# Extending Linux Network Stack for Persistent Memory

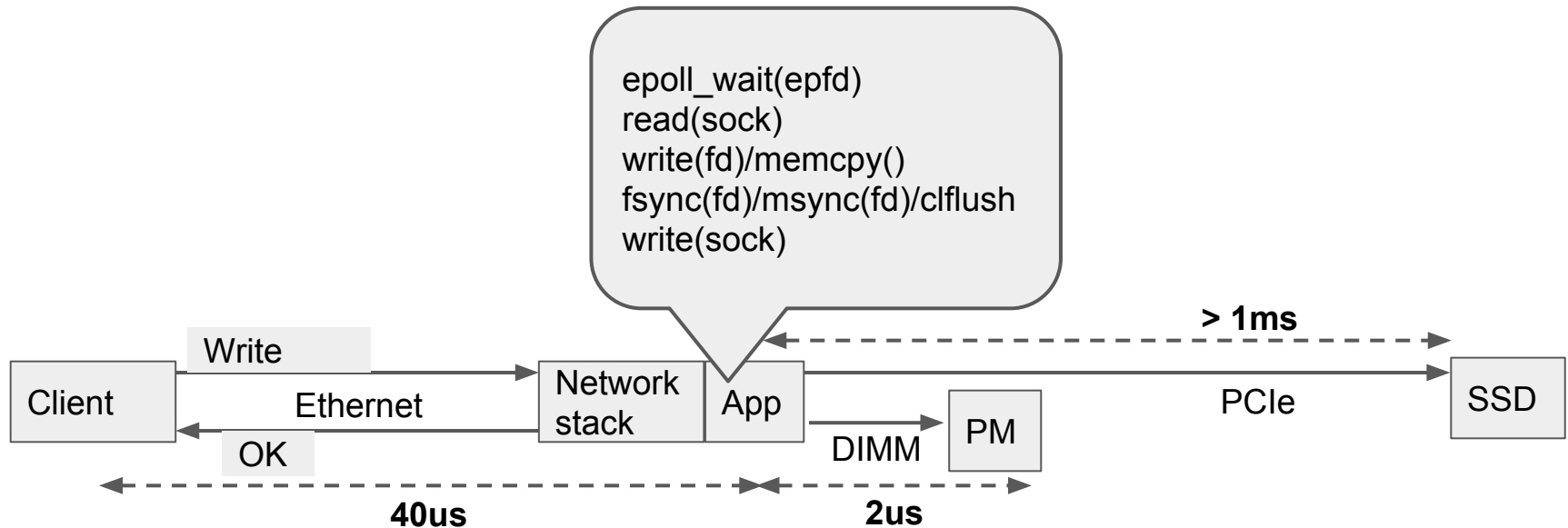**Michio Honda (NEC Laboratories Europe)**
Giuseppe Lettieri (Universita di Pisa)
Lars Eggert (NetApp)
Douglas Santry (NetApp)

# Problem

- Persistent memory (PM) is emerging
  - HPE NVDIMM, Intel 3D-Xpoint etc.
- **End-to-end latencies are now dominated by networking**

epoll_wait(epfd)
read(sock)
write(fd)/memcpy()
fsync(fd)/msync(fd)/clflush
write(sock)

Write

Ethernet

OK

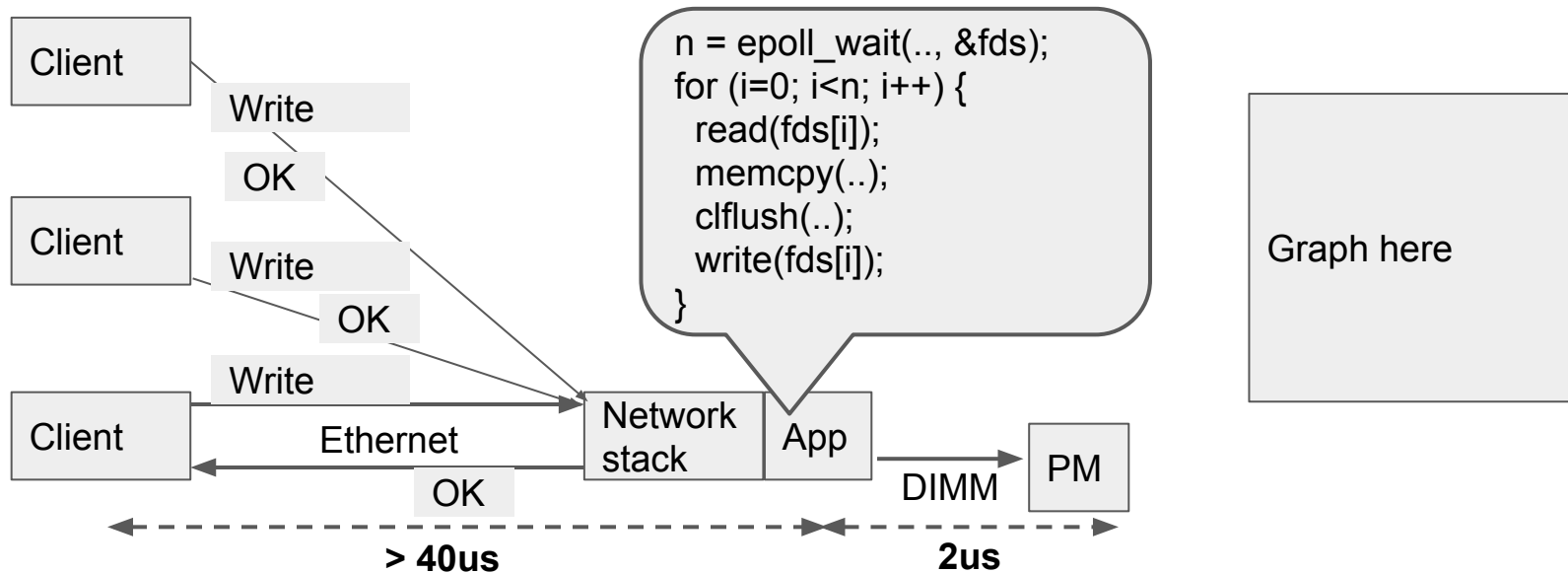Client — Network stack — App — DIMM — PM — PCIe — SSD

> 1ms

40us          2us

# Problem

- Persistent memory (PM) is emerging
  - HPE NVDIMM, Intel 3D-Xpoint etc.
- **End-to-end latencies are now dominated by networking**



```
n = epoll_wait(.., &fds);
for (i=0; i<n; i++) {
  read(fds[i]);
  memcpy(..);
  clflush(..);
  write(fds[i]);
}
```

# Status Quo

|  | Zero copy | Syscall batching | Synchronous I/O | Linux TCP/IP | Named packet buffers |
|---|---|---|---|---|---|
| MSG_ZEROCOPY | ✔ | x | x | ✔ | x |
| KCM | x | ✔ | x | ✔ | x |
| DPDK | ✔ | ✔ | ✔ | x | x |
| PASTE | ✔ | ✔ | ✔ | ✔ | ✔ |

# PASTE Design

- Zero-copy packet I/O to and from PM-backed file
- All the best practices for high-performance network stacks



App

socket(), bind(), listen(), accept(), connect()

nmd =
   nm_open("eth1", va);
poll(nmd->fd)

open(
   "**/mnt/pm/foo**");
va = mmap();

user
kernel

RX    TX

Network stack
kernel_sendpage()
sk->sk_data_ready()
netif_receive_skb()
dev->ndo_start_xmit()

PASTE

RX    TX

File system
w/ DAX

get_user_pages(va)

**/mnt/pm/foo**

Packet buffers
(fixed size (e.g., 2K) each)

NIC

# In Action



Application
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()

netmap API

Kernel

metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

**/mnt/nvmm/myapp_metadata**

packet buffers
(static)

**/mnt/nvmm/pktbufs**

TCP/IP
input
and
output

NIC ring

# In Action

Unread

Read or written

**Application**
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()

netmap API

Kernel

metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

TCP/IP input and output

metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

packet buffers (static)

NIC ring

**/mnt/nvmm/myapp_metadata**

**/mnt/nvmm/pktbufs**

# In Action

Unread

Read or written

**Application**
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()

netmap API

Kernel

metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

TCP/IP input and output

packet buffers (static)

NIC ring

**/mnt/nvmm/myapp_metadata**

**/mnt/nvmm/pktbufs**

# In Action

# In Action

Unread ┄┄┄┄

Read or written ▭

Flushed ▭

**Application**
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()

netmap API

Kernel

## metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

## metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

**/mnt/nvmm/myapp_metadata**

TCP/IP input and output

packet buffers (static)

NIC ring

**/mnt/nvmm/pktbufs**

# In Action



Unread

Read or written

Flushed

Application
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()
netmap API
Kernel

metadata header
```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

**/mnt/nvmm/myapp_metadata**

TCP/IP input and output

packet buffers (static)

NIC ring

**/mnt/nvmm/pktbufs**

# In Action



Unread

Read or written

Flushed

**Application**
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()

netmap API

Kernel

metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

**/mnt/nvmm/myapp_metadata**

TCP/IP input and output

packet buffers (static)

NIC ring

**/mnt/nvmm/pktbufs**

# In Action

Unread

Read or written

Flushed

**Application**
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()

netmap API

Kernel

metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

**/mnt/nvmm/myapp_metadata**

TCP/IP input and output

NIC ring

packet buffers (static)

**/mnt/nvmm/pktbufs**

# In Action

Unread

Read or written

Flushed

Application
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()

netmap API

Kernel

metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

**/mnt/nvmm/myapp_metadata**

TCP/IP input and output

packet buffers (static)

NIC ring

**/mnt/nvmm/pktbufs**

# In Action

Unread

Read or written

Flushed

Application
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

Idempotent request

User

mmap()

netmap API

Kernel

metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

**/mnt/nvmm/myapp_metadata**

TCP/IP input and output

packet buffers (static)

NIC ring

**/mnt/nvmm/pktbufs**

# In Action

Unread

Read or written

Flushed

**Application**
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()

netmap API

Kernel

### metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

### metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

**/mnt/nvmm/myapp_metadata**

TCP/IP input and output

packet buffers (static)

NIC ring

**/mnt/nvmm/pktbufs**

# In Action

Unread

Read or written

Flushed

Application
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()                    netmap API          Kernel

metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

metadata entries

| buf_idx | off | len  |
|---------|-----|------|
| 0       | 100 | 1135 |
| 1       | 100 | 932  |
| 3       | 100 | 1024 |

**/mnt/nvmm/myapp_metadata**

packet buffers (static)

**/mnt/nvmm/pktbufs**

TCP/IP input and output

NIC ring

# In Action

Unread

Read or written

Flushed

**Application**
**(1) Read data (zero copy)**
**(2) Write metadata entry**
**(3) Flush (buffer and metadata)**

User

mmap()

netmap API

Kernel

## metadata header

```
/mnt/nvmm/pktbufs
buf_ofs: 123
```

## metadata entries

| buf_idx | off | len |
|---------|-----|------|
| 0 | 100 | 1135 |
| 1 | 100 | 932 |
| 3 | 100 | 1024 |

`/mnt/nvmm/myapp_metadata`

TCP/IP input and output

packet buffers (static)

NIC ring

`/mnt/nvmm/pktbufs`

# In Action

poll(nmd->fd)                                      App

head
tail
9
10
11  RX
12

Network stack                    PASTE

sk->sk_data_ready()

1
2
netif_receive_skb()              3  RX
4

NIC

# In Action

App

```
struct netmap_ring *rxr = nmd->rx_ring;
poll(nmd->fd)
for (i = ring->head; i < ring->tail; i++) {
    char *p = nmb(ring, slot[i].buf_idx);
    if (is_write(p)) {
        clflush(p);
        swap_buf_idx(slot, nmd->extra);
    }
}
}
```

head

1
2
3 RX
4

tail

Network stack

sk->sk_data_ready()

netif_receive_skb()

PASTE

9
10
11 RX
12

NIC

# Challenges

- Moving packets to and from the TCP/IP stack
    - Two-level skb destructors
    - sk's callbacks
- Avoiding HoL
    - Buffer swapping to drain sender or NIC queue
- Getting kernel pages
    - get_user_pages()

# Performance

# Conclusion