# PASTE: A Network Programming Interface for Non-Volatile Main Memory

**Michio Honda** (NEC Laboratories Europe)
Giuseppe Lettieri (Università di Pisa)
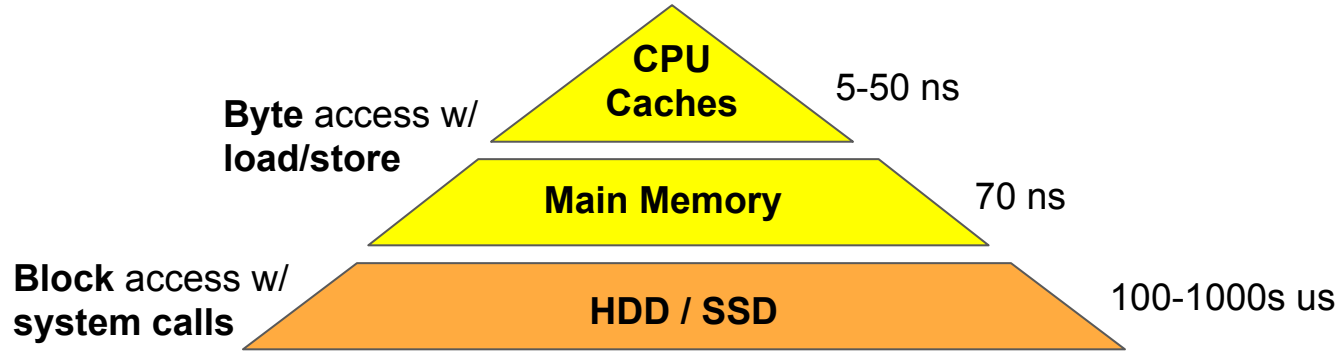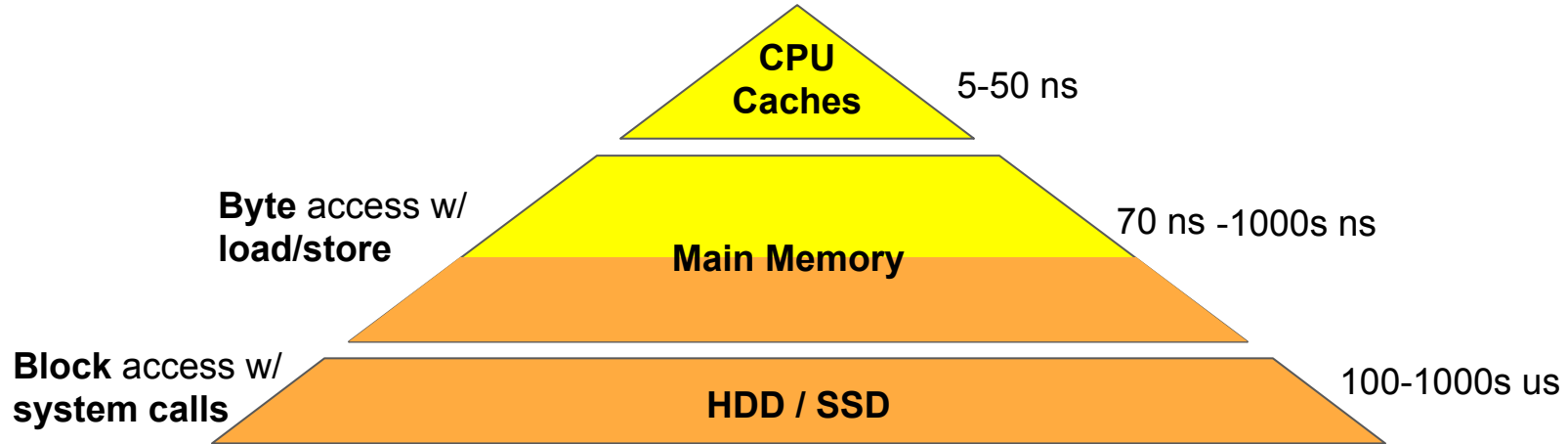Lars Eggert and Douglas Santry (NetApp)

USENIX NSDI 2018

SSICLOPS
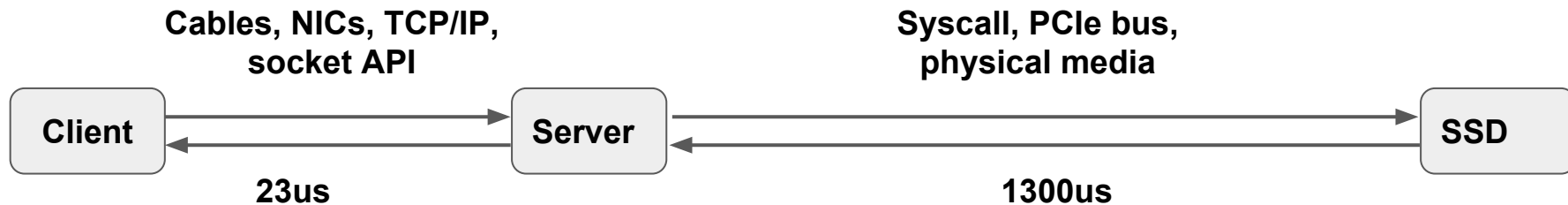
# Review: Memory Hierarchy

Slow, block-oriented persistence

**Byte** access w/ **load/store**

**Block** access w/ **system calls**

**CPU Caches** — 5-50 ns

**Main Memory** — 70 ns

**HDD / SSD** — 100-1000s us

# Review: Memory Hierarchy

*Fast, byte-addressable persistence*



**CPU Caches** — 5-50 ns

**Byte** access w/ **load/store**

**Main Memory** — 70 ns -1000s ns

**Block** access w/ **system calls**

**HDD / SSD** — 100-1000s us

# Networking is faster than disks/SSDs

## 1.2KB durable write over TCP/HTTP

**Cables, NICs, TCP/IP, socket API**

**Syscall, PCIe bus, physical media**

Client ⇄ Server ⇄ SSD

**23us**

**1300us**

# Networking is slower than NVMM

1.2KB durable write over TCP/HTTP

Cables, NICs, TCP/IP,     Memcpy, memory bus,
     socket API              physical media

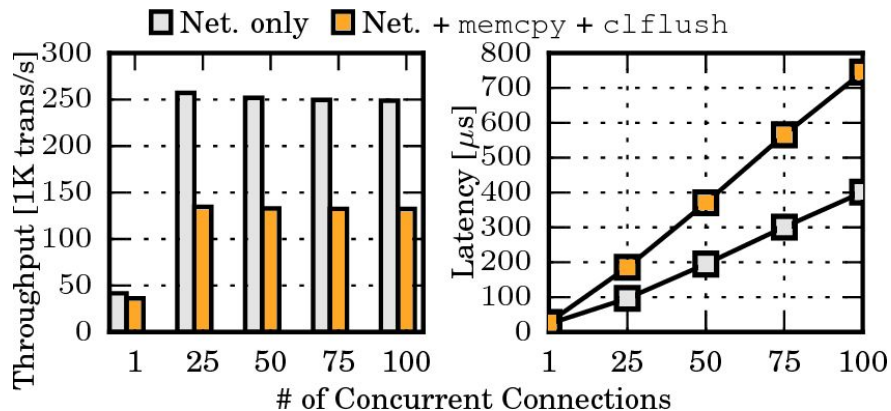**Client** ⟶ **Server** ⟶ **NVMM**

23us                     2us

# Networking is slower than NVMM

## 1.2KB durable write over TCP/HTTP

```
nevts = epoll_wait(fds)
for (i =0; i < nevts; i++) {
    read(fds[i], buf);
    ...
    memcpy(nvmm, buf);
    ...
    write(fds[i], reply)
}
```

**Cables, NICs, TCP/IP, socket API**

**Memcpy, memory bus, physical media**

Client → Server → NVMM

# Innovations at both stacks



**Network stack**

MegaPipe [OSDI'12]
Seastar
mTCP [NSDI'14]
IX [OSDI'14]
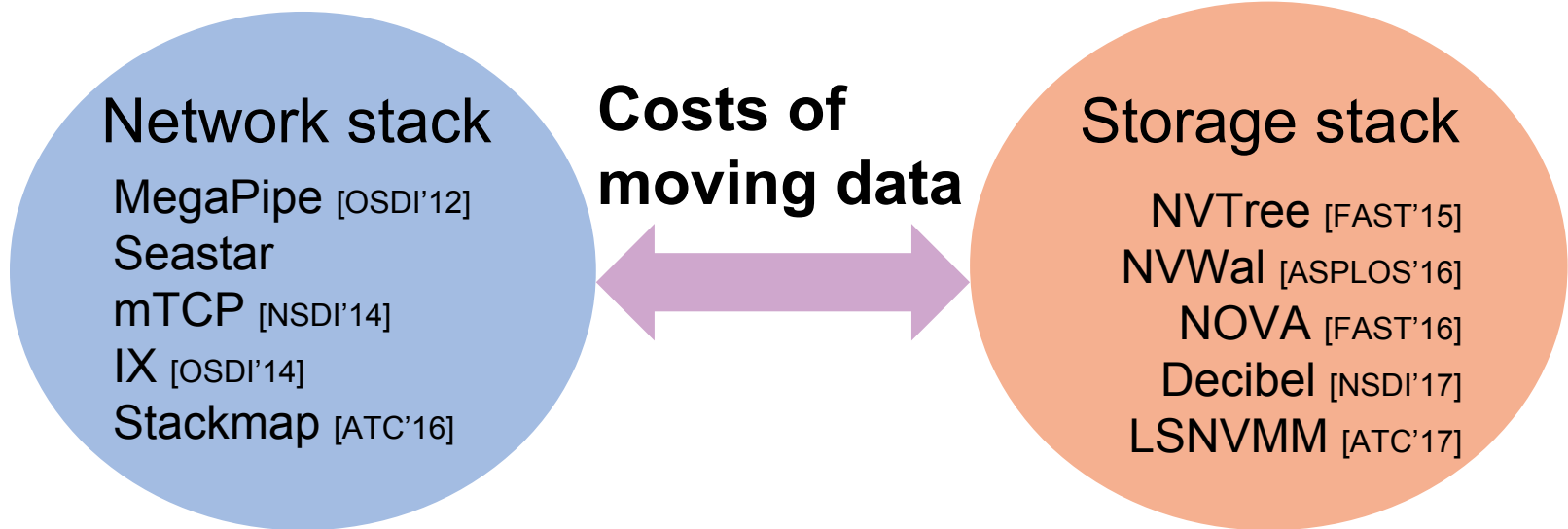Stackmap [ATC'16]

**Storage stack**

NVTree [FAST'15]
NVWal [ASPLOS'16]
NOVA [FAST'16]
Decibel [NSDI'17]
LSNVMM [ATC'17]

# Stacks are isolated



**Network stack**

MegaPipe [OSDI'12]
Seastar
mTCP [NSDI'14]
IX [OSDI'14]
Stackmap [ATC'16]

**Costs of moving data**

**Storage stack**

NVTree [FAST'15]
NVWal [ASPLOS'16]
NOVA [FAST'16]
Decibel [NSDI'17]
LSNVMM [ATC'17]

# Bridging the gap



Network stack

MegaPipe [OSDI'12]
Seastar
mTCP [NSDI'14]
IX [OSDI'14]
Stackmap [ATC'16]

**PASTE**

Storage stack
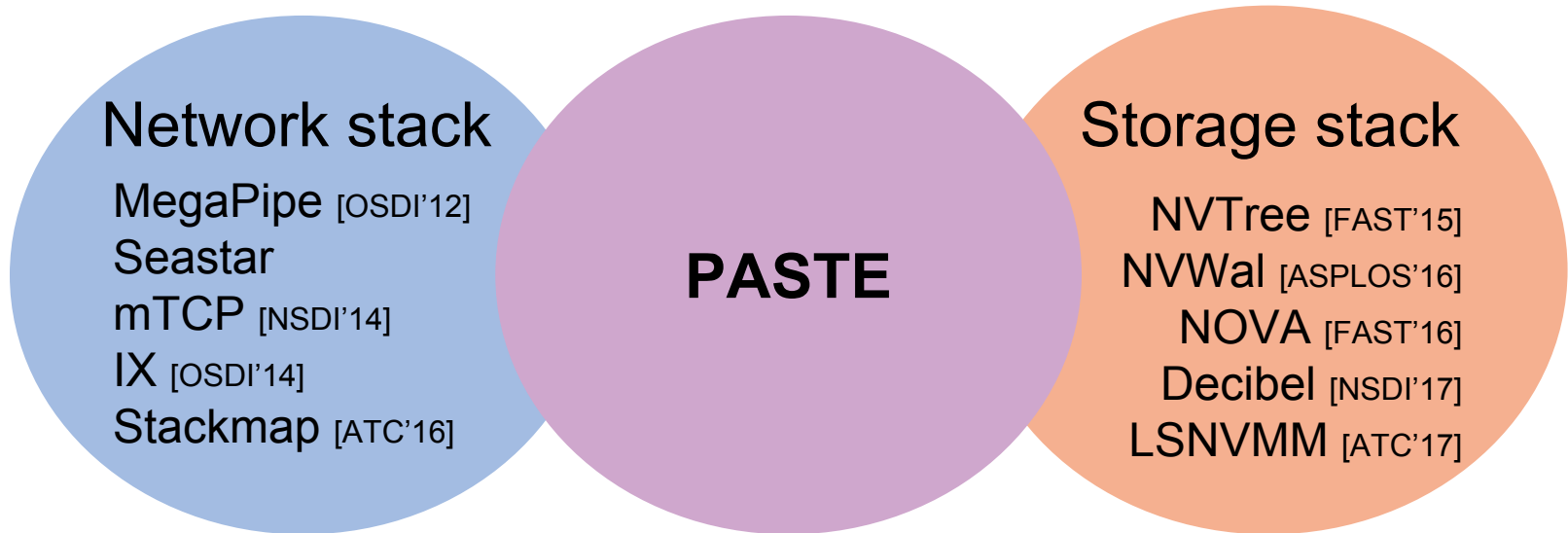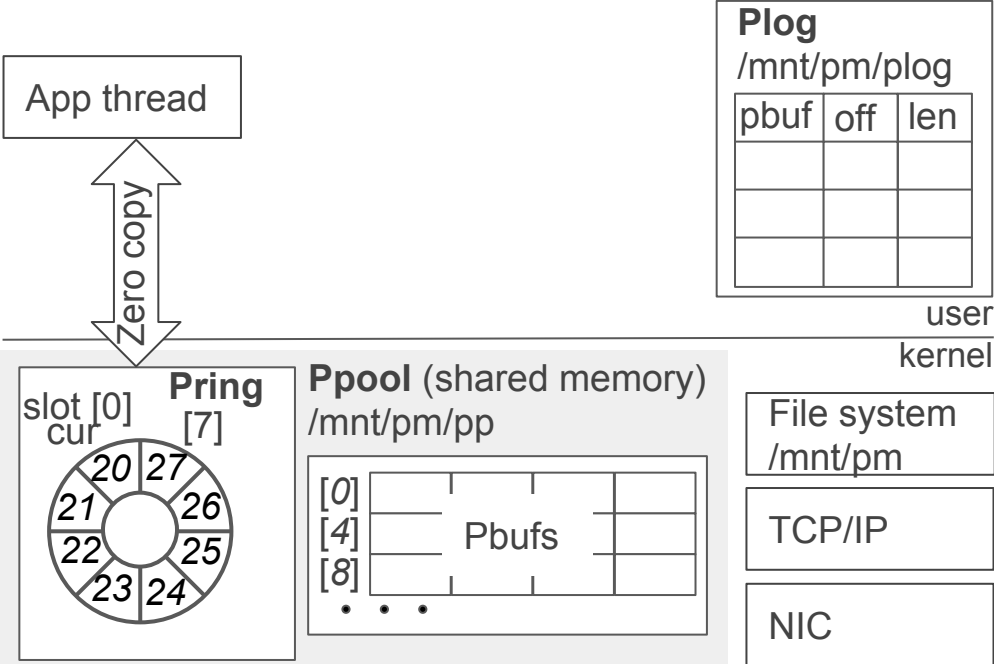
NVTree [FAST'15]
NVWal [ASPLOS'16]
NOVA [FAST'16]
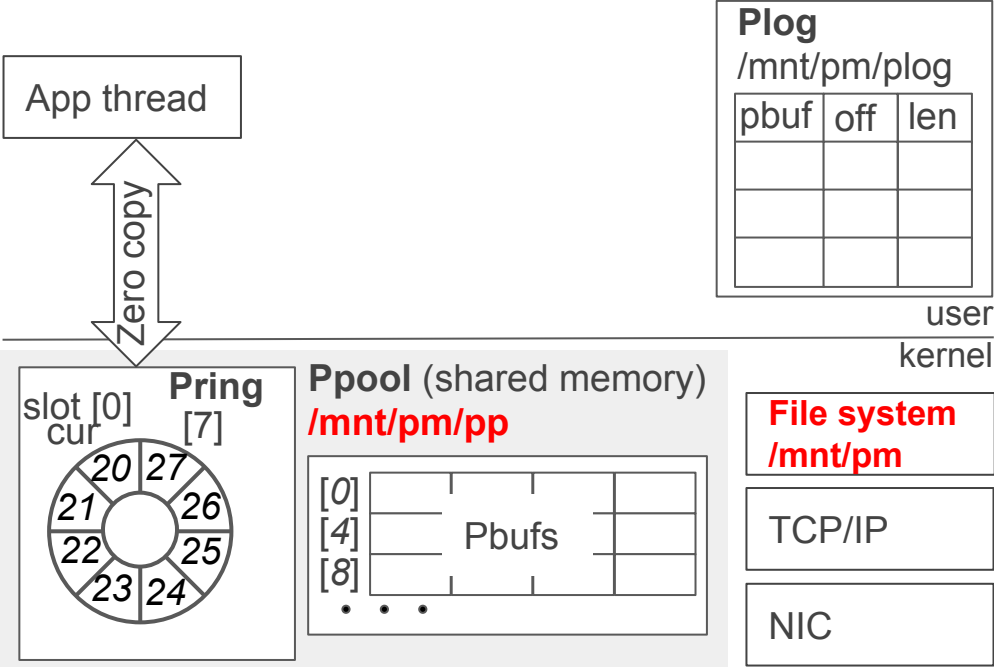Decibel [NSDI'17]
LSNVMM [ATC'17]

# PASTE Design Goals

- Durable zero copy
  - **DMA to NVMM**
- Selective persistence
  - **Exploit modern NIC's DMA to L3 cache**
- Persistent data structures
  - **Indexed, named packet buffers backed fy a file**
- Generality and safety
  - **TCP/IP** in the kernel **and netmap API**
- Best practices from modern network stacks
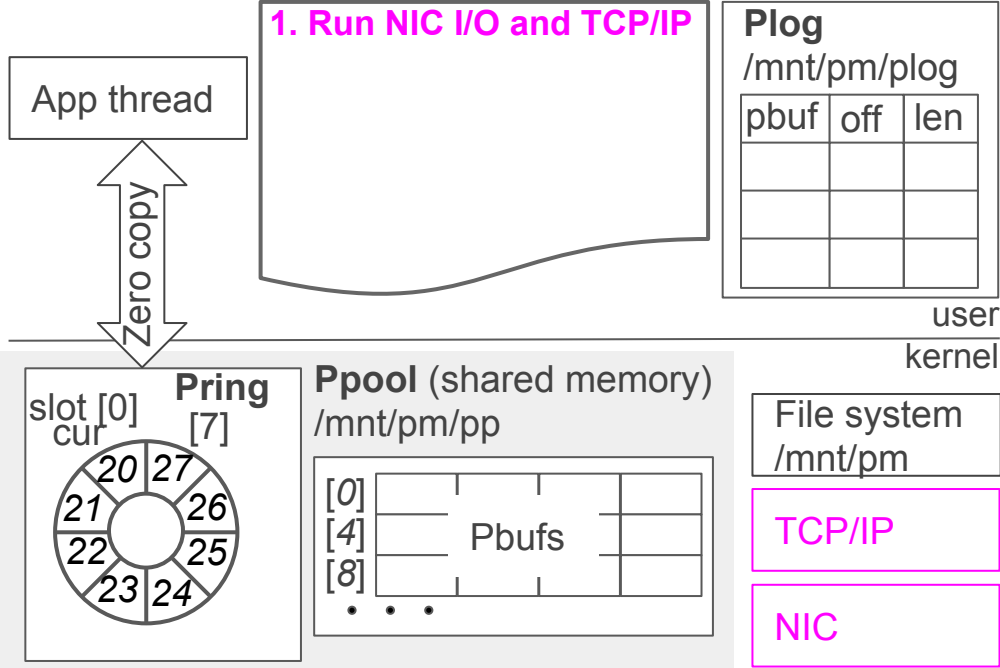  - Run-to-completion, blocking, busy-polling, batching etc

# PASTE in Action

App thread

Zero copy

**Plog**
/mnt/pm/plog

| pbuf | off | len |
|------|-----|-----|
|      |     |     |
|      |     |     |
|      |     |     |

user
kernel

**Pring**

slot [0]
cur

[7]

20 27
21 26
22 25
23 24

**Ppool** (shared memory)
/mnt/pm/pp

[0]
[4]    Pbufs
[8]

• • •

File system
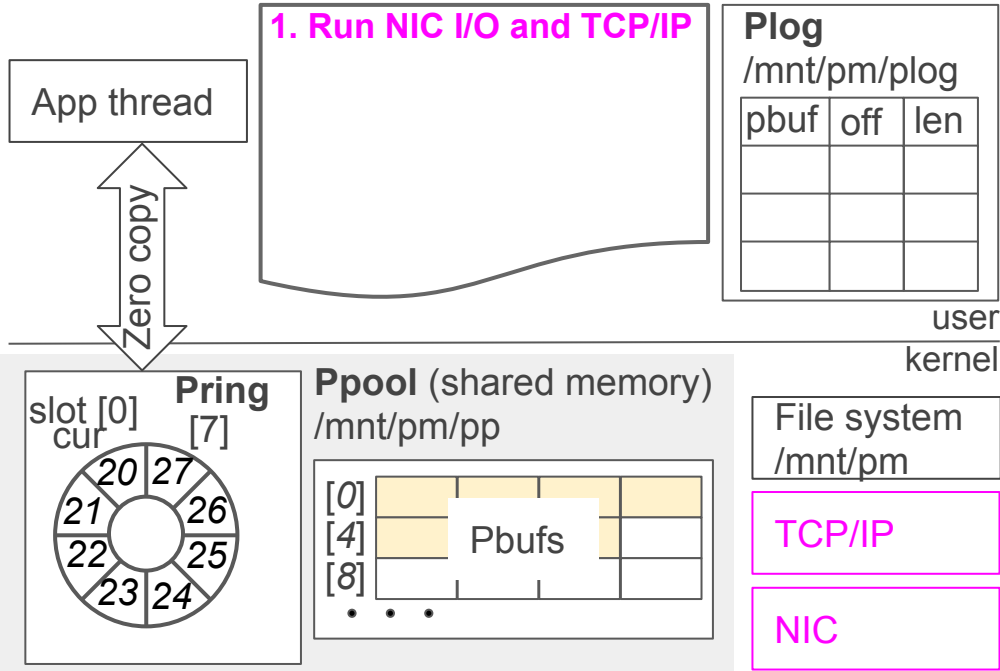/mnt/pm

TCP/IP
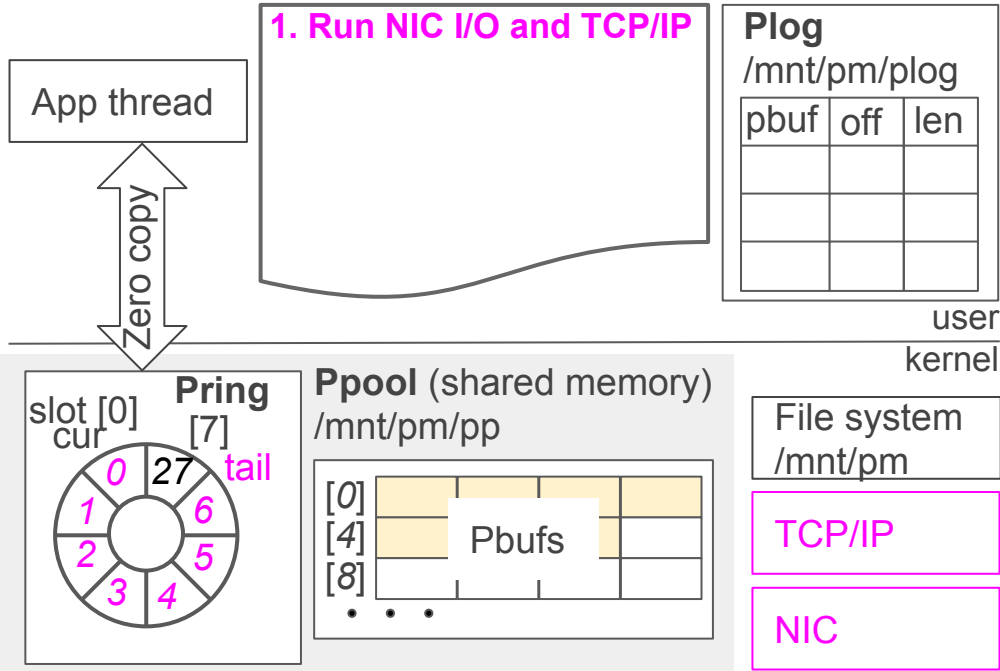
NIC

# PASTE in Action

# PASTE in Action
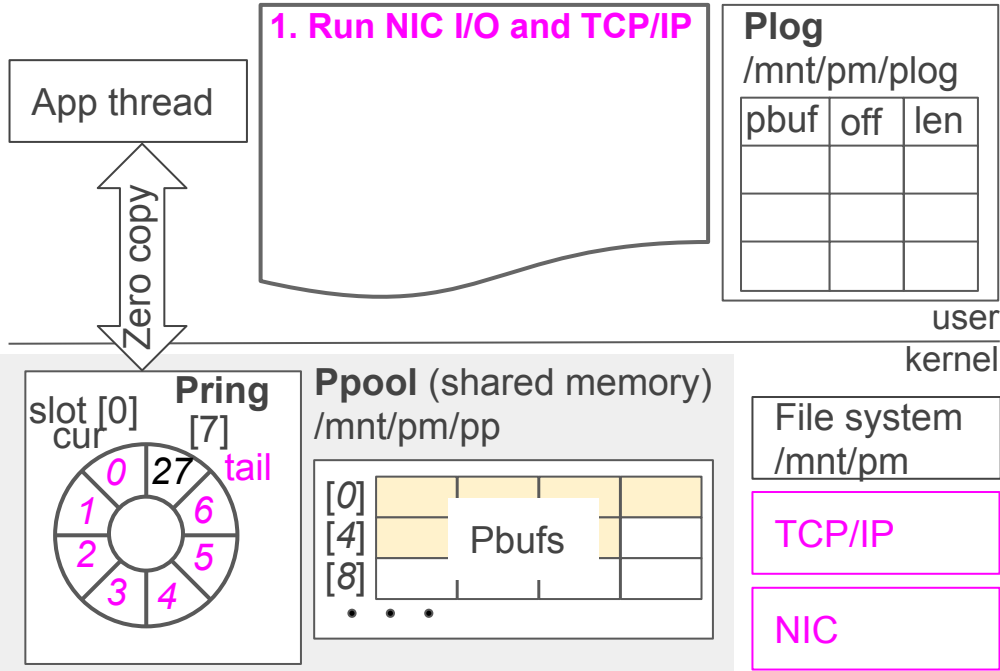


- poll() system call

# PASTE in Action



- poll() system call
  - Got 6 in-order TCP segments

# PASTE in Action

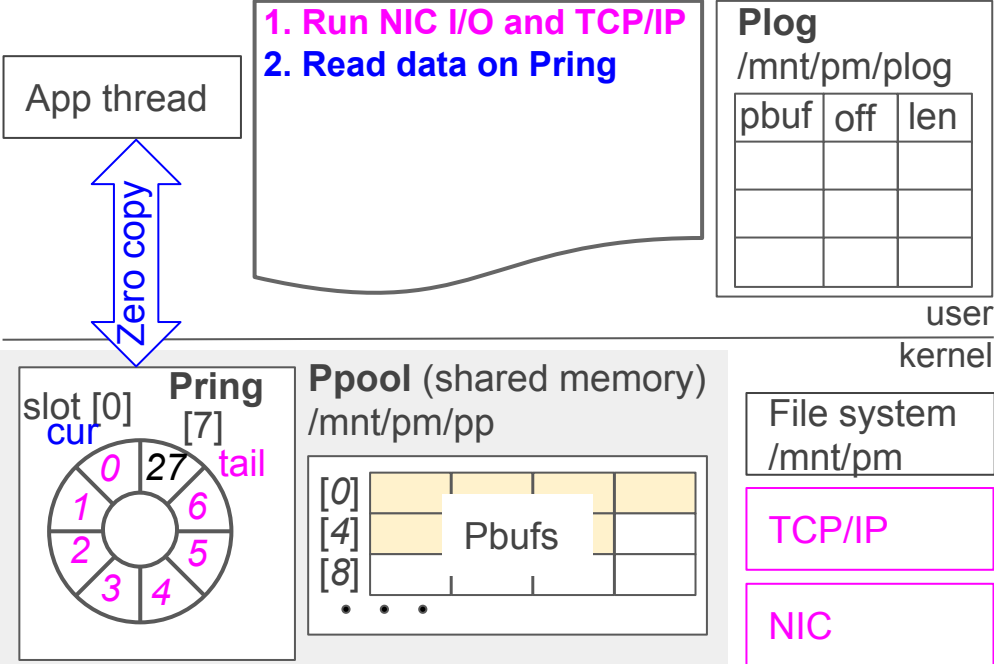

- poll() system call
  - They are set to Pring slots

# PASTE in Action



- Return from poll()

# PASTE in Action

App thread

**1. Run NIC I/O and TCP/IP**
**2. Read data on Pring**

**Plog**
/mnt/pm/plog

| pbuf | off | len |
|------|-----|-----|
|      |     |     |
|      |     |     |
|      |     |     |

Zero copy

user
kernel

**Pring**

slot [0]   [7]
cur       tail

0   27
1   6
2   5
3   4

**Ppool** (shared memory)
/mnt/pm/pp

[0]
[4]   Pbufs
[8]

• • •

File system
/mnt/pm

TCP/IP

NIC

# PASTE in Action

App thread

Zero copy

1. **Run NIC I/O and TCP/IP**
2. **Read data on Pring**
3. **Flush Pbuf(s)**

**Plog**
/mnt/pm/plog

| pbuf | off | len |
|------|-----|-----|
|      |     |     |
|      |     |     |
|      |     |     |

user
kernel

**Pring**

slot [0]      [7]

cur

tail

0    27

1         6

2         5

3    4

**Ppool** (shared memory)
/mnt/pm/pp

[0]
[4]    Pbufs
[8]

• • •

File system
/mnt/pm

TCP/IP

NIC

- **flush Pbuf data from CPU cache to DIMM**
  - clflush(opt) instruction

# PASTE in Action



**App thread**

1. **Run NIC I/O and TCP/IP**
2. **Read data on Pring**
3. **Flush Pbuf(s)**
4. **Flush Plog entry(ies)**

**Plog**
/mnt/pm/plog

| pbuf | off | len |
|------|-----|-----|
| 1 | 96 | 120 |
| | | |
| | | |

user / kernel

Zero copy

**Pring**
slot [0]
cur
[7]
0  27  tail
1  6
2  5
3  4

**Ppool** (shared memory)
/mnt/pm/pp

[0]
[4]  Pbufs
[8]
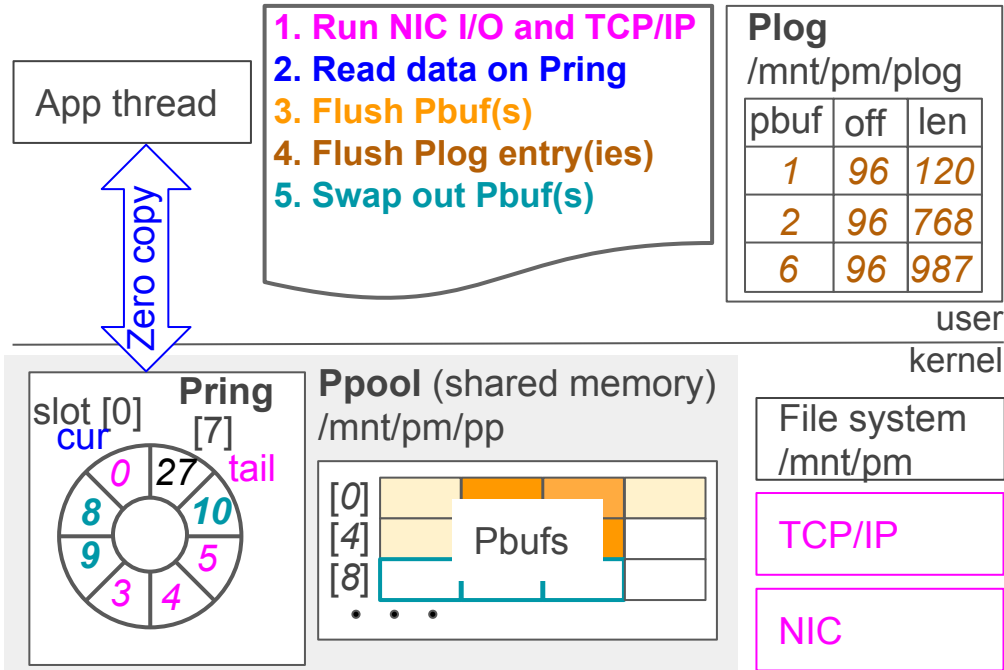• • •

**File system**
/mnt/pm

**TCP/IP**

**NIC**

- **Pbuf is persistent data representation**
  - Base address is static i.e., file (/mnt/pm/pp)
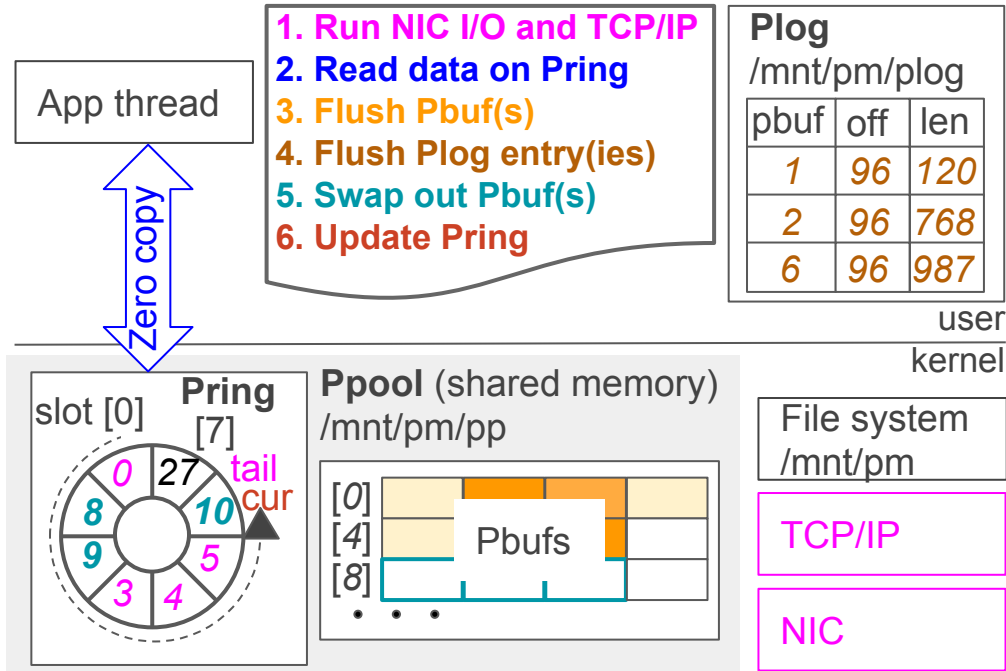  - Buffers can be recovered after reboot

# PASTE in Action



App thread

Zero copy

1. **Run NIC I/O and TCP/IP**
2. **Read data on Pring**
3. **Flush Pbuf(s)**
4. **Flush Plog entry(ies)**
5. **Swap out Pbuf(s)**

**Plog**
/mnt/pm/plog

| pbuf | off | len |
|------|-----|-----|
| 1 | 96 | 120 |
| | | |
| | | |

user
kernel

slot [0]
**Pring**
[7]
cur
0    27    tail
8         6
2         5
3    4

**Ppool** (shared memory)
/mnt/pm/pp

[0]
[4]    Pbufs
[8]
• • •

File system
/mnt/pm

TCP/IP

NIC

- **Prevent the kernel from recycling the buffer**

# PASTE in Action



1. **Run NIC I/O and TCP/IP**
2. **Read data on Pring**
3. **Flush Pbuf(s)**
4. **Flush Plog entry(ies)**
5. **Swap out Pbuf(s)**

App thread

Zero copy

**Plog**
/mnt/pm/plog

| pbuf | off | len |
|------|-----|-----|
| 1 | 96 | 120 |
| 2 | 96 | 768 |
| 6 | 96 | 987 |

user
kernel

slot [0]
**Pring**
[7]
cur
0  27  tail
8  10
9  5
3  4

**Ppool** (shared memory)
/mnt/pm/pp

[0]
[4]  Pbufs
[8]
• • •

File system
/mnt/pm

TCP/IP

NIC

- Same for Pbuf 2 and 6

# PASTE in Action

App thread

Zero copy

1. **Run NIC I/O and TCP/IP**
2. **Read data on Pring**
3. **Flush Pbuf(s)**
4. **Flush Plog entry(ies)**
5. **Swap out Pbuf(s)**
6. **Update Pring**

**Plog**
/mnt/pm/plog

| pbuf | off | len |
|------|-----|-----|
| 1 | 96 | 120 |
| 2 | 96 | 768 |
| 6 | 96 | 987 |

user

kernel

**Pring**

slot [0]    [7]

0   27   tail
8    10   cur
9    5
3   4

**Ppool** (shared memory)
/mnt/pm/pp

[0]
[4]   Pbufs
[8]
• • •

File system
/mnt/pm

TCP/IP

NIC

● Advance cur
  ○ Return buffers in slot 0-6 to the kernel at next poll()

# PASTE in Action

# PASTE in Action

**Plog**
/mnt/pm/plog

| 3 | 5 | |
|---|---|---|
|   |   | |

| 0 | | |   | 5 | | | | 7 |
|---|---|---|---|---|---|---|---|---|

**B+tree**

1. **Run NIC I/O and TCP/IP**
2. **Read data on Pring**
3. **Flush Pbuf(s)**
4. **Flush Plog entry(ies)**
5. **Swap out Pbuf(s)**
6. **Update Pring**

(*1*, 96, 120)

(2, 96, 987)

(6, 96, 512)

App thread

Zero copy

user
kernel

slot [0]  **Pring**  **Ppool** (shared memory)
[7]       /mnt/pm/pp

0   27   tail
8        cur
          10
9        5
3   4

[0]
[4]   Pbufs
[8]

• • •

File system
/mnt/pm

TCP/IP

NIC

- We can organize various data structures in **Plog**

# Evaluation

1. How does PASTE outperform existing systems?
2. Is PASTE applicable to existing applications?
3. Is PASTE useful for systems other than file/DB storage?
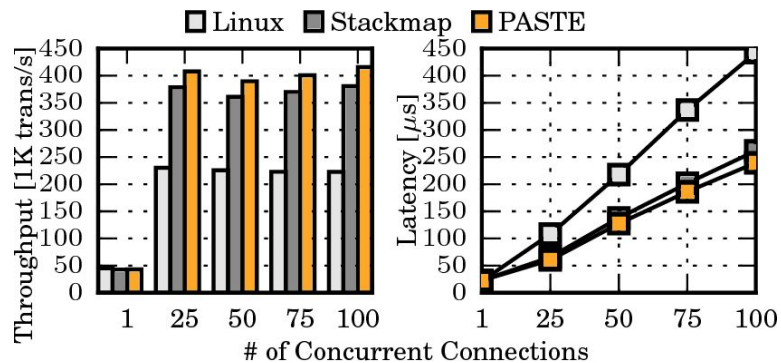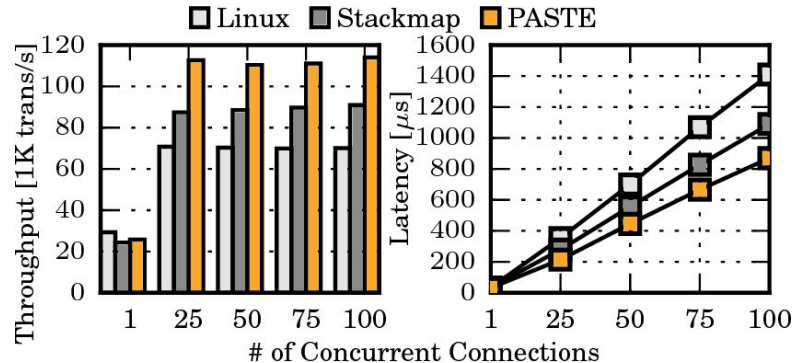
# How does PASTE outperform existing systems?



**64B**
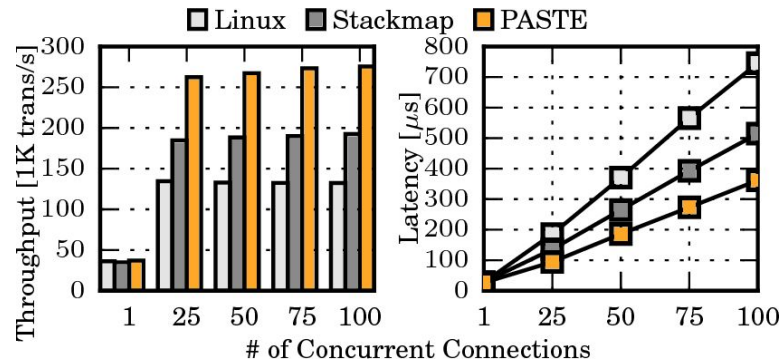
**1280B**

**WAL**

**What if we use more complex data structures?**

# How does PASTE outperform existing systems?



**64B**
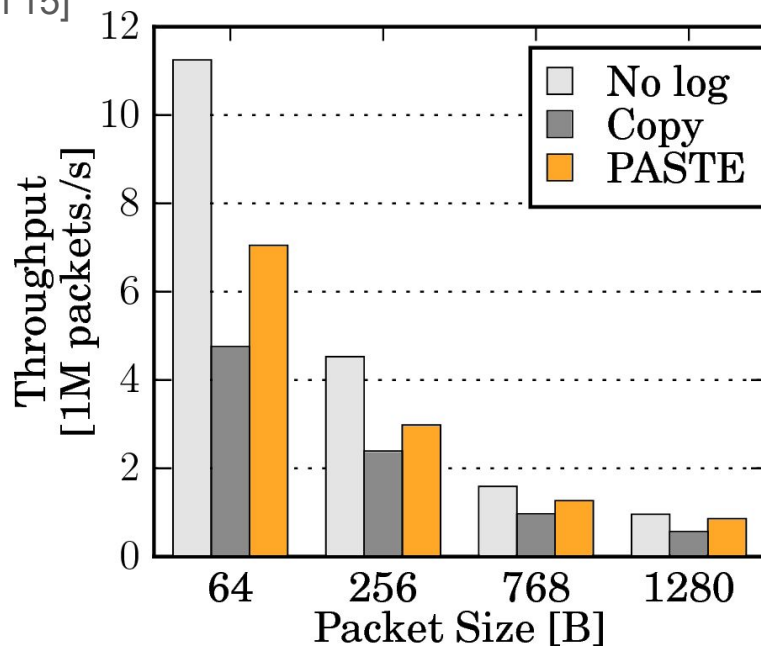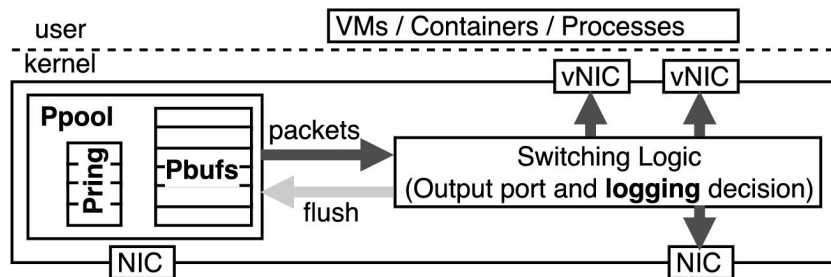
**1280B**

**WAL**

**B+tree (all writes)**

# Is PASTE applicable to existing applications?

- Redis

# Is PASTE useful for systems other than DB/file storage?

- Packet logging *prior to forwarding*
  - Fault-tolerant middlebox [Sigcomm'15]
  - Traffic recording
- Extend mSwitch [SOSR'15]
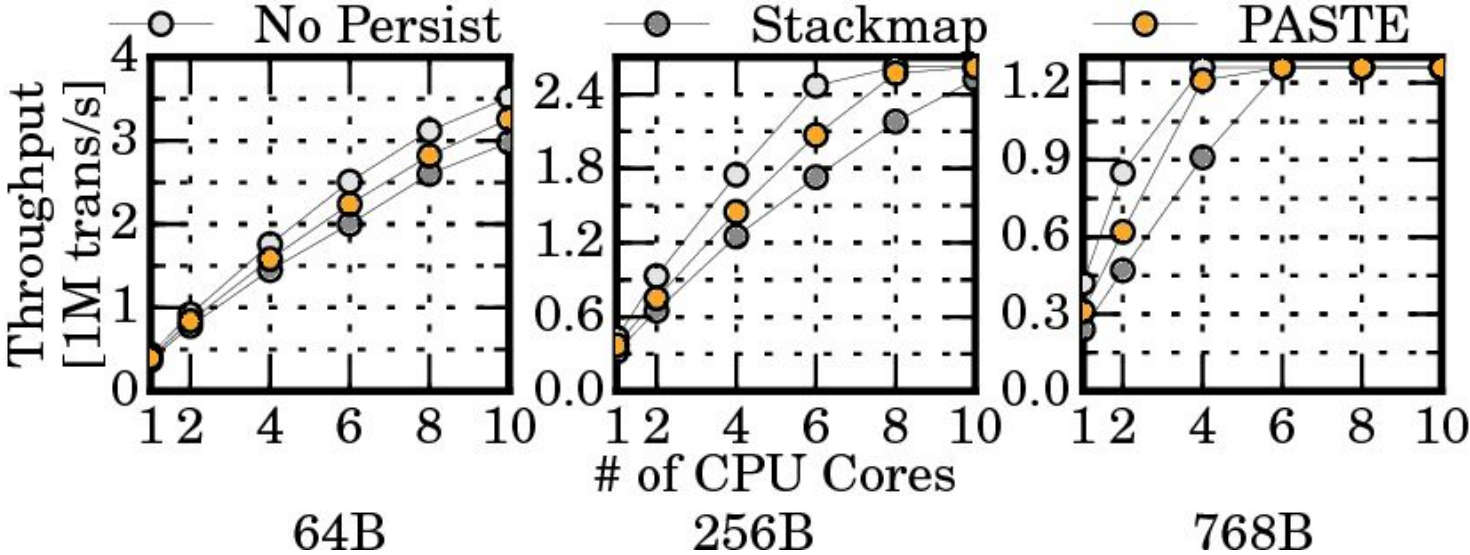  - Scalable NFV backend switch

# Conclusion

- PASTE is a network programming interface that:
  - Enables durable zero copy to NVMM
  - Helps apps organize persistent data structures on NVMM
  - Lets apps use TCP/IP and be protected
  - Offers high-performance network stack even w/o NVMM

https://github.com/luigirizzo/netmap/tree/paste
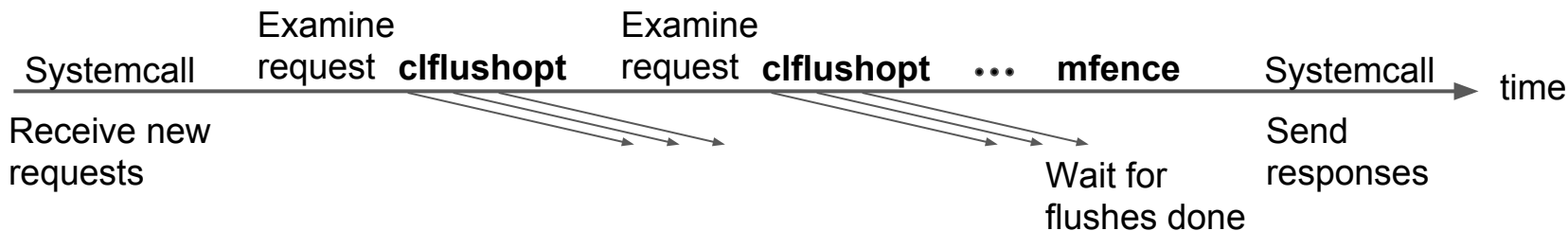micchie@sfc.wide.ad.jp or @michioh

# Multicore Scalability

● WAL throughput



64B       256B       768B

# Further Opportunity with Co-designed Stacks

- What if we use higher access latency NVMM?
  - e.g., 3D-Xpoint
- Overlap flushes and processing with clflushopt and mfence before system call (triggers packet I/O)
  - See the paper for results

Examine request **clflushopt**   Examine request   **clflushopt**   ···   **mfence**   Systemcall

Systemcall

Receive new requests

Wait for flushes done

Send responses

time

# Experiment Setup

- Intel Xeon E5-2640v4 (2.4 Ghz)
- HPE 8GB NVDIMM (NVDIMM-N)
- Intel X540 10 GbE NIC
- Comparison
  - Linux and Stackmap [ATC'15] (current state-of-the art)
  - Fair to use the same kernel TCP/IP implementation