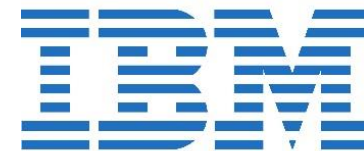


Iron: Isolating Network-based CPU in Container Environments

Junaid Khalid¹, Eric Rozner², Wesley Felter^{2,3}, Cong Xu²,
Karthick Rajamani², Alexandre Ferreira^{2,4}, Aditya Akella¹

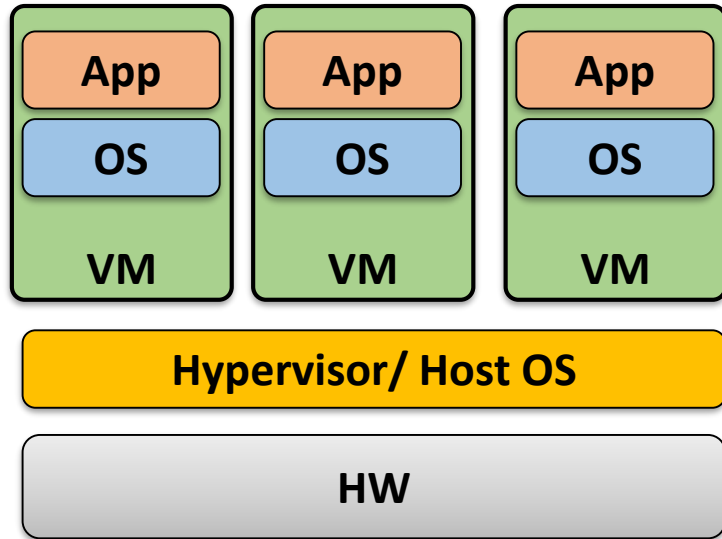


WISCONSIN
UNIVERSITY OF WISCONSIN-MADISON

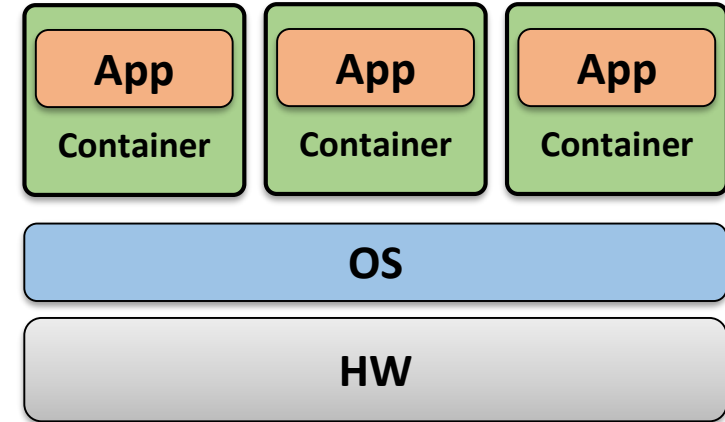


Virtualization

Virtual machines

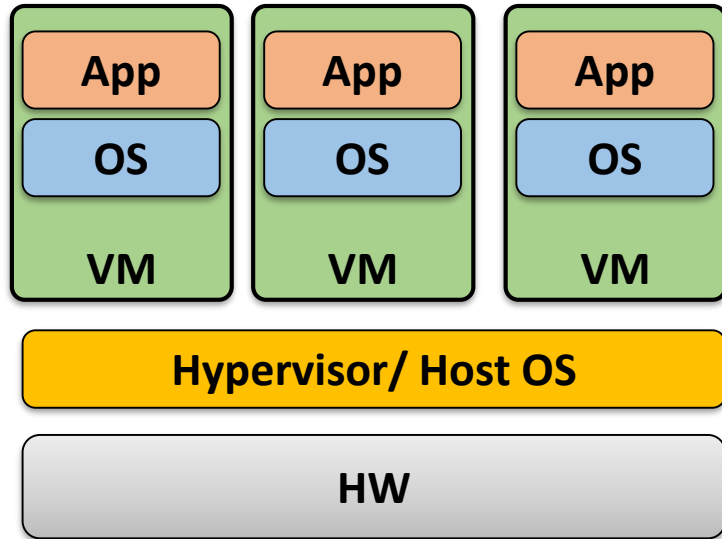


Containers



Virtualization

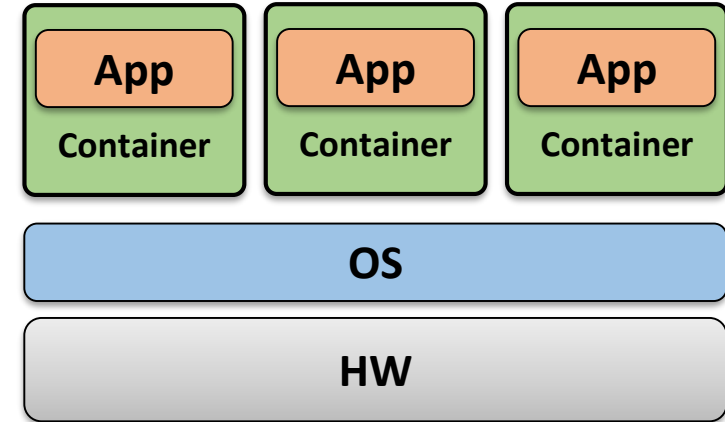
Virtual machines



✗ Heavy weight

✗ Slow

Containers

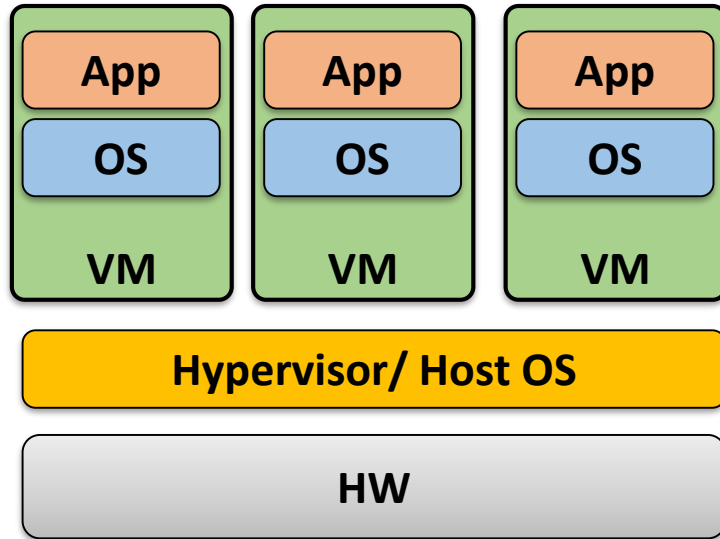


✓ Light weight

✓ Fast

Virtualization

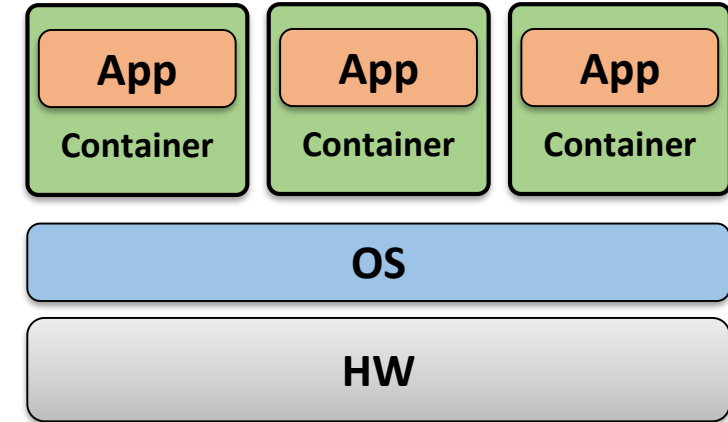
Virtual machines



✗ Heavy weight

✗ Slow

Containers



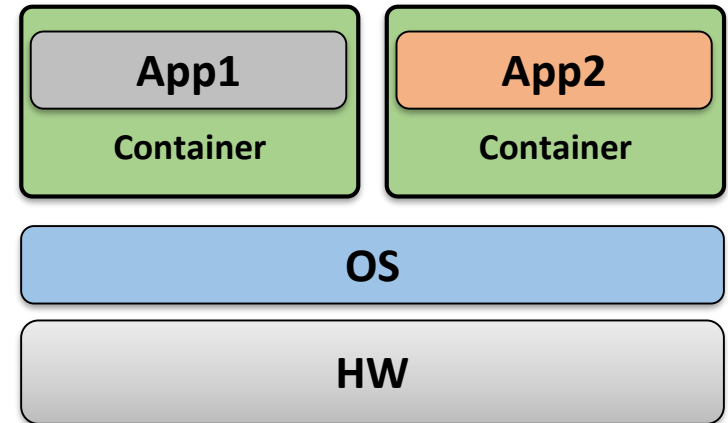
✓ Light weight

✓ Fast

Supports new workloads such as microservices and serverless

Isolation

Containers

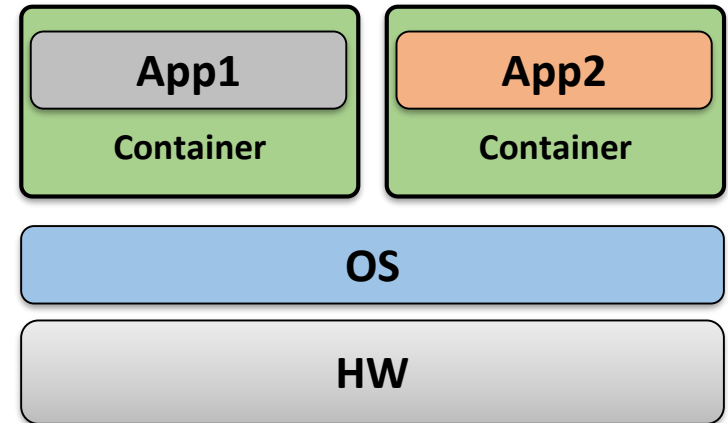


Isolation

Containers require strong resource isolation

- Memory
- Network
- CPU

Containers



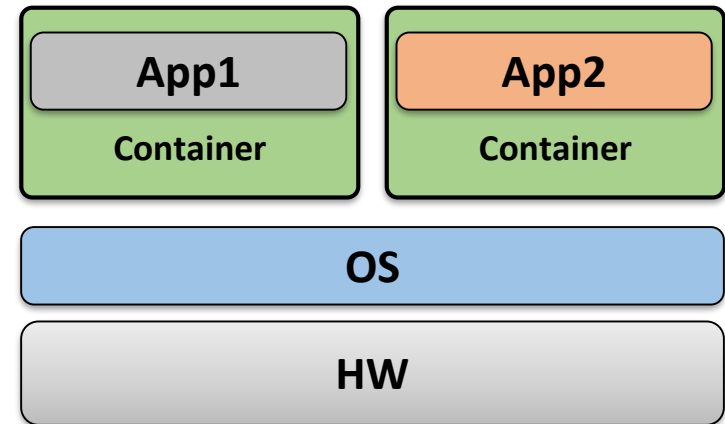
Isolation

Containers require strong resource isolation

- Memory
- Network
- CPU

Administrators want to strongly control resource allocation in multi-tenant environments

Containers



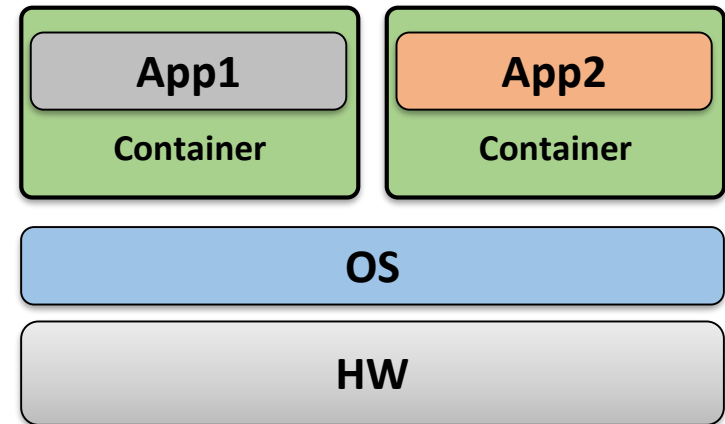
Isolation

Containers

Containers require strong resource isolation

- Memory
- Network
- CPU

Administrators want to strongly control resource allocation in multi-tenant environments

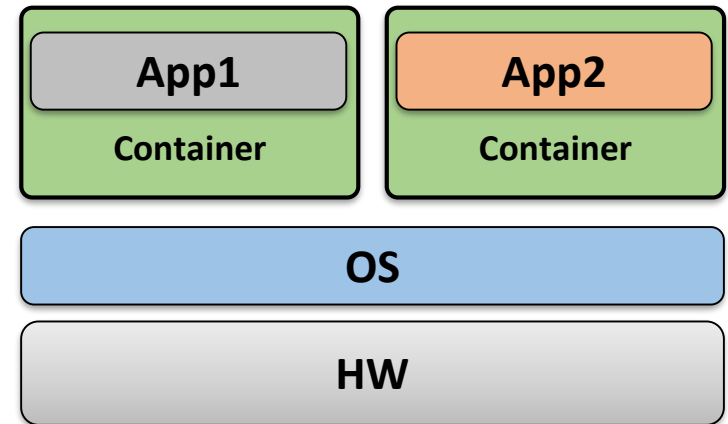


Strong isolation is important for ***performance***,
predictability and ***efficiency***

Isolation

Isolation: A container *shouldn't consume* more than its assigned share of *resources*

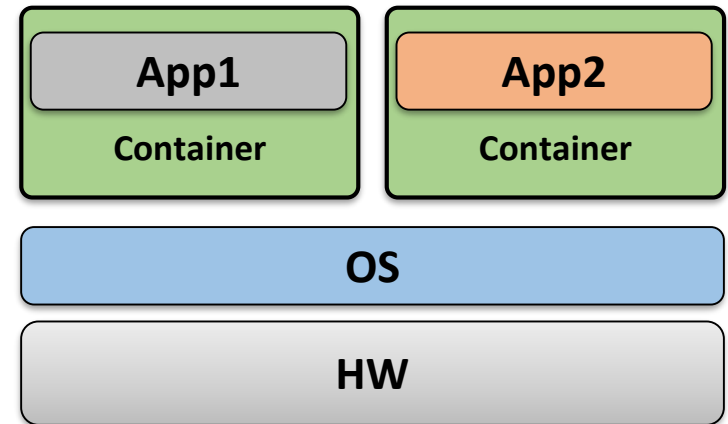
Containers



Isolation

Isolation: A container *shouldn't consume* more than its assigned share of *resources*

Containers



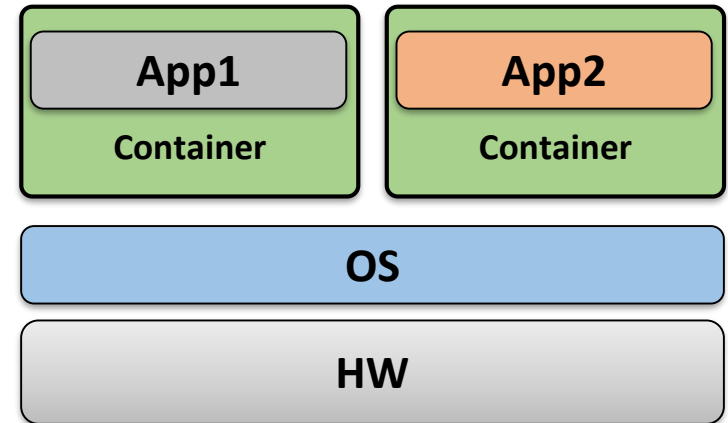
allocated share



Isolation

Isolation: A container *shouldn't consume* more than its assigned share of *resources*

Containers



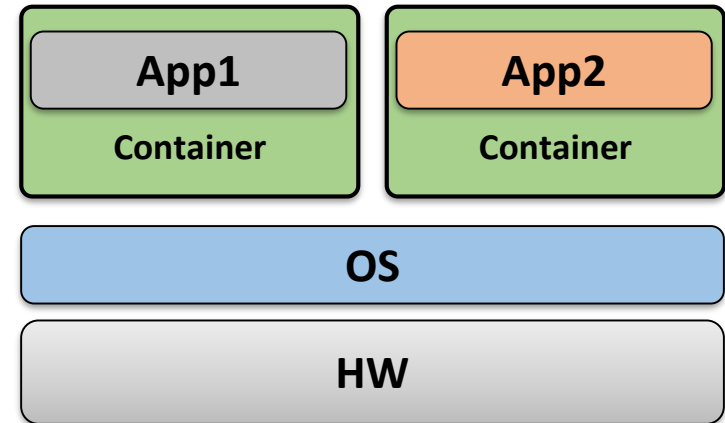
allocated share	50%	50%
actual usage	50%	50%

Isolation

Isolation: A container *shouldn't consume* more than its assigned share of *resources*

cgroups ensures **CPU isolation** by *allocating*, *metering*, and *enforcing* resource usage in the kernel

Containers



allocated share
actual usage



Isolation

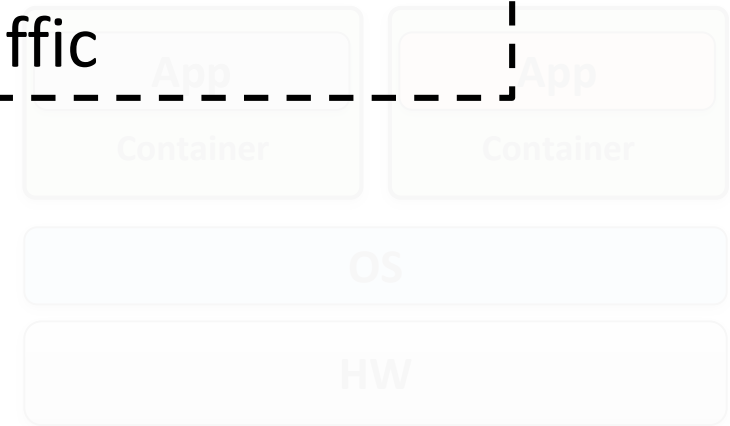
Containers

Isolation: A container *shouldn't consume* more than its assigned share of resources

CPU isolation provided by Linux breaks down while handling the network traffic

network
intensive

compute
intensive



cGroup ensures resource isolation by *allocating, metering, and enforcing* resource usage in the kernel

allocated share
actual usage



Outline

- How and by how much is isolation broken

Outline

- How and by how much is isolation broken
- Iron's design
 - Accounting of per-packet processing cost
 - Ensuring isolation via enforcement
 - Integration with Linux scheduler
 - Hardware-based packet dropping

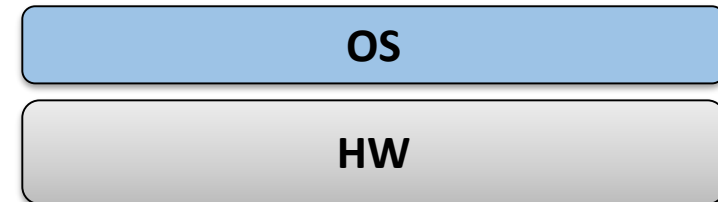
Outline

- How and by how much is isolation broken
- Iron's design
 - Accounting of per-packet processing cost
 - Ensuring isolation via enforcement
 - Integration with Linux scheduler
 - Hardware-based packet dropping
- Evaluation
 - Controlled workload
 - Realistic workload

Isolation is Broken

Containers

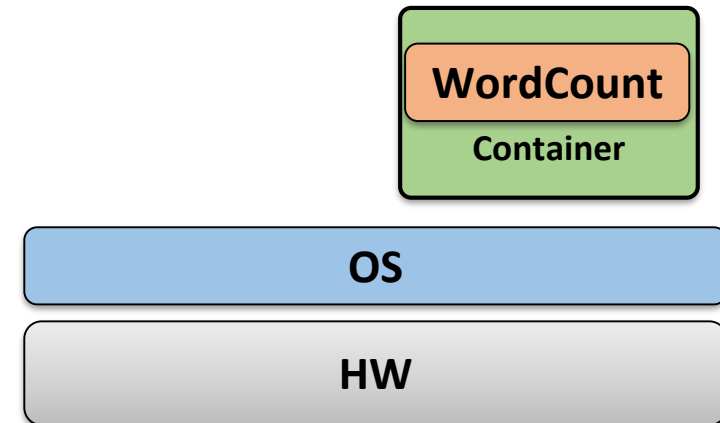
$$\text{Penalty factor} = \frac{\text{Time that job takes when competing with traffic}}{\text{Time that job takes when competing with compute}}$$



Isolation is Broken

Containers

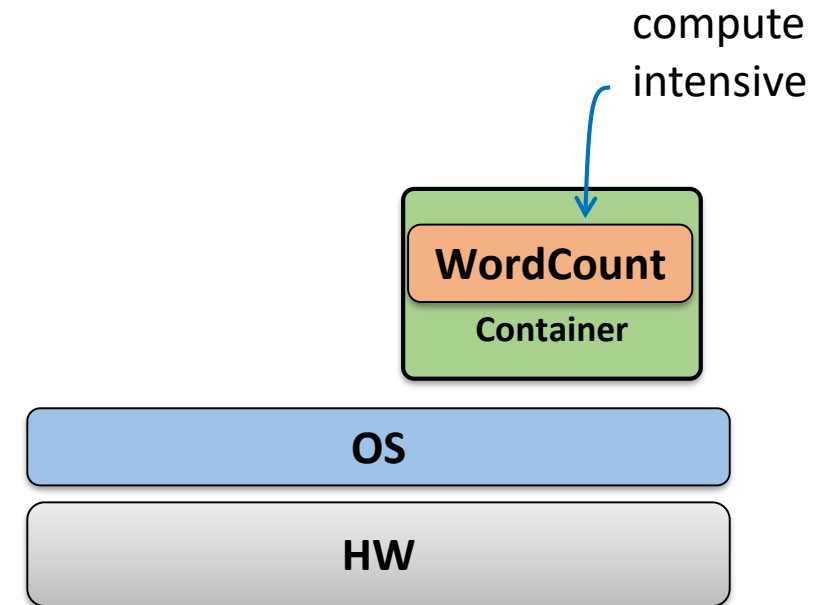
$$\text{Penalty factor} = \frac{\text{Time that job takes when competing with traffic}}{\text{Time that job takes when competing with compute}}$$



Isolation is Broken

Containers

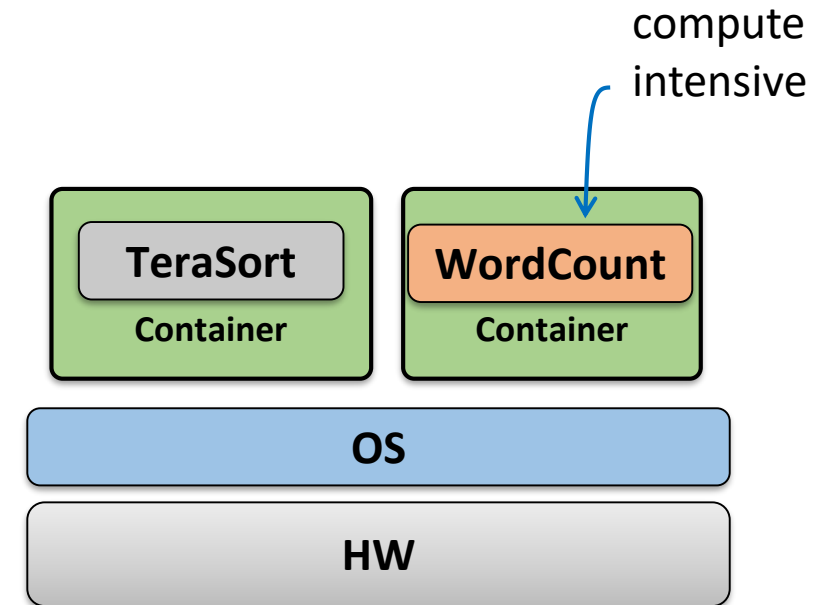
$$\text{Penalty factor} = \frac{\text{Time that job takes when competing with traffic}}{\text{Time that job takes when competing with compute}}$$



Isolation is Broken

Containers

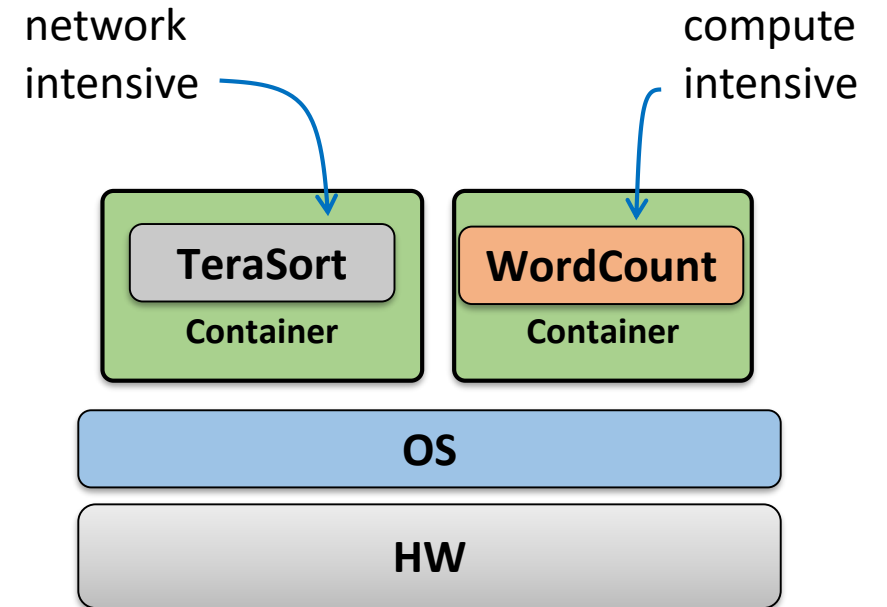
$$\text{Penalty factor} = \frac{\text{Time that job takes when competing with traffic}}{\text{Time that job takes when competing with compute}}$$



Isolation is Broken

Containers

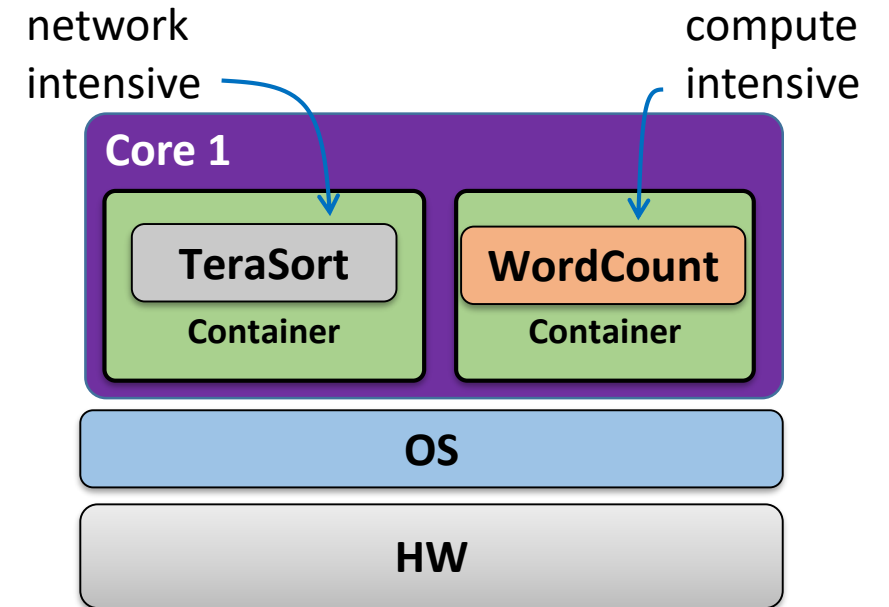
$$\text{Penalty factor} = \frac{\text{Time that job takes when competing with traffic}}{\text{Time that job takes when competing with compute}}$$



Isolation is Broken

Containers

$$\text{Penalty factor} = \frac{\text{Time that job takes when competing with traffic}}{\text{Time that job takes when competing with compute}}$$

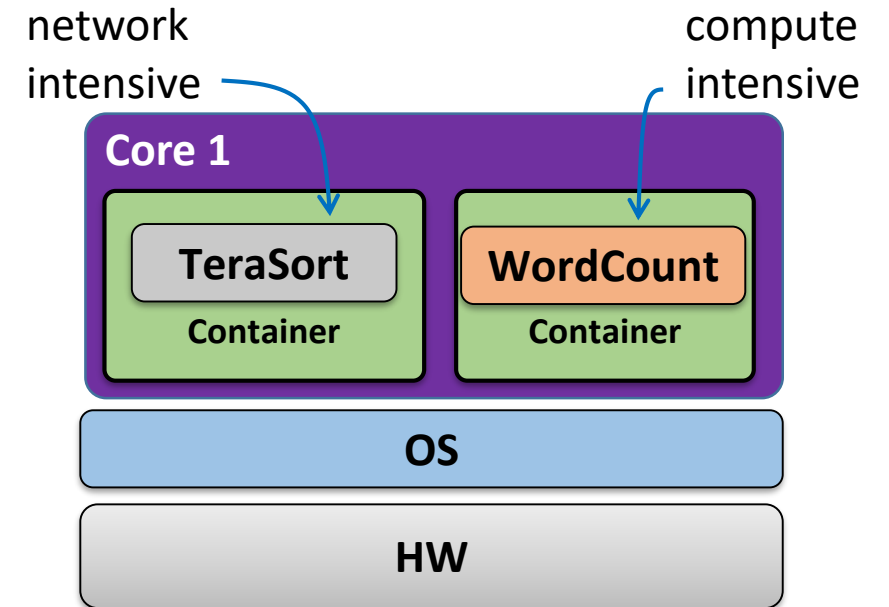


Isolation is Broken

Containers

$$\text{Penalty factor} = \frac{\text{Time that job takes when competing with traffic}}{\text{Time that job takes when competing with compute}}$$

Wordcount can take **1.5x** longer when it shares the core with **TeraSort**

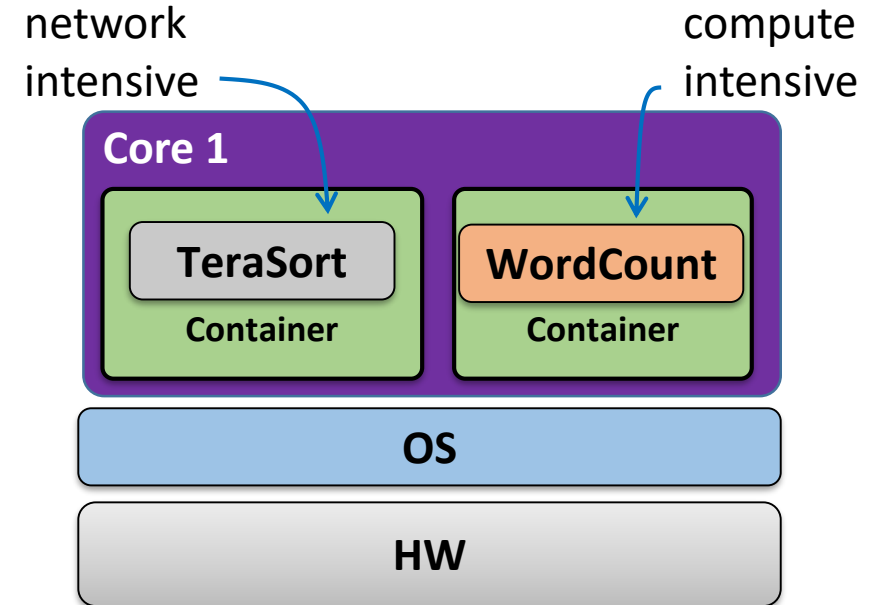


Isolation is Broken

Containers

$$\text{Penalty factor} = \frac{\text{Time that job takes when competing with traffic}}{\text{Time that job takes when competing with compute}}$$

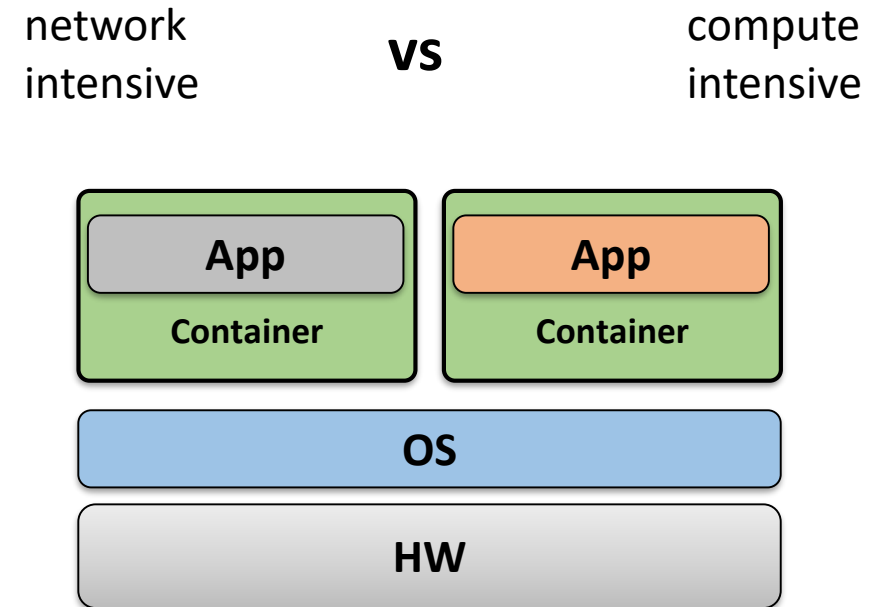
Wordcount can take **1.5x** longer when it shares the core with **TeraSort**
vs
running alone



Practical Implications

Containers

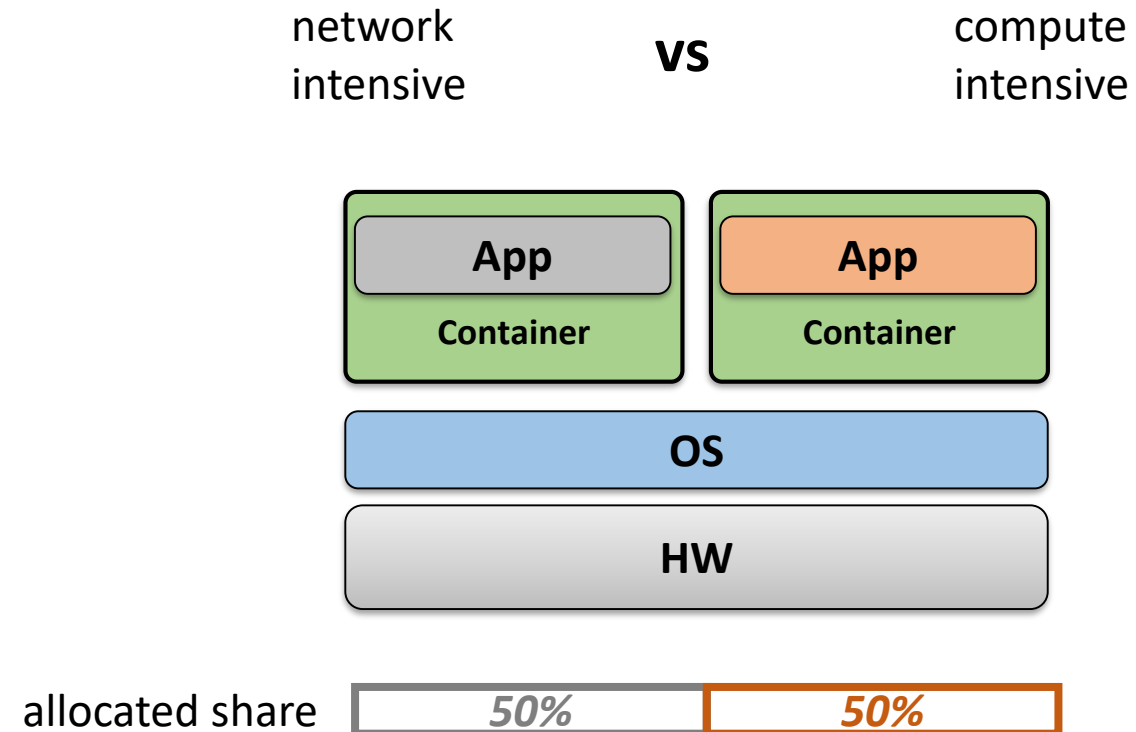
1) Insufficient isolation →
overcharging & high variance
in the performance



Practical Implications

Containers

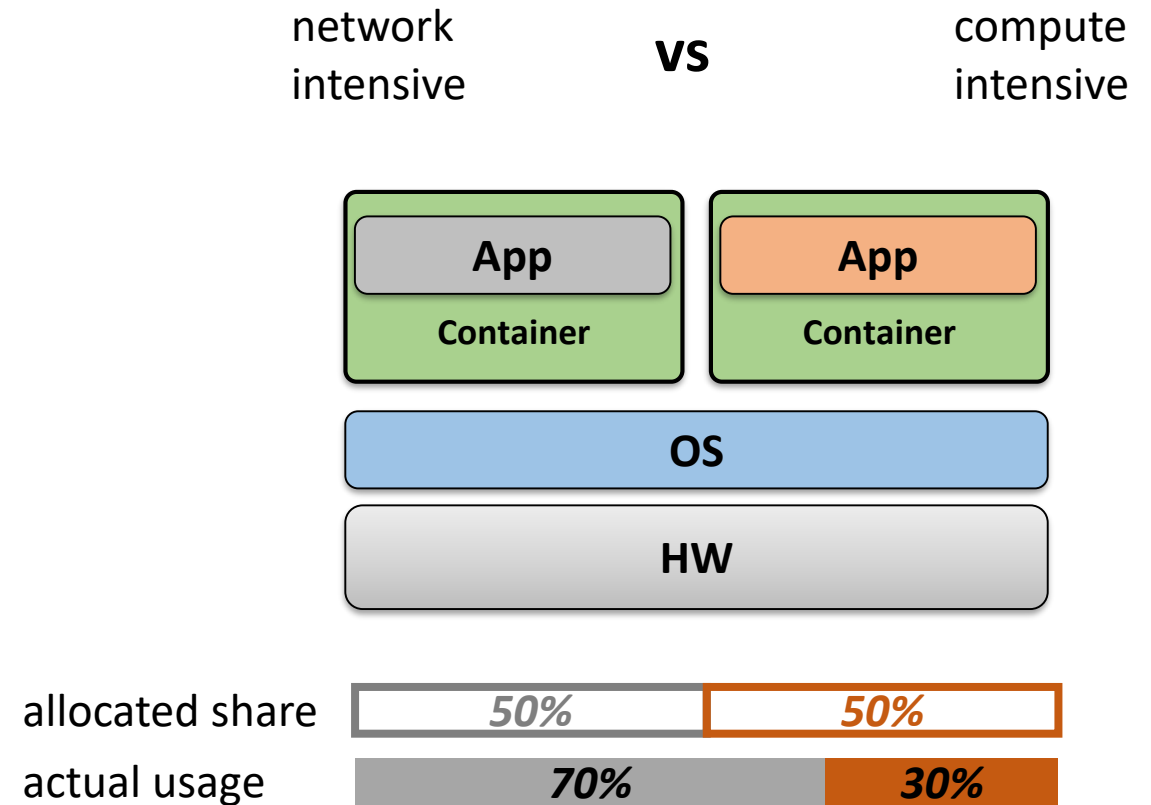
1) Insufficient isolation →
overcharging & high variance
in the performance



Practical Implications

Containers

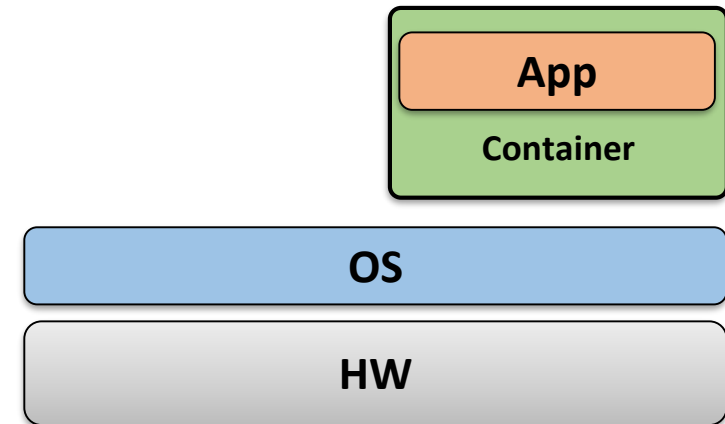
1) Insufficient isolation →
overcharging & high variance
in the performance



Practical Implications

Containers

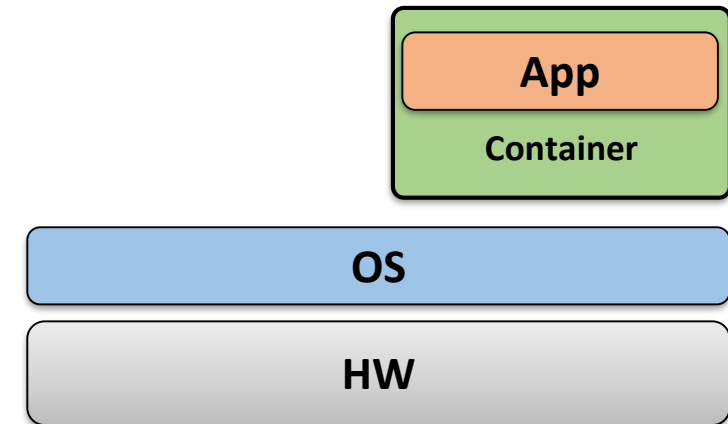
- 1) Insufficient isolation →
overcharging & high variance
in the performance
- 2) Under provisioning →
waste of potential revenue



Practical Implications

Containers

- 1) Insufficient isolation →
overcharging & high variance
in the performance
- 2) Under provisioning →
waste of potential revenue



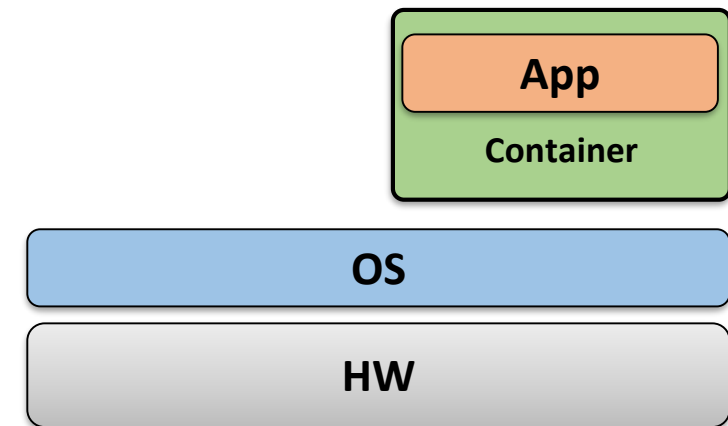
allocated share

100%

Practical Implications

Containers

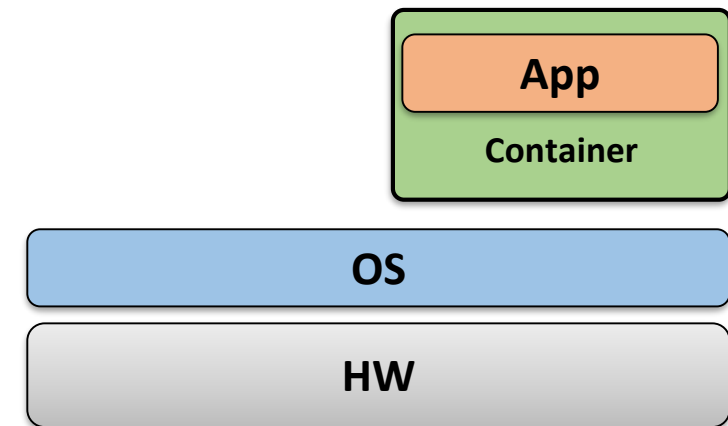
- 1) Insufficient isolation →
overcharging & high variance
in the performance
- 2) Under provisioning →
waste of potential revenue



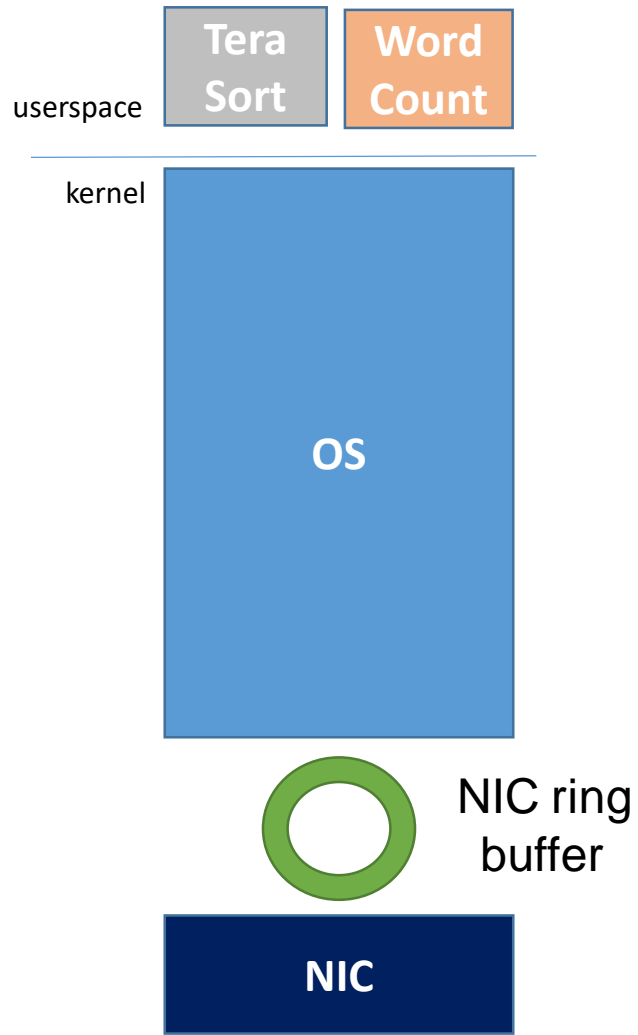
Practical Implications

Containers

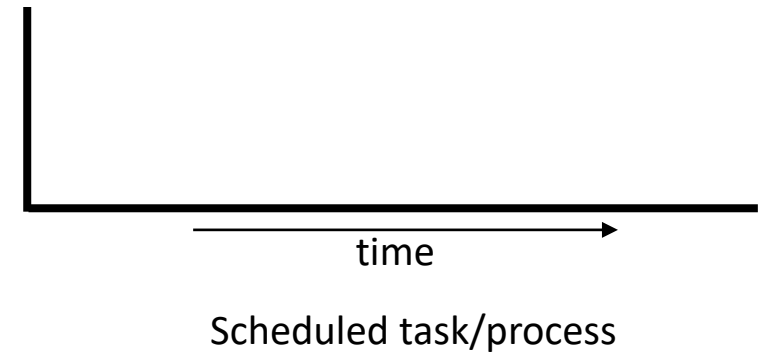
- 1) Insufficient isolation →
overcharging & high variance
in the performance
- 2) Under provisioning →
waste of potential revenue



How is Isolation Broken?

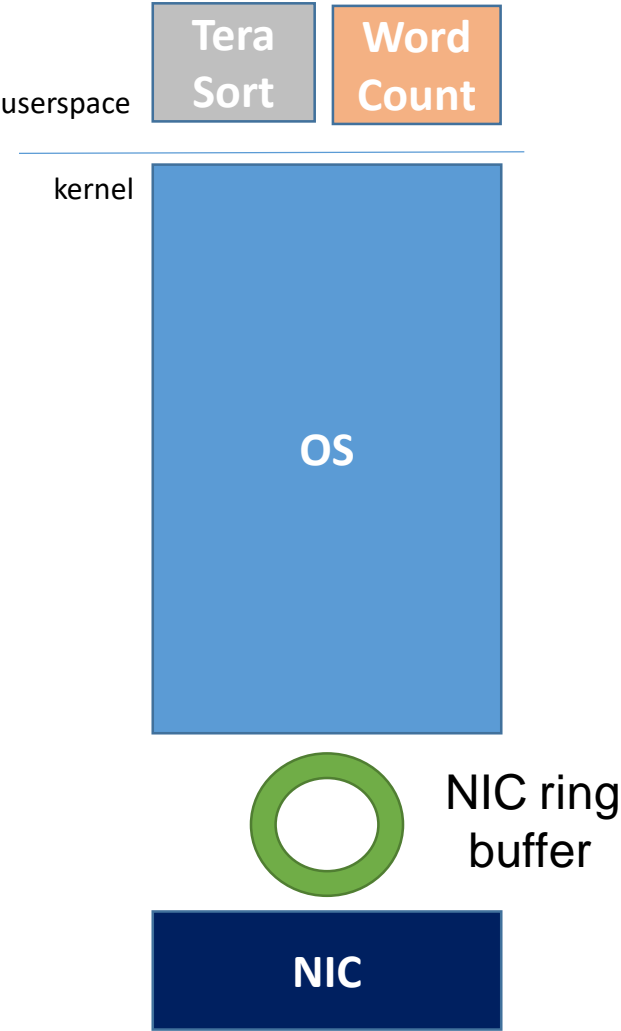


Interrupt handler
wordcount

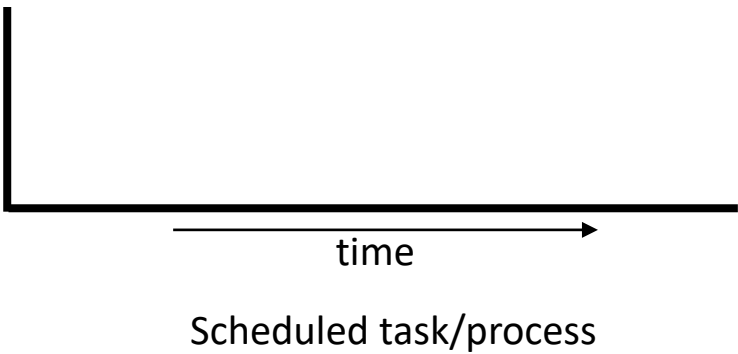


How is Isolation Broken?

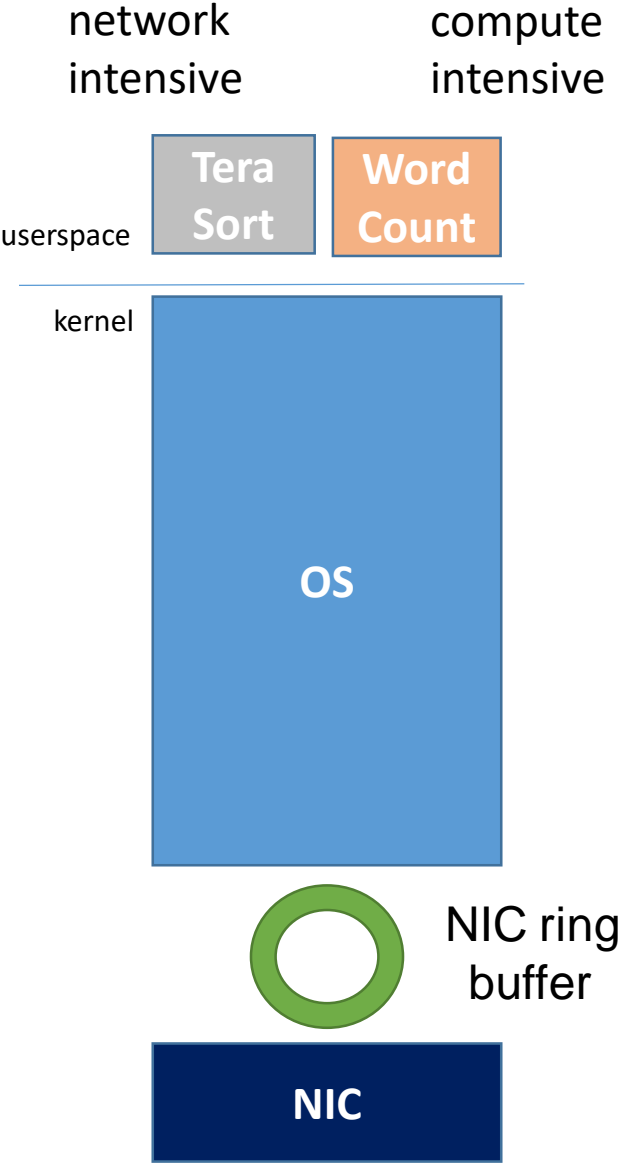
compute
intensive



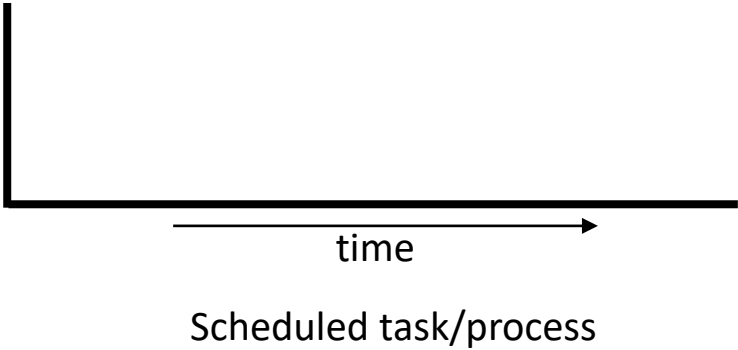
Interrupt handler
wordcount



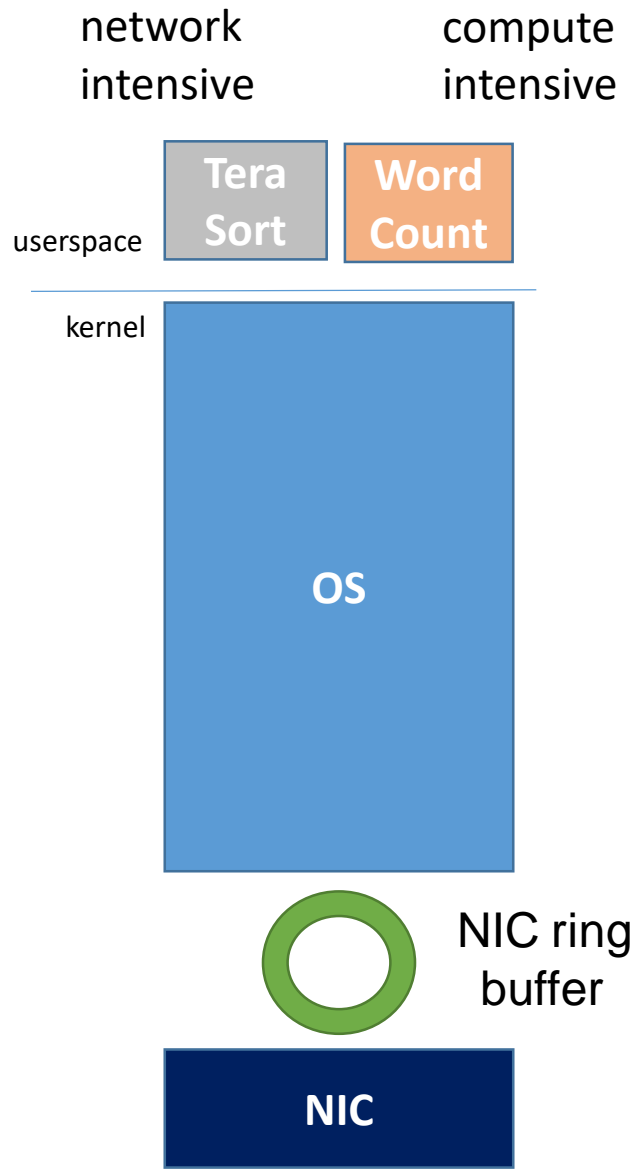
How is Isolation Broken?



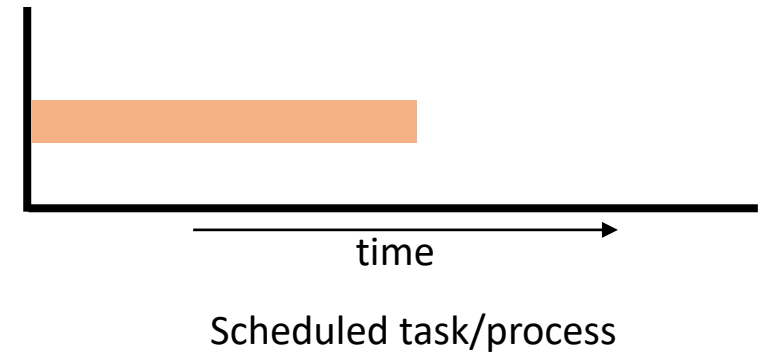
Interrupt handler
wordcount



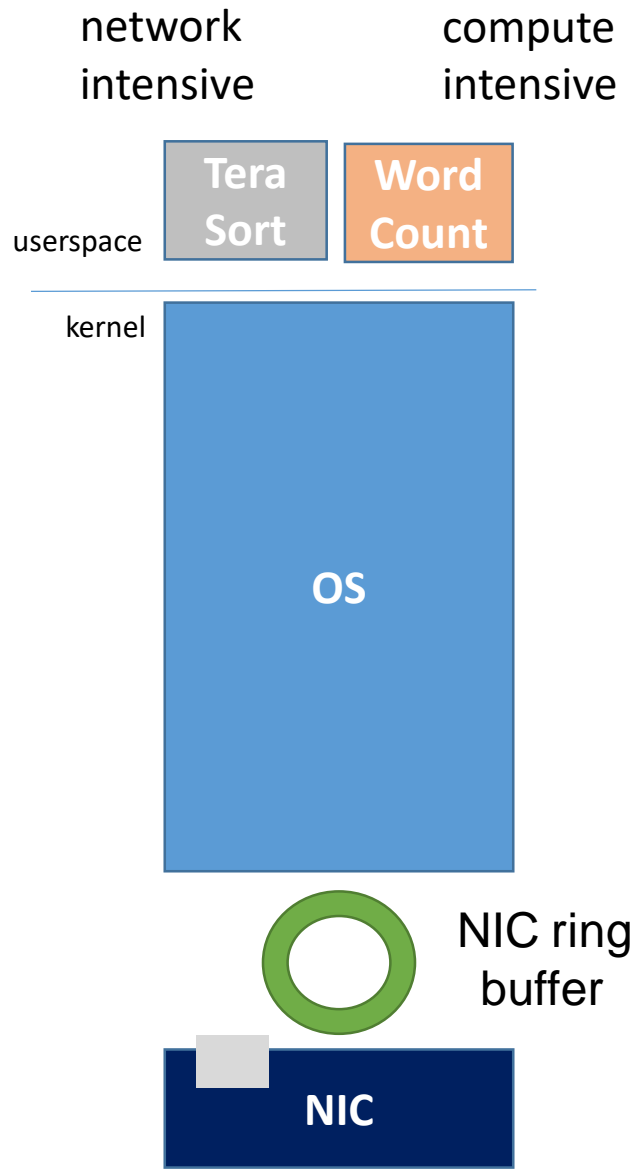
How is Isolation Broken?



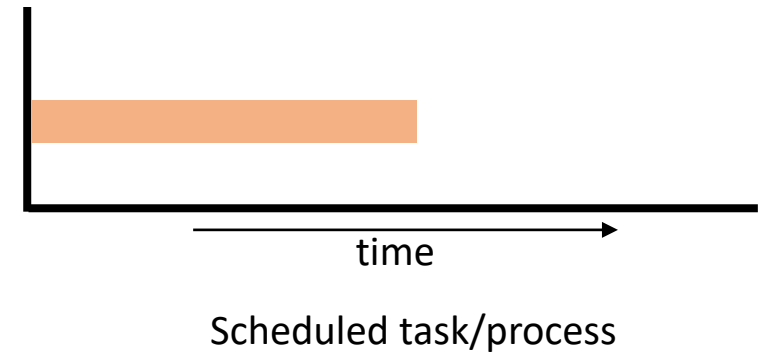
Interrupt handler
wordcount



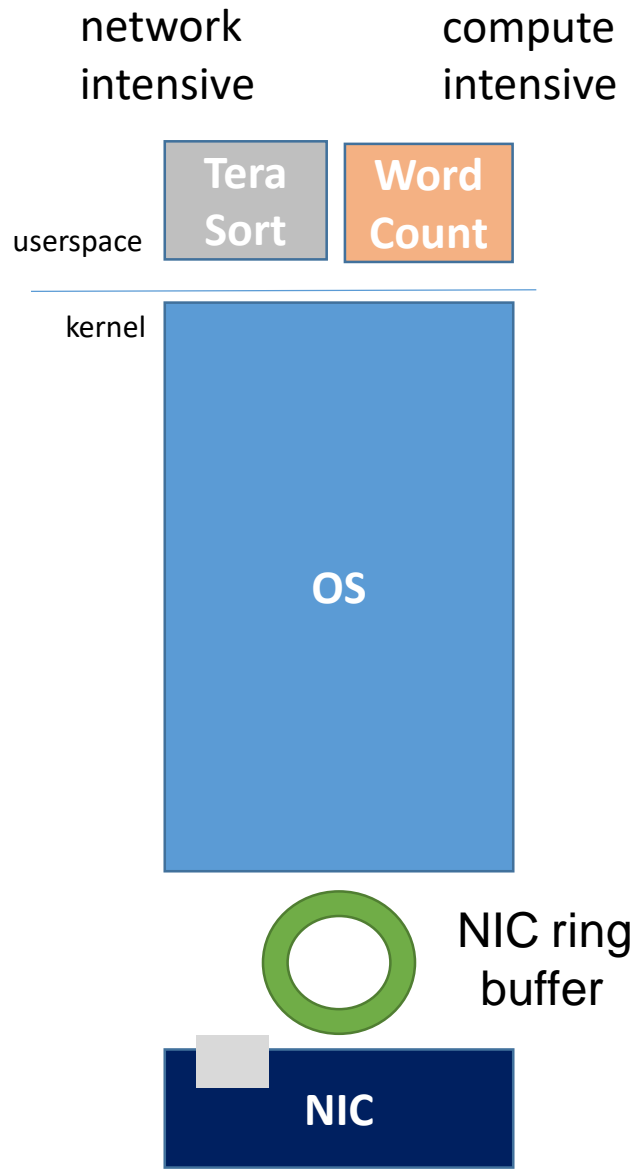
How is Isolation Broken?



Interrupt handler
wordcount

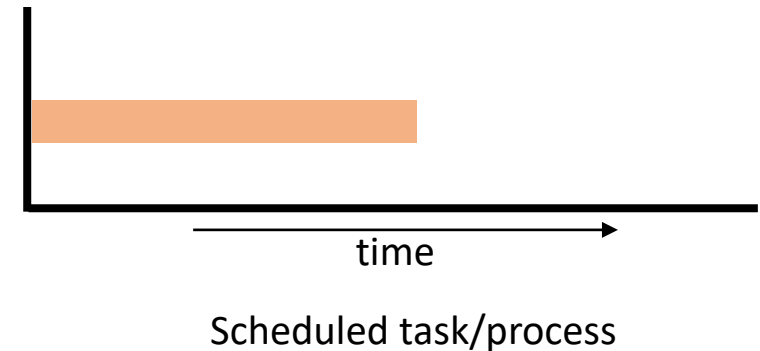


How is Isolation Broken?

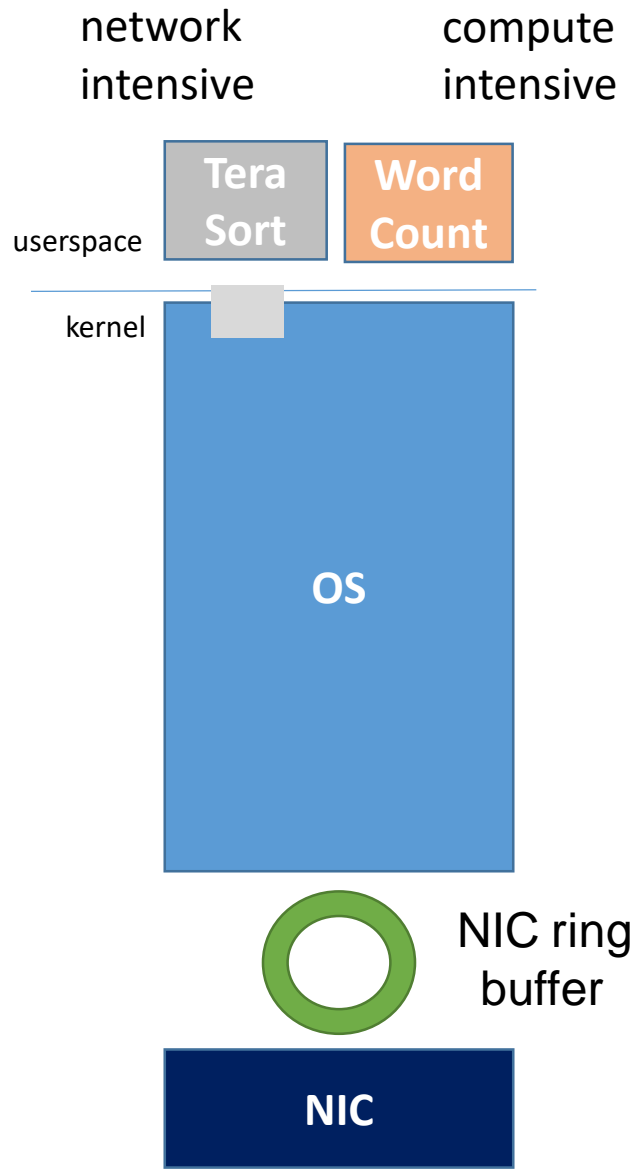


- Kernel processes network traffic via interrupts

Interrupt handler
wordcount

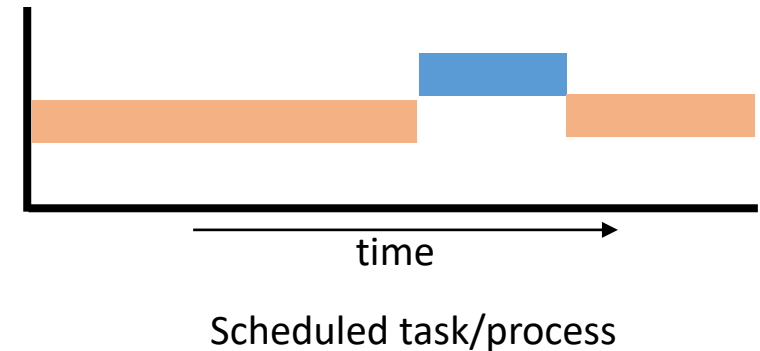


How is Isolation Broken?

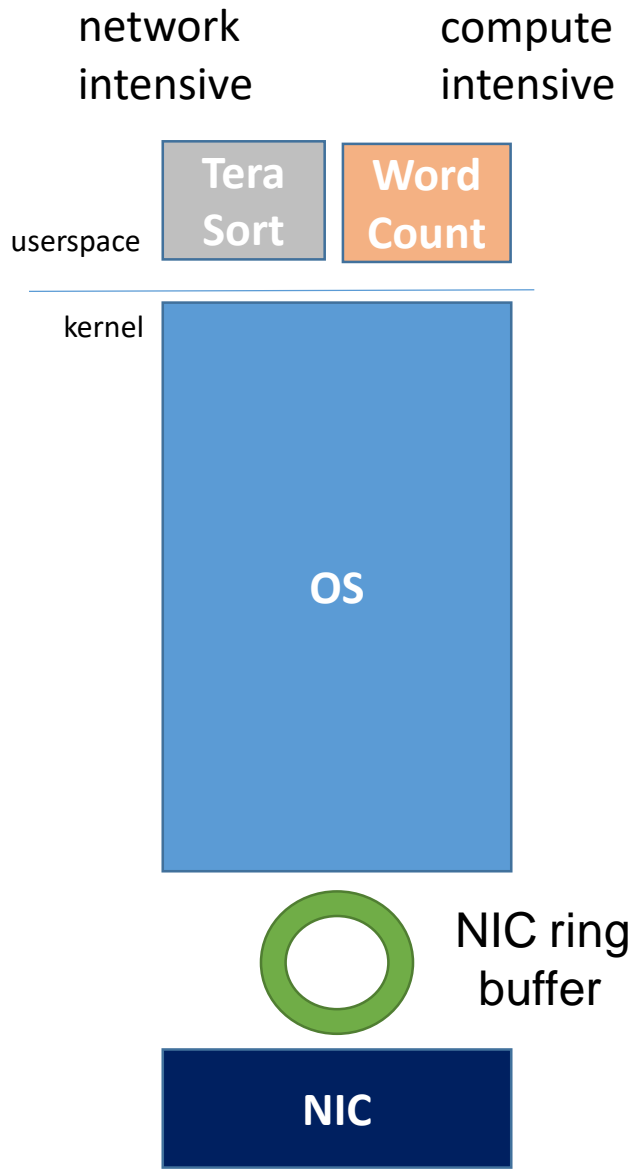


- Kernel processes network traffic via interrupts

Interrupt handler
wordcount



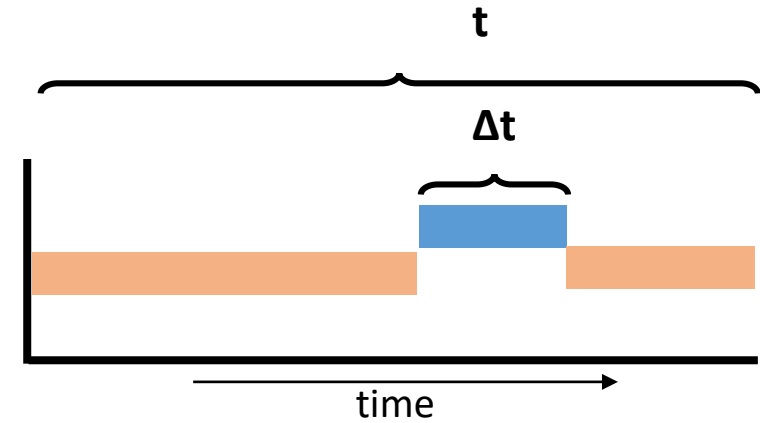
How is Isolation Broken?



- Kernel processes network traffic via interrupts

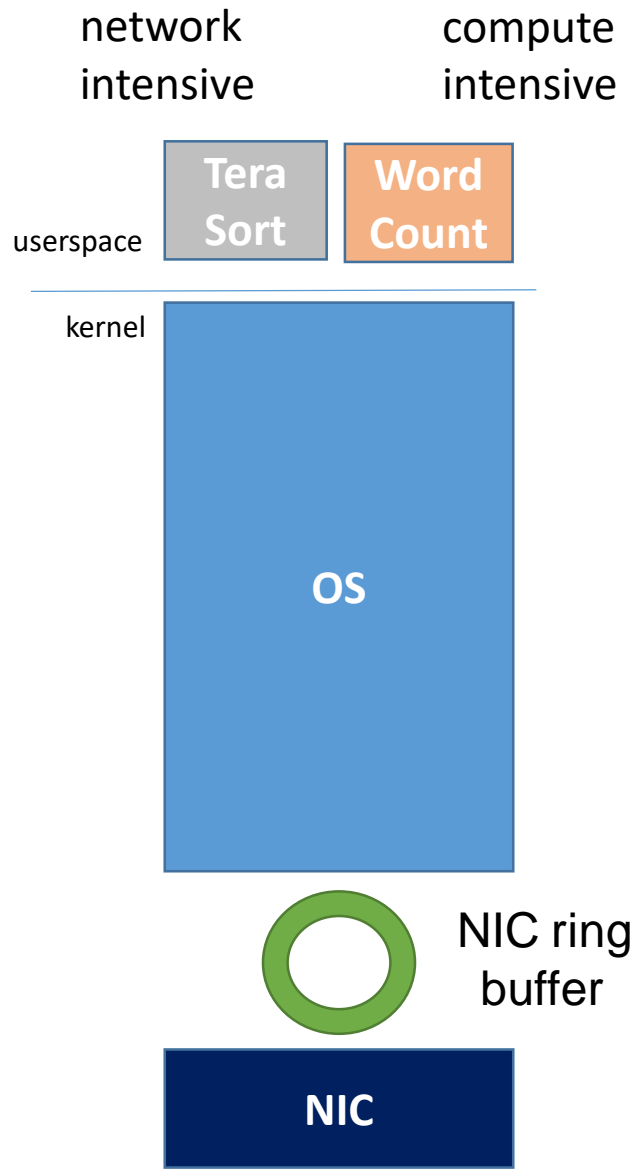
Interrupt handler
wordcount

Scheduler **charges** wordcount for "t", instead of " $t - \Delta t$ "



Scheduled task/process

How is Isolation Broken?

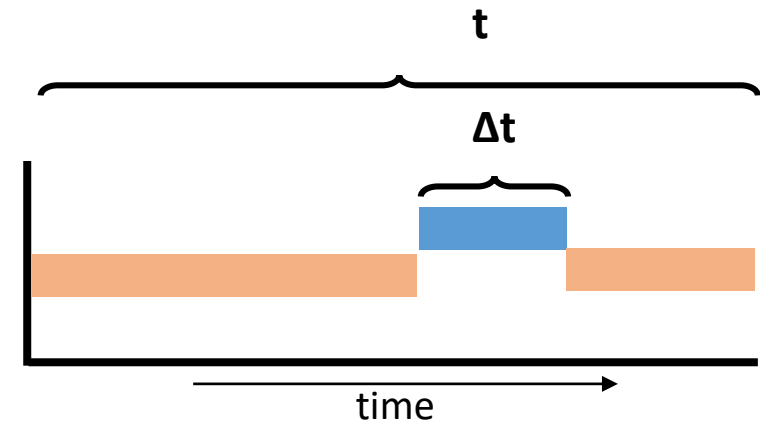


- Kernel processes network traffic via interrupts

Charging: Reduction in the runtime of a container

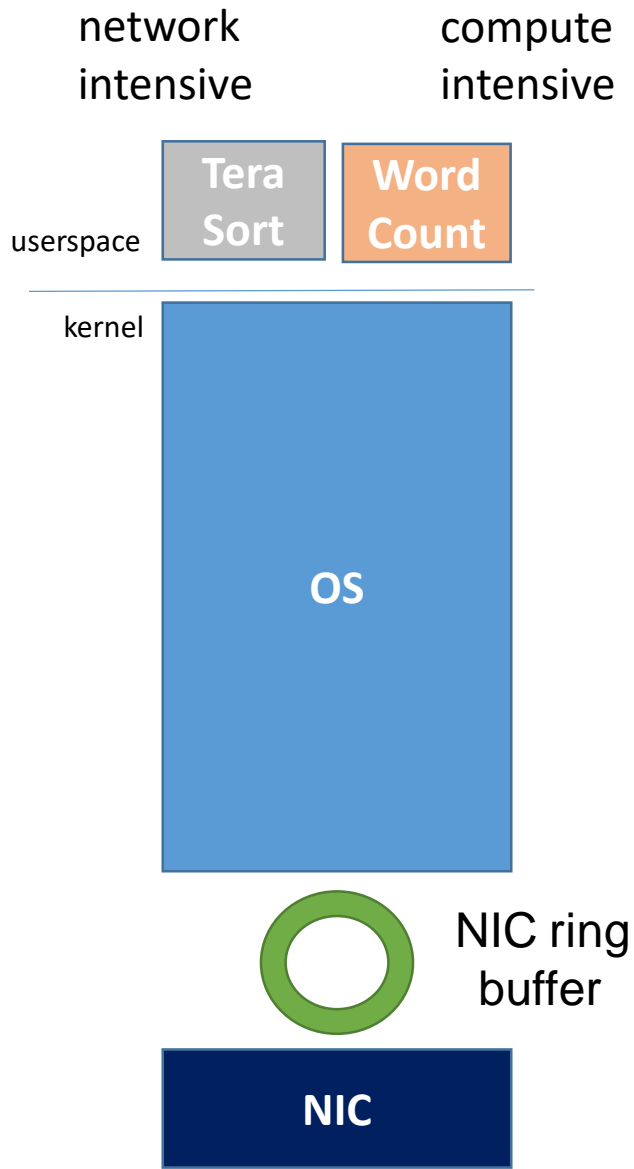
Scheduler **charges** wordcount for "t", instead of " $t - \Delta t$ "

Interrupt handler
wordcount



Scheduled task/process

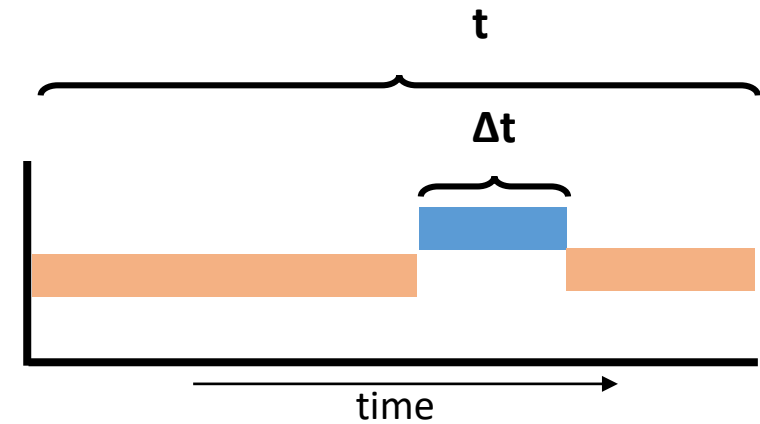
How is Isolation Broken?



- Kernel processes network traffic via interrupts
- Time spend in servicing interrupts is incorrectly charged

Charging: Reduction in the runtime of a container

Scheduler **charges** wordcount for "**t**", instead of "**t - Δt** "



Scheduled task/process

How did Linux get here?

“... [software interrupts] a conglomerate of mostly unrelated jobs, which run in the context of a **randomly chosen victim** w/o the ability to put any control on them.”

--Thomas Gleixner (Linux developer)

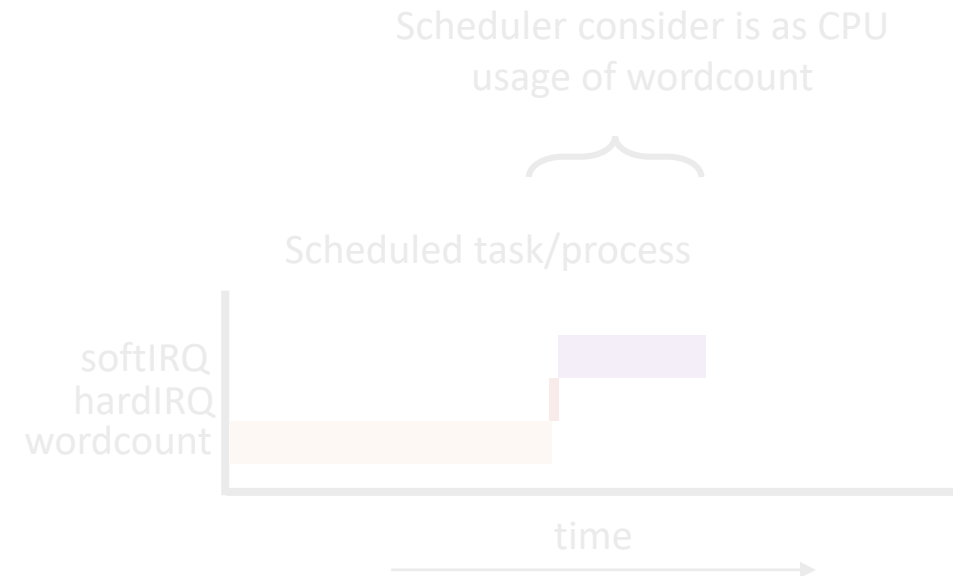
- Kernel processes network traffic via interrupts
- Time spend in servicing interrupts is not accounted properly

os



NIC ring buffer

NIC



How did Linux get here?

“... [software interrupts] a conglomerate of mostly unrelated jobs, which run in the context of a **randomly chosen victim** w/o the ability to put any control on them.”

--Thomas Gleixner (Linux developer)

- Kernel processes network traffic via interrupts
- Time spend in servicing interrupts is not accounted properly

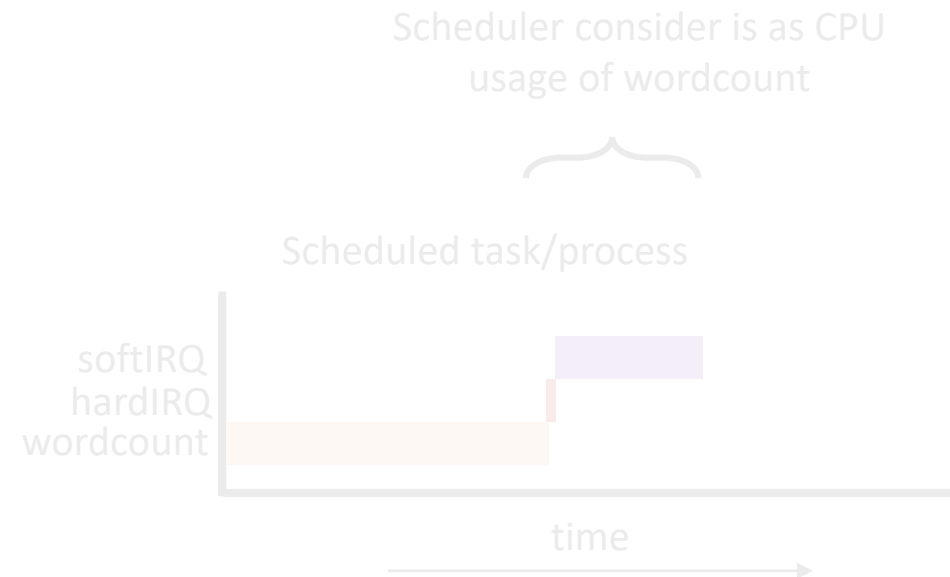
os



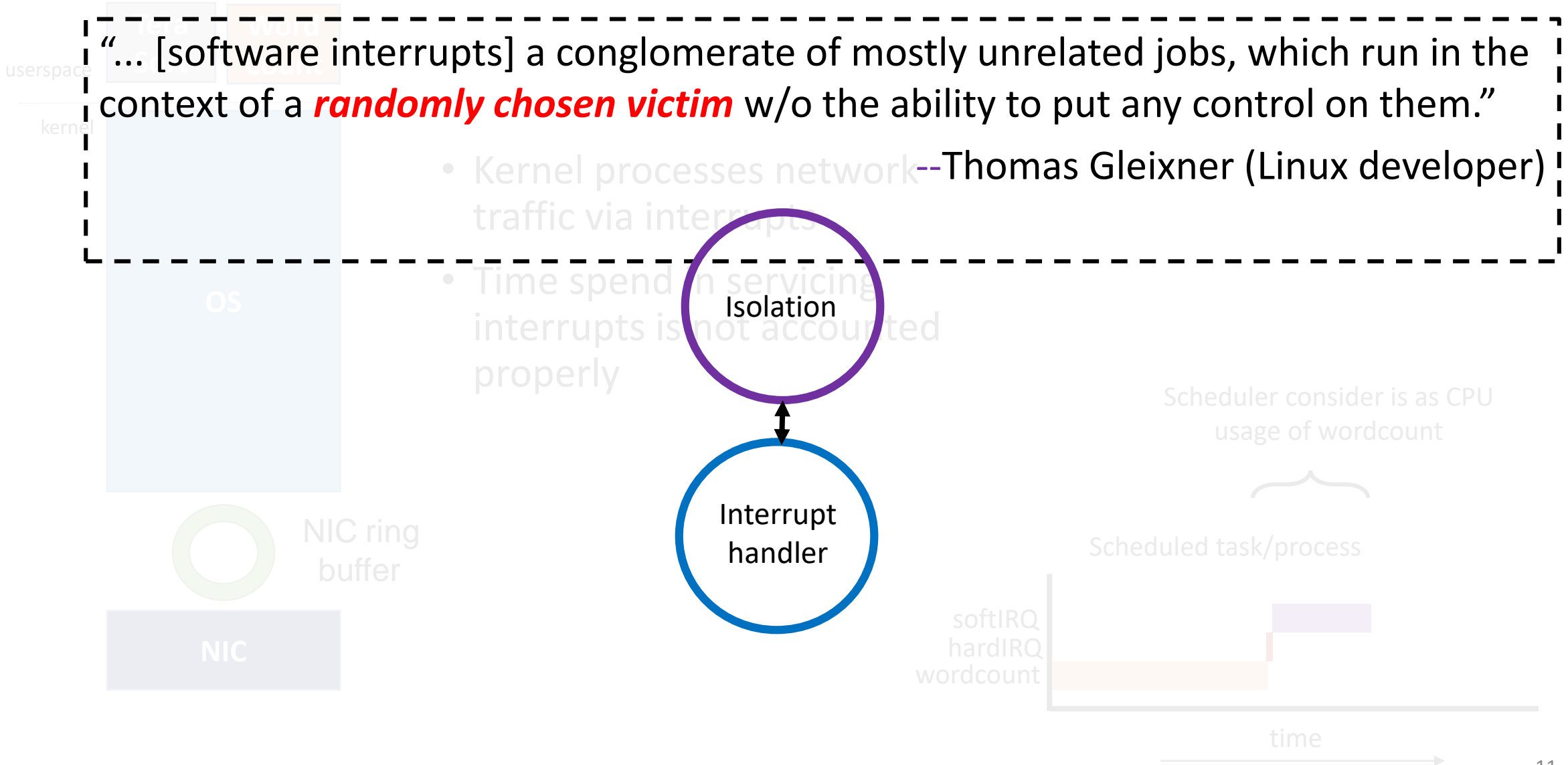
NIC ring buffer

NIC

Interrupt handler



How did Linux get here?

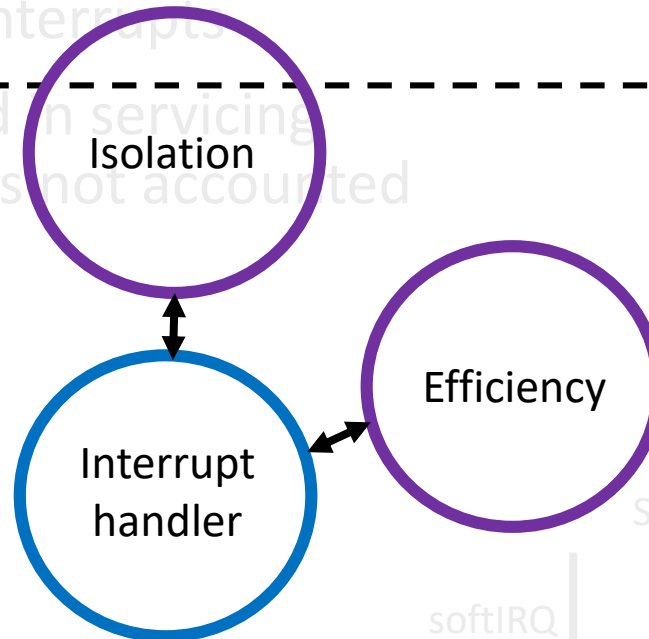


How did Linux get here?

“... [software interrupts] a conglomerate of mostly unrelated jobs, which run in the context of a **randomly chosen victim** w/o the ability to put any control on them.”

--Thomas Gleixner (Linux developer)

- Kernel processes network traffic via interrupts
- Time spend in servicing interrupts is not accounted properly



Scheduler consider is as CPU usage of wordcount

Scheduled task/process

softIRQ
hardIRQ
wordcount

time

How did Linux get here?

“... [software interrupts] a conglomerate of mostly unrelated jobs, which run in the context of a **randomly chosen victim** w/o the ability to put any control on them.”

--Thomas Gleixner (Linux developer)

- Kernel processes network traffic via interrupts

- Time spend in servicing interrupts is not accounted properly

Isolation

Efficiency

Interrupt handler

Responsiveness

Scheduler consider is as CPU usage of wordcount

Scheduled task/process

SOFTIRQ handling wordcount

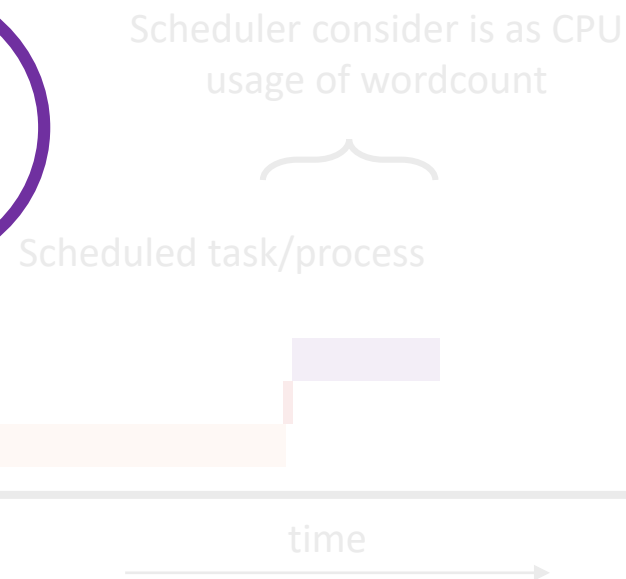
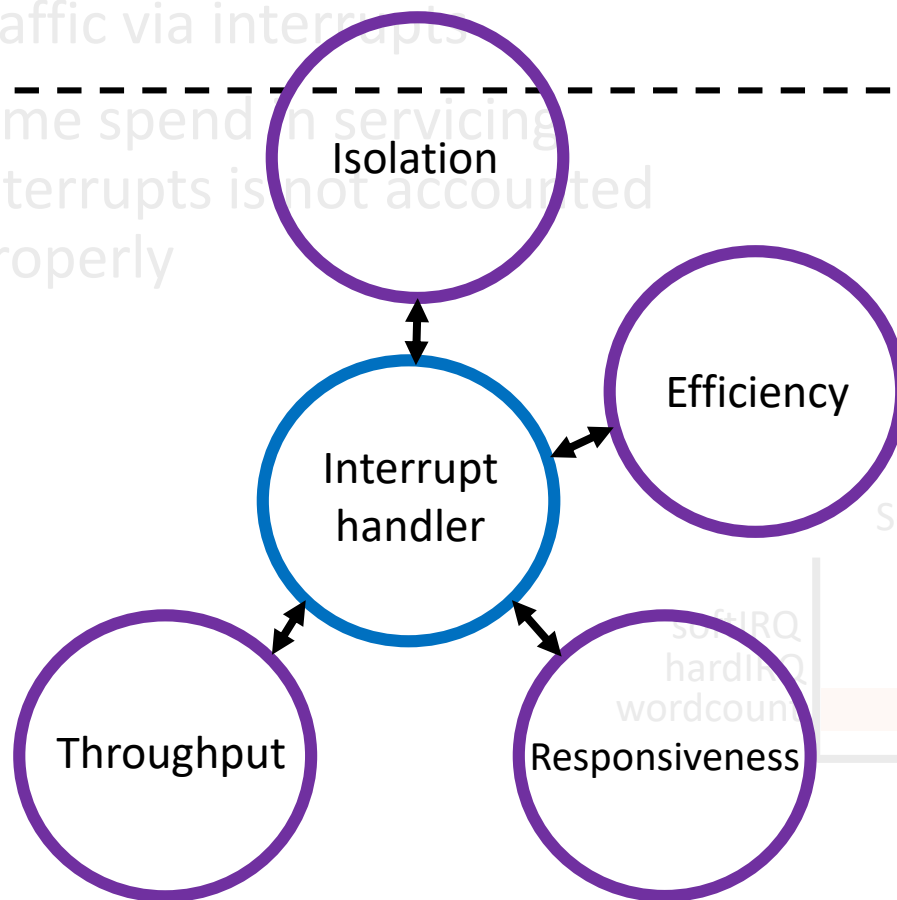
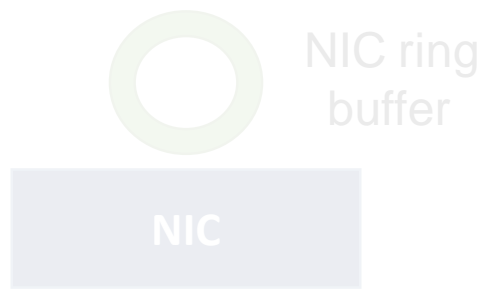
time

How did Linux get here?

“... [software interrupts] a conglomerate of mostly unrelated jobs, which run in the context of a **randomly chosen victim** w/o the ability to put any control on them.”

--Thomas Gleixner (Linux developer)

- Kernel processes network traffic via interrupts
- Time spend in servicing interrupts is not accounted properly

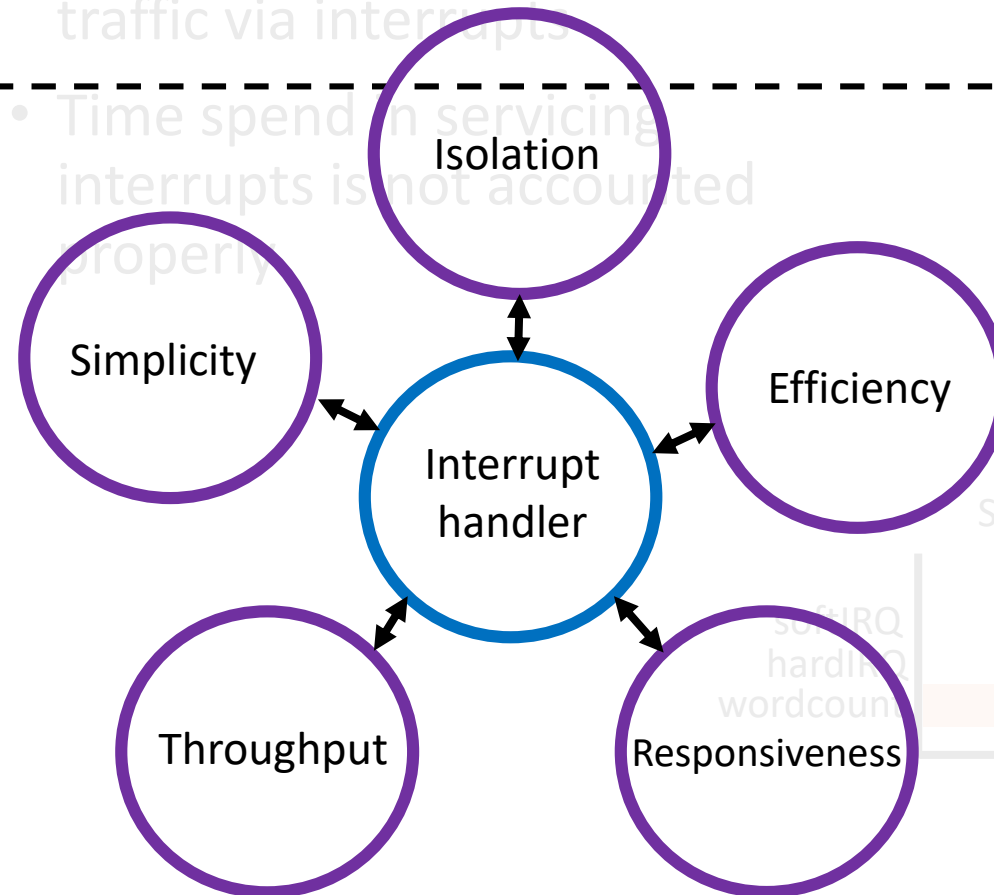


How did Linux get here?

“... [software interrupts] a conglomerate of mostly unrelated jobs, which run in the context of a **randomly chosen victim** w/o the ability to put any control on them.”

--Thomas Gleixner (Linux developer)

- Kernel processes network traffic via interrupts
- Time spend in servicing interrupts is not accounted properly



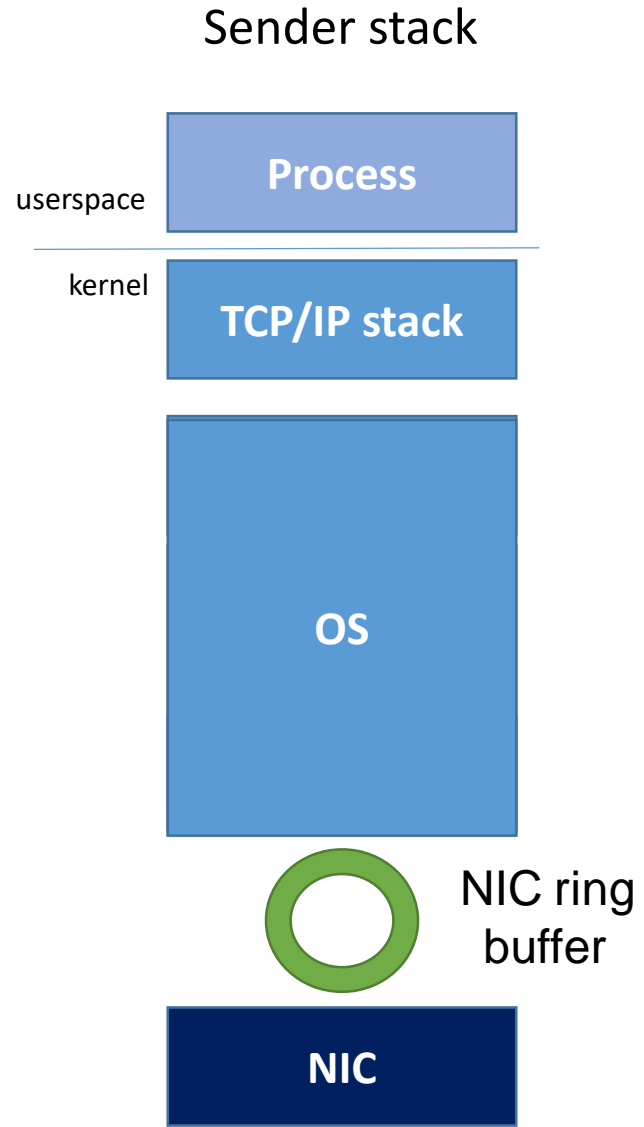
Scheduler consider is as CPU usage of wordcount

Scheduled task/process

SO_IRQ handling wordcount

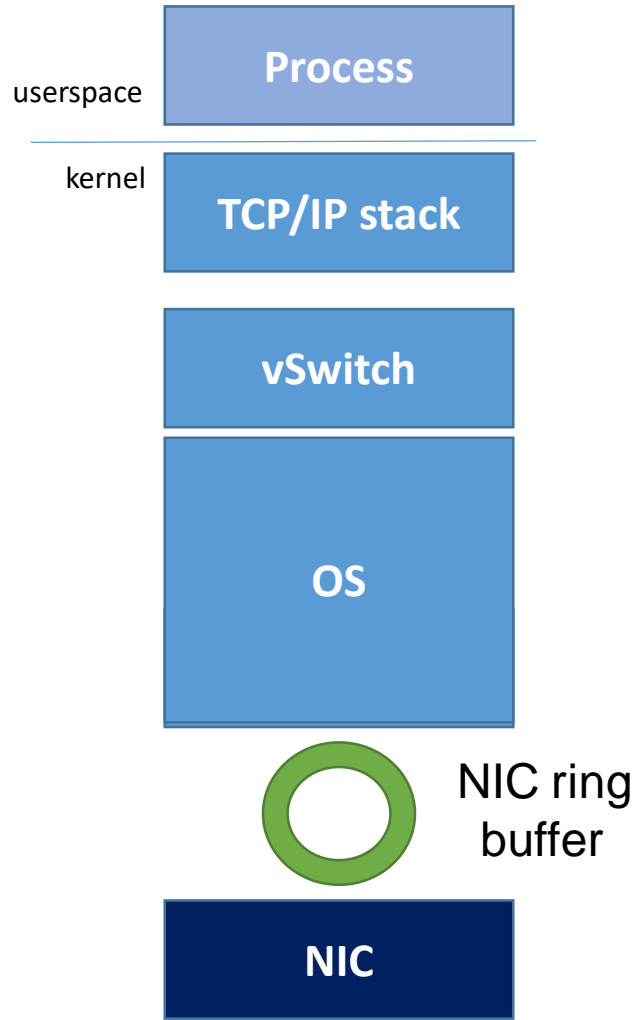
time

Sender Side



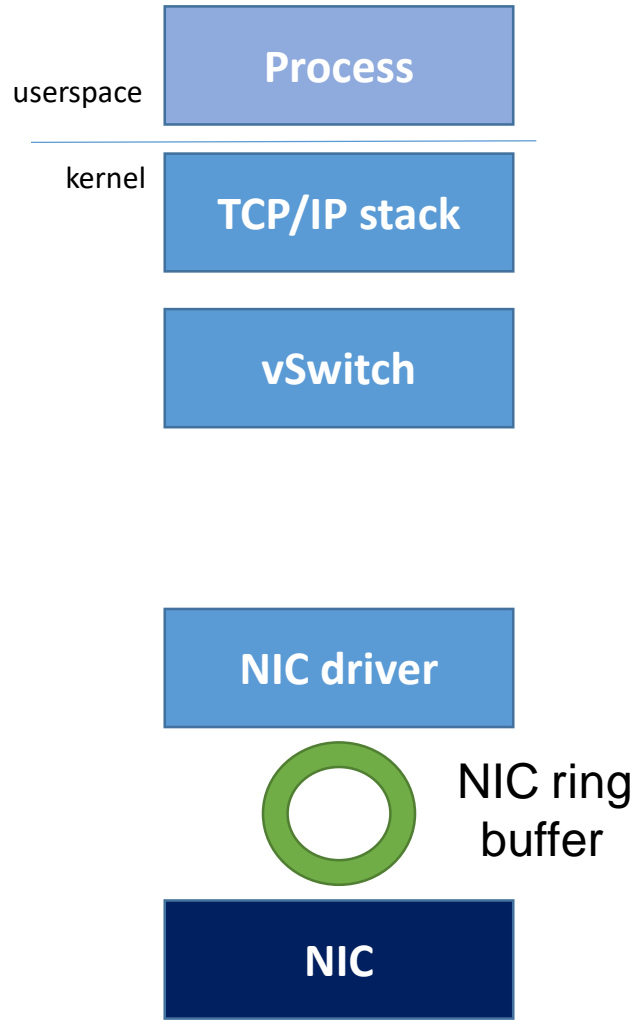
Sender Side

Sender stack



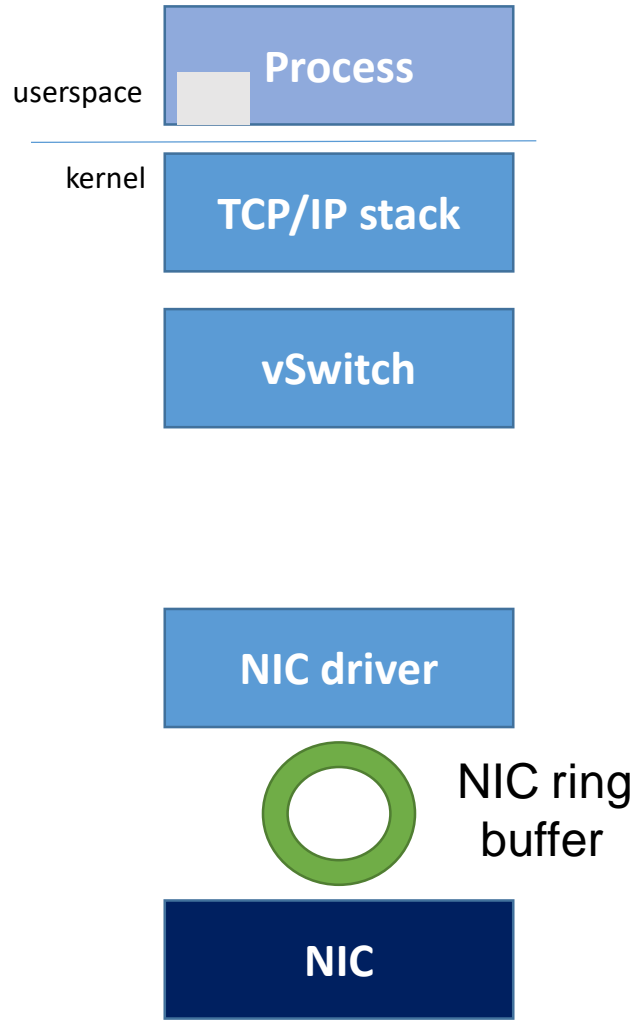
Sender Side

Sender stack



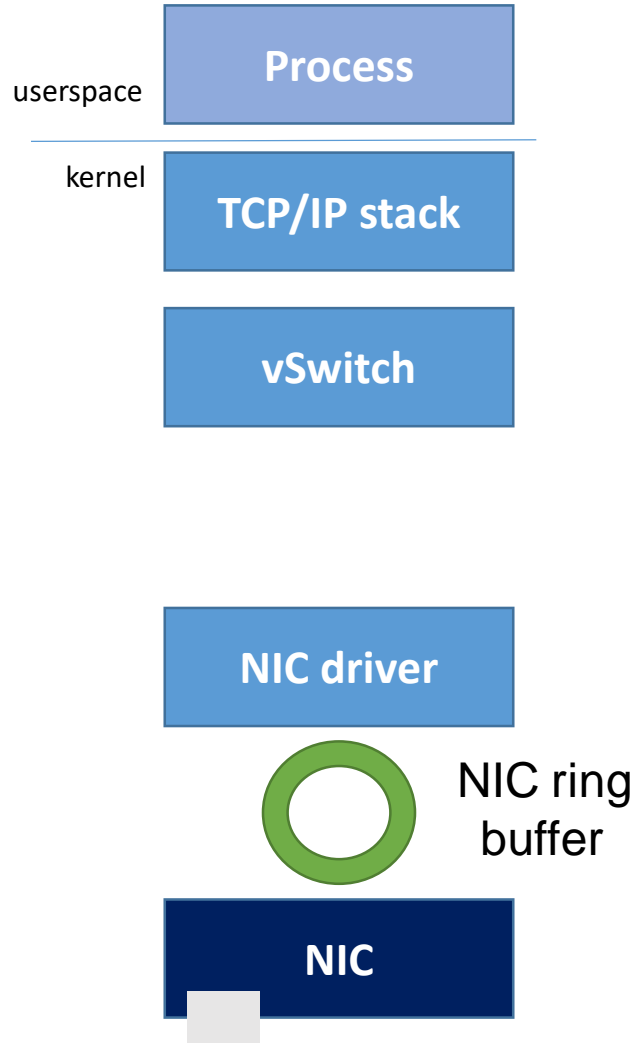
Sender Side

Sender stack



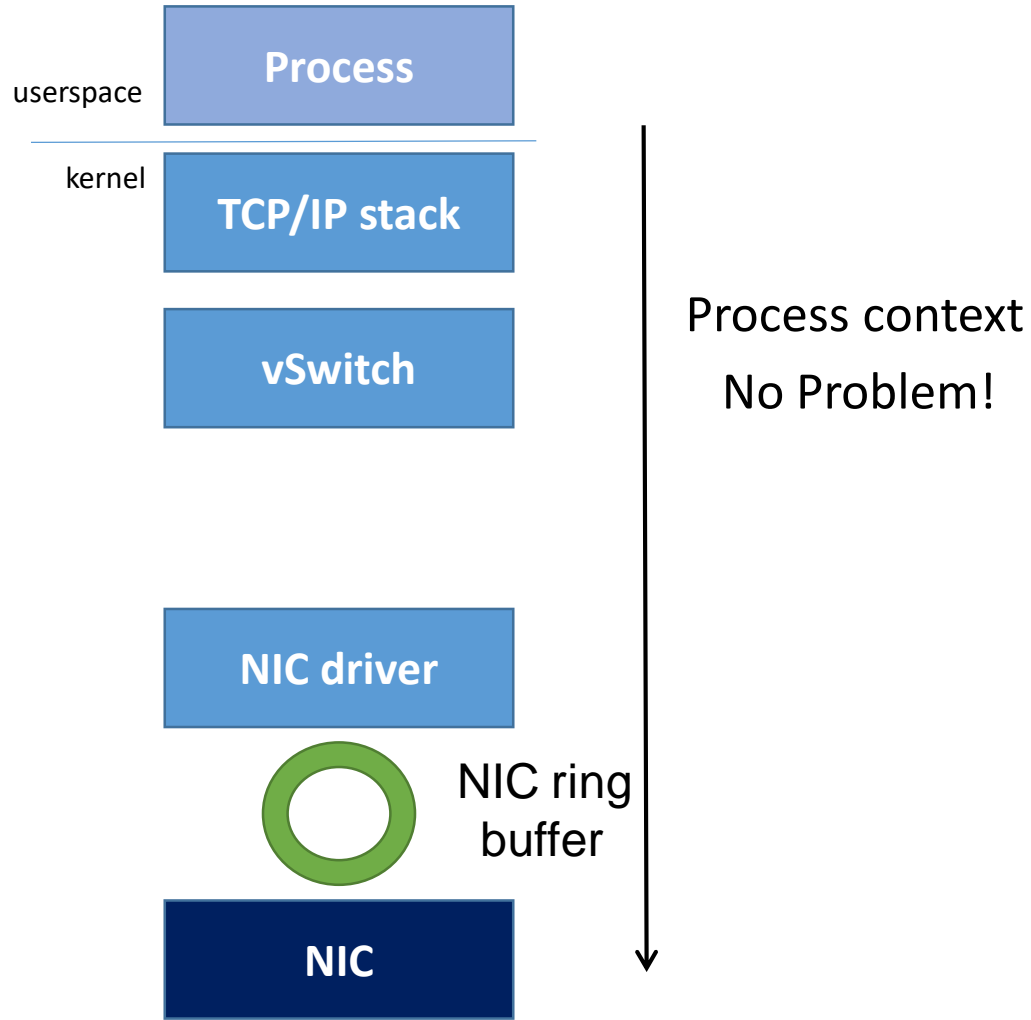
Sender Side

Sender stack

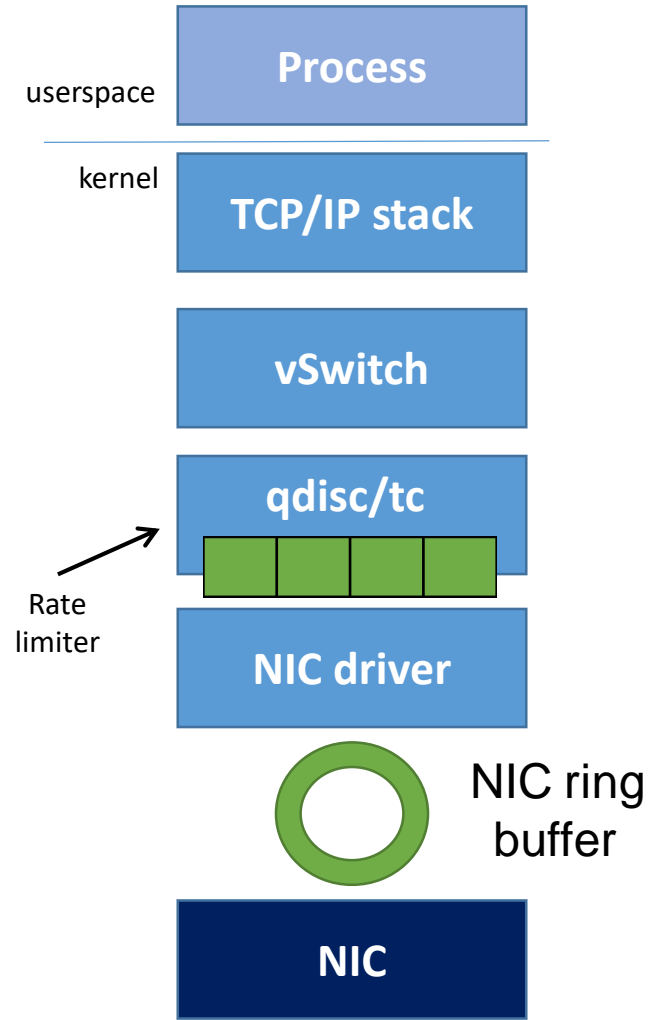


Sender Side

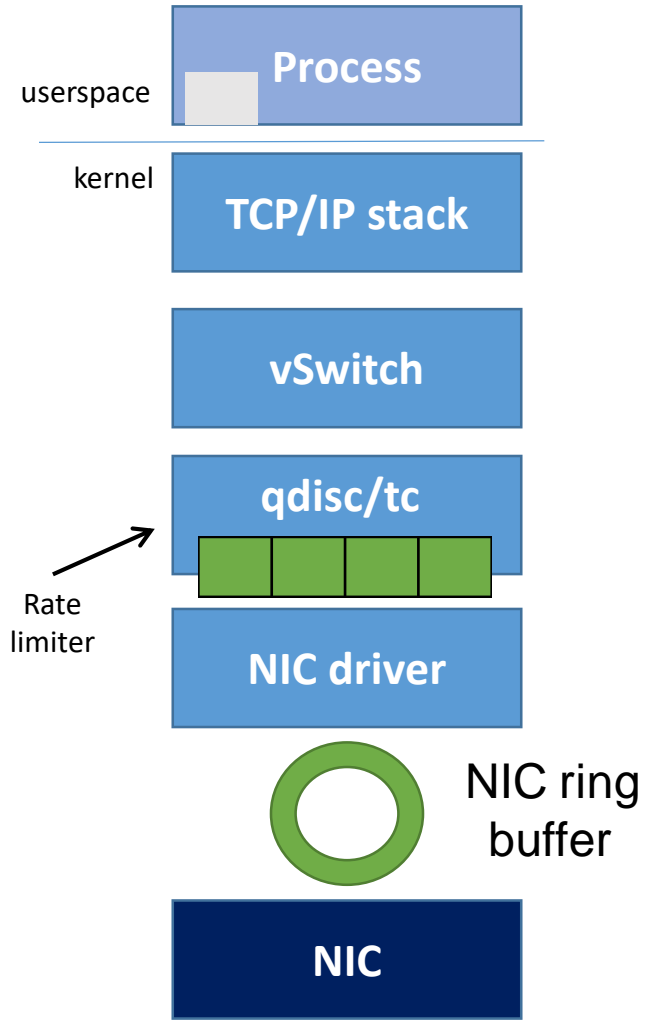
Sender stack



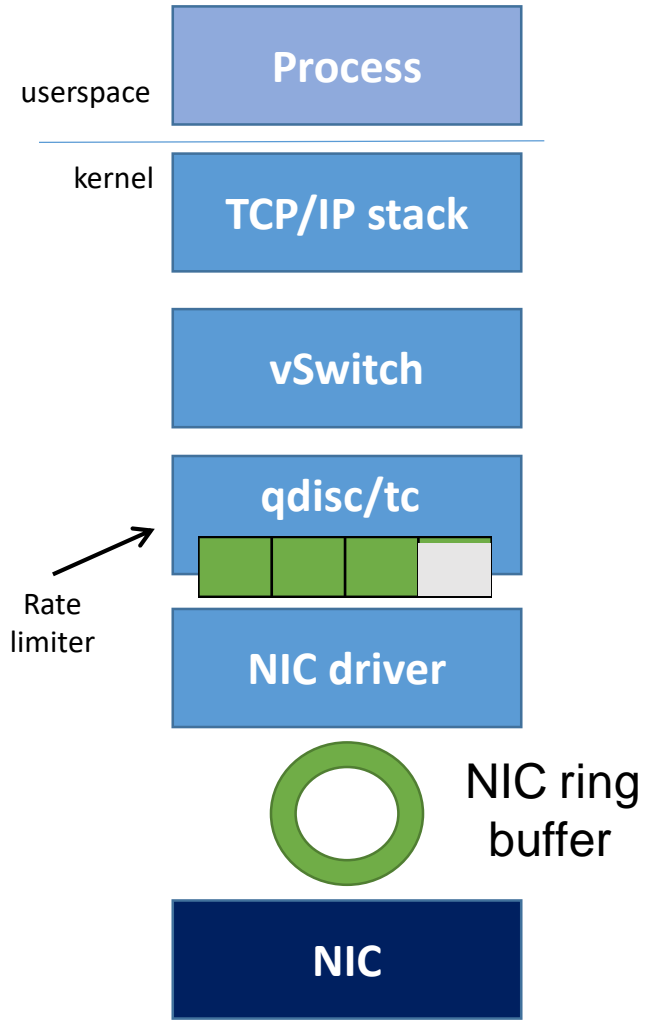
Sender Side



Sender Side



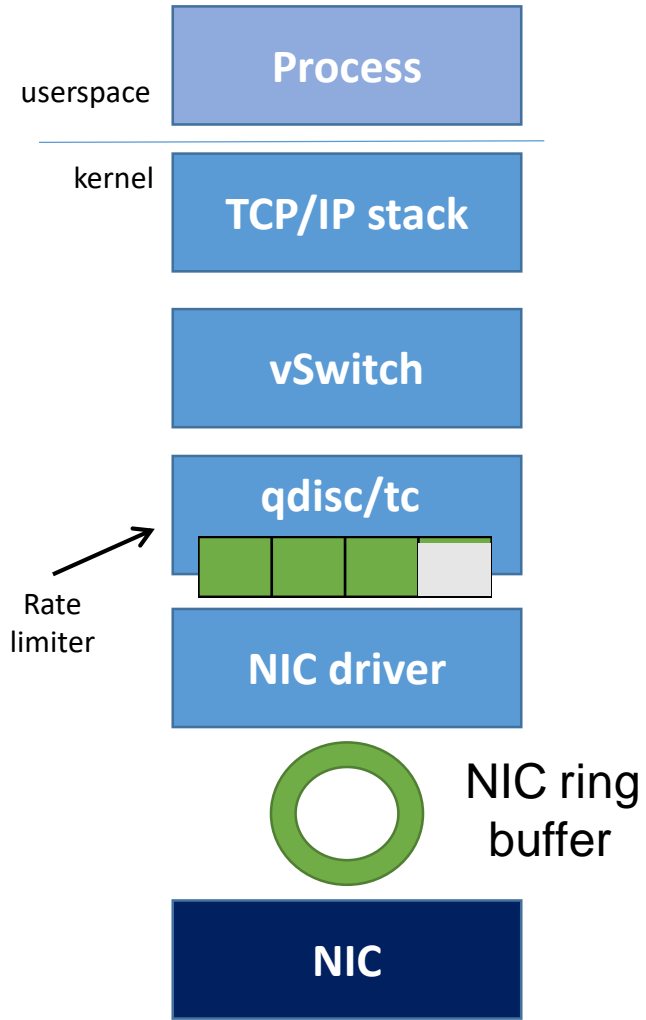
Sender Side



Process context
No Problem!

- Packet is enqueued in the process context

Sender Side

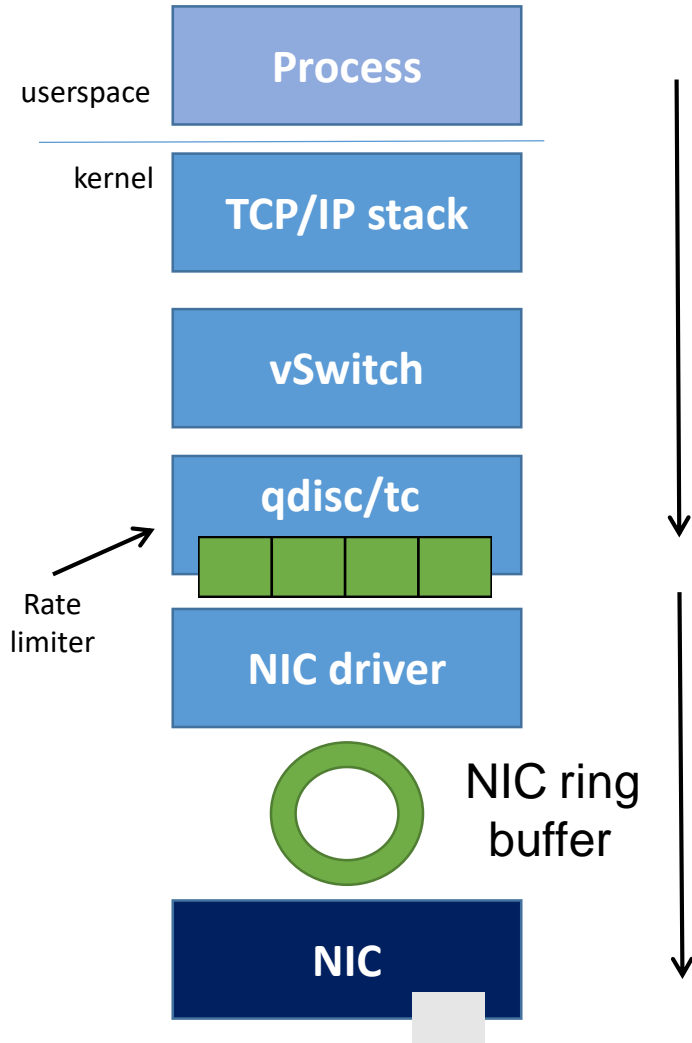


Process context
No Problem!

- Packet is enqueued in the process context
- System call exits after enqueueing the packet

Sender Side

Sender stack



Process context

No Problem!

Non-process
context

Process may not
get charged

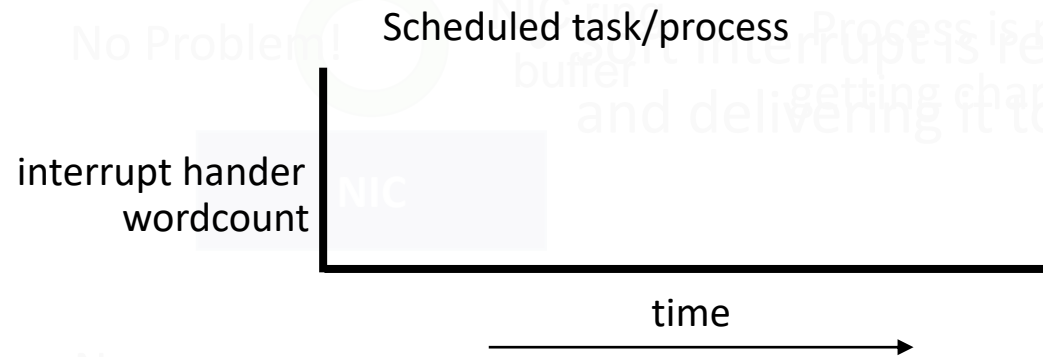
- Packet is enqueued in the process context
- System call exits after enqueueing the packet
- Soft interrupt is responsible for dequeuing and delivering it to the NIC

Sender Side

Sender side stack

Linux services a softirq

- 1) at the end of hardware interrupt processing, in the context of the *currently scheduled process*

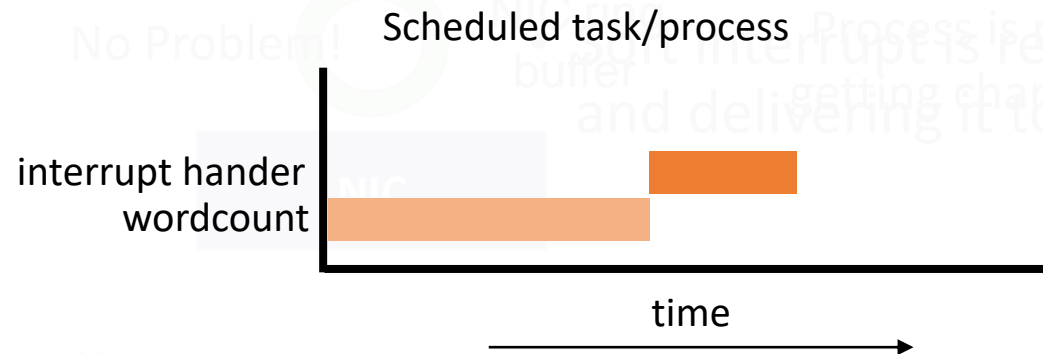


Sender Side

Sender side stack

Linux services a softirq

- 1) at the end of hardware interrupt processing, in the context of the *currently scheduled process*

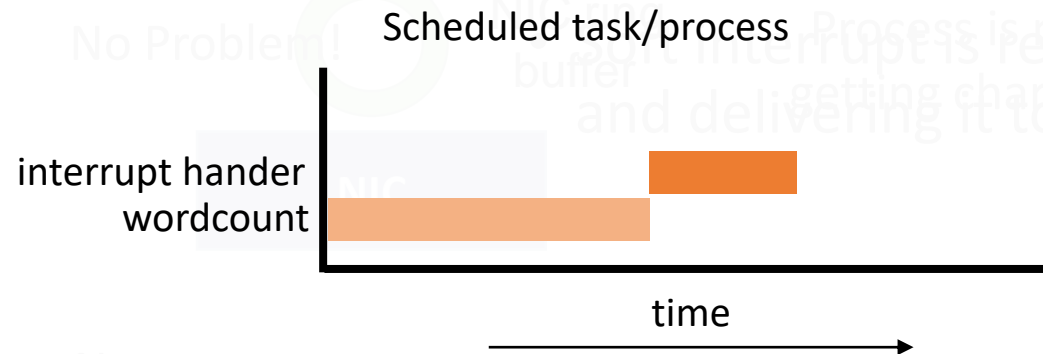


Sender Side

Sender side stack

Linux services a softirq

- 1) at the end of hardware interrupt processing, in the context of the *currently scheduled process*



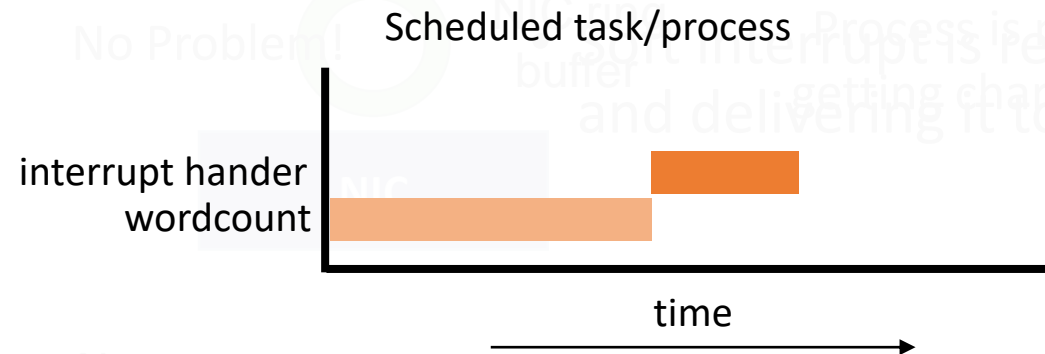
- 2) through *ksoftirqd thread* (a per core kernel thread)

Sender Side

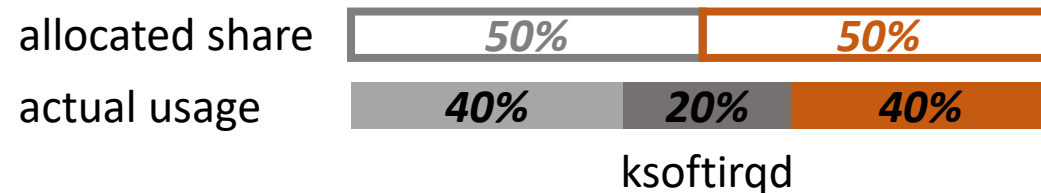
Sender side stack

Linux services a softirq

- 1) at the end of hardware interrupt processing, in the context of the *currently scheduled process*



- 2) through *ksoftirqd thread* (a per core kernel thread)



Sender Side

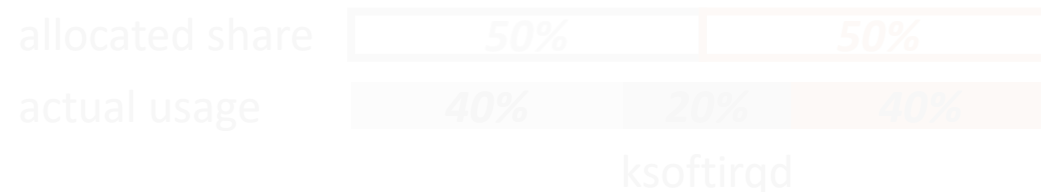
Sender side stack

Linux services the softirq

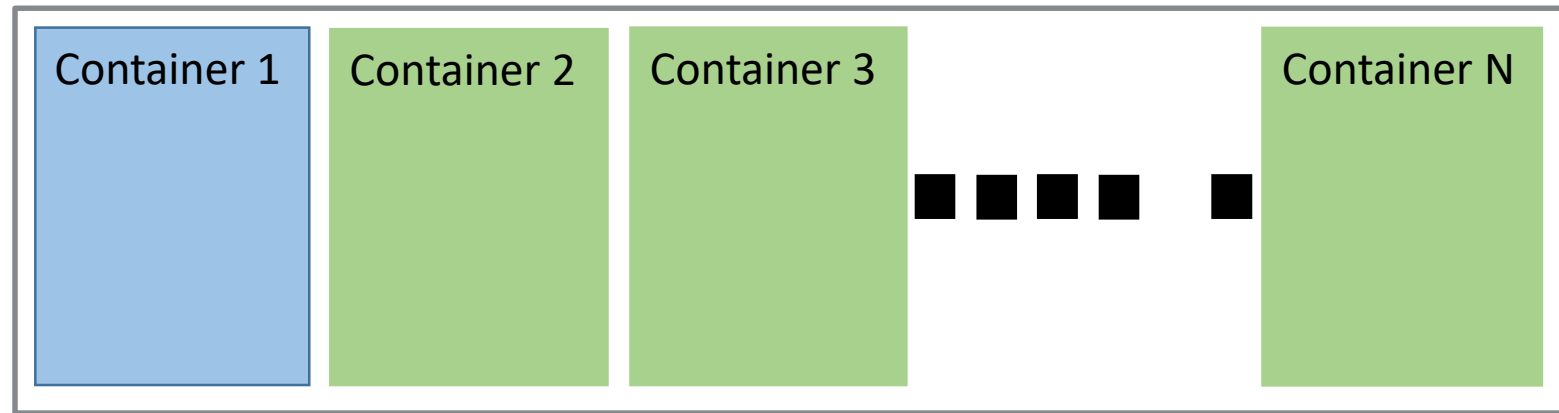
- 1) at the end of hardware interrupt processing, in the context of *one of* the *current scheduled process*.

Softirq processing can be ***charged incorrectly*** or ***not charged*** at all to any container

- 2) through *ksoftirqd thread* (a per core kernel thread)

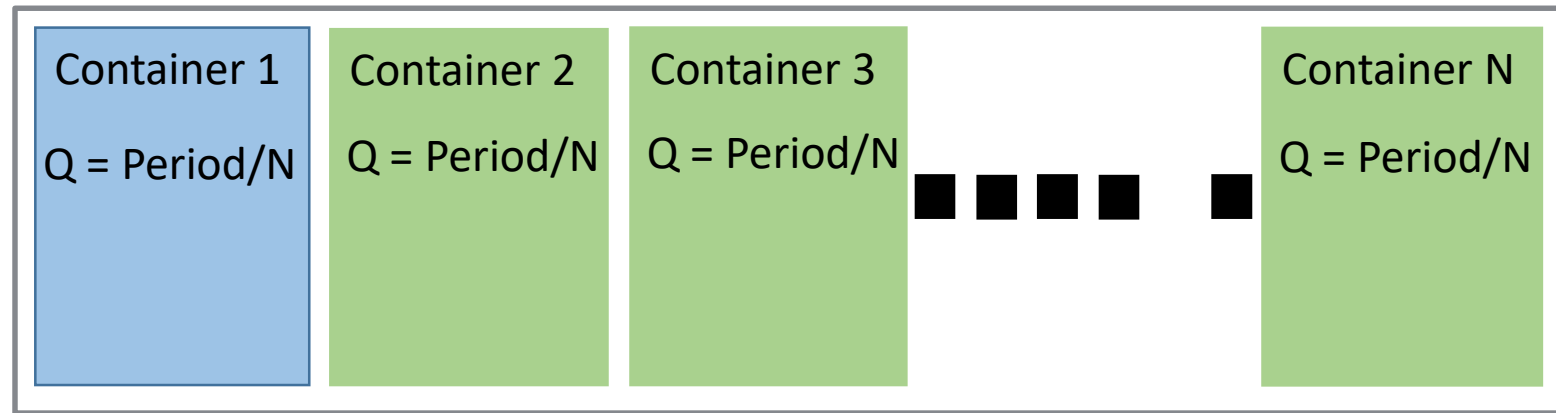


Experiment Setup



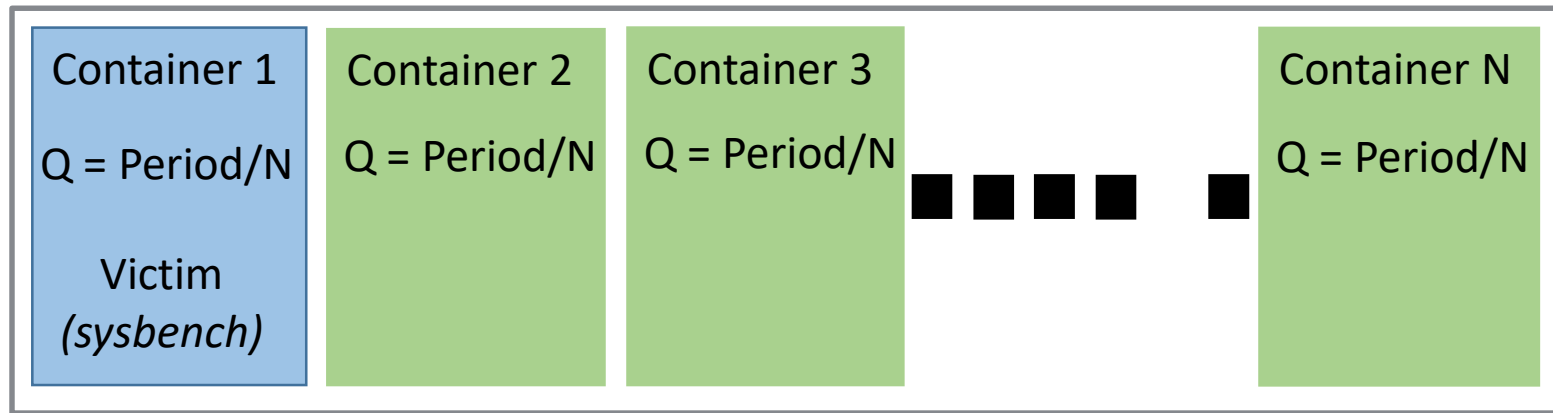
Core

Experiment Setup



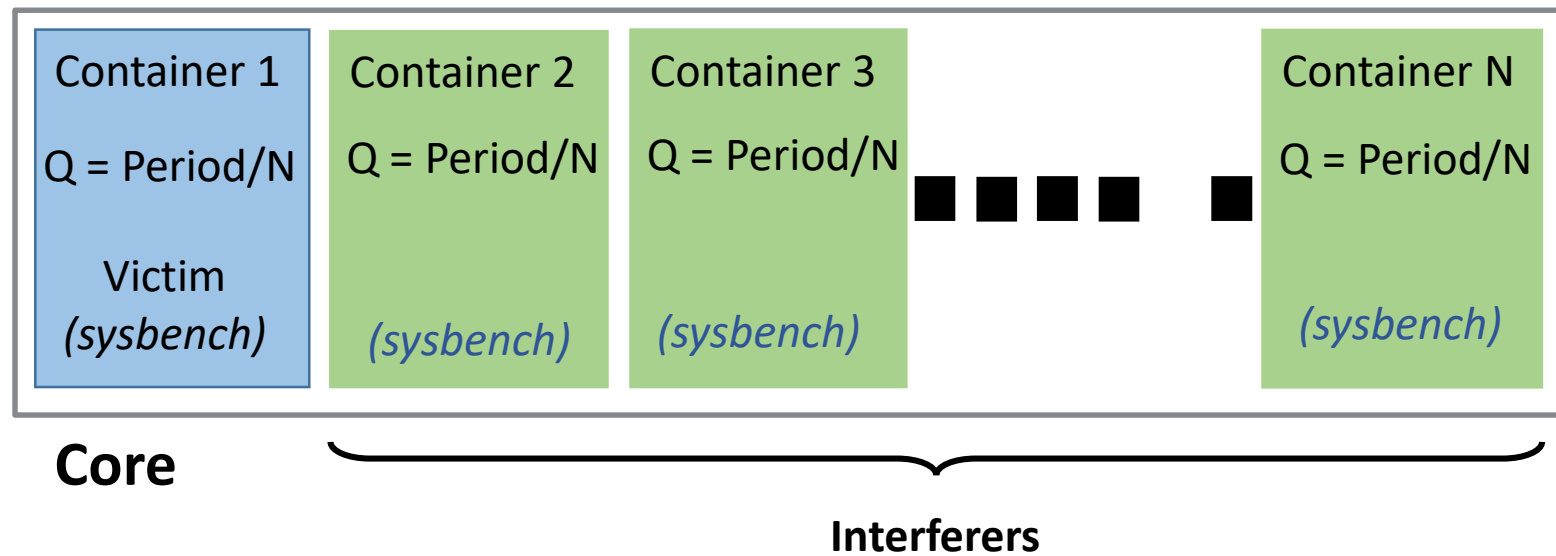
Core

Experiment Setup



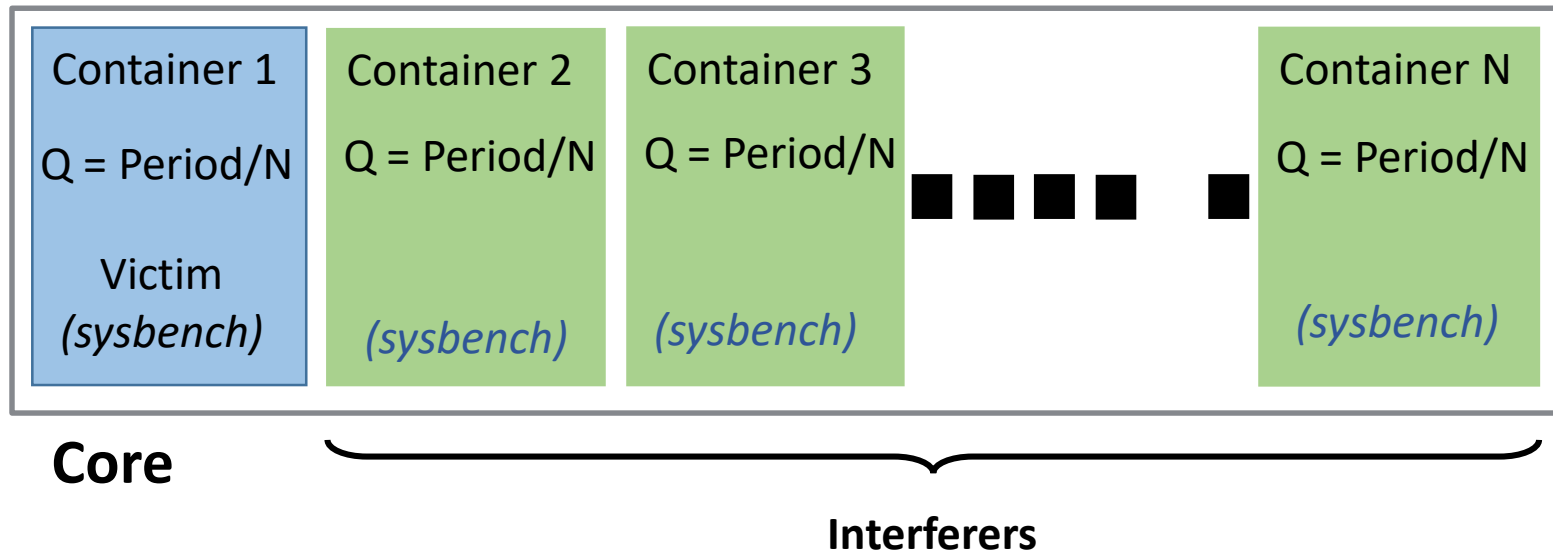
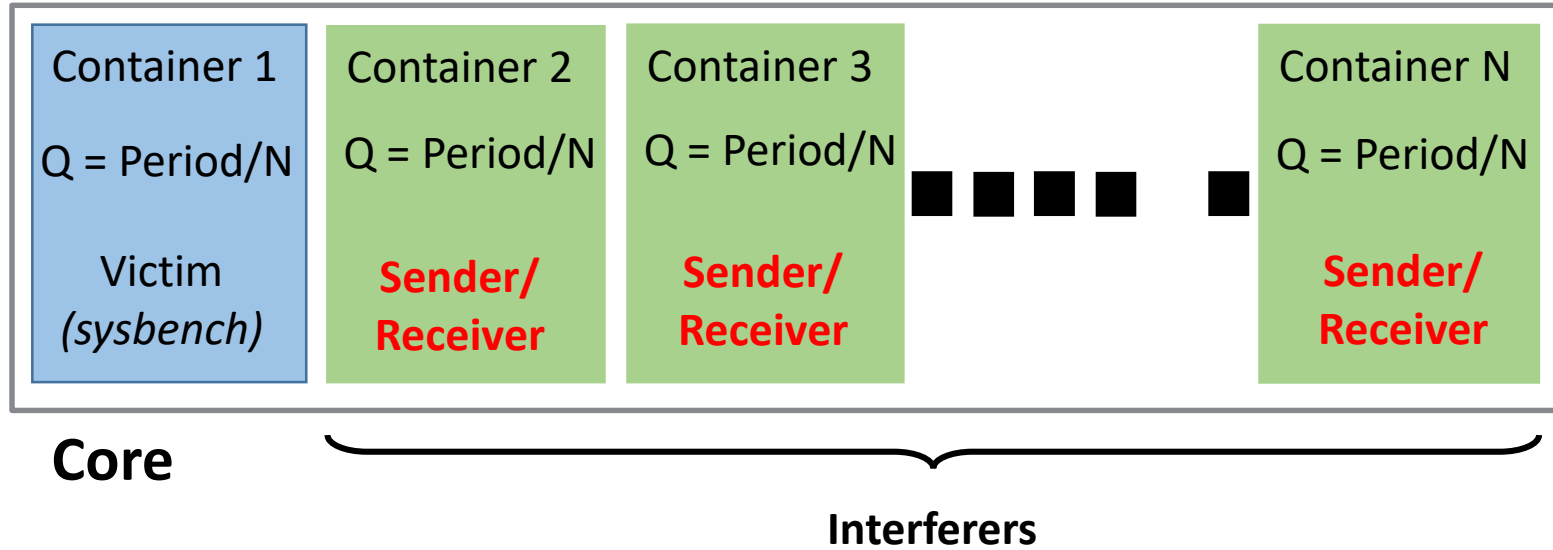
Core

Experiment Setup



Experiment Setup

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$

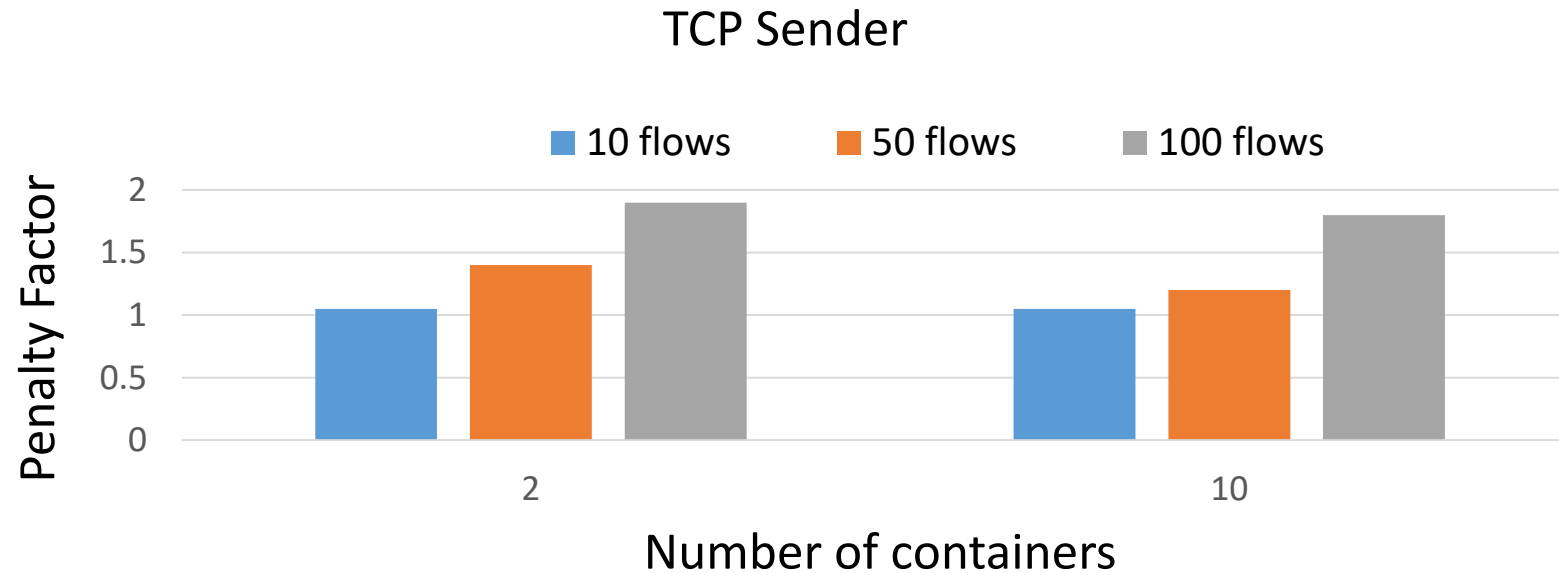


Impact Of Network Traffic

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$

Impact Of Network Traffic

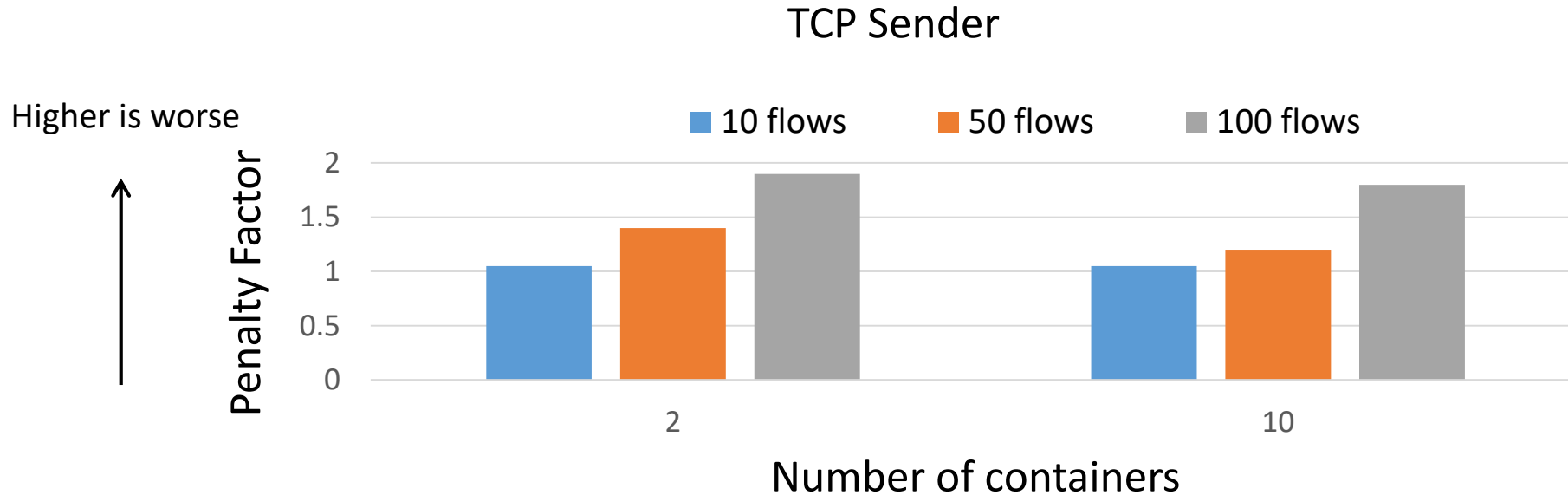
$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



HTB is used for traffic shaping @ 5Gbps

Impact Of Network Traffic

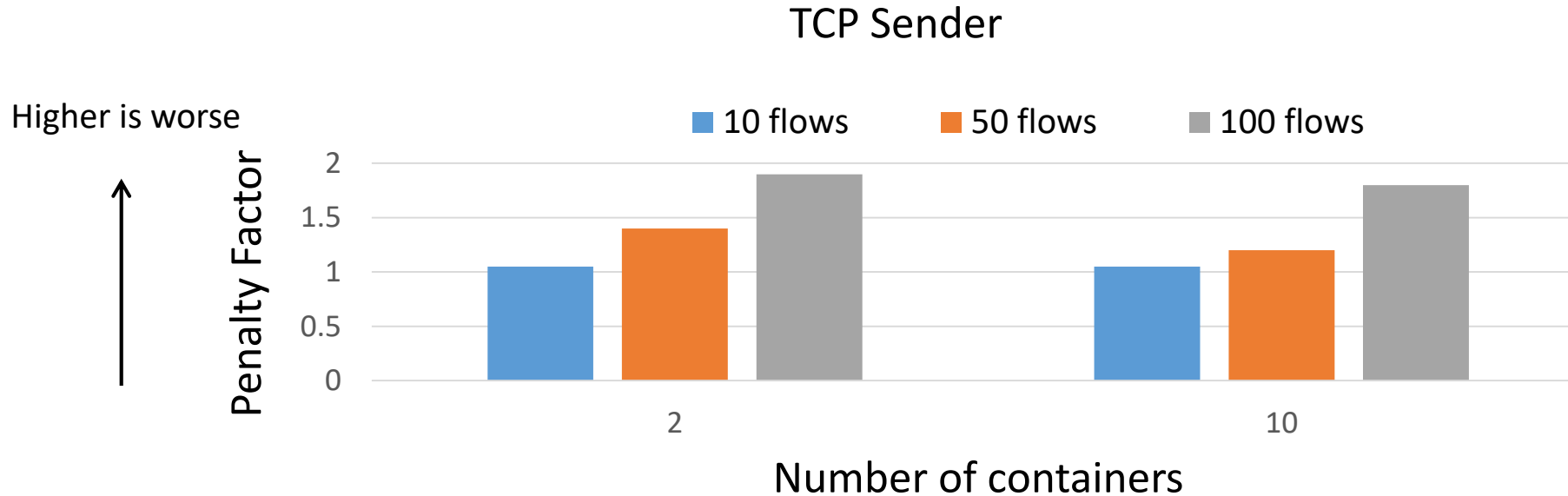
$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



HTB is used for traffic shaping @ 5Gbps

Impact Of Network Traffic

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$

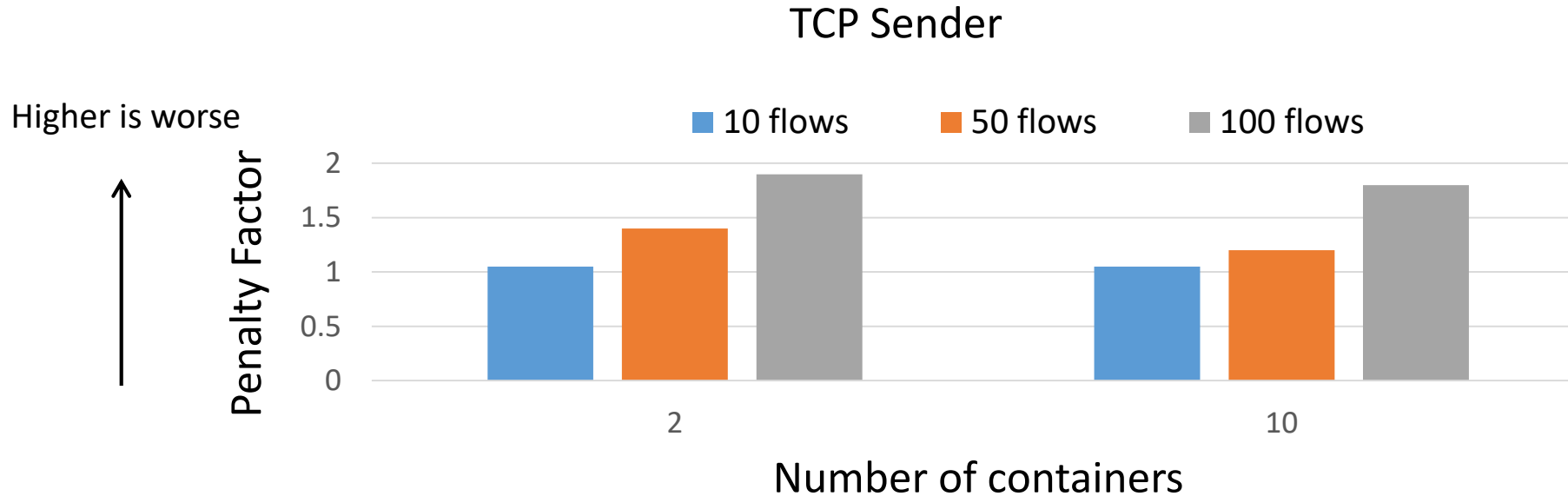


HTB is used for traffic shaping @ 5Gbps

Maximum penalty factor is around **1.85**

Impact Of Network Traffic

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



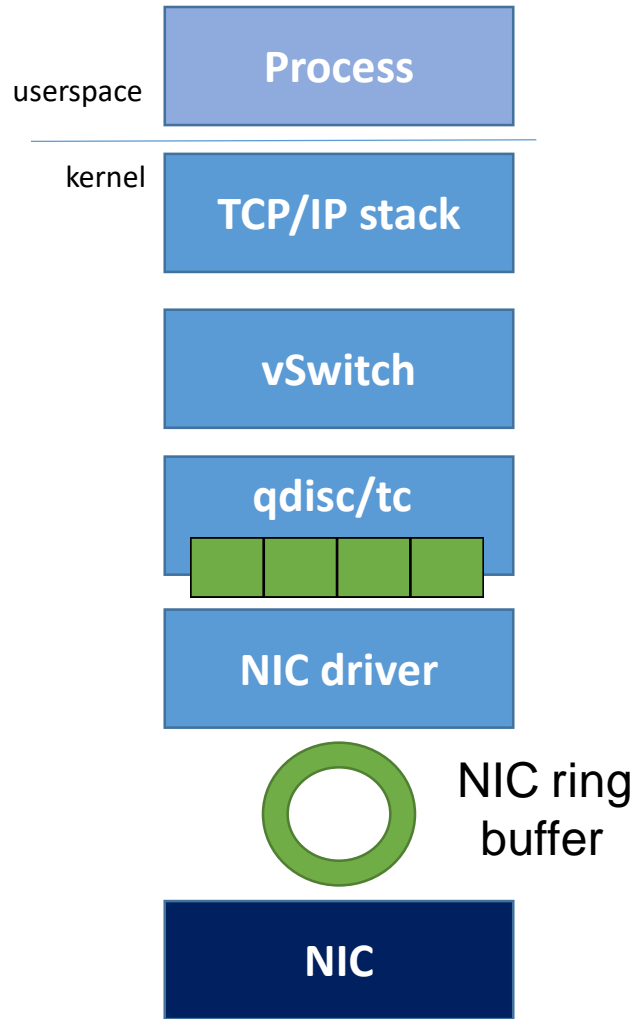
HTB is used for traffic shaping @ 5Gbps

Maximum penalty factor is around **1.85**

Look at our paper for the impact of UDP traffic

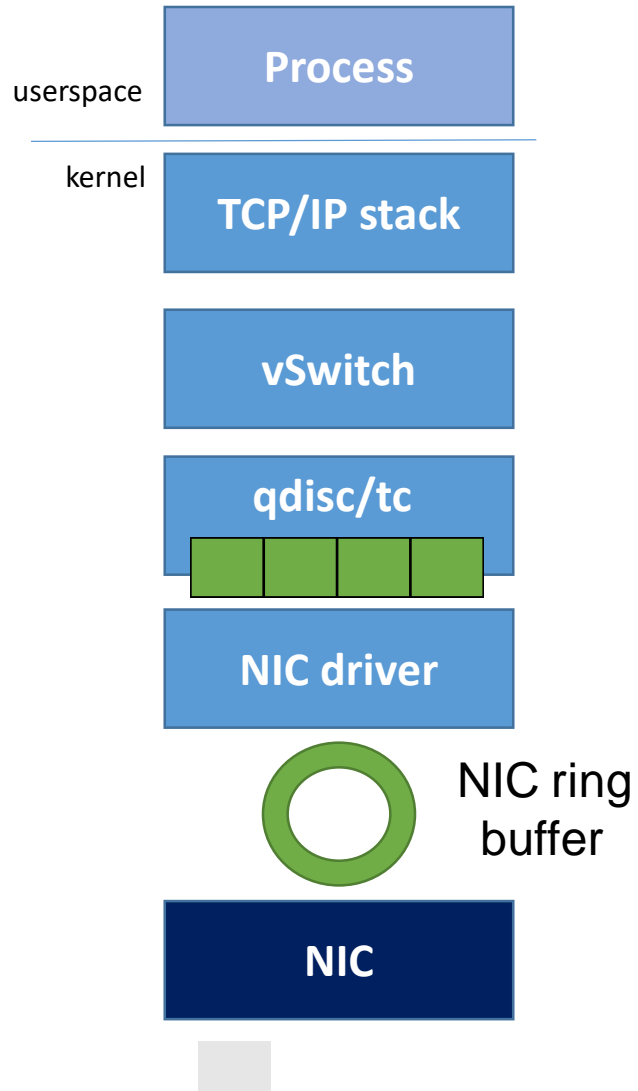
Receiver Side

Receiver stack



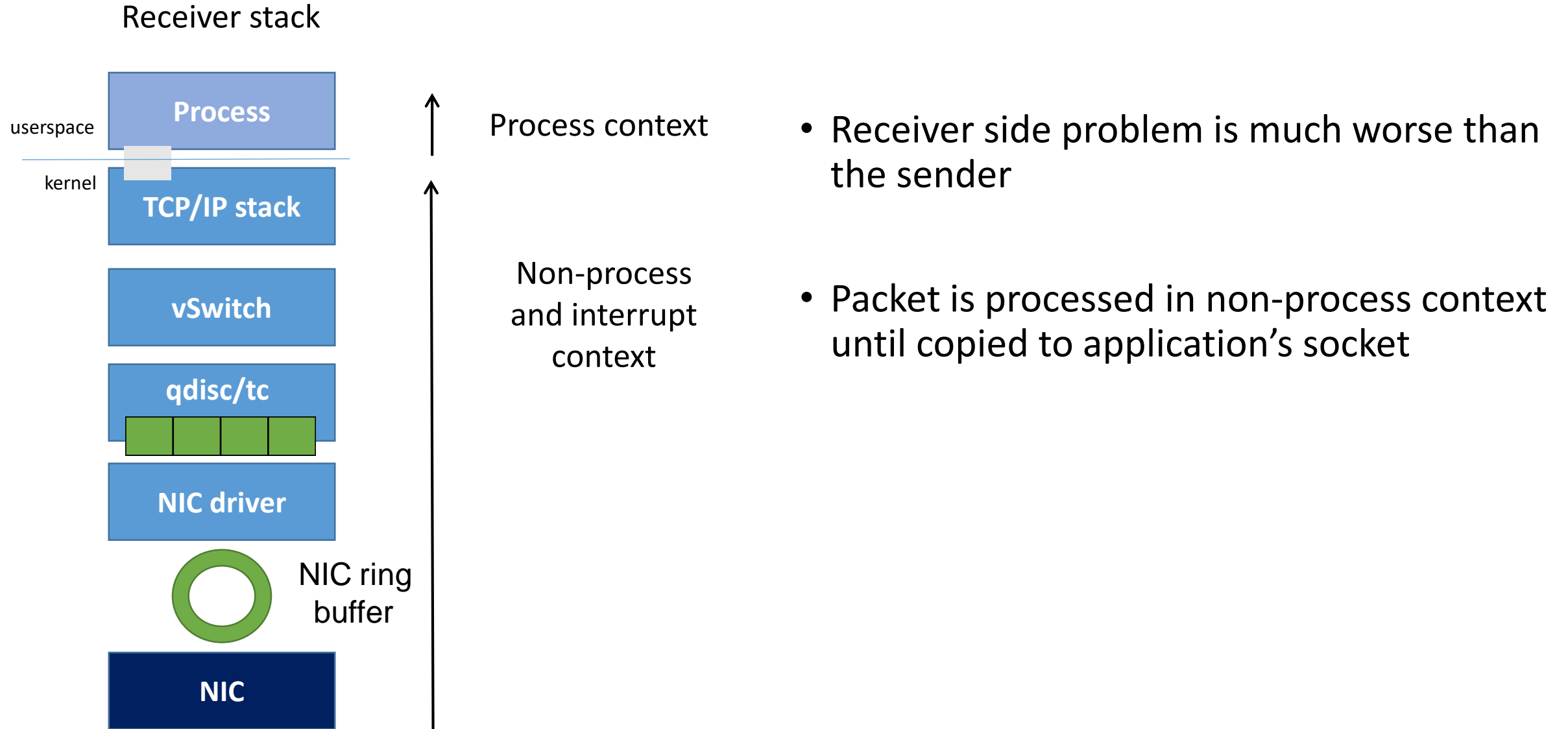
Receiver Side

Receiver stack



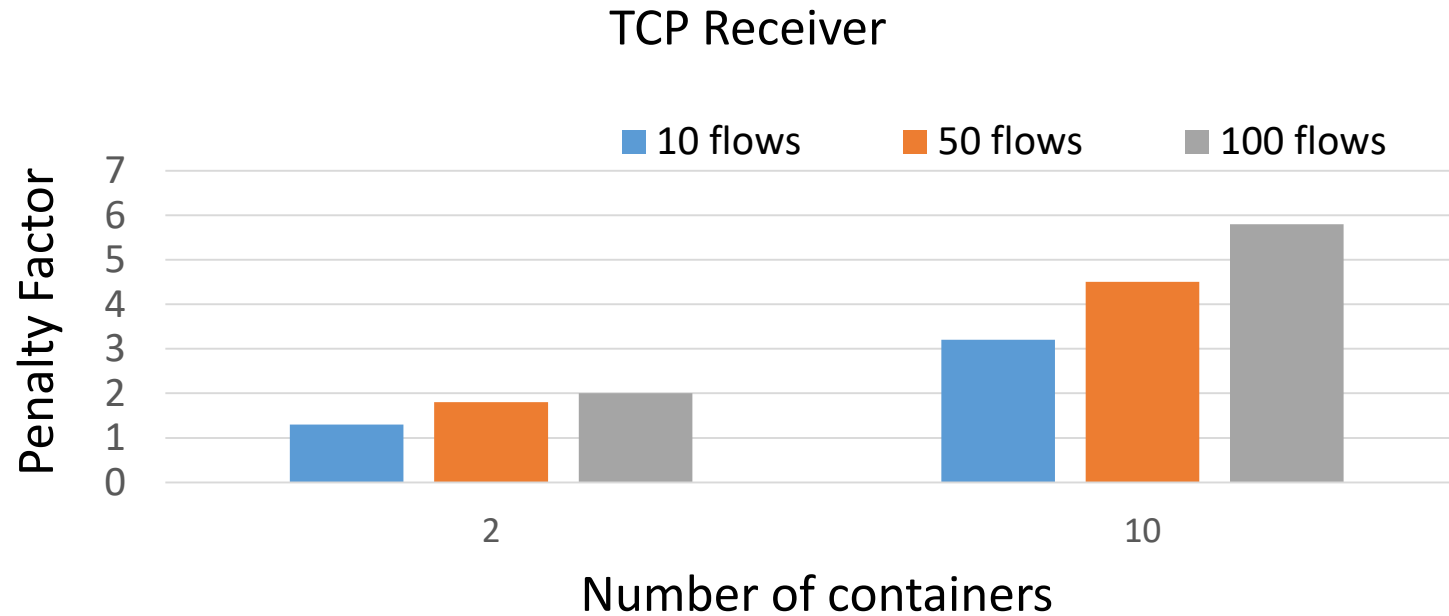
- Receiver side problem is much worse than the sender

Receiver Side



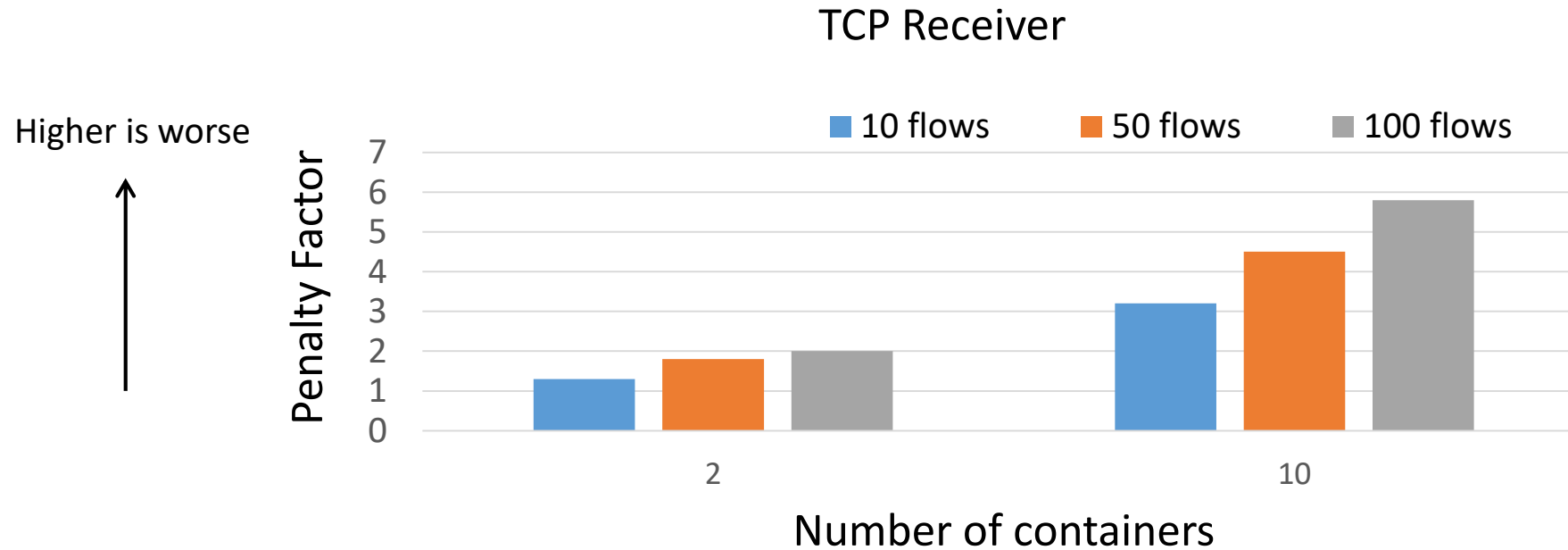
Impact Of Network Traffic

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



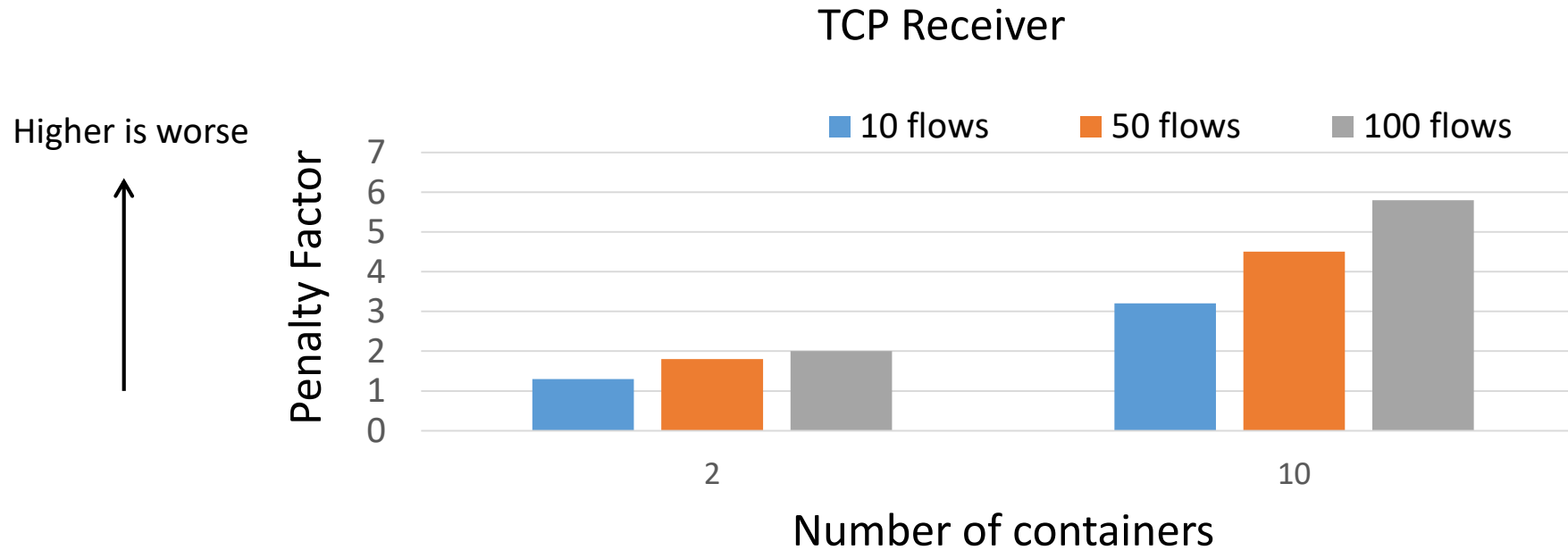
Impact Of Network Traffic

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



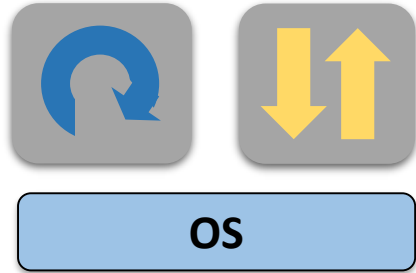
Impact Of Network Traffic

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$

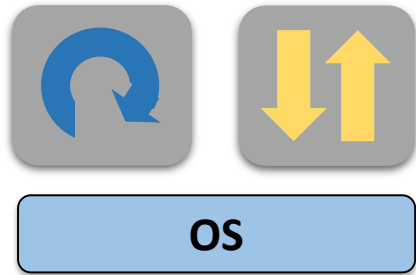


Maximum penalty factor is around **6**

Scenarios When Isolation Breaks



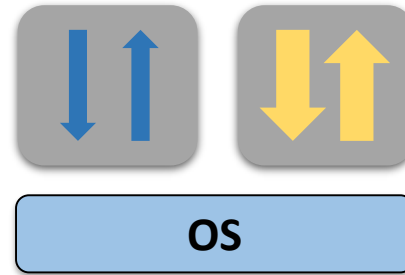
Scenarios When Isolation Breaks



Compute intensive

vs

Network intensive

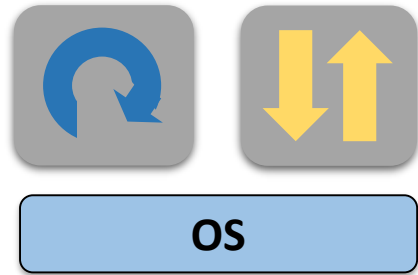


Low network workload

vs

High network workload

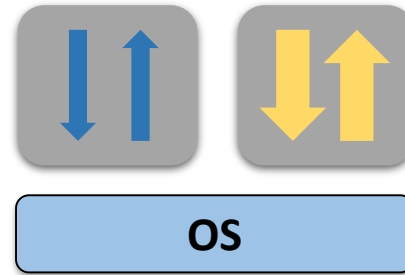
Scenarios When Isolation Breaks



Compute intensive

vs

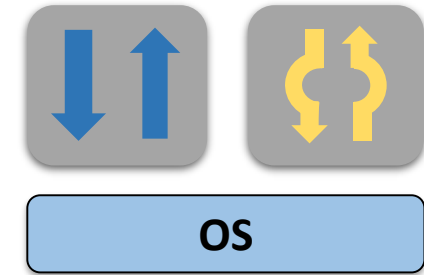
Network intensive



Low network workload

vs

High network workload



Network intensive

vs

Network intensive with
kernel bypass

Iron

A scheme that ensures and enforces accounting of *network-based CPU* consumed in the kernel on the *behalf of a container*.

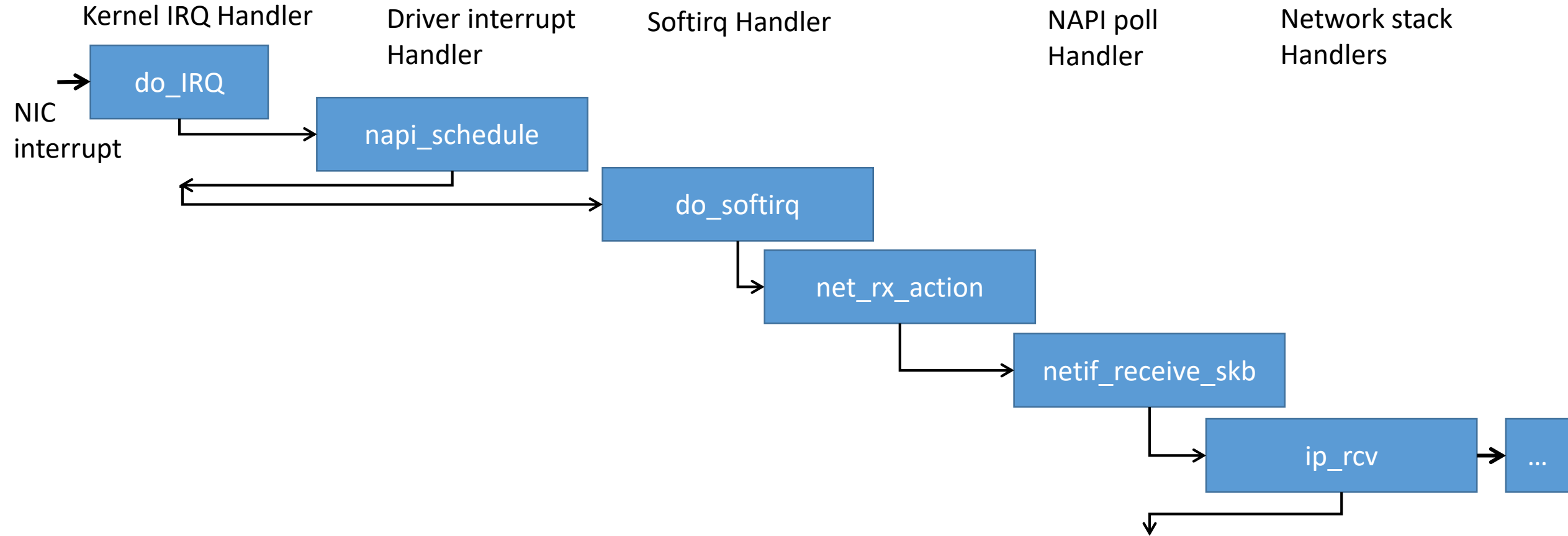
Outline

- How and by how much is isolation broken
- Iron's design
 - **Accounting of per-packet processing cost**
 - Ensuring isolation via enforcement
 - Integration with Linux scheduler
 - Hardware-based packet dropping
- Evaluation
 - Controlled workload
 - Realistic workload

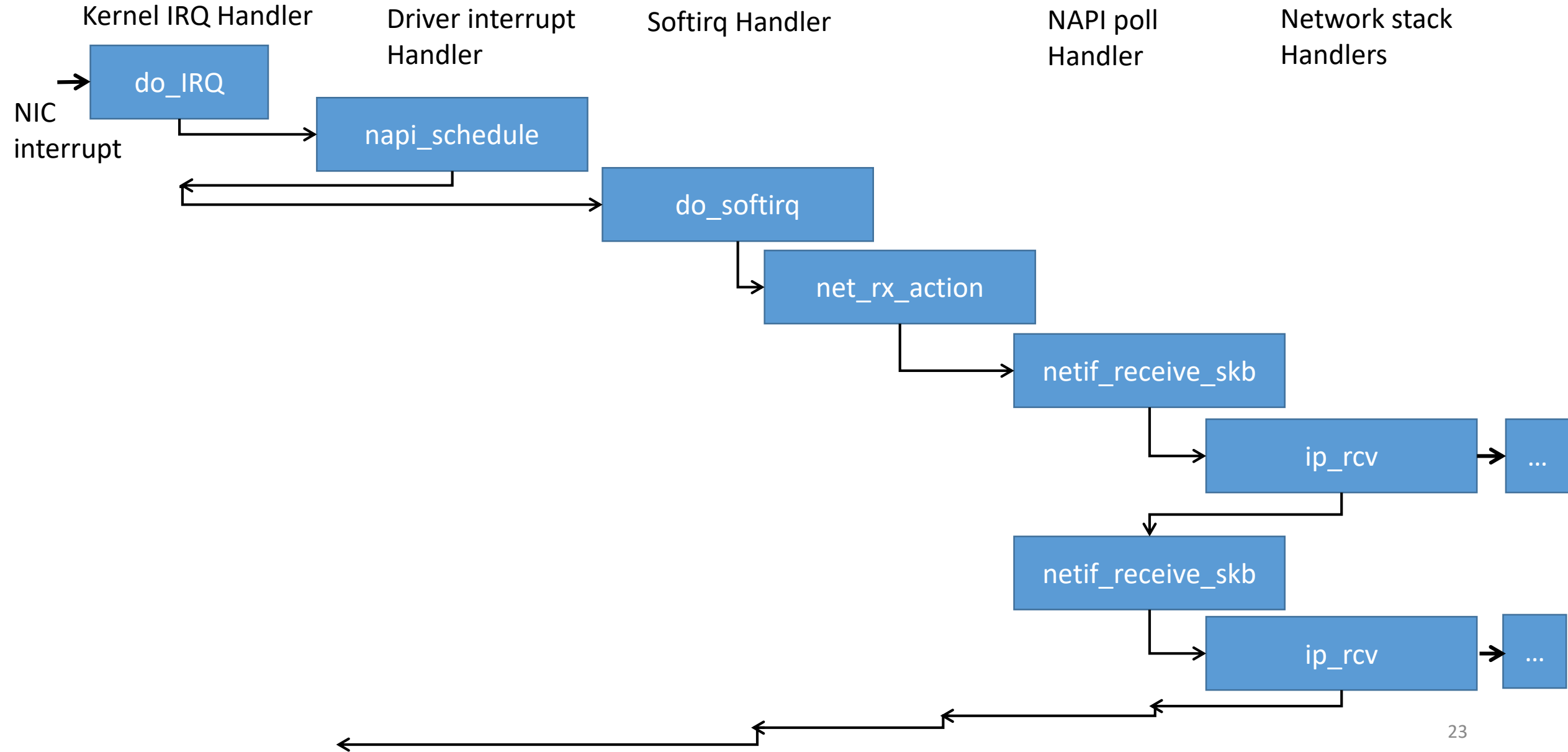
Network Call Stack – Background

NIC
interrupt

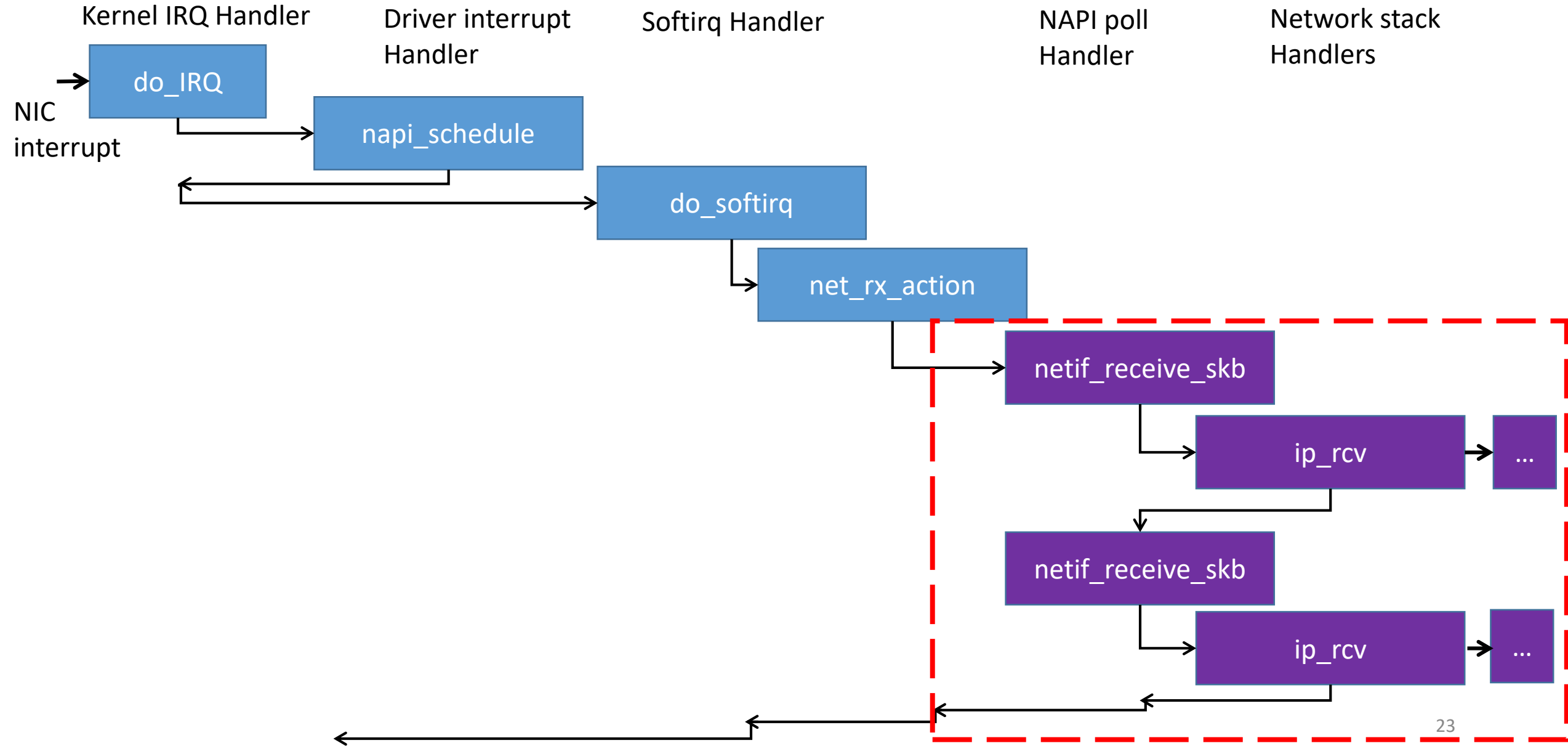
Network Call Stack – Background



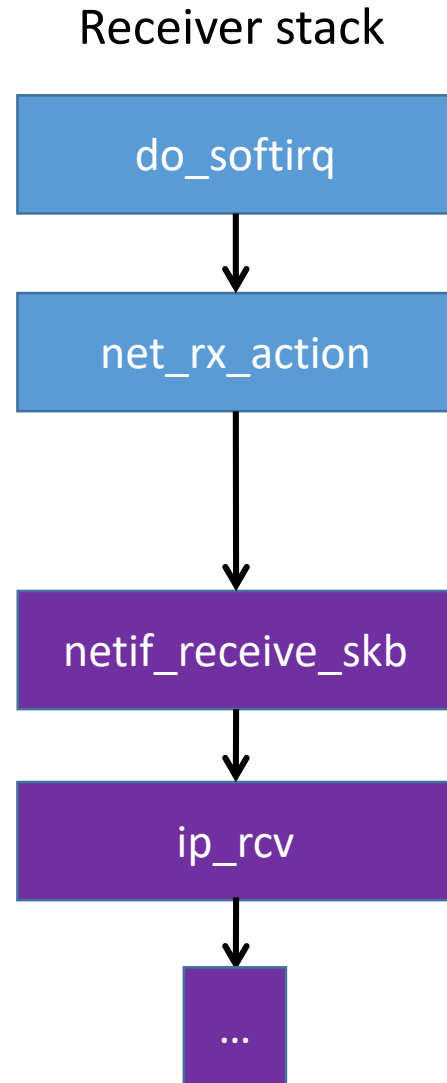
Network Call Stack – Background



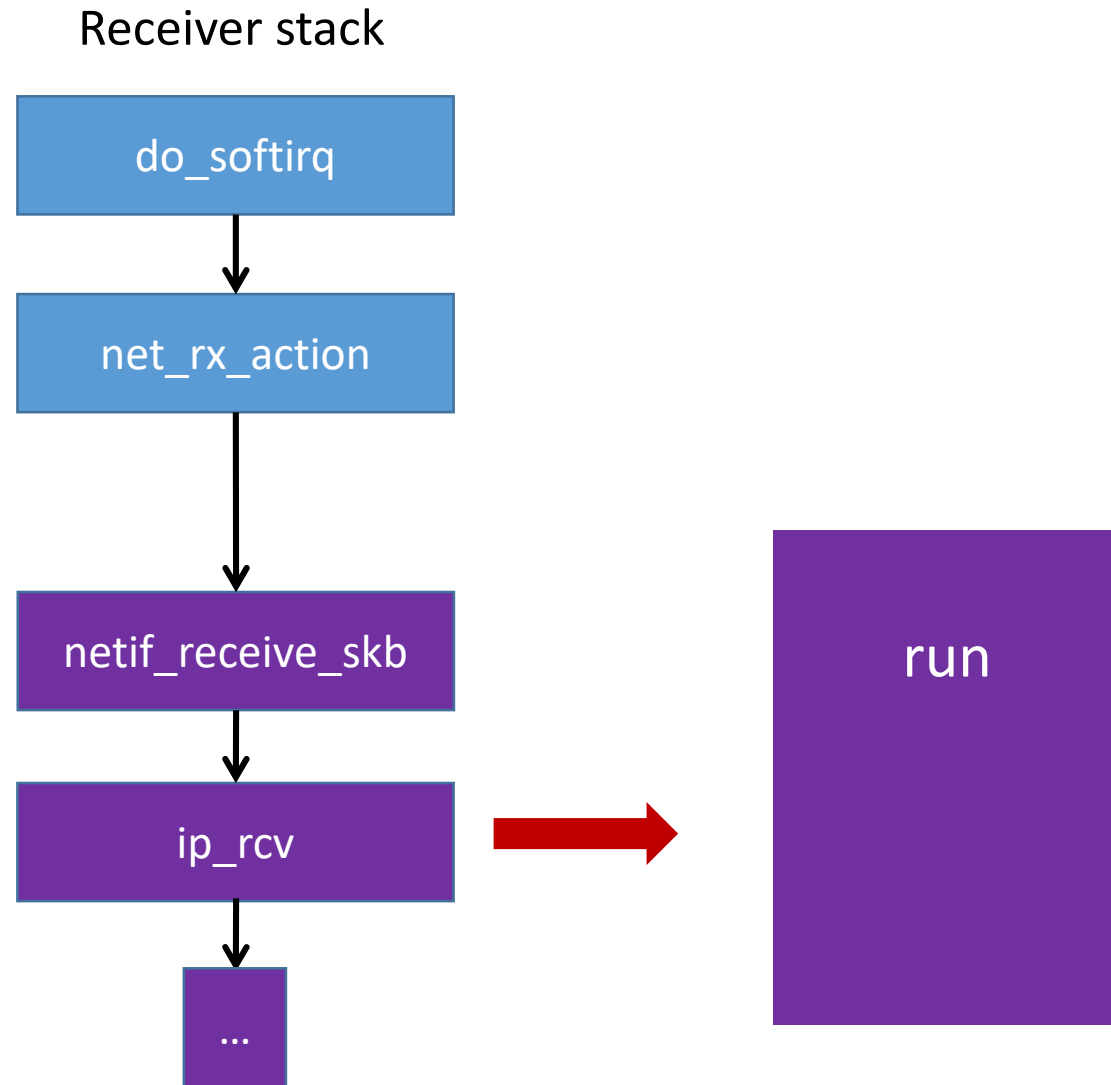
Network Call Stack – Background



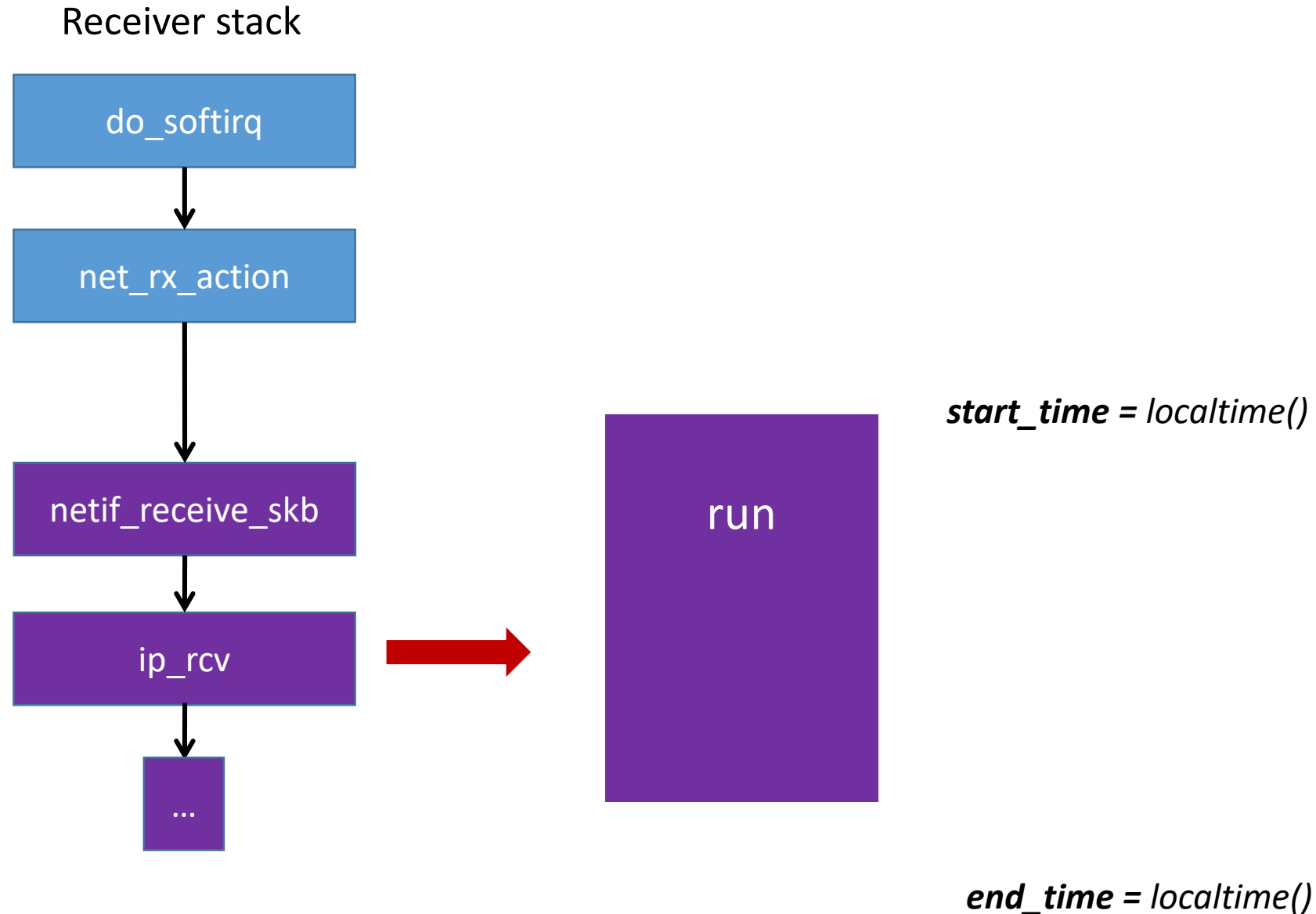
Iron – Accounting



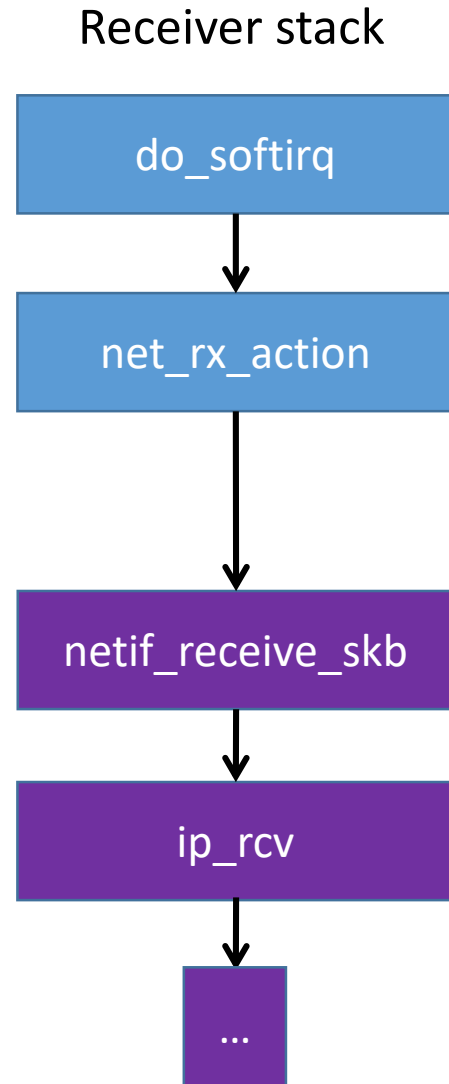
Iron – Accounting



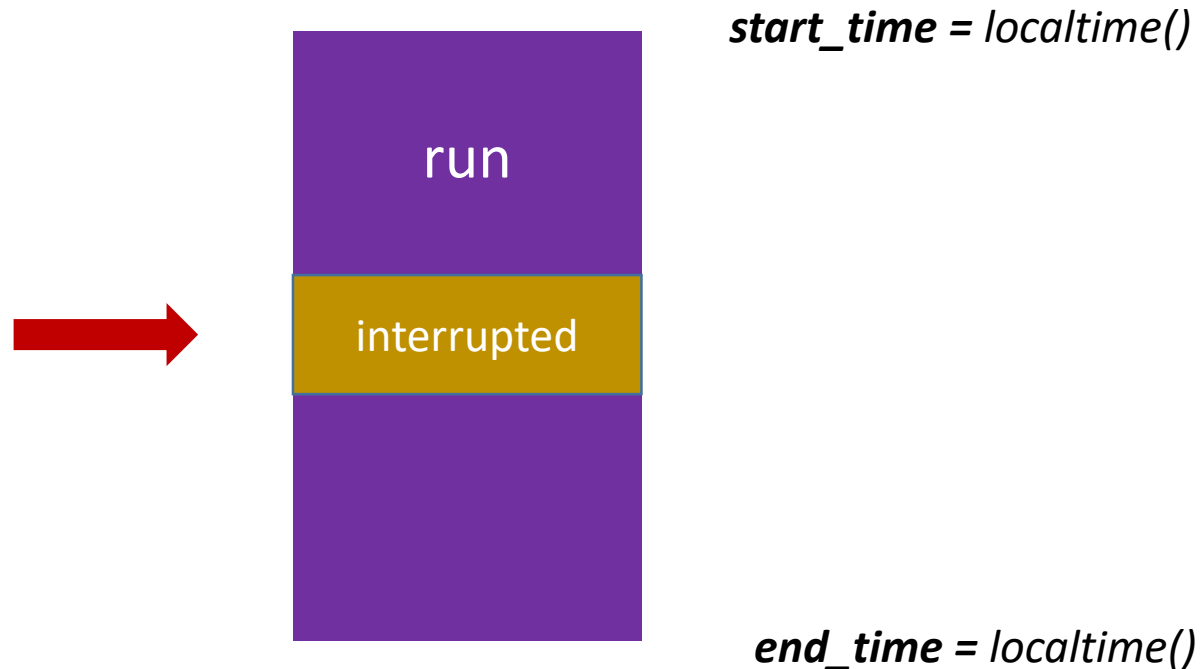
Iron – Accounting



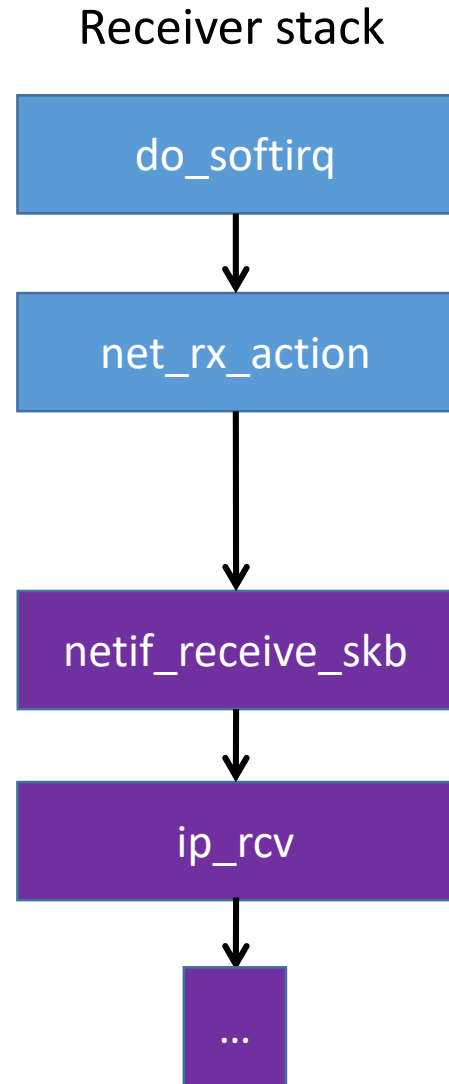
Iron – Accounting



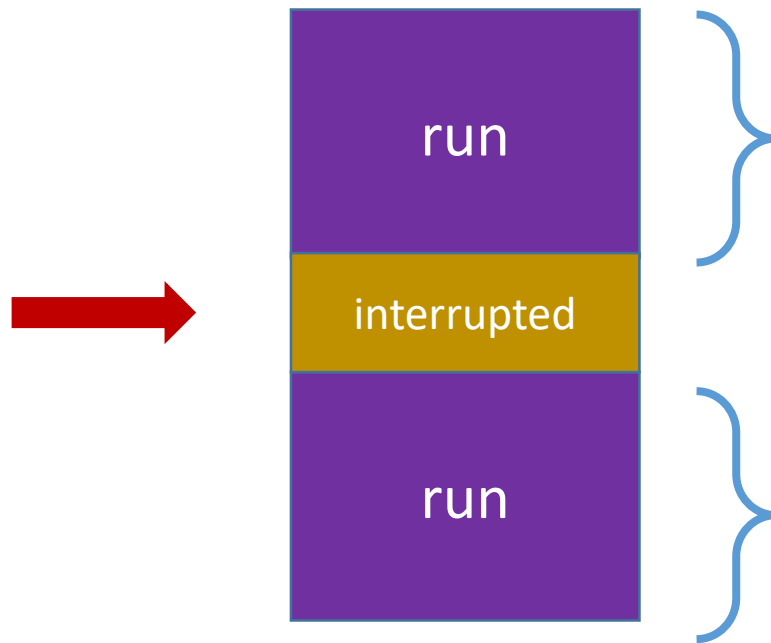
- Measuring time difference is non-trivial
 - Kernel is preemptable
 - Function in the call stack can be interrupted at any time



Iron – Accounting

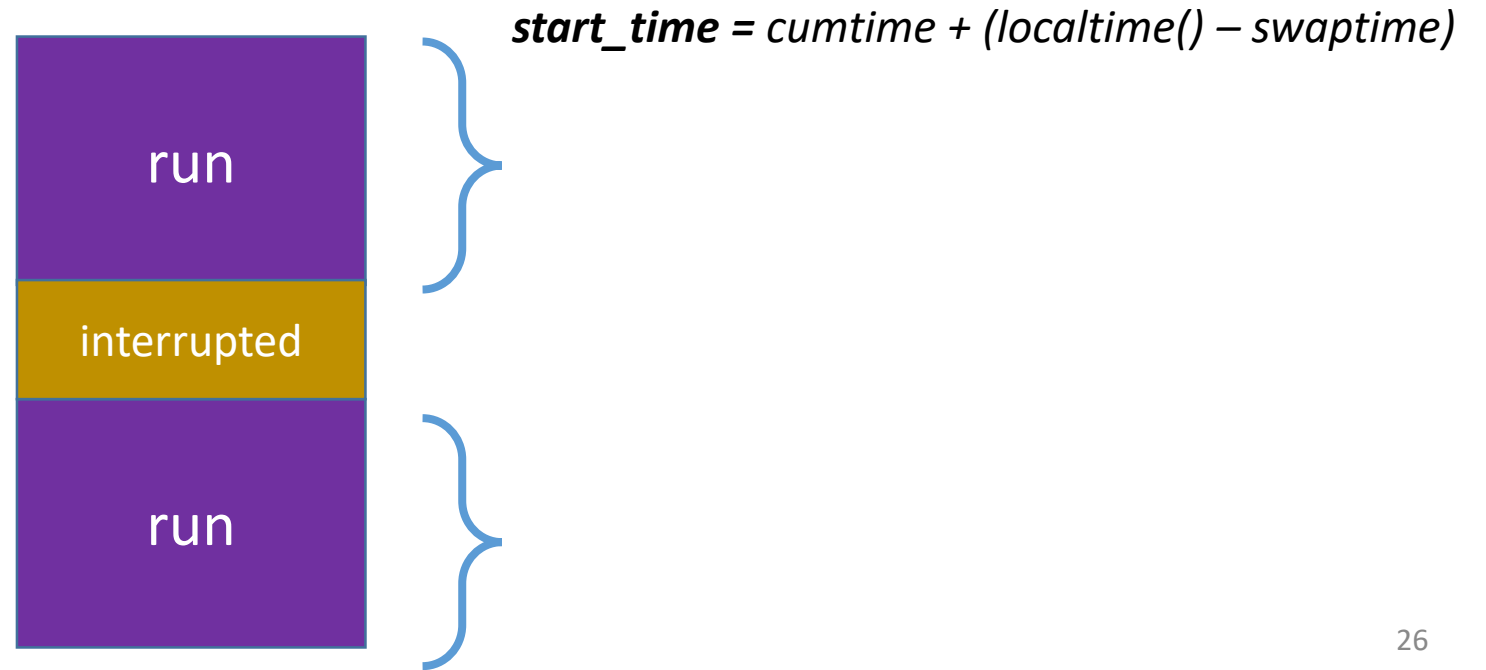
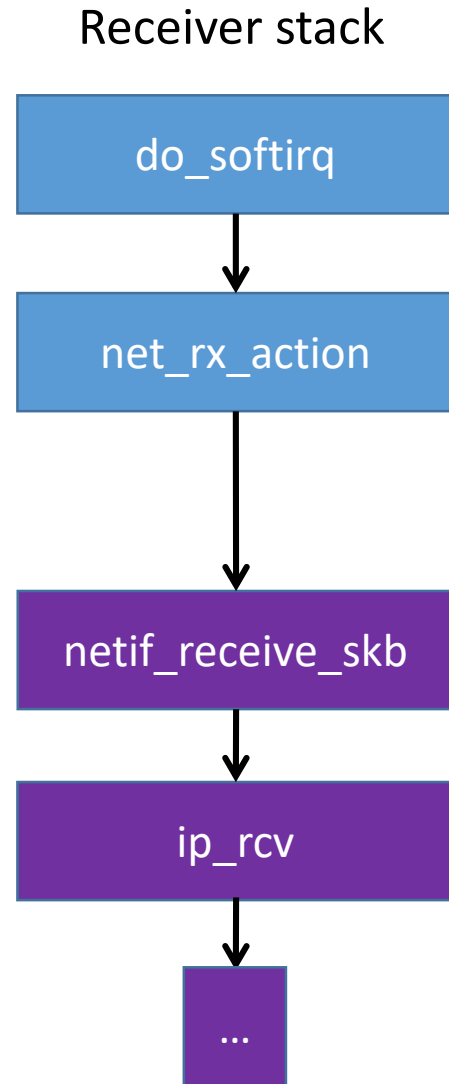


- Measuring time difference is non-trivial
 - Kernel is preemptable
 - Function in the call stack can be interrupted at any time



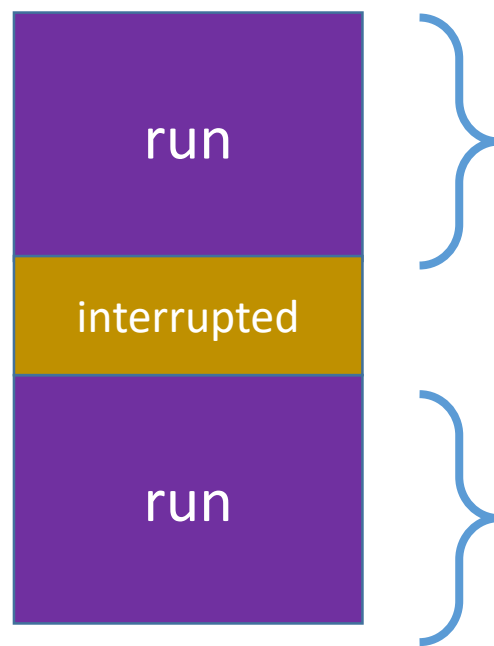
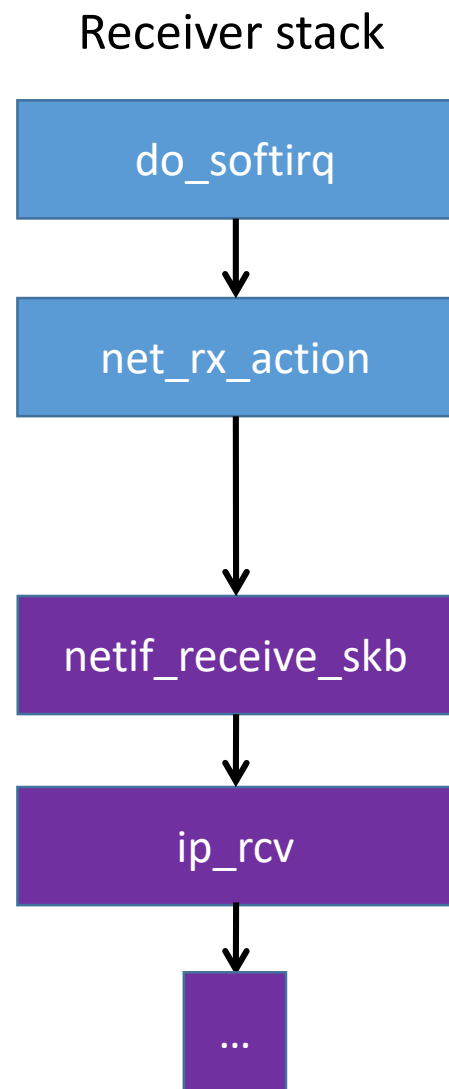
Iron – Accounting

- Measuring time difference is non-trivial
 - Kernel is preemptable
 - Function in the call stack can be interrupted at any time



Iron – Accounting

- Measuring time difference is non-trivial
 - Kernel is preemptable
 - Function in the call stack can be interrupted at any time

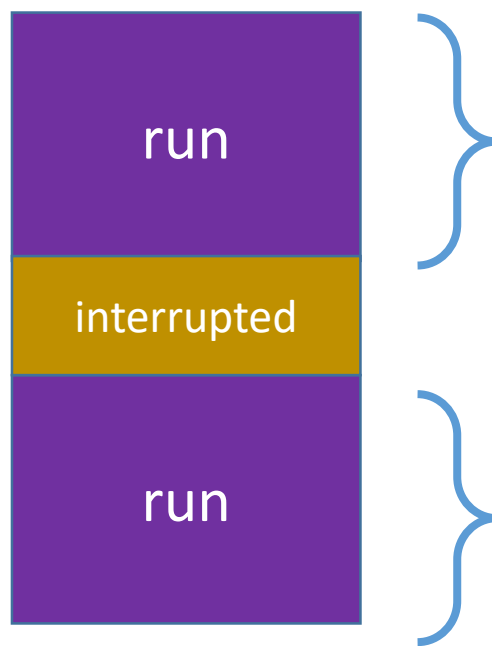
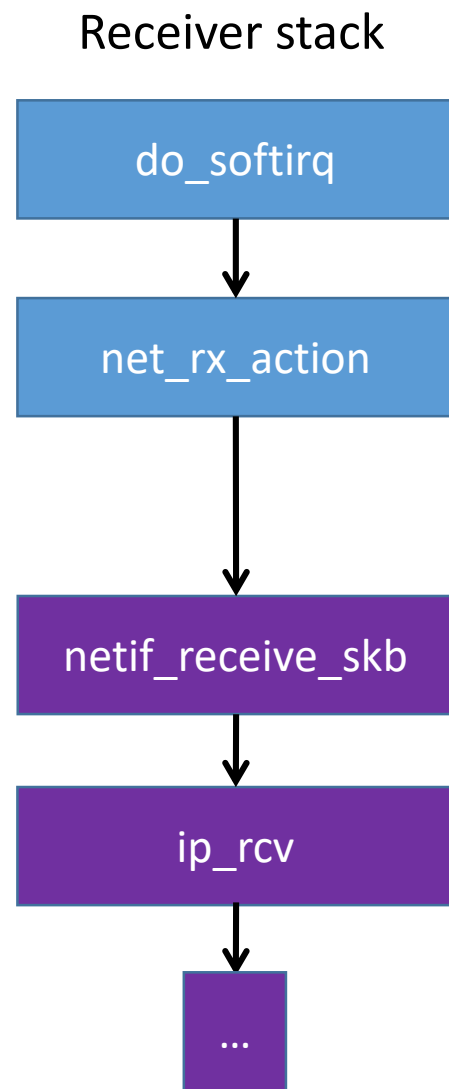


cumulative execution time

$$\text{start_time} = \text{cumtime} + (\text{localtime}() - \text{swaptime})$$

Iron – Accounting

- Measuring time difference is non-trivial
 - Kernel is preemptable
 - Function in the call stack can be interrupted at any time



cumulative execution time

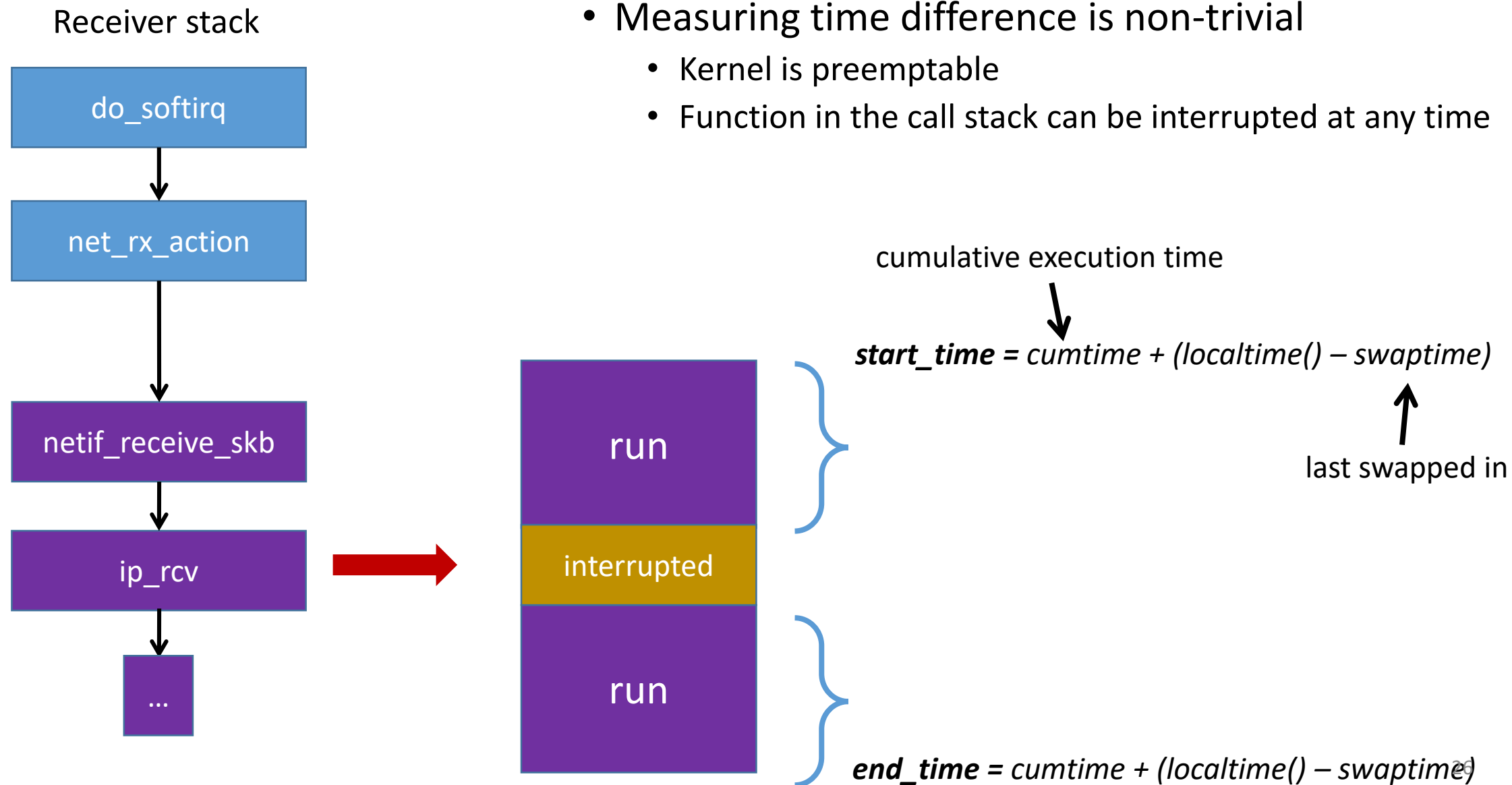
$$\text{start_time} = \text{cumtime} + (\text{localtime}() - \text{swaptime})$$

last swapped in

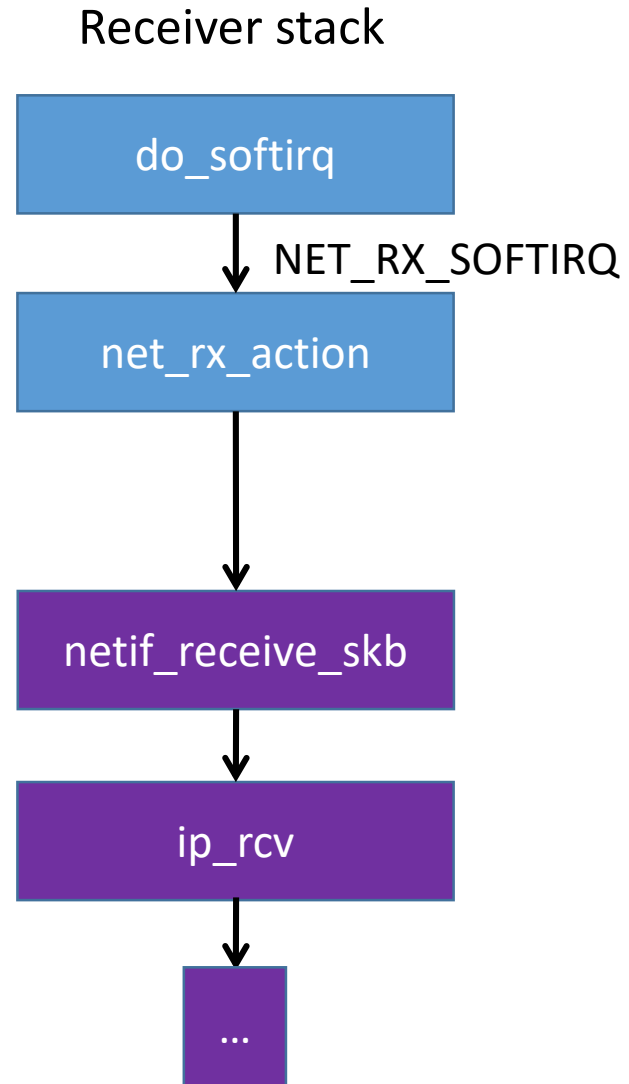
The diagram illustrates the calculation of `start_time`. It shows the equation `start_time = cumtime + (localtime() - swaptime)`. An arrow points from the text "cumulative execution time" to `cumtime`. Another arrow points from the text "last swapped in" to `swaptime`.

Iron – Accounting

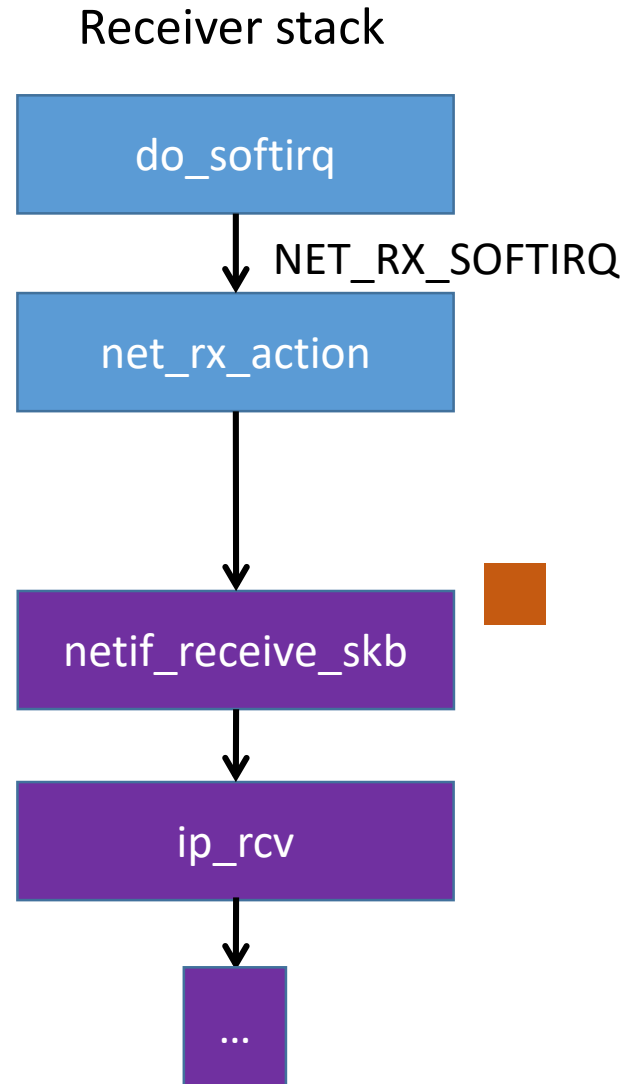
- Measuring time difference is non-trivial
 - Kernel is preemptable
 - Function in the call stack can be interrupted at any time



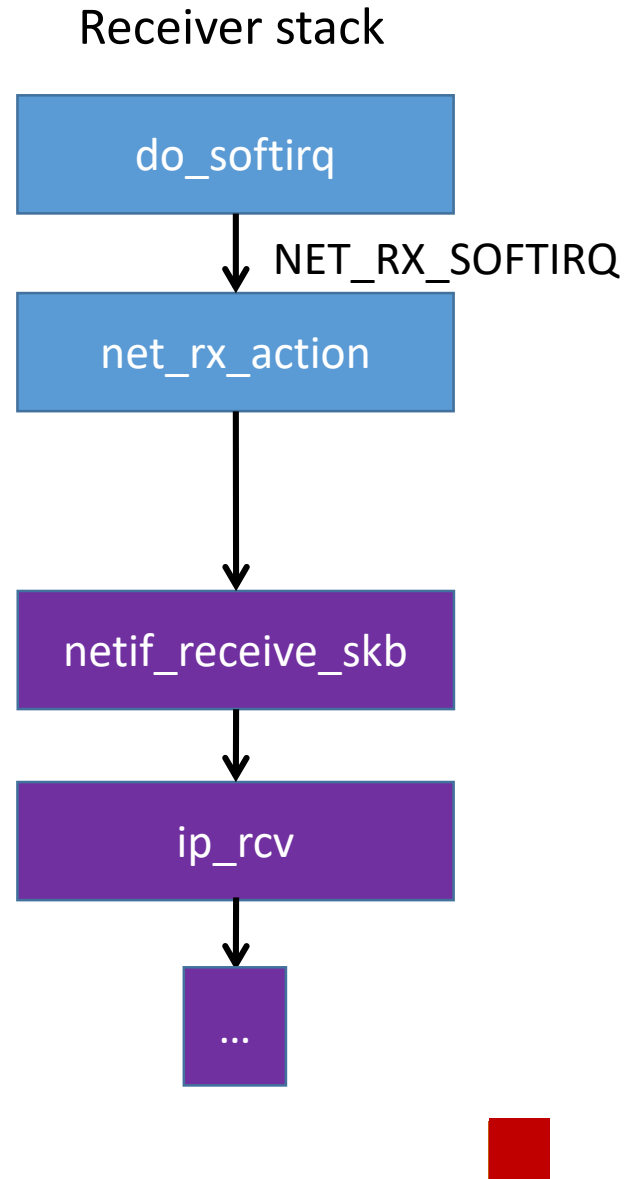
Iron – Accounting



Iron – Accounting



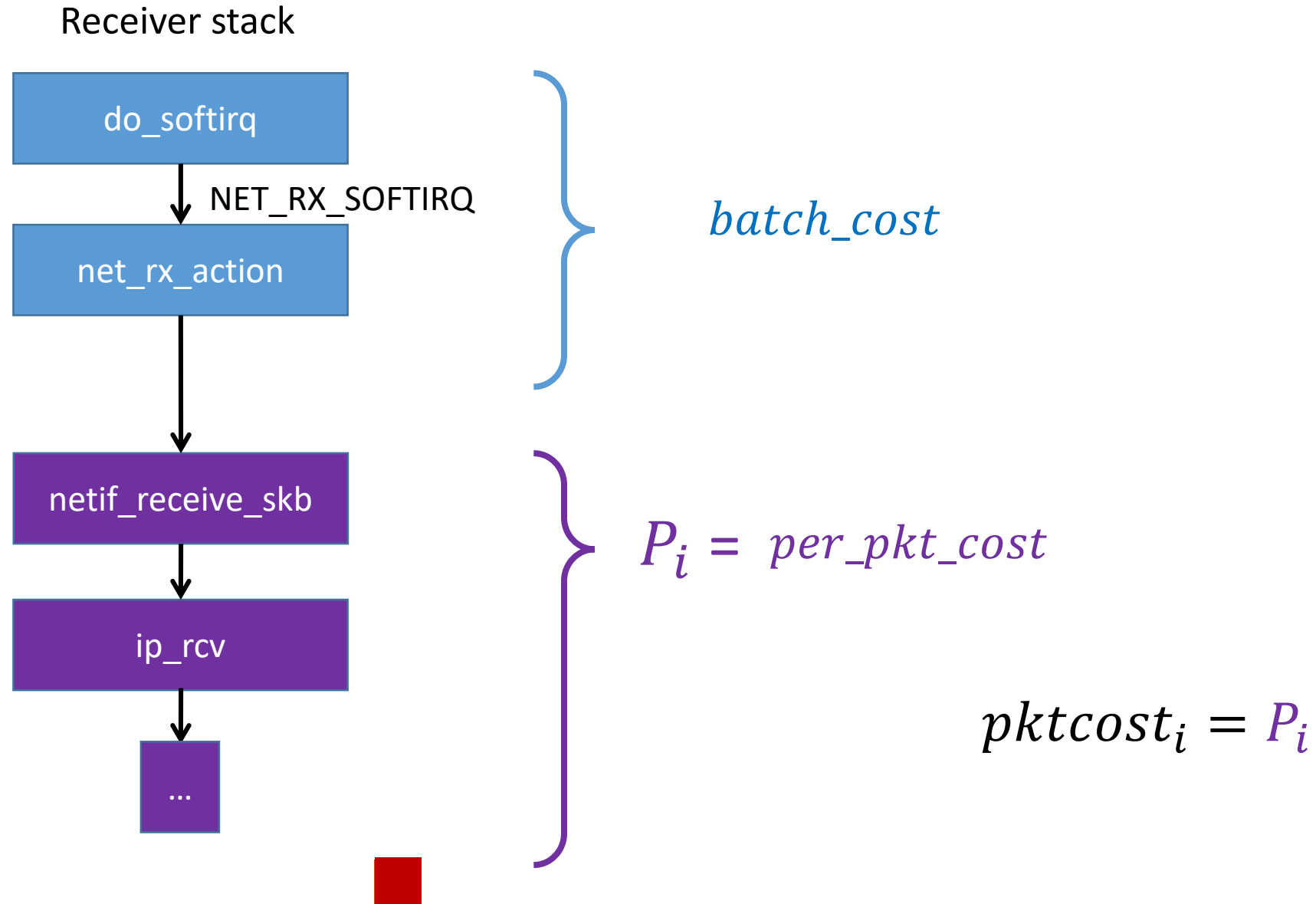
Iron – Accounting



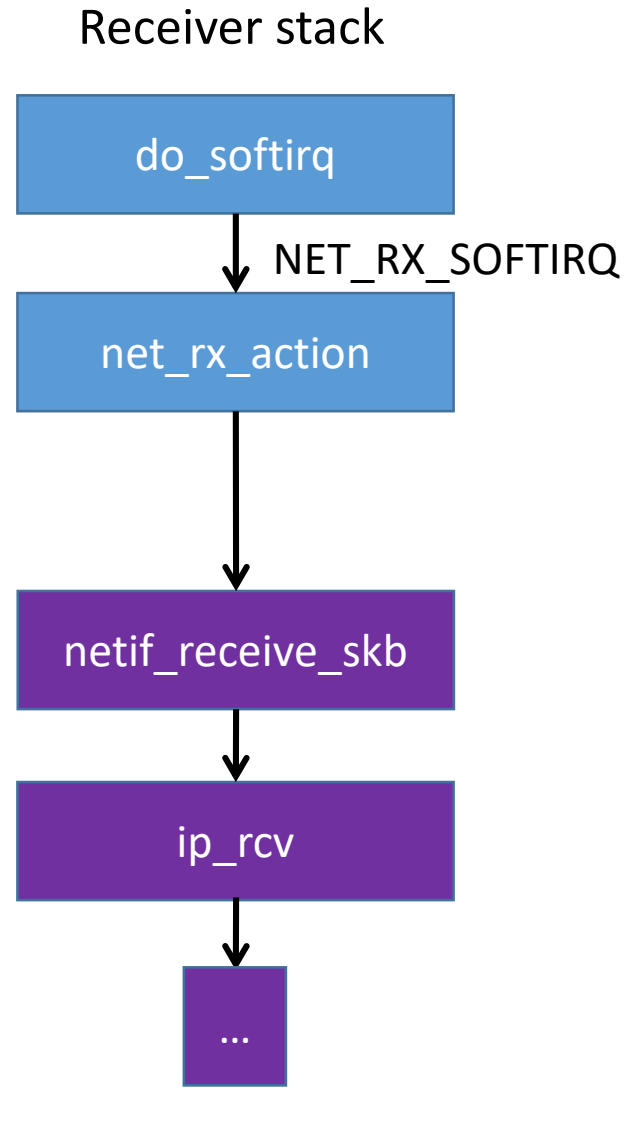
$$\left. \begin{array}{l} \text{netif_receive_skb} \\ \text{ip_rcv} \end{array} \right\} P_i = \text{per_pkt_cost}$$

$$\text{pktnetcost}_i = P_i$$

Iron – Accounting



Iron – Accounting

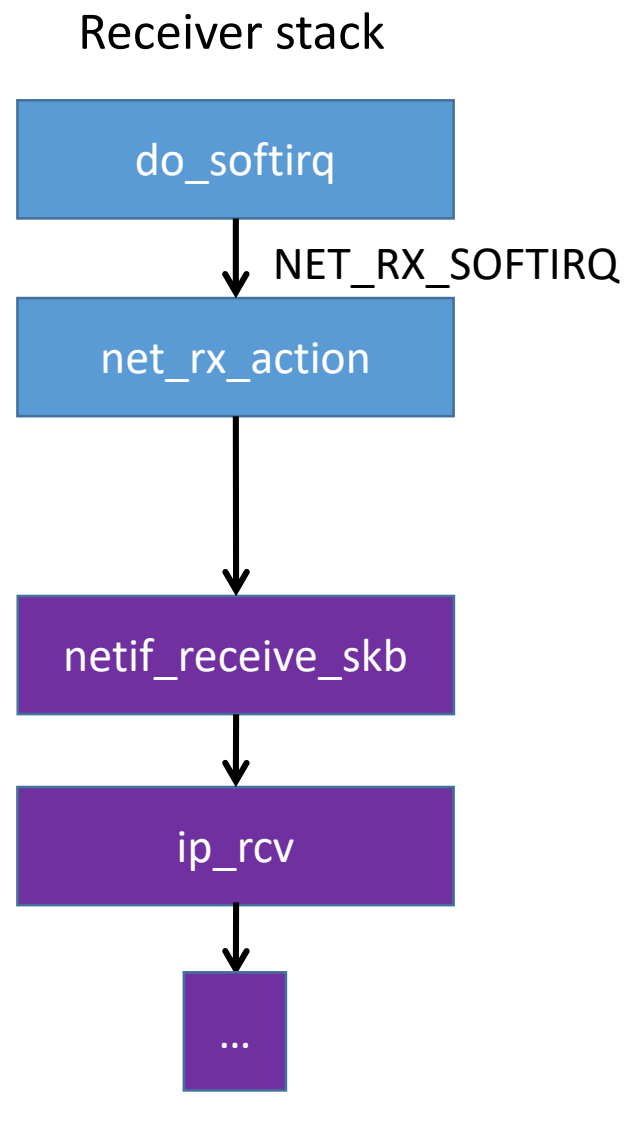


batch_cost

$P_i = \text{per_pkt_cost}$

$$pktnost_i = P_i + \frac{\text{batch_cost}}{|P|}$$

Iron – Accounting



HI, TX, RX, TIMER, SCSI & TASKLET

$$batch_cost = \left(do_softirq_cost - \sum S_i \right) * \frac{S_{RX}}{\sum S_i}$$

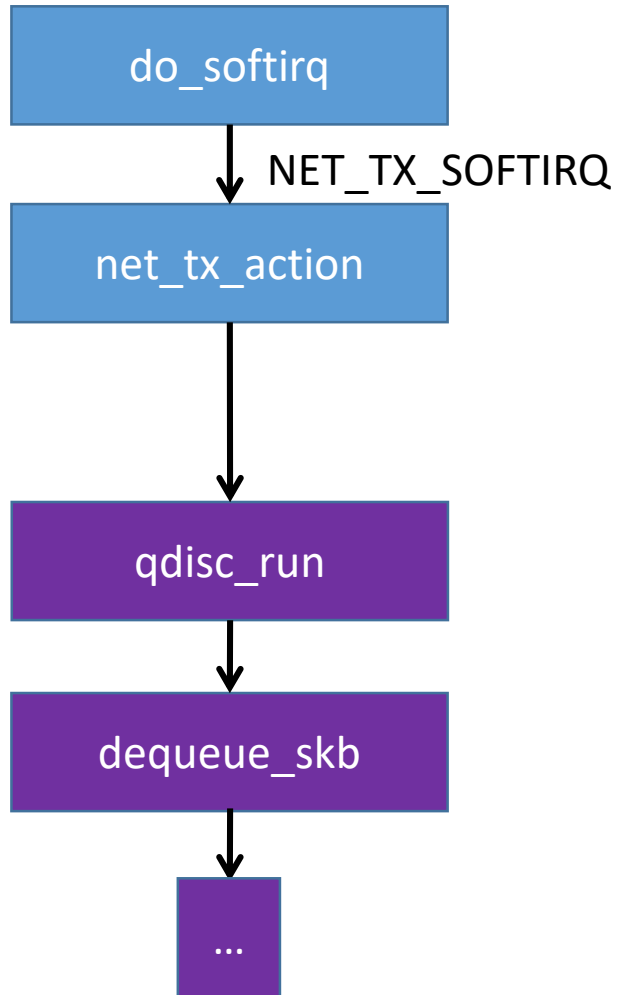
cost of all the softirqs

$$P_i = per_pkt_cost$$

$$pktpcost_i = P_i + \frac{batch_cost}{|P|}$$

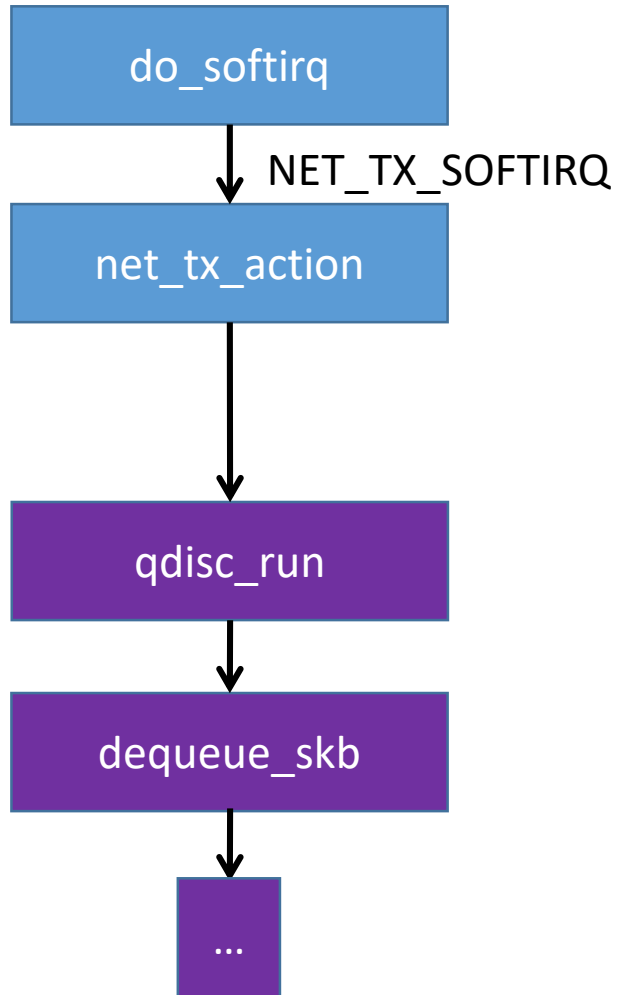
Iron – Accounting

Transmitter stack



Iron – Accounting

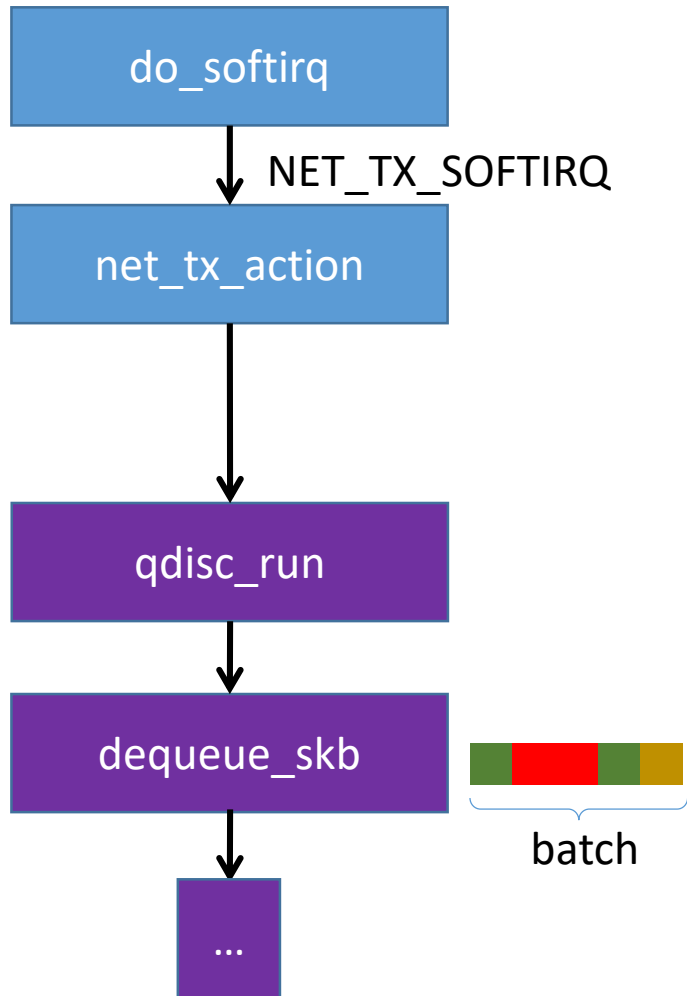
Transmitter stack



- Linux batches packets for transmission

Iron – Accounting

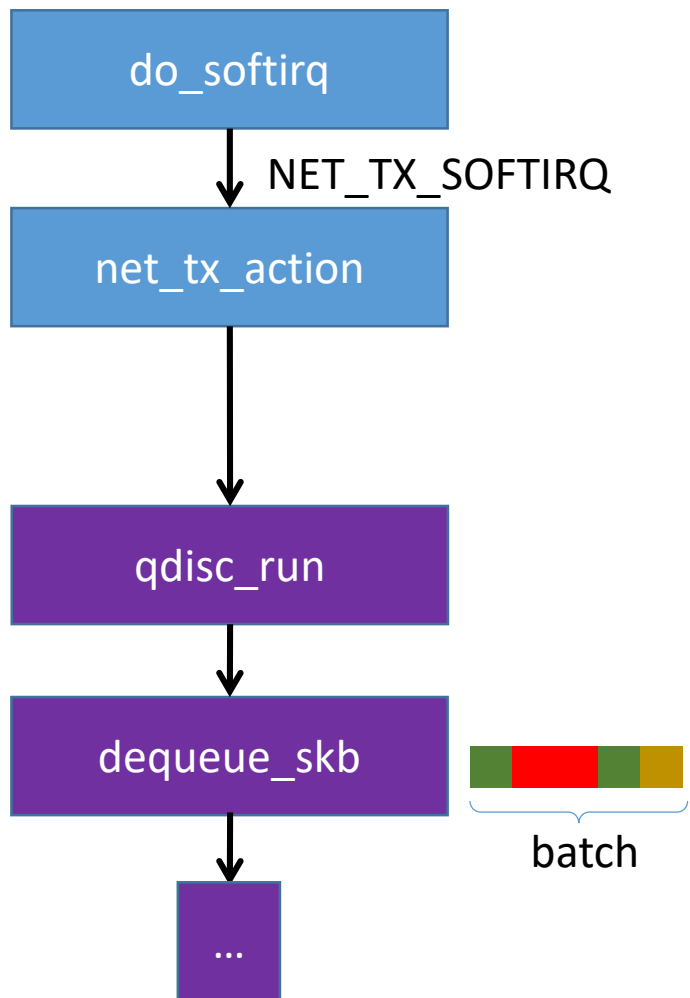
Transmitter stack



- Linux batches packets for transmission

Iron – Accounting

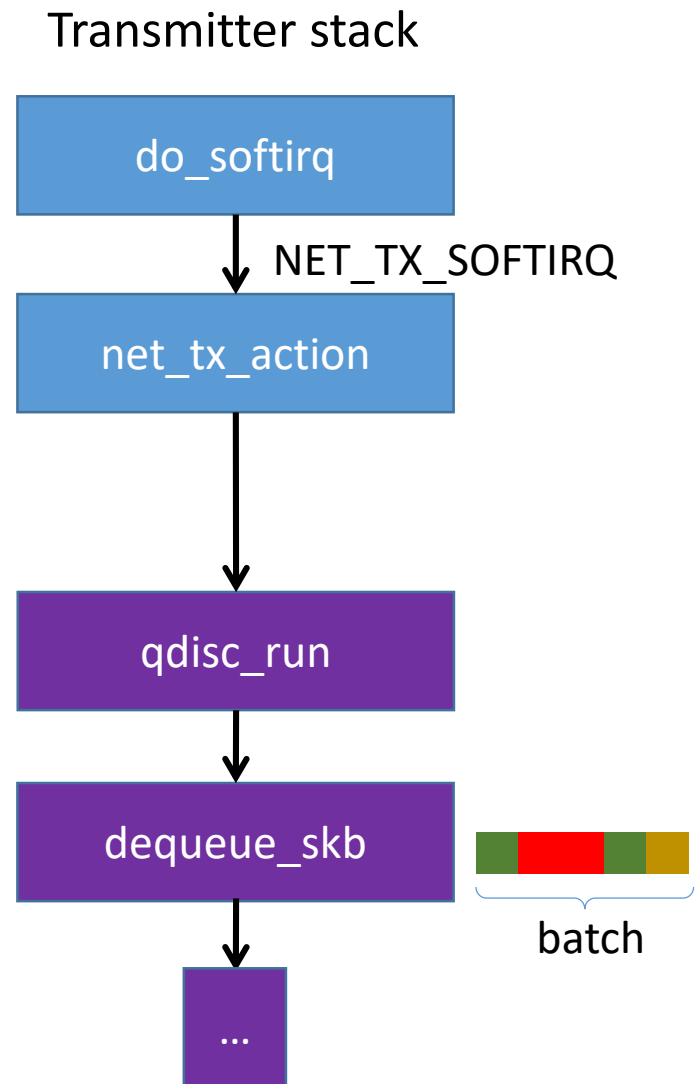
Transmitter stack



- Linux batches packets for transmission
- We measure the cost of the batch and charge each packet within the batch an equal share

$$pktnost = \frac{batch_cost}{batch_size}$$

Iron – Accounting



- Linux batches packets for transmission
- We measure the cost of the batch and charge each packet within the batch an equal share
- To identify the container to charge at dequeue
 - We encode the container information in the skb while enqueueing the packet

$$pktnost = \frac{batch_cost}{batch_size}$$

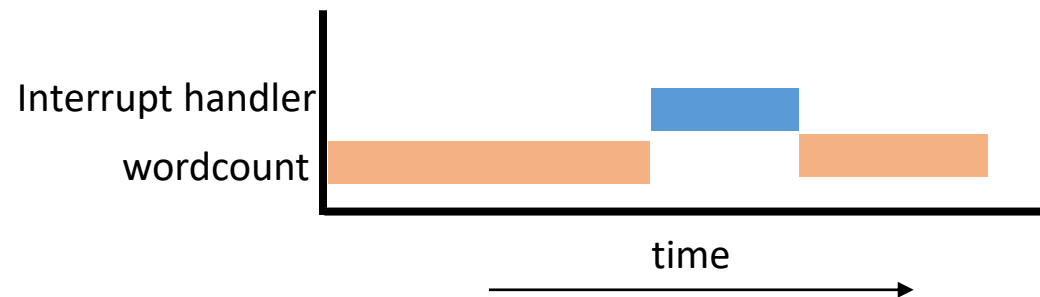
Outline

- How and by how much is isolation broken
- Iron's Design
 - Accounting of per-packet processing cost
 - **Ensuring isolation via enforcement**
 - **Integration with Linux scheduler**
 - Hardware-based packet dropping
- Evaluation
 - Controlled workload
 - Realistic workload

Iron – Enforcement

Scheduler Integration

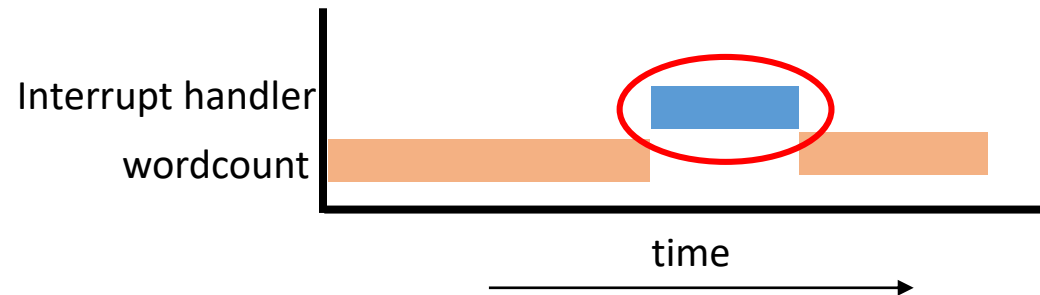
- ***Return the accounted time*** to the container which was ***incorrectly*** charged



Iron – Enforcement

Scheduler Integration

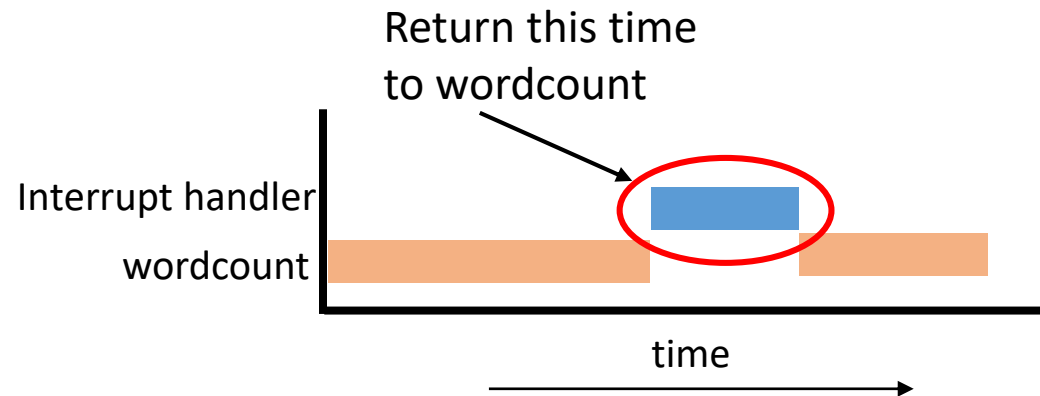
- ***Return the accounted time*** to the container which was ***incorrectly*** charged



Iron – Enforcement

Scheduler Integration

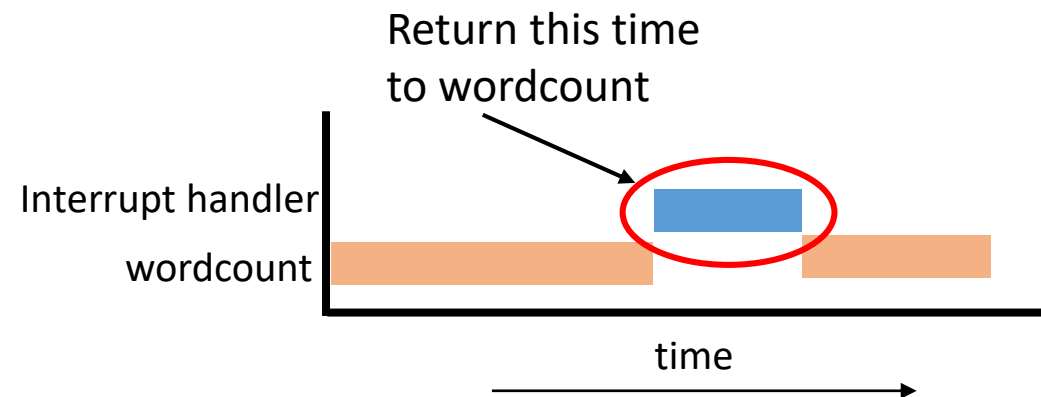
- ***Return the accounted time*** to the container which was ***incorrectly*** charged



Iron – Enforcement

Scheduler Integration

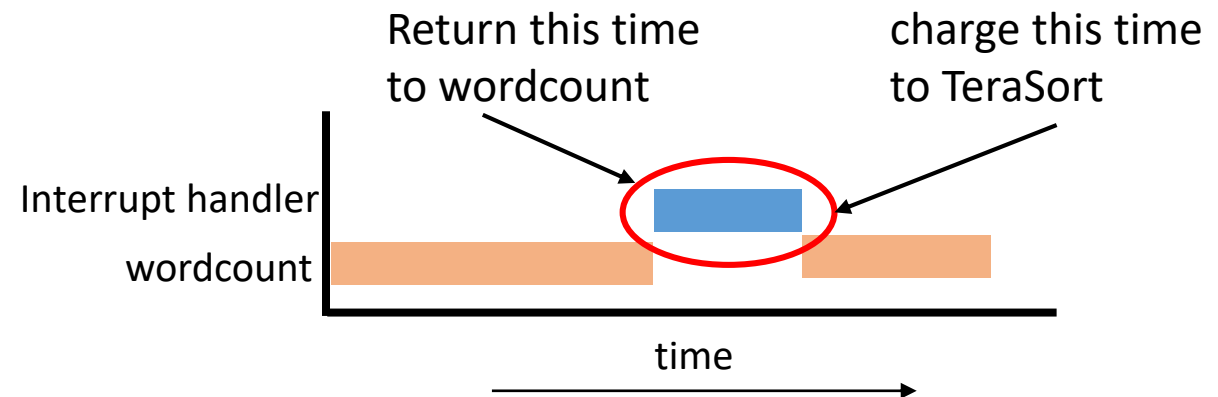
- ***Return the accounted time*** to the container which was ***incorrectly*** charged
- ***Charge the accounted time*** to the container which was ***responsible*** for the network traffic



Iron – Enforcement

Scheduler Integration

- ***Return the accounted time*** to the container which was ***incorrectly*** charged
- ***Charge the accounted time*** to the container which was ***responsible*** for the network traffic



Iron – Enforcement

Scheduler Integration

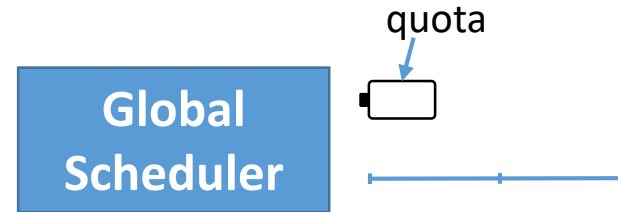
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

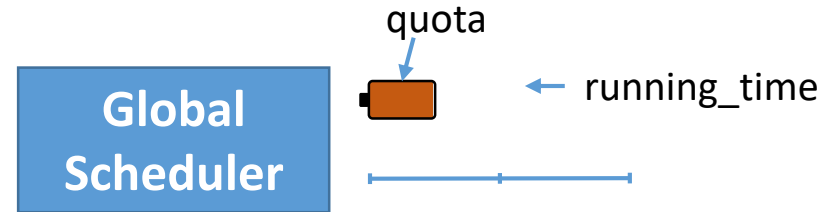
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

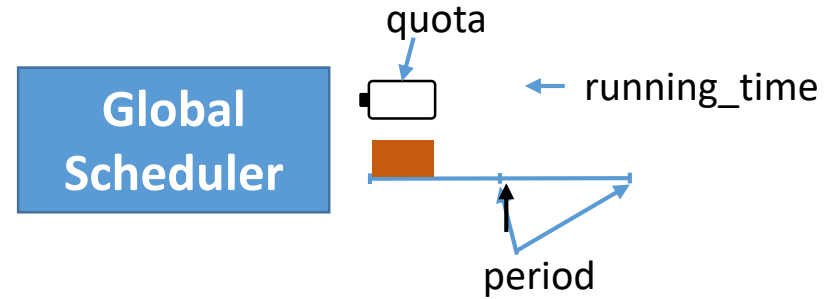
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

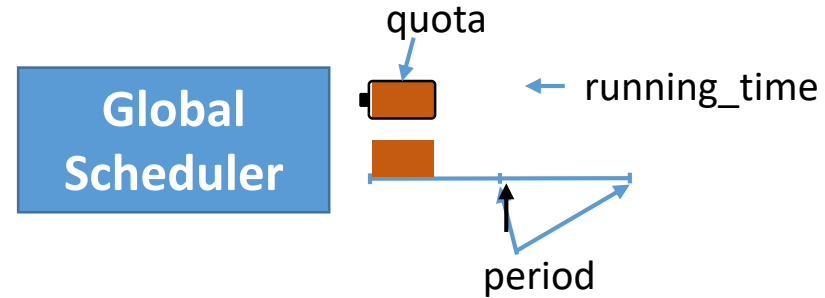
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

Reuse infrastructure from cgroup and Linux scheduler

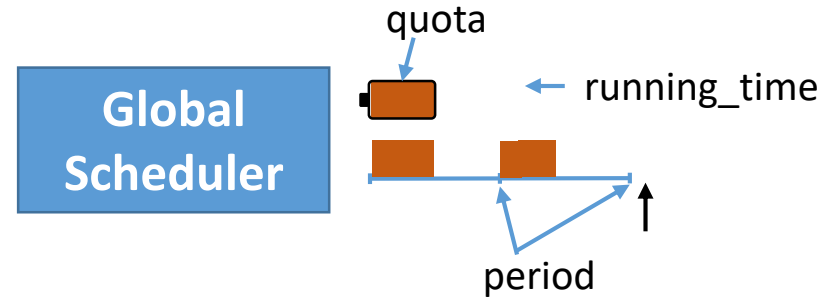


At the end of the period,
running_time is refilled by *quota*.

Iron – Enforcement

Scheduler Integration

Reuse infrastructure from cgroup and Linux scheduler

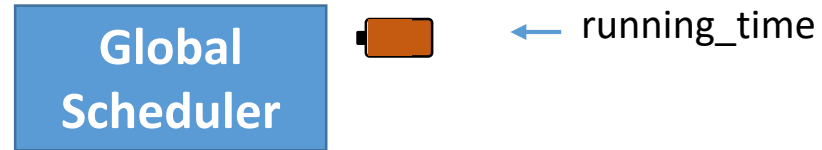


At the end of the period,
running_time is refilled by *quota*.

Iron – Enforcement

Scheduler Integration

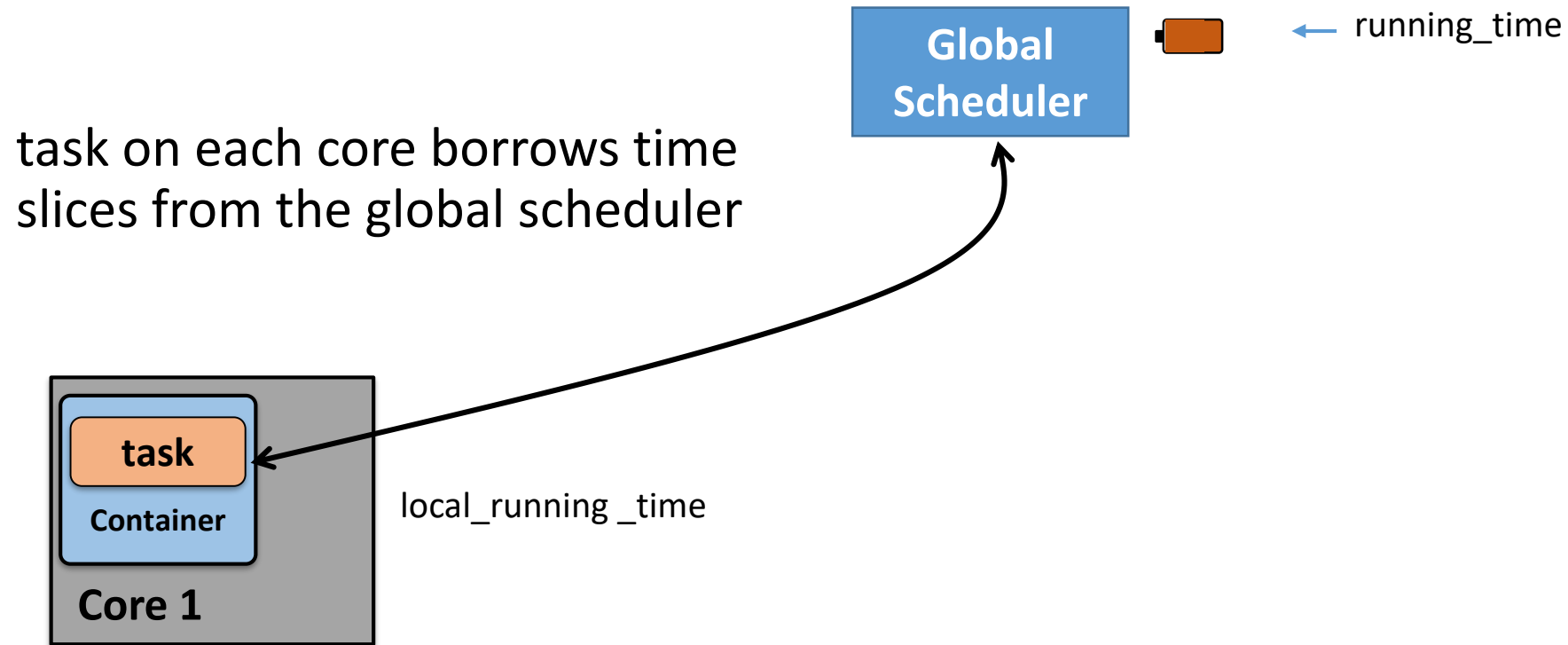
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

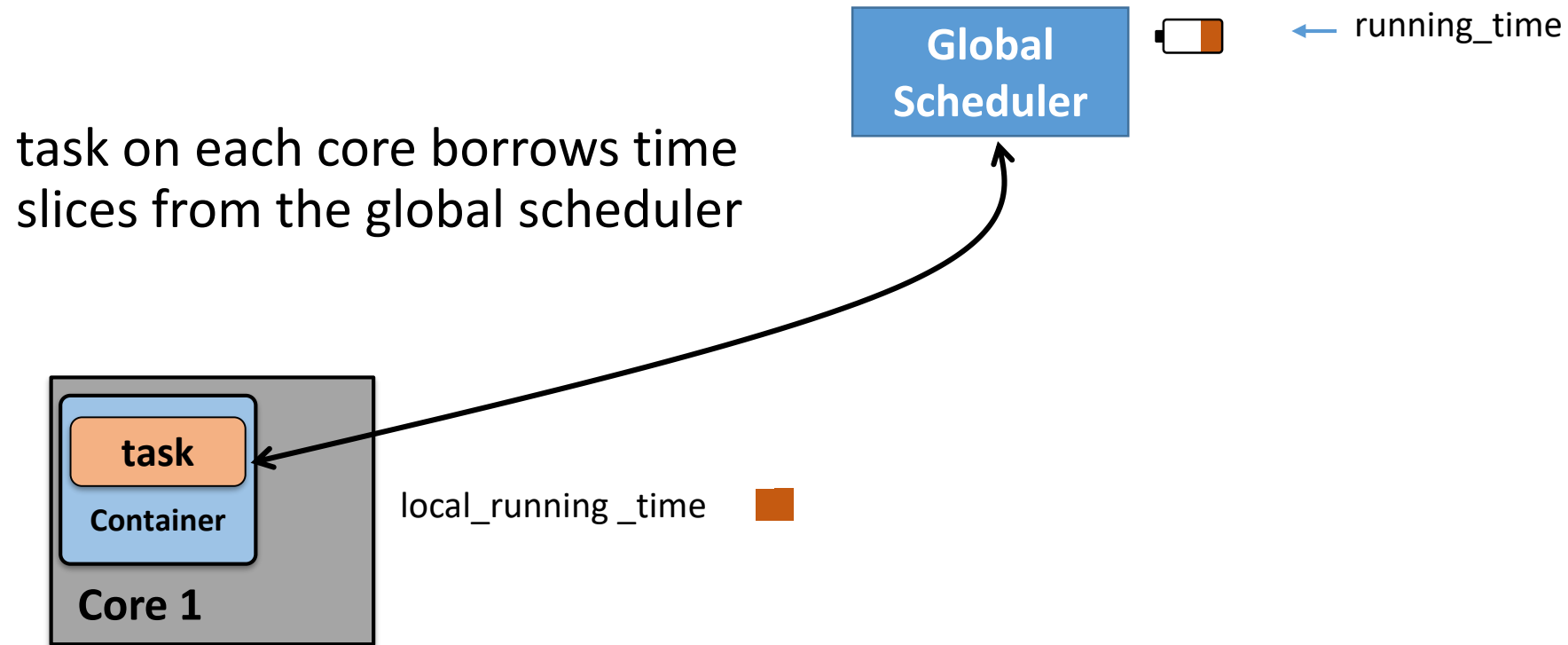
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

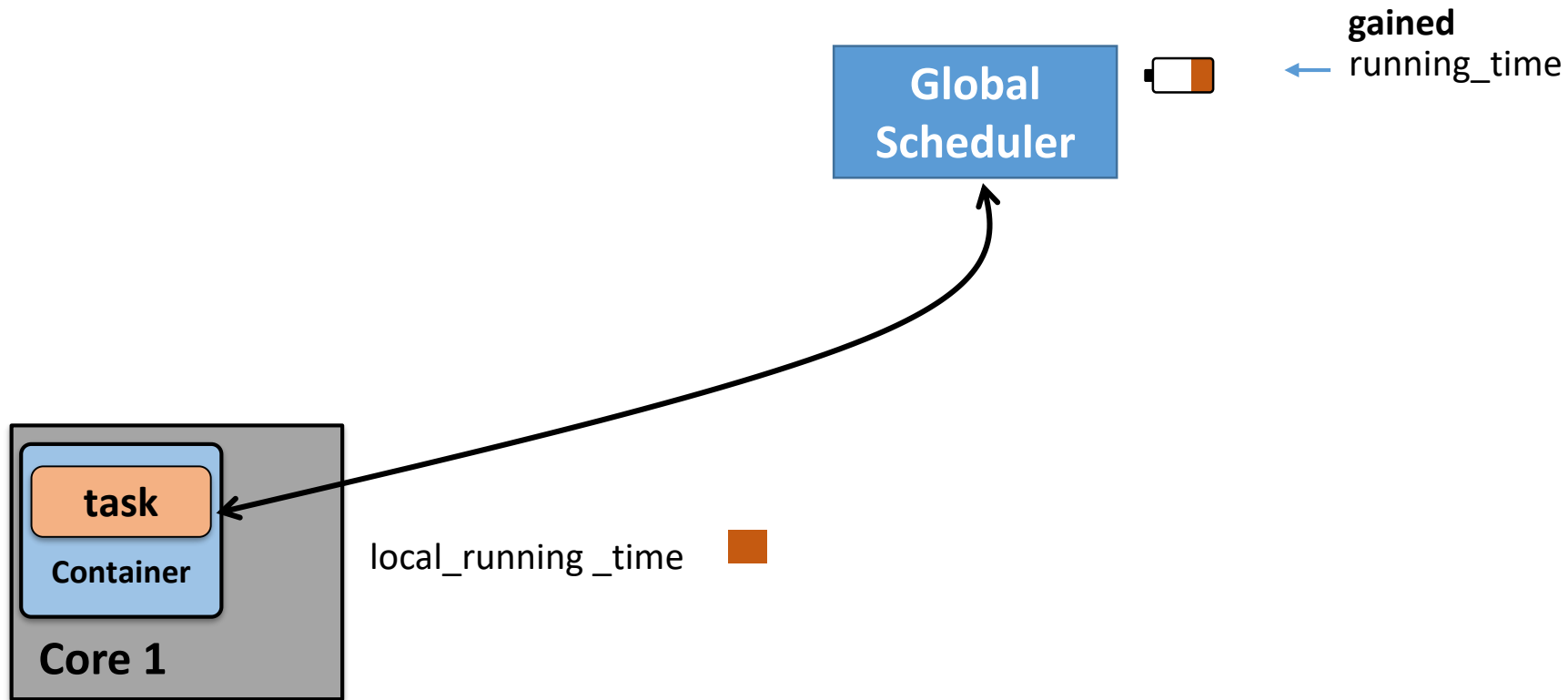
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

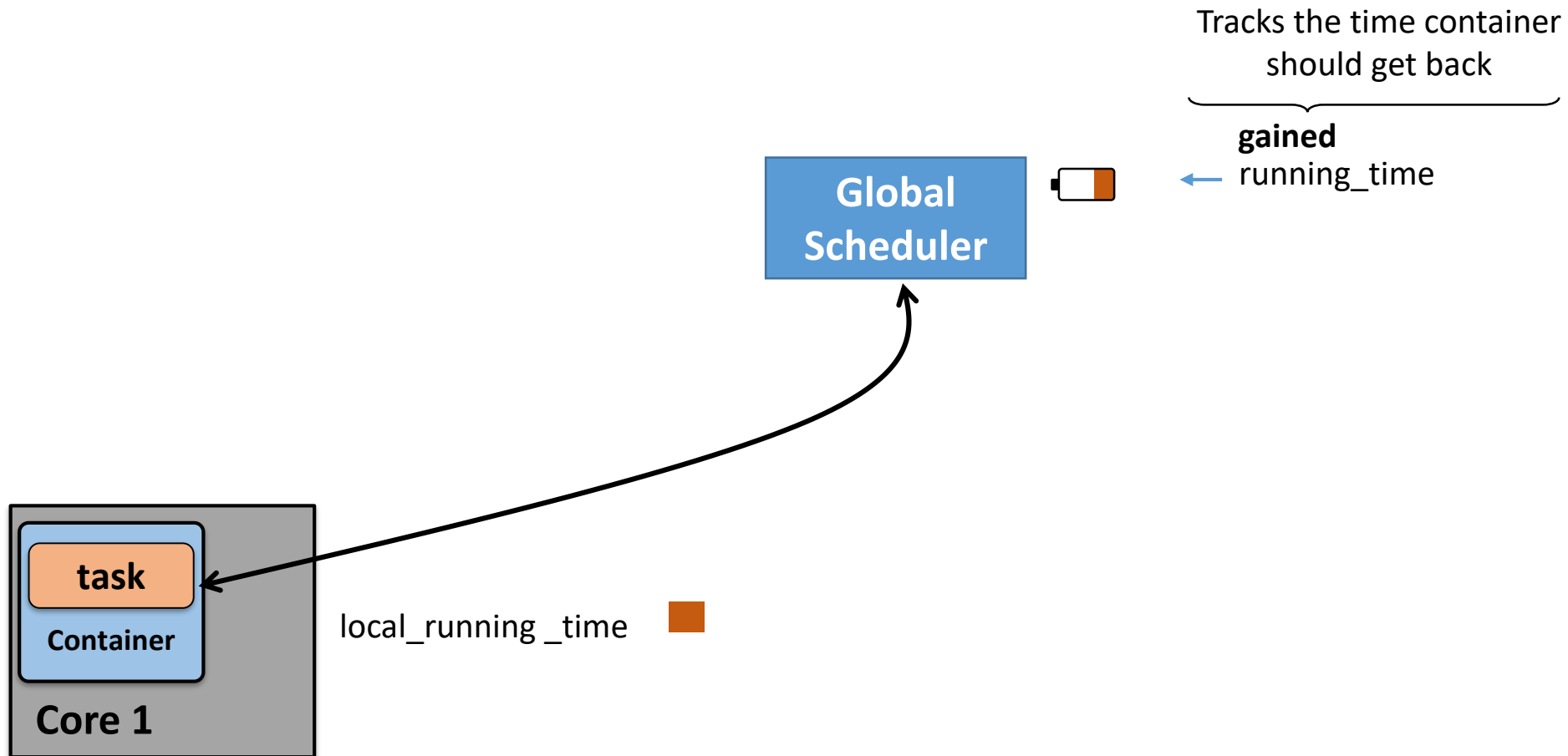
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

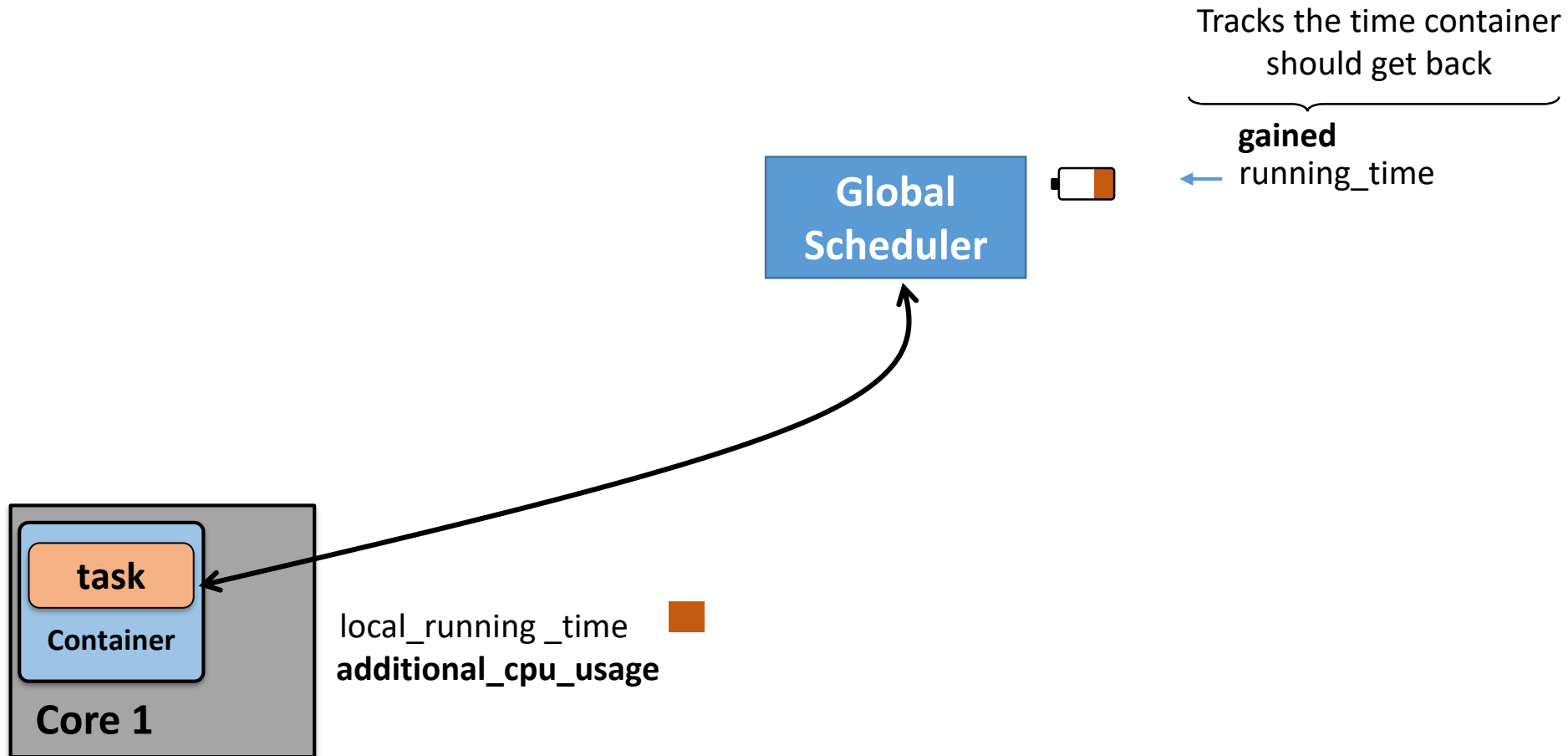
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

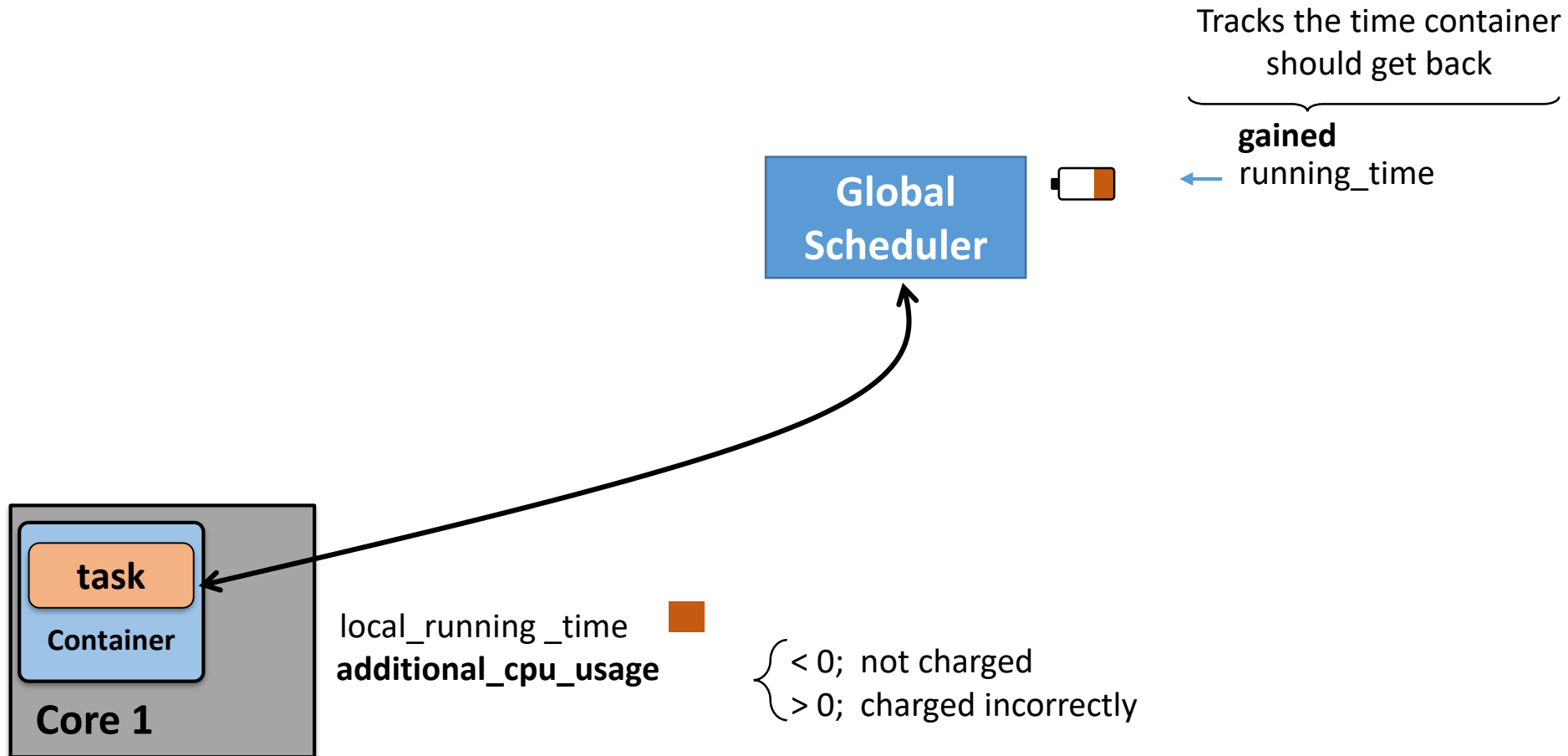
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

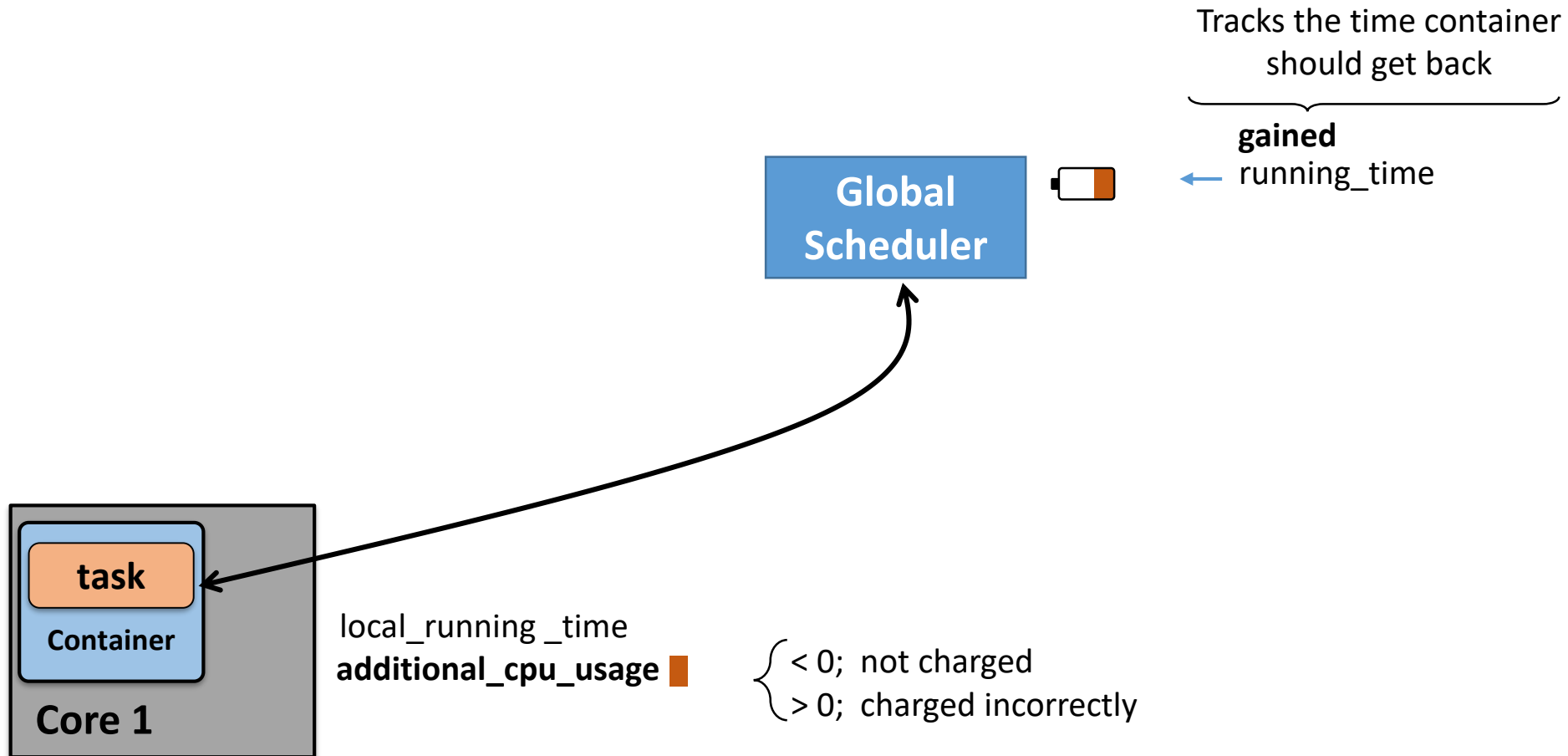
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

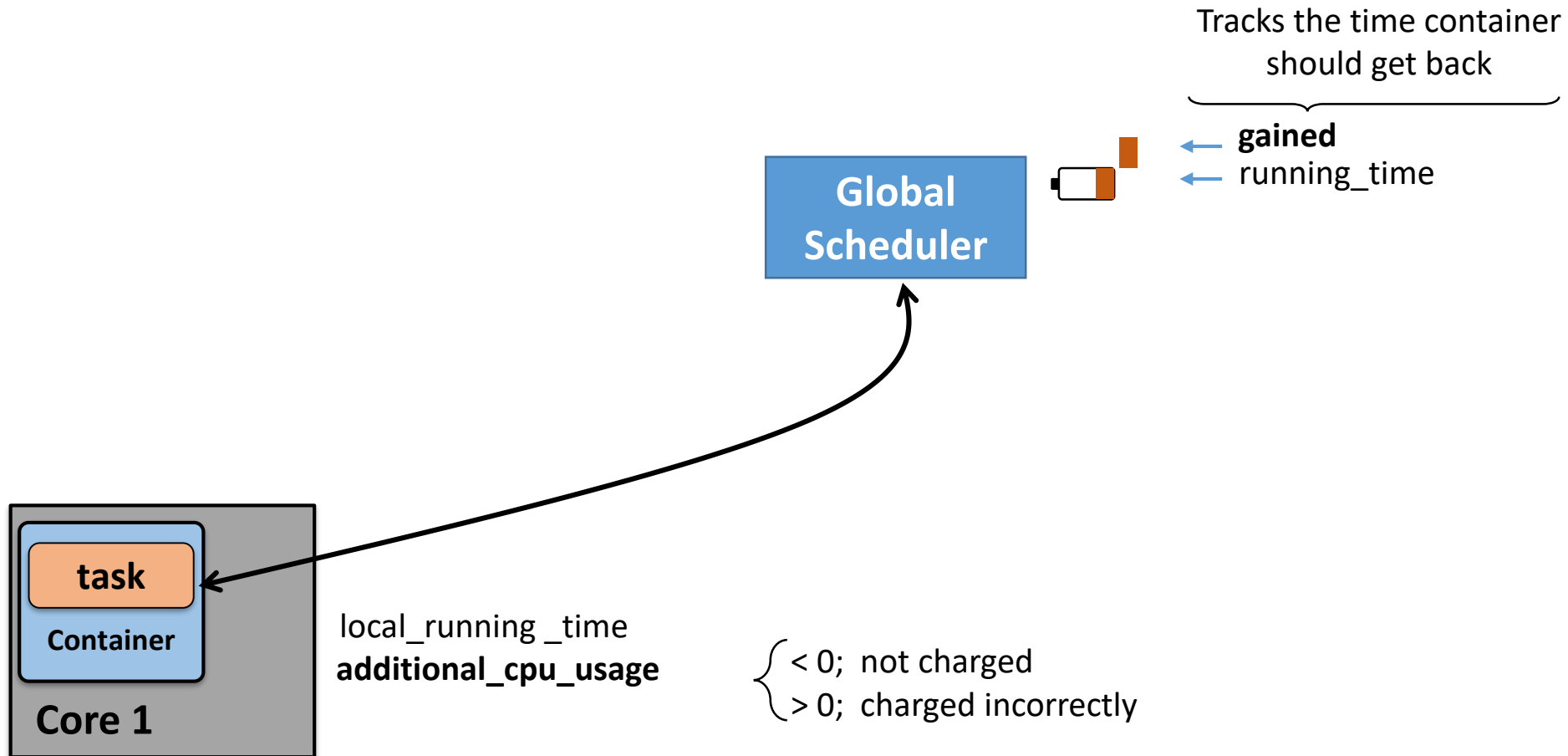
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

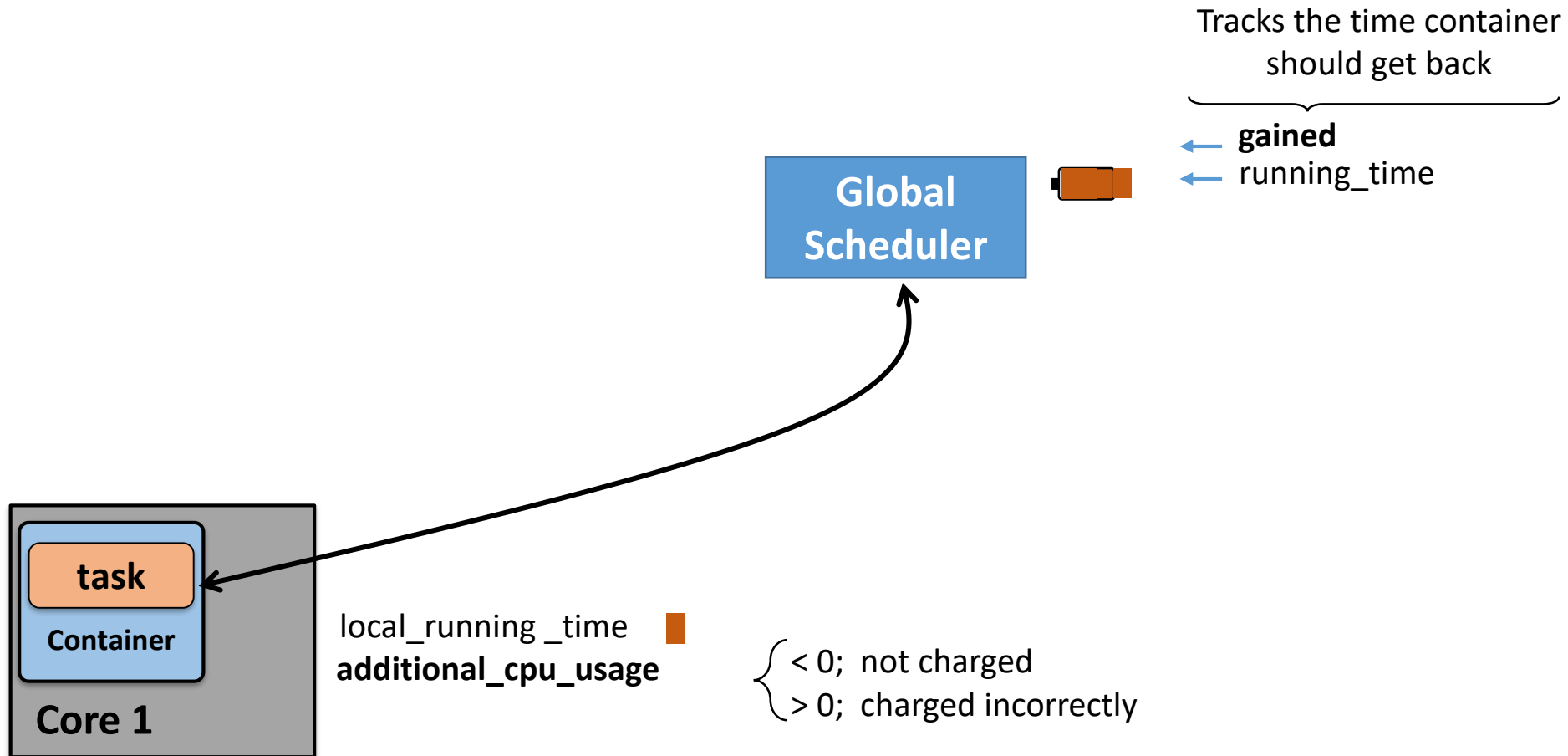
Reuse infrastructure from cgroup and Linux scheduler



Iron – Enforcement

Scheduler Integration

Reuse infrastructure from cgroup and Linux scheduler

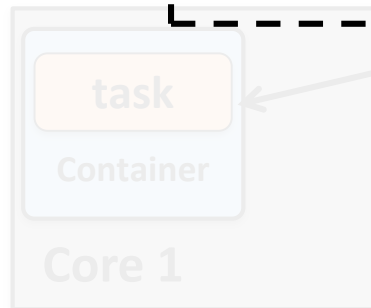


Iron – Enforcement

Scheduler Integration

Reuse infrastructure from cgroup and Linux scheduler

Throttling a sender ensures isolation! Because throttled sender ($\text{runtime} < 0$) ***cannot generate*** outgoing traffic.



■ $\text{rt_remain} = \min(\text{runtime}, \text{min_amount})$
■ $\text{additional_cpu_usage}$

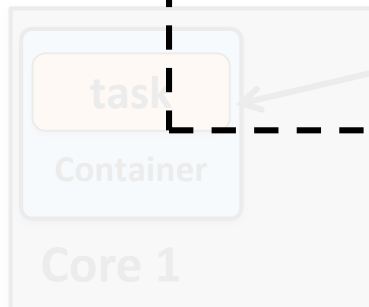
Iron – Enforcement

Scheduler Integration

Reuse infrastructure from cgroup and Linux scheduler

Throttling a sender ensures isolation! Because throttled sender ($\text{runtime} < 0$) ***cannot generate*** outgoing traffic.

If the receiver is throttled, incoming traffic can still arrive and consume CPU.



$\text{rt_remain} = \min(\text{runtime}, \text{min_amount})$
additional_cpu_usage

Outline

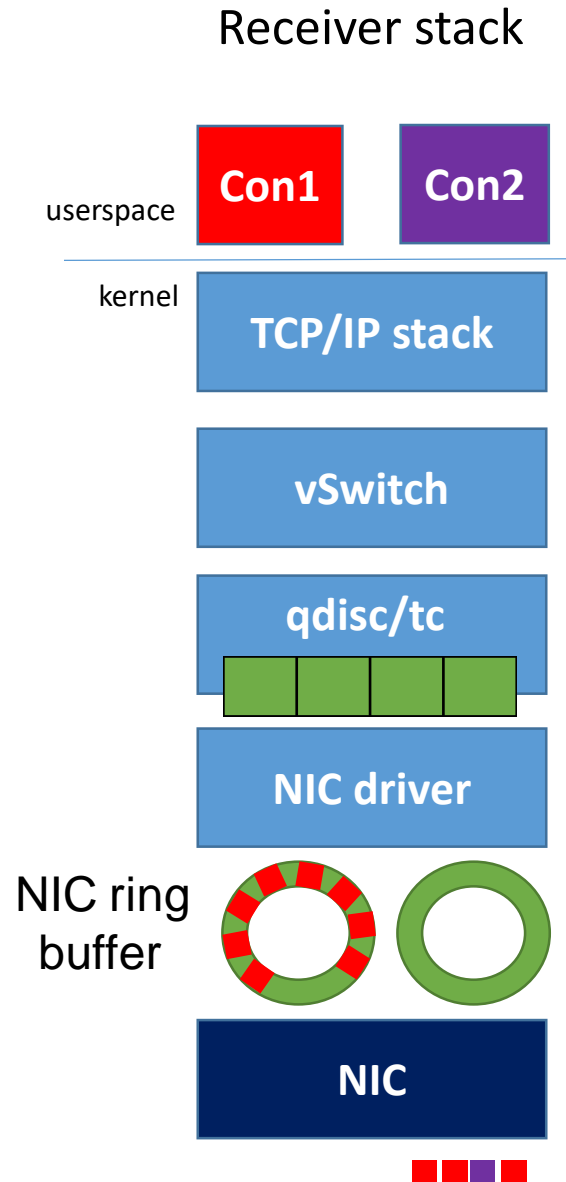
- How and by how much is isolation broken
- Iron's Design
 - Accounting of per-packet processing cost
 - **Ensuring isolation via enforcement**
 - Integration with Linux scheduler
 - **Hardware-based packet dropping**
- Evaluation
 - Controlled workload
 - Realistic workload

Iron – Enforcement

Dropping Packets

Modifies the Linux's polling mechanism (NAPI)

- Assigns a queue (ring buffer) to a container

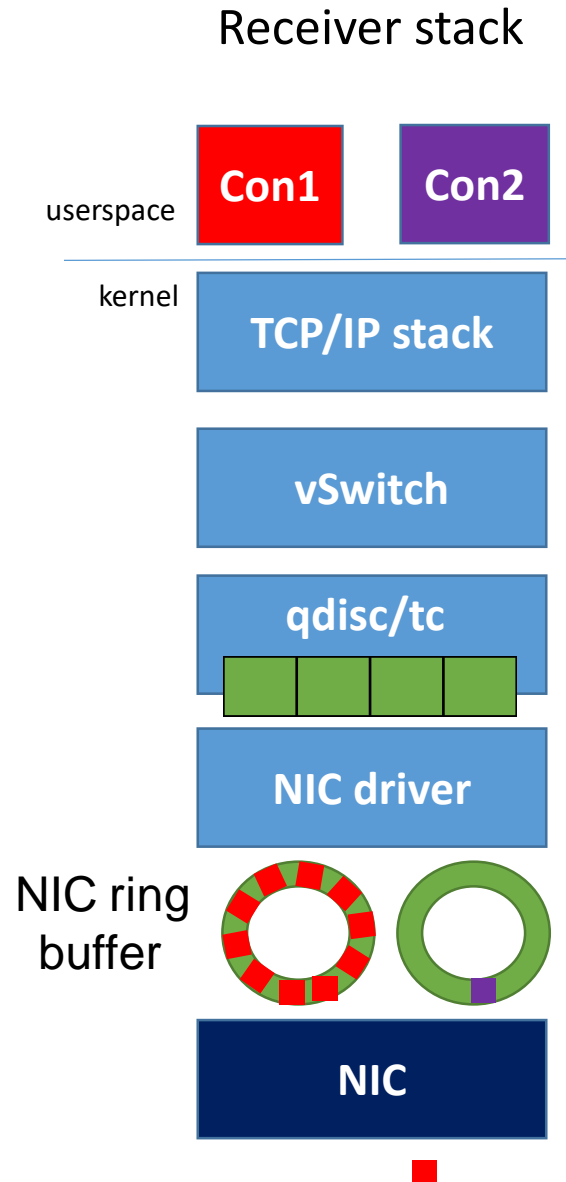


Iron – Enforcement

Dropping Packets

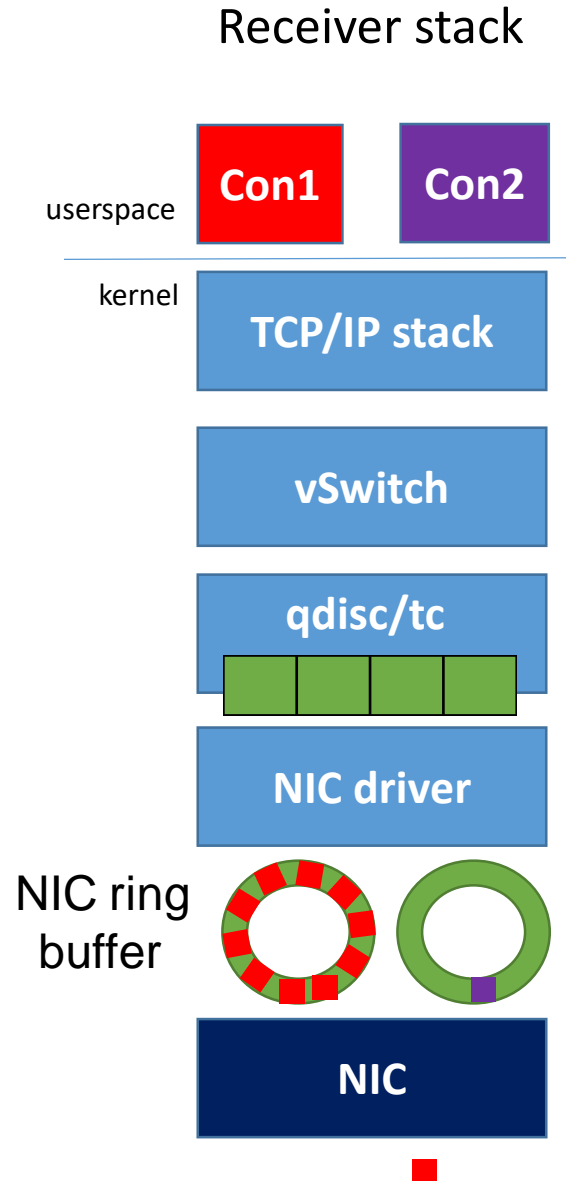
Modifies the Linux's polling mechanism (NAPI)

- Assigns a queue (ring buffer) to a container
- Iron strips the throttled queue from the polling list



Iron – Enforcement

Dropping Packets

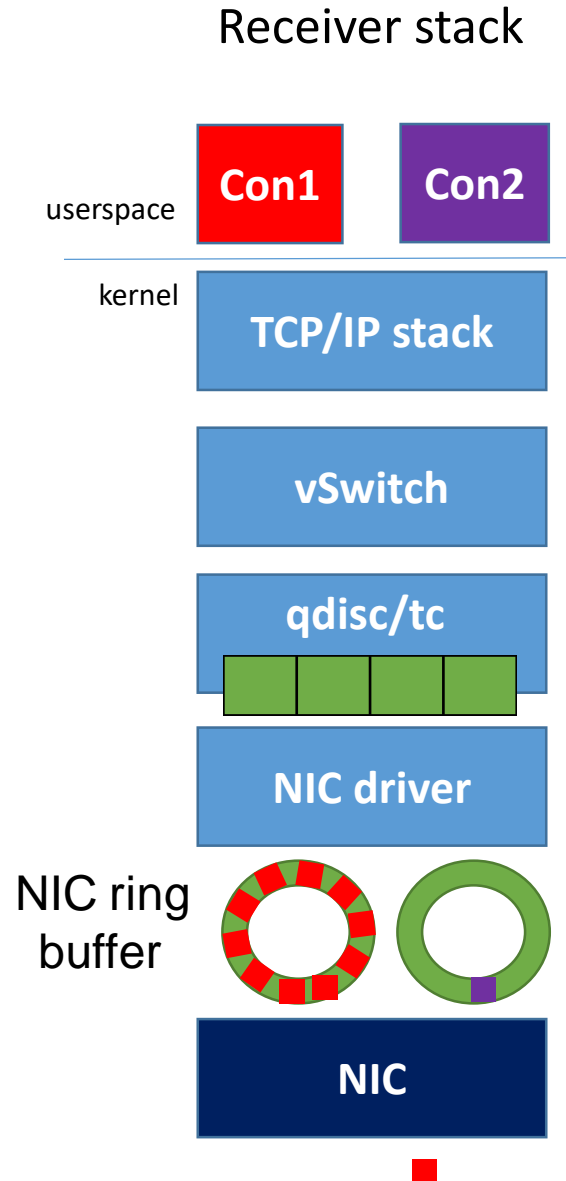


Modifies the Linux's polling mechanism (NAPI)

- Assigns a queue (ring buffer) to a container
- Iron strips the throttled queue from the polling list
 - From kernel's point of view, there are no more interrupts – no packets

Iron – Enforcement

Dropping Packets



Modifies the Linux's polling mechanism (NAPI)

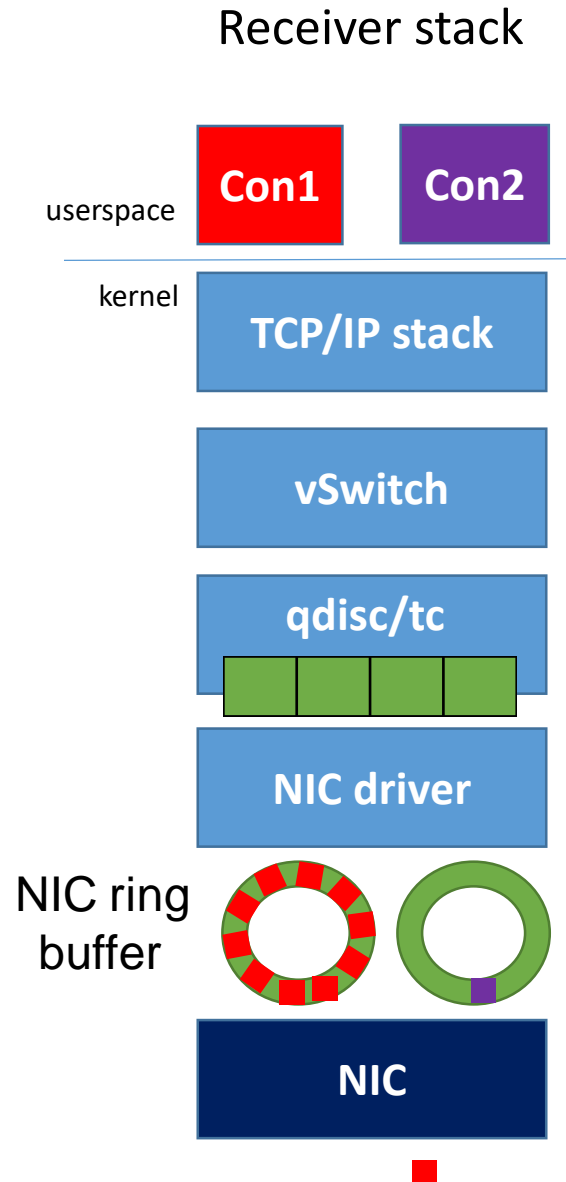
- Assigns a queue (ring buffer) to a container
- Iron strips the throttled queue from the polling list
 - From kernel's point of view, there are no more interrupts – no packets
 - From NIC's point of view, kernel is busy and is not polling packets from the queue, so it stays in the polling mode.

Iron – Enforcement

Dropping Packets

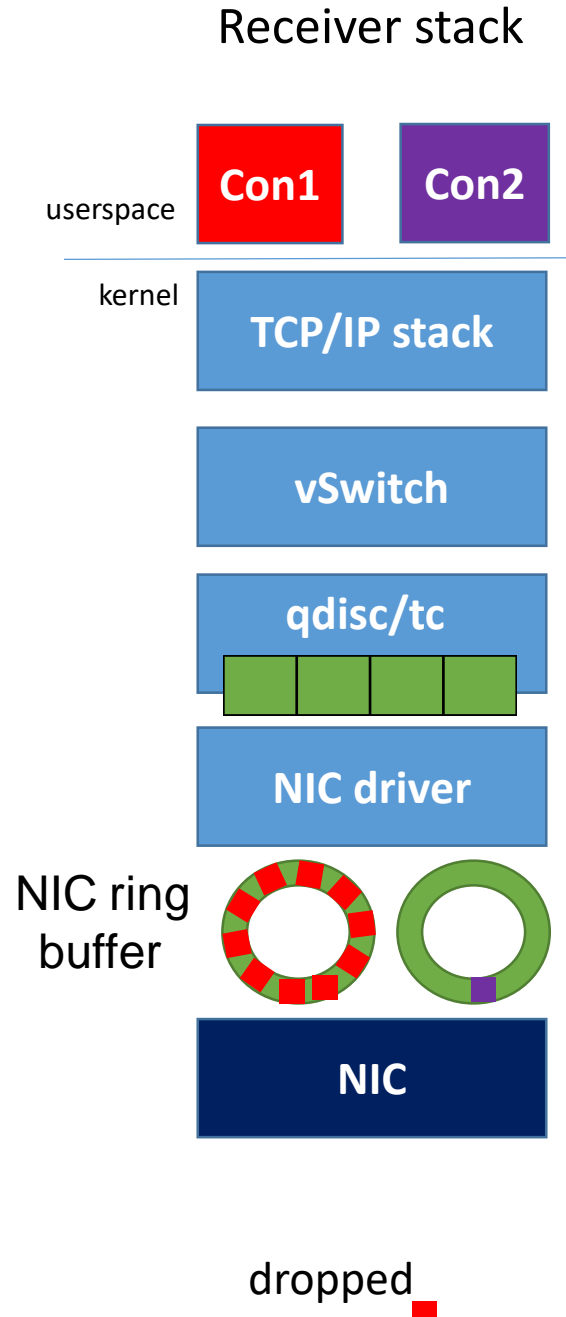
Modifies the Linux's polling mechanism (NAPI)

- Assigns a queue (ring buffer) to a container
- Iron strips the throttled queue from the polling list
 - From kernel's point of view, there are no more interrupts – no packets
 - From NIC's point of view, kernel is busy and is not polling packets from the queue, so it stays in the polling mode.
 - If a new packet arrives and the ring buffer is full, it gets dropped



Iron – Enforcement

Dropping Packets



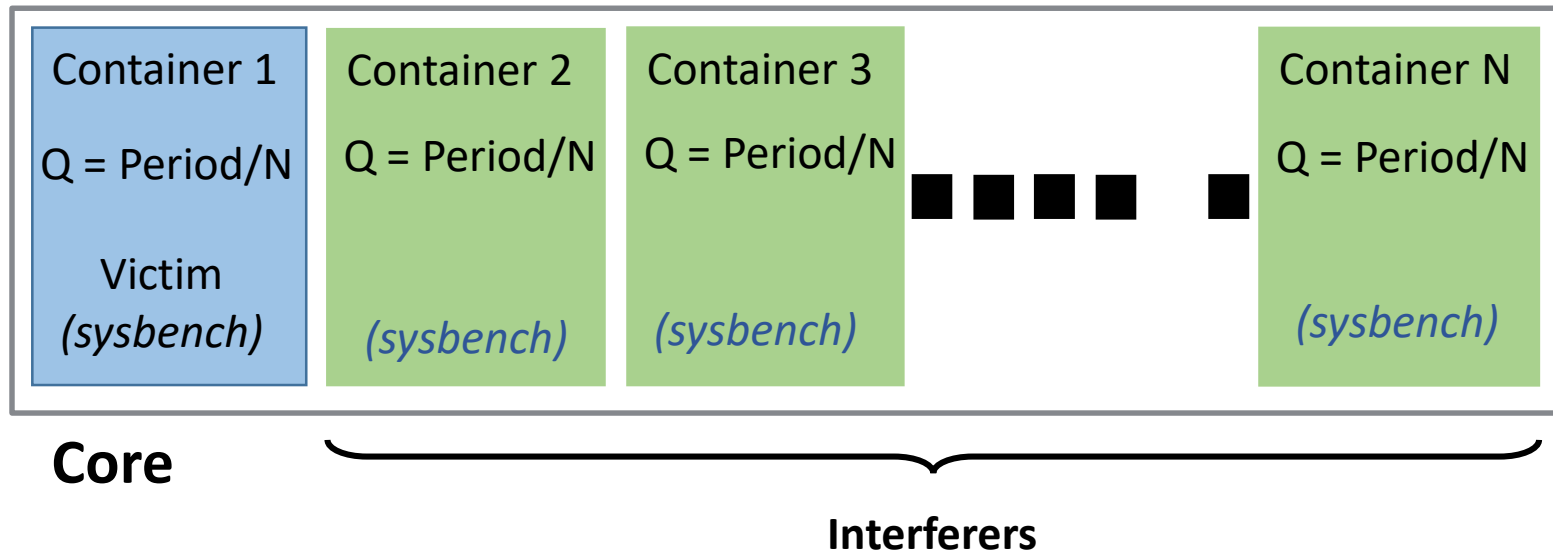
Modifies the Linux's polling mechanism (NAPI)

- Assigns a queue (ring buffer) to a container
- Iron strips the throttled queue from the polling list
 - From kernel's point of view, there are no more interrupts – no packets
 - From NIC's point of view, kernel is busy and is not polling packets from the queue, so it stays in the polling mode.
 - If a new packet arrives and the ring buffer is full, it gets dropped

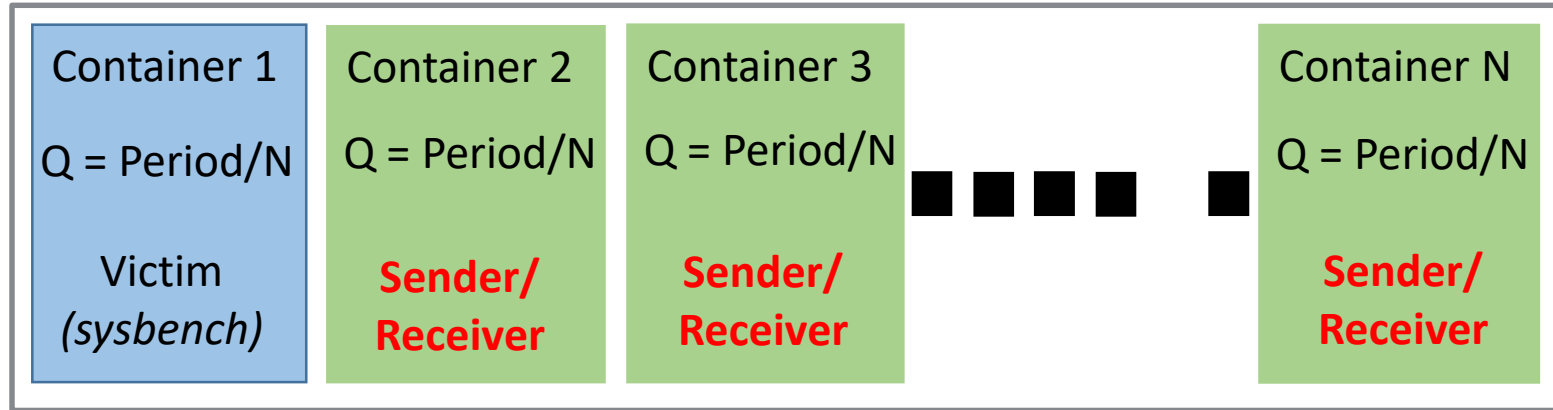
Outline

- How and by how much is isolation broken
- Iron's design
 - Accounting of per-packet processing cost
 - Ensuring isolation via enforcement
 - Integration with Linux scheduler
 - Hardware-based packet dropping
- **Evaluation**
 - **Controlled workload**
 - Realistic workload

Experiment Setup

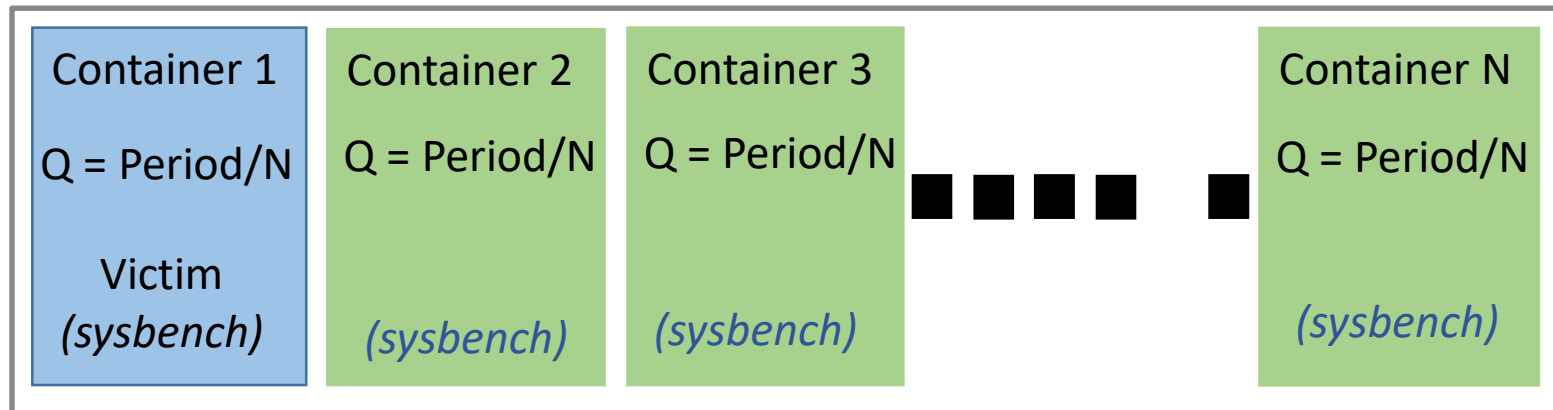


Experiment Setup



Core

Interferers

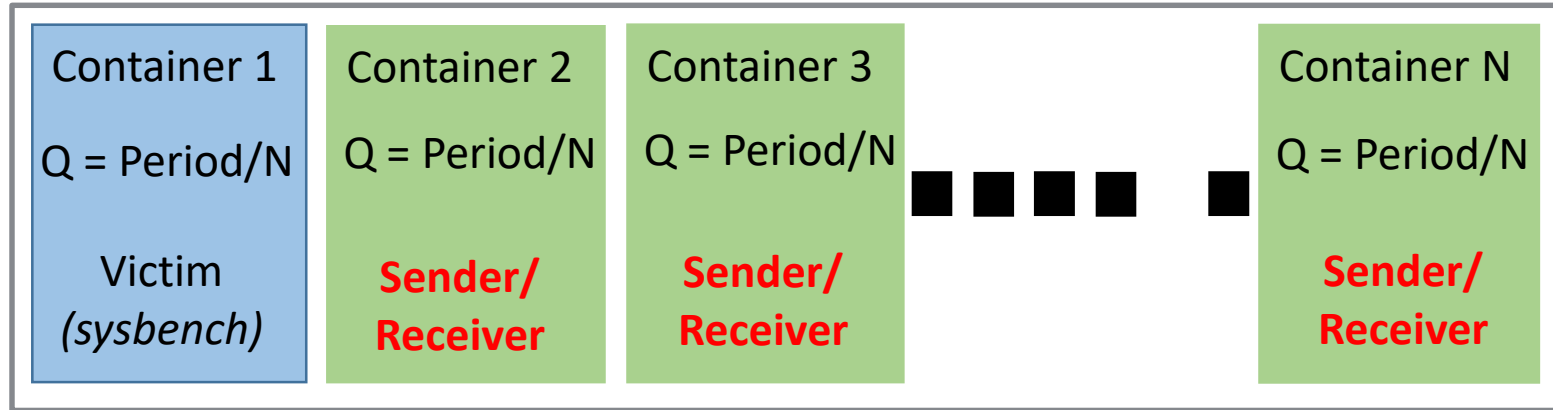


Core

Interferers

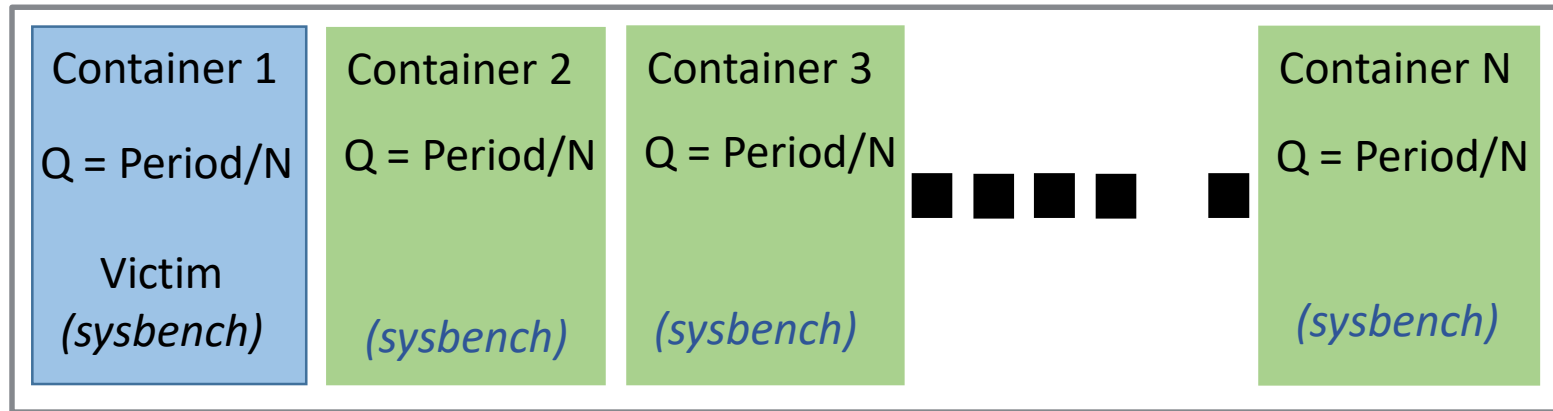
Experiment Setup

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



Core

Interferers

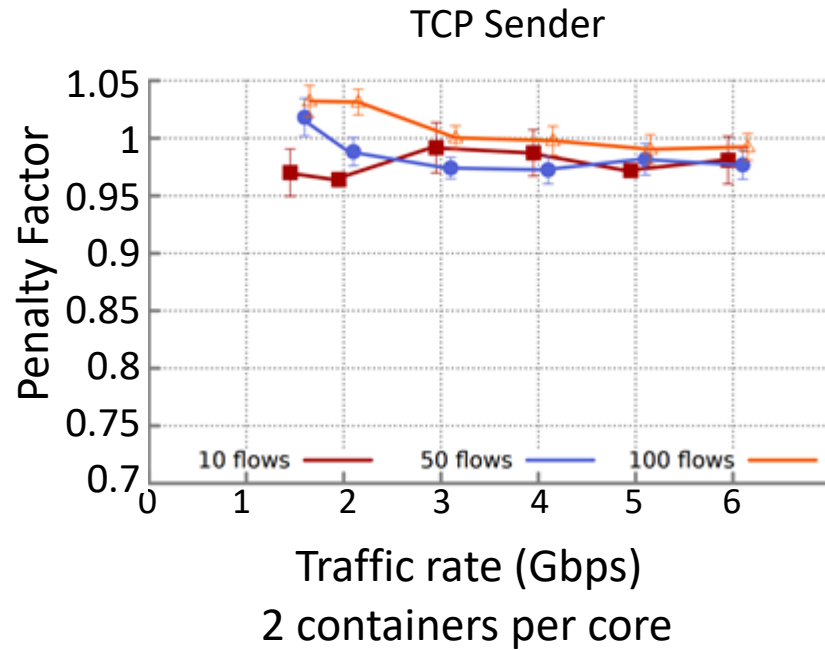


Core

Interferers

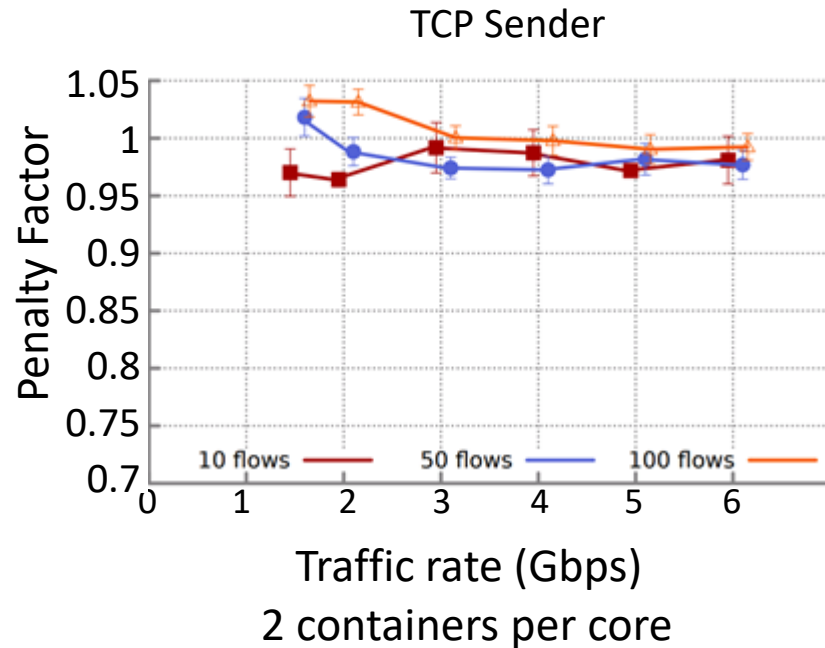
Sender Interference With Iron

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



Sender Interference With Iron

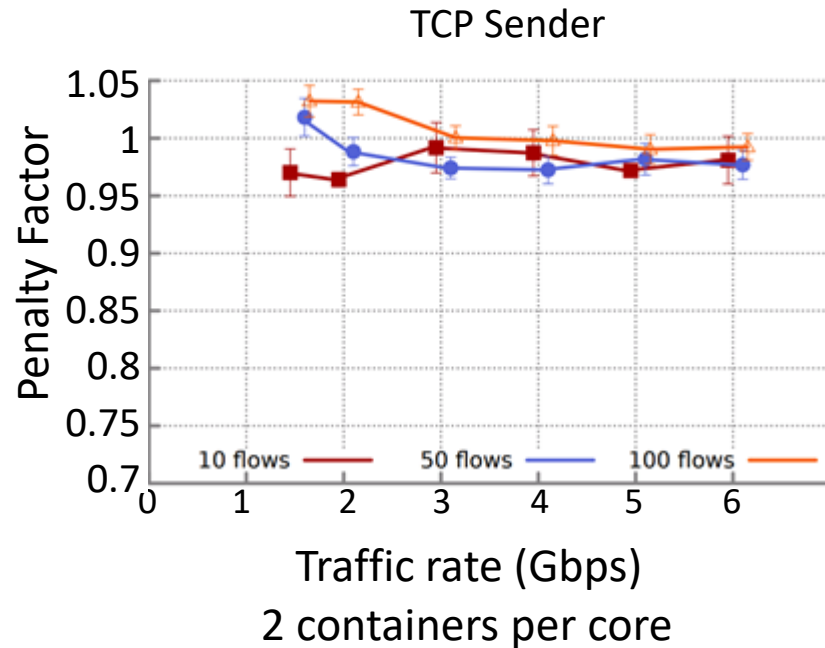
$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



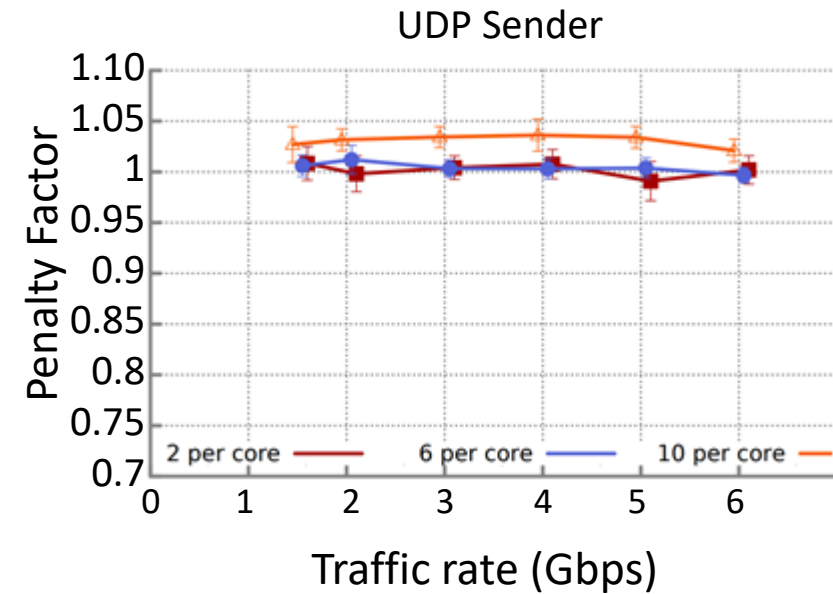
Penalty factor remains below **1.04**,
significant decrease from **1.85**

Sender Interference With Iron

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



Penalty factor remains below **1.04**,
significant decrease from **1.85**



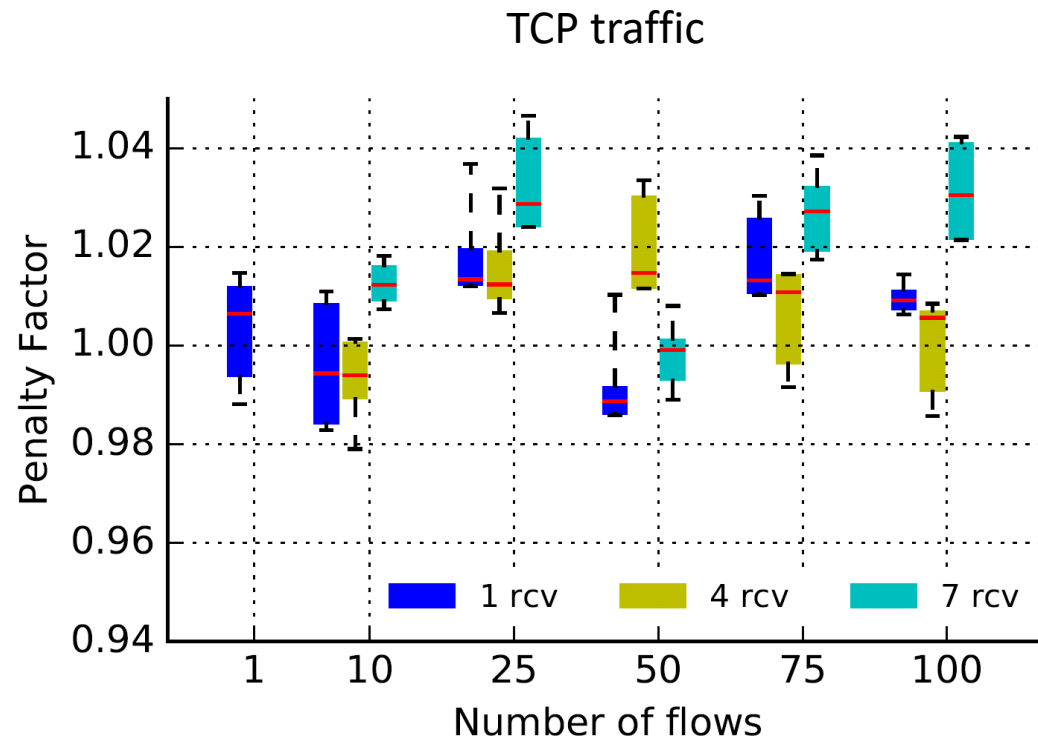
Penalty factor remains below **1.04**,
significant decrease from **1.18**

Receiver Interference With Iron

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$

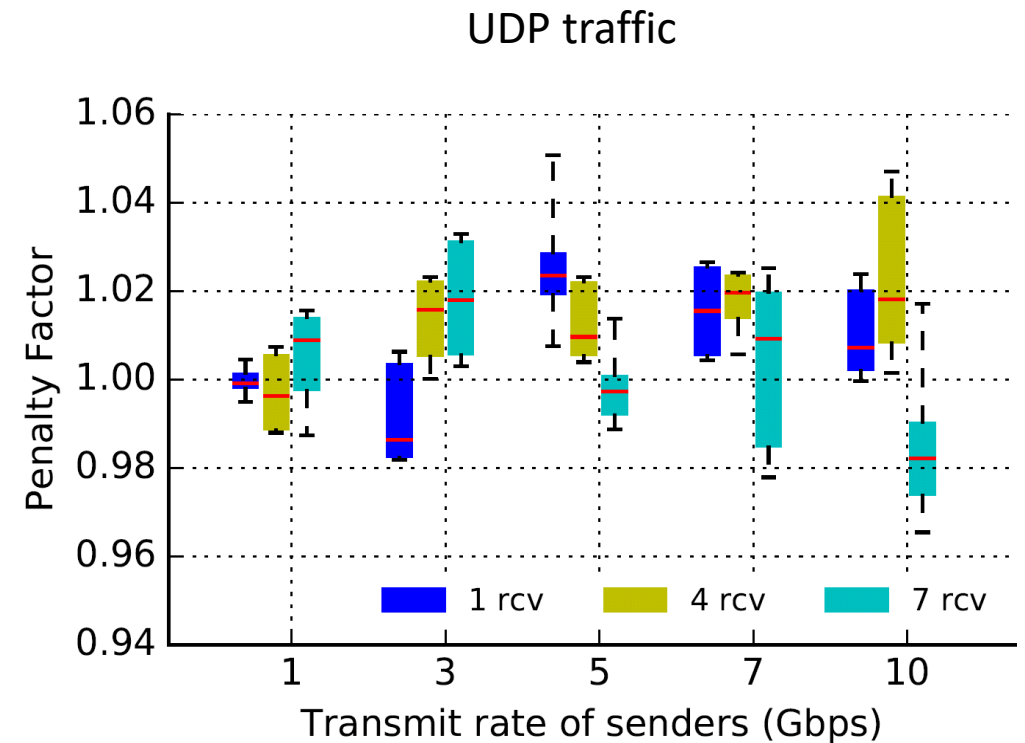
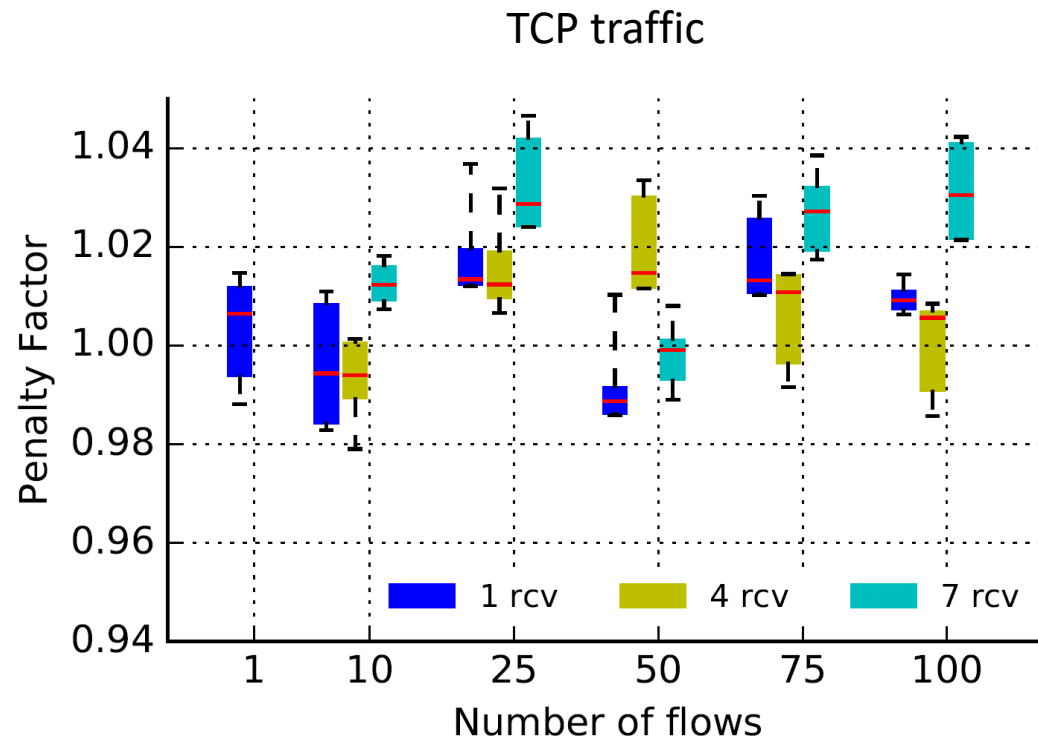
Receiver Interference With Iron

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



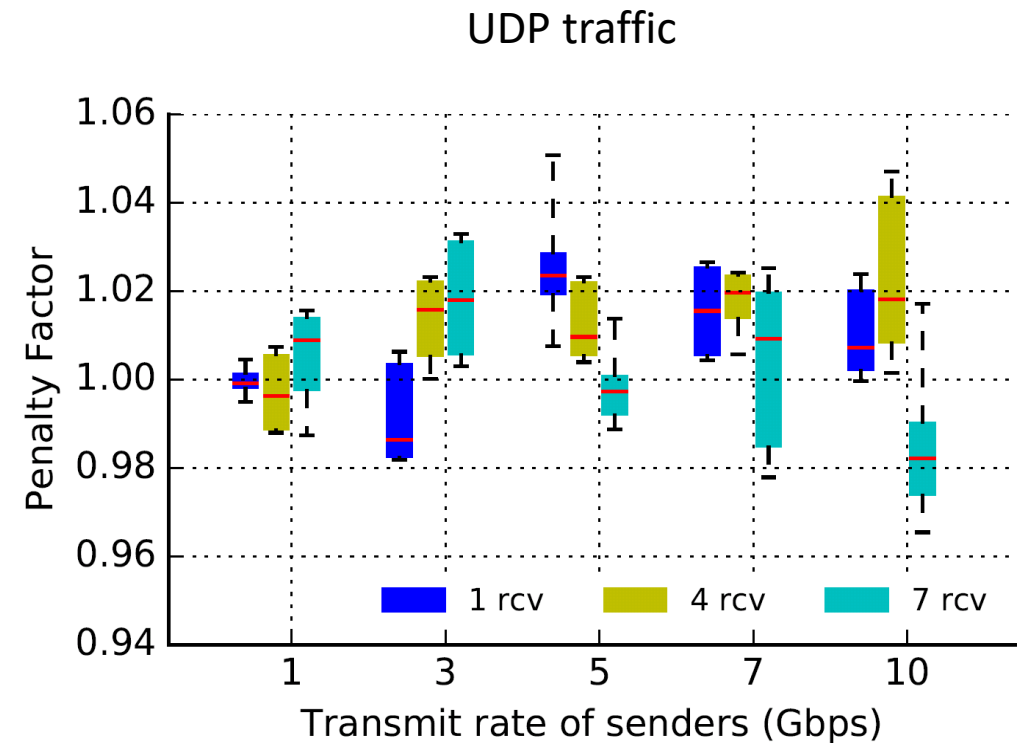
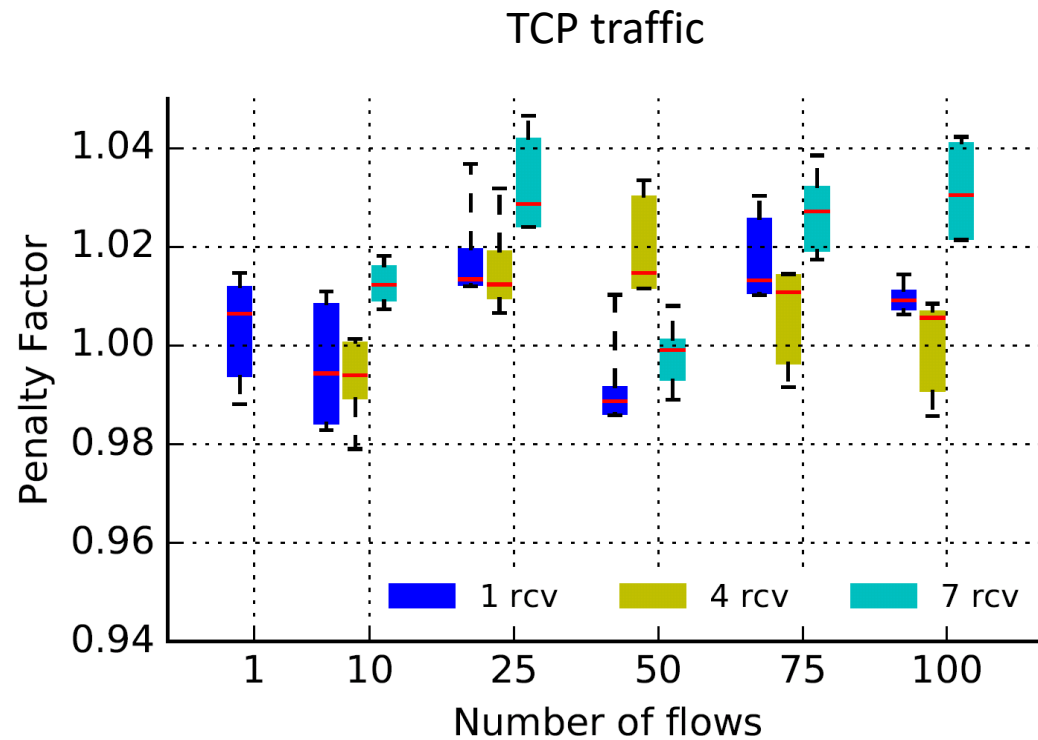
Receiver Interference With Iron

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



Receiver Interference With Iron

$$\text{Penalty factor} = \frac{\text{Time that victim takes when competing with traffic}}{\text{Time that victim takes when competing with sysbench}}$$



Penalty factor never exceeds **1.05**, significant decrease from **6 for TCP** and **4.45 for UDP**

Outline

- How and by how much is isolation broken
- Iron's Design
 - Accounting of per-packet processing cost
 - Ensuring isolation via enforcement
 - Integration with Linux scheduler
 - Hardware-based packet dropping
- **Evaluation**
 - Controlled workload
 - **Realistic workload**

Impact on Big Data Applications

Setup

- 48 containers spread over 6 machines
- Each job runs over 24 containers

Impact on Big Data Applications

Setup

- 48 containers spread over 6 machines
- Each job runs over 24 containers

MapReduce jobs as victim:

- **wordcount**: counts word frequency
- **pi**: computes the value of pi
- **grep**: searches for a given word

Trace based Interferer:

- Shuffle phase of TeraSort job with 115GB input file

Impact on Big Data Applications

Setup

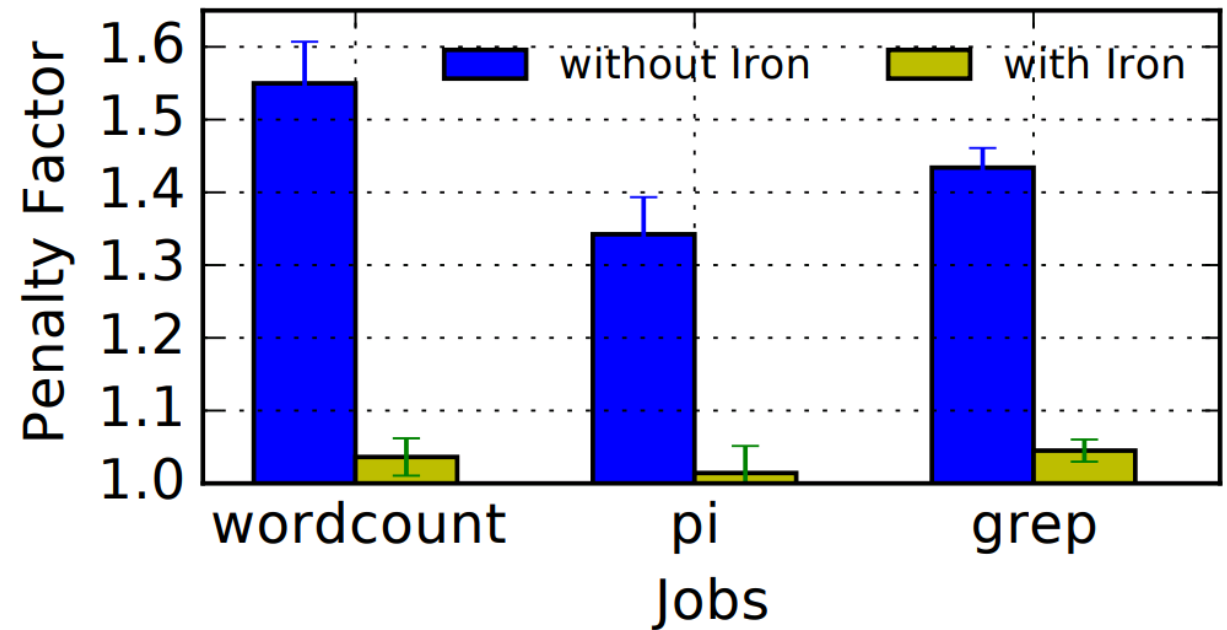
- 48 containers spread over 6 machines
- Each job runs over 24 containers

MapReduce jobs as victim:

- **wordcount**: counts word frequency
- **pi**: computes the value of pi
- **grep**: searches for a given word

Trace based Interferer:

- Shuffle phase of TeraSort job with 115GB input file



Impact on Big Data Applications

Setup

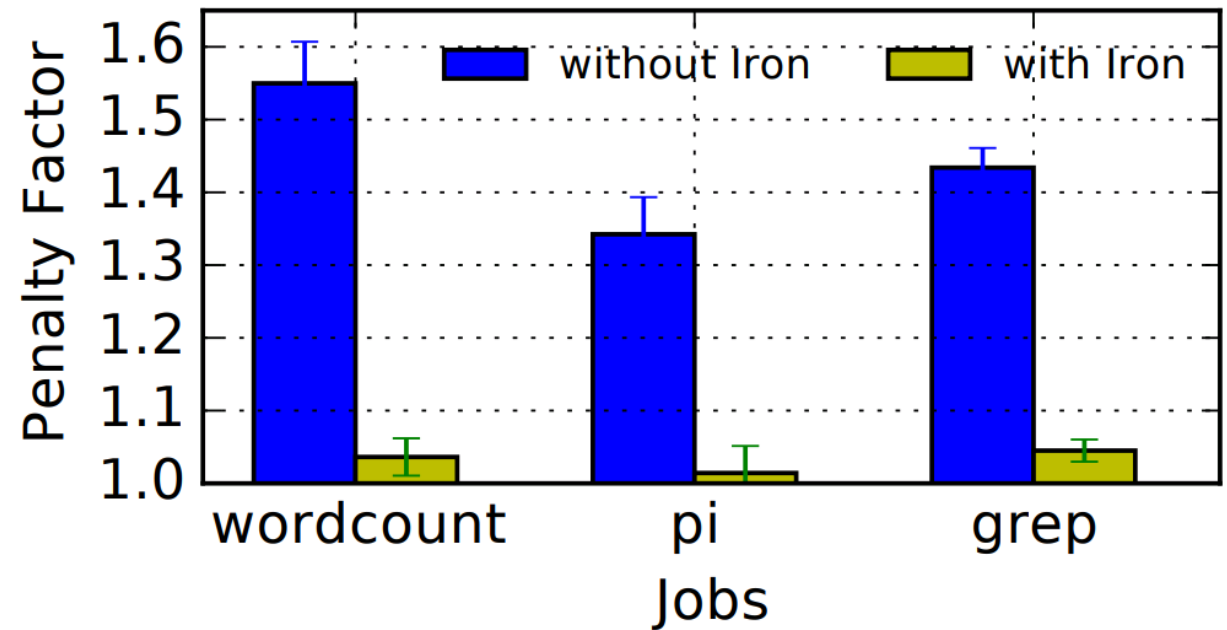
- 48 containers spread over 6 machines
- Each job runs over 24 containers

MapReduce jobs as victim:

- **wordcount**: counts word frequency
- **pi**: computes the value of pi
- **grep**: searches for a given word

Trace based Interferer:

- Shuffle phase of TeraSort job with 115GB input file



Penalty factor never exceeds **1.04**

Summary

- Evaluated the interference caused by network-based containers.
- Provided hardened isolation for network-based processing in containerized environment.
- Ensures accurate accounting of the time spent processing network traffic in softirq.
- Integrated with Linux scheduler with minimal changes.
- Novel packet dropping mechanism to limit the interference.