# High Performance Networking
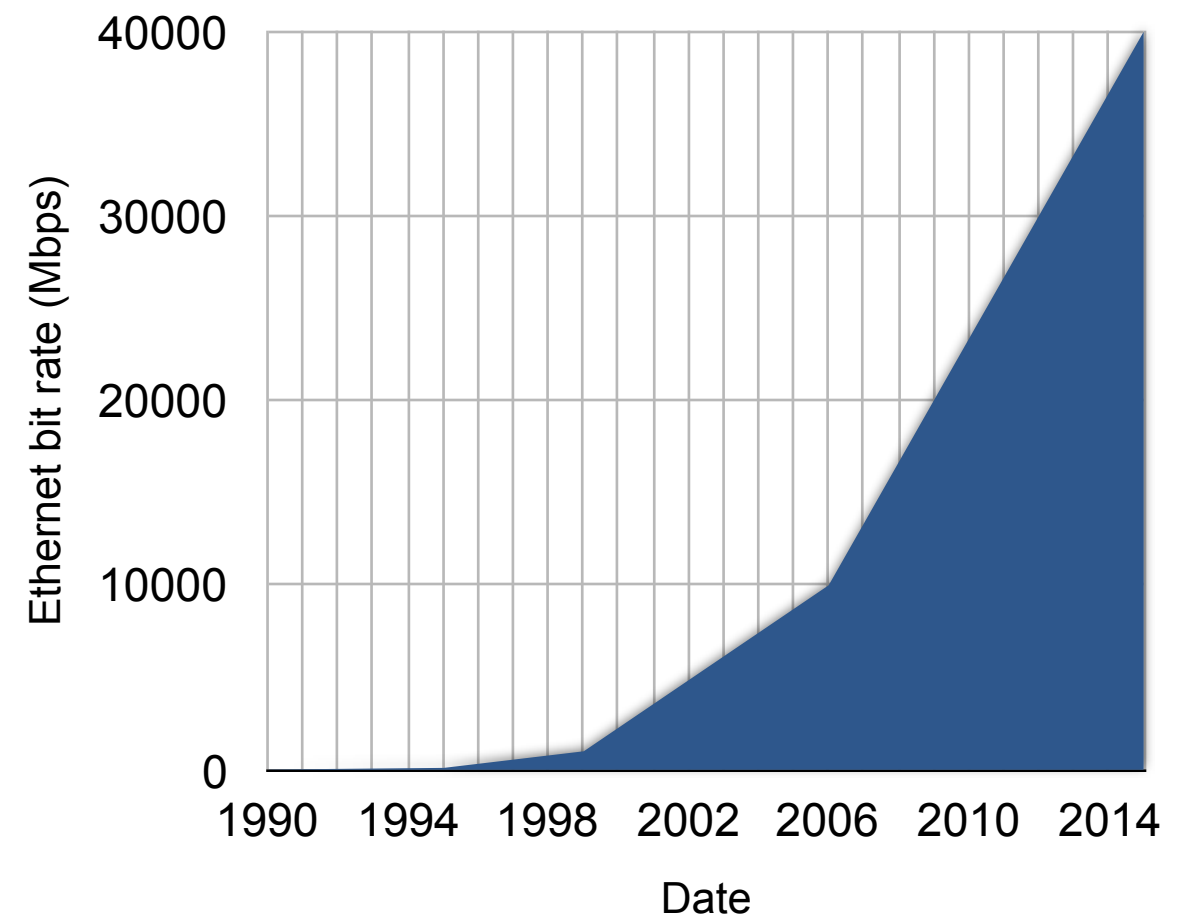
Advanced Operating Systems

Lecture 14

# Lecture Outline

- Limitations of the kernel protocol stack

- Alternative network stacks

- Accelerating TCP via API improvements

- Some cautionary remarks

# Network Performance Growth

- Network performance is increasing faster than CPU performance

  - Chart shows Ethernet bit rate over time – wireless links follow a similar curve

  - Closely tracking exponential growth over time – unlike CPU speed, which stopped growing significantly mid-2000s

- MTU remains constant → packet rate increases

  - Maximum 14,880,952 packets/second on 10Gbps Ethernet (scales linearly with link rate)

  - Minimum size packet is 60 bytes data, with 8 byte preamble, 4 byte CRC; 12 byte inter-frame gap (silent period) between packets

- CPU cycles available to process each packet decreasing

# Limitations of the Kernel Protocol Stack

- Why does the traditional kernel protocol stack offer sub-optimal performance?

- Designed when CPUs were faster than networks

  - Allocates memory for buffers on per-packet basis

  - Copies data multiple times, from NIC ("network interface card") to kernel to application

  - System call to send/receive each packet

    - Layered architecture offers clean design, but not efficient packet processing

    - Example on right: timing of a sendto() system call on FreeBSD: 950ns total; overheads at each layer boundary due to system call, copies, etc.

- How to redesign the protocol stack to reduce overheads?

| File | Function/description | time ns | delta ns |
|---|---|---|---|
| user program | `sendto` system call | 8 | 96 |
| uipc_syscalls.c | `sys_sendto` | 104 | |
| uipc_syscalls.c | `sendit` | 111 | |
| uipc_syscalls.c | `kern_sendit` | 118 | |
| uipc_socket.c | `sosend` | — | |
| uipc_socket.c | `sosend_dgram` sockbuf locking, mbuf allocation, copyin | 146 | 137 |
| udp_usrreq.c | `udp_send` | 273 | |
| udp_usrreq.c | `udp_output` | 273 | 57 |
| ip_output.c | `ip_output` route lookup, ip header setup | 330 | 198 |
| if_ethersubr.c | `ether_output` MAC header lookup and copy, loopback | 528 | 162 |
| if_ethersubr.c | `ether_output_frame` | 690 | |
| ixgbe.c | `ixgbe_mq_start` | 698 | |
| ixgbe.c | `ixgbe_mq_start_locked` | 720 | |
| ixgbe.c | `ixgbe_xmit` mbuf mangling, device programming | 730 | 220 |
| – | on wire | 950 | |

Source: L. Rizzo. netmap: a novel framework for fast packet I/O. In Proceedings of the USENIX Annual Technical Conference, Boston, MA, USA, June 2012.
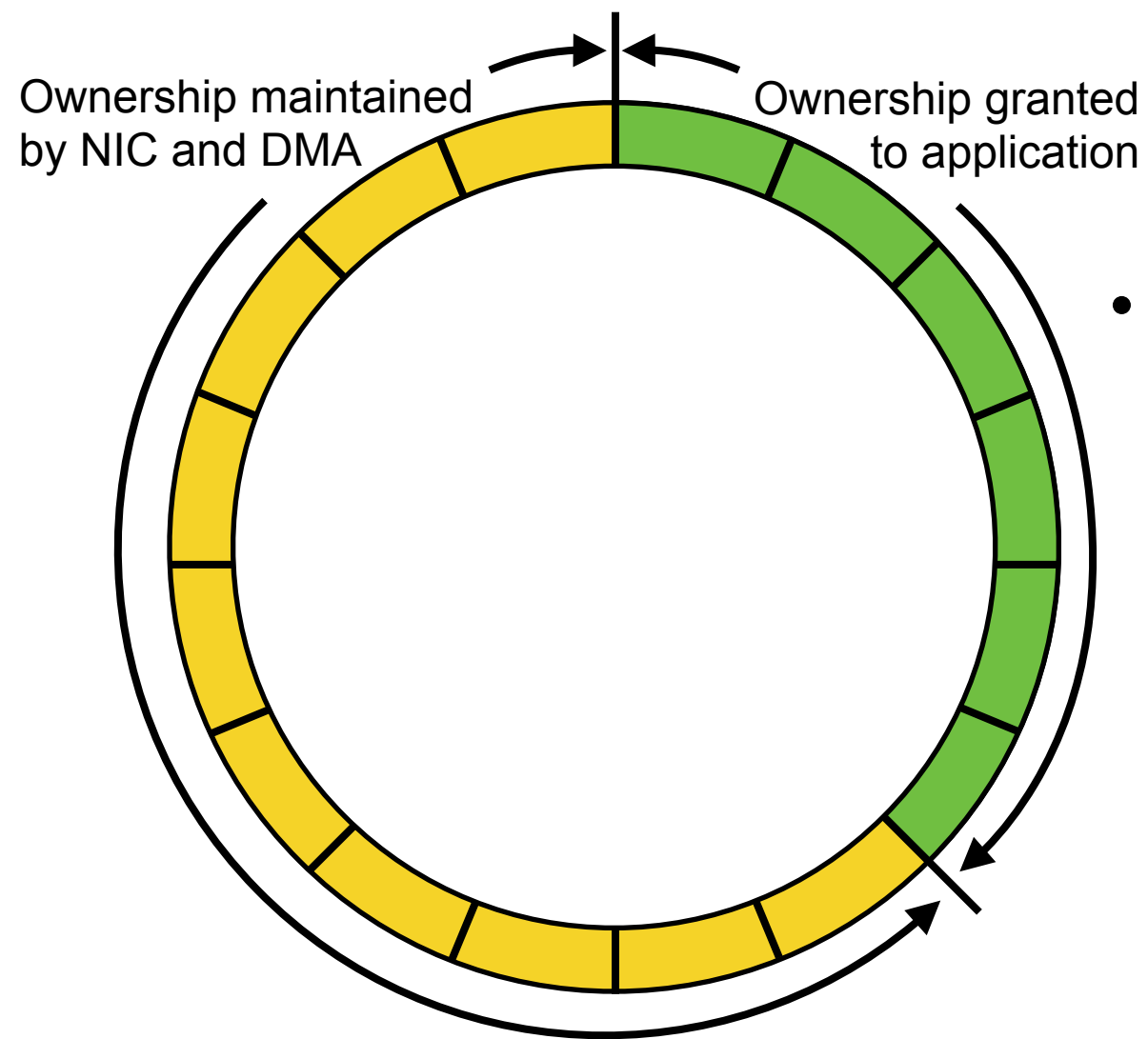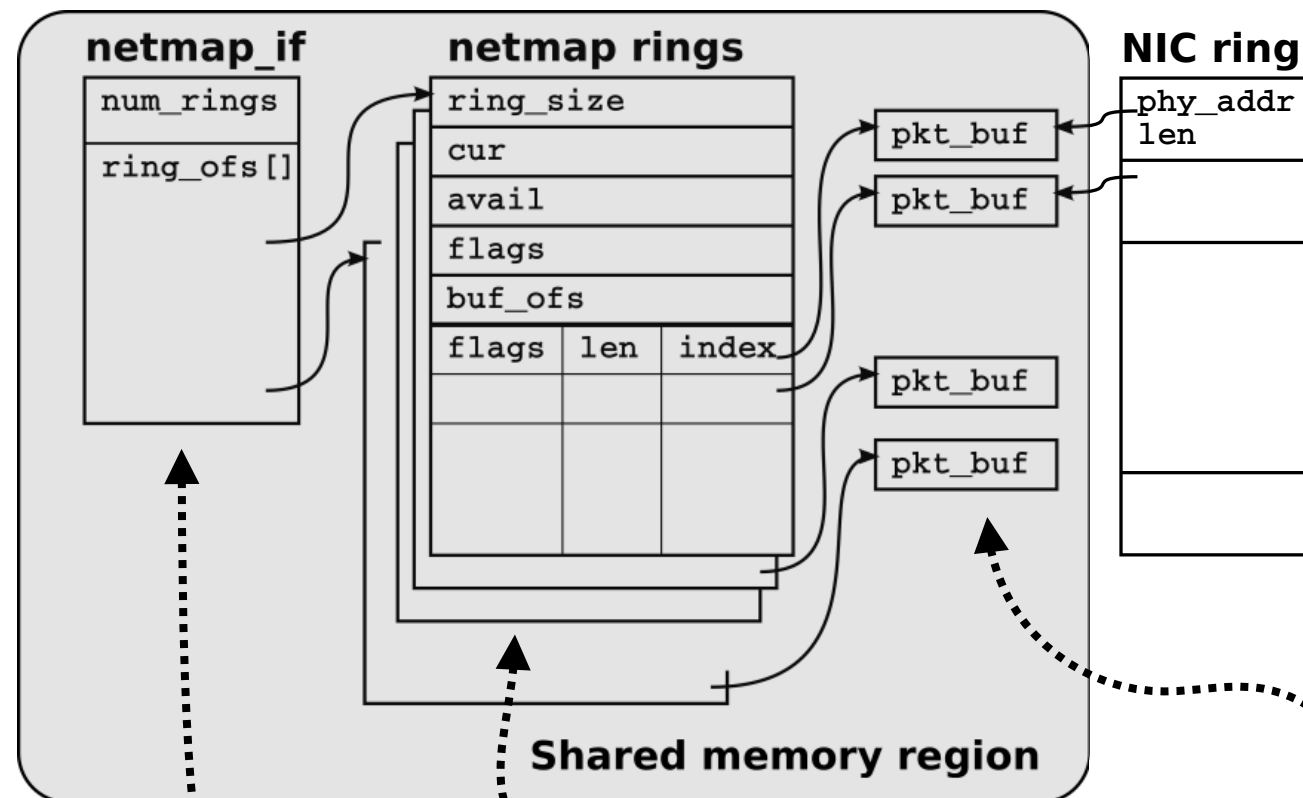
# An Alternative Network Stack: netmap

- Changes network API – a mechanism to give an application dedicated control of a NIC

  - A replacement for `libpcap`, not the Sockets API

  - Basis for fast packet capture applications; software router; user-space protocol stack – not general purpose

- Pre-allocate buffers, that are shared between OS and user application; coordinate buffer ownership

  - No memory allocation at time the packets are sent/received

  - No data copies – DMA direct to application accessible memory

  - Fewer system calls – one system call can transfer ownership of multiple buffers between application and kernel

- NIC uses efficient DMA to transfer packets to and from memory; kernel manages synchronisation and memory protection

  - Kernel is the control plane

  - NIC and DMA transfer manage the data plane

L. Rizzo. netmap: a novel framework for fast packet I/O. In Proceedings of the USENIX Annual Technical Conference, Boston, MA, USA, June 2012. https://www.usenix.org/conference/atc12/technical-sessions/presentation/rizzo (paper and video of presentation)

# Shared Ownership of NIC Buffers (1)

Ownership maintained
by NIC and DMA

Ownership granted
to application

- Modern NICs maintain circular buffers sized to hold queues of full size packets

  - NIC writes incoming packets direct to one segment of ring buffer via DMA

  - Operating system copies from other segment into lower layer of protocol stack – first of several copies

  - If using netmap → OS disconnected; ownership of ring buffer segment is temporarily granted to application to process packets in place

  - (Analogous for outgoing packets)

# Shared Ownership of NIC Buffers (2)



Source: L. Rizzo. netmap: a novel framework for fast packet I/O. In Proceedings of the USENIX Annual Technical Conference, Boston, MA, USA, June 2012.

`pkt_buf` objects form the circular buffer, and are in memory shared between NIC and application

`netmap_ring` provides index into the circular buffer; tracks current ownership of each `pkt_buf`

`netmap_if` provides metadata about the interface

# Sample netmap Code

```
struct netmap_if *nifp;
struct nmreq req;
int i, len;
char *buf;

fd = open("/dev/netmap", 0);
strcpy(req.nr_name, "ix0"); // register the interface
ioctl(fd, NIOCREG, &req); // offset of the structure
mem = mmap(NULL, req.nr_memsize, PROT_READ|PROT_WRITE, 0, fd, 0);
nifp = NETMAP_IF(mem, req.nr_offset);
for (;;) {
    struct pollfd x[1];
    struct netmap_ring *ring = NETMAP_RX_RING(nifp, 0);

    x[0].fd = fd;
    x[0].events = POLLIN;
    poll(x, 1, 1000);
    for ( ; ring->avail > 0 ; ring->avail--) {
        i = ring->cur;
        buf = NETMAP_BUF(ring, i);
        use_data(buf, ring->slot[i].len);
        ring->cur = NETMAP_NEXT(ring, i);
    }
}
```

Updates `netmap_ring` structure, based on the received data (`ring->cur` and `ring->avail` only, no data copied, no synchronisation)

Gets pointer to shared `pkf_buf`

# Benefits and Limitations of netmap

- Memory shared between NIC and application

    - Misbehaving applications can read memory owned by NIC – see unpredictable contents since DMA active in this region

    - Kernel data structures are protected – cannot crash the kernel or see other kernel data

- Operates on granularity of network interface

    - Application has access to all traffic on netmap interface

    - Limited applicability – not a replacement for the TCP/IP stack, but well suited to network monitoring or software router implementation

- Performance excellent – saturates 10Gbps Ethernet sending minimum size packets on 900MHz CPU

    - Using minimum size packets requires highest packet rate → highest overhead

    - Similar performance receiving packets



Source: L. Rizzo. netmap: a novel framework for fast packet I/O. In Proceedings of the USENIX Annual Technical Conference, Boston, MA, USA, June 2012.

# StackMap: Accelerating TCP with netmap

- Key insight: the TCP/IP stack processing is not expensive – inefficiencies are primarily system call overheads, copying data, and Socket API limitations with large number of file descriptors

- The netmap framework avoids the copies and reduces number of system calls, and gains in performance – but without the TCP/IP stack

- StackMap integrates the Linux kernel TCP/IP stack with netmap

  - Uses kernel TCP/IP stack for the control place

  - Uses netmap for the data plane – new API

  - Like netmap, requires a dedicated network interface for each StackMap enabled application



K. Yasukata, M. Honda, D. Santry, and L. Eggert. StackMap: Low-latency networking with the OS stack and dedicated NICs. In Proceedings of the USENIX Annual Technical Conference, Denver, CO, USA, June 2016. https://www.usenix.org/conference/atc16/technical-sessions/presentation/yasukata

# StackMap Architecture



Source: K. Yasukata, M. Honda, D. Santry, and L. Eggert. StackMap: Low-latency networking with the OS stack and dedicated NICs. In Proceedings of the USENIX Annual Technical Conference, Denver, CO, USA, June 2016.

- Socket API for control: `socket()`, `bind()`, `listen()`, `accept()`, etc.

- Netmap API used for data

  - `STACKMAP_BUF()` updates netmap's circular buffer, and passes data through the TCP/IP stack for processing

- StackMap manages the buffer pool

  - Circular buffer used by netmap

  - If packet must be stored for retransmission, its buffer is swapped out of the netmap ring and replaced by another from the pool, until needed

    - Zero copy – just swaps pointers; buffer is already shared between netmap, kernel, and application

  - Also manages *scratchpad* data structure, to simplify iteration over data from on multiple connections

# StackMap API

**Stack port RX**

```
 1    struct sockaddr_in sin = { AF_INET, "10.0.0.1", INADDR_ANY };
 2    int sd = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP);
 3    bind(sd, &sin);
 4    // prefix "stack" opens stack port for given interface
 5    struct nm_desc *nmd = nm_open("stack:ix0");
 6    connect(sd, dst_addr); /* non-blocking */
 7    // transmit using ring 0 only, for this example
 8    struct netmap_ring *ring = NETMAP_TXRING(nmd->nifp, 0);
 9    uint32_t cur = ring->cur;
10    while (app_has_data && cur != ring->tail) {
11        struct netmap_slot *slot = &ring->slot[cur];
12        char *buf = STACKMAP_BUF(ring, slot->buf_index);
13        // place payload in buf, then
14        slot->fd = sd;
15        cur = nm_ring_next(ring, cur);
16    }
17    ring->head = ring->cur = cur;
18    ioctl(nmd->fd, NIOCTXSYNC);
```

Sockets API used to initiate connection

The netmap API is used to send and receive data

The `STACKMAP_BUF()` call is an extension that passes the data to the kernel TCP/IP stack, handles ACKs, retransmission, congestion control, etc.

Source: K. Yasukata, M. Honda, D. Santry, and L. Eggert. StackMap: Low-latency networking with the OS stack and dedicated NICs. In Proceedings of the USENIX Annual Technical Conference, Denver, CO, USA, June 2016.

# StackMap Performance



64 B Response Size      512 B Response Size      1280 B Response Size

- Linux with StackMap extensions outperforms standard Linux

- Primary benefits are avoiding copies and better scaling with concurrent flows: performance benefit increases with response size and number of connections

- Note: uses the full kernel TCP/IP stack – not a cut-down version – benefits due to new API, not reduced functionality

# Accelerating TCP: Other Approaches

- Numerous other attempts to accelerate TCP/IP stack processing exist:
  - Example: sandstorm builds on netmap; combines application, TCP/IP, and ethernet processing into a user-space library – builds highly optimised, special purpose, protocol stack for each application (see "Network Stack Specialization for Performance")
  - Example: Google's QUIC protocol aims to provide an alternative to TCP, with better performance, in a user space protocol running over UDP
  - …

- No clear consensus on the right approach
  - API changes to batch packet processing and reduce copies are clearly beneficial – how general purpose the resulting API needs to be is unclear
  - Pervasive encryption may reduce the benefits – you need to make a copy while decrypting, and many of these approaches benefit from zero-copy stacks

I. Marinos, R. N. M. Watson, and M. Handley. Network stack specialization for performance. Proceedings of the SIGCOMM Conference, Chicago, IL, USA, August 2014. ACM. DOI: 10.1145/2619239.2626311 (Open access via http://dl.acm.org/authorize?N71202)

# Conclusions

- Improvements to network performance relative to CPU performance highlight limitations of Sockets API for high performance – ongoing work to find good new APIs

- Beware: it's very hard to beat TCP – incredibly well tuned, secure, and well maintained

  - TCP highly optimised over many years – handles edge cases in security and congestion control that are non-obvious

    - IETF RFC 7414 ("A Roadmap for TCP Specification Documents") is 55 pages, and references 150 other documents – textbook explanations of TCP omit much important detail

  - Likely better to optimise kernel TCP stack and API, than end up stuck on a poorly maintained user-space stack

# Further Reading





L. Rizzo. netmap: a novel framework for fast packet I/O. In Proceedings of the USENIX Annual Technical Conference, Boston, MA, USA, June 2012.
https://www.usenix.org/conference/atc12/technical-sessions/presentation/rizzo (paper and video of presentation)

K. Yasukata, M. Honda, D. Santry, and L. Eggert. StackMap: Low-latency networking with the OS stack and dedicated NICs. In Proceedings of the USENIX Annual Technical Conference, Denver, CO, USA, June 2016.
https://www.usenix.org/conference/atc16/technical-sessions/presentation/yasukata