# iip: an integratable TCP/IP stack
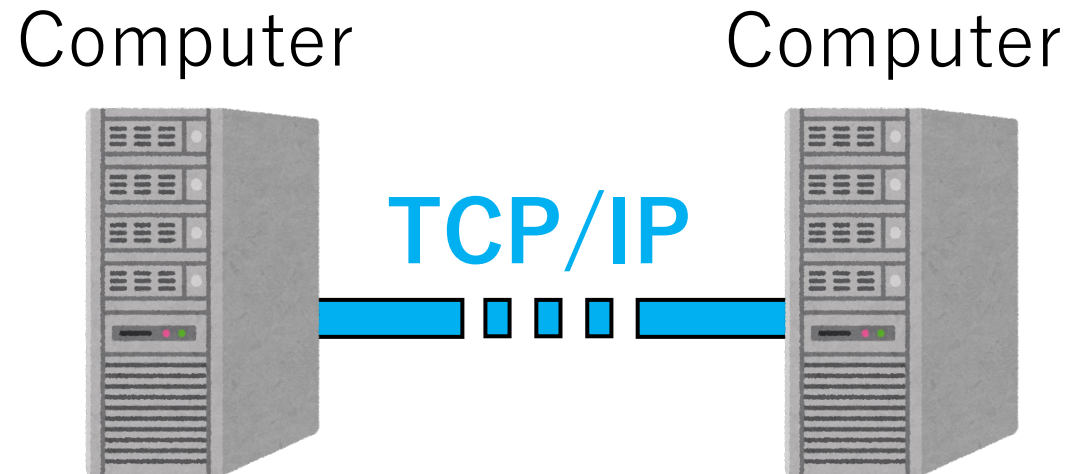
Kenichi Yasukata

Internet Initiative Japan

# TCP/IP and TCP/IP Stacks
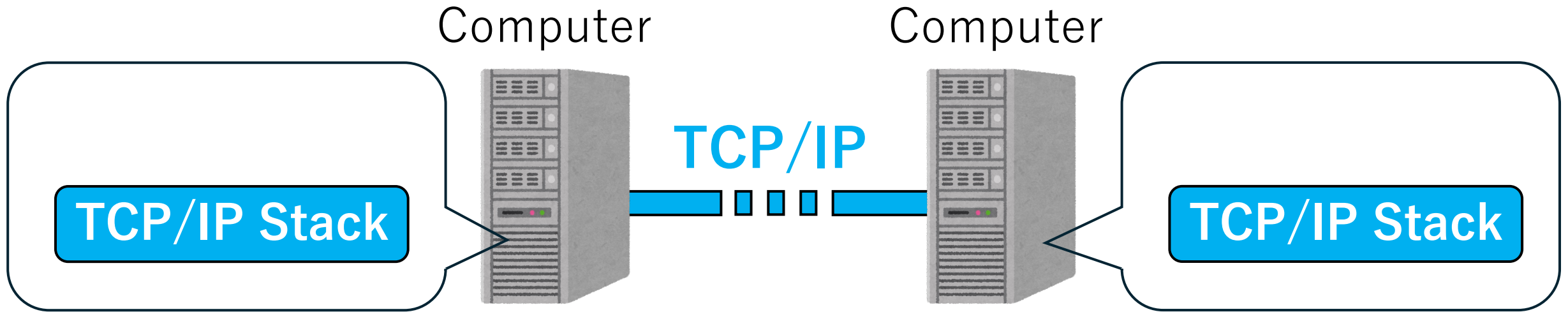
- TCP/IP is a standardized protocol suite commonly used for communication in computer networks

Computer                    Computer
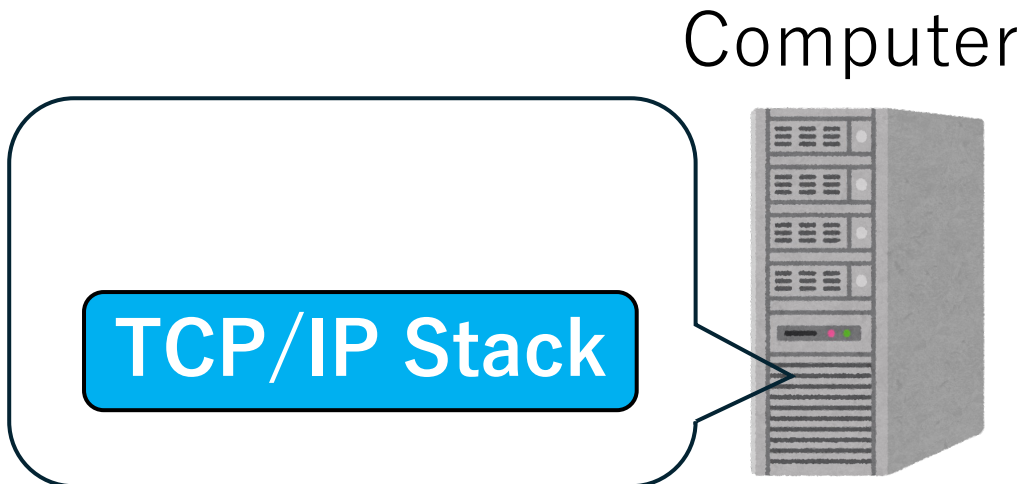
**TCP/IP**

https://github.com/yasukata/iip

# TCP/IP and TCP/IP Stacks

- TCP/IP is a standardized protocol suite commonly used for communication in computer networks
- TCP/IP stacks are typically software that implements procedures to comply with the TCP/IP standard

Computer          Computer

**TCP/IP**

**TCP/IP Stack**          **TCP/IP Stack**

https://github.com/yasukata/iip

# TCP/IP Stacks in Legacy OS Kernels

- TCP/IP stacks are typically maintained as part of OS kernels
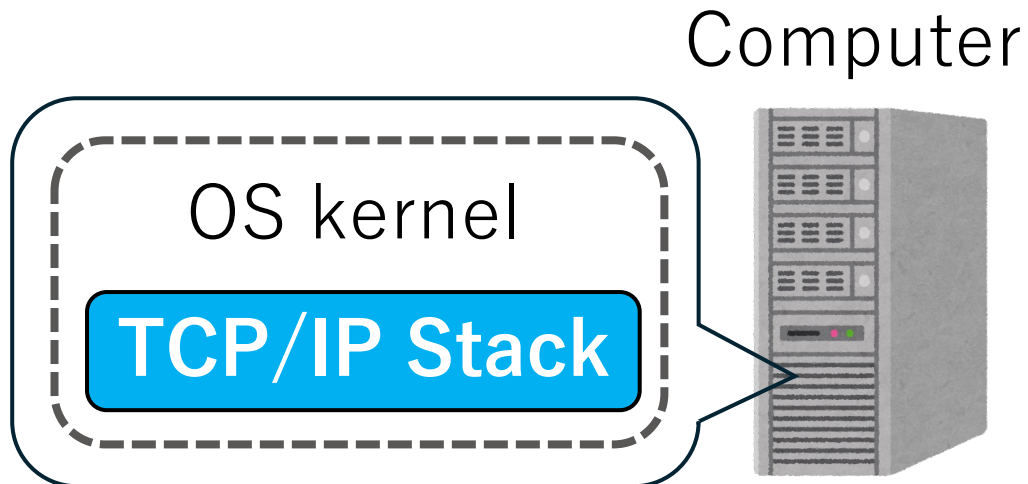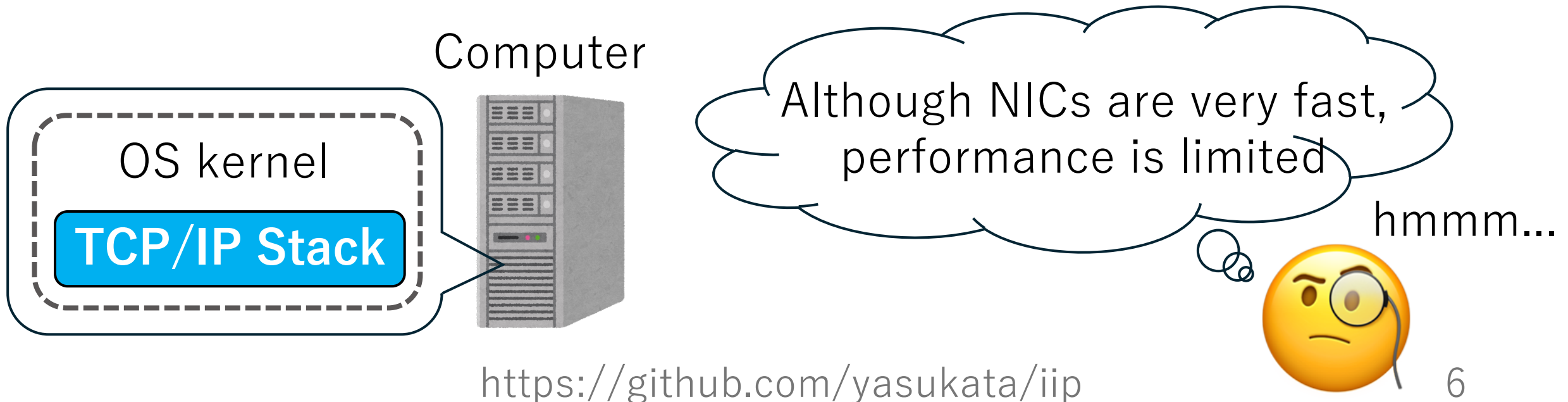
Computer

**TCP/IP Stack**

https://github.com/yasukata/iip

# TCP/IP Stacks in Legacy OS Kernels

- TCP/IP stacks are typically maintained as part of OS kernels

Computer

OS kernel

**TCP/IP Stack**

https://github.com/yasukata/iip

# TCP/IP Stacks in Legacy OS Kernels

- TCP/IP stacks are typically maintained as part of OS kernels
- People found it is hard for TCP/IP stacks in legacy OS kernels to effectively utilize the benefits of high-speed NICs

Computer

OS kernel

**TCP/IP Stack**

Although NICs are very fast, performance is limited

hmmm...

https://github.com/yasukata/iip

6

# Performance-optimized TCP/IP Stacks

- To address this issue, research and industry communities have invented many <u>performance-optimized TCP/IP stacks</u>
  - e.g., Sandstorm (SIGCOMM'14), mTCP (NSDI'14)

Computer

**Performance-optimized TCP/IP Stack**

Power of NICs is unleashed by fast TCP/IP stacks!

https://github.com/yasukata/iip

# Performance-optimized TCP/IP Stacks

😄 Their performance is excellent
e.g., Sandstorm (SIGCOMM'14), mTCP (NSDI'14)

# Performance-optimized TCP/IP Stacks

😄 Their performance is excellent
  e.g., Sandstorm (SIGCOMM'14), mTCP (NSDI'14)

😢 They often incur high integration complexity

| Category | Performance | Integration |
|---|---|---|
| Performance-optimized | ✔ | |

# Portability-aware TCP/IP Stacks

😄 There are TCP/IP stacks that allow for easy integration
  e.g., lwIP, FNET, picoTCP

| Category | Performance | Integration |
|---|:---:|:---:|
| Performance-optimized | ✓ | |
| Portability-aware | | ✓ |

https://github.com/yasukata/iip

# Portability-aware TCP/IP Stacks

😄 There are TCP/IP stacks that allow for easy integration
   e.g., lwIP, FNET, picoTCP

😢 They often lack the care for performance-critical factors

| Category | Performance | Integration |
|---|---|---|
| Performance-optimized | ✓ | |
| Portability-aware | | ✓ |

https://github.com/yasukata/iip

# Problem

- None of previous TCP/IP stack implementations allow for <u>easy integration</u> and <u>good performance</u> simultaneously

| Category | Performance | Integration |
|---|---|---|
| Performance-optimized | ✔ | |
| Portability-aware | | ✔ |

# Problem

- As a result, developers only had limited and laborious options

| Category | Performance | Integration |
|---|---|---|
| Performance-optimized | ✓ | |
| Portability-aware | | ✓ |

# Problem

- As a result, developers only had limited and laborious options
  - intensively modifying one of the existing TCP/IP stacks

| Category | Performance | Integration |
|---|:---:|:---:|
| Performance-optimized | ✔ | |
| Portability-aware | | ✔ |

https://github.com/yasukata/iip

# Problem

- As a result, developers only had limited and laborious options
  - intensively modifying one of the existing TCP/IP stacks
  - building a new TCP/IP stack from scratch

| Category | Performance | Integration |
|---|---|---|
| Performance-optimized | ✓ | |
| Portability-aware | | ✓ |

https://github.com/yasukata/iip

# Problem

- As a result, developers only had limited and laborious options
  - intensively modifying one of the existing TCP/IP stacks
  - building a new TCP/IP stack from scratch
  - accepting performance limitations of an applicable TCP/IP stack

| Category | Performance | Integration |
|---|---|---|
| Performance-optimized | ✔ | |
| Portability-aware | | ✔ |

https://github.com/yasukata/iip

# Problem

- As a result, developers only had limited and laborious options
  - intensively modifying one of the existing TCP/IP stacks
  - building a new TCP/IP stack from scratch
  - accepting performance limitations of an applicable TCP/IP stack
  - giving up the integration

| Category | Performance | Integration |
|---|---|---|
| Performance-optimized | ✔ | |
| Portability-aware | | ✔ |

https://github.com/yasukata/iip

# This Work

- We develop iip, an integratable TCP/IP stack, that allows for <u>easy integration</u> and <u>good performance</u> simultaneously

| Category | Performance | Integration |
|---|:---:|:---:|
| Performance-optimized | ✔ | |
| Portability-aware | | ✔ |
| This work | ✔ | ✔ |

https://github.com/yasukata/iip

# Issues of Existing TCP/IP Stacks

- Performance-optimized TCP/IP stacks
  - Dependencies on other components
  - Functionality conflicts
  - Limited choices for CPU core assignment models


- Portability-aware TCP/IP stacks
  - Unaware of NIC offloading features
  - Lack of zero-copy I/O capability
  - Lack of multi-core scalability

https://github.com/yasukata/iip

# Issues of Existing TCP/IP Stacks

Performance-optimized
TCP/IP stack

Performance-optimized
TCP/IP stacks often consist of
various components

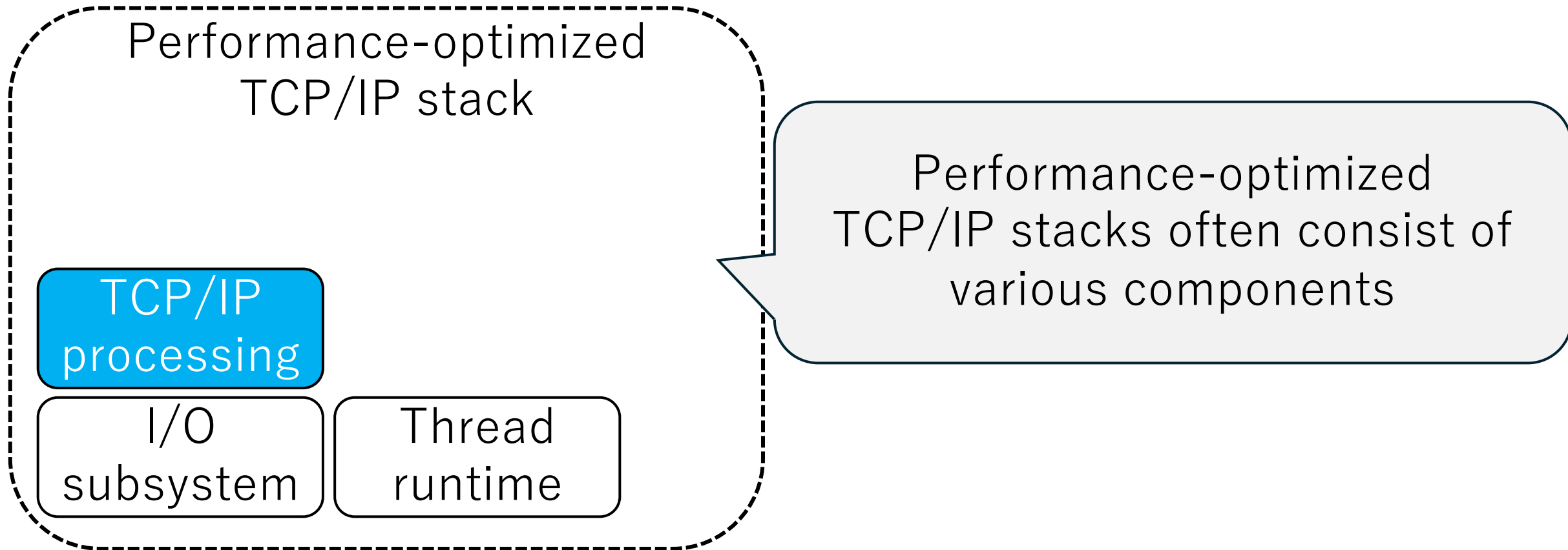# Issues of Existing TCP/IP Stacks

Performance-optimized
TCP/IP stack

TCP/IP processing

Performance-optimized
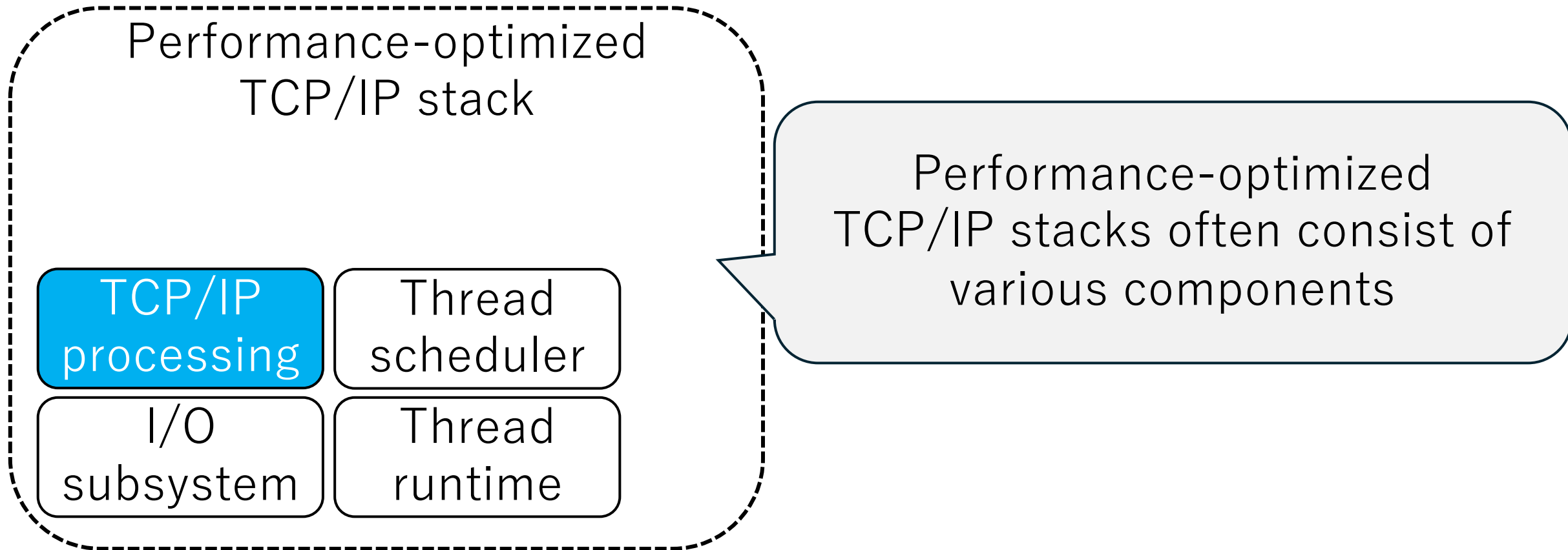TCP/IP stacks often consist of
various components

https://github.com/yasukata/iip
21

# Issues of Existing TCP/IP Stacks

Performance-optimized
TCP/IP stack

TCP/IP
processing

I/O
subsystem

Performance-optimized
TCP/IP stacks often consist of
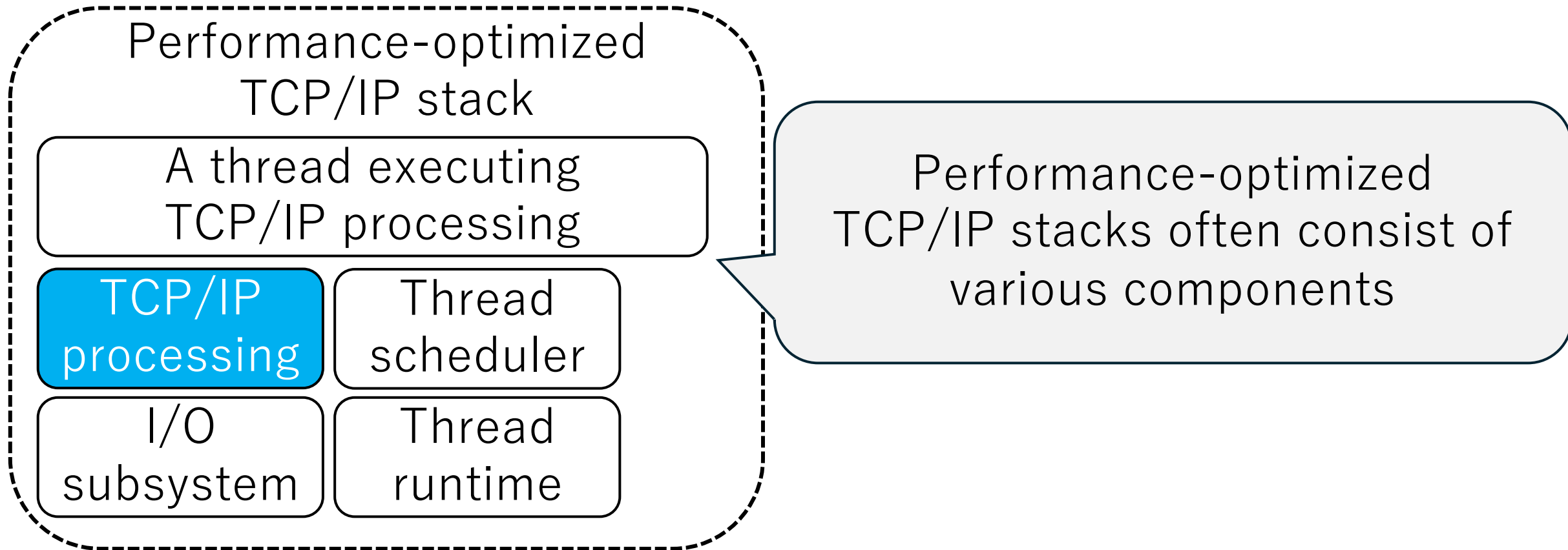various components

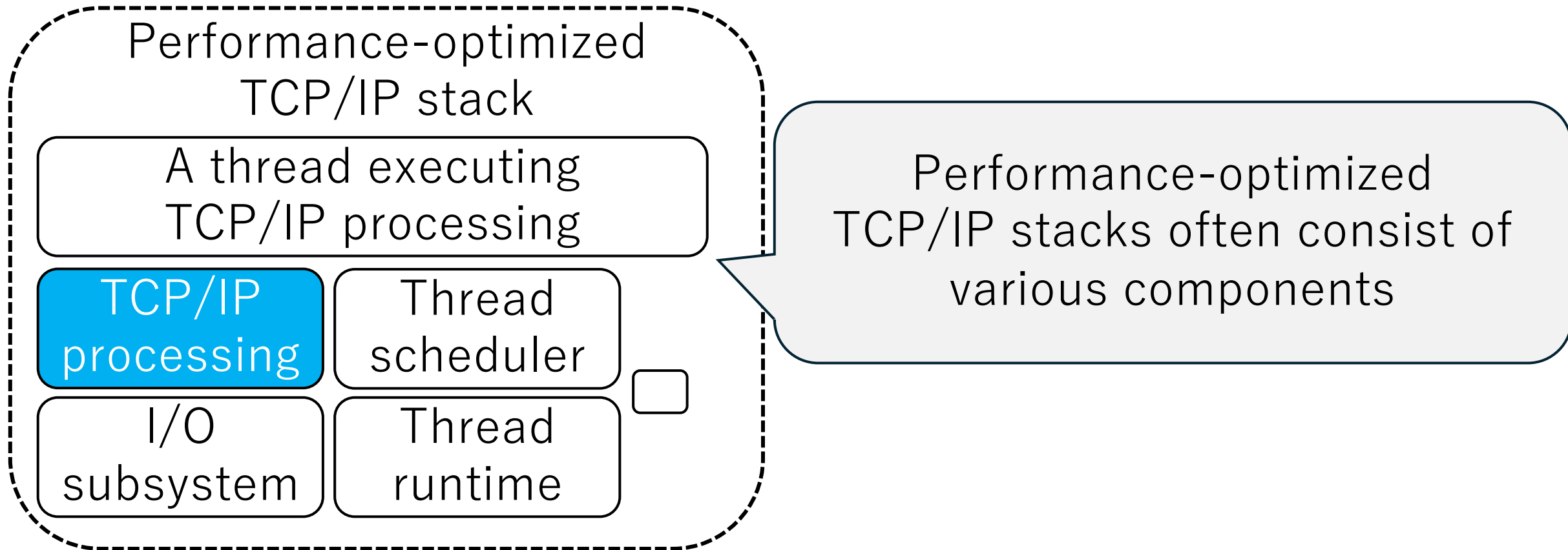# Issues of Existing TCP/IP Stacks

Performance-optimized
TCP/IP stack

TCP/IP
processing

I/O
subsystem

Thread
runtime

Performance-optimized
TCP/IP stacks often consist of
various components

# Issues of Existing TCP/IP Stacks

Performance-optimized
TCP/IP stack

| TCP/IP processing | Thread scheduler |
|---|---|
| I/O subsystem | Thread runtime |

Performance-optimized
TCP/IP stacks often consist of
various components

# Issues of Existing TCP/IP Stacks

Performance-optimized
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

Performance-optimized
TCP/IP stacks often consist of
various components

# Issues of Existing TCP/IP Stacks

Performance-optimized
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

Performance-optimized
TCP/IP stacks often consist of
various components

# Issues of Existing TCP/IP Stacks

Performance-optimized
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

Performance-optimized
TCP/IP stacks often consist of
various components

# Issues of Existing TCP/IP Stacks

## Dependencies on Other Components

Performance-optimized
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

# Issues of Existing TCP/IP Stacks

## Dependencies on Other Components

The I/O subsystem depends on Library-A

TCP/IP processing

| TCP/IP processing | Thread scheduler |
| I/O subsystem | Thread runtime |

Library-A

# Issues of Existing TCP/IP Stacks

## Dependencies on Other Components

Library-A depends on
OS-B version X

TCP/IP processing

**TCP/IP processing**

Thread scheduler

I/O subsystem

Thread runtime

Library-A

OS-B version X

# Issues of Existing TCP/IP Stacks

## Dependencies on Other Components

Performance-optimized
TCP/IP stack

Implemented by TCP/IP stack users
(application developers)

A thread executing
TCP/IP processing

Application
logic

TCP/IP
processing

Thread
scheduler

Library-A

I/O
subsystem

Thread
runtime

OS-B
version X

https://github.com/yasukata/iip

31

# Issues of Existing TCP/IP Stacks

## Dependencies on Other Components

Implemented by TCP/IP stack users
(application developers)

Developers wish to integrate this TCP/IP stack with their app

TCP/IP processing

| TCP/IP processing | Thread scheduler |
| I/O subsystem | Thread runtime |

Application logic

Library-A

OS-B version X

# Issues of Existing TCP/IP Stacks

## Dependencies on Other Components

This app depends on
Library-C

Implemented by TCP/IP stack users
(application developers)

TCP/IP processing

| TCP/IP processing | Thread scheduler |
| I/O subsystem | Thread runtime |

Application logic

Library-A

Library-C

OS-B version X

# Issues of Existing TCP/IP Stacks

## Dependencies on Other Components

Implemented by TCP/IP stack users
(application developers)

Library-C depends on
OS-B version Y

TCP/IP processing

**TCP/IP processing**

Thread scheduler

I/O subsystem

Thread runtime

Application logic

Library-A

Library-C

OS-B version X

OS-B version Y

# Issues of Existing TCP/IP Stacks

## Dependencies on Other Components

Normally, two different versions of an OS cannot coexist

TCP/IP processing

TCP/IP processing

Thread scheduler

I/O subsystem

Thread runtime

Implemented by TCP/IP stack users
(application developers)

Application logic

Library-A

Library-C

**OS-B version X**

**OS-B version Y**

# Issues of Existing TCP/IP Stacks

## Dependencies on Other Components

As a result, it is hard to use this TCP/IP stack for this app

Implemented by TCP/IP stack users
(application developers)

TCP/IP processing

TCP/IP processing

Thread scheduler

I/O subsystem

Thread runtime

Application logic

Library-A

Library-C

**OS-B version X**

**OS-B version Y**

https://github.com/yasukata/iip

# Issues of Existing TCP/IP Stacks

## Dependencies on Other Components

**Dependencies often increase integration complexity**

Implemented by TCP/IP stack users
(application developers)

TCP/IP processing

| TCP/IP processing | Thread scheduler |
| --- | --- |
| I/O subsystem | Thread runtime |

Application logic

Library-A

Library-C

**OS-B version X**

**OS-B version Y**

https://github.com/yasukata/iip

37

# Issues of Existing TCP/IP Stacks

## Functionality Conflicts

Performance-optimized
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

# Issues of Existing TCP/IP Stacks

## Functionality Conflicts

**Performance-optimized TCP/IP stack**

- A thread executing TCP/IP processing
- TCP/IP processing
- Thread scheduler
- I/O subsystem
- Thread runtime

**A new OS specialized for performance**

- A thread executing TCP/IP processing
- TCP/IP processing
- Thread scheduler
- I/O subsystem
- Thread runtime

# Issues of Existing TCP/IP Stacks

## Functionality Conflicts



Let's say we wish to use this TCP/IP component on this new OS

A new OS specialized for performance

A thread executing TCP/IP processing

TCP/IP processing

Thread scheduler

I/O subsystem

Thread runtime

# Issues of Existing TCP/IP Stacks

## Functionality Conflicts

This TCP/IP component depends on its specific thread runtime

TCP/IP processing

TCP/IP processing

Th

scheduler

I/O subsystem

Thread runtime

A new OS specialized for performance

A thread executing TCP/IP processing

TCP/IP processing

Thread scheduler

I/O subsystem

Thread runtime

https://github.com/yasukata/iip

# Issues of Existing TCP/IP Stacks

## Functionality Conflicts



We also need to port
the thread runtime to this new OS

A new OS specialized
for performance

A thread executing
TCP/IP processing

TCP/IP processing

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

https://github.com/yasukata/iip

42

# Issues of Existing TCP/IP Stacks

## Functionality Conflicts



A new OS specialized for performance

A thread executing TCP/IP processing

However, this new OS has its own thread runtime

TCP/IP processing

TCP/IP processing

Thread scheduler

Thread scheduler

I/O subsystem

I/O subsystem

Thread runtime

Thread runtime

# Issues of Existing TCP/IP Stacks

## Functionality Conflicts

Two independent thread runtimes normally cannot coexist

A new OS specialized for performance

TCP/IP processing

A thread executing TCP/IP processing

TCP/IP processing

Thread scheduler

I/O subsystem

Thread runtime

TCP/IP processing

Thread scheduler

I/O subsystem

Thread runtime

Thread runtime

https://github.com/yasukata/iip

44

# Issues of Existing TCP/IP Stacks

## Functionality Conflicts



Here, the thread runtime functionality is conflicting

A new OS specialized for performance

A thread executing TCP/IP processing

TCP/IP processing

TCP/IP processing

Thread scheduler

Thread scheduler

I/O subsystem

I/O subsystem

Thread runtime

Thread runtime

https://github.com/yasukata/iip

45

# Issues of Existing TCP/IP Stacks

## Functionality Conflicts

A new OS specialized
for performance

A thread executing
TCP/IP processing

As a result, it is hard to use this
TCP/IP component on this new OS

| TCP/IP processing | Thread scheduler |
| I/O subsystem | Thread runtime |

| TCP/IP processing | Thread scheduler |
| I/O subsystem | Thread runtime |

# Issues of Existing TCP/IP Stacks

## Functionality Conflicts

**Functionality conflicts increase integration complexity**

A new OS specialized for performance

A thread executing TCP/IP processing

TCP/IP processing

TCP/IP processing | Thread scheduler

I/O subsystem | Thread runtime

TCP/IP processing | Thread scheduler

I/O subsystem | Thread runtime

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

Performance-optimized
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

Performance-optimized
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP

Thr

TCP/IP stacks often include
threads for TCP/IP processing

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

Performance-optimized
TCP/IP stack

Implemented by TCP/IP stack users
（application developers）

A thread executing
TCP/IP processing

A thread executing
application logic

TCP/IP          Thread

TCP/IP stacks often assume
app logic is run by another thread

# Issues of Existing TCP/IP Stacks

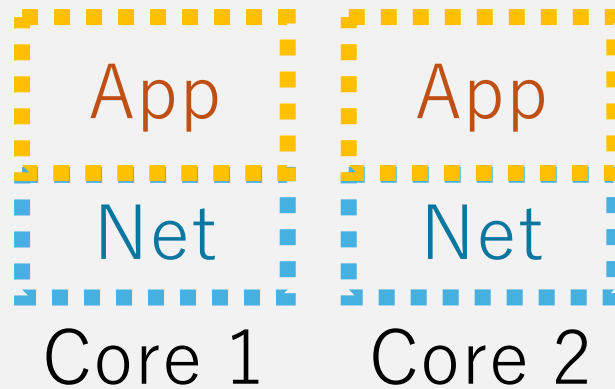## Limited Choices for CPU Core Assignment Models

Three potential CPU core assignment models

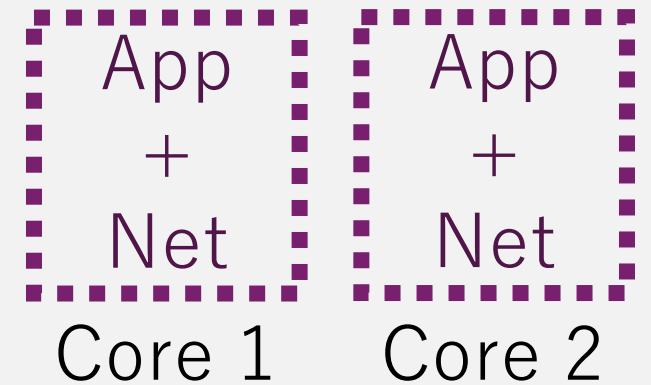# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models
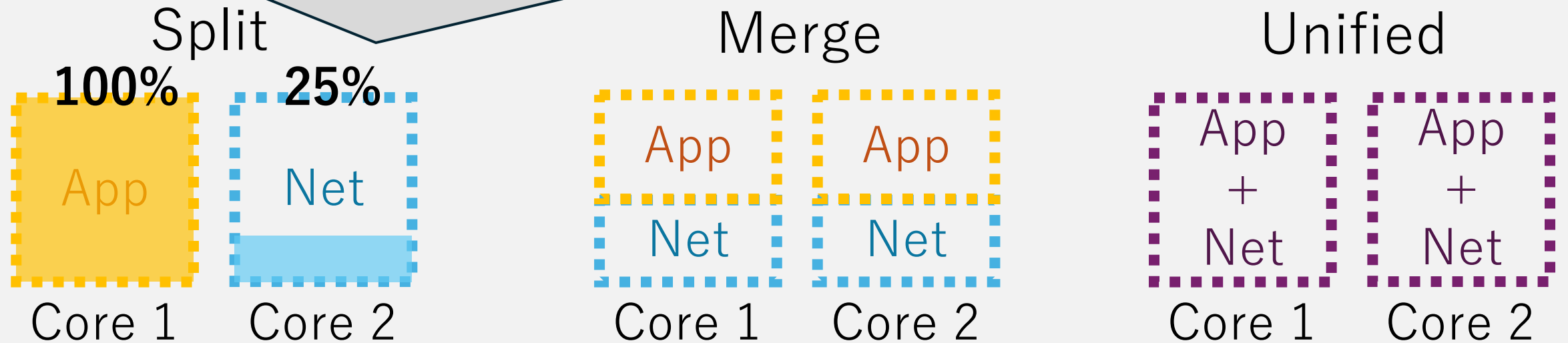
Three potential CPU core assignment models

Split

App

Net

Core 1     Core 2

App and Net threads run
on different CPU cores

# Issues of Existing TCP/IP Stacks

Limited ~~[App]~~ ment Models

Three ~~ployment~~ ment models

Split

App and Net threads run on the same CPU core

App

Net

Core 1  Core 2

Merge

App

Net

Core 1

# Issues of Existing TCP/IP Stacks

## Limited [Thread Deployment] Models

Three [thread deployment] models

**Split**

App | Net

Core 1 | Core 2

**Merge**

> Duplicate the same setup to available CPU cores

App | App

Net | Net

Core 1 | Core 2

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

A thread executes both App and Net logic

Split

Merge

Unified

App

Net

App

Net

App

Net

App
+
Net

Core 1

Core 2

Core 1

Core 2

Core 1

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

Duplicate the setup to available CPU cores

Split

Merge

Unified

App

Net

Core 1

Core 2

App

Net

App

Net

Core 1

Core 2

App + Net

App + Net

Core 1

Core 2

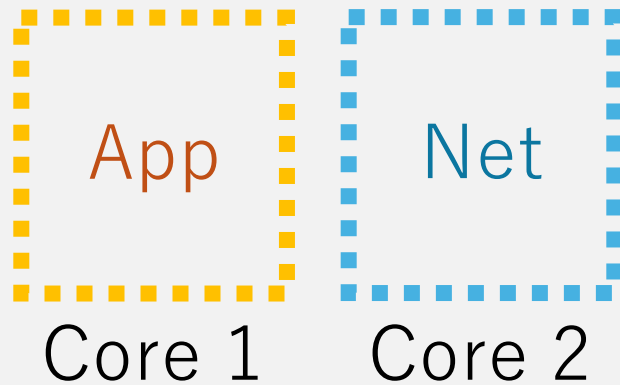# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

Three potential CPU core assignment models

Split

Merge

Unified

App

Net

Core 1　Core 2

App

App

Net

Net

Core 1　Core 2

App + Net

App + Net

Core 1　Core 2

Each of them has different performance characteristics

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models
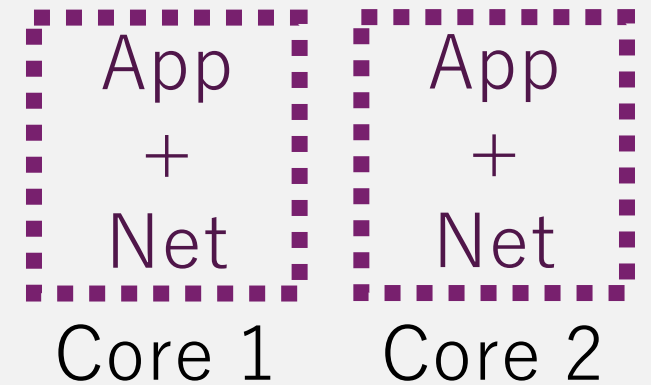
Three potential CPU core assignment models

Split

App — Core 1

Net — Core 2

Merge

App / Net — Core 1

App / Net — Core 2

Unified

App + Net — Core 1

App + Net — Core 2

CPU utilization

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

A busy thread cannot use other CPU cores' unused cycles

### Split

**100%** **25%**

App | Net

Core 1 | Core 2

### Merge

App | App
Net | Net

Core 1 | Core 2

### Unified

App + Net | App + Net

Core 1 | Core 2

CPU utilization

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

A busy thread cannot use other CPU cores' unused cycles

### Split

**100%**     **25%**

App     Net

Core 1     Core 2

### Merge

App     App

Net     Net

Core 1     Core 2

### Unified

App + Net     App + Net

Core 1     Core 2

Low CPU utilization     CPU utilization

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

The merge and unified models can fully use the CPU cycles

### Split

**100%**   **25%**

App   Net

Core 1   Core 2

Low CPU utilization

### Merge

**100%**   **100%**

App   App

Net   Net

Core 1   Core 2

### Unified

**100%**   **100%**

App + Net   App + Net

Core 1   Core 2

CPU utilization

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

Three potential CPU core assignment models

Split

App | Net

Core 1 | Core 2

Merge

App / Net | App / Net

Core 1 | Core 2

Unified

App + Net | App + Net

Core 1 | Core 2

Low CPU utilization

Context switch

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

A thread switching is necessary for every app/net transition

### Split

App
Net

Core 1   Core 2

Low CPU utilization

### Merge

App
Net

App
Net

Core 1   Core 2

### Unified

App
+
Net

App
+
Net

Core 1   Core 2

Context switch

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

A thread switching is necessary for every app/net transition

**Split**

App — Core 1

Net — Core 2

Low CPU utilization

**Merge**

App ⟷ Net — Core 1

App / Net — Core 2

High transition cost

**Unified**

App + Net — Core 1

App + Net — Core 2

Context switch

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

The split and unified models are free from this transition cost

### Split

App | Net

Core 1 | Core 2

Low CPU utilization

### Merge

App ⟲ Net | App Net

Core 1 | Core 2

High transition cost

Context switch

### Unified

App + Net | App + Net

Core 1 | Core 2

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

Performance-optimized
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

Implemented by TCP/IP stack users
（application developers）

A thread executing
application logic

https://github.com/yasukata/iip

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

Performance-optimized
TCP/IP stack

Implemented by TCP/IP stack users
(application developers)

A thread executing
TCP/IP processing

A thread executing
application logic

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

TCP/IP stack design often includes
the assignment model selection

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

Performance-optimized
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

Implemented by TCP/IP stack users
(application developers)

A thread executing
application logic

Developers cannot choose
a desired assignment model

https://github.com/yasukata/iip

68

# Issues of Existing TCP/IP Stacks

## Limited Choices for CPU Core Assignment Models

Performance-optimized
TCP/IP stack

A thread executing
TCP/IP processing

| TCP/IP processing | Thread scheduler |
| --- | --- |
| I/O subsystem | Thread runtime |

Implemented by TCP/IP stack users
(application developers)

A thread executing
application logic

**This issue makes it hard to build performance-optimal systems**

# Issues of Existing TCP/IP Stacks

Portability-aware
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP
processing

Thread
scheduler

I/O
subsystem

Thread
runtime

https://github.com/yasukata/iip

# Issues of Existing TCP/IP Stacks

Portability-aware
TCP/IP stack

A thread executing
TCP/IP processing

**TCP/IP processing**

Thread scheduler

I/O subsystem

Thread runtime

Portability-aware TCP/IP
TCP/IP stacks usually maintain
limited functionalities

https://github.com/yasukata/iip

71

# Issues of Existing TCP/IP Stacks

Portability-aware
TCP/IP stack

A thread executing
TCP/IP processing

TCP/IP
processing

Thread
scheduler

Developers are assumed to provide
platform-dependent components

Provided by TCP/IP stack users
(developers)

A thread executing
TCP/IP processing

Thread
scheduler

I/O
subsystem

Thread
runtime

https://github.com/yasukata/iip

72

# Issues of Existing TCP/IP Stacks

# Issues of Existing TCP/IP Stacks

Not sufficiently aware of performance-critical factors

Portability-aware
TCP/IP stack

Provided by TCP/IP stack users
(developers)

TCP/IP
processing

A thread executing
TCP/IP processing

Thread
scheduler

I/O
subsystem

Thread
runtime

# Issues of Existing TCP/IP Stacks

## Not sufficiently aware of performance-critical factors

Portability-aware
TCP/IP stack

Provided by TCP/IP stack users
(developers)

A thread executing
TCP/IP processing

TCP/IP
processing

- Unaware of NIC offloading features

# Issues of Existing TCP/IP Stacks

## Not sufficiently aware of performance-critical factors

Portability-aware
TCP/IP stack

Provided by TCP/IP stack users
(developers)

A thread executing
TCP/IP processing

**TCP/IP processing**

- Unaware of NIC offloading features
- Lack of zero-copy I/O capability

er

https://github.com/yasukata/iip

76

# Issues of Existing TCP/IP Stacks

## Not sufficiently aware of performance-critical factors

Portability-aware
TCP/IP stack

Provided by TCP/IP stack users
(developers)

A thread executing
TCP/IP processing

TCP/IP
processing

- Unaware of NIC offloading features
- Lack of zero-copy I/O capability
- Lack of multi-core scalability

https://github.com/yasukata/iip

77

# Issues of Existing TCP/IP Stacks

Not sufficiently aware of performance-critical factors

Portability

Provided by TCP/IP stack users

**Portability-aware TCP/IP stacks cannot achieve good performance in various workloads**

TCP/IP processing

TCP/IP processing

- Unaware of NIC offloading features
- Lack of zero-copy I/O capability
- Lack of multi-core scalability

# iip

iip

TCP/IP processing

Provided by TCP/IP stack users
(developers)

A thread executing
TCP/IP processing

Thread
scheduler

I/O
subsystem

Thread
runtime

# iip

iip only implements the TCP/IP processing functionality to minimize the chance for causing functionality conflicts

iip

TCP/IP processing

Provided by TCP/IP stack users (developers)

A thread executing TCP/IP processing

Thread scheduler

I/O subsystem

Thread runtime

https://github.com/yasukata/iip

80

iip

Platform-dependent functionalities are assumed
to be provided by developers through the API

nality
nflicts

iip

Provided by TCP/IP stack users
（developers）

A thread executing
TCP/IP processing

TCP/IP
processing

API

Thread
scheduler

I/O
subsystem

Thread
runtime

# iip

Platform-dependent functionalities are assumed to be provided by developers through the API

...nality ...nflicts

iip

No dependency on CPUs, NICs, OSes, libraries, and compilers

Provided by TCP/IP stack users （developers）

**TCP/IP processing**

API

A thread executing TCP/IP processing

Thread scheduler

I/O subsystem

Thread runtime

# iip

Platform-dependent functionalities are assumed to be provided by developers through the API

nality nflicts

iip

No dependency on CPUs, NICs, OSes, libraries, and compilers

TCP/IP processing

C89 / C++98 compliant for old and future compilers

API

Provided by TCP/IP stack users (developers)

A thread executing TCP/IP processing

Thread scheduler

I/O subsystem

Thread runtime

# iip

Platform-dependent functionalities are assumed to be provided by developers through the API

nality nflicts

Provided by TCP/IP stack users (developers)

iip

No dependency on CPUs, NICs, OSes, libraries, and compilers

A thread executing

iip pays attention to performance-critical factors

TCP/IP processing

API

C89 / C++98 compliant for old and future compilers

# iip

Platform-dependent functionalities are assumed to be provided by developers through the API

nality nflicts

## iip

No dependency on CPUs, NICs, OSes, libraries, and compilers

Provided by TCP/IP stack users (developers)

A thread executing

TCP/IP processing

API

iip pays attention to performance-critical factors

- uses NIC offloading features

C89 / C++98 compliant for old and future compilers

# iip

Platform-dependent functionalities are assumed
to be provided by developers through the API

nality
nflicts

iip

No dependency on CPUs, NICs,
OSes, libraries, and compilers

Provided by TCP/IP stack users
(developers)

A thread executing

iip pays attention to
performance-critical factors

TCP/IP
processing

API

- uses NIC offloading features

- supports zero-copy I/O

C89 / C++98 compliant
for old and future compilers

https://github.com/yasukata/iip

86

# iip

Platform-dependent functionalities are assumed to be provided by developers through the API

nality nflicts

iip

No dependency on CPUs, NICs, OSes, libraries, and compilers

Provided by TCP/IP stack users (developers)

A thread executing

iip pays attention to performance-critical factors

- uses NIC offloading features

TCP/IP processing

API

- supports zero-copy I/O

C89 / C++98 compliant for old and future compilers

- aware of multi-core scalability

https://github.com/yasukata/iip

87

# Integratability Benefits

- iip can be integrated into the ns-3 simulator written in C++
  - https://github.com/yasukata/iip-ns

ns-3

A thread executing
TCP/IP processing

iip

Thread
scheduler

I/O
subsystem

Thread
runtime

https://github.com/yasukata/iip

# Integratability Benefits

- iip can be integrated into the ns-3 simulator written in C++
  - https://github.com/yasukata/iip-ns

ns-3

A thread executing
TCP/IP processing

iip

Thread
scheduler

I/O
subsystem

Thread
runtime

These functionalities are
specific to ns-3

# Integratability Benefits

- iip can be integrated into the ns-3 simulator written in C++
  - https://github.com/yasukata/iip-ns

ns-3

A thread executing
TCP/IP processing

iip

Thread
scheduler

I/O
subsystem

Thread
runtime

These functionalities are
specific to ns-3

These functionalities makes it hard
for TCP/IP stacks depending on
common OS features to run on ns-3

https://github.com/yasukata/iip

# Integratability Benefits

- iip can be integrated into the ns-3 simulator written in C++
  - https://github.com/yasukata/iip-ns

ns-3

A thread executing
TCP/IP processing

iip

Thread
scheduler

I/O
subsystem

These functionalities are
specific to ns-3

These functionalities makes it hard
for TCP/IP stacks depending on

iip can be easily integrated into ns-3 because
it does not depend on platform-specific functionalities

https://github.com/yasukata/iip

# Integratability Benefits

- iip runs on various I/O backends, including but not limited to:
    - DPDK (Linux): https://github.com/yasukata/iip-dpdk
    - AF_XDP (Linux): https://github.com/yasukata/iip-af_xdp
    - netmap (Linux/FreeBSD): https://github.com/yasukata/iip-netmap



iip does not depend on specific I/O subsystems

# Evaluation: Small Message Exchange

- TCP ping-pong workload

32 CPU cores

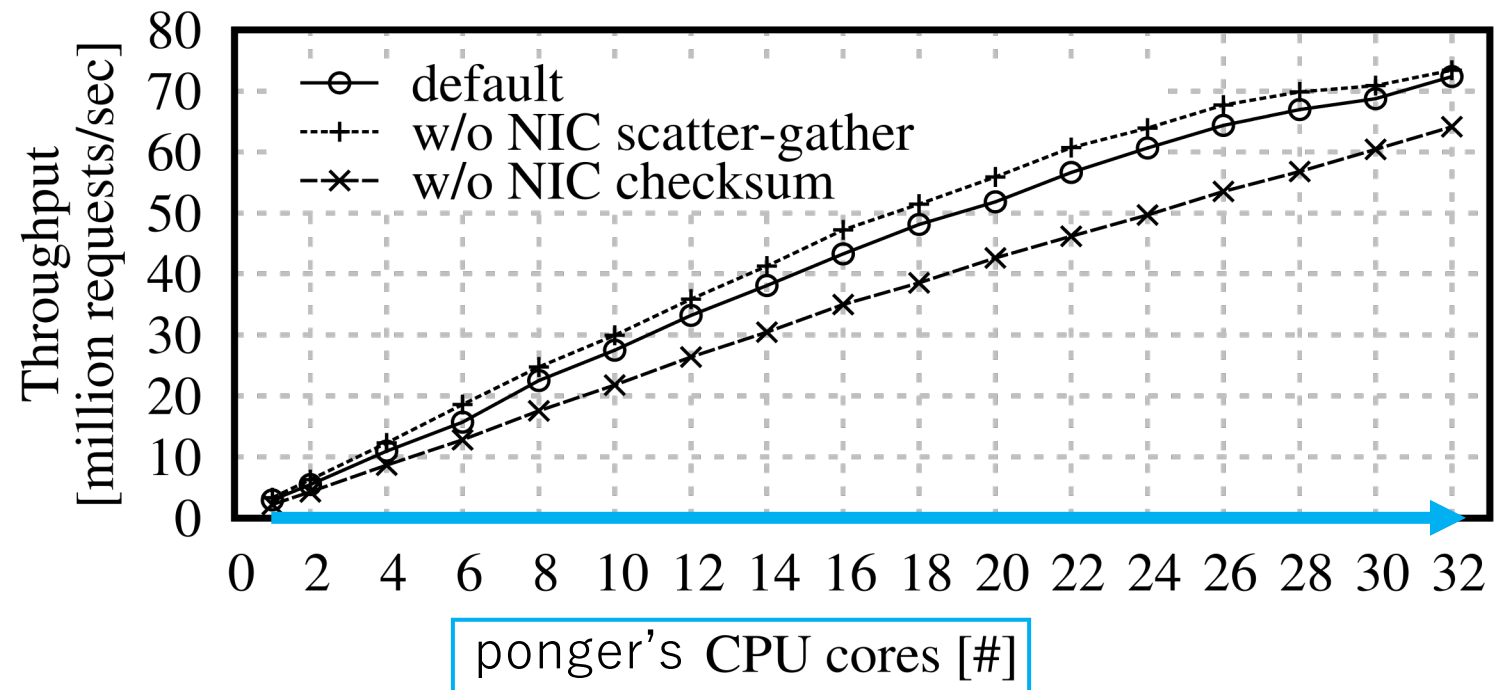We have tested while changing the number of CPU cores on the ponger side

The pinger and ponger exchange 1-byte TCP payloads

A ponger thread handles 32 concurrent TCP connections

1~32 CPU cores

| pinger app |
| iip |
| DPDK |

| ponger app |
| iip |
| DPDK |

100 Gbps

# Evaluation: Small Message Exchange

- TCP ping-pong workload

We have tested while changing
the number of CPU cores
on the ponger side

The pinger and ponger exchange
1-byte TCP payloads

32 CPU
cores

A ponger thread handles

1~32 CPU
cores

| pinger app |
| iip |
| DPDK |

- These two perform ping-pong with 1-byte TCP payloads
- Pinger always uses 32 CPU cores
- Ponger uses 1 ~ 32 CPU cores

| ponger app |
| iip |
| DPDK |

https://github.com/yasukata/iip

# Evaluation: Small Message Exchange

- The pinger and ponger apps exchange 1-byte TCP payloads

# Evaluation: Small Message Exchange

- The pinger and ponger apps exchange 1-byte TCP payloads

# Evaluation: Small Message Exchange

- The pinger and ponger apps exchange 1-byte TCP payloads

# Evaluation: Small Message Exchange

- The pinger and ponger apps exchange 1-byte TCP payloads

- Performance factors
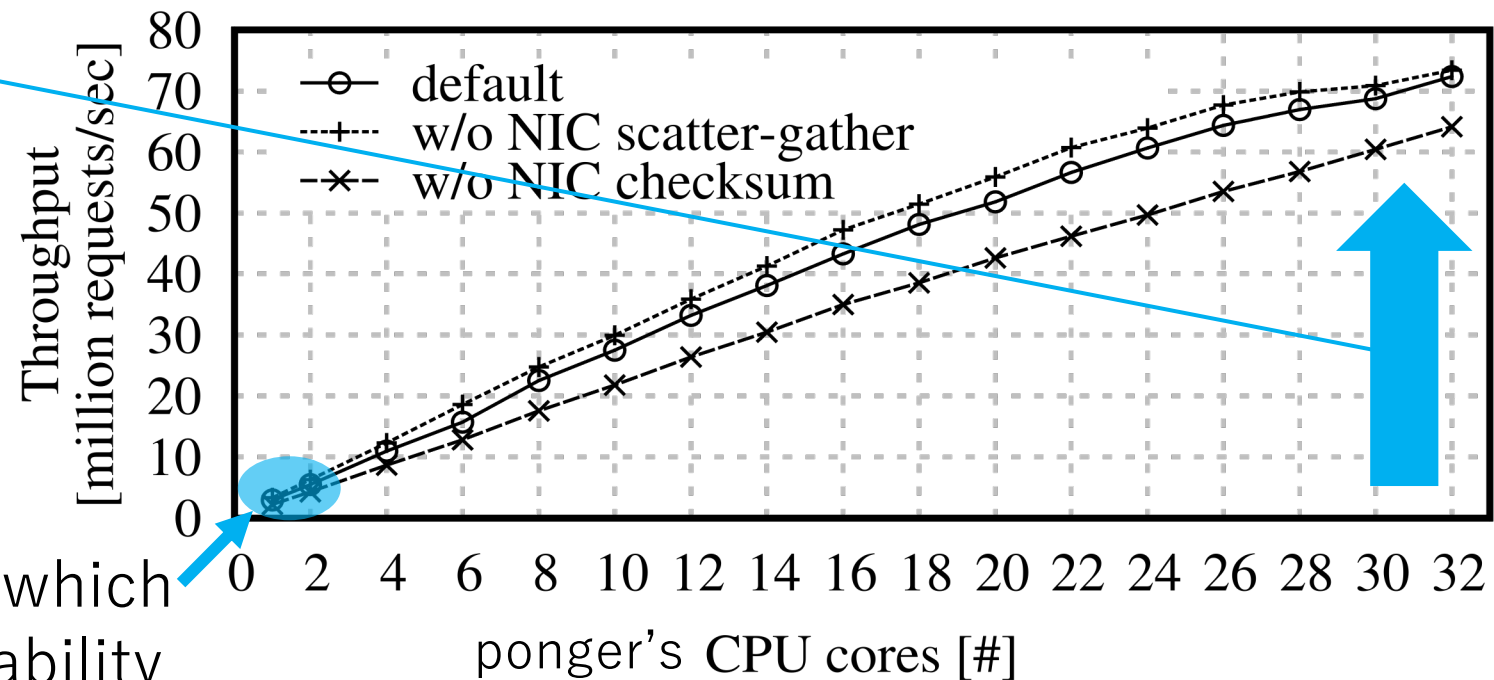  - Multi-core scalability

# Evaluation: Small Message Exchange

- The pinger and ponger apps exchange 1-byte TCP payloads

- Performance factors
  - Multi-core scalability

Throughput [million requests/sec] vs ponger's CPU cores [#]

- default
- w/o NIC scatter-gather
- w/o NIC checksum

Performance of TCP/IP stacks which are unaware of multi-core scalability
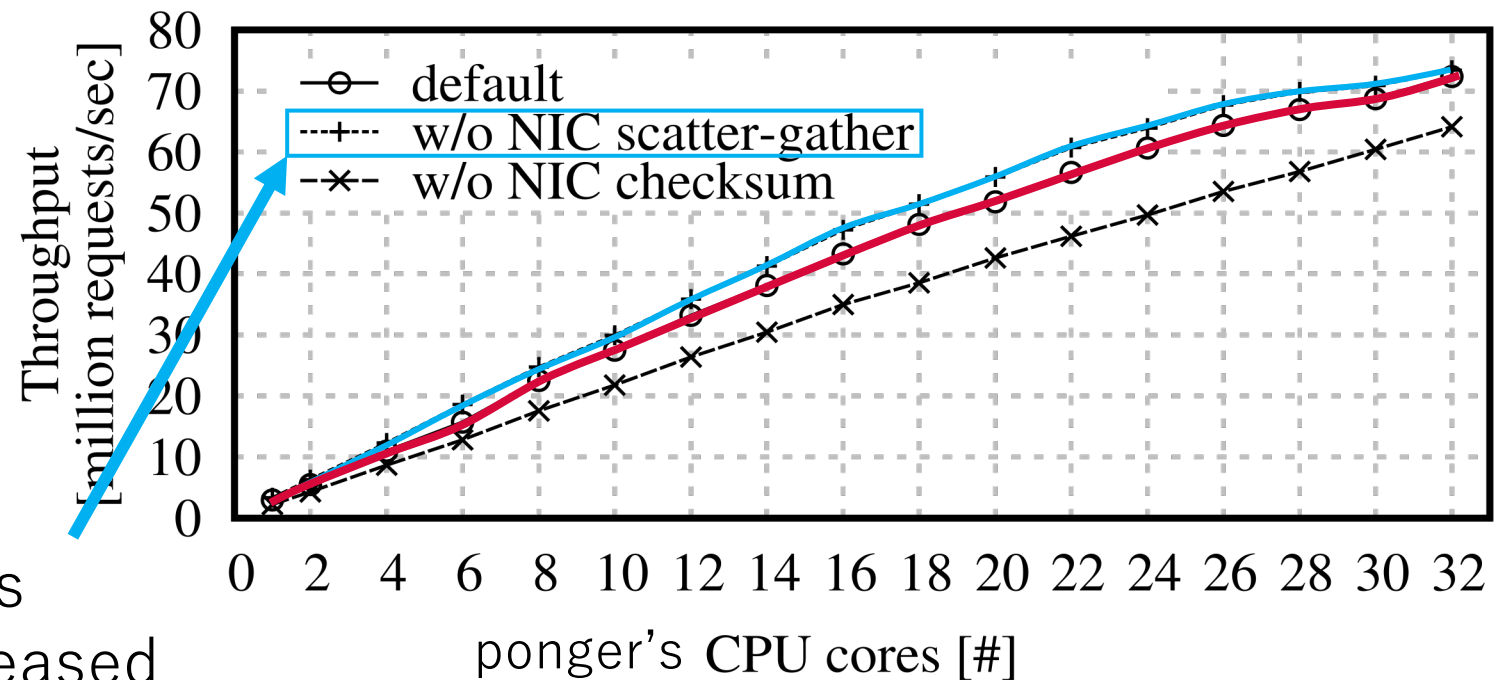
https://github.com/yasukata/iip

# Evaluation: Small Message Exchange

- The pinger and ponger apps exchange 1-byte TCP payloads

- Performance factors
  - Multi-core scalability

Performance of TCP/IP stacks which are unaware of multi-core scalability



Throughput [million requests/sec] vs ponger's CPU cores [#]

- default
- w/o NIC scatter-gather
- w/o NIC checksum

# Evaluation: Small Message Exchange

- The pinger and ponger apps exchange 1-byte TCP payloads

- Performance factors
  - Multi-core scalability

All NIC offloading features and
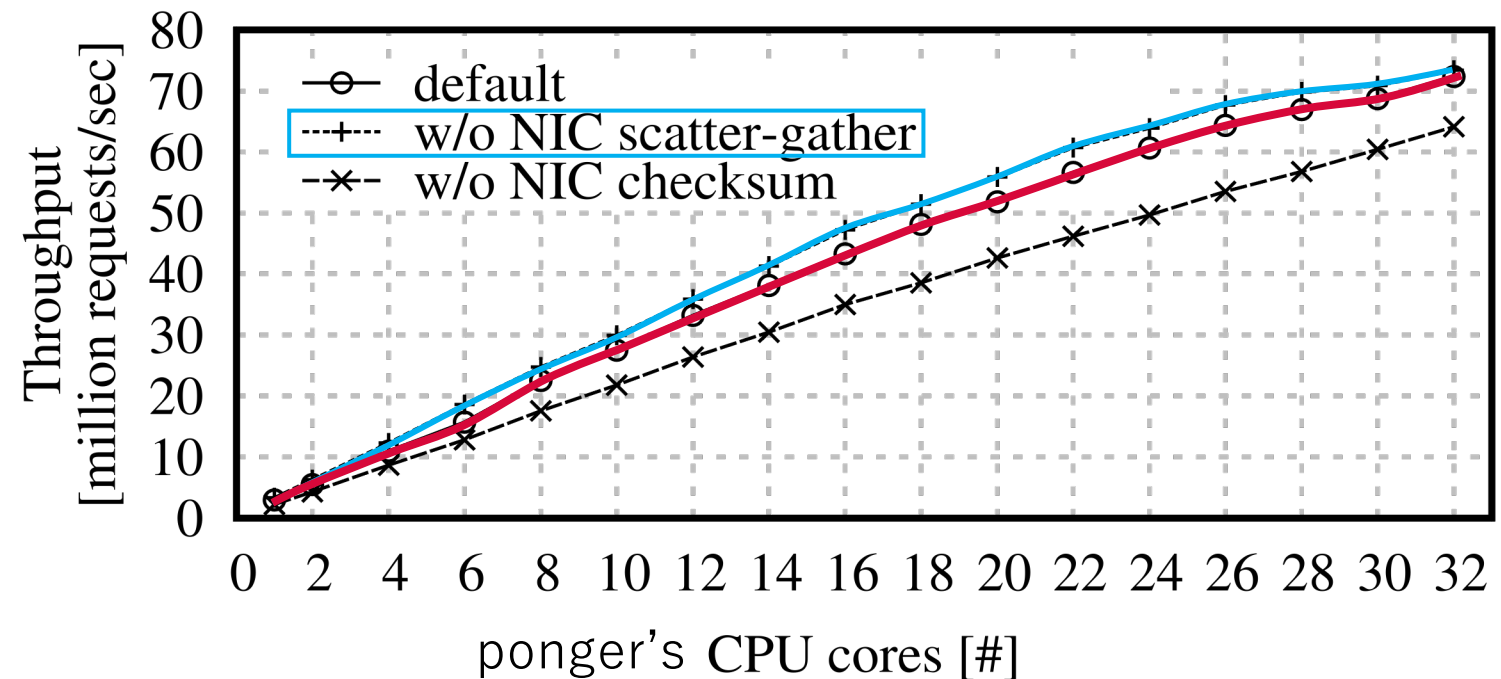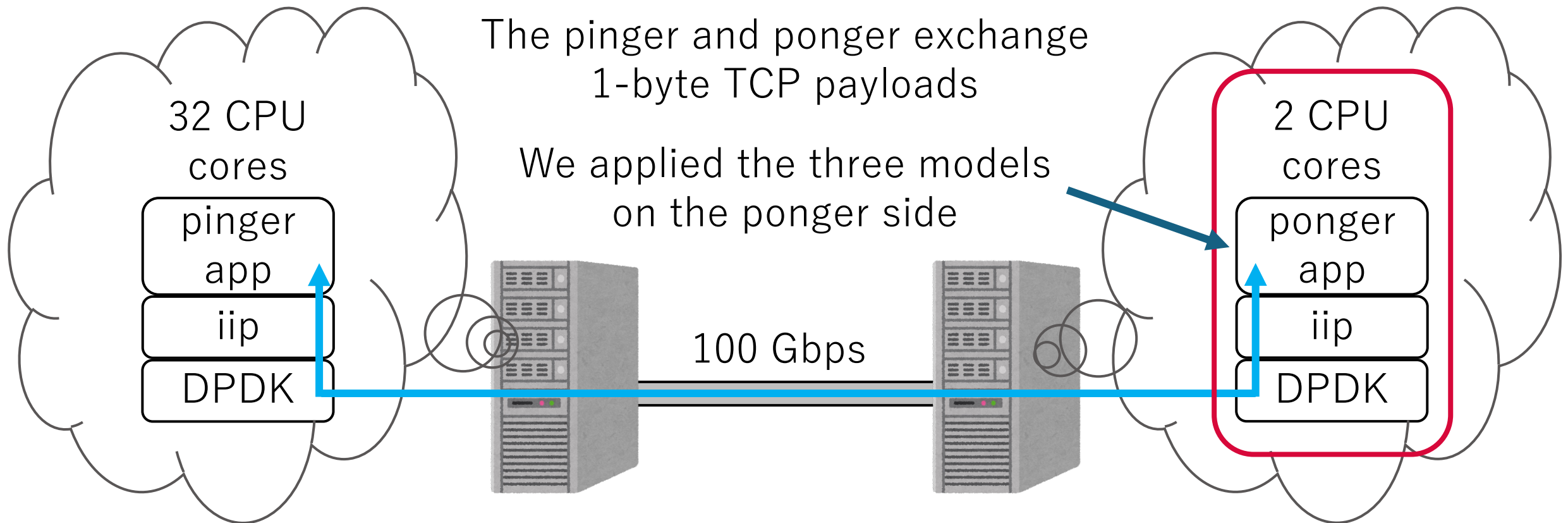zero-copy transmission are activated

# Evaluation: Small Message Exchange

- The pinger and ponger apps exchange 1-byte TCP payloads

- Performance factors
  - Multi-core scalability
  - Checksum offloading

When NIC checksum offloading is deactivated, throughput is degraded



Throughput [million requests/sec] vs ponger's CPU cores [#]

Legend:
- default
- w/o NIC scatter-gather
- w/o NIC checksum

# Evaluation: Small Message Exchange

- The pinger and ponger apps exchange 1-byte TCP payloads

- Performance factors
  - Multi-core scalability
  - Checksum offloading



Throughput [million requests/sec] vs ponger's CPU cores [#]

Legend: default, w/o NIC scatter-gather, w/o NIC checksum

When zero-copy transmission is deactivated, throughput is increased

# Evaluation: Small Message Exchange

- The pinger and ponger apps exchange 1-byte TCP payloads

- Performance factors
  - Multi-core scalability
  - Checksum offloading
- Tips
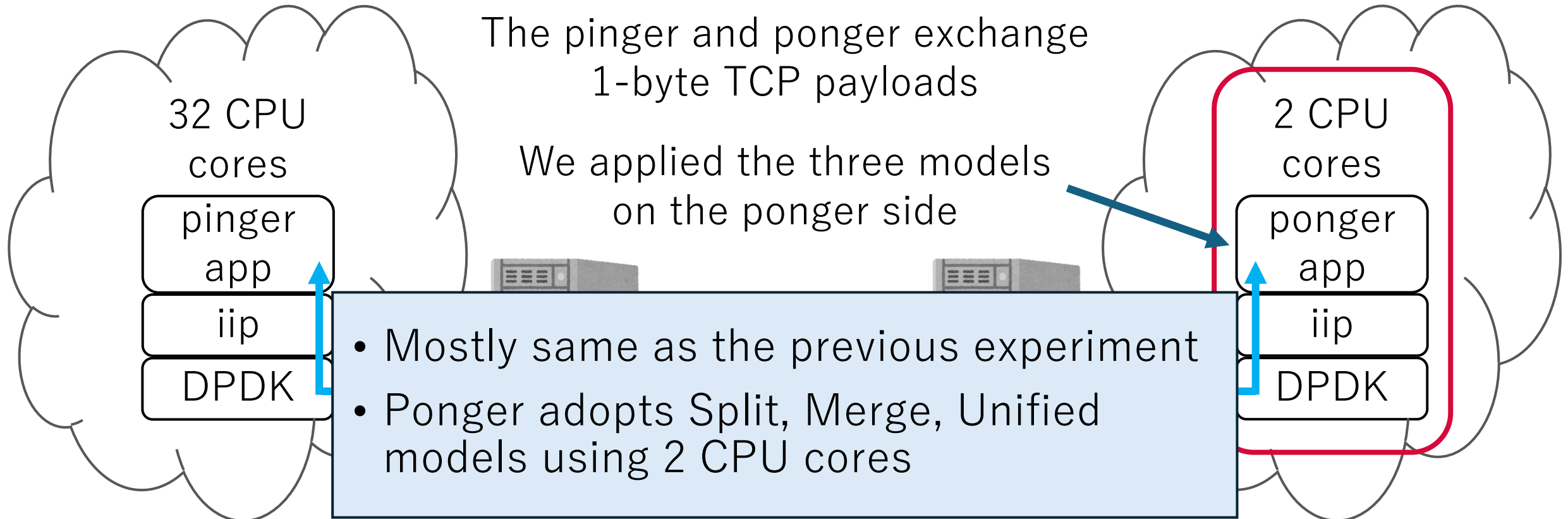  - For small messages, copy is faster than zero-copy transmission



Throughput [million requests/sec] vs ponger's CPU cores [#]
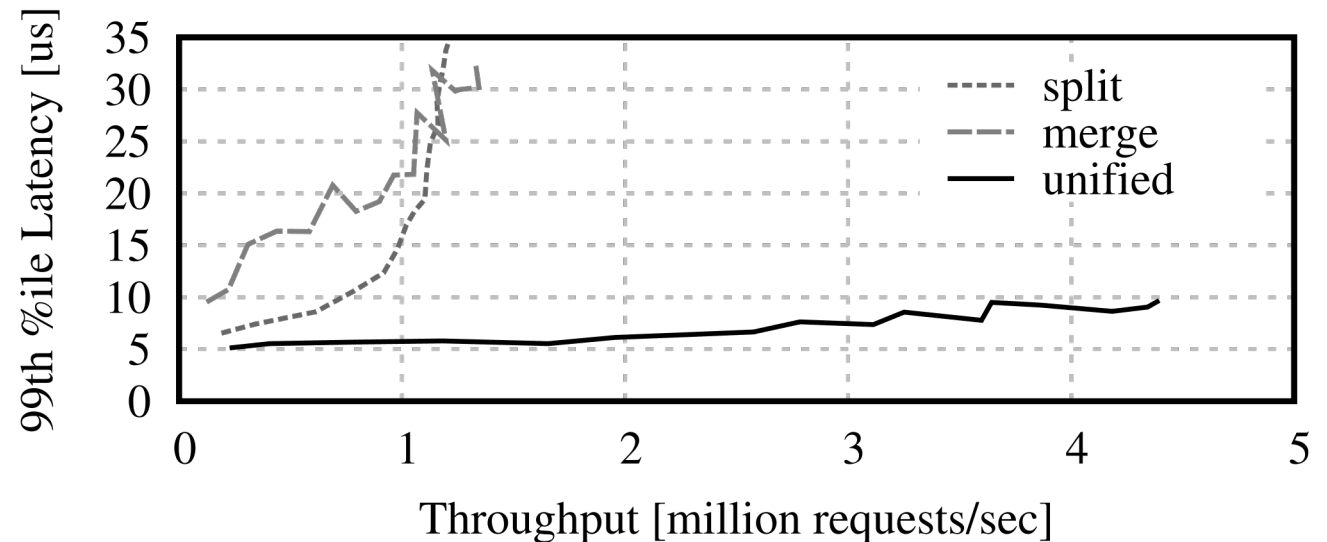
Legend: default, w/o NIC scatter-gather, w/o NIC checksum

# Evaluation: CPU Core Assignment Models

- TCP ping-pong workload



32 CPU cores

pinger app

iip

DPDK

The pinger and ponger exchange
1-byte TCP payloads

We applied the three models
on the ponger side

100 Gbps

2 CPU cores

ponger app

iip

DPDK

https://github.com/yasukata/iip

105

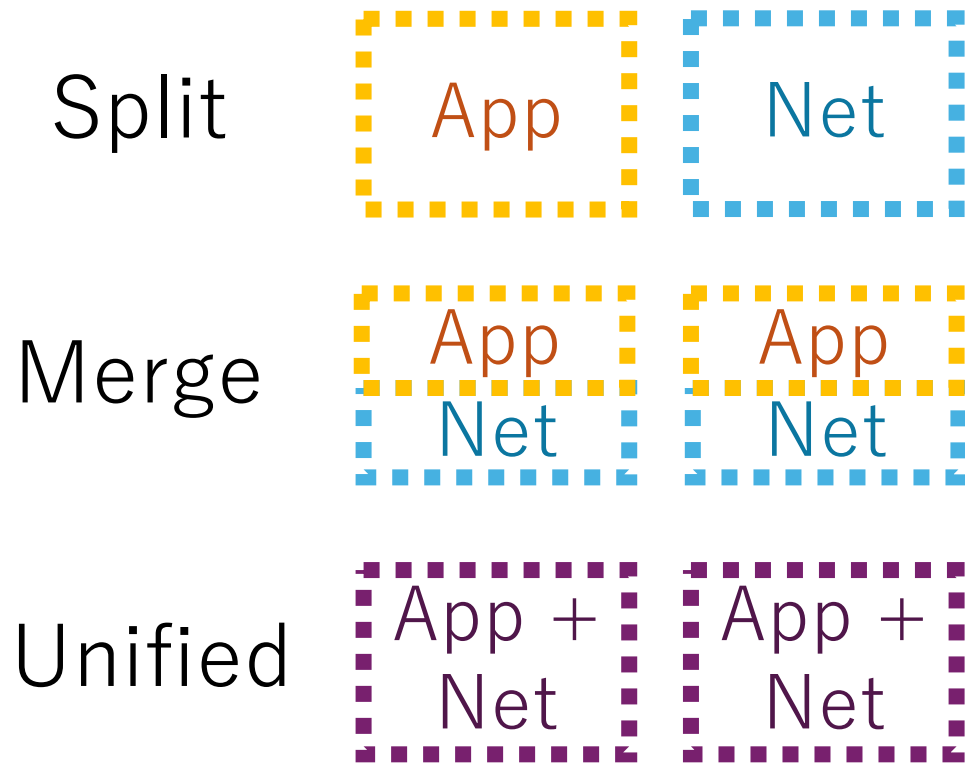# Evaluation: CPU Core Assignment Models

- TCP ping-pong workload

The pinger and ponger exchange
1-byte TCP payloads

We applied the three models
on the ponger side

32 CPU
cores

pinger
app

iip

DPDK

2 CPU
cores

ponger
app

iip

DPDK

- Mostly same as the previous experiment
- Ponger adopts Split, Merge, Unified
  models using 2 CPU cores

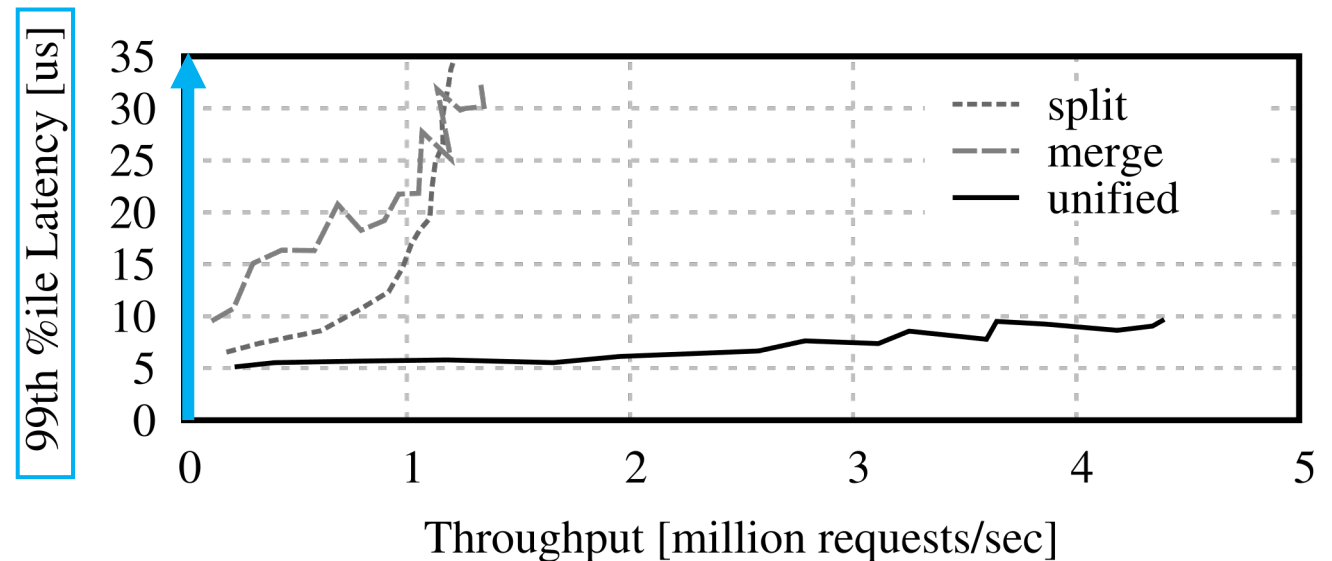https://github.com/yasukata/iip

# Evaluation: CPU Core Assignment Models

- The pinger and ponger apps exchange 1-byte TCP payloads

Split

App  Net

Merge

App  App
Net  Net

Unified

App +  App +
Net    Net

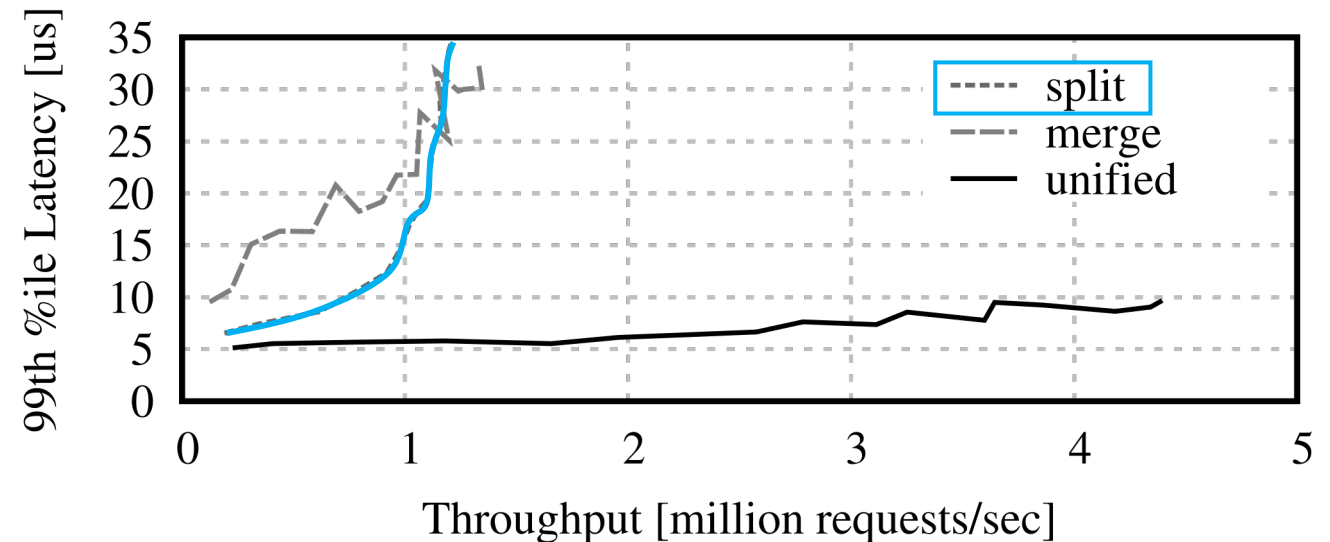# Evaluation: CPU Core Assignment Models

- The pinger and ponger apps exchange 1-byte TCP payloads

Split

App    Net

Merge

App    App
Net    Net

Unified

App +    App +
Net      Net



99th %ile Latency [us]

35
30
25
20
15
10
5
0

- - - - split
— - merge
—— unified

0    1    2    3    4    5

Throughput [million requests/sec]

# Evaluation: CPU Core Assignment Models

- The pinger and ponger apps exchange 1-byte TCP payloads

Split

App    Net

Merge

App    App
Net    Net

Unified

App +    App +
Net      Net



99th %ile Latency [us] vs Throughput [million requests/sec]

- - - - split
— - — merge
——— unified

# Evaluation: CPU Core Assignment Models

- The pinger and ponger apps exchange 1-byte TCP payloads

| | | |
|---|---|---|
| **Split** | App | Net |
| **Merge** | App<br>Net | App<br>Net |
| **Unified** | App +<br>Net | App +<br>Net |



https://github.com/yasukata/iip

110

# Evaluation: CPU Core Assignment Models

- The pinger and ponger apps exchange 1-byte TCP payloads

Split

| App | Net |

Merge

| App | App |
| Net | Net |

Unified

| App + Net | App + Net |



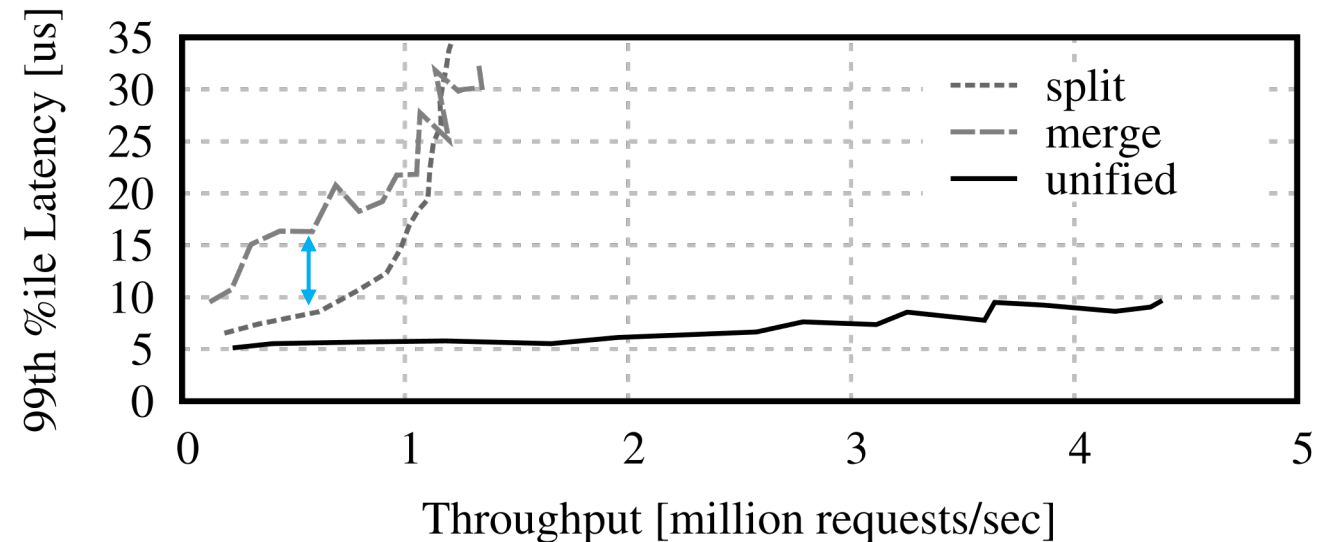https://github.com/yasukata/iip                                    111

# Evaluation: CPU Core Assignment Models

- The pinger and ponger apps exchange 1-byte TCP payloads

Split · App · Net

The latency gap comes from the app/net transition cost of the merge model

Merge · App / Net · App / Net

Unified · App + Net · App + Net

# Evaluation: CPU Core Assignment Models

- The pinger and ponger apps exchange 1-byte TCP payloads

Split

App    Net

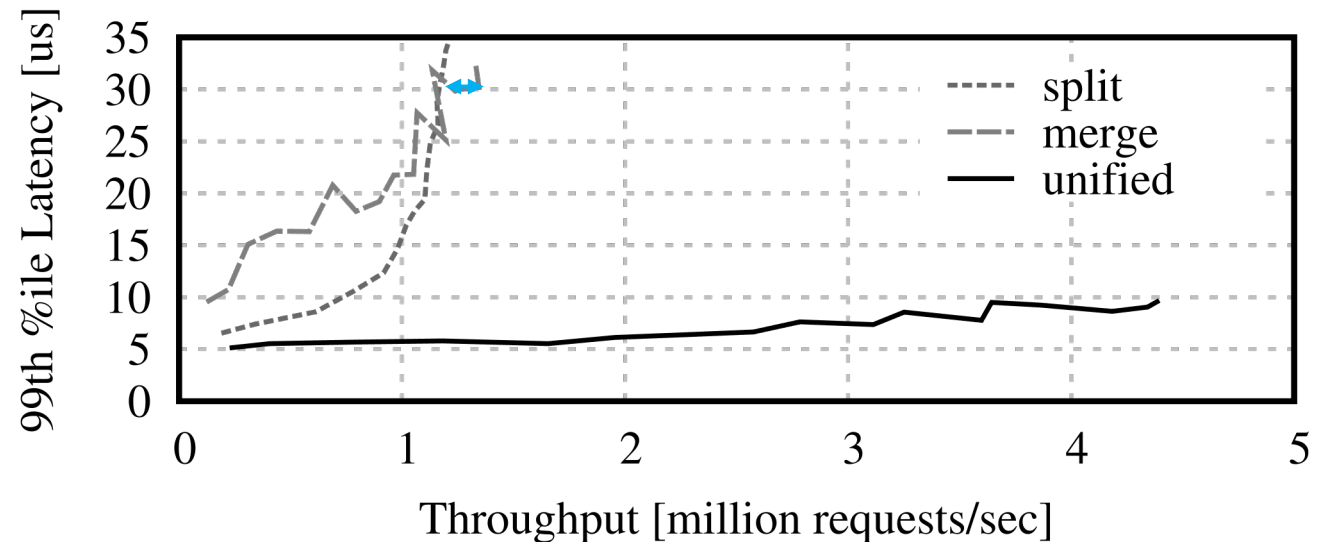The merge model exhibits higher throughput
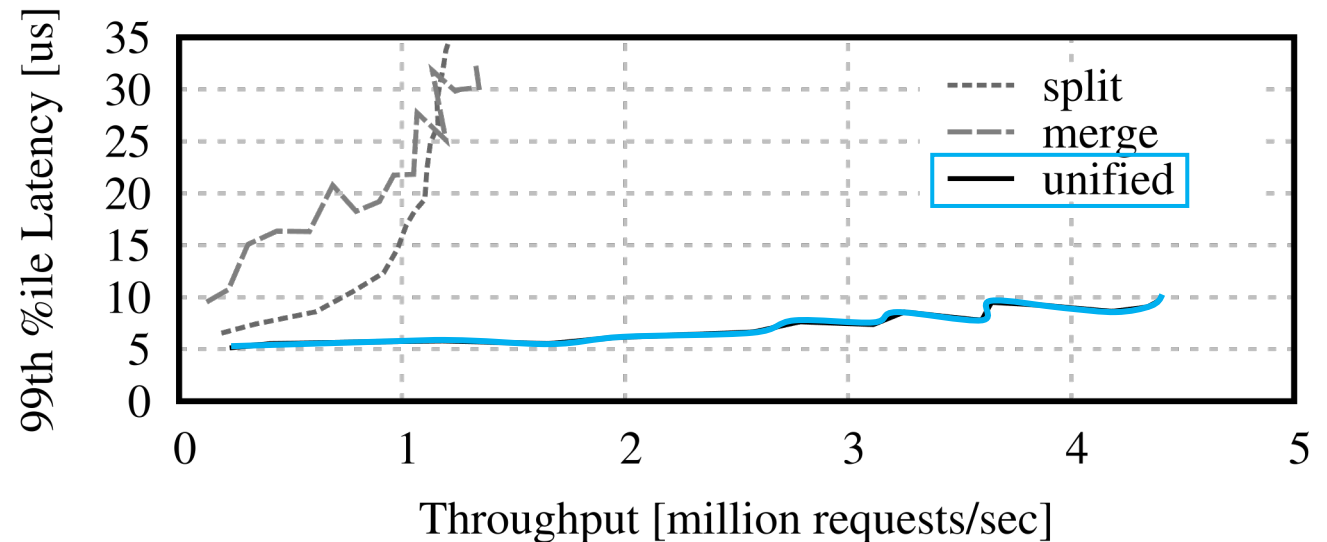than the split model because of CPU utilization

Merge

App    App
Net    Net

Unified

App +    App +
Net      Net

# Evaluation: CPU Core Assignment Models

- The pinger and ponger apps exchange 1-byte TCP payloads

Split

App    Net

The unified model achieves the best speed because it is free from the issues of the others

Merge

App    App
Net    Net

Unified

App +    App +
Net      Net

# Evaluation: Bulk Data Transfer

- Data transfer workload

Each uses 1 CPU core

The sender sends the same data
to the receiver repeatedly

These two communicate
over 1 TCP connection

1 CPU
cores

| sender app |
| iip |
| DPDK |

100 Gbps

1 CPU
core

| receiver app |
| iip |
| DPDK |

# Evaluation: Bulk Data Transfer

Each uses 1 CPU core

- Data transfer workload

The sender sends the same data
to the receiver repeatedly

1 CPU
cores

1 CPU
core

These two communicate
over 1 TCP connection

| sender app |
| iip |
| DPDK |

| receiver app |
| iip |
| DPDK |

- The sender sends the same data to the receiver over 1 TCP connection
- Each uses 1 CPU core

https://github.com/yasukata/iip

116

# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver

https://github.com/yasukata/iip

# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver

Legend:
- default
- sender w/o NIC scatter-gather
- sender w/o NIC TSO
- sender w/o NIC TSO and NIC checksum
- receiver w/o NIC LRO
- receiver w/o NIC checksum



https://github.com/yasukata/iip

# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver



https://github.com/yasukata/iip

# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver

All NIC offloading features and
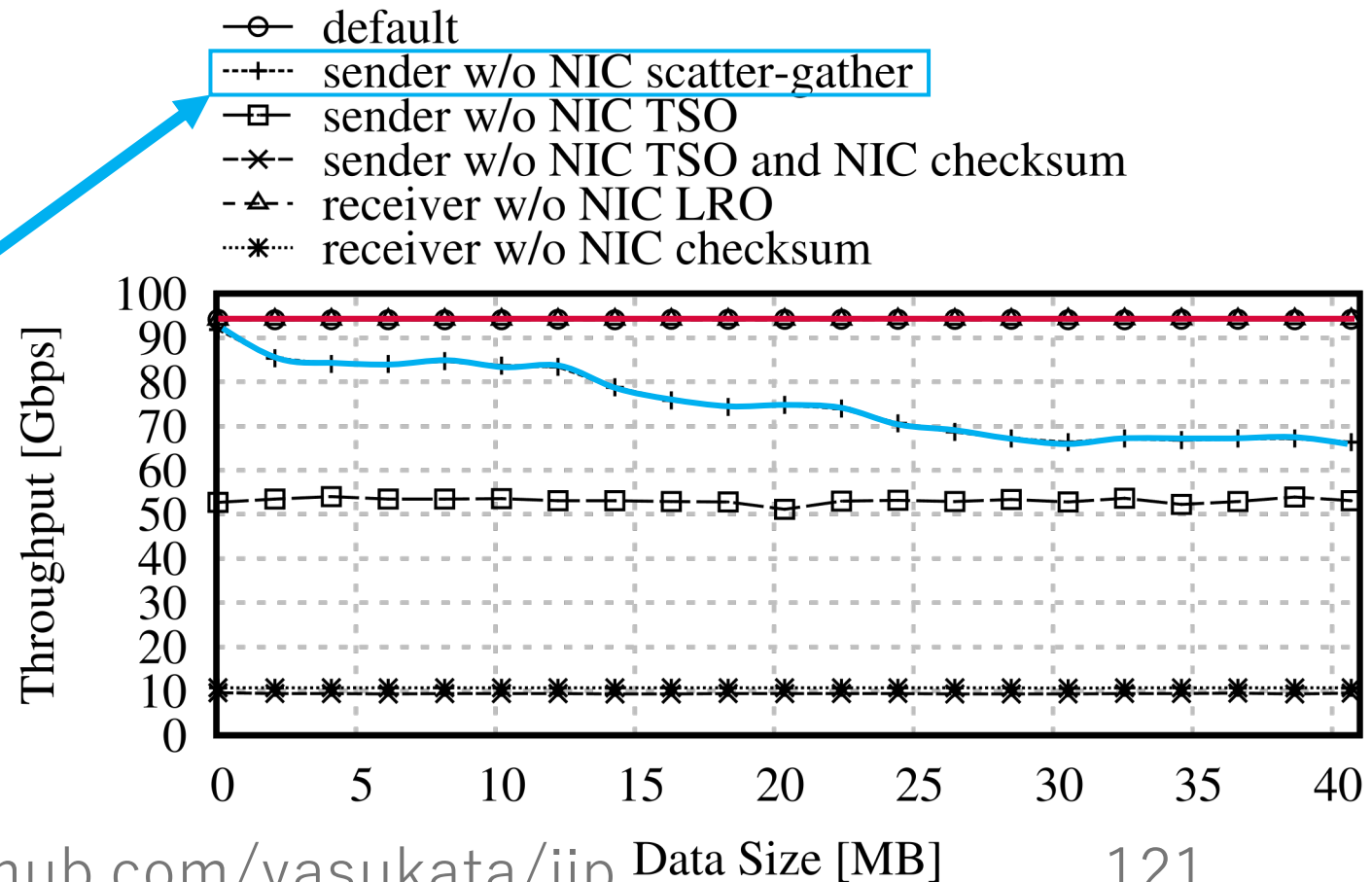zero-copy transmission are activated

# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver

- Performance factors
  - Zero-copy transmission

When the sender deactivates zero-copy transmission, throughput is degraded according to the size of the data

Legend:
- default
- sender w/o NIC scatter-gather
- sender w/o NIC TSO
- sender w/o NIC TSO and NIC checksum
- receiver w/o NIC LRO
- receiver w/o NIC checksum

Chart: Throughput [Gbps] vs Data Size [MB]

# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver

- Performance factors
  - Zero-copy transmission

When the sender deactivates zero-copy transmission, throughput is degraded according to the size of the data
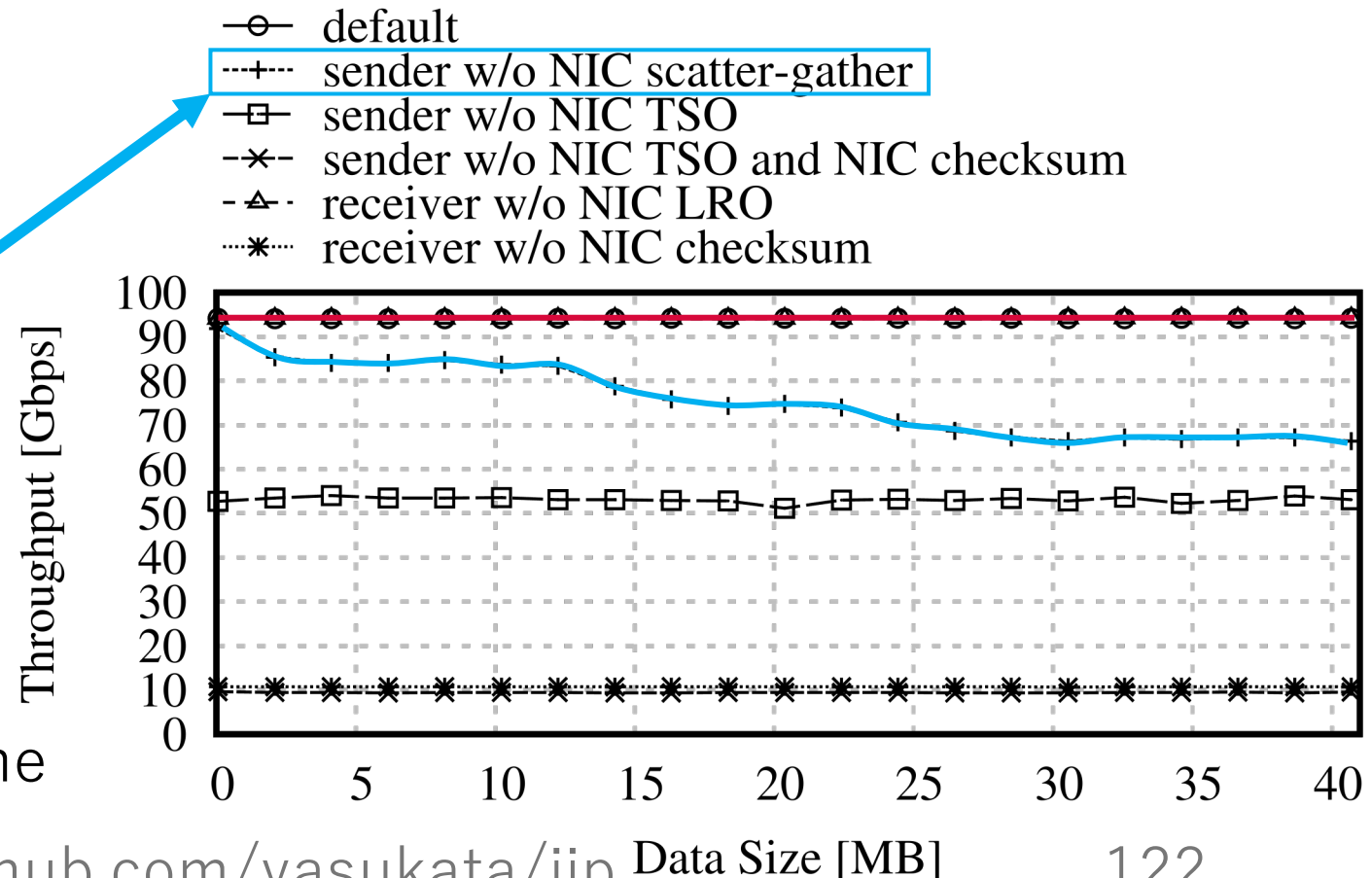
We consider this happens because the payload occupies the CPU cache



Legend:
- default
- sender w/o NIC scatter-gather
- sender w/o NIC TSO
- sender w/o NIC TSO and NIC checksum
- receiver w/o NIC LRO
- receiver w/o NIC checksum

Y-axis: Throughput [Gbps] (0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100)
X-axis: Data Size [MB] (0, 5, 10, 15, 20, 25, 30, 35, 40)

# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver

- Performance factors
  - Zero-copy transmission
  - TSO

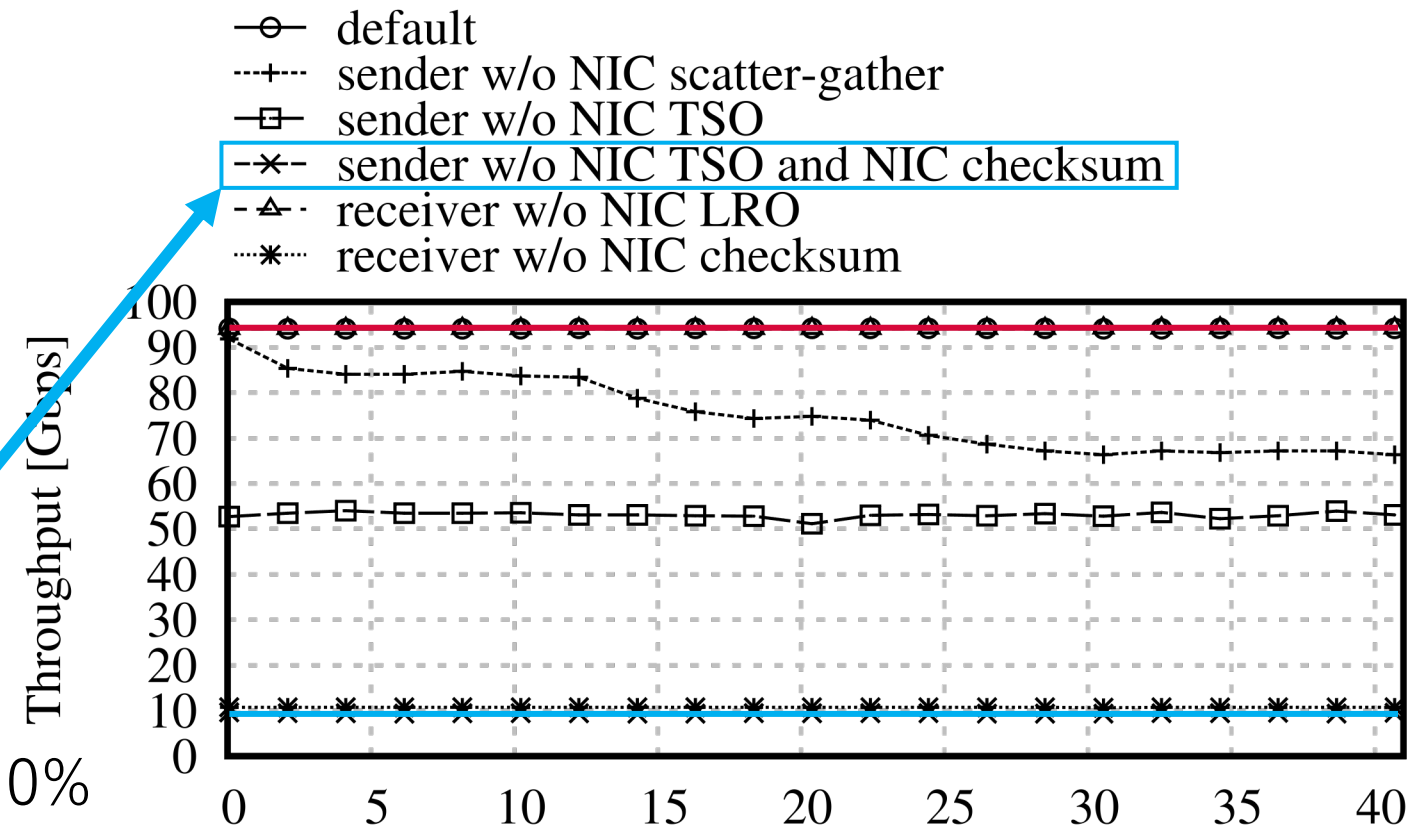When the sender deactivates TSO, throughput is almost halved

# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver

- Performance factors
  - Zero-copy transmission
  - TSO
  - Checksum offloading

When the sender deactivates
TSO and checksum offloading,
throughput goes down to around 10%

Legend:
- default
- sender w/o NIC scatter-gather
- sender w/o NIC TSO
- sender w/o NIC TSO and NIC checksum
- receiver w/o NIC LRO
- receiver w/o NIC checksum

Y-axis: Throughput [Gbps]
X-axis: Data Size [MB]

# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver

- Performance factors
  - Zero-copy transmission
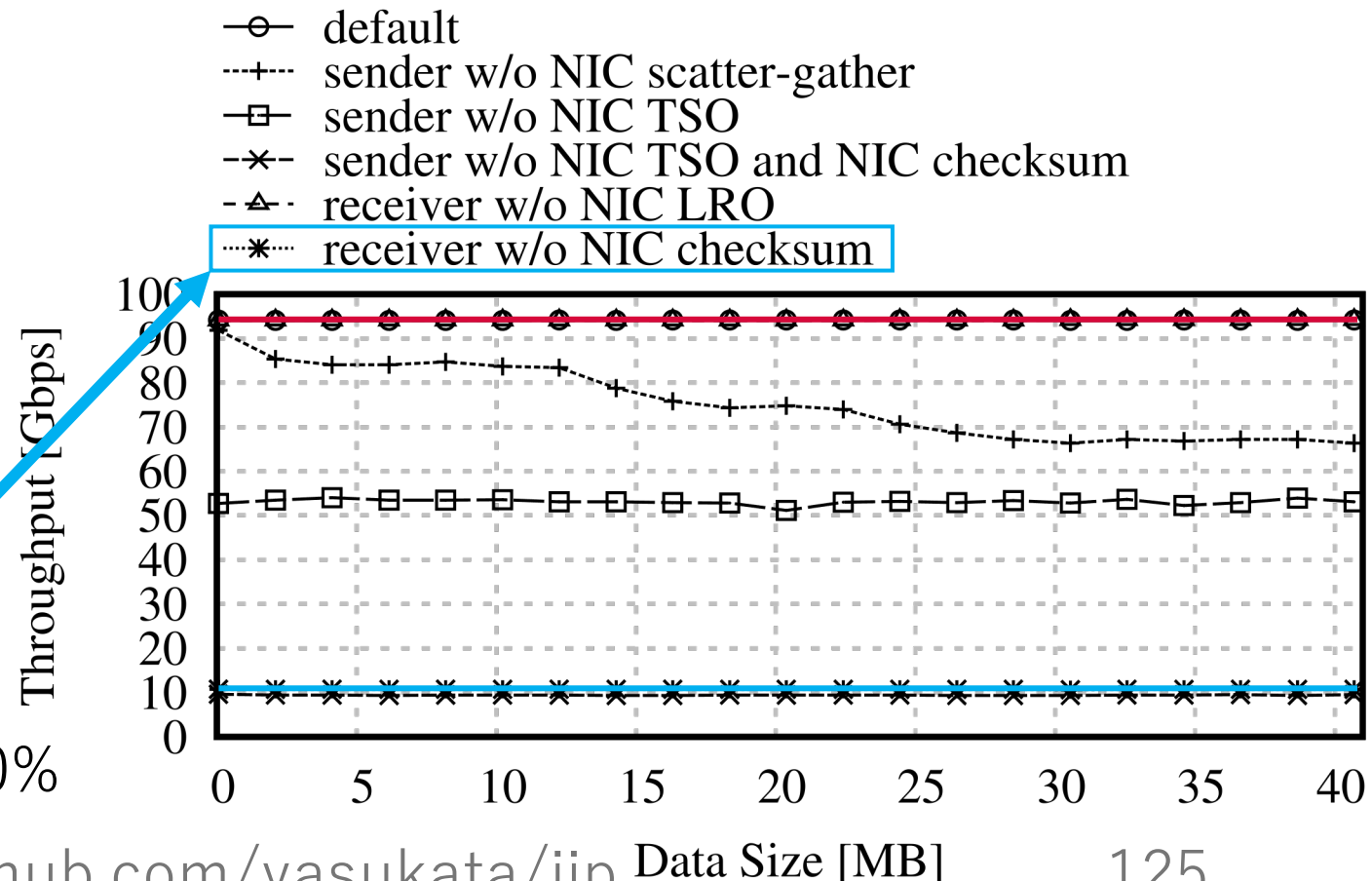  - TSO
  - Checksum offloading

When the receiver deactivates checksum offloading, throughput goes down to around 10%



Legend:
- default
- sender w/o NIC scatter-gather
- sender w/o NIC TSO
- sender w/o NIC TSO and NIC checksum
- receiver w/o NIC LRO
- receiver w/o NIC checksum

Y-axis: Throughput [Gbps] (0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100)
X-axis: Data Size [MB] (0, 5, 10, 15, 20, 25, 30, 35, 40)

https://github.com/yasukata/iip

125

# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver

- Performance factors
  - Zero-copy transmission
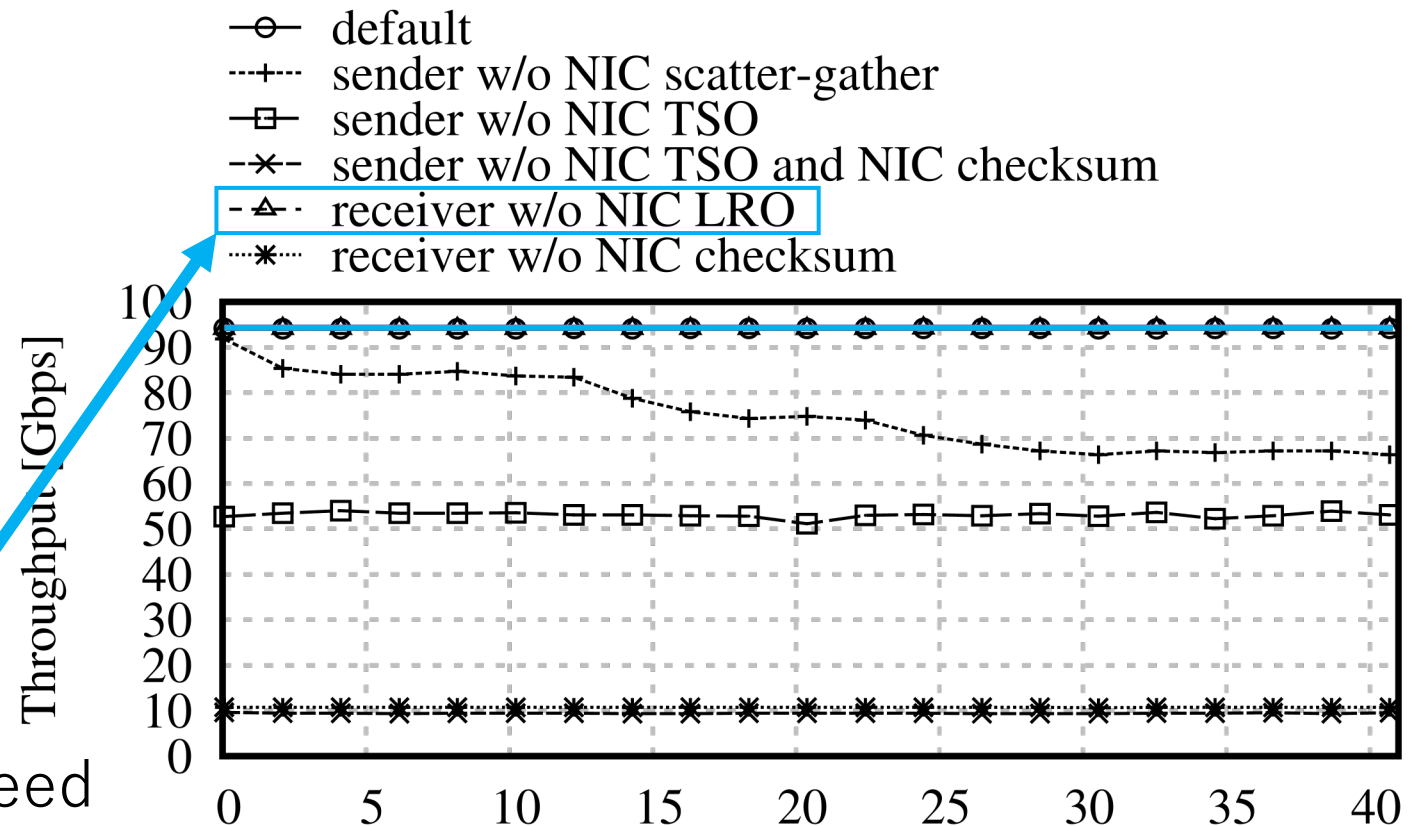  - TSO
  - Checksum offloading

When the receiver disables LRO,
it can still catch up with the link speed

Legend:
- default
- sender w/o NIC scatter-gather
- sender w/o NIC TSO
- sender w/o NIC TSO and NIC checksum
- receiver w/o NIC LRO
- receiver w/o NIC checksum

Throughput [Gbps] vs Data Size [MB]
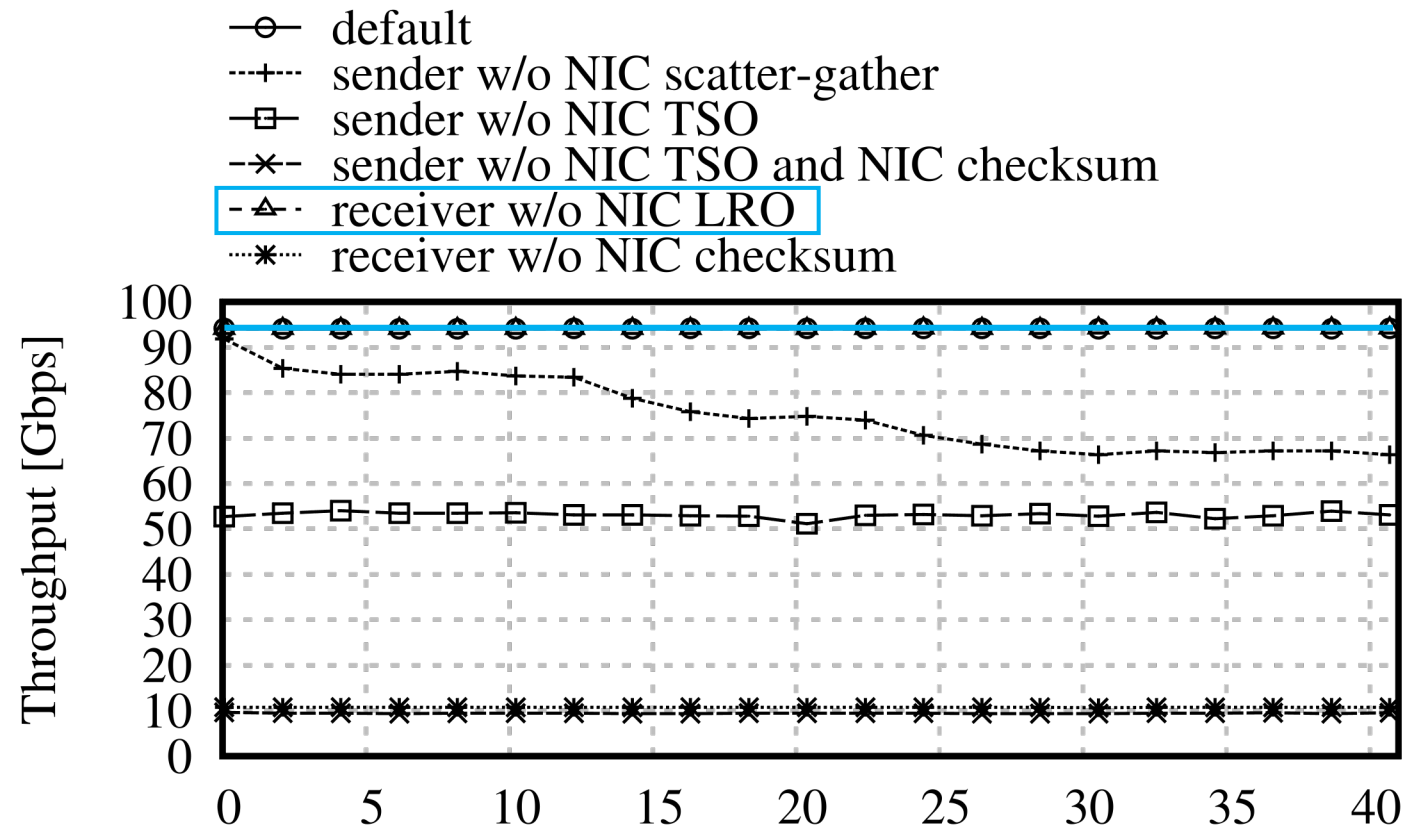
# Evaluation: Bulk Data Transfer

- The sender repeatedly sends the same data to the receiver

- Performance factors
  - Zero-copy transmission
  - TSO
  - Checksum offloading

- Note
  - This does not mean LRO is not necessary
    - Just we did not see differences in this workload



Legend:
- default
- sender w/o NIC scatter-gather
- sender w/o NIC TSO
- sender w/o NIC TSO and NIC checksum
- receiver w/o NIC LRO
- receiver w/o NIC checksum

Y-axis: Throughput [Gbps]
X-axis: Data Size [MB]

# Summary

- iip is a TCP/IP stack implementation that aims to allow for <u>easy integration</u> and <u>good performance</u> simultaneously

Please try it if you are interested

- Main page: https://github.com/yasukata/iip
- Assets used in the paper: https://github.com/yasukata/bench-iip/tree/9cf2488ec93ae51f4bd7b18923a5d1a233852f66