# Much Faster Networking

SOLARFLARE®

David Riddoch

driddoch@solarflare.com

# What is kernel bypass?

SOLARFLARE®

1. DMA

2. Interrupt

3. Protocol Stack

4. recv()

OS

App

1. DMA

OS

App

2. Doorbell

1. Form message

OS

App

1. PIO

OS

App

# What does all of this cleverness achieve?

# Better performance!

- Fewer CPU instructions for network operations

- Better cache locality
  - Faster response (lower latency)
  - Higher throughput (higher bandwidth/message rate)

- Reduced contention between threads
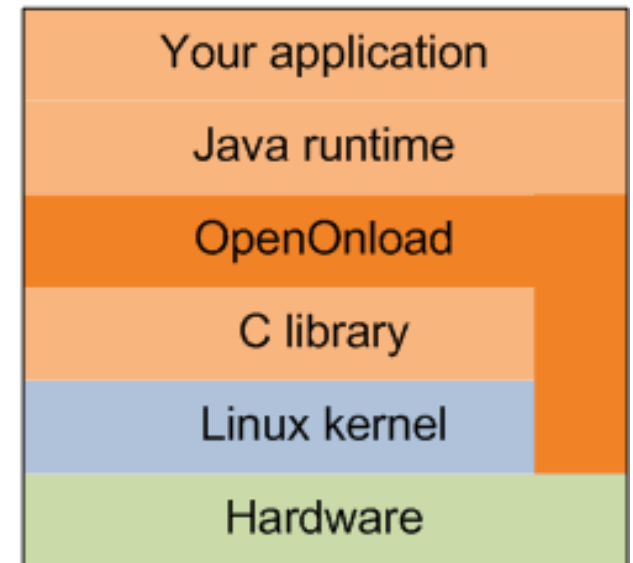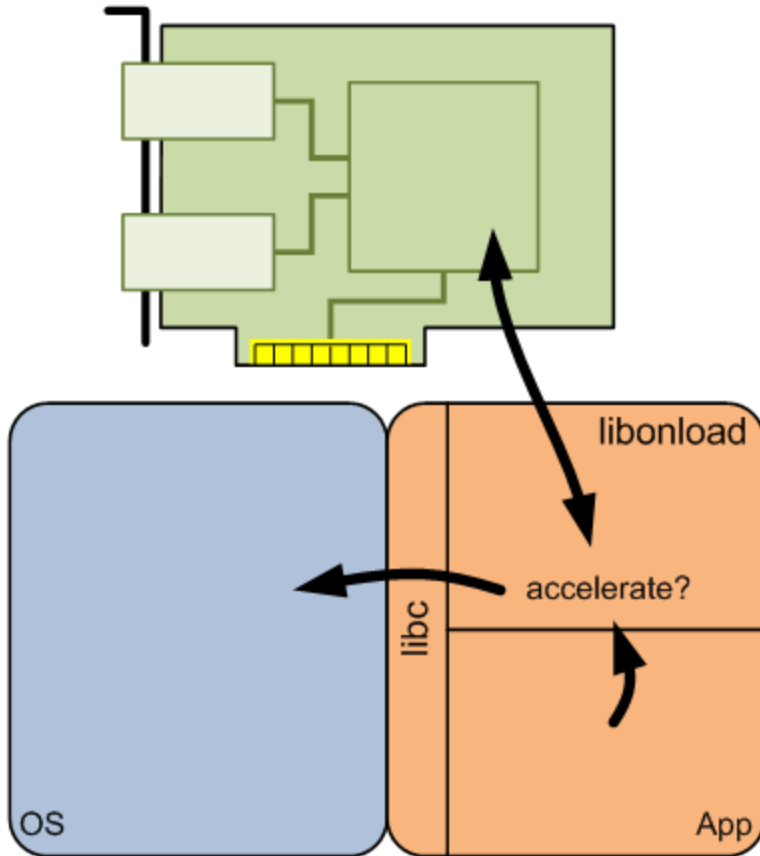  - Better core scaling
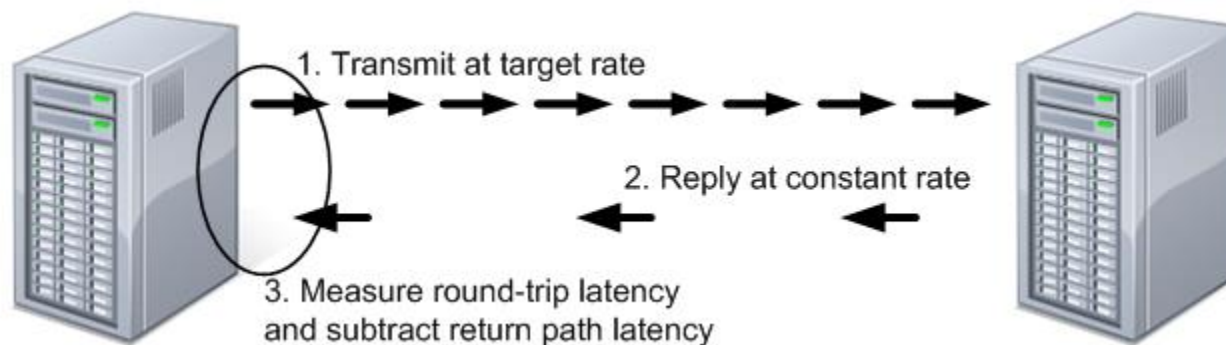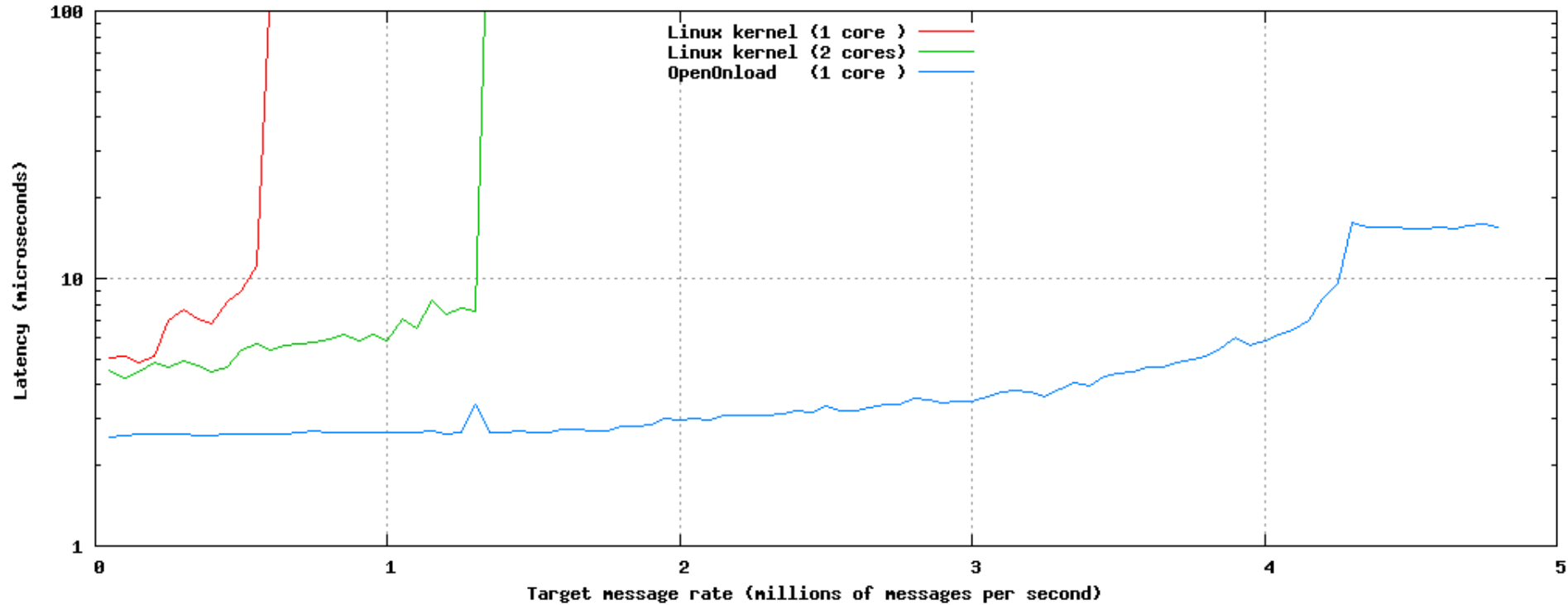  - Reduced latency jitter

SOLARFLARE®

- Sockets acceleration using kernel bypass

- Standard Ethernet, IP, TCP and UDP

- Standard BSD sockets API

- Binary compatible with existing applications

libonload

libc

accelerate?

OS

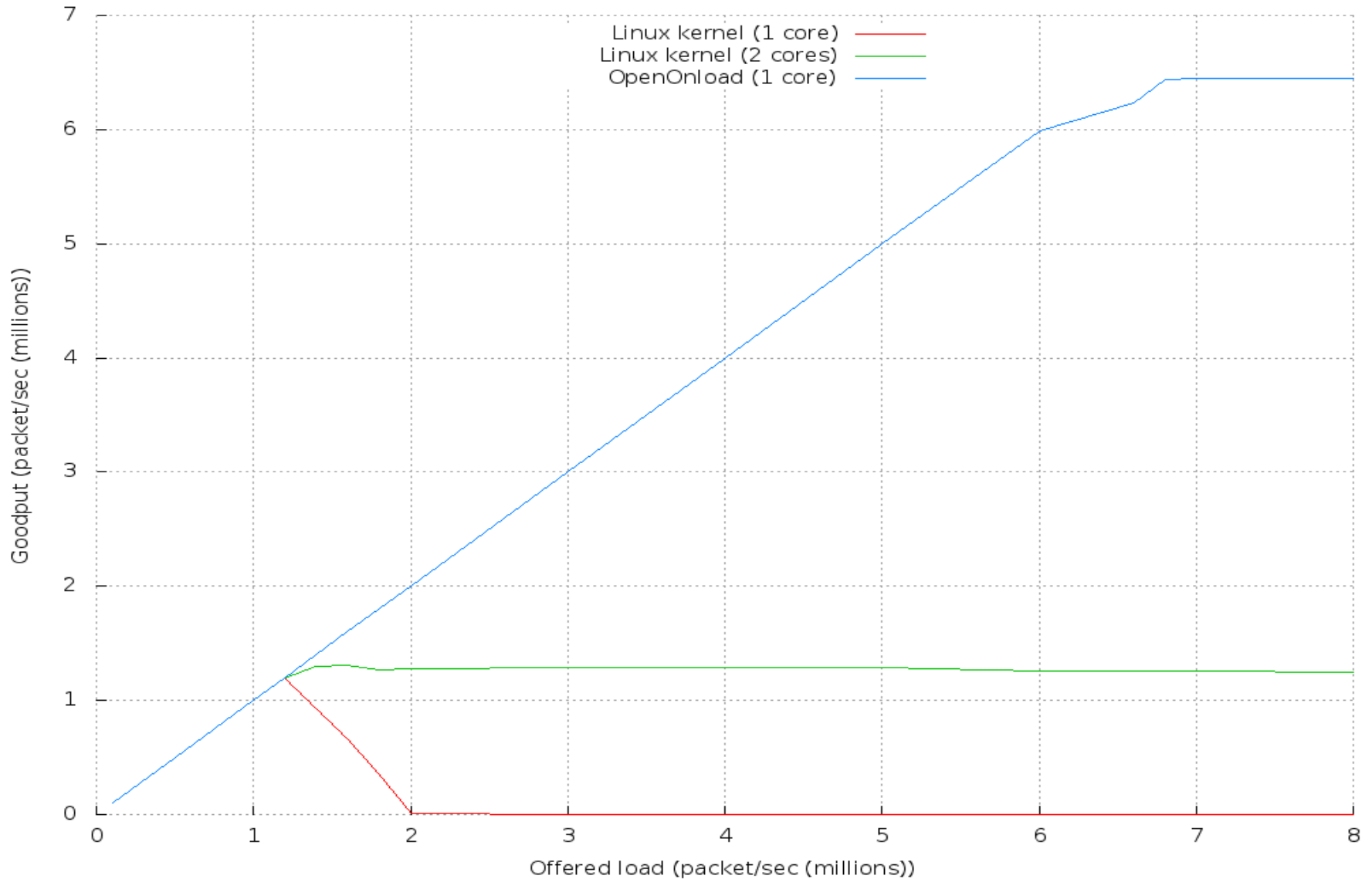App

Your application

Java runtime

OpenOnload

C library

Linux kernel

Hardware

# UDP receive throughput (small messages)

HAProxy

# HAProxy performance and scaling

SOLARFLARE

## 1 KiB message size



## 100 KiB message size

# Why doesn't performance scale when using the kernel stack?

# Better question:

## *How come it scales as well as it does?*

CPU cores

# Received packet is delivered into memory (or L3 cache)

Interrupt triggers packet
handling on this CPU core

Application calls recv()

Multiple receive channels
(up to one per core)

OS    App

channel_id = hash(4tuple) % n_cores;

Multiple flows

Hopefully!

Usually

```
channel_id, n = lookup(4tuple);
if( n > 1 )
    channel_id += hash(4tuple) % n_cores;
```

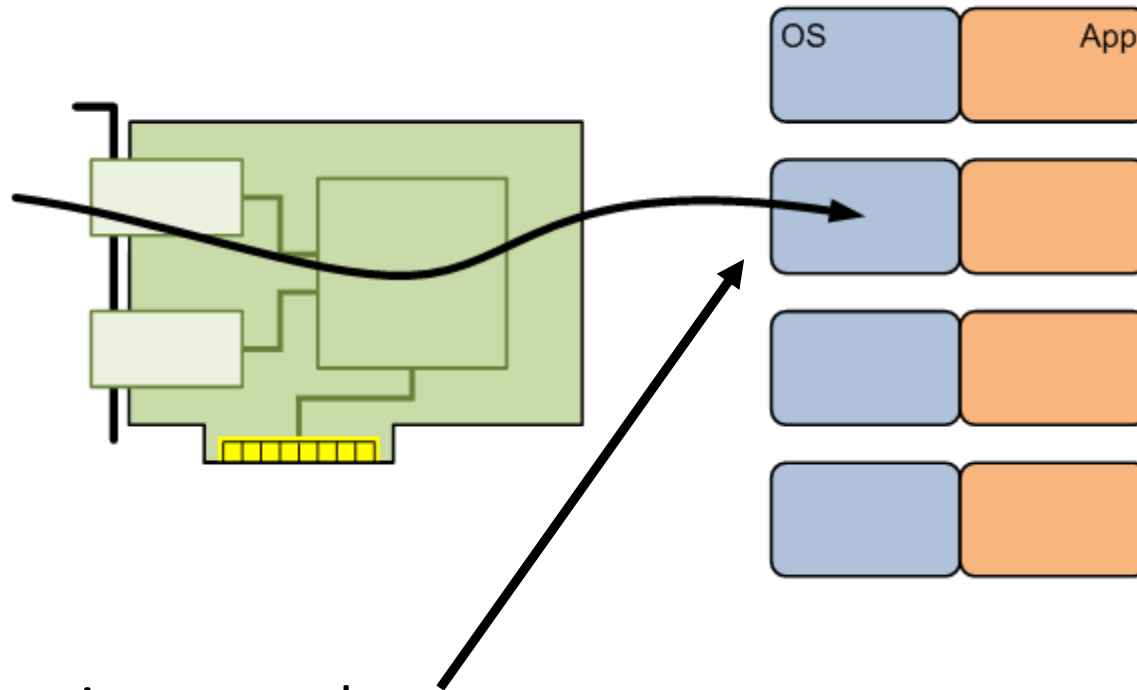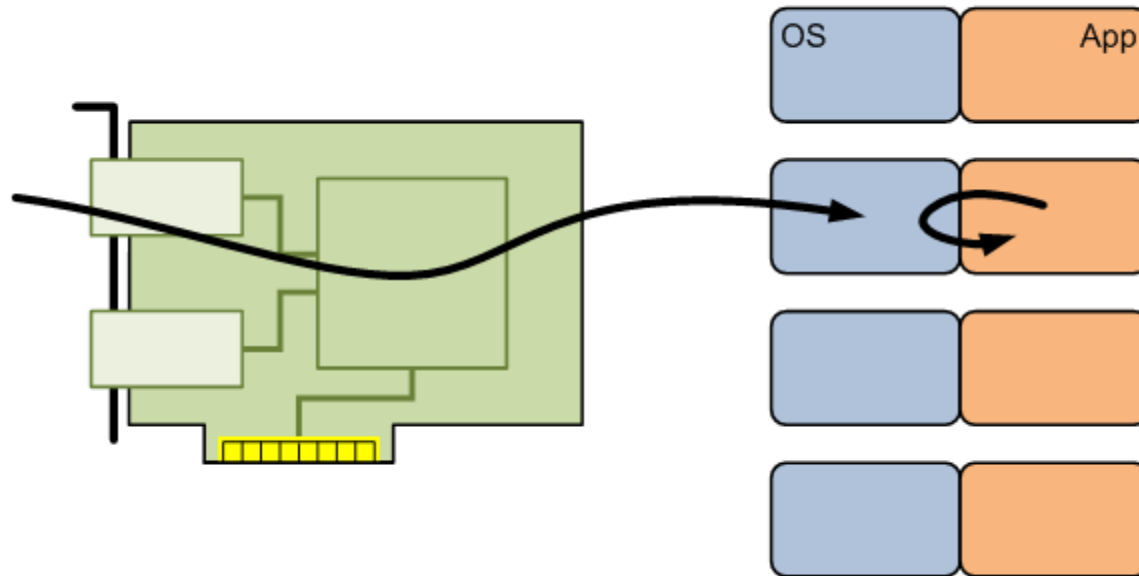- Multiple listening sockets on the same TCP port
  - One listening socket per worker thread
  - Each gets a subset of incoming connection requests
  - New connections go to the worker running on the core that the flow hashes to

- Connection establishment scales with the number of workers

- Received packets are delivered to the 'right' core

Problem solved?

QS

1 KiB message size

Number of CPU cores

| | |
|---|---|
| 5-7 | Application |
| | Middleware |
| 4 | TCP, UDP, ... |
| 3 | IPv4, IPv6 |
| 1-2 | Ethernet |

← Sockets

| 5-7 | Application |
| | Middleware |
| 4 | TCP, UDP, ... |
| 3 | IPv4, IPv6 |
| 1-2 | Ethernet |

← Sockets

1-6 Mpps/core
Many protocols
End-host applications

← Layer-2 APIs

1-**60** Mpps/core
**All** protocols
**Any** network function

# 60 million pkt/s?

# 17 ns

Some tips for achieving really fast networking...

Tip 1. The faster your application is already, the more speedup you'll get from kernel bypass

(This is just Ahmdal's law)

# Tip 2. NUMA locality applies doubly to I/O devices

- DMA transfers will use the L3 cache if:
  - The targeted cache line is resident
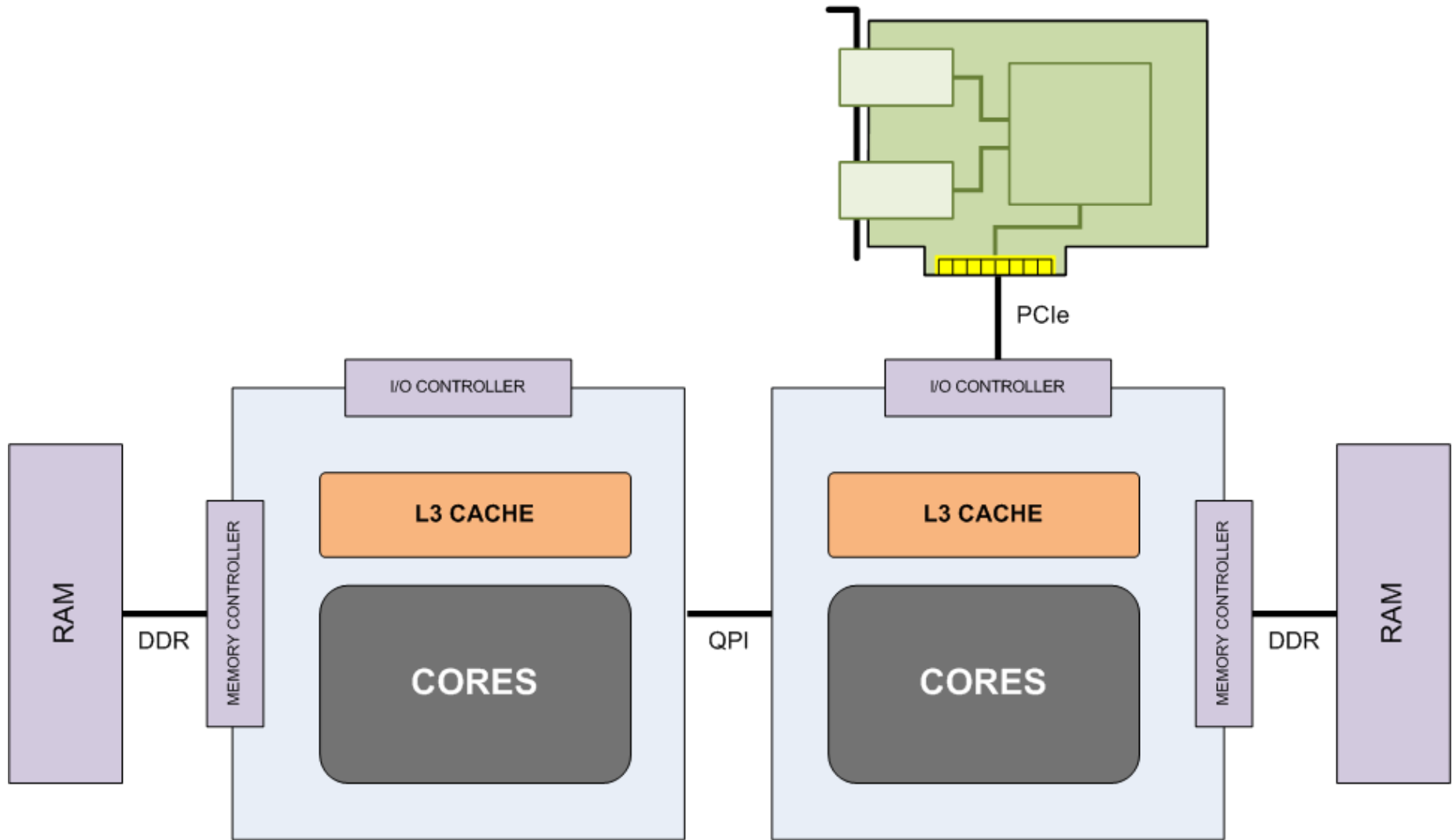  - Or if not then up to 10% of L3 is available for write-allocate

- Therefore
  - If you want consistent high performance, DMA buffers must be resident in L3 cache

- To achieve that
  - Small set of DMA buffers recycled quickly
  - (Even if that means doing an extra copy)

# Tip 3.
# Queue management is *critical*

Queues exist mostly to handle mismatch between *arrival rate* and *service rate*

- Buffers in switches and routers

- Descriptor rings in network adapters

- Socket send and receive buffers
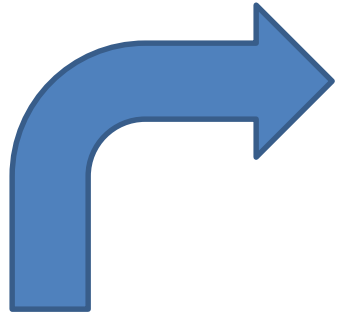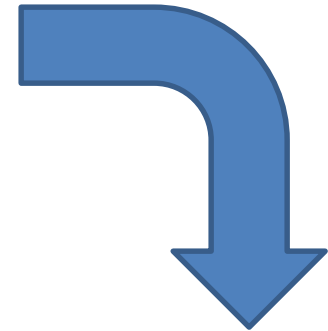
- Shared memory queues, locks etc.

- Run queue in the kernel task scheduler

# What happens when queues start to fill?

Service rate < arrival rate

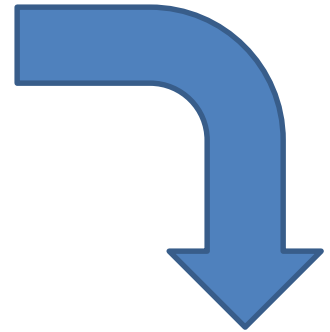Queue fill level increases
(latency++)

Service rate drops

Working-set size
increases
(efficiency--)

Service rate < arrival rate

Queue fills

DROP!

Drops are bad, m'kay

Make  SO_RCVBUF  bigger!

# Dilemma!

Small buffers:

Necessary for stable performance when overloaded

Large buffers:
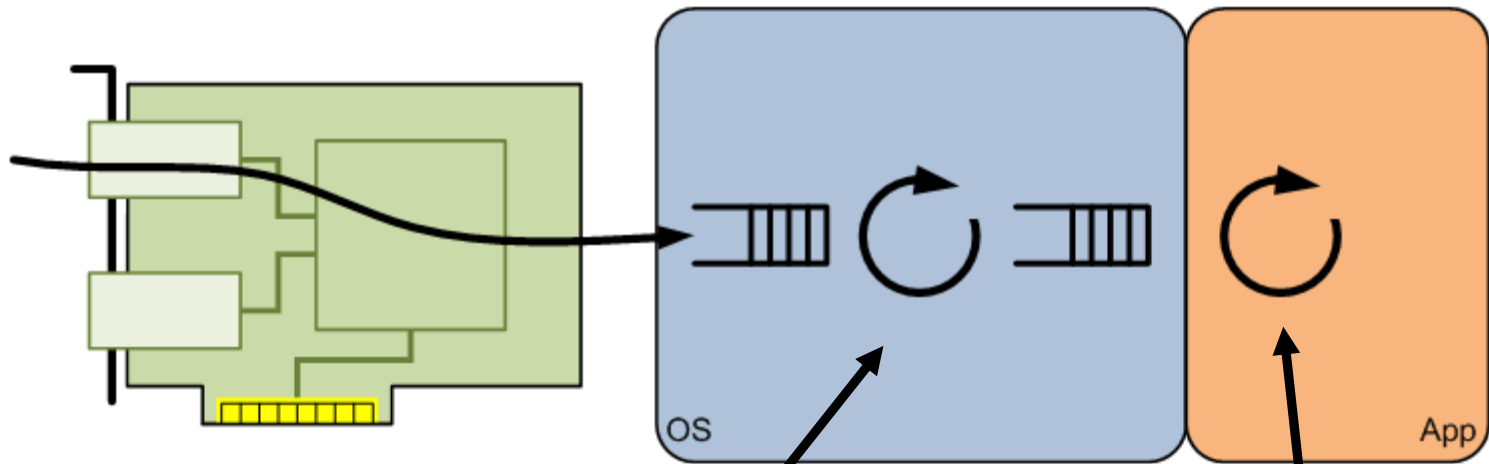
Necessary for absorbing bursts without loss

Limit working set size

- Limit the sizes of pools, queues, socket buffers etc.

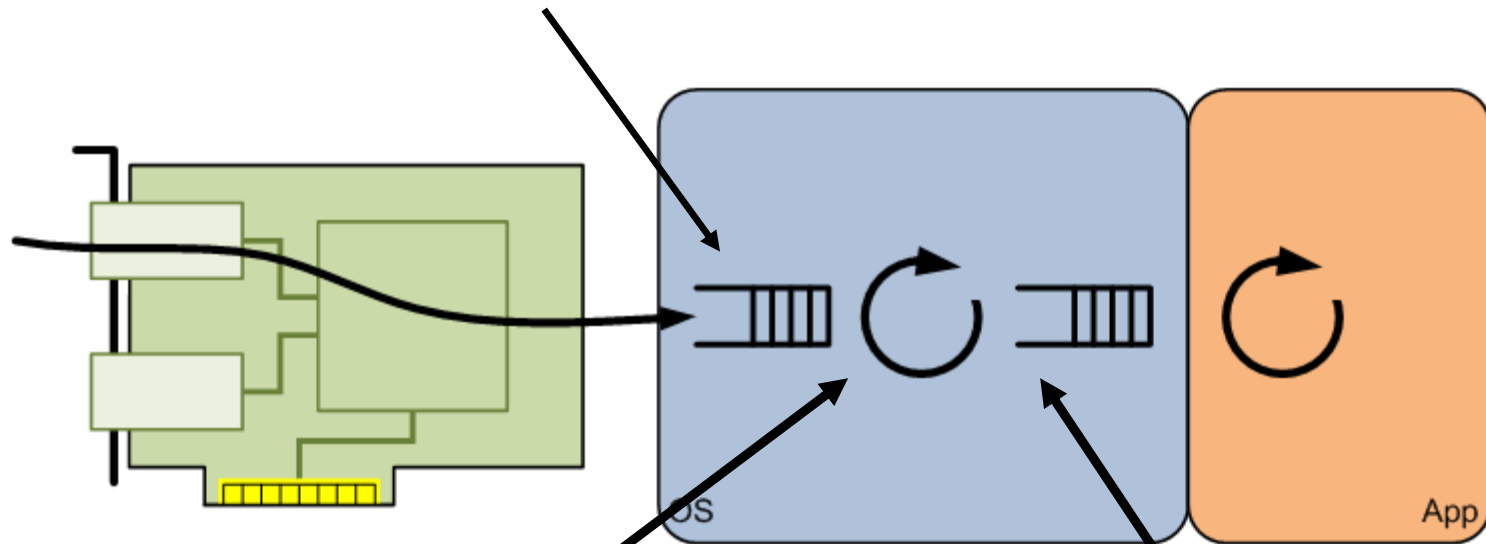Shed excess load early  (Tip 3.1)

- Before you've wasted time on requests you're going to have to drop anyway

Interrupt moves packets from descriptor ring to sockets
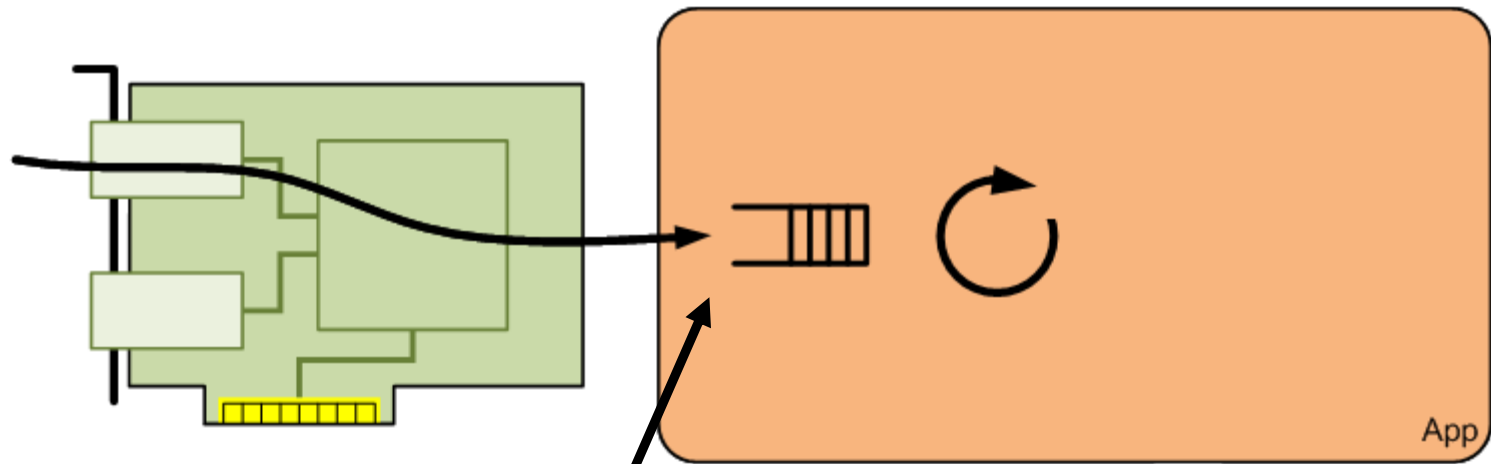
App thread consumes from socket buffer

Drop newest data
(only at very high rates)
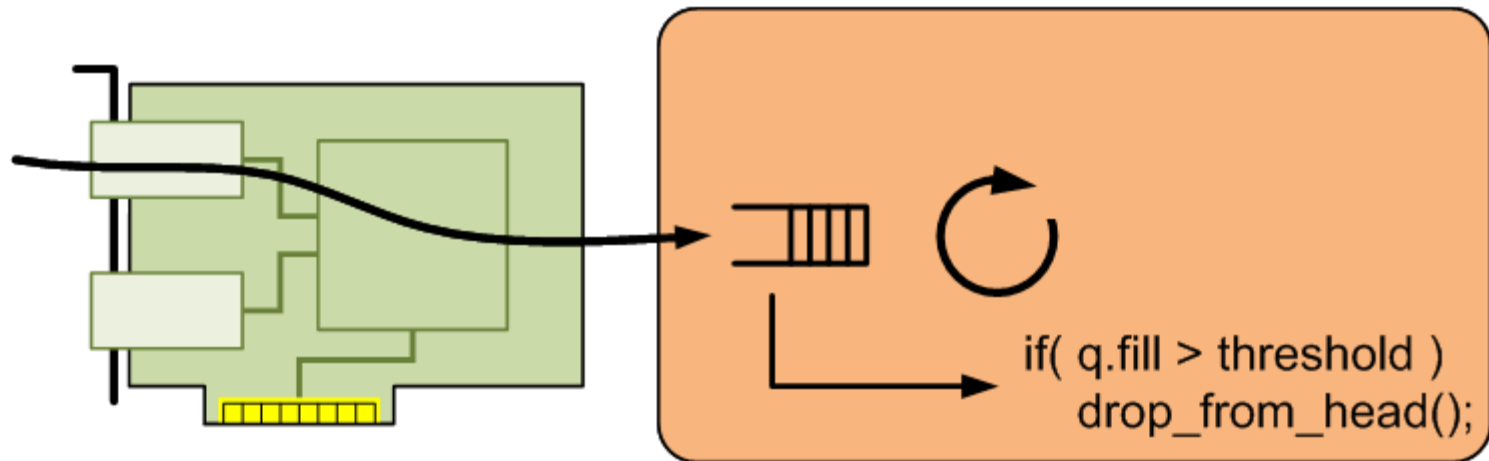
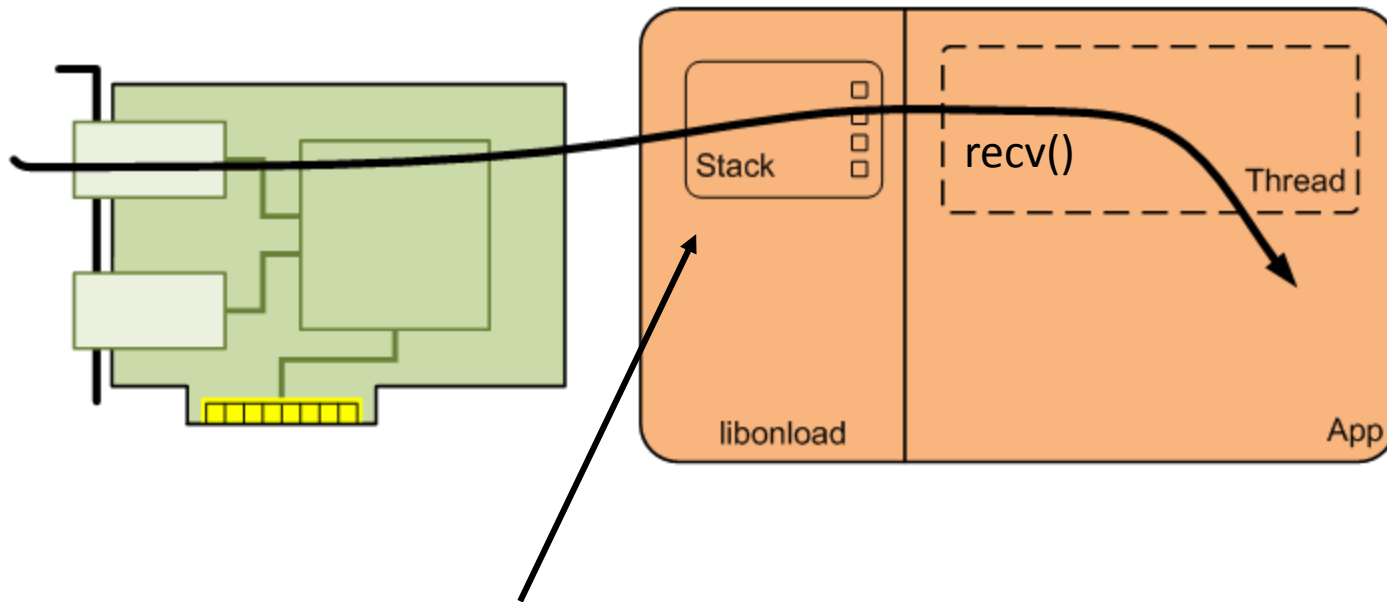Do work for every packet,
whether dropped or not

Drop newest data

Drop newest data

```
if( q.fill > threshold )
    drop_from_head();
```
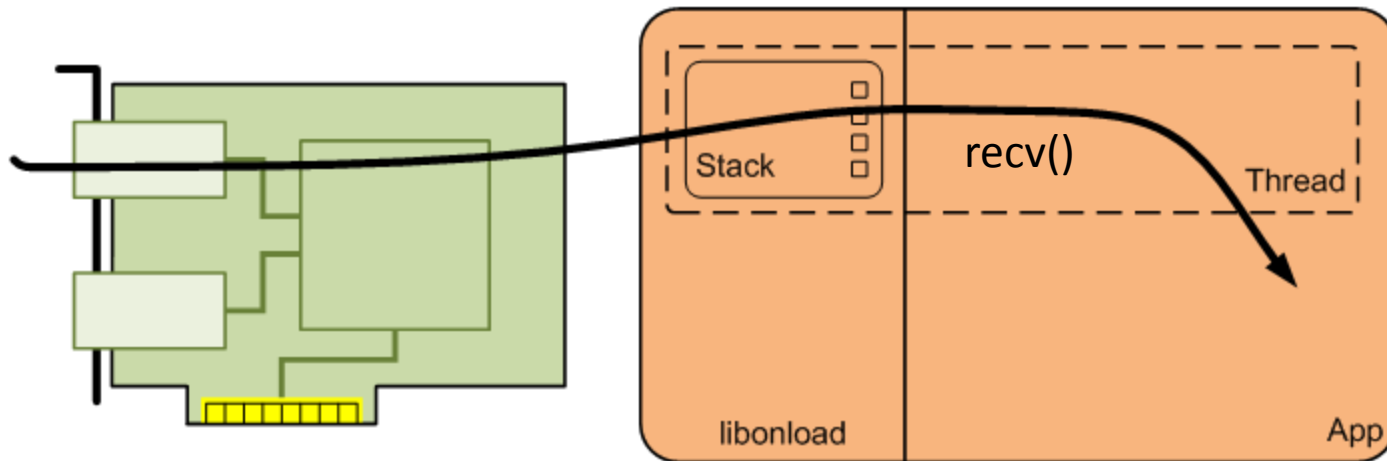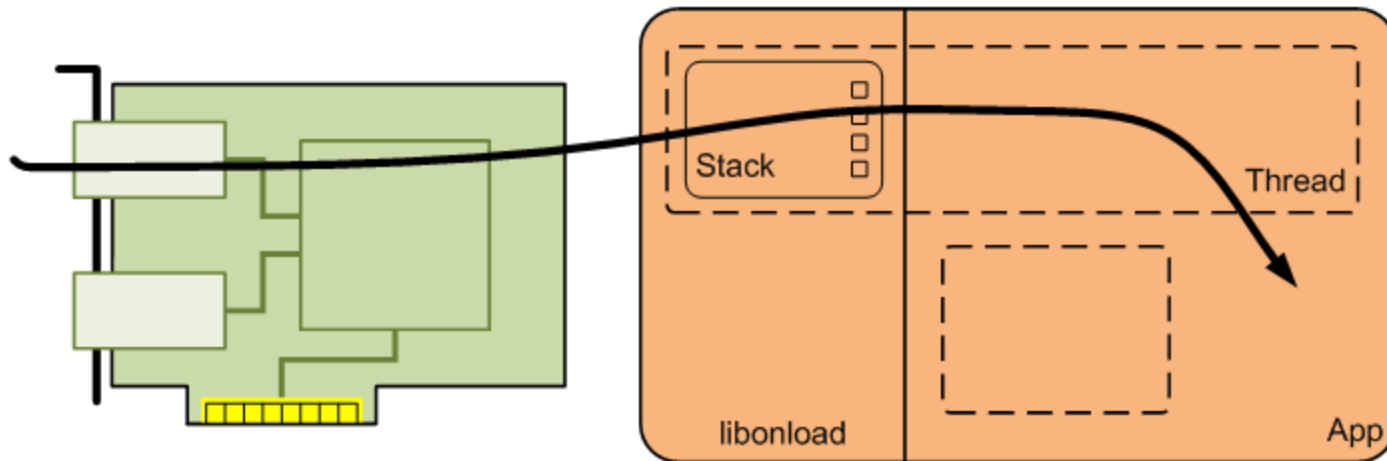
# Last example: Cache locality

OpenOnload stack; includes sockets, DMA buffers, TX and RX rings, control plane etc.
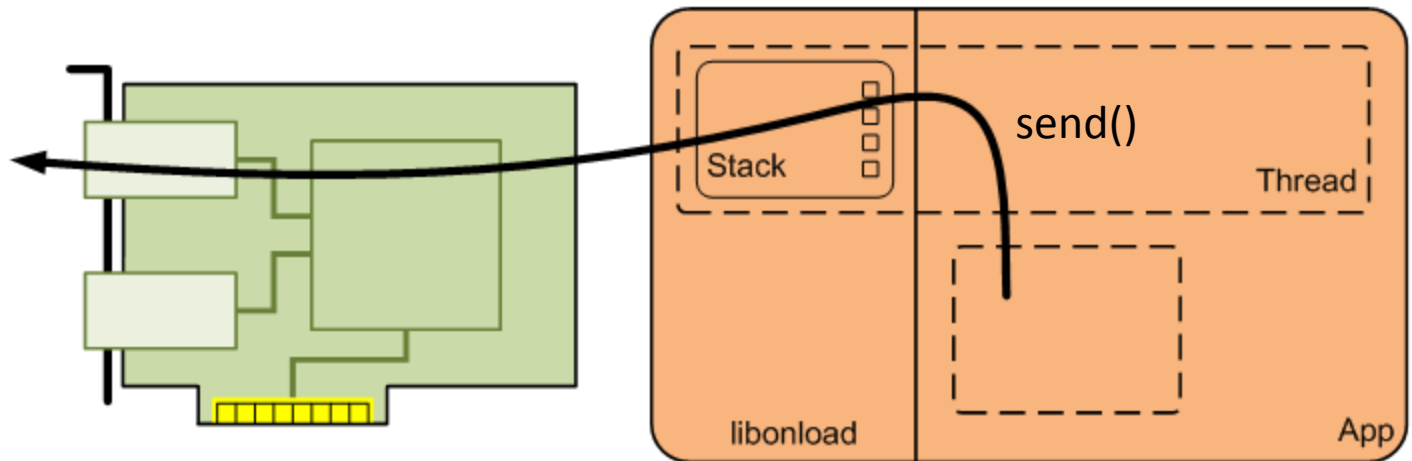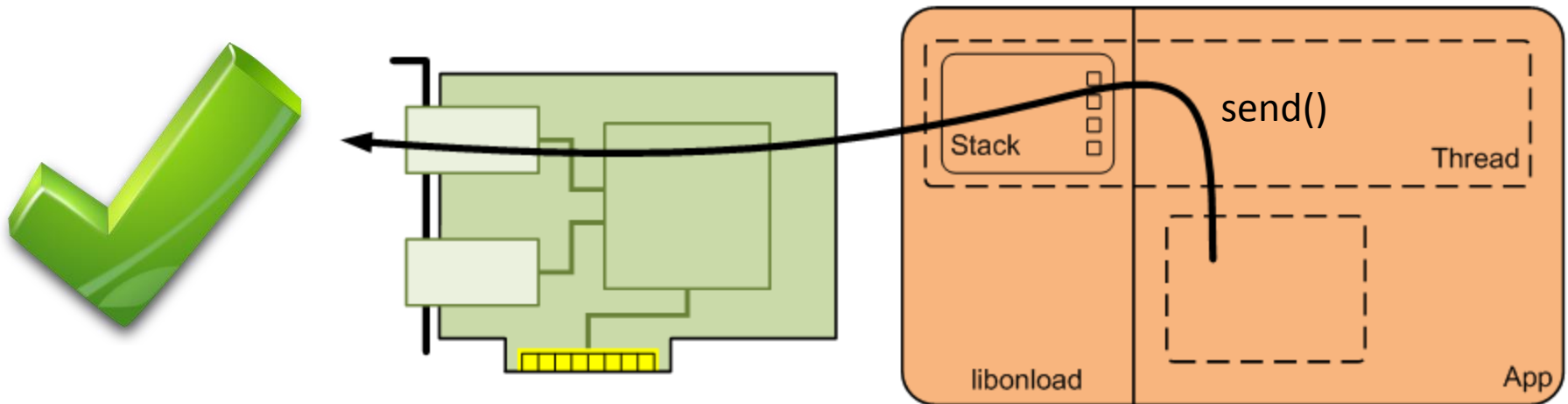
Stack

recv()

Thread

libonload

App

Problem: Send a small (eg. 200 bytes) reply from a different thread

Or…

- Passing a small message to another thread:
  - A few cache misses

- send() on a socket last accessed on another core:
  - Dozens of cache misses

# Thank you!