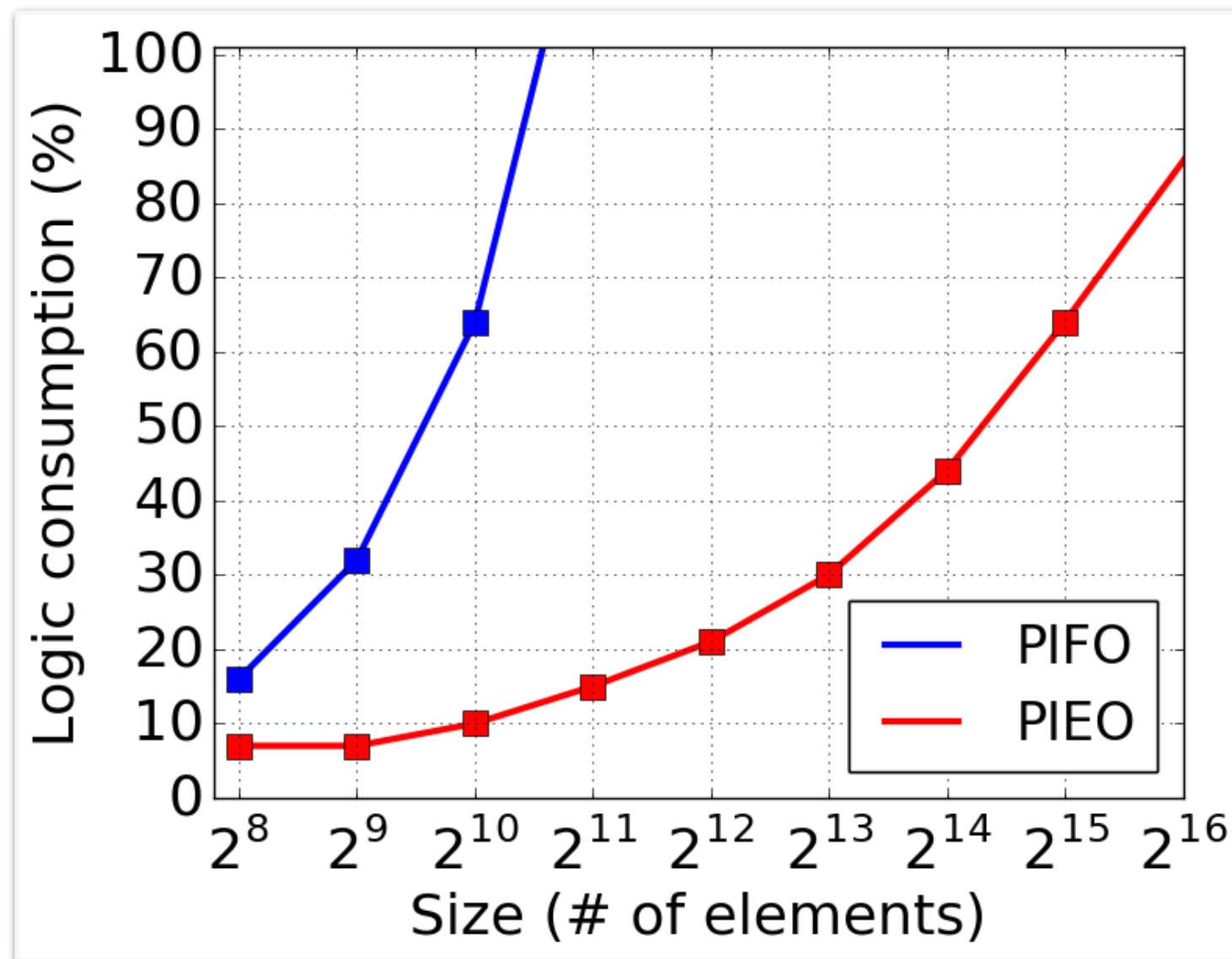
PIEO on FPGA : 80 MHz (for 32K flows) \longrightarrow 50ns per primitive operation

~200 Gbps scheduling rate for 1500B (MTU) packets

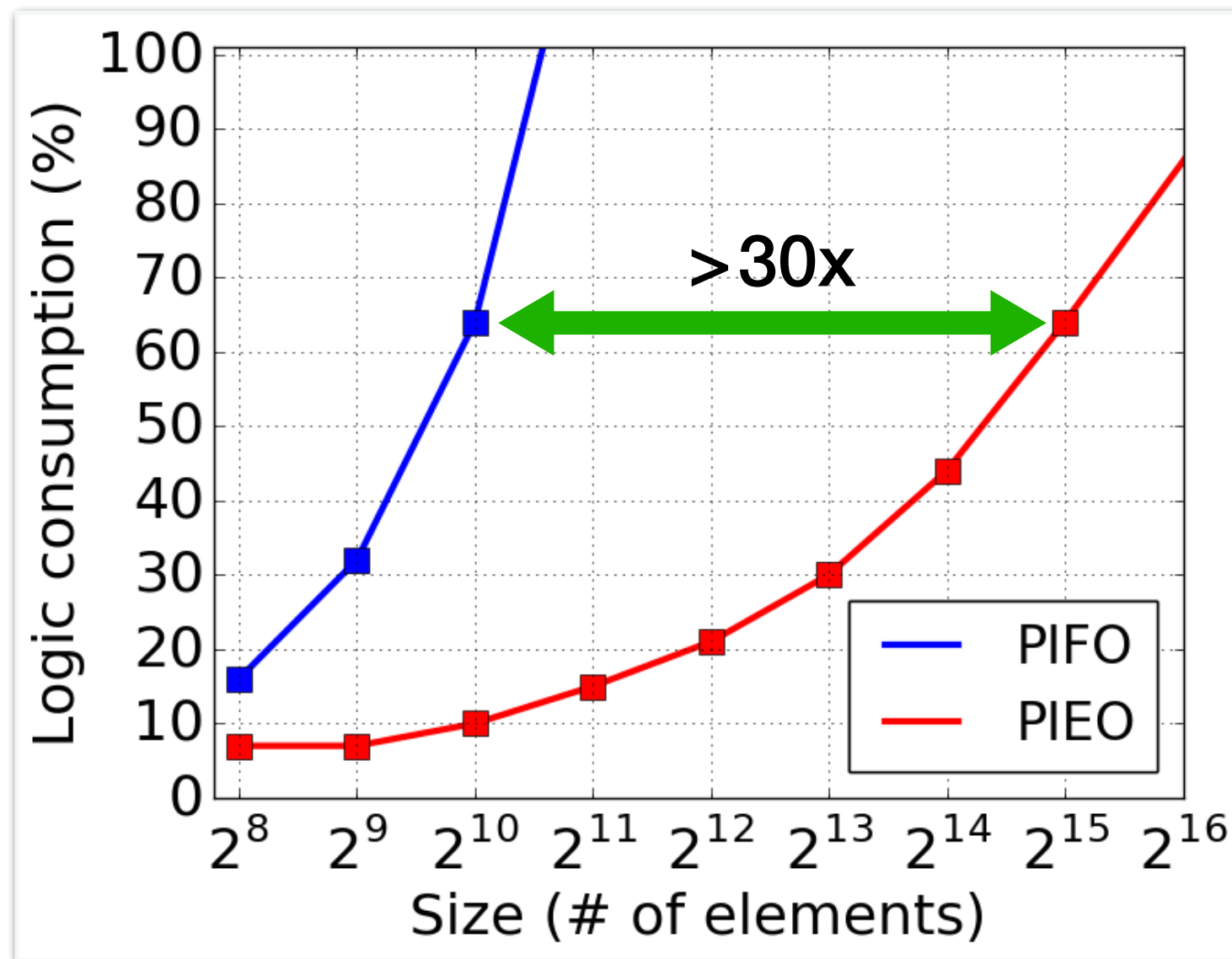
PIEO on ASIC : >1 GHz (expected) \longrightarrow <4ns per primitive operation

>3 Tbps scheduling rate for 1500B (MTU) packets

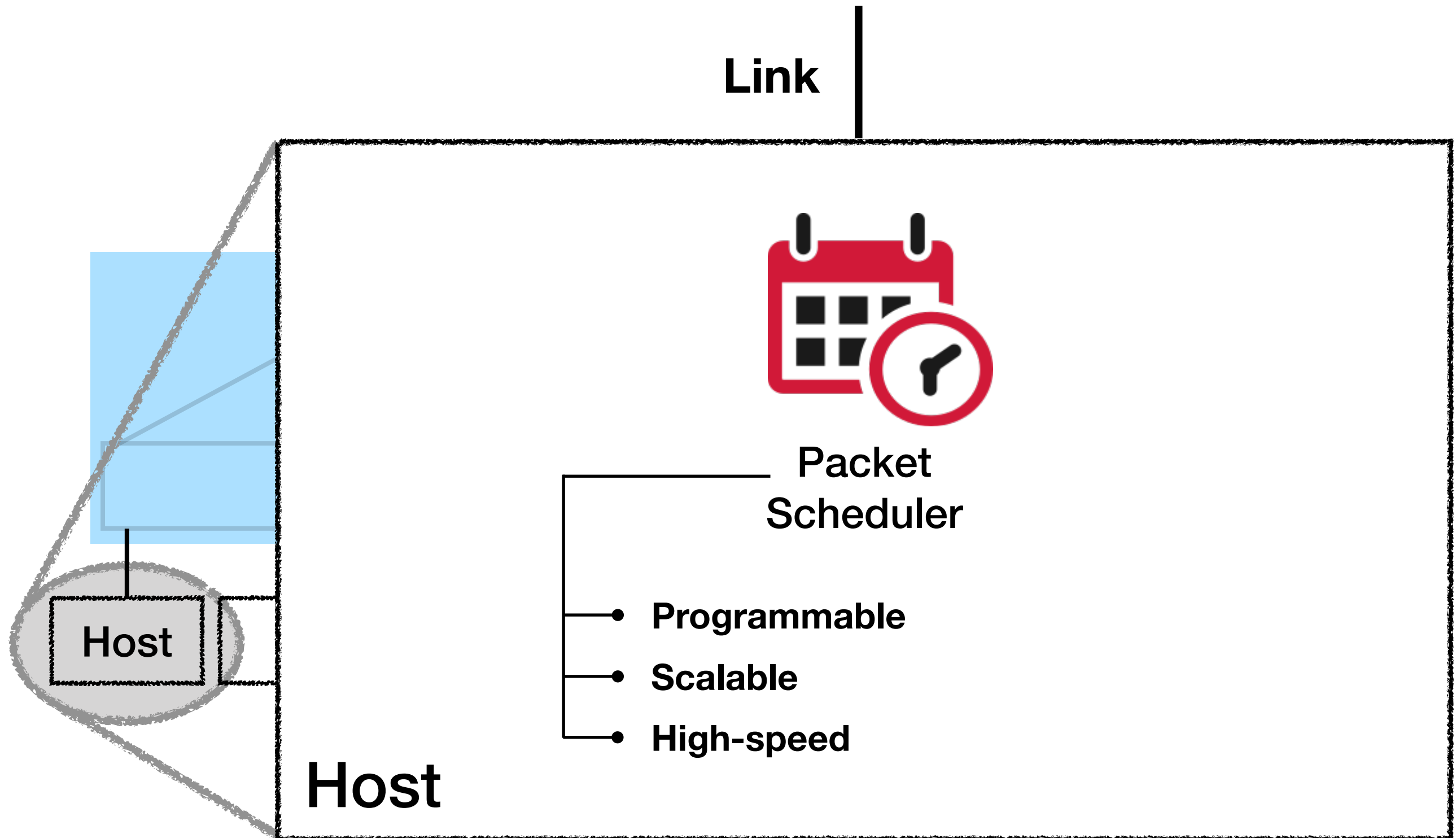
Scalability



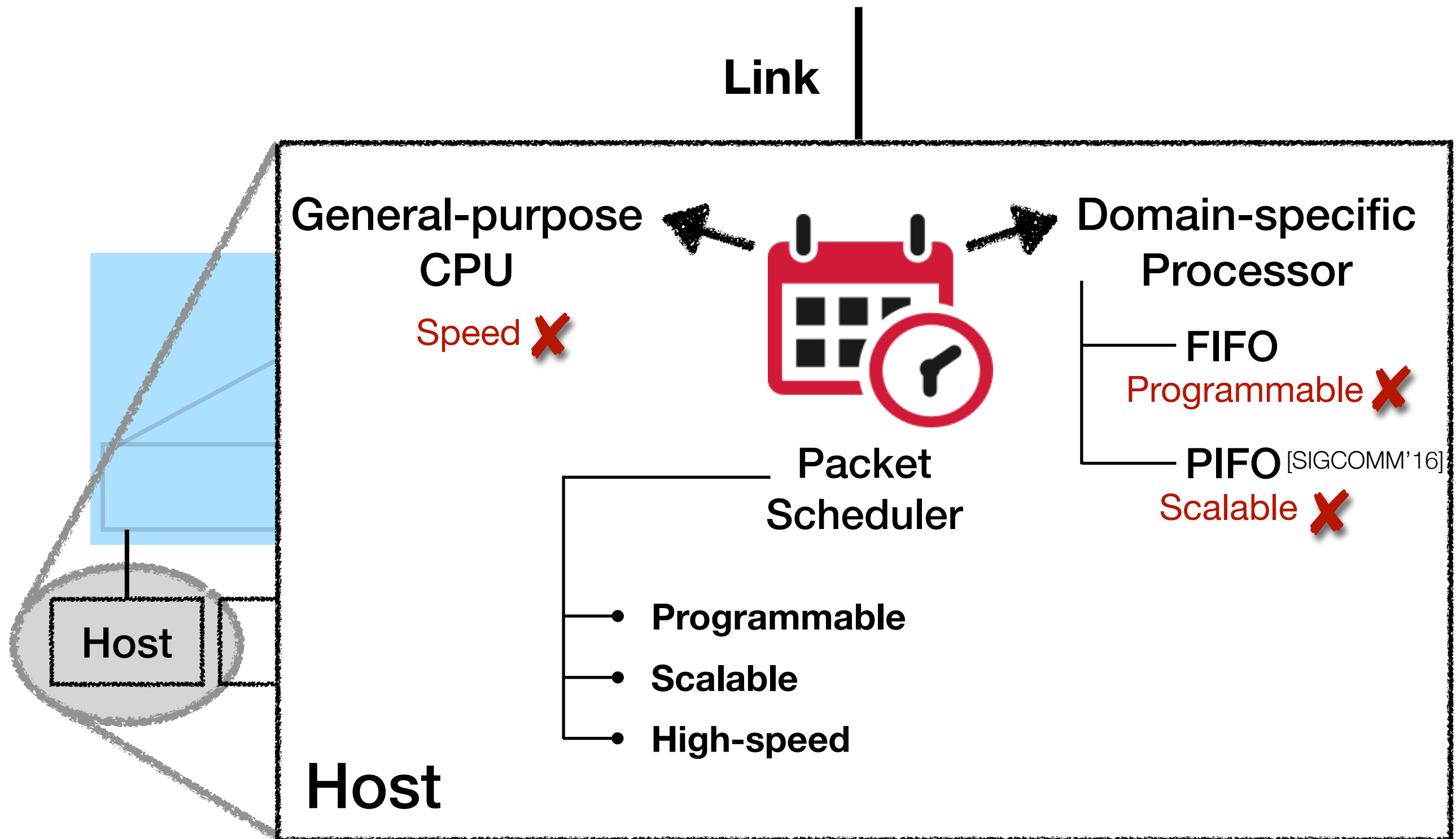
Scalability



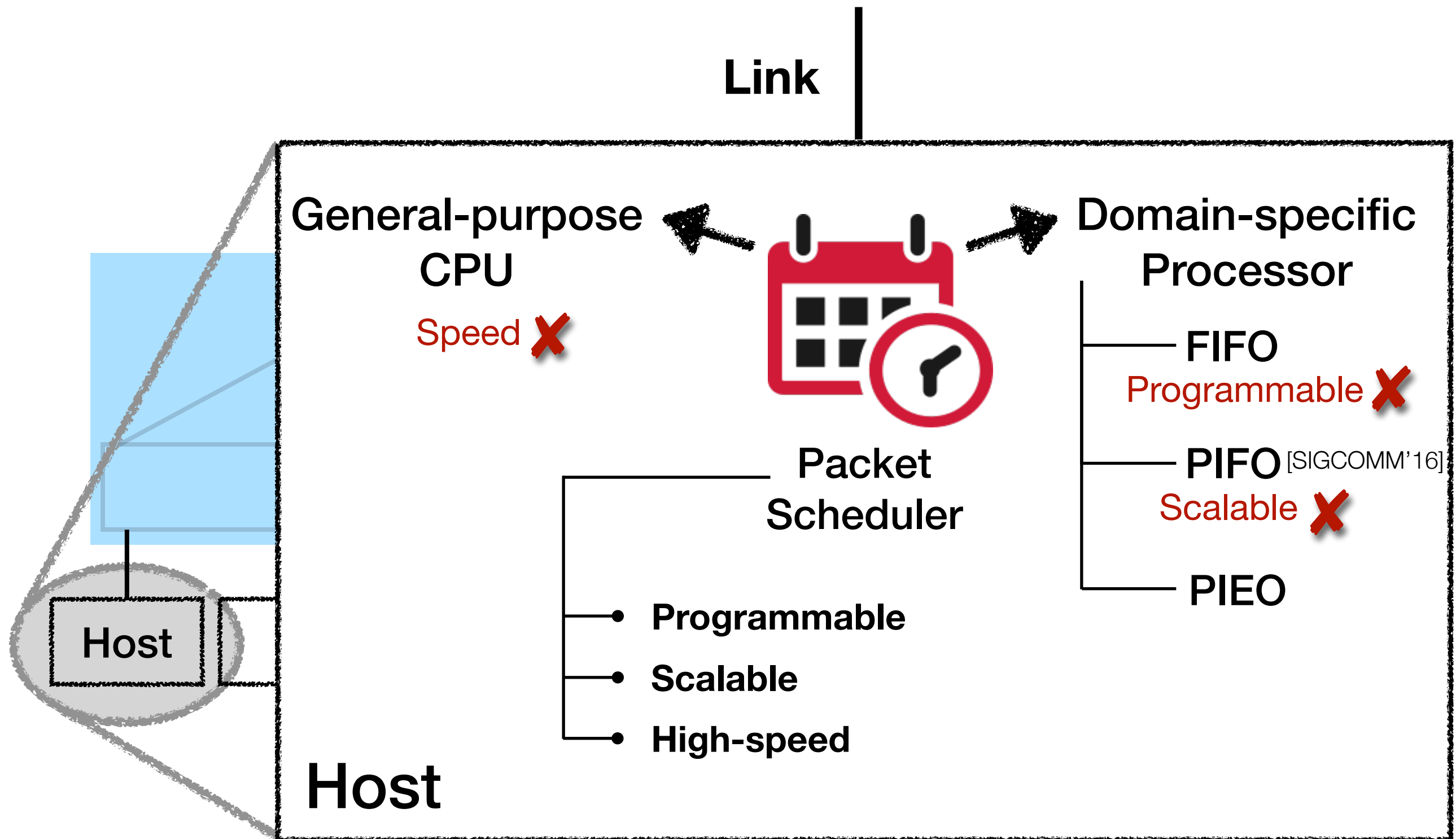
Part I : Summary



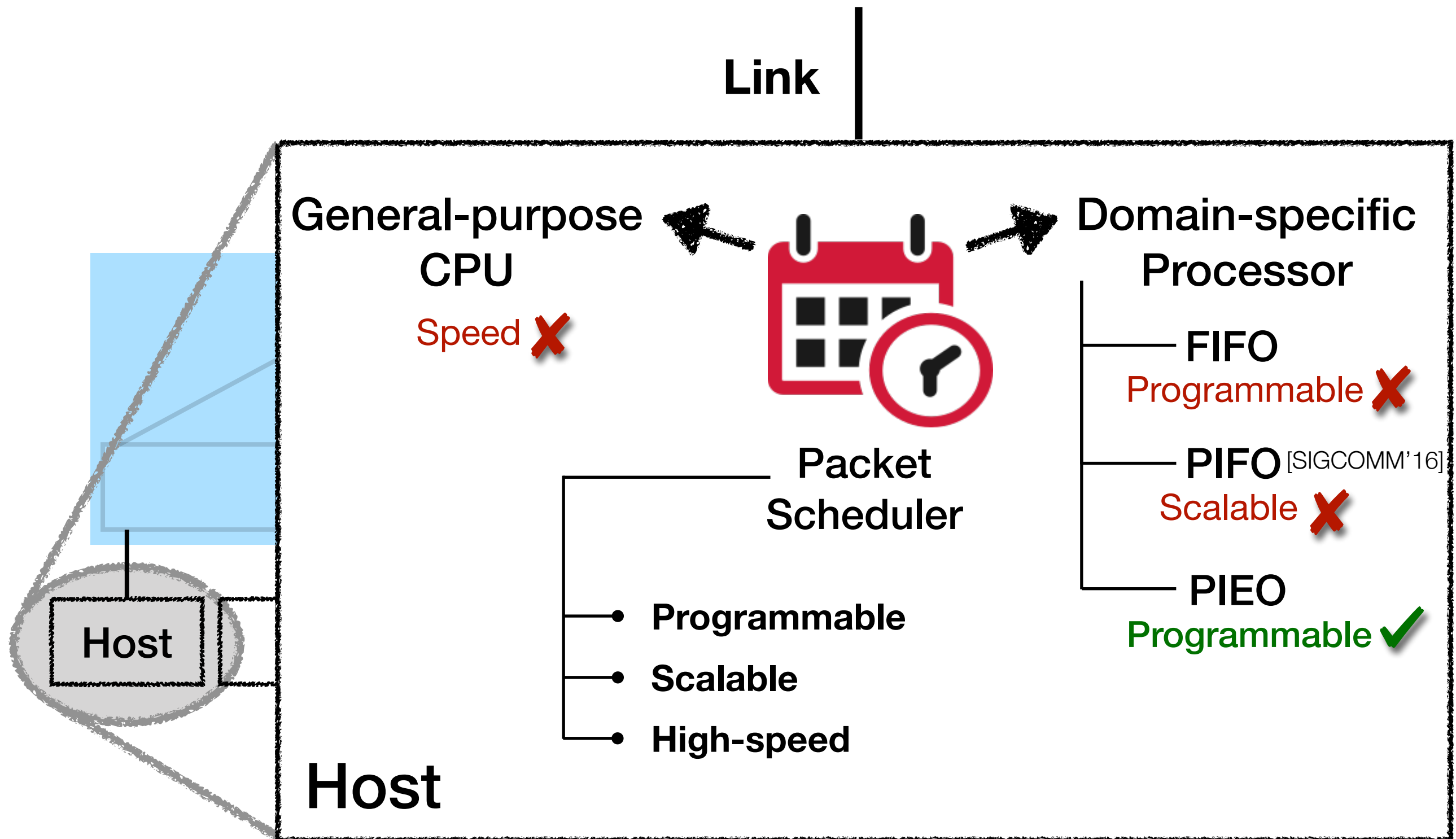
Part I : Summary



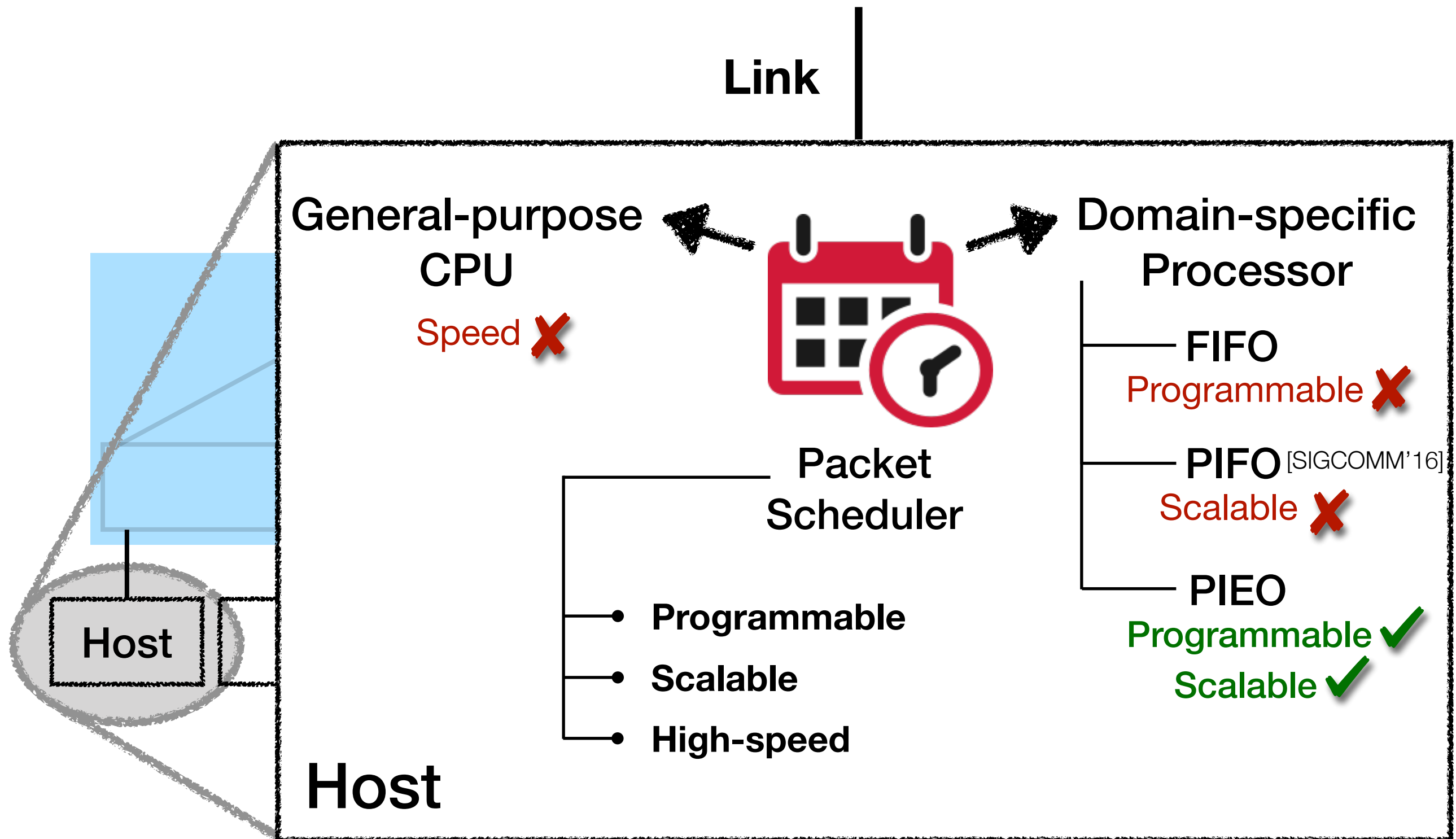
Part I : Summary



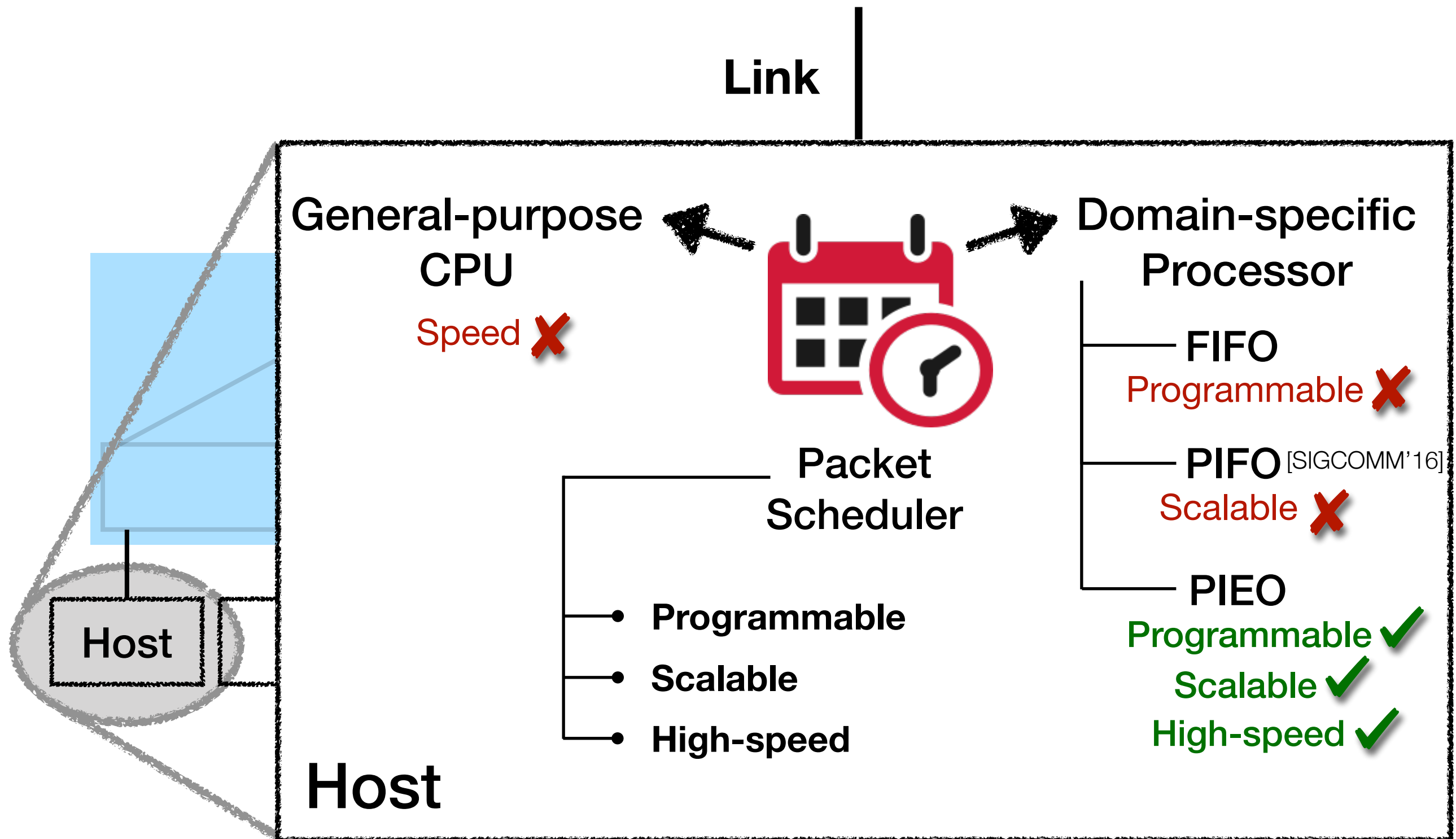
Part I : Summary



Part I : Summary

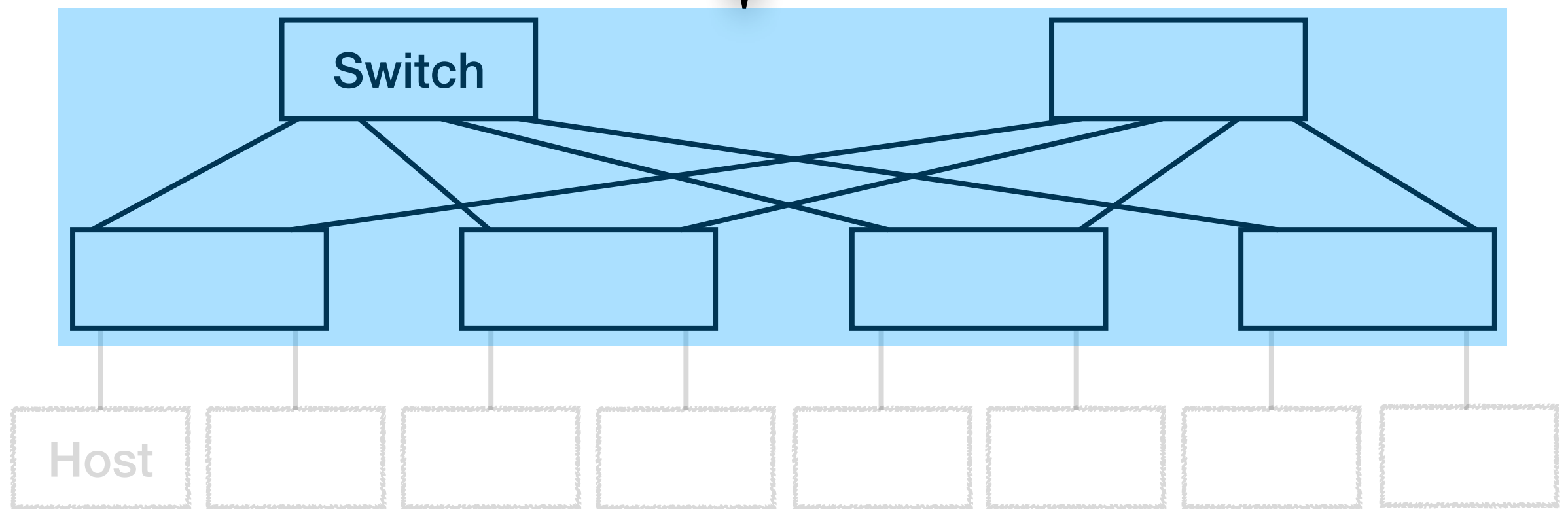


Part I : Summary



Part II : Switching Fabric

How to build high-speed switching fabric?



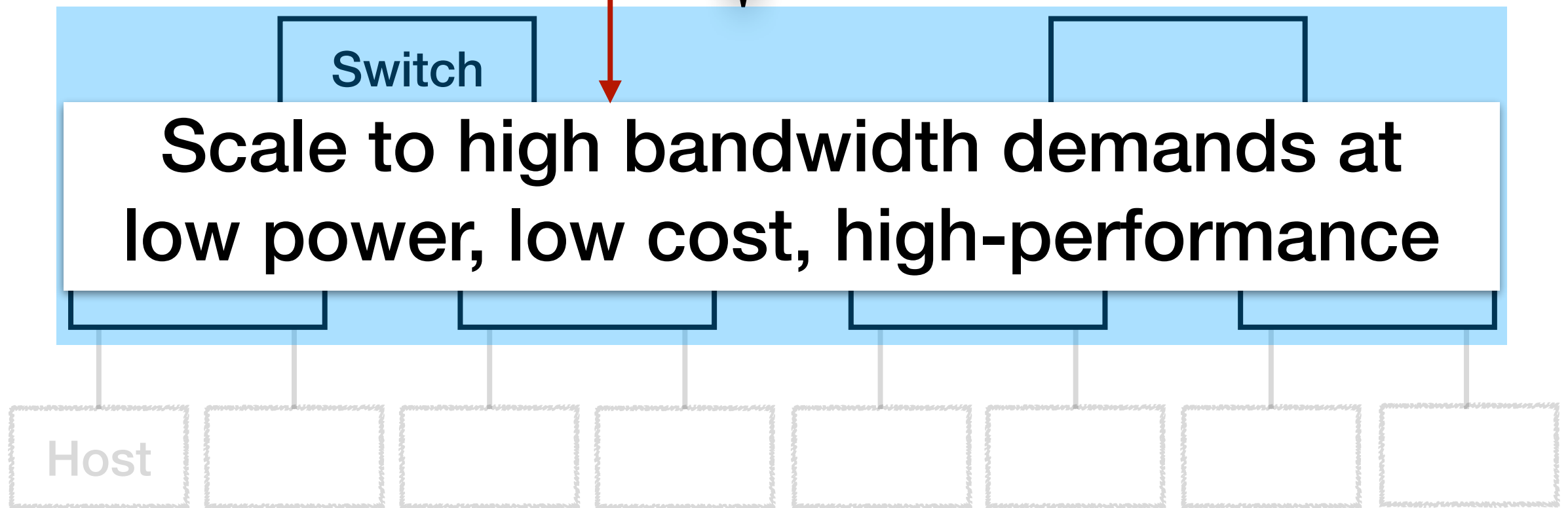
Part II : Switching Fabric

How to build high-speed switching fabric?

Switch

**Scale to high bandwidth demands at
low power, low cost, high-performance**

Host



Part II : Switching Fabric

How to build high-speed switching fabric?

Switch

Scale to high bandwidth demands at
low power, low cost, high-performance

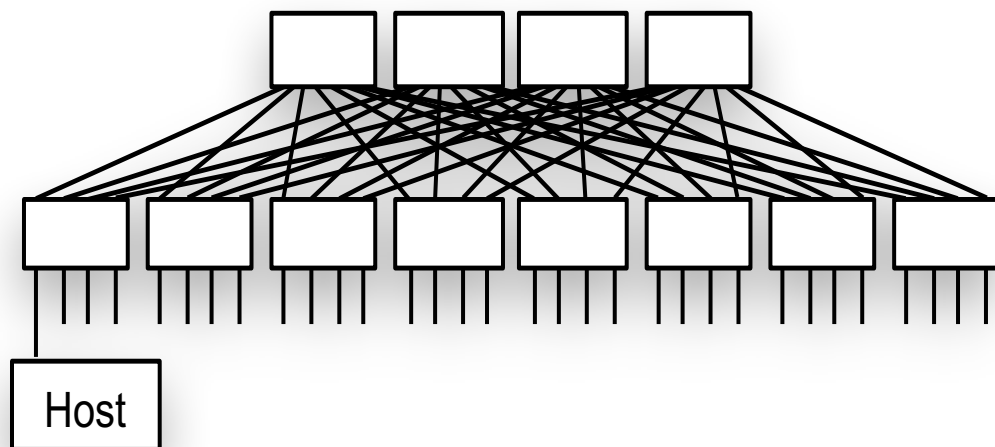
Host

Circuit Switching



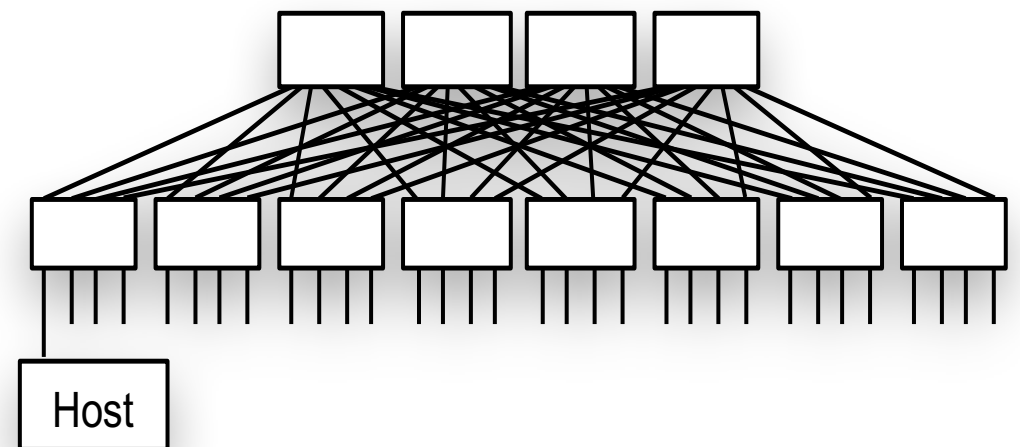
Promise of Circuit Switching

Packet Switching



VS.

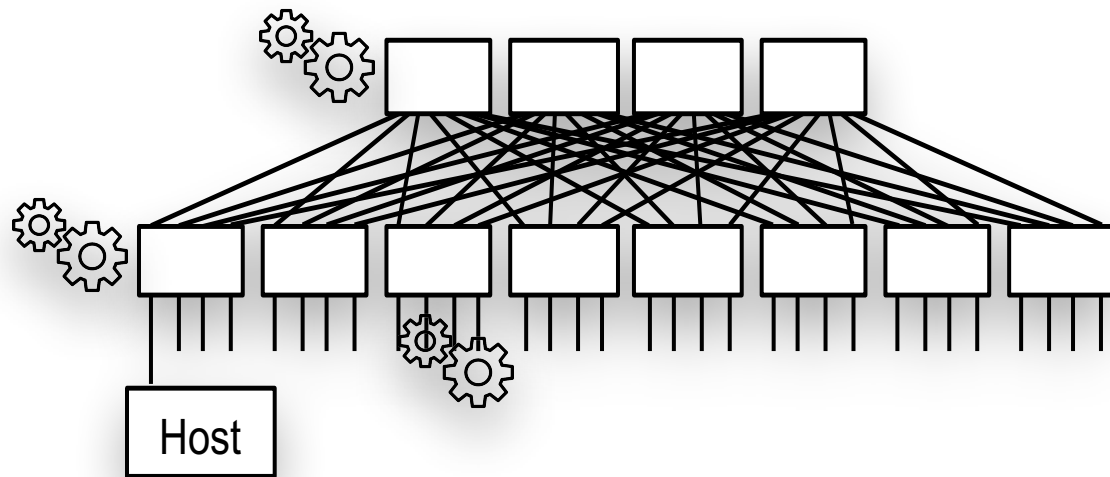
Circuit Switching



Promise of Circuit Switching

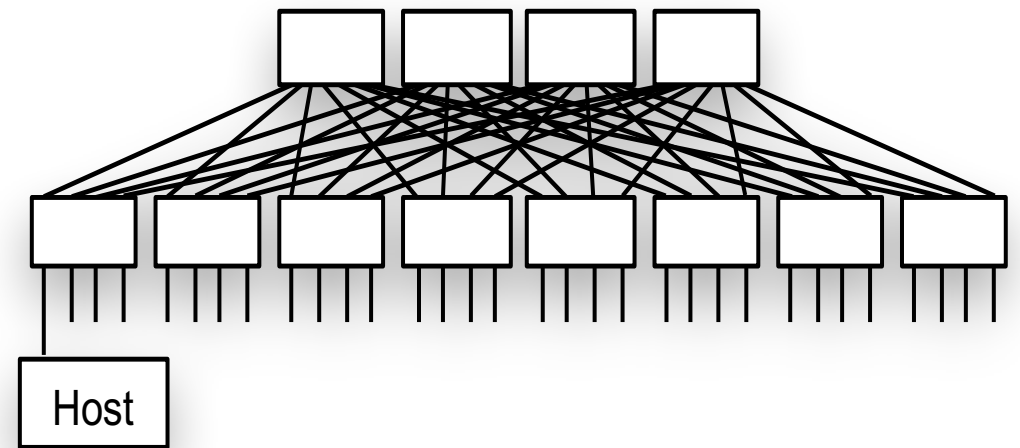
Packet Switching

(Per packet, Per hop Decisions)



VS.

Circuit Switching



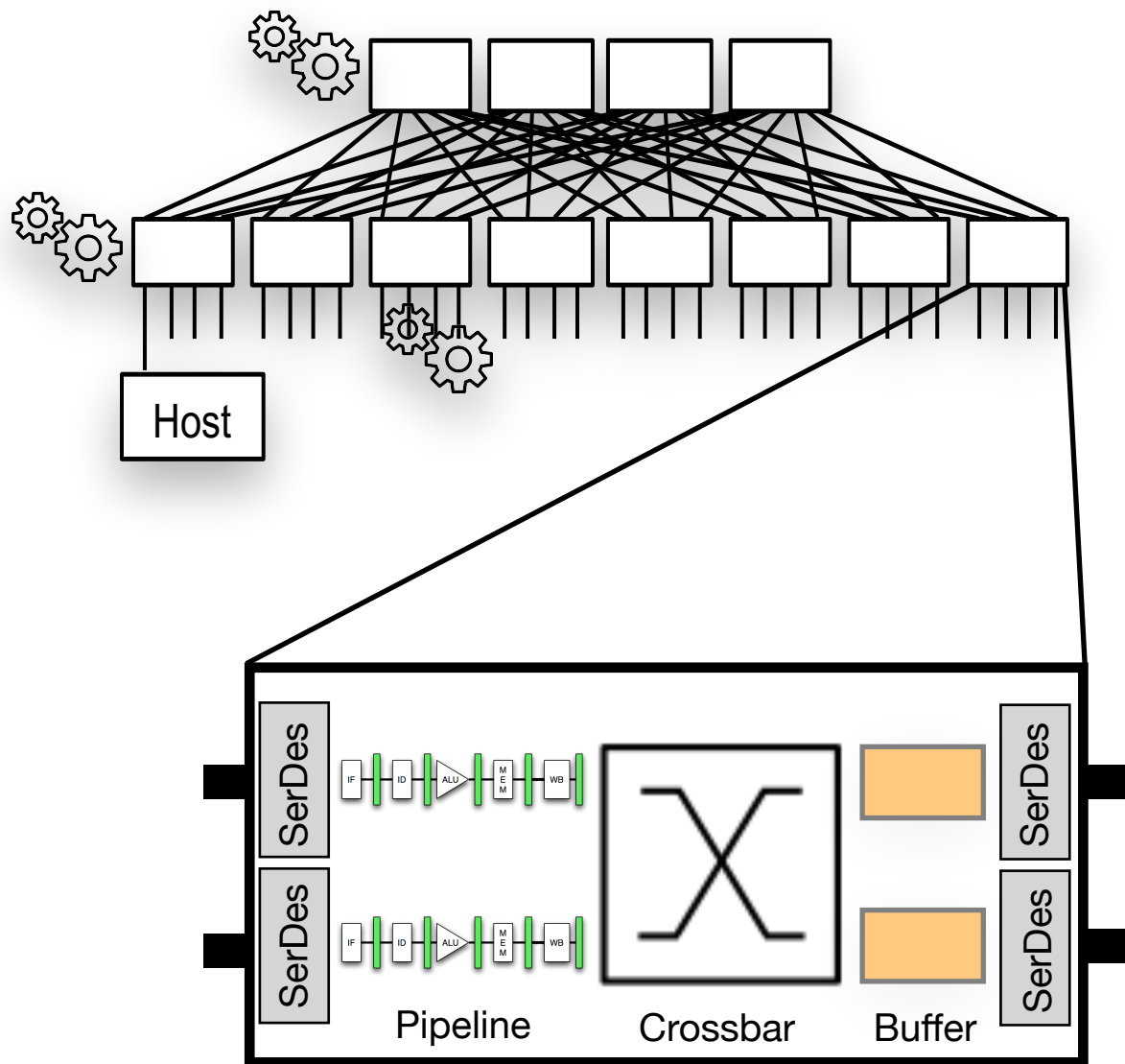
Promise of Circuit Switching

Packet Switching

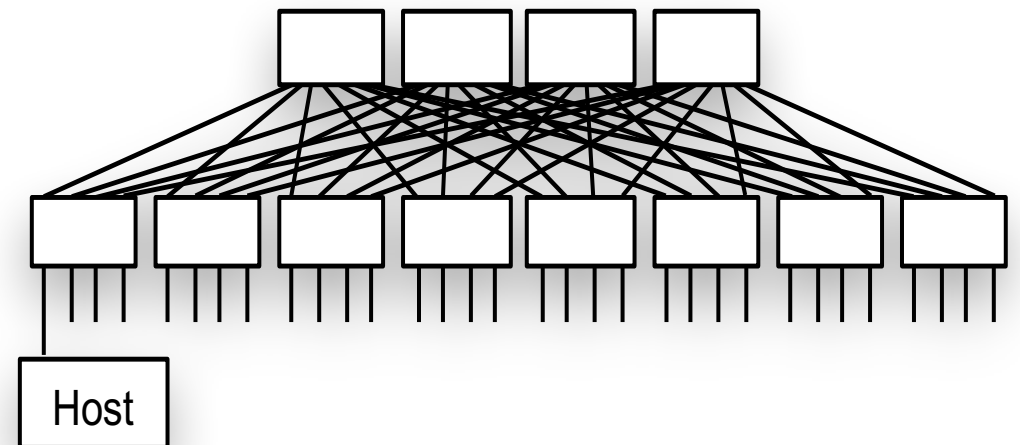
(Per packet, Per hop Decisions)

VS.

Circuit Switching



Packet Switch



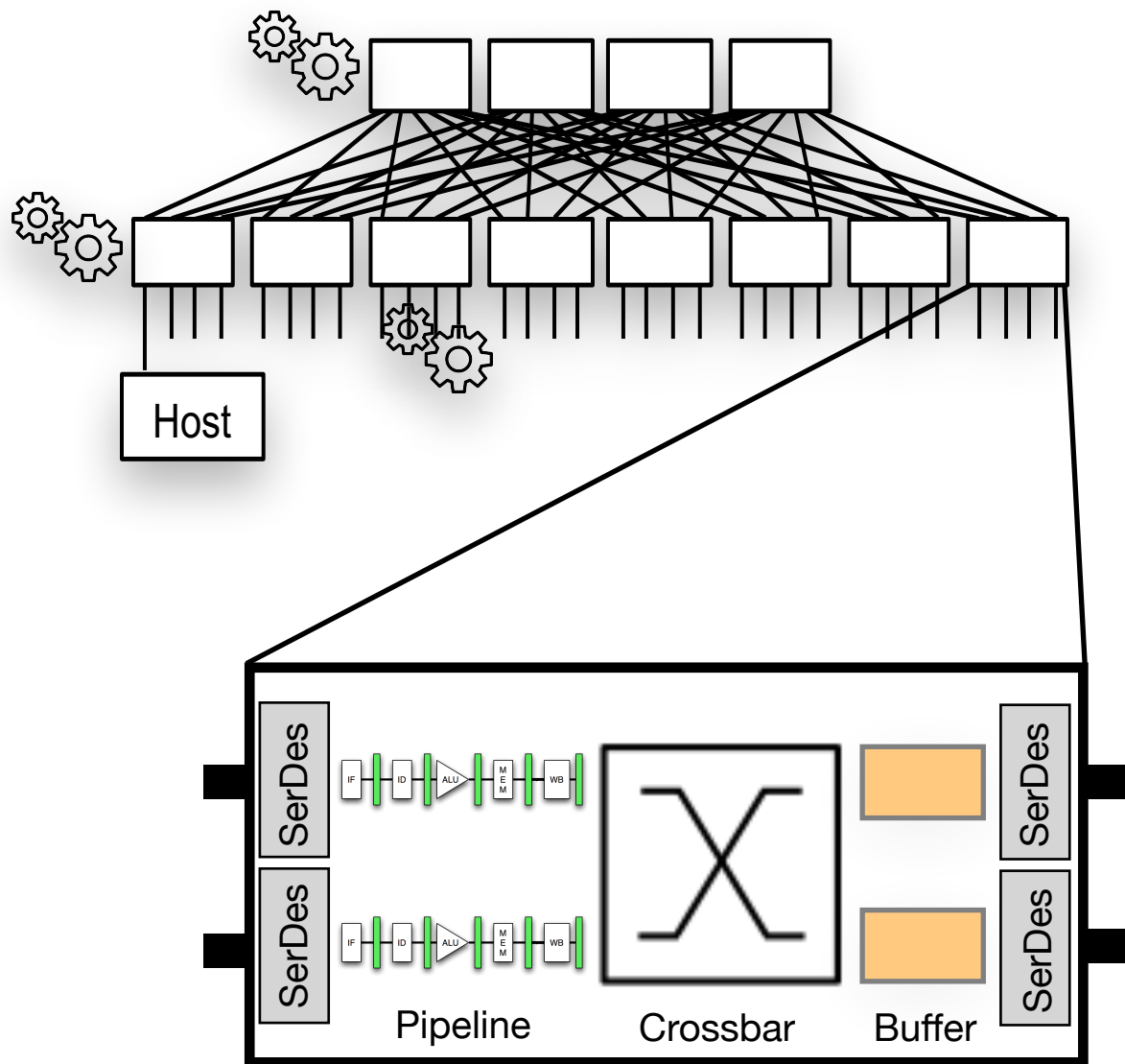
Promise of Circuit Switching

Packet Switching

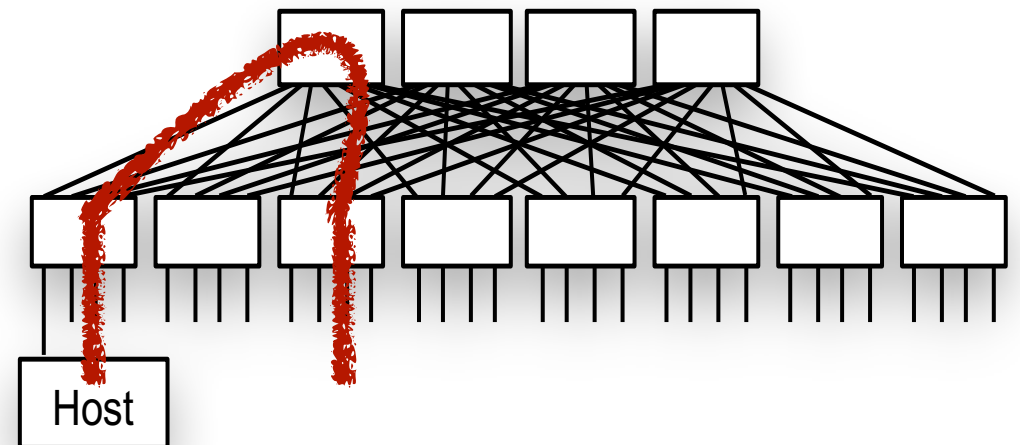
(Per packet, Per hop Decisions)

VS.

Circuit Switching



Packet Switch



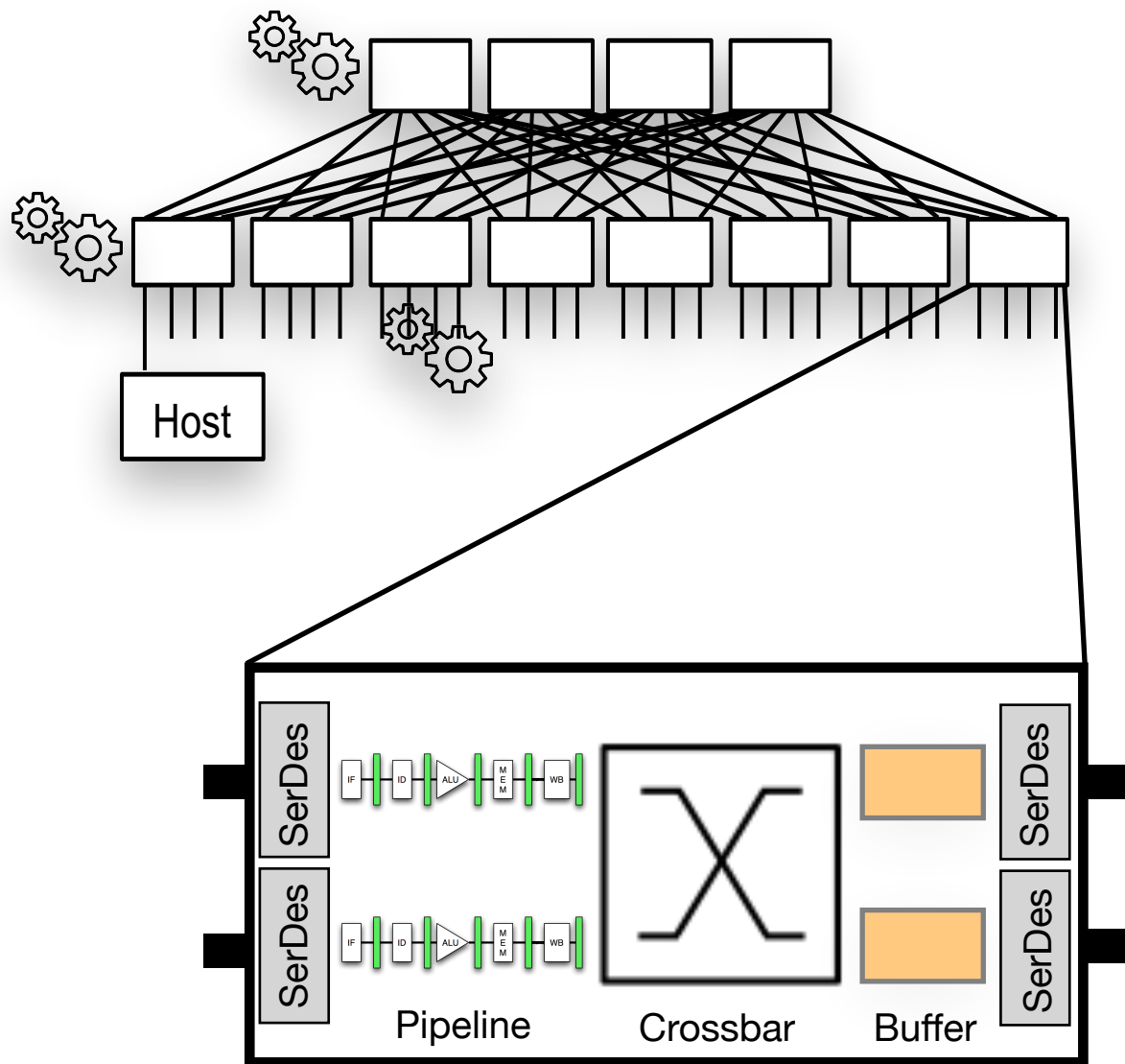
Promise of Circuit Switching

Packet Switching

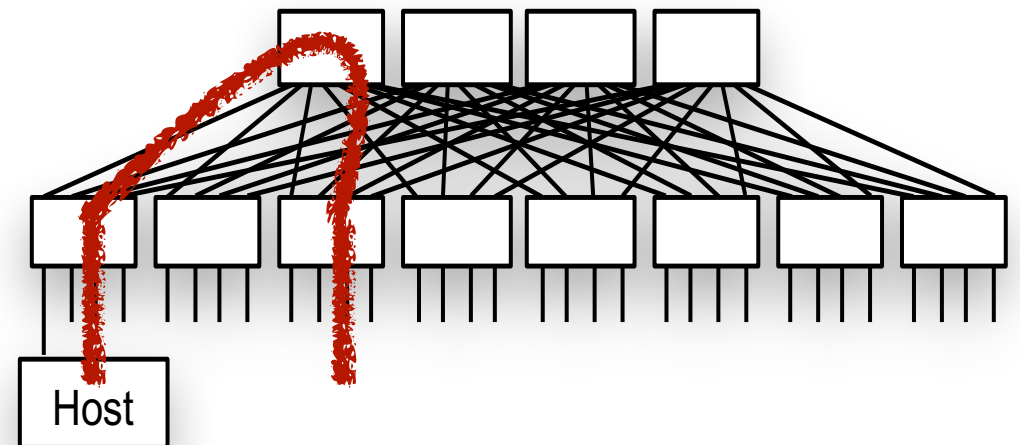
(Per packet, Per hop Decisions)

VS.

Circuit Switching



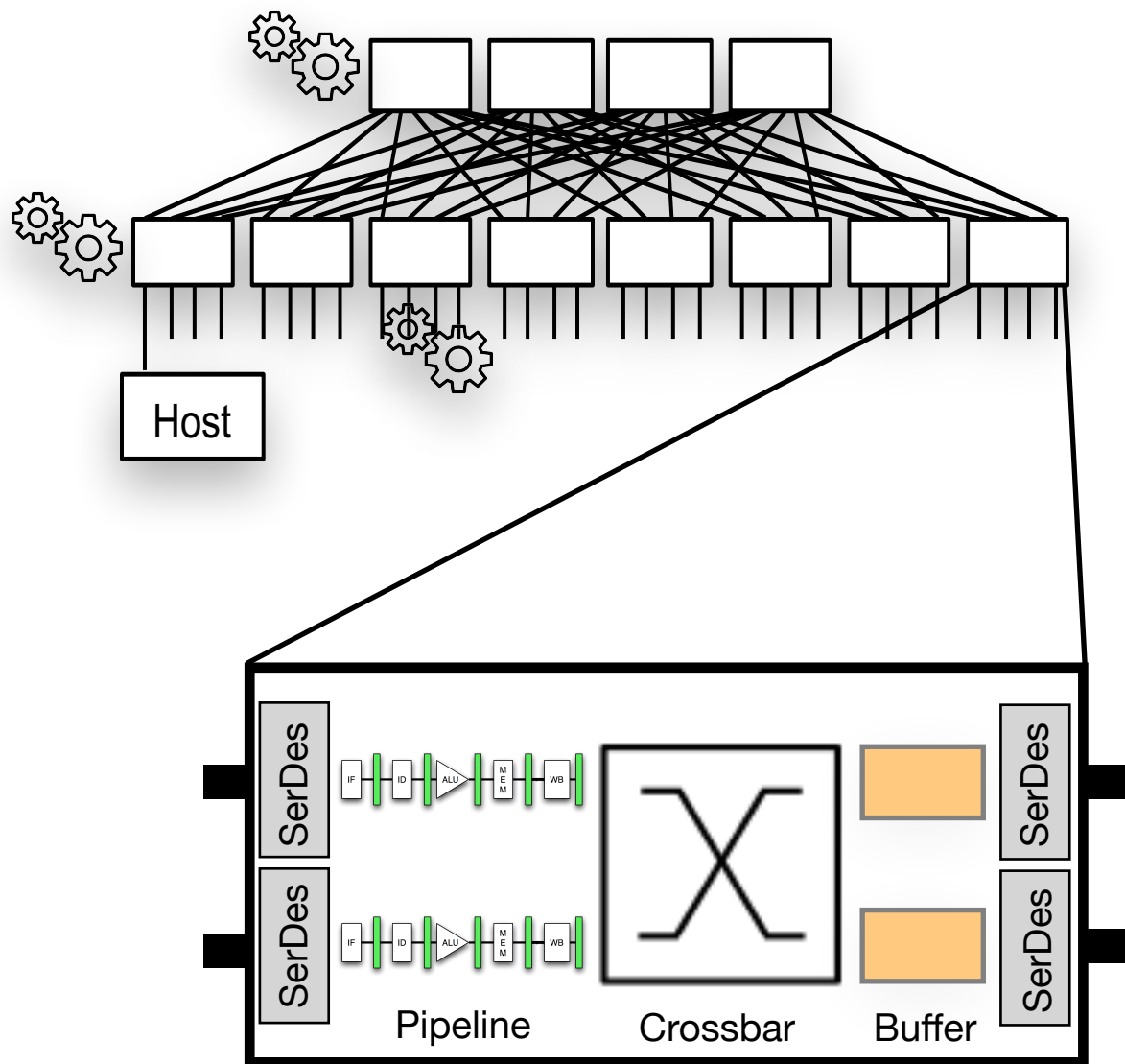
Packet Switch



Promise of Circuit Switching

Packet Switching

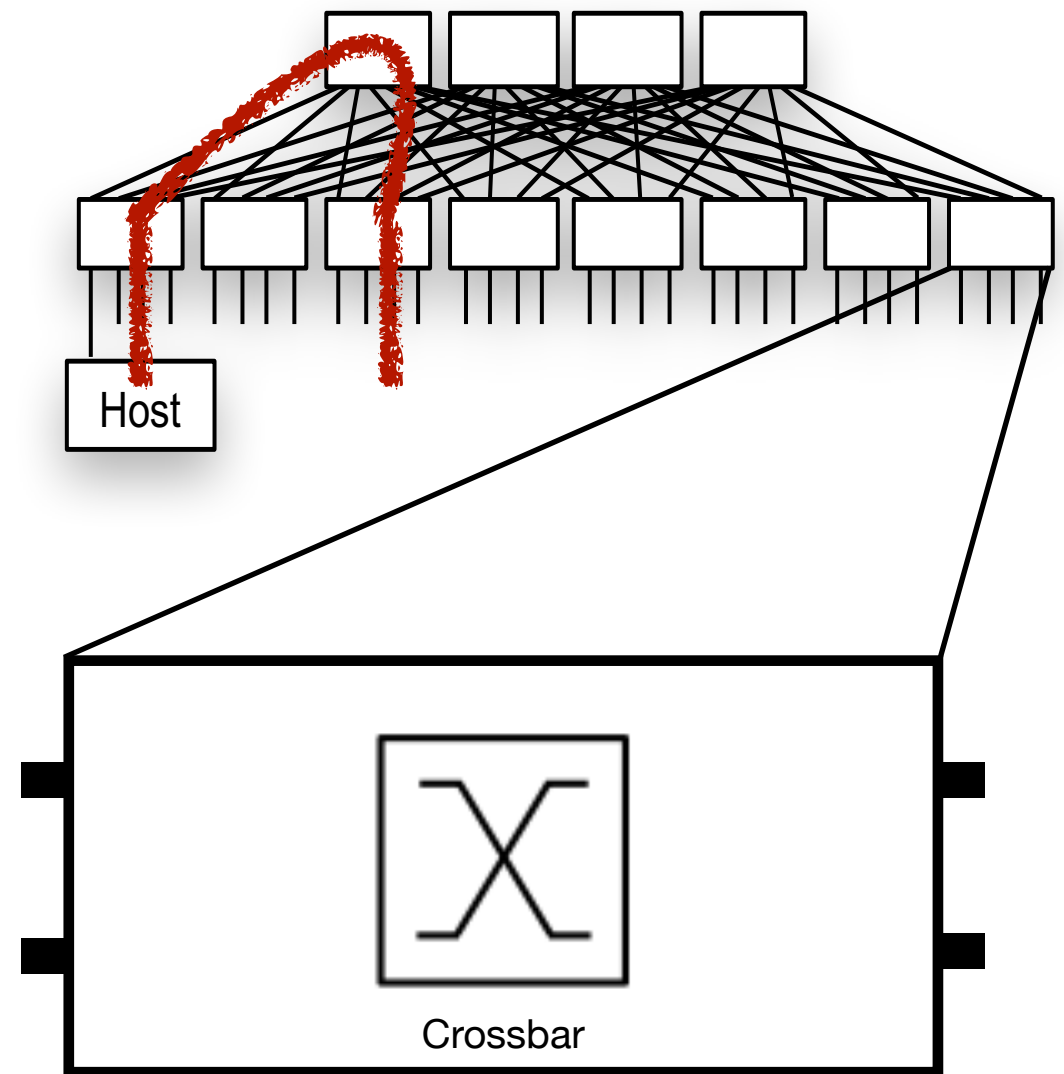
(Per packet, Per hop Decisions)



Packet Switch

VS.

Circuit Switching

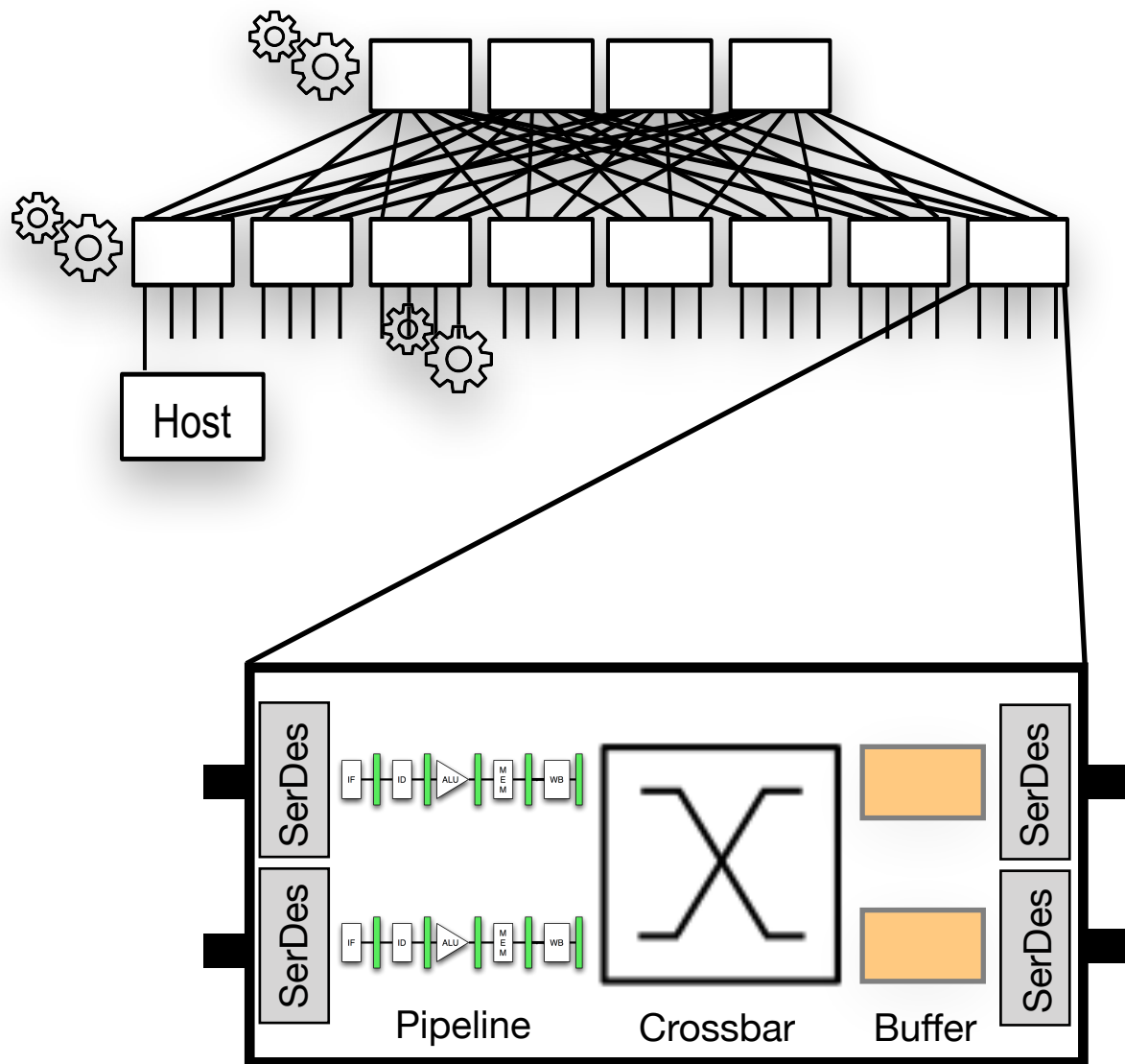


Circuit Switch

Promise of Circuit Switching

Packet Switching

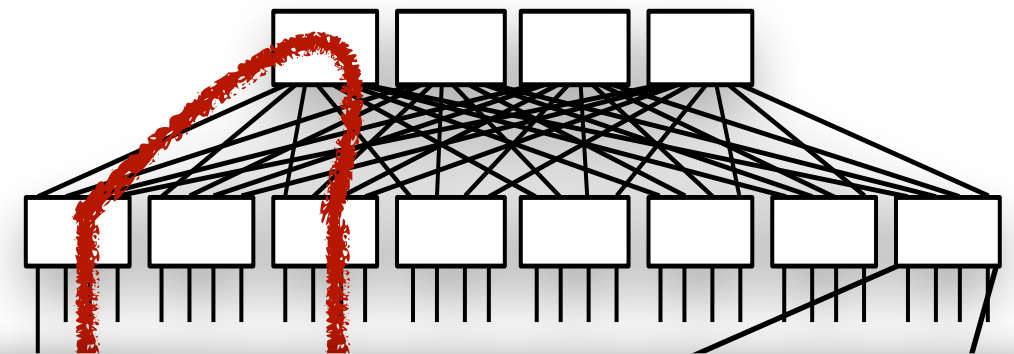
(Per packet, Per hop Decisions)



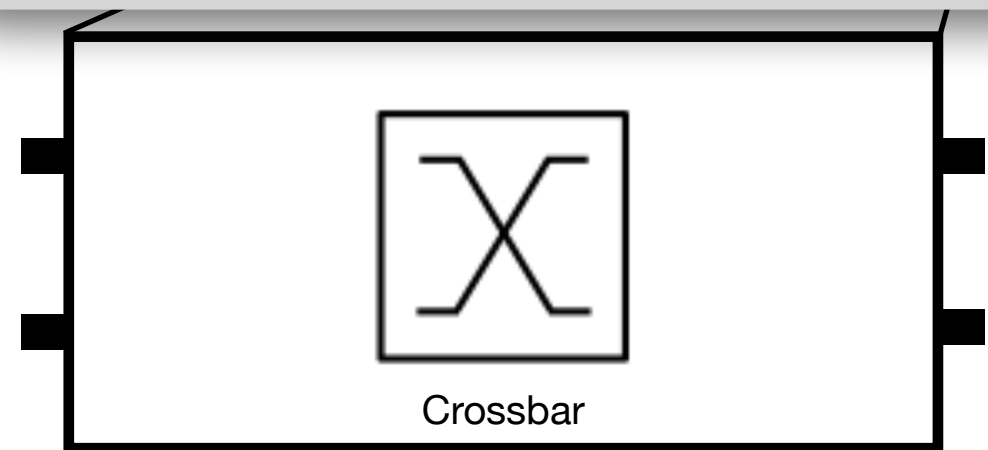
Packet Switch

VS.

Circuit Switching

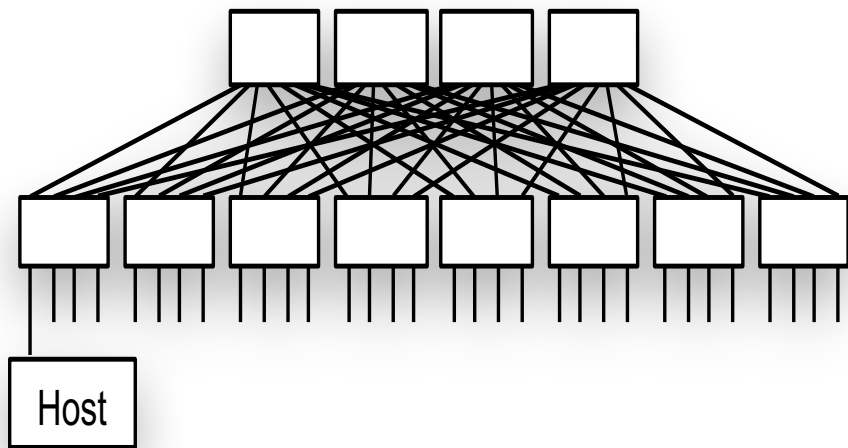


Lower Power
Lower \$\$
Unlimited Bandwidth Scaling (w/ optics)

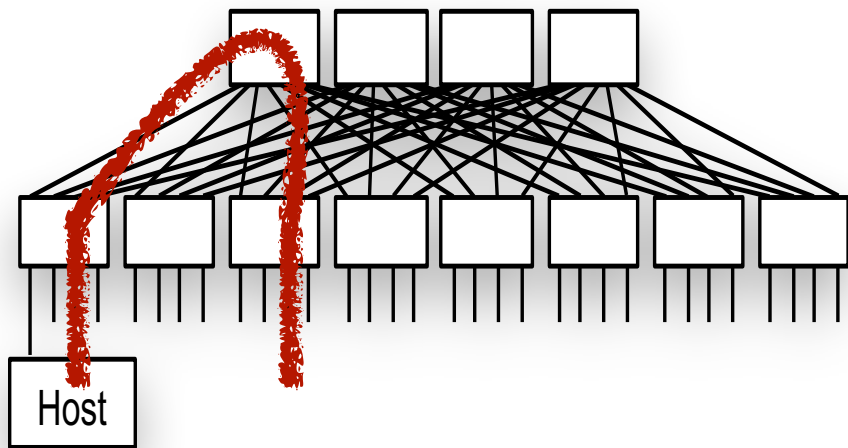


Circuit Switch

Challenge of Circuit Switching

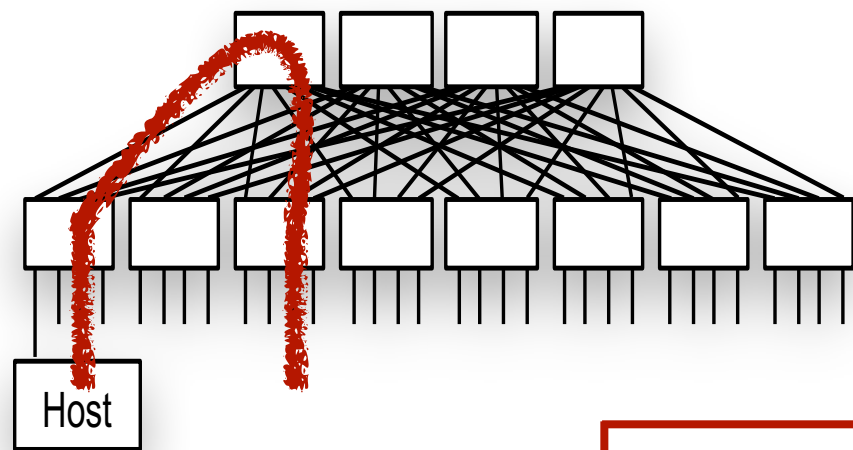


Challenge of Circuit Switching



👹 Circuit
Configuration

Challenge of Circuit Switching



 Circuit Configuration

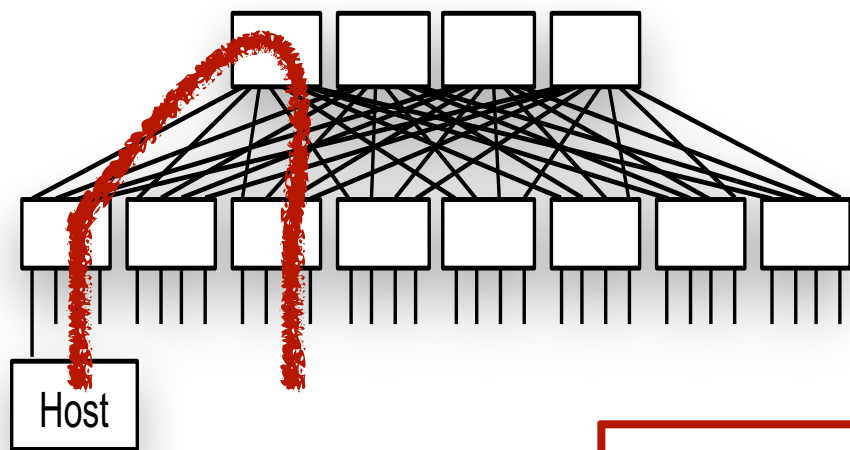
Data Plane

creating physical circuits

Control Plane

circuit scheduling

Challenge of Circuit Switching



 Circuit Configuration

Data Plane

creating physical circuits

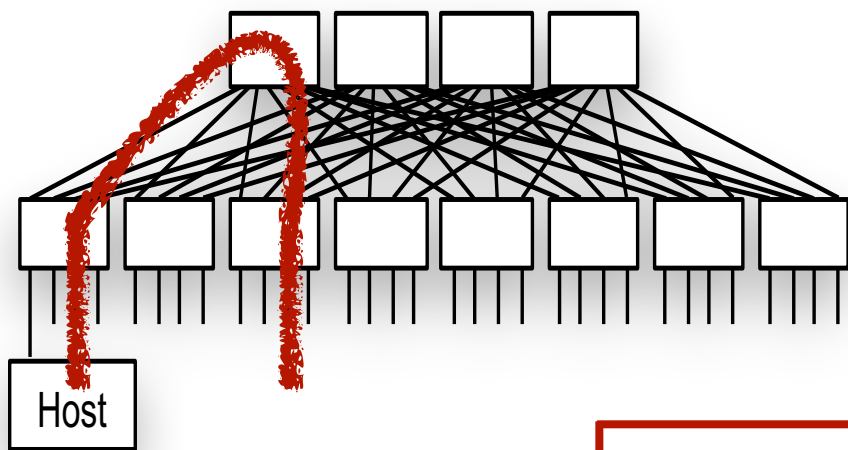
Control Plane

circuit scheduling

~2010

milliseconds

Challenge of Circuit Switching



 Circuit Configuration

Data Plane

creating physical circuits

Control Plane

circuit scheduling

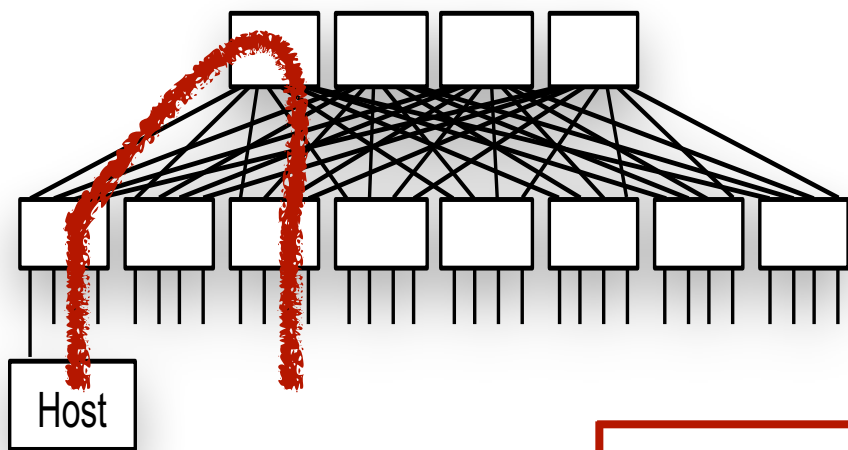
~2010

milliseconds

Today

nanoseconds

Challenge of Circuit Switching



 Circuit Configuration

Data Plane

creating physical circuits

Control Plane

circuit scheduling



Centralized

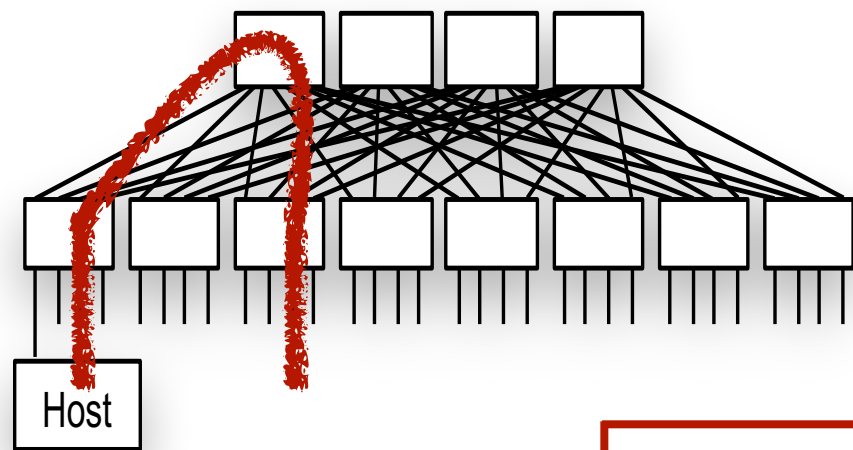
~2010

milliseconds

Today

nanoseconds

Challenge of Circuit Switching



 **Circuit Configuration**

Data Plane

creating physical circuits

Control Plane

circuit scheduling

~2010

milliseconds

Traffic demand



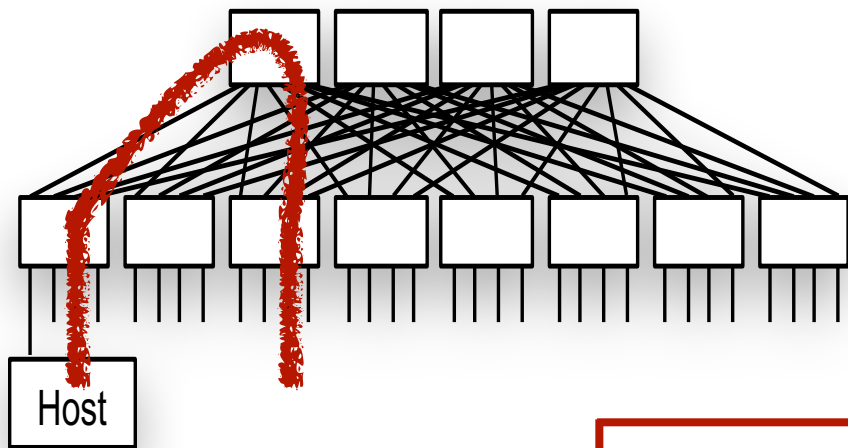
Circuits

Centralized

Today

nanoseconds

Challenge of Circuit Switching



 **Circuit Configuration**

Data Plane

creating physical circuits

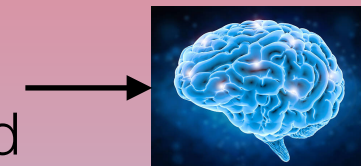
Control Plane

circuit scheduling

~2010

milliseconds

Traffic demand



Circuits

Centralized

Today

nanoseconds

???

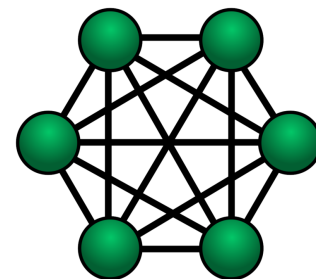
How to build a fast control plane for circuit switching?

nanosecond-scale circuit scheduling with high performance

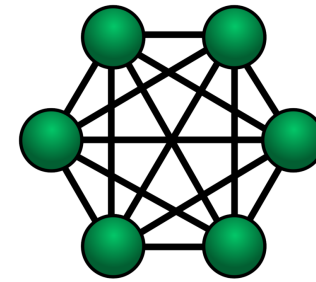
How to build a fast control plane for circuit switching?

nanosecond-scale circuit scheduling with high performance

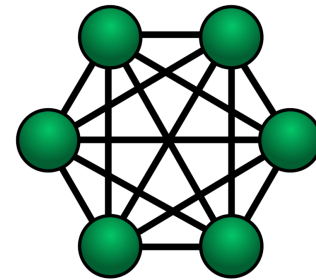
Shoal



Shoal

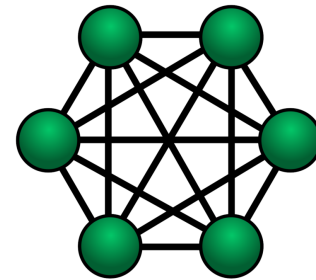


Shoal



1. **Physical Layer** : Fast circuit scheduling mechanism
2. **Routing** : Bounded worst-case throughput
3. **Congestion Control** : Bounded worst-case queuing

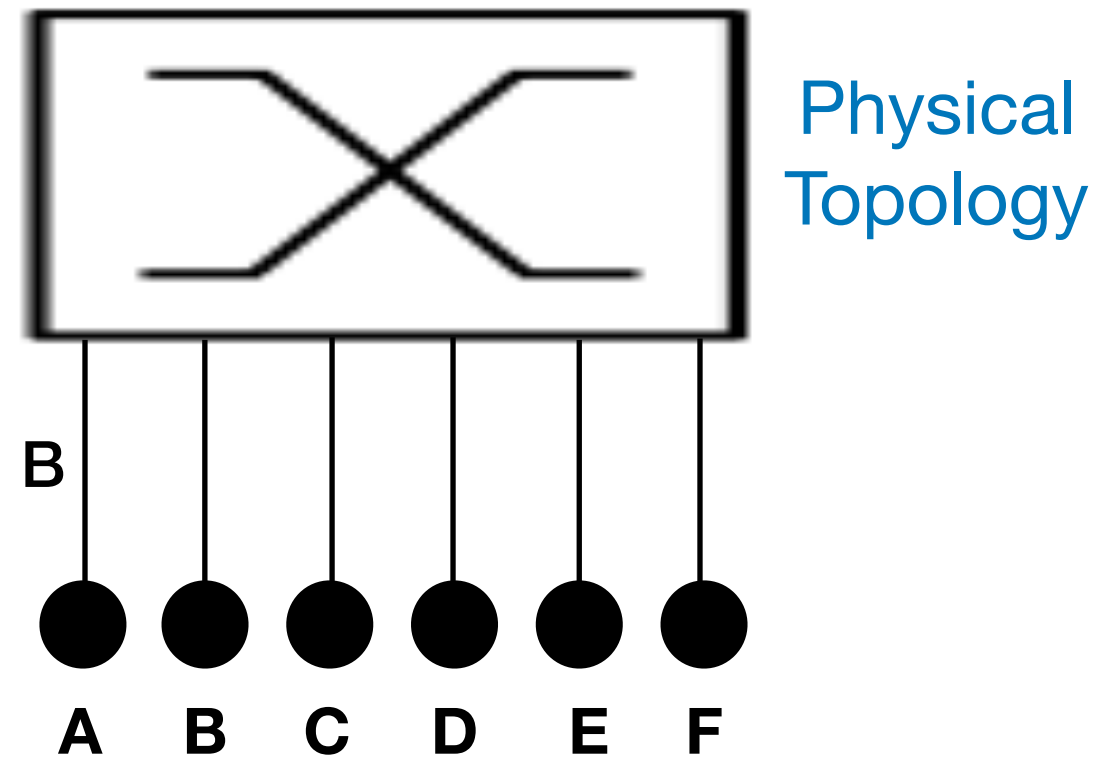
Shoal



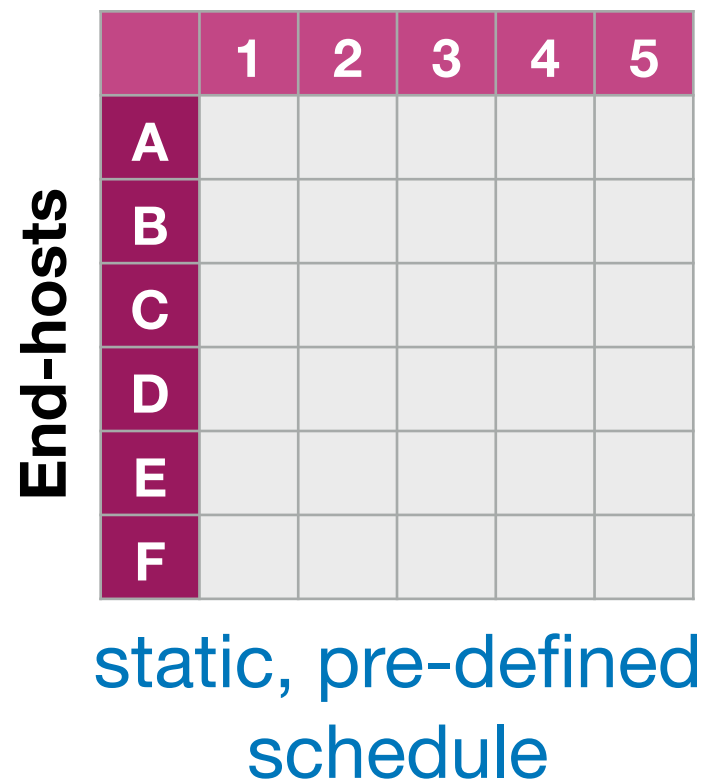
1. **Physical Layer** : Fast circuit scheduling mechanism
2. **Routing** : Bounded worst-case throughput
3. **Congestion Control** : Bounded worst-case queuing

Achieves comparable or better performance than several recent packet-switched network designs

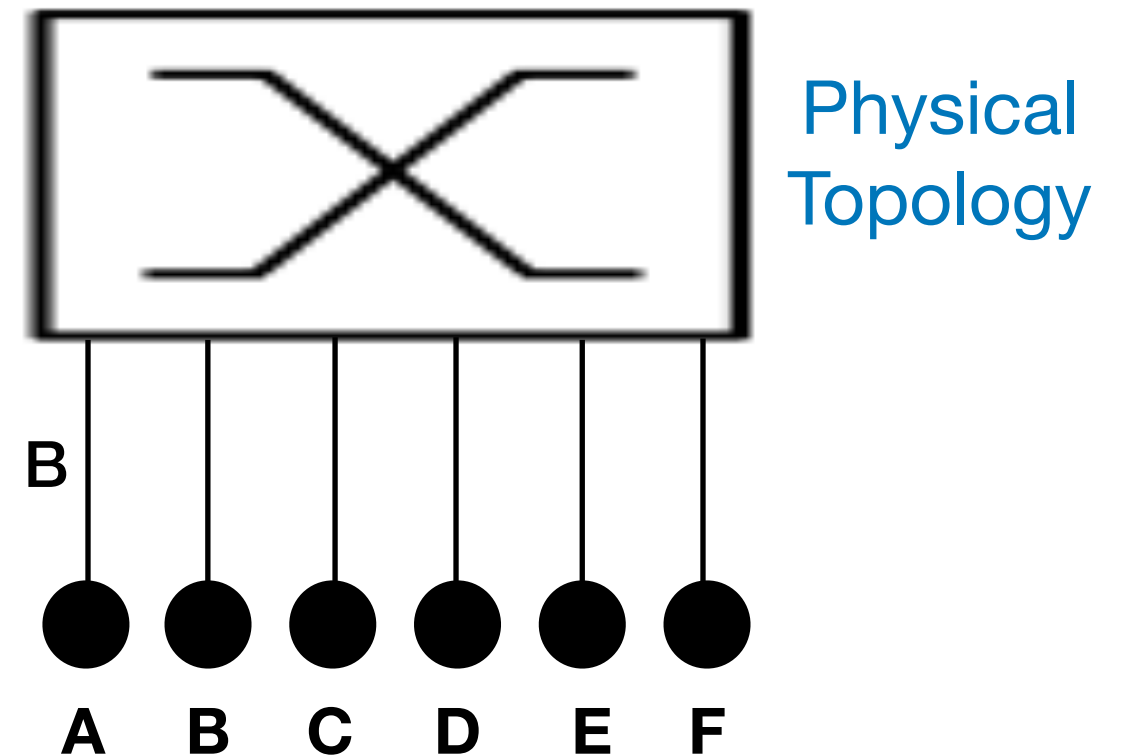
Physical Layer



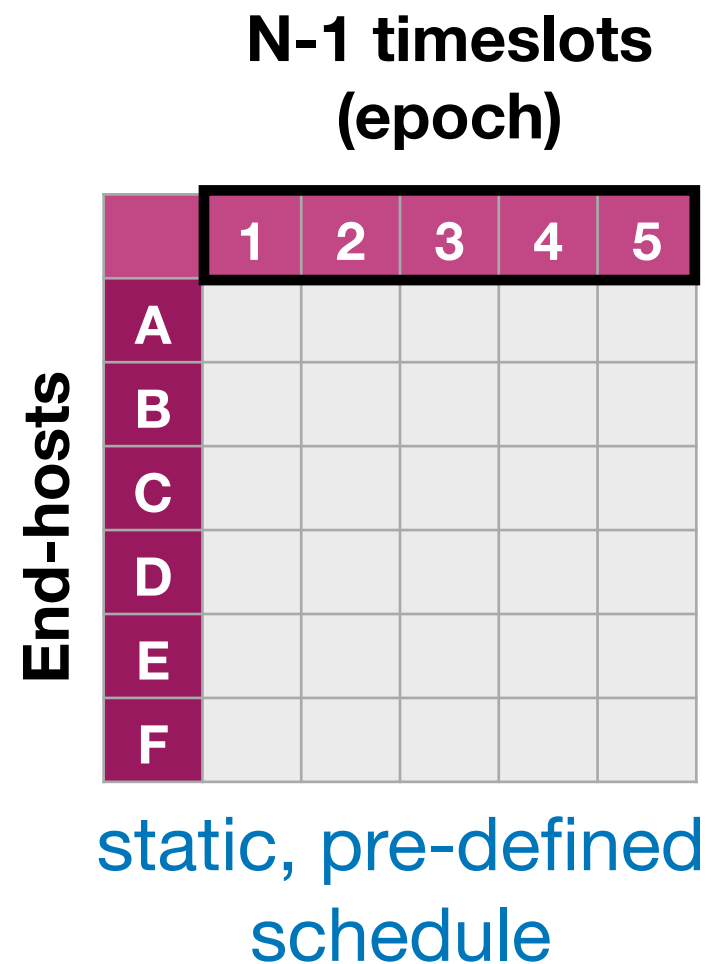
Physical Layer



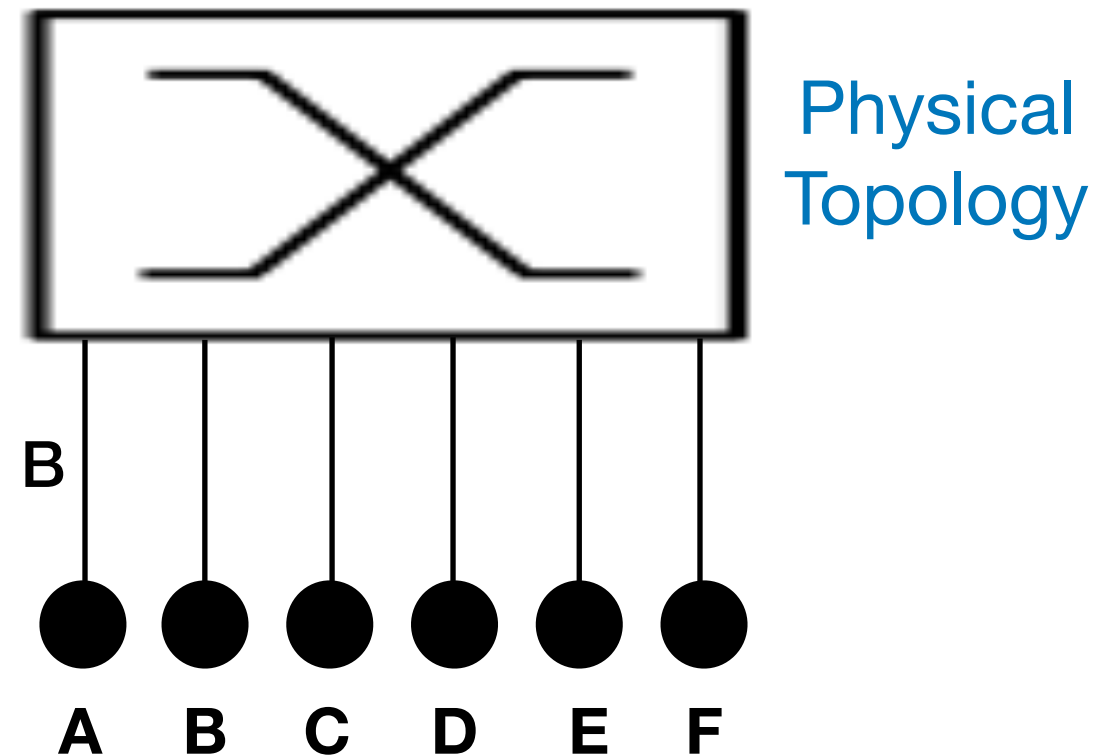
+



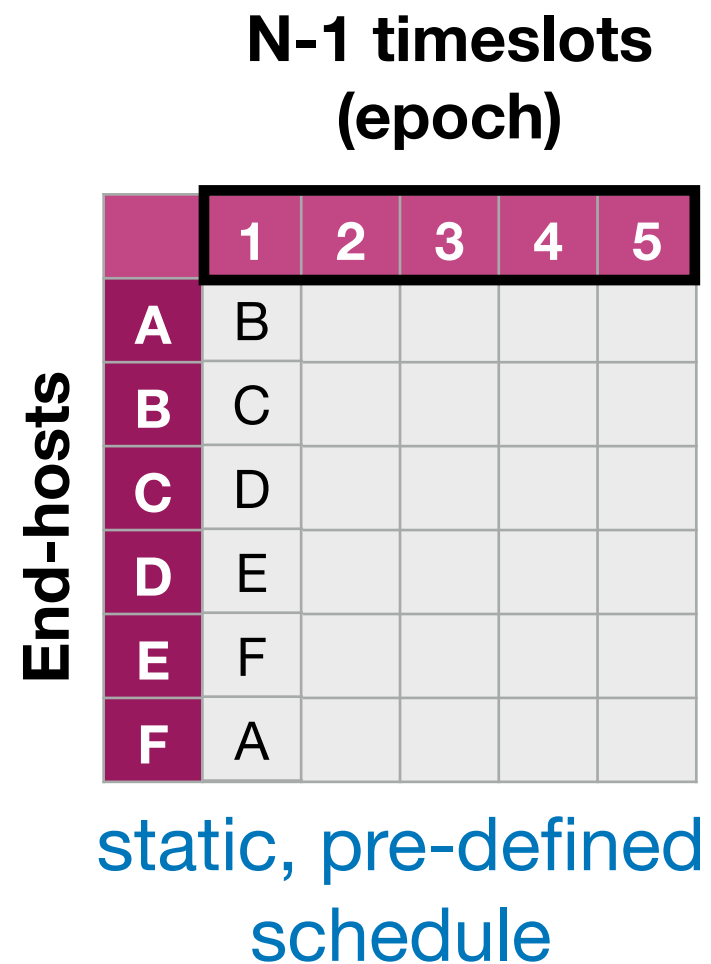
Physical Layer



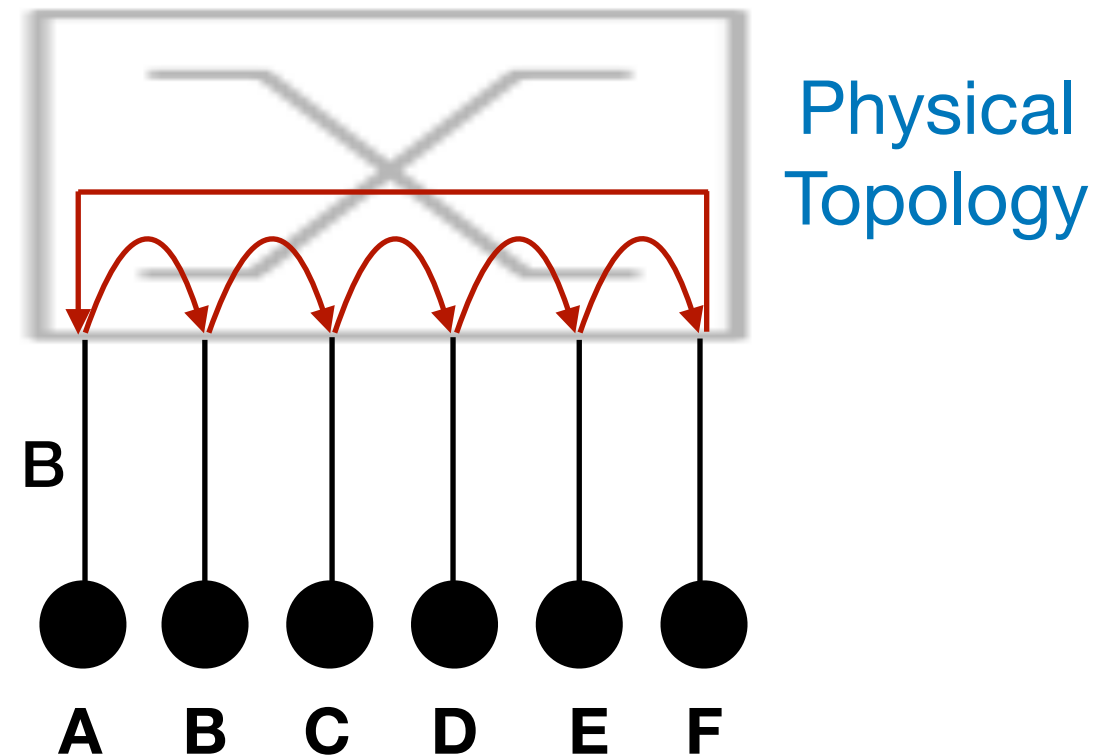
+



Physical Layer



+

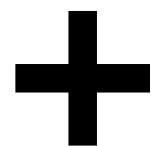


Physical Layer

N-1 timeslots
(epoch)

End-hosts		1	2	3	4	5
	A	B				
	B	C				
	C	D				
	D	E				
	E	F				
	F	A				

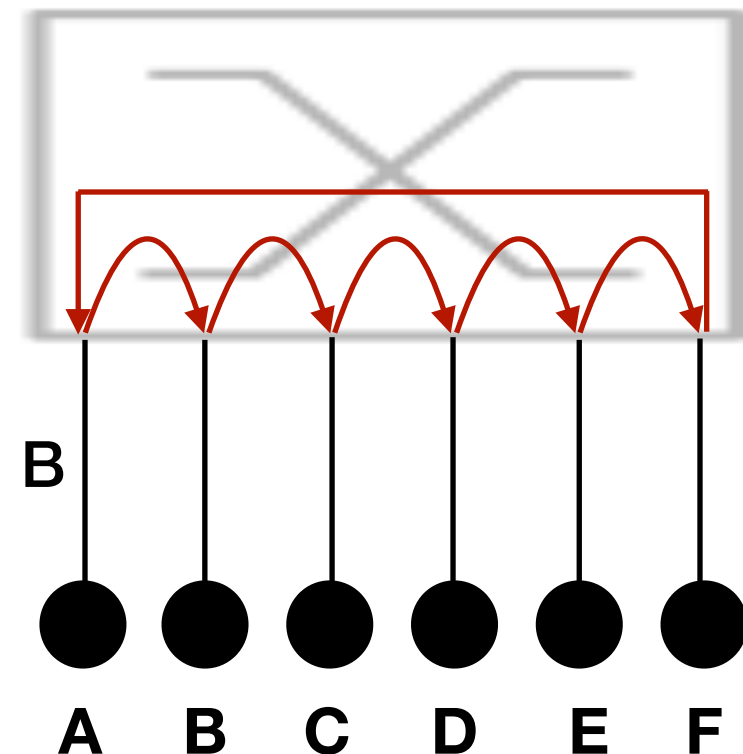
static, pre-defined
schedule



Synchronous
System

ns-precision time sync

DTP [SIGCOMM'16]



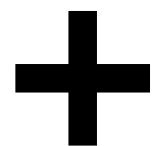
Physical
Topology

Physical Layer

**N-1 timeslots
(epoch)**

End-hosts		1	2	3	4	5
	A	B	C	D	E	F
	B	C	D	E	F	A
	C	D	E	F	A	B
	D	E	F	A	B	C
	E	F	A	B	C	D
	F	A	B	C	D	E

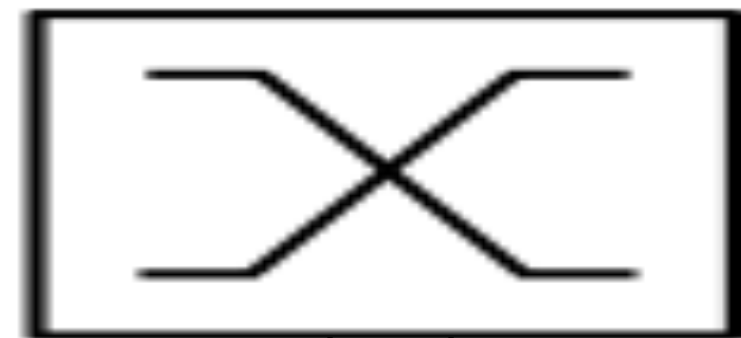
static, pre-defined
schedule



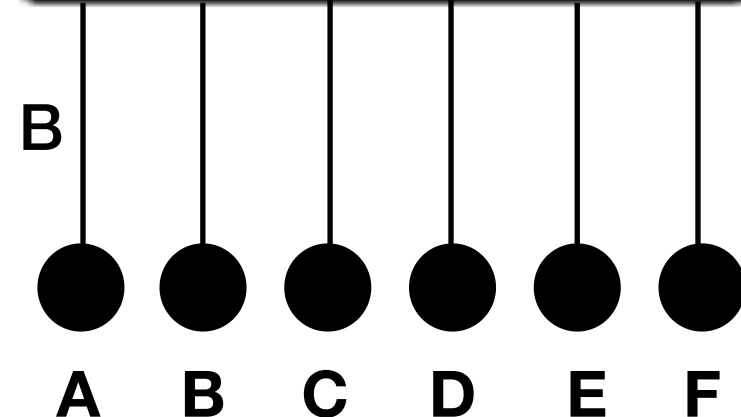
**Synchronous
System**

ns-precision time sync

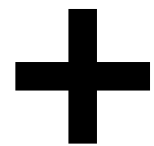
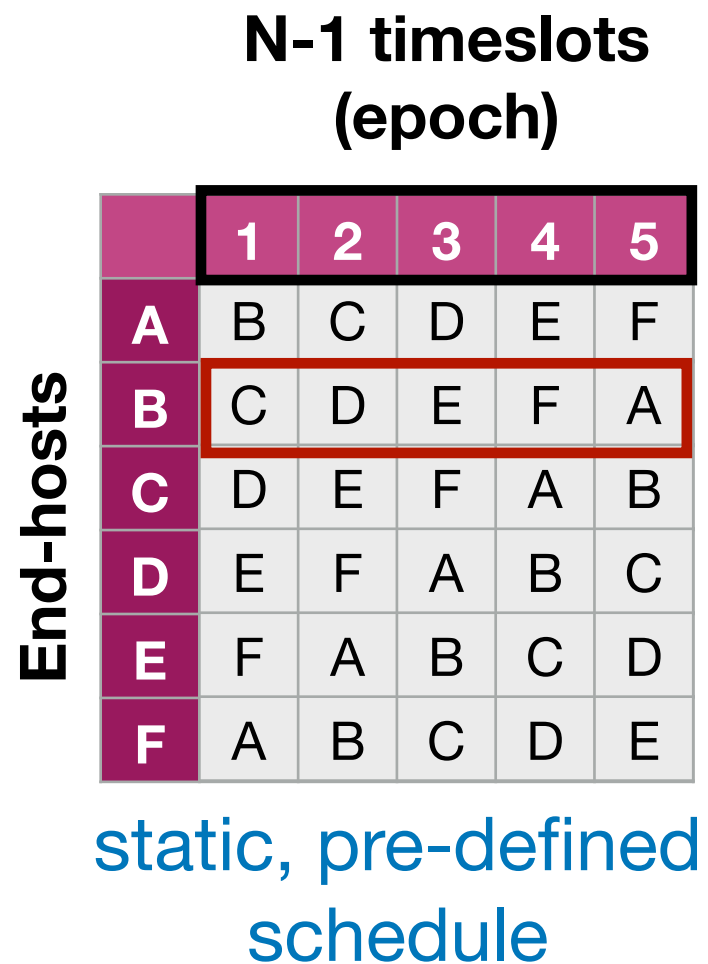
DTP [SIGCOMM'16]



Physical
Topology



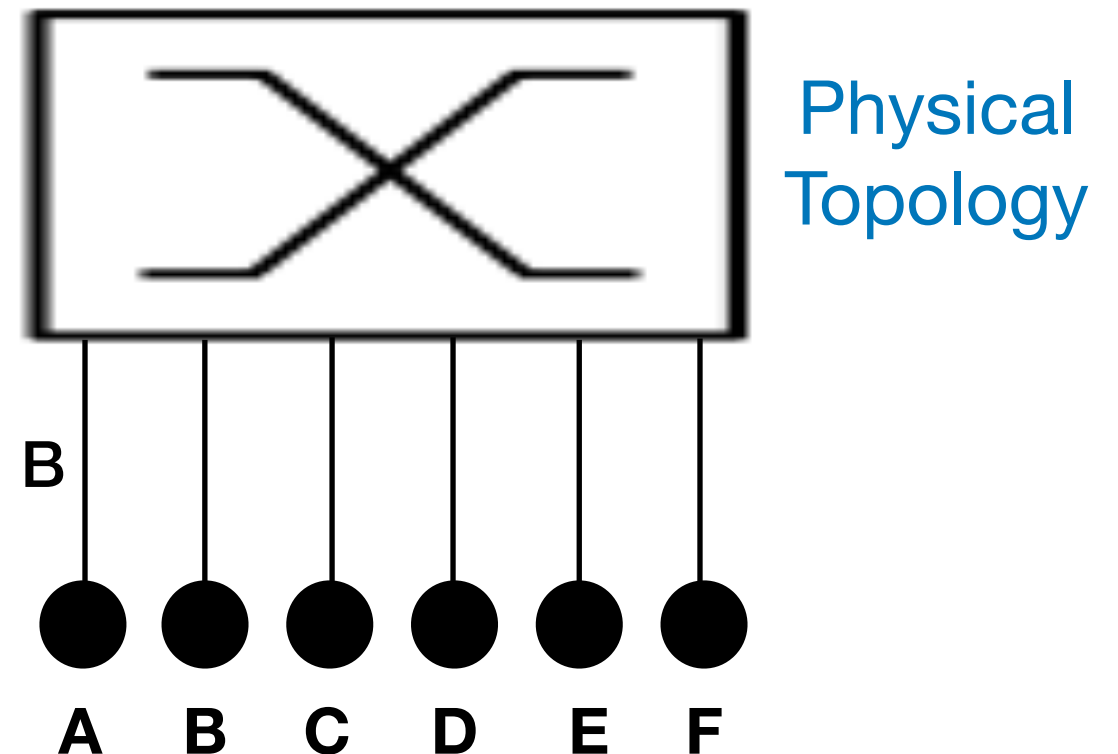
Physical Layer



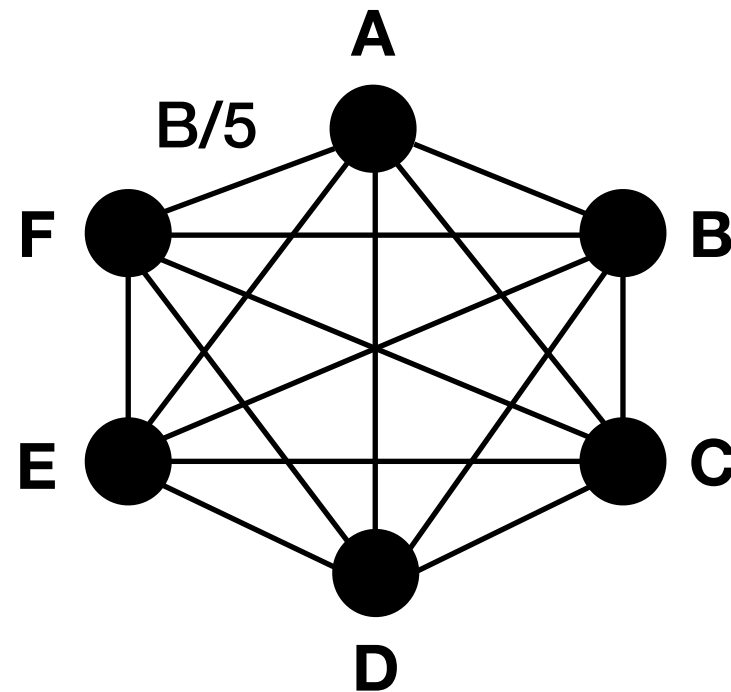
**Synchronous
System**

ns-precision time sync

DTP [SIGCOMM'16]



Physical Layer

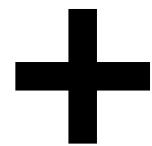


Full Mesh
Virtual Topology

N-1 timeslots
(epoch)

End-hosts		1	2	3	4	5
	A	B	C	D	E	F
	B	C	D	E	F	A
	C	D	E	F	A	B
	D	E	F	A	B	C
	E	F	A	B	C	D
	F	A	B	C	D	E

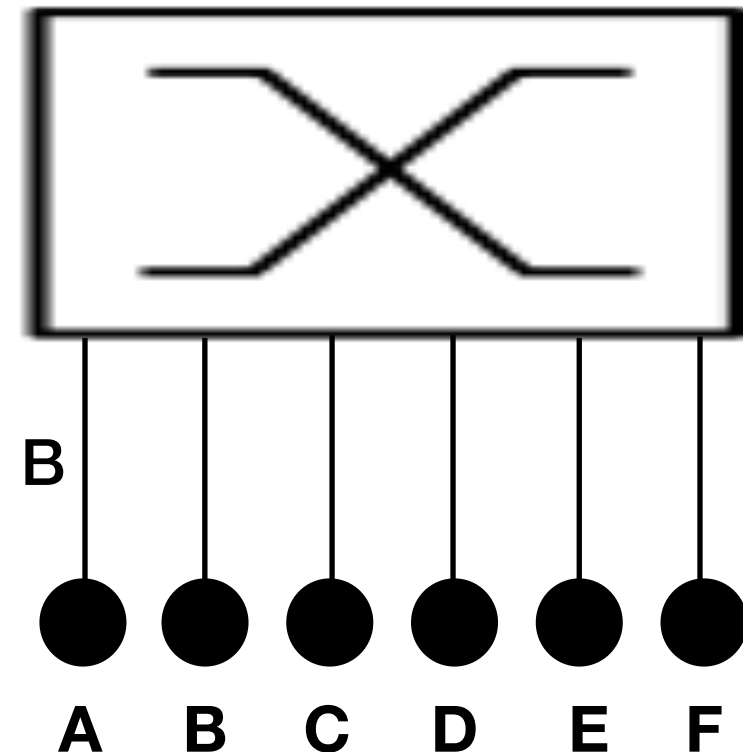
static, pre-defined
schedule



Synchronous
System

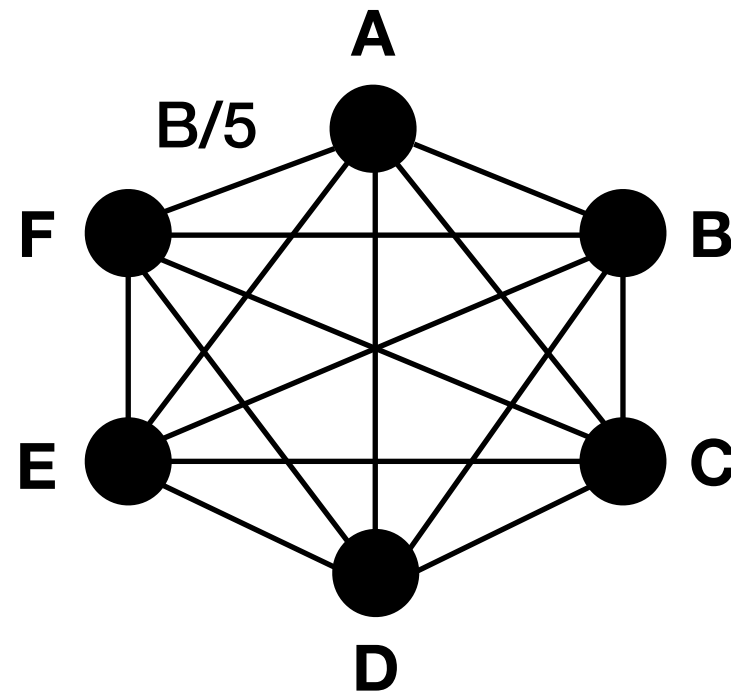
ns-precision time sync

DTP [SIGCOMM'16]



Physical
Topology

Physical Layer

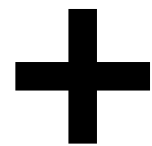


Full Mesh
Virtual Topology

N-1 timeslots
(epoch)

End-hosts		1	2	3	4	5
	A	B	C	D	E	F
	B	C	D	E	F	A
	C	D	E	F	A	B
	D	E	F	A	B	C
	E	F	A	B	C	D
	F	A	B	C	D	E

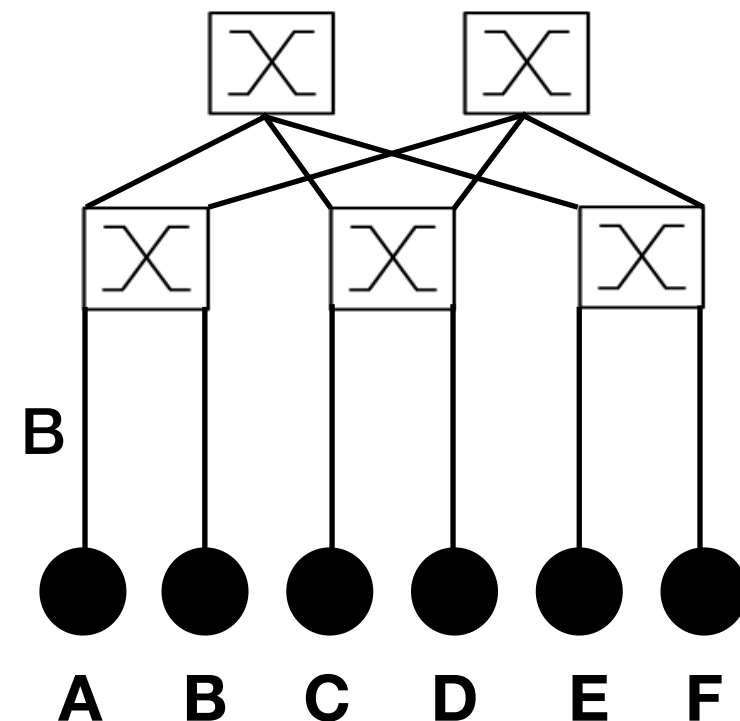
static, pre-defined
schedule



Synchronous
System

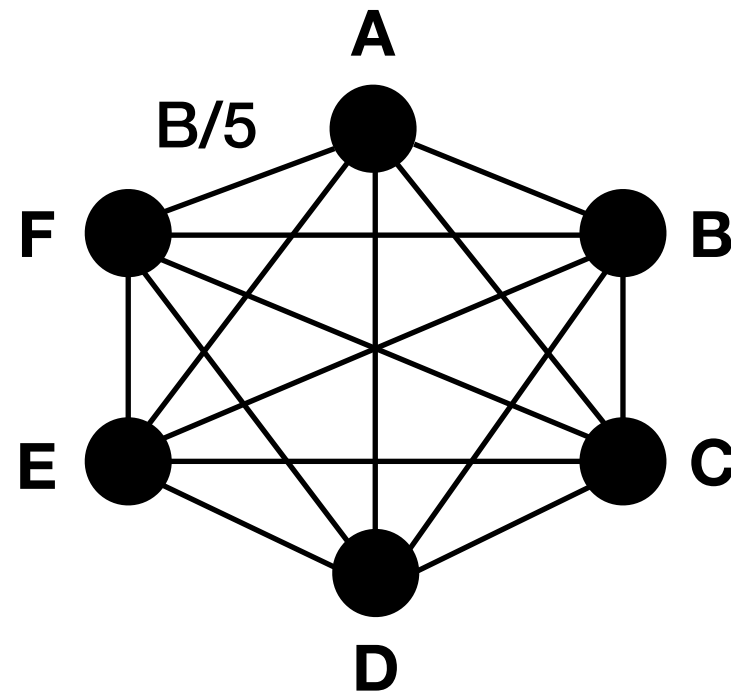
ns-precision time sync

DTP [SIGCOMM'16]



Physical
Topology

Physical Layer

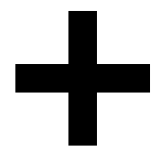


Full Mesh
Virtual Topology

N-1 timeslots
(epoch)

End-hosts		1	2	3	4	5
	A	B	C	D	E	F
	B	C	D	E	F	A
	C	D	E	F	A	B
	D	E	F	A	B	C
	E	F	A	B	C	D
	F	A	B	C	D	E

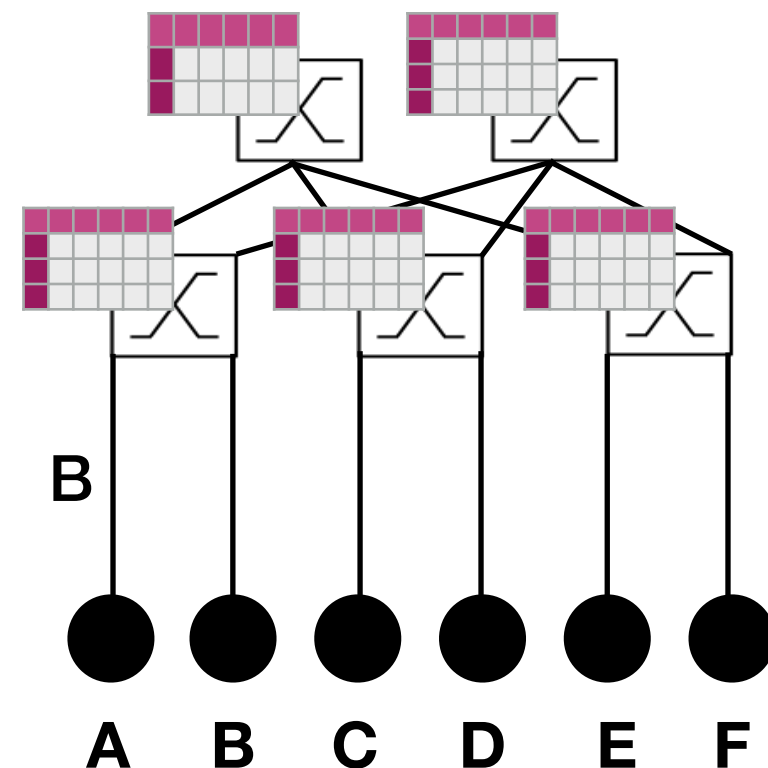
static, pre-defined
schedule



Synchronous
System

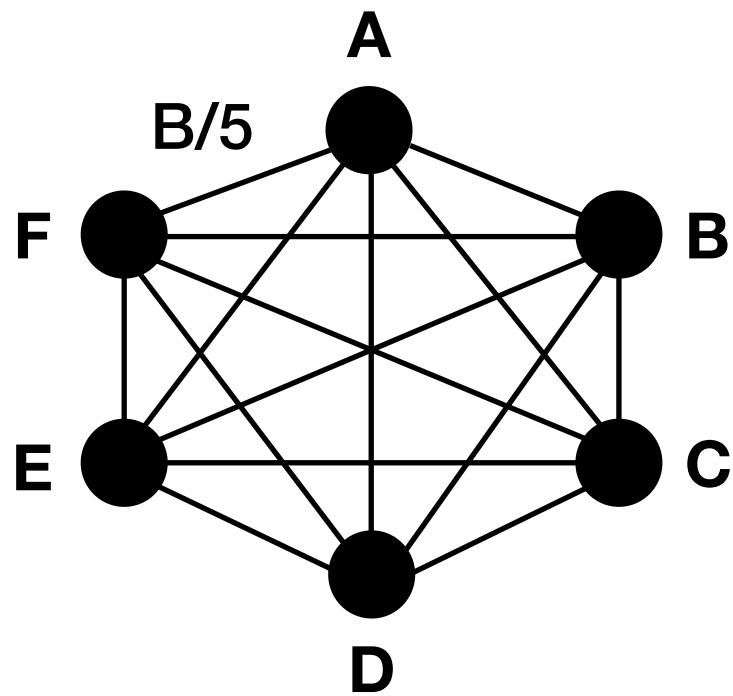
ns-precision time sync

DTP [SIGCOMM'16]

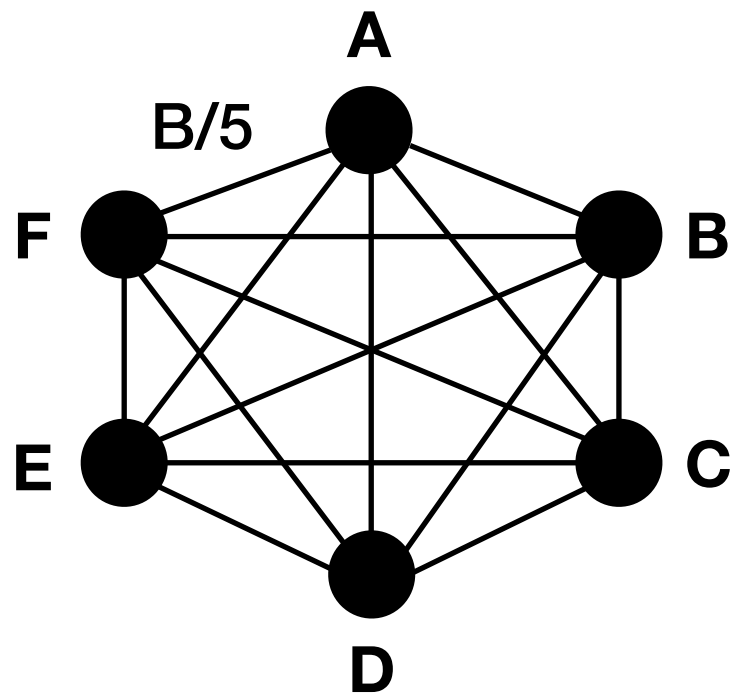


Physical
Topology

Routing & Congestion Control

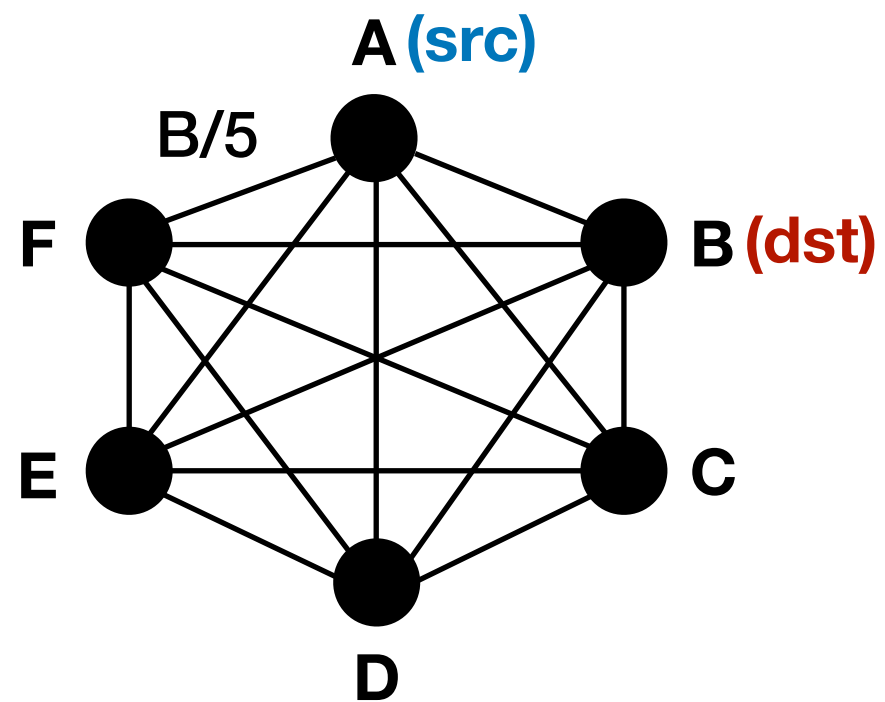


Routing & Congestion Control



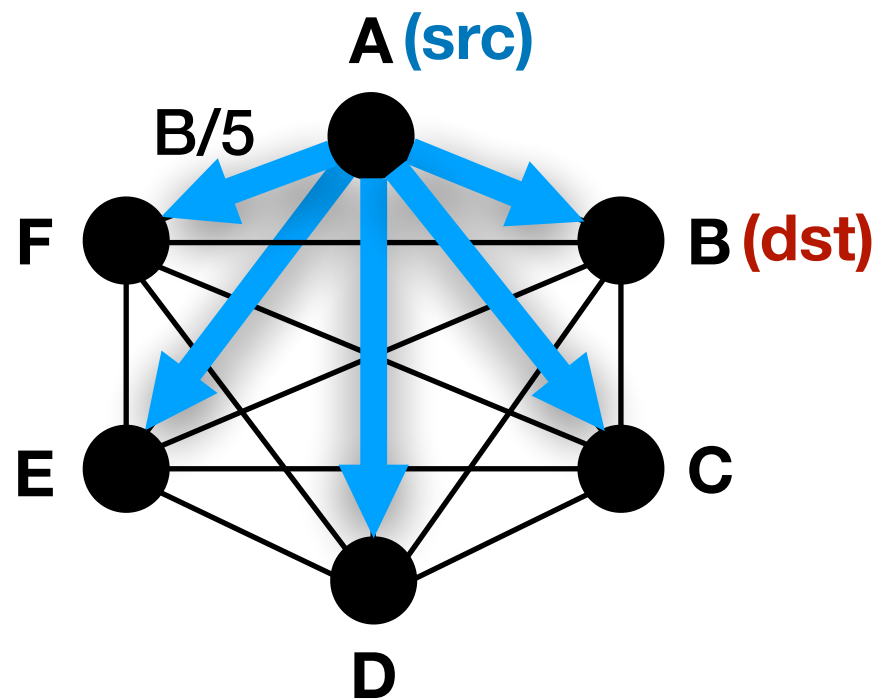
Valiant
Load
Balancing

Routing & Congestion Control



Valiant
Load
Balancing

Routing & Congestion Control

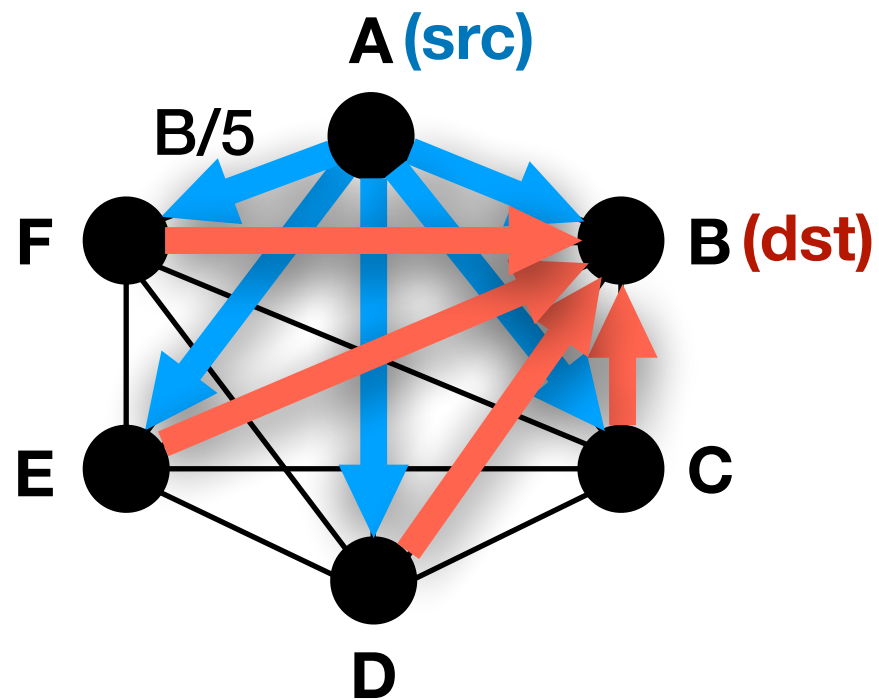


Valiant
Load
Balancing

Uniform load
balancing

Routing in Shoal

Routing & Congestion Control



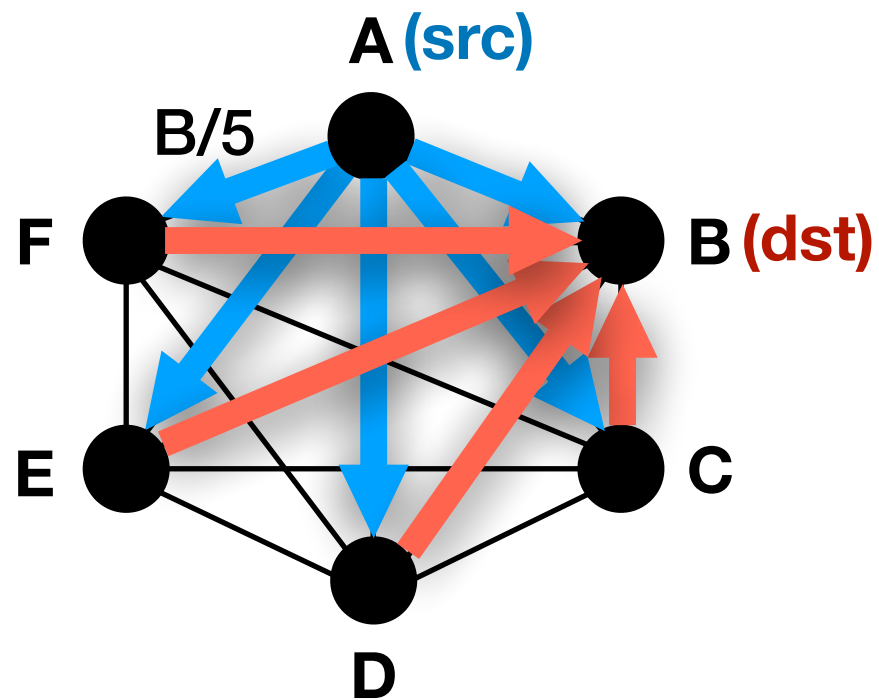
Valiant
Load
Balancing

Uniform load
balancing

Forward to
destination

Routing in Shoal

Routing & Congestion Control



Valiant
Load
Balancing

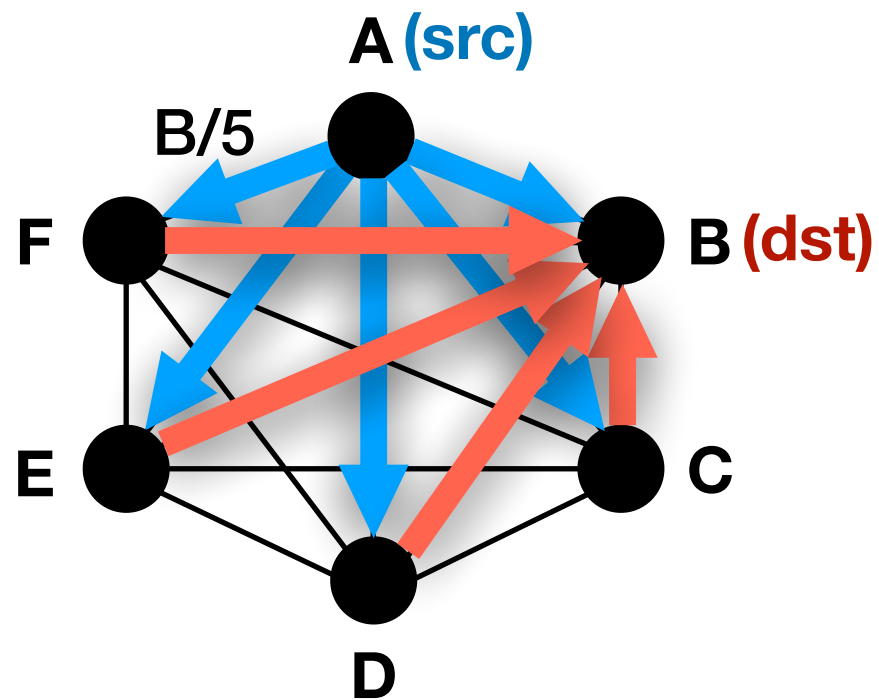
Uniform load
balancing

Forward to
destination

Routing in Shoal

Bounded worst-case throughput of 50%
compared to an ideal packet-switched network

Routing & Congestion Control



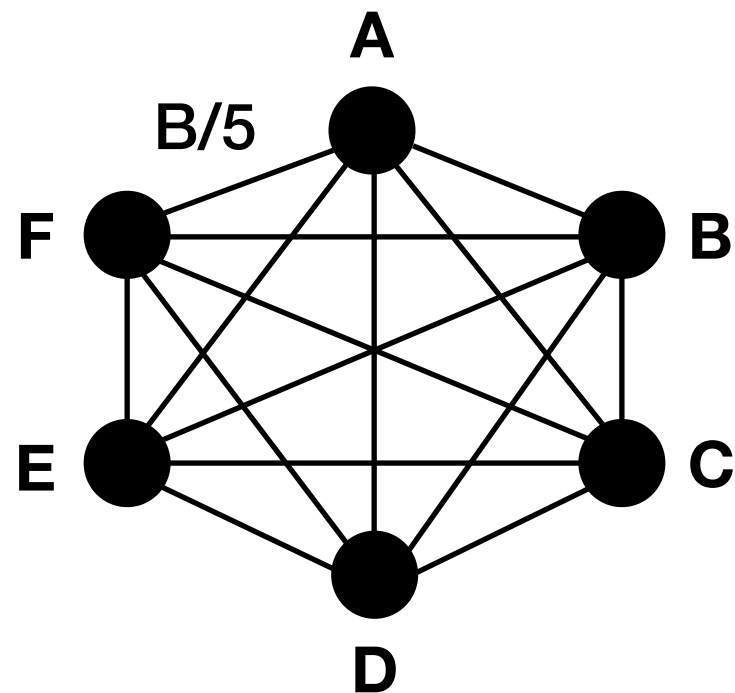
Valiant
Load
Balancing

Uniform load
balancing

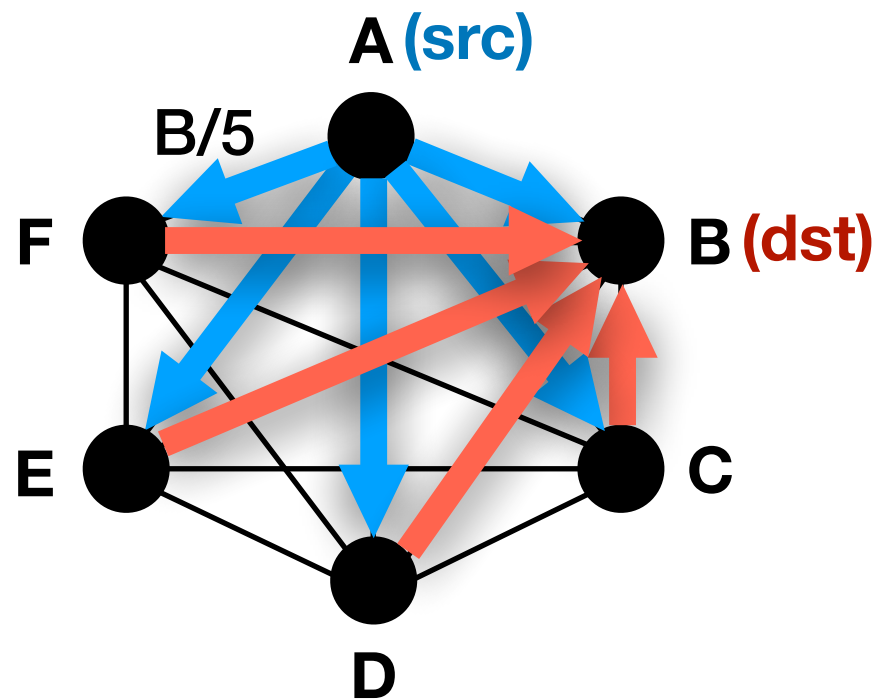
Forward to
destination

Routing in Shoal

Bounded worst-case throughput of 50%
compared to an ideal packet-switched network



Routing & Congestion Control



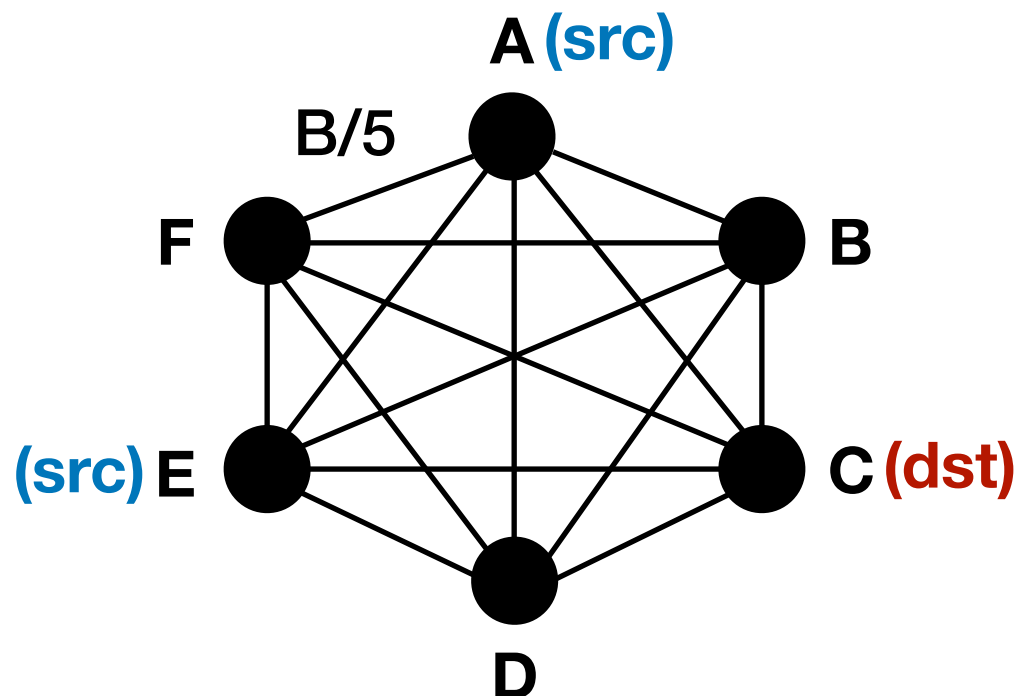
Valiant
Load
Balancing

Uniform load
balancing

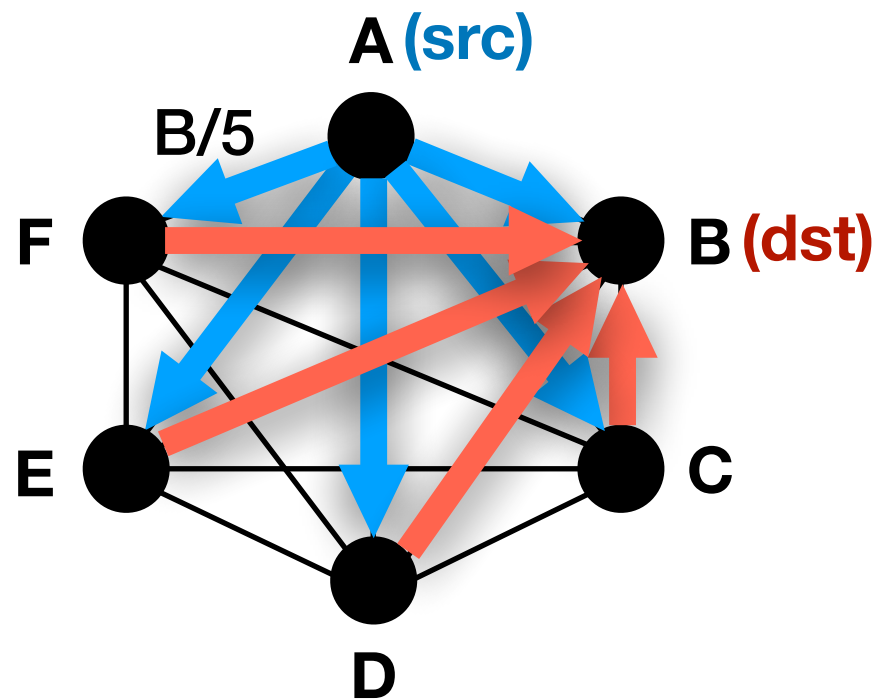
Forward to
destination

Routing in Shoal

Bounded worst-case throughput of 50%
compared to an ideal packet-switched network



Routing & Congestion Control



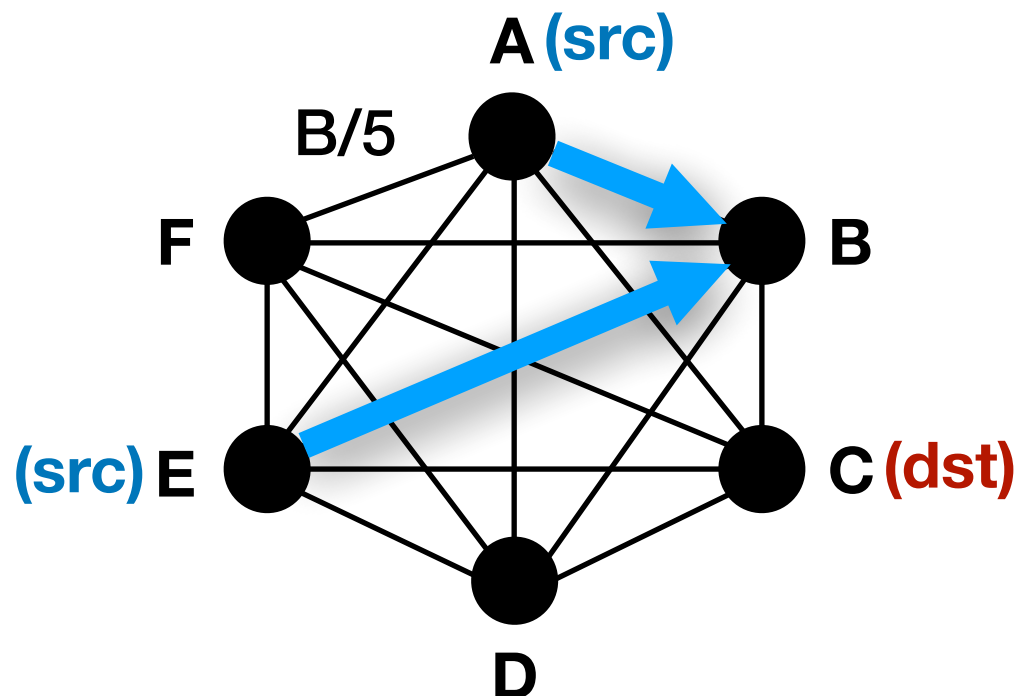
Valiant
Load
Balancing

Uniform load
balancing

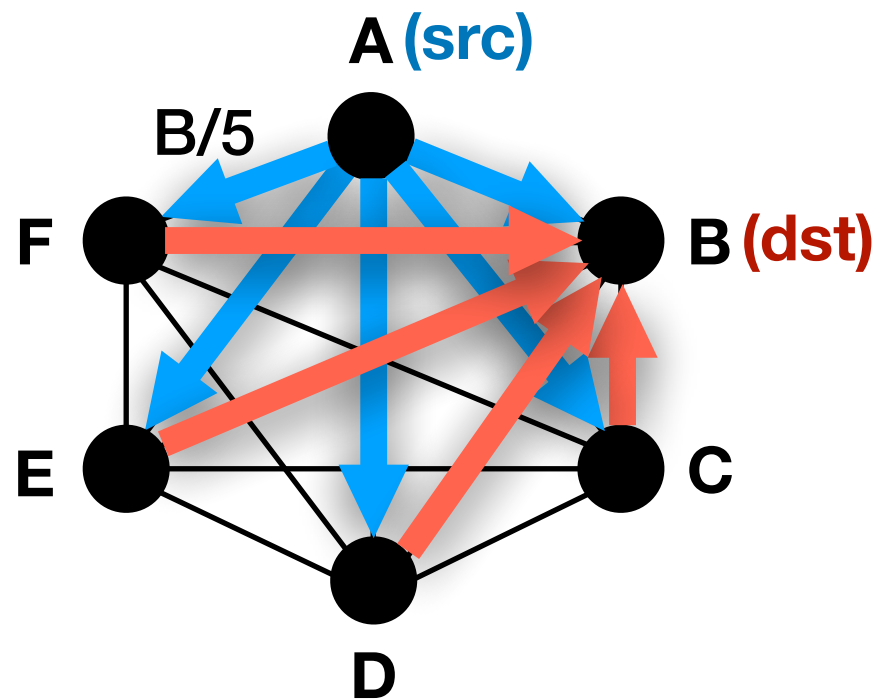
Forward to
destination

Routing in Shoal

Bounded worst-case throughput of 50%
compared to an ideal packet-switched network



Routing & Congestion Control



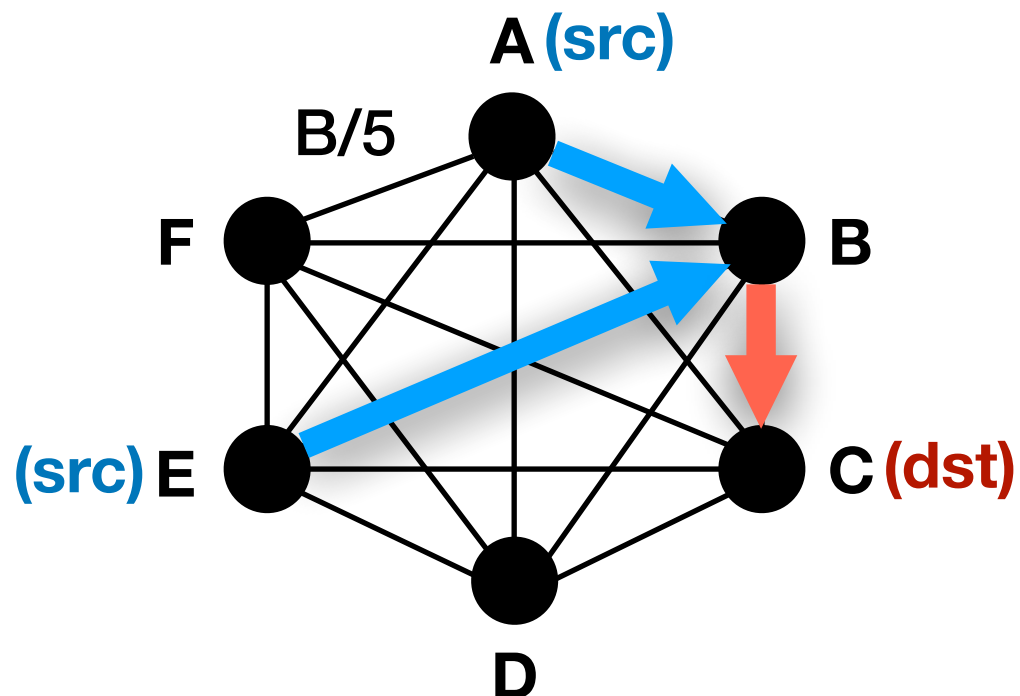
Valiant
Load
Balancing

Uniform load
balancing

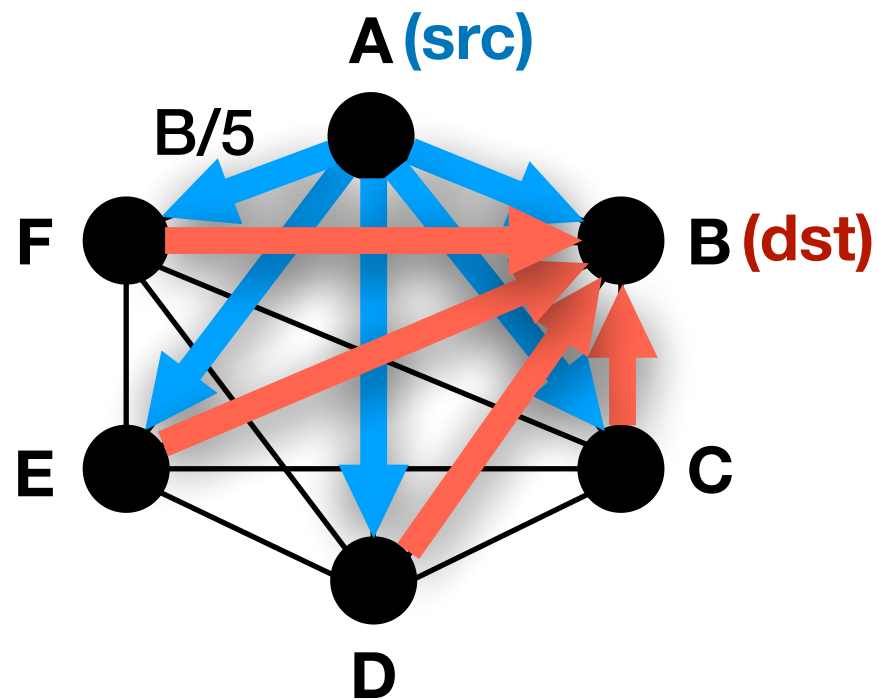
Forward to
destination

Routing in Shoal

Bounded worst-case throughput of 50%
compared to an ideal packet-switched network



Routing & Congestion Control



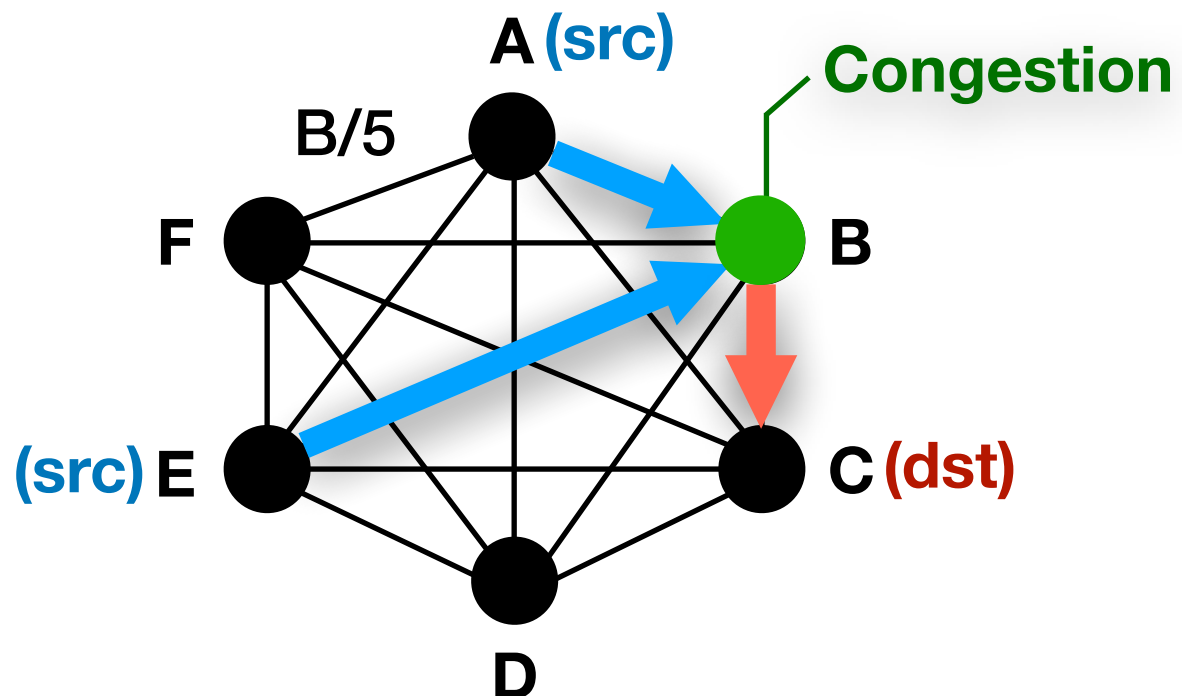
Valiant
Load
Balancing

Uniform load
balancing

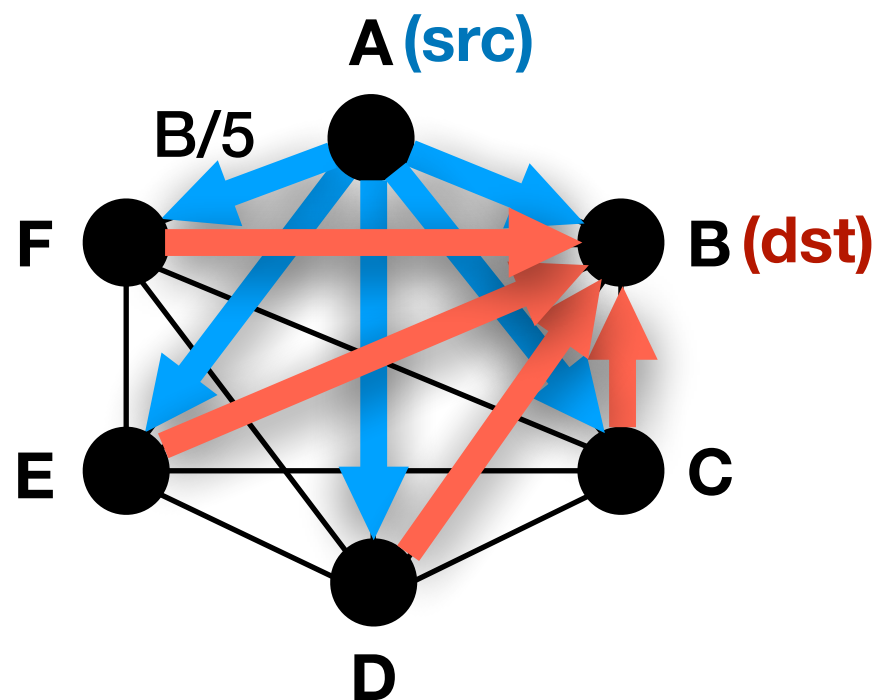
Forward to
destination

Routing in Shoal

Bounded worst-case throughput of 50%
compared to an ideal packet-switched network



Routing & Congestion Control



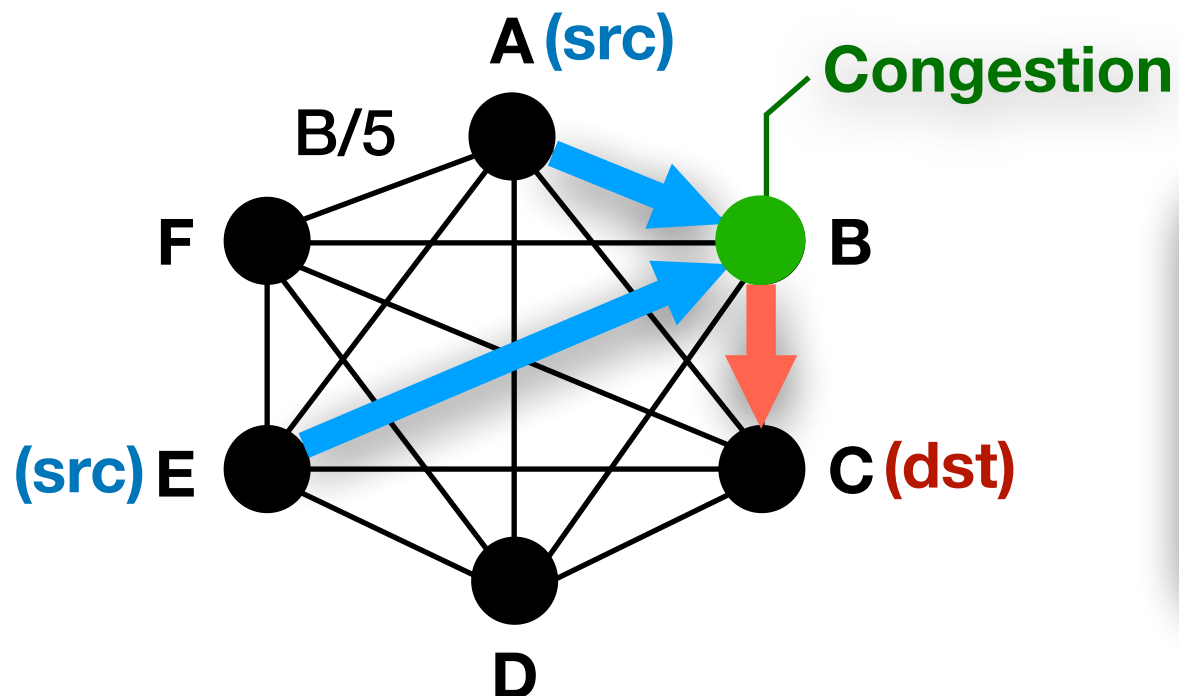
Valiant
Load
Balancing

Uniform load
balancing

Forward to
destination

Routing in Shoal

Bounded worst-case throughput of 50%
compared to an ideal packet-switched network



Congestion Control

Each per-destination queue Q_i
for each destination i is bounded!
 $len(Q_i) \leq 1 + incast_degree(i)$ packets

Key Properties of Shoal

- Fast, de-centralized, and traffic-agnostic **circuit scheduling**
- **Routing** with bounded worst-case throughput
- **Congestion Control** with bounded worst-case queuing

Key Properties of Shoal

- Fast, de-centralized, and traffic-agnostic **circuit scheduling**
- **Routing** with bounded worst-case throughput
- **Congestion Control** with bounded worst-case queuing
- **Worst-case network throughput of 50%**
(compared to an ideal packet-switched network)

Key Properties of Shoal

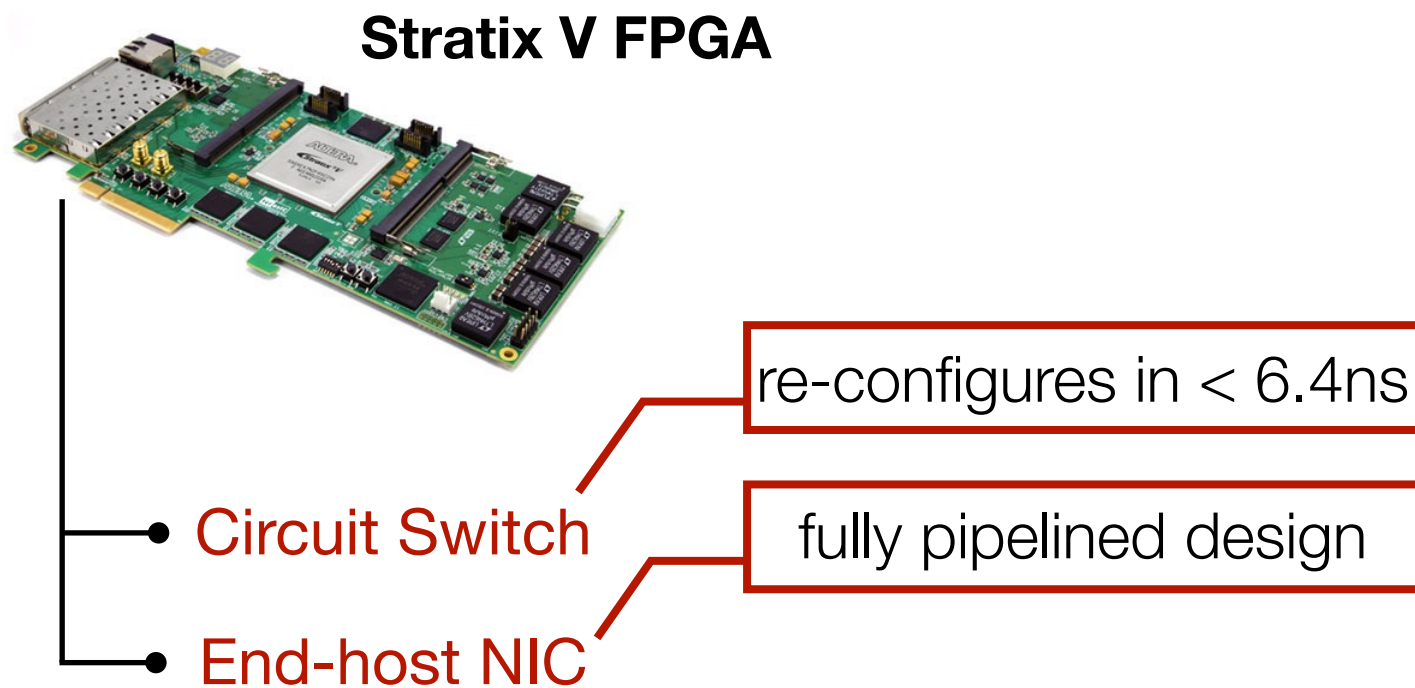
- Fast, de-centralized, and traffic-agnostic **circuit scheduling**
- **Routing** with bounded worst-case throughput
- **Congestion Control** with bounded worst-case queuing
- **Worst-case network throughput of 50%**
(compared to an ideal packet-switched network)
 - ▶ Equip each node with 2x bandwidth

Key Properties of Shoal

- Fast, de-centralized, and traffic-agnostic **circuit scheduling**
- **Routing** with bounded worst-case throughput
- **Congestion Control** with bounded worst-case queuing
- **Worst-case network throughput of 50%**
(compared to an ideal packet-switched network)
 - ▶ Equip each node with 2x bandwidth
 - ▶ $power(\text{Shoal}) < power(\text{packet-switched network with } 1/2 \text{ bandwidth of Shoal})$
 - ▶ $cost(\text{Shoal}) < cost(\text{packet-switched network with } 1/2 \text{ bandwidth of Shoal})$

Implementation

Implementation



FPGA code for the implementation of Shoal is available at:
<https://github.com/vishal1303/Shoal>

Implementation

- Routing
- Congestion Control
- Scheduling
- Synchronization



Stratix V FPGA

- Circuit Switch
- End-host NIC

re-configures in $< 6.4\text{ns}$

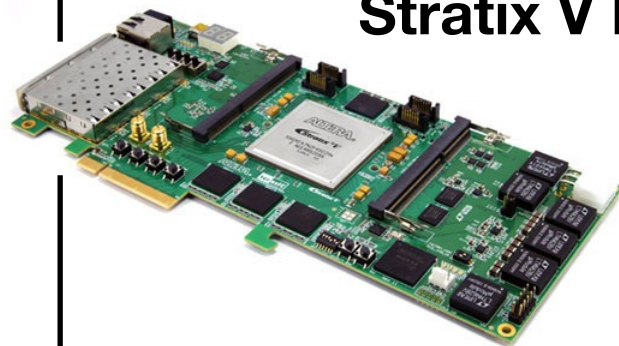
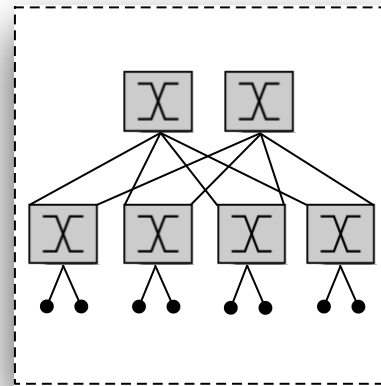
fully pipelined design

FPGA code for the implementation of Shoal is available at:

<https://github.com/vishal1303/Shoal>

Implementation

- Routing
- Congestion Control
- Scheduling
- Synchronization

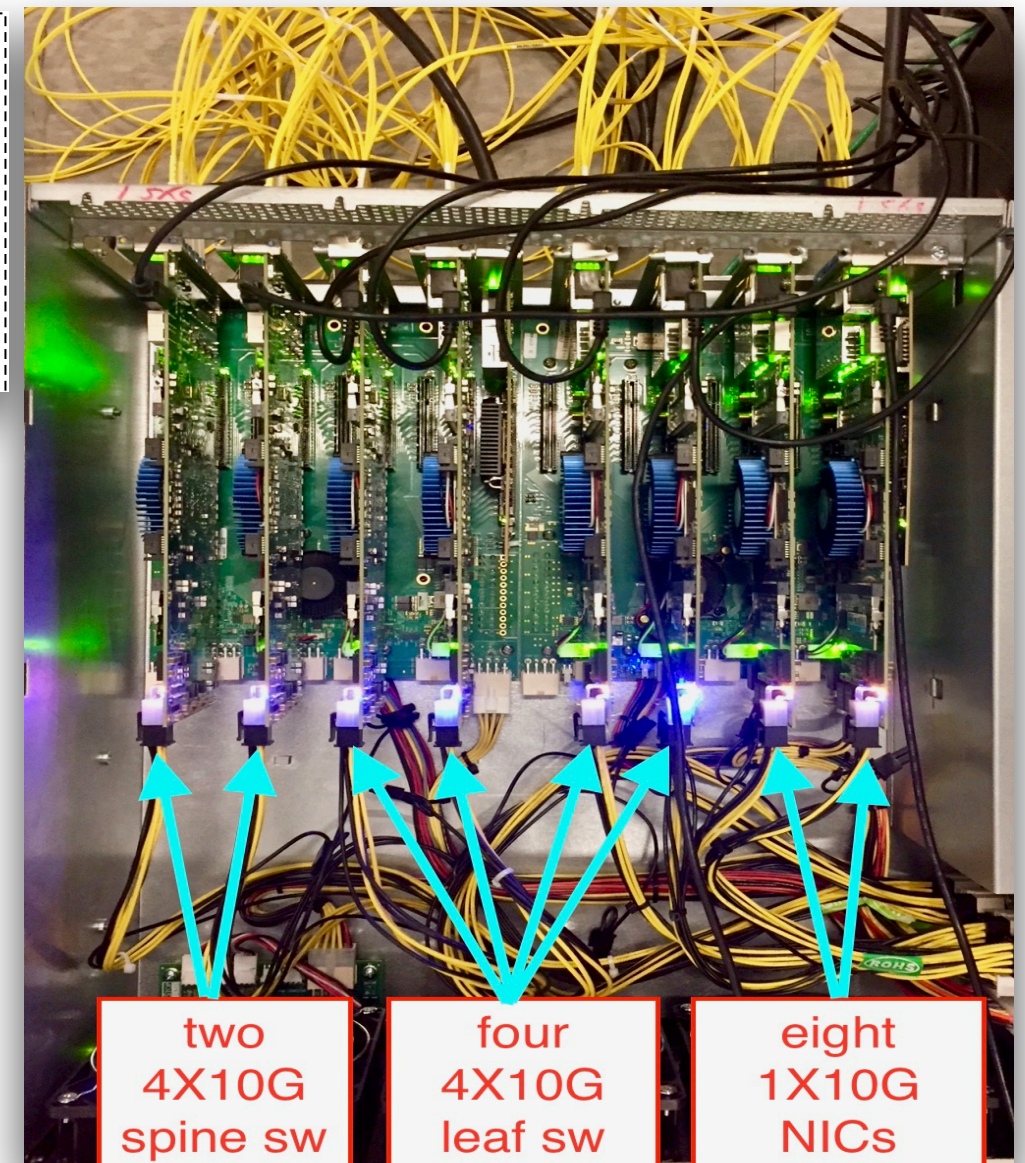


Stratix V FPGA

- **Circuit Switch**
- **End-host NIC**

re-configures in < 6.4ns

fully pipelined design



two
4X10G
spine sw

four
4X10G
leaf sw

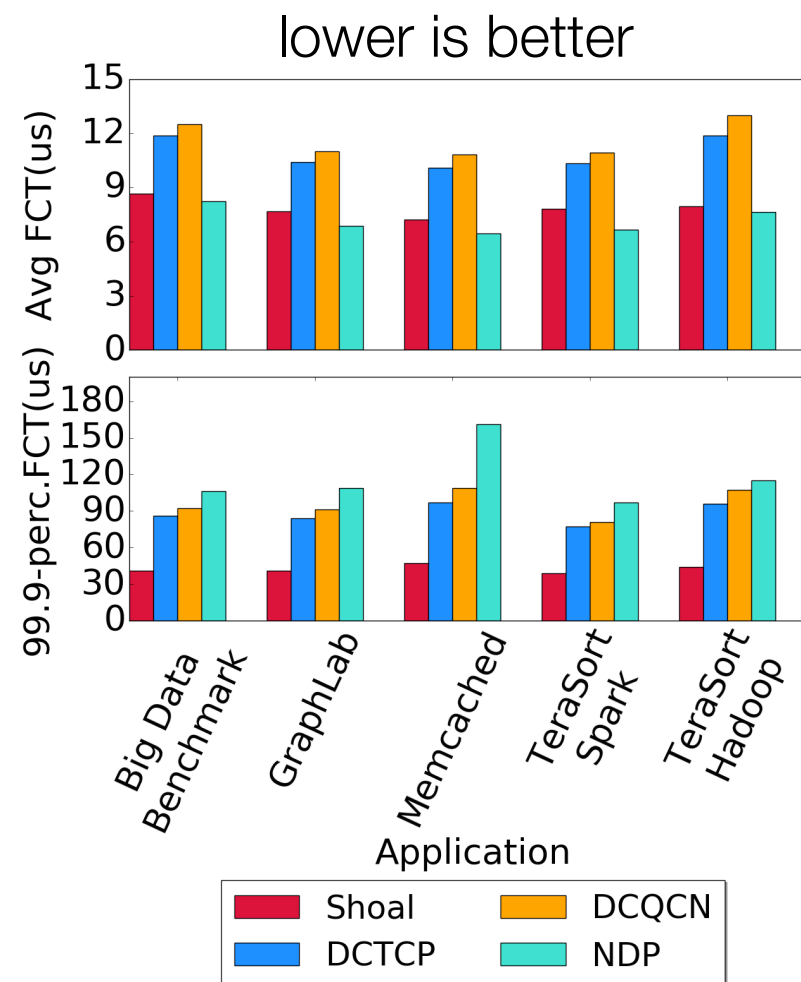
eight
1X10G
NICs

FPGA code for the implementation of Shoal is available at:

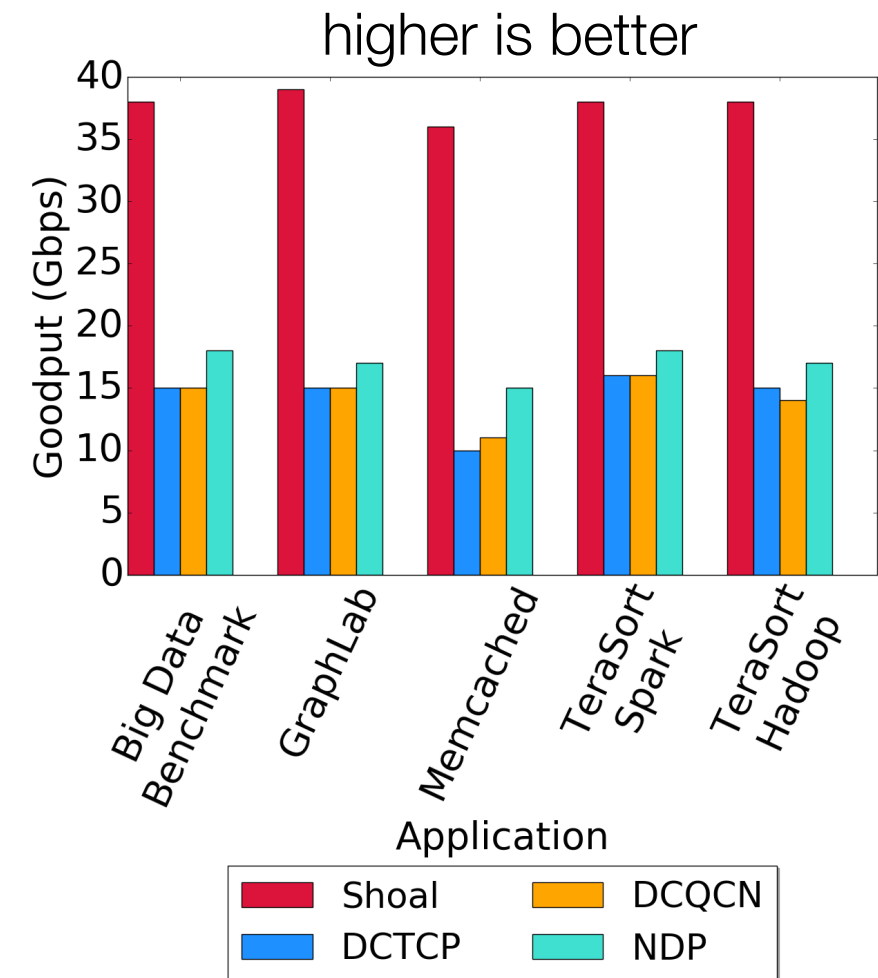
<https://github.com/vishal1303/Shoal>

Evaluation

- Packet-level simulator in C
- 512 end-hosts
- Shoal has 2x bandwidth (100Gbps vs. 50Gbps)
at lower cost and power!



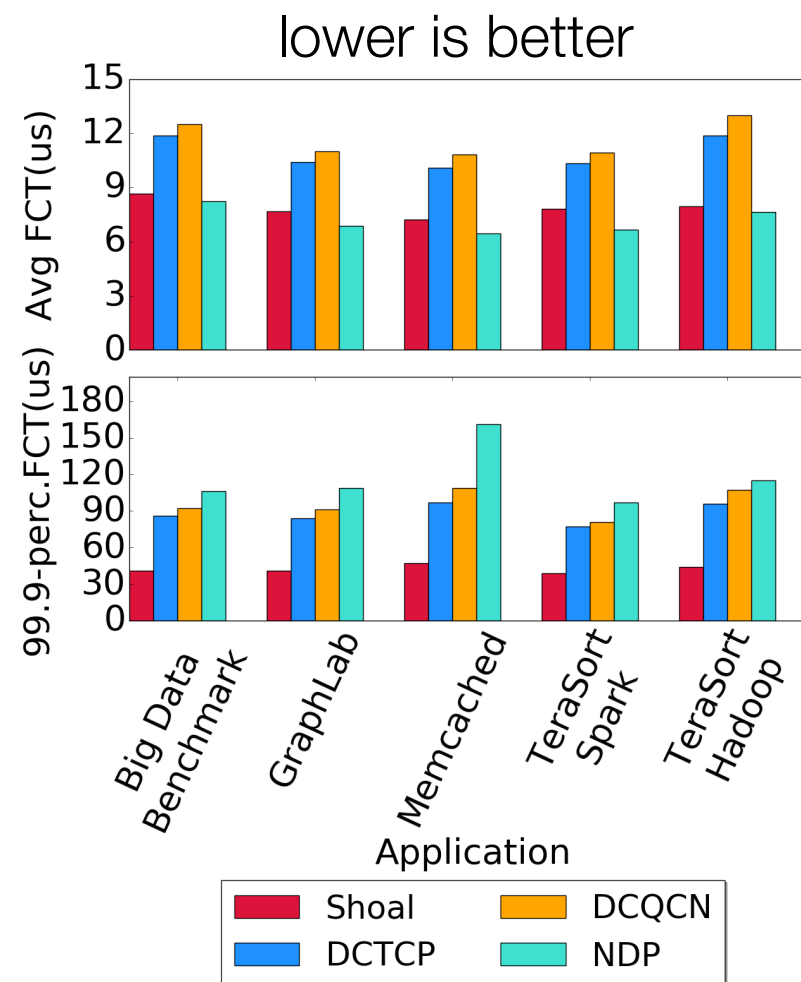
Short Flows (<100KB)



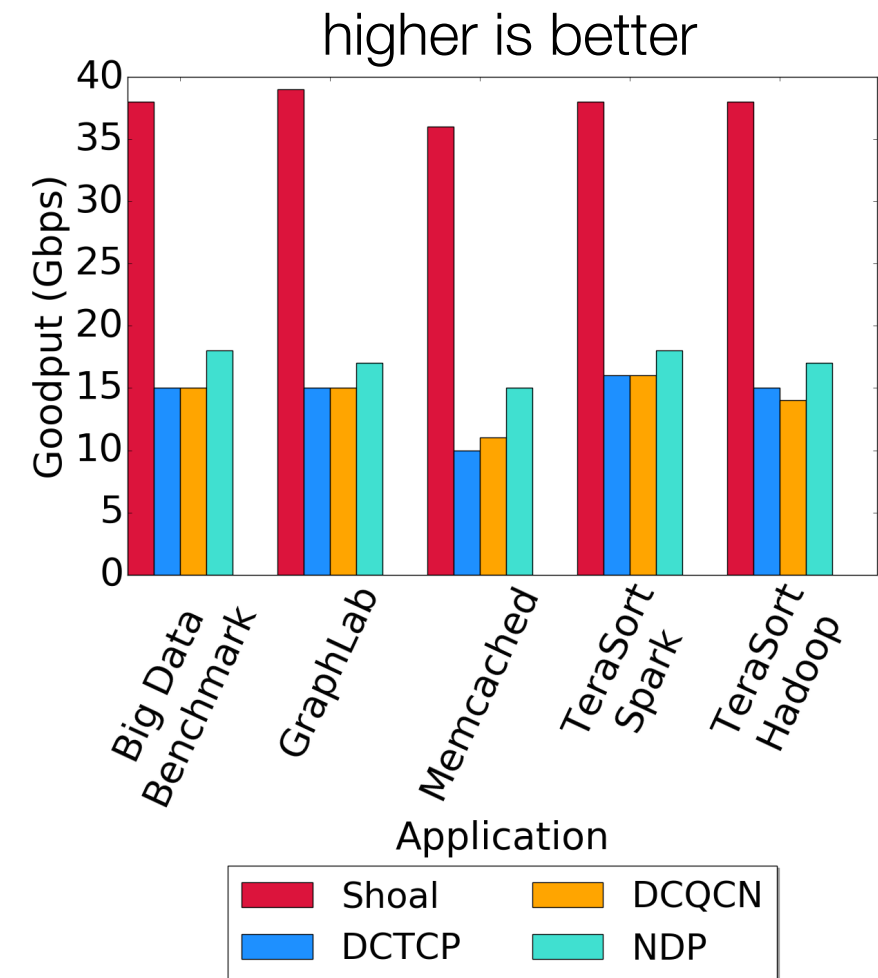
Long Flows (>1MB)

Evaluation

- Packet-level simulator in C
- 512 end-hosts
- Shoal has 2x bandwidth (100Gbps vs. 50Gbps)
at lower cost and power!



Short Flows (<100KB)



Long Flows (>1MB)

Shoal performs comparable or better than several recent packet-switched network designs

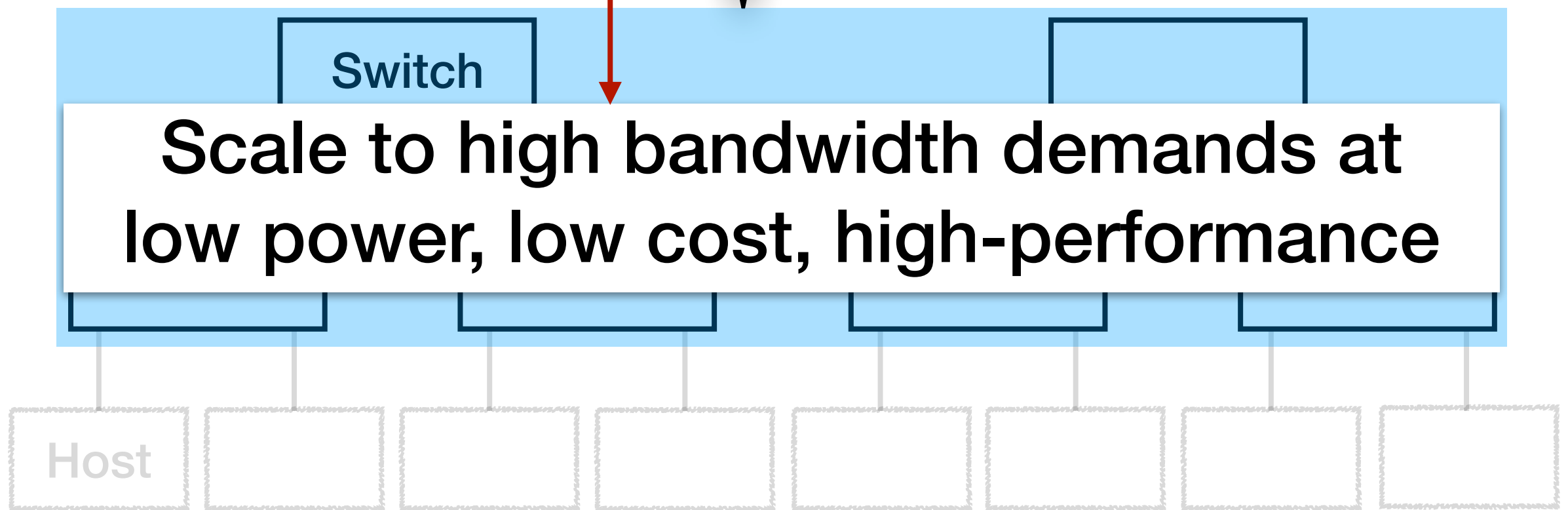
Part II : Summary

How to build high-speed switching fabric?

Switch

**Scale to high bandwidth demands at
low power, low cost, high-performance**

Host



Part II : Summary

How to build high-speed switching fabric?

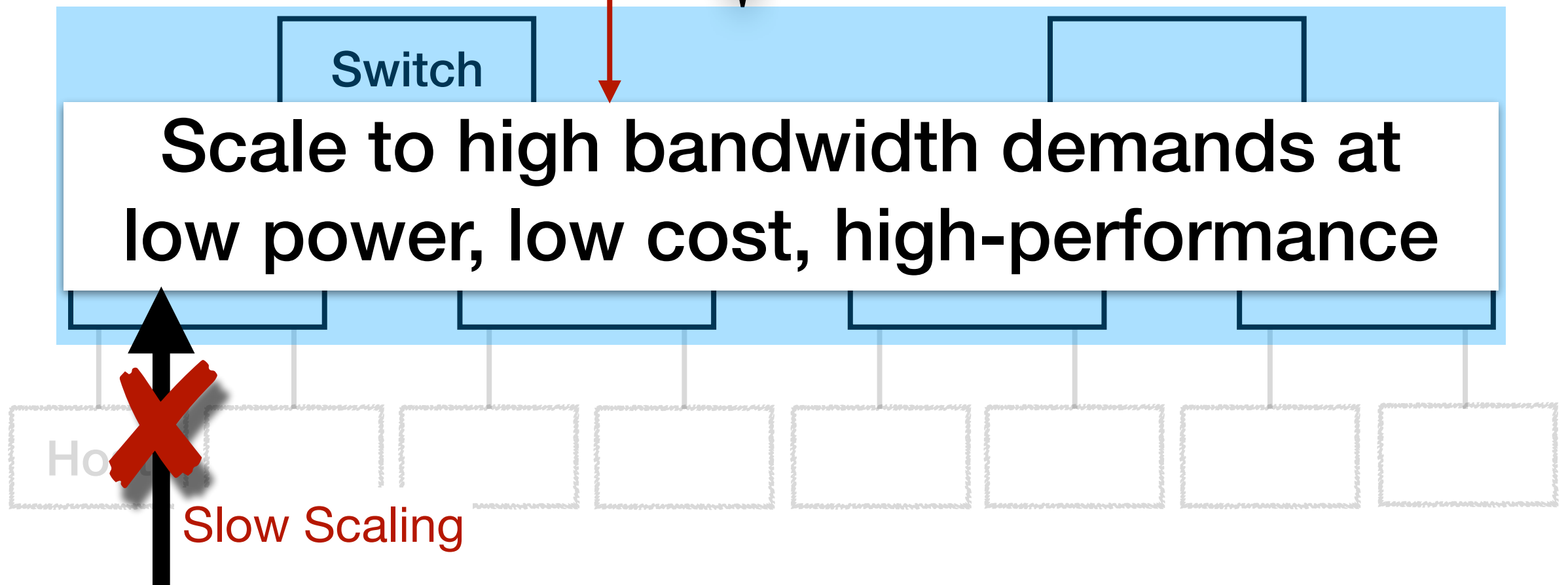
Switch

Scale to high bandwidth demands at
low power, low cost, high-performance

Ho

Slow Scaling

Packet Switches



Part II : Summary

How to build high-speed switching fabric?

Switch

Scale to high bandwidth demands at
low power, low cost, high-performance

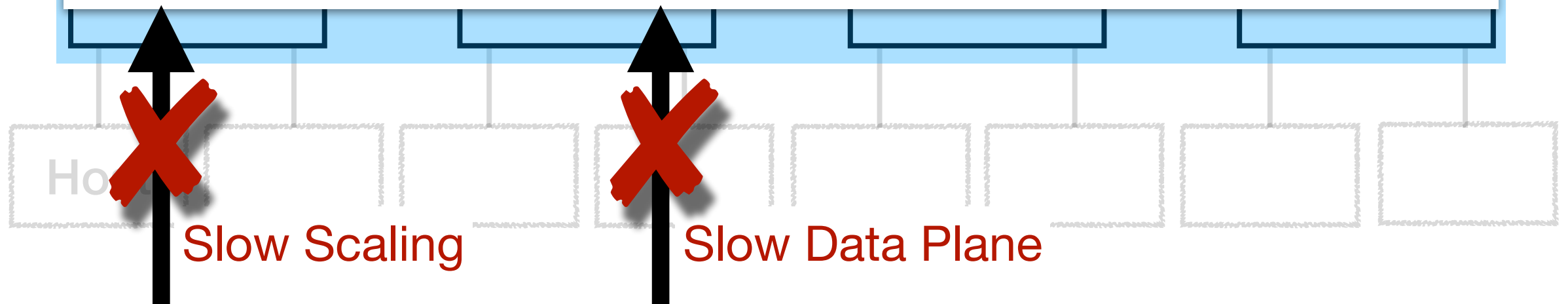
Ho

Slow Scaling

Slow Data Plane

Packet Switches

Traditional
Circuit Switches



Part II : Summary

How to build high-speed switching fabric?

Switch

Scale to high bandwidth demands at
low power, low cost, high-performance

Ho

Slow Scaling

Slow Data Plane

Fast Data Plane ✓

Packet Switches

Traditional
Circuit Switches

Fast
Circuit Switches

Part II : Summary

How to build high-speed switching fabric?

Switch

Scale to high bandwidth demands at
low power, low cost, high-performance

Ho

Slow Scaling

Slow Data Plane

Shoal
Fast Control Plane
Fast Data Plane

Packet Switches

Traditional
Circuit Switches

Fast
Circuit Switches

Future Directions

Future Directions

High-speed and Programmable Processors for
(stateful) Network Functions

e.g., load balancers, firewalls, deep packet inspectors

Future Directions

High-speed and Programmable Processors for
(stateful) Network Functions

e.g., load balancers, firewalls, deep packet inspectors

Low latency Optical Circuit Switching

Conclusion

Conclusion



Moore's Law
Dennard Scaling

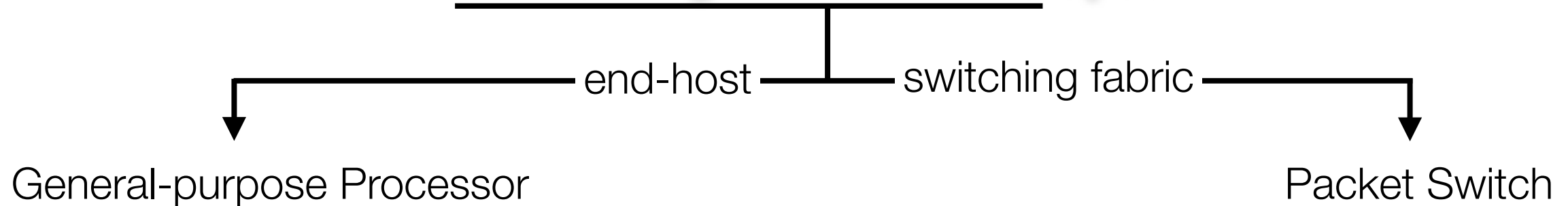
Conclusion



Moore's Law
Dennard Scaling



Networking Infrastructure Speed

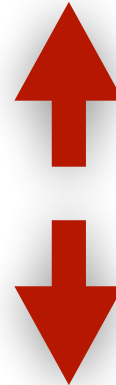


Conclusion

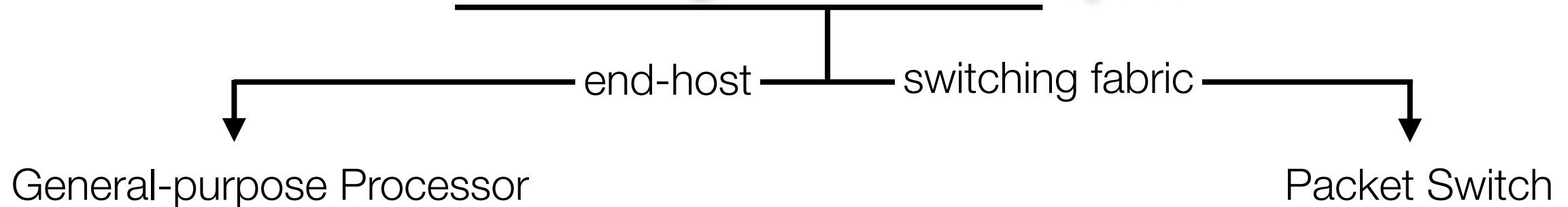


X
Moore's Law
Dennard Scaling

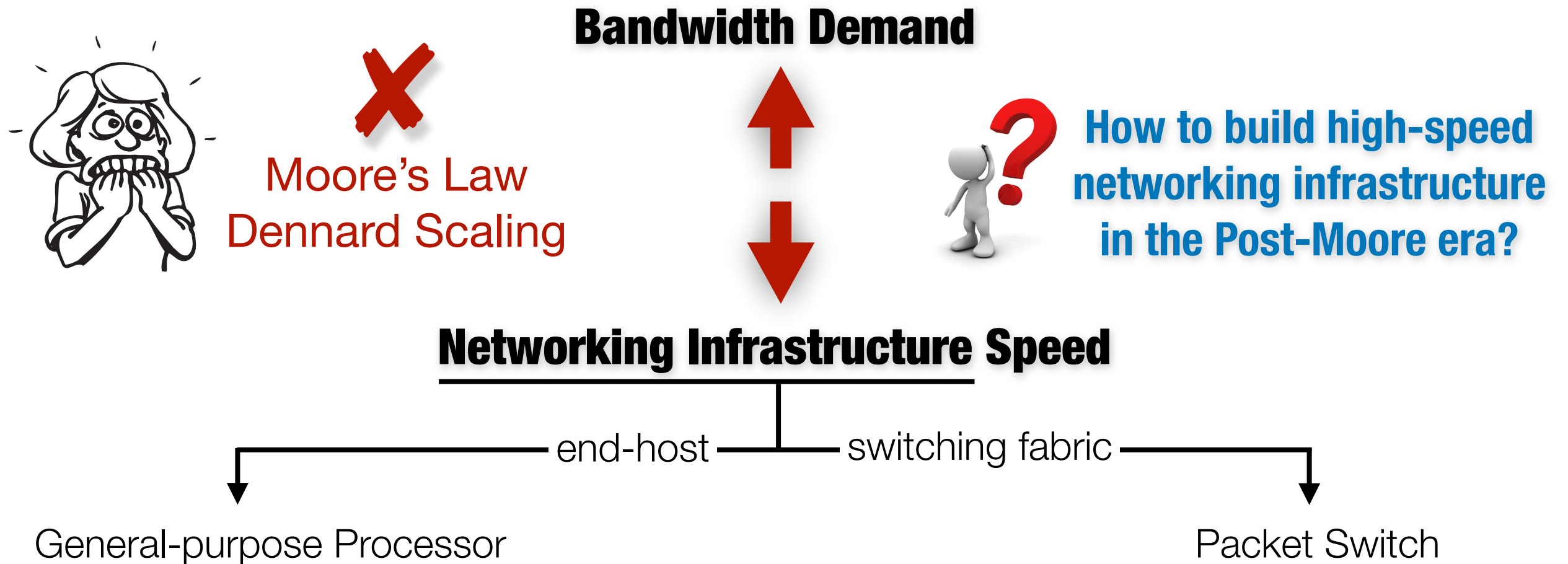
Bandwidth Demand



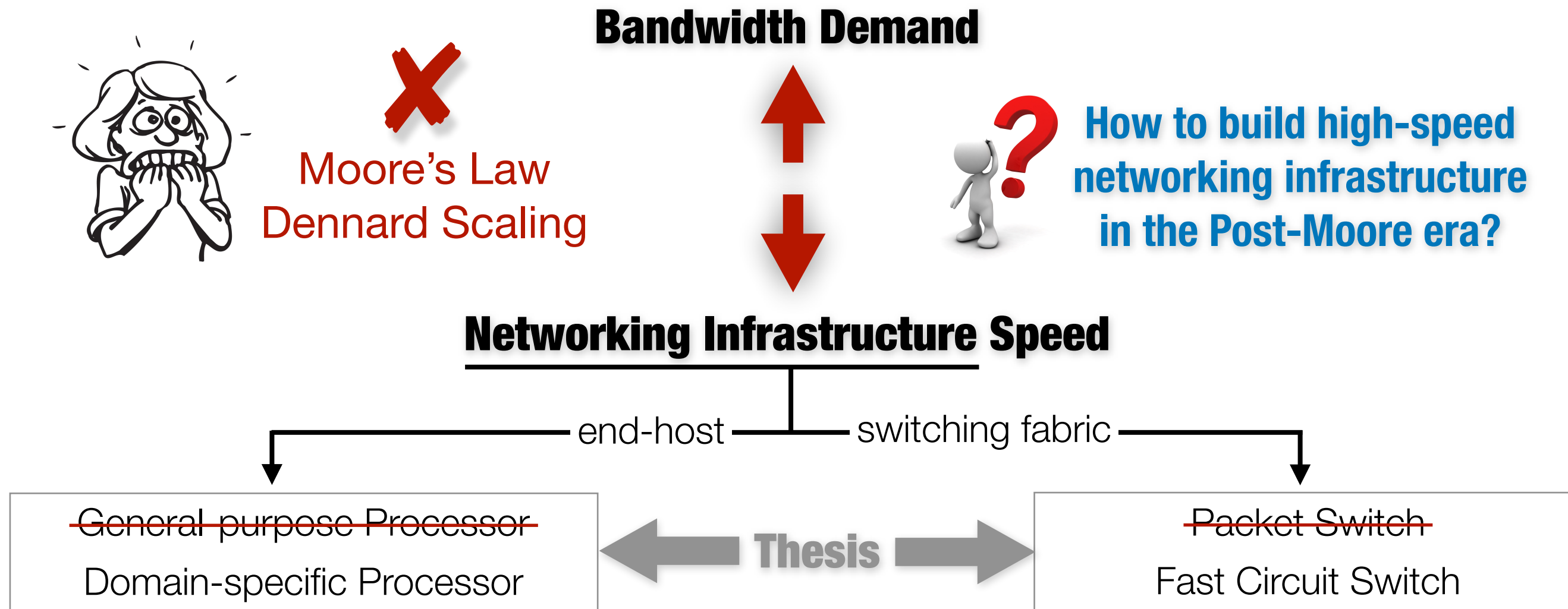
Networking Infrastructure Speed



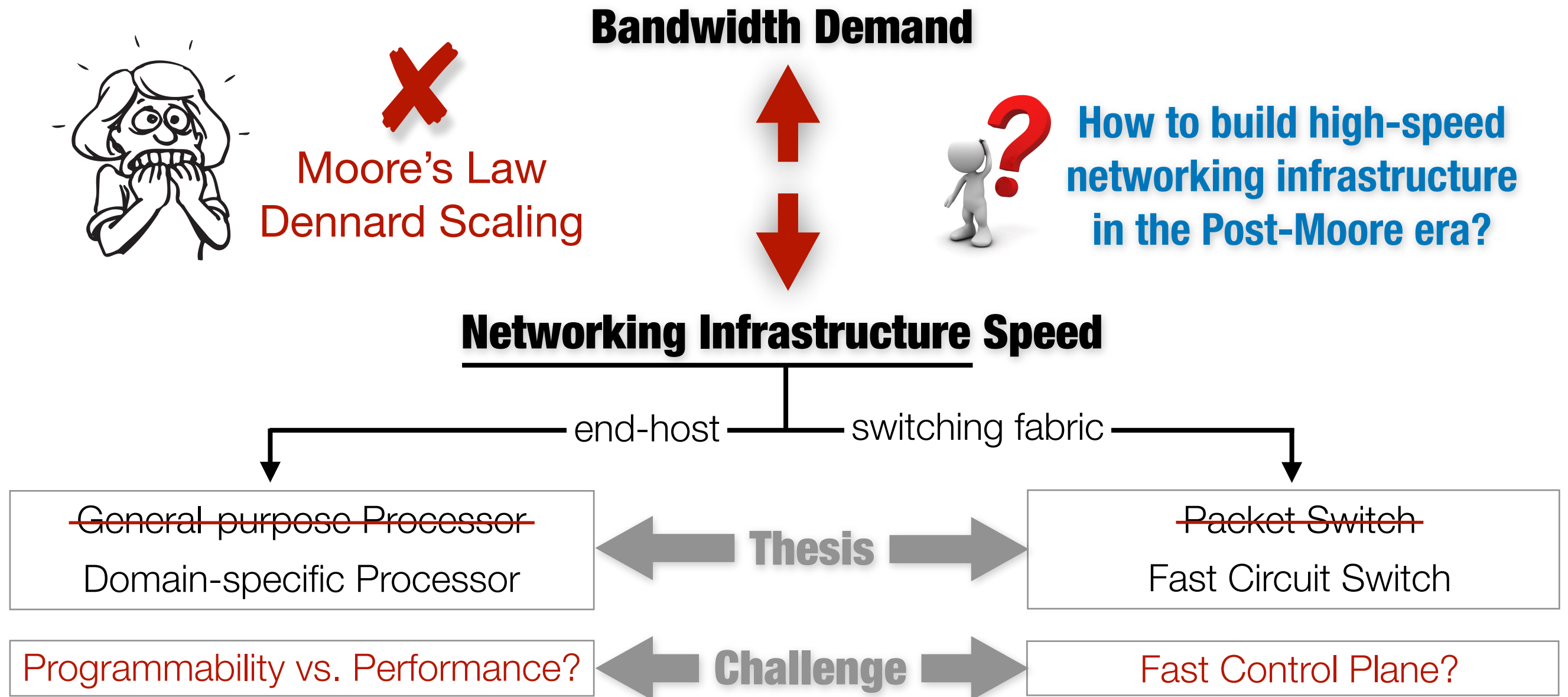
Conclusion



Conclusion



Conclusion

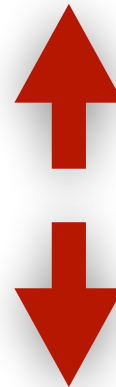


Conclusion

Bandwidth Demand

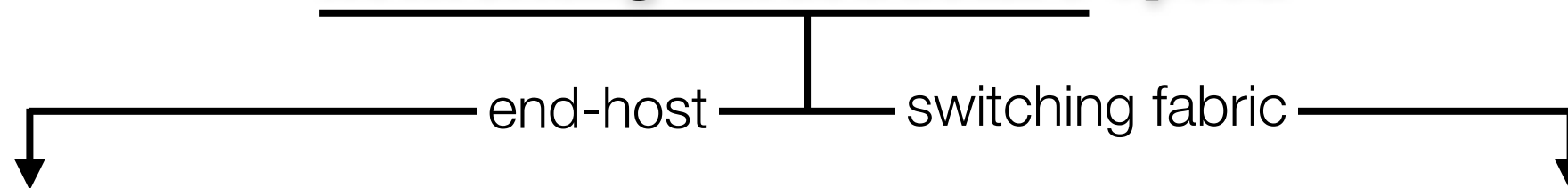


Moore's Law
Dennard Scaling



How to build high-speed
networking infrastructure
in the Post-Moore era?

Networking Infrastructure Speed



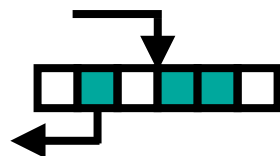
~~General purpose Processor~~
Domain-specific Processor

~~Packet Switch~~
Fast Circuit Switch

Programmability vs. Performance?

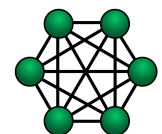
Fast Control Plane?

PIEO



Programmable and high-performance
domain-specific processor for
packet scheduling

Shoal



High-performance, ns-granularity
control plane for circuit switching

Contributions

Thesis

Challenge

Thank you!