

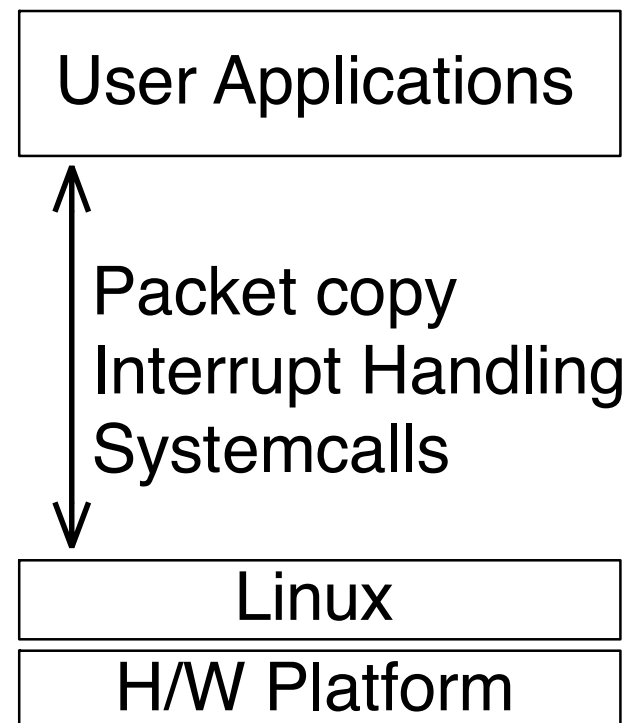
Software-Based Networking

Getting started with DPDK

Linux Packet Processing

Traditional networking:

- NIC uses DMA to copy data into kernel buffer
- Interrupt when packets arrive
- Copy packet data from kernel space to user space
- Use system call to send data from user space

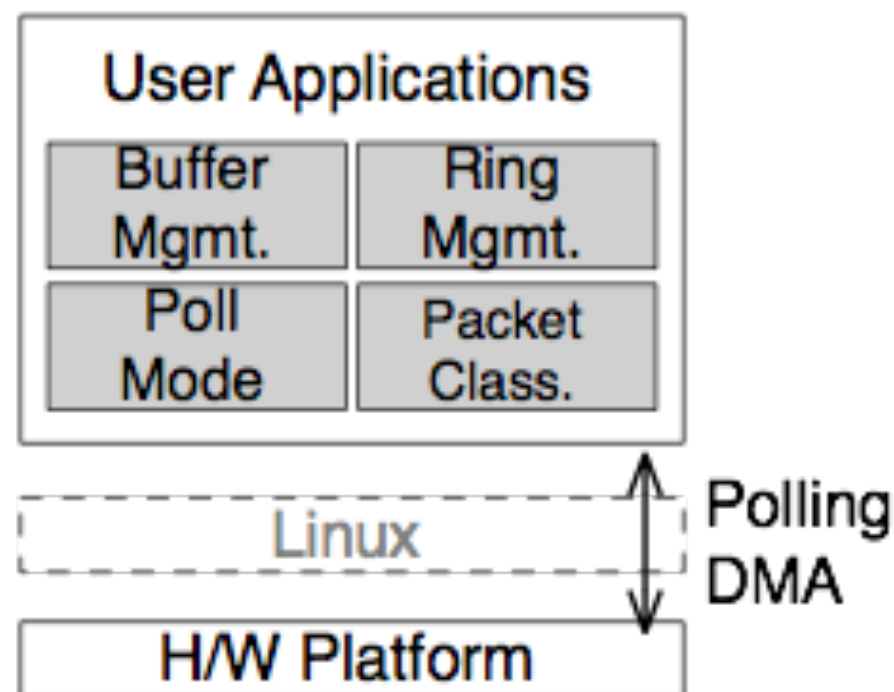


Can you handle being interrupted 60 million times per second?

User Space Packet Processing

Recent NICs and OS support allow user space apps to directly access packet data

- NIC uses DMA to copy data into ~~kernel~~ **user space** buffer
- ~~Interrupt~~ **use polling to find** when packets arrive
- ~~Copy packet data from kernel space to user space~~
- Use ~~system~~ **regular function** call to send data from user space



Data Plane Development Kit

High performance I/O library

Poll mode driver reads packets from NIC

Packets bypass the OS and are copied directly into user space memory

Low level library... does not provide:

- Support for multiple network functions
- SDN-based control
- Interrupt-driven NFs
- State management
- TCP stack

Data Plane Development Kit

Where to find it:

- <http://dpdk.org/>

What to use it for:

- Applications that need high speed access to low-level packet data

Why try it:

- One of the best documented open source projects I've ever seen

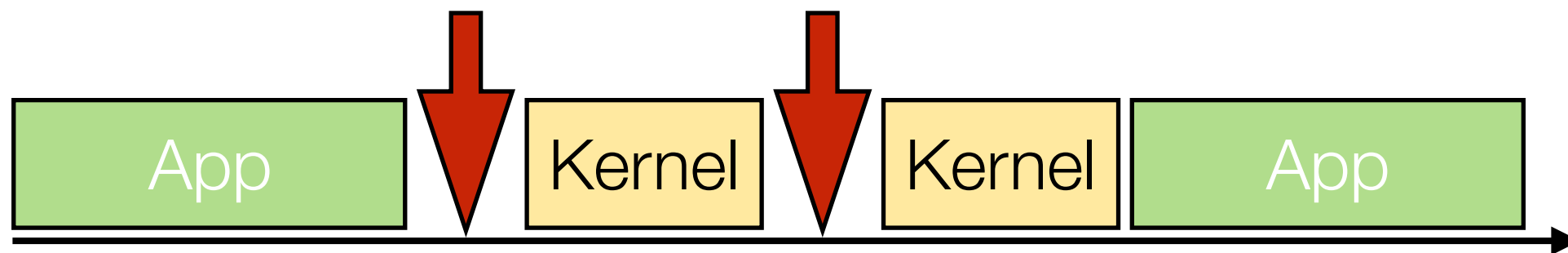
Alternatives:

- netmap
- PF_RING

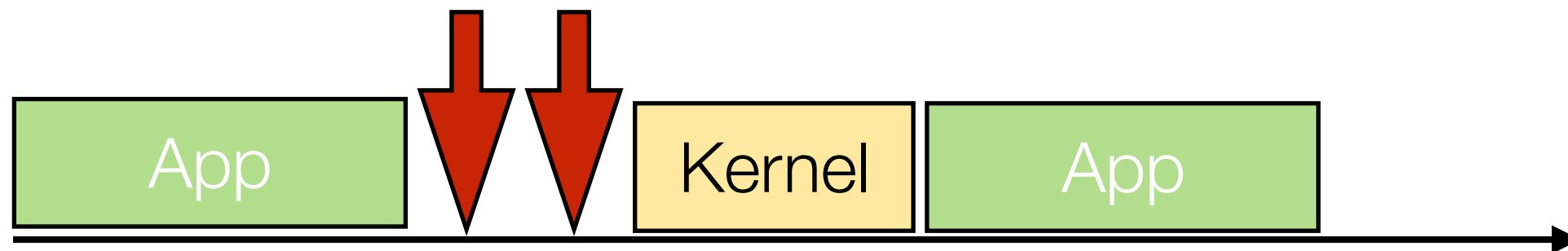
Network Interrupts

~~Interrupt
Context
Switch
Overhead~~

Very distracting! Have to stop doing useful work to handle incoming packets



Coalescing interrupts helps, but still causes problems



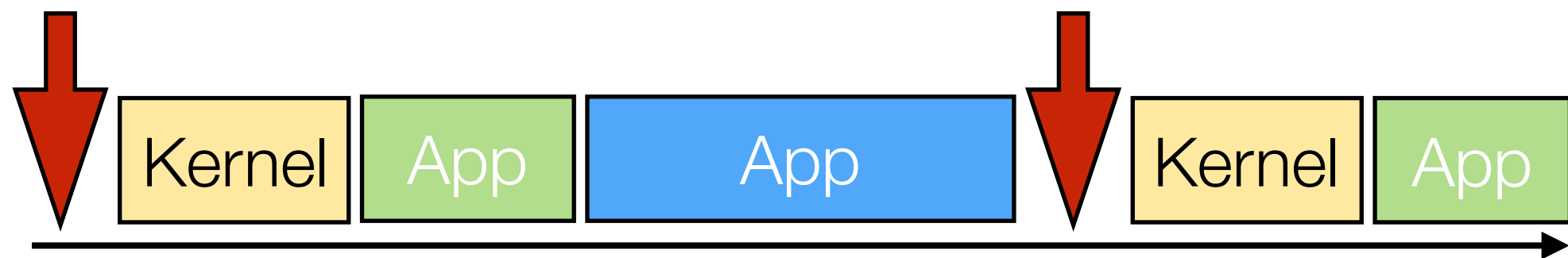
- Interrupts can arrive during critical sections!
- Interrupts can be delivered to the wrong CPU core!
- Still must pay context switch cost

Polling

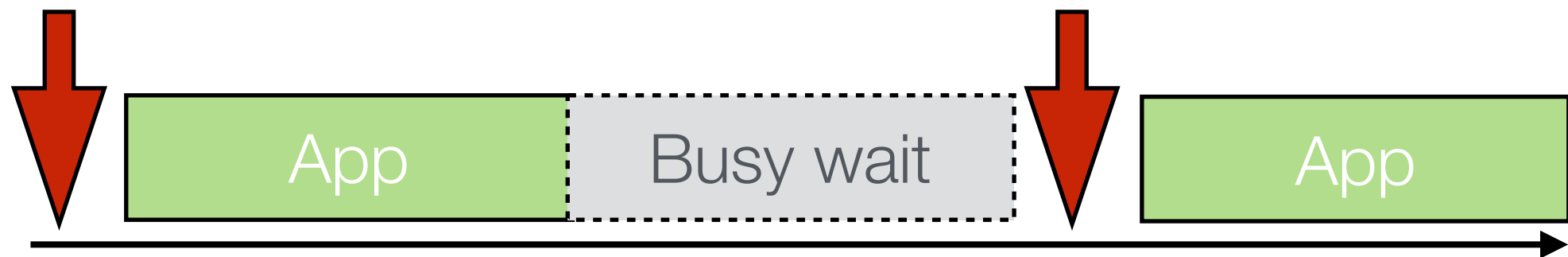
~~Interrupt
Context
Switch
Overhead~~

Continuously loop looking for new packet arrivals

Trade-off?



Interrupts help share the CPU



Polling can be wasteful

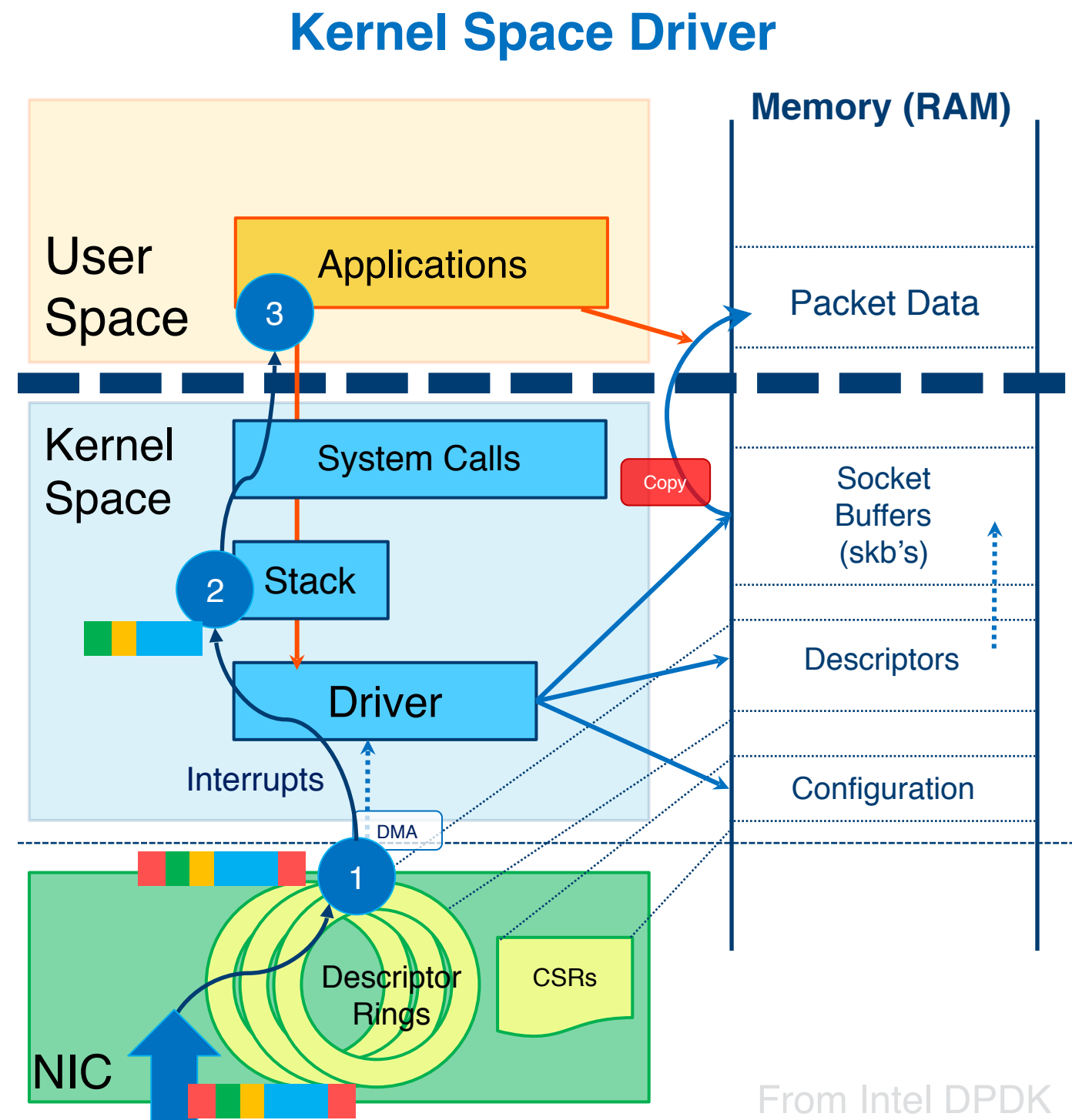
Kernel-User Overhead

Kernel
User
Overhead

NIC Driver operates in kernel mode

- Reads packets into kernel memory
- Stack pulls data out of packets
- Data is copied into user space for application
- Application uses system calls to interface with OS

Why is copying so bad?



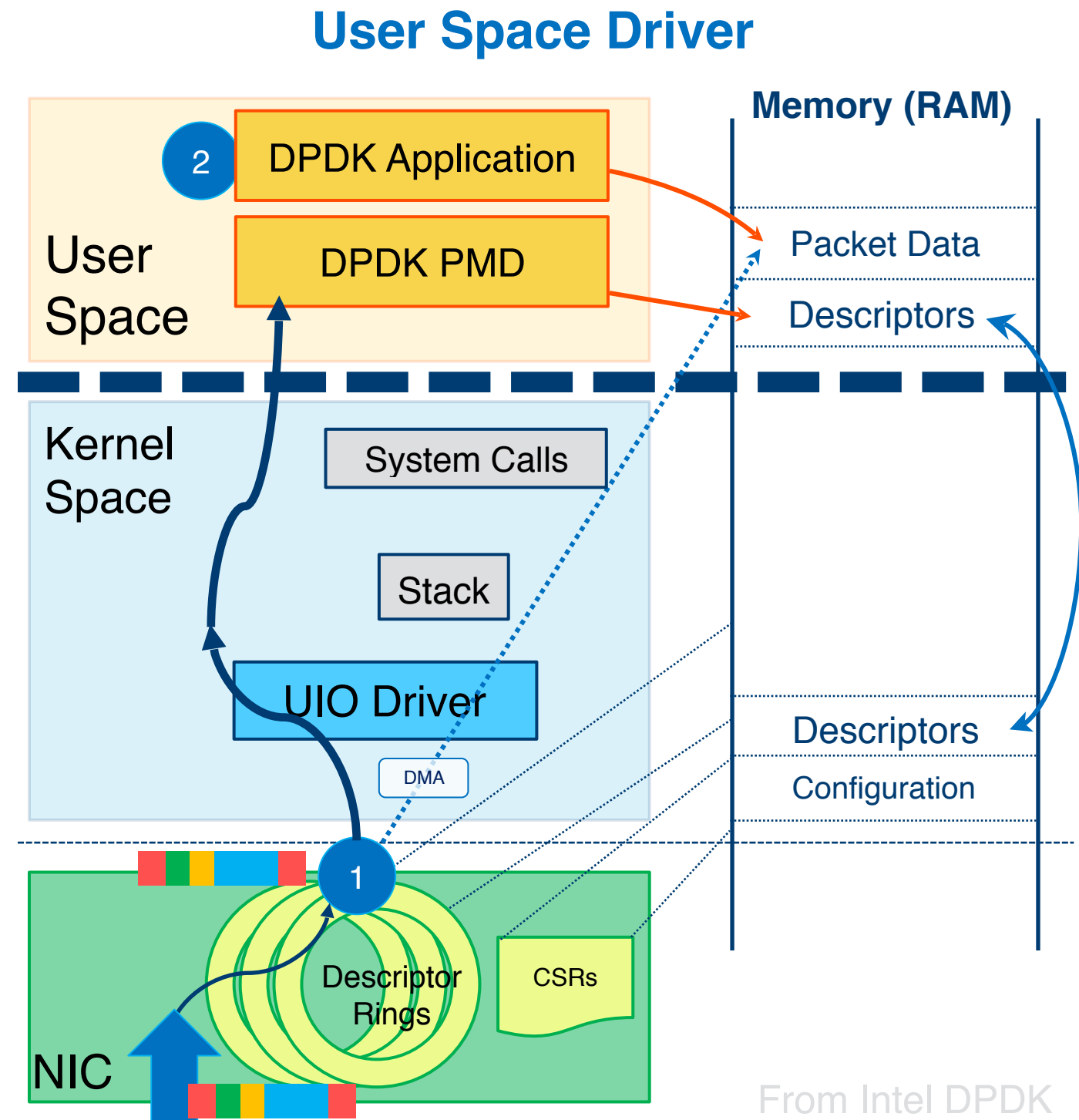
From Intel DPDK
University Lecture

Kernel Bypass

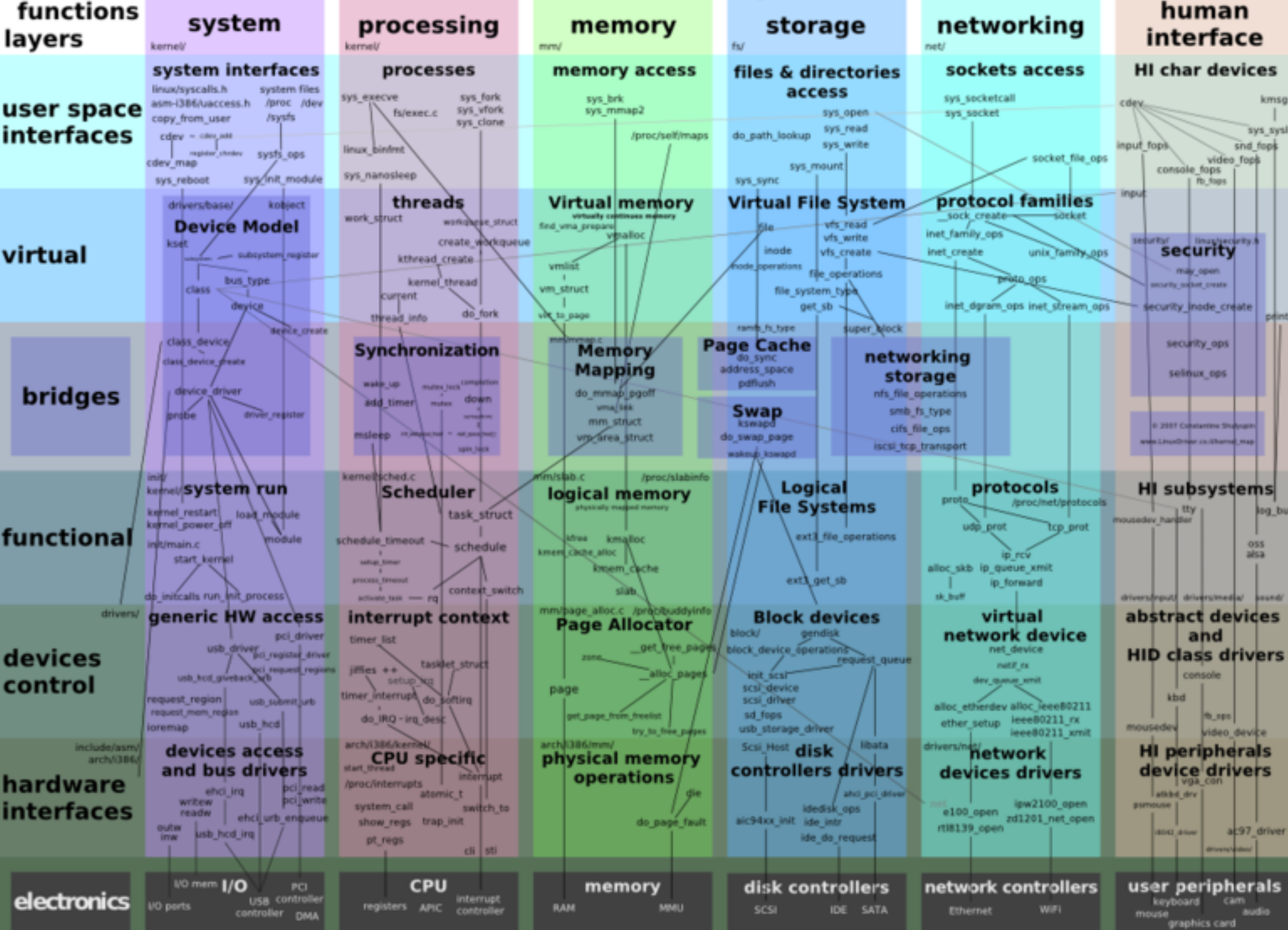
~~Kernel
User
Overhead~~

User-mode Driver

- Kernel only sets up basic access to NIC
- User-space driver tells NIC to DMA data directly into user-space memory
- No extra copies
- No in-kernel processing
- No context switching



From Intel DPDK
University Lecture



Networking

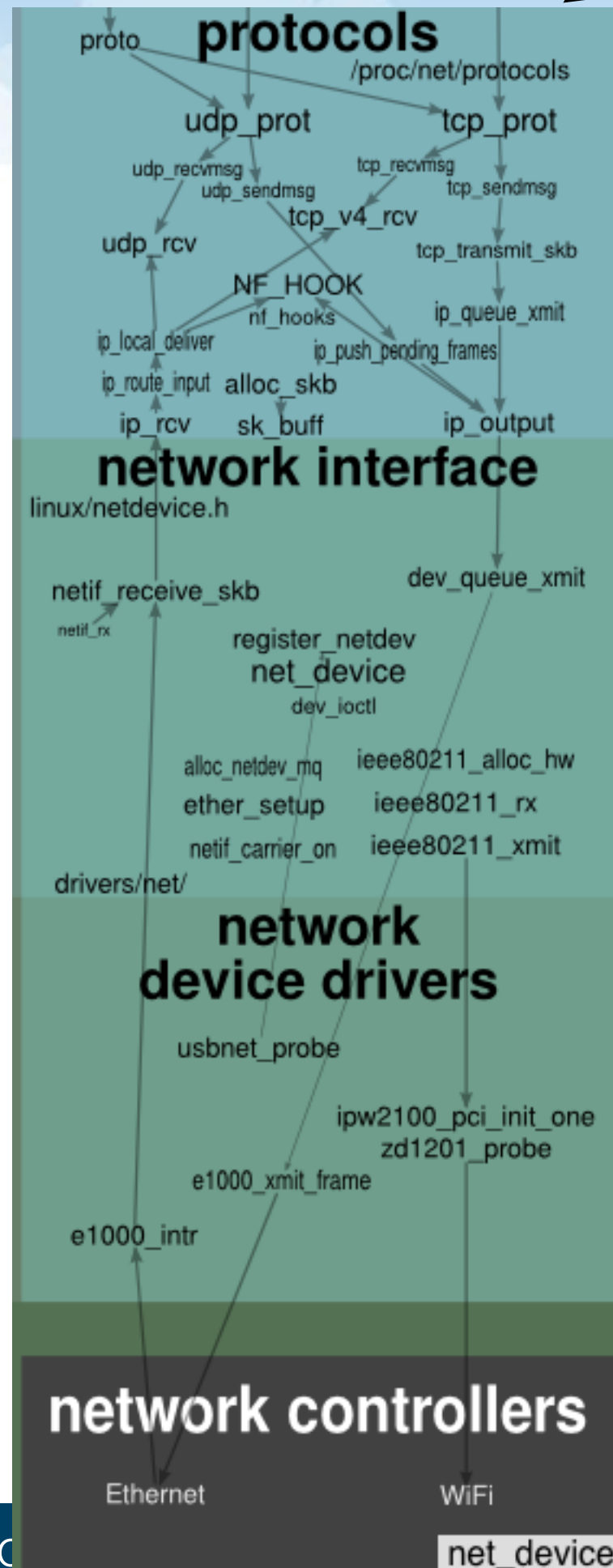
Linux networking stack has a lot of extra components

For NFV middlebox we don't use all of this:

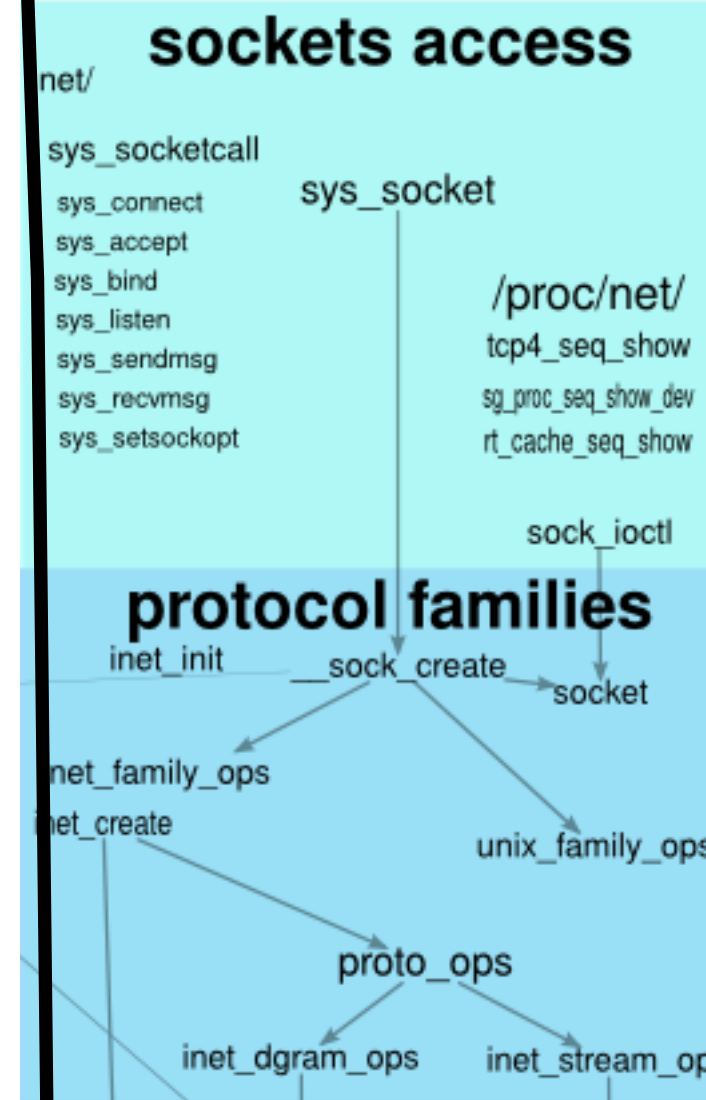
- TCP, UDP, sockets

NFV middle boxes just need packet data

- Need it fast!



Application



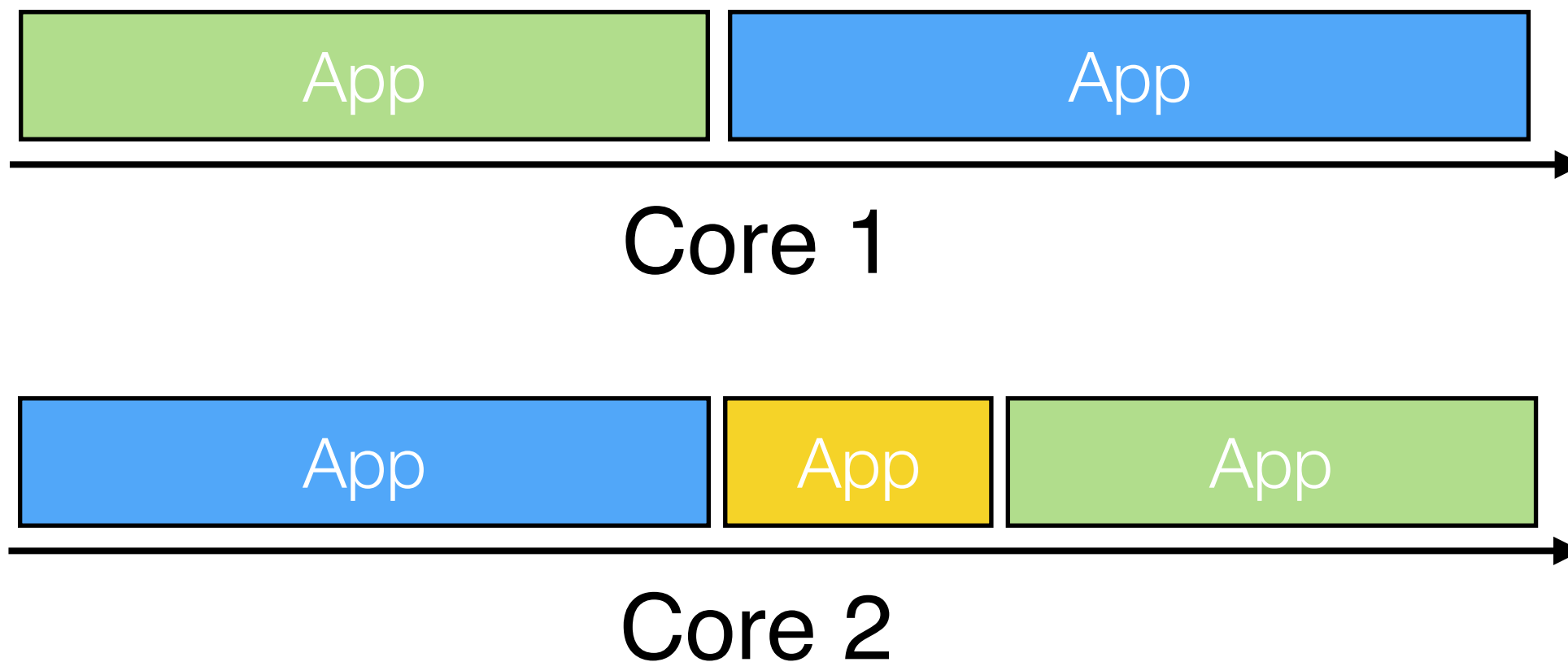
~~Kernel~~
~~User~~
~~Overhead~~

CPU Core Affinity

Core To
Thread
Scheduling
Overhead

Linux Scheduler can move threads between cores

- Context switches :(
- Cache locality :(

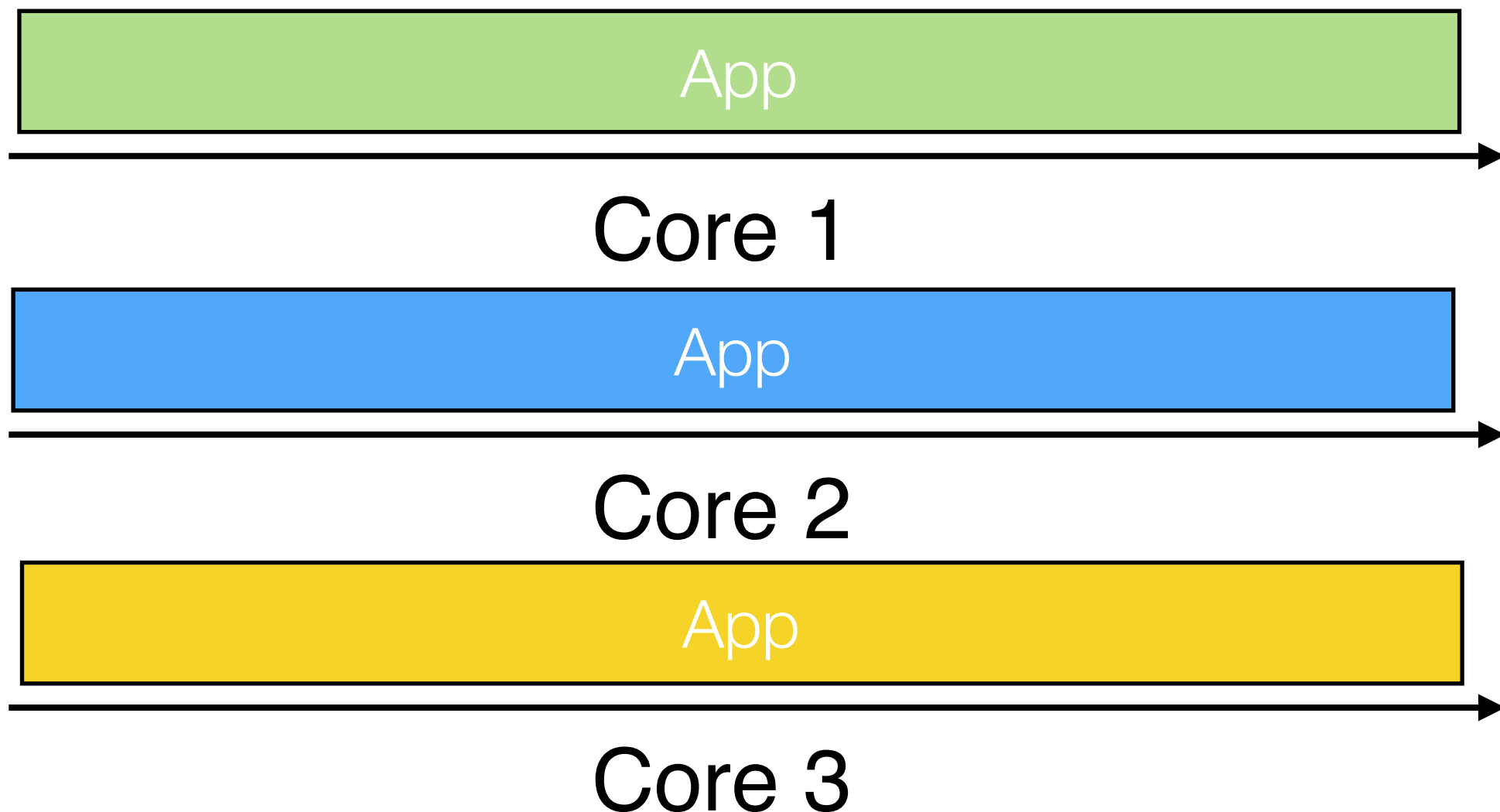


CPU Core Affinity

Core To
Thread
Scheduling
Overhead

Pin threads and dedicate cores

- **Trade-offs?**



Paging Overhead

4K
Paging
Overhead

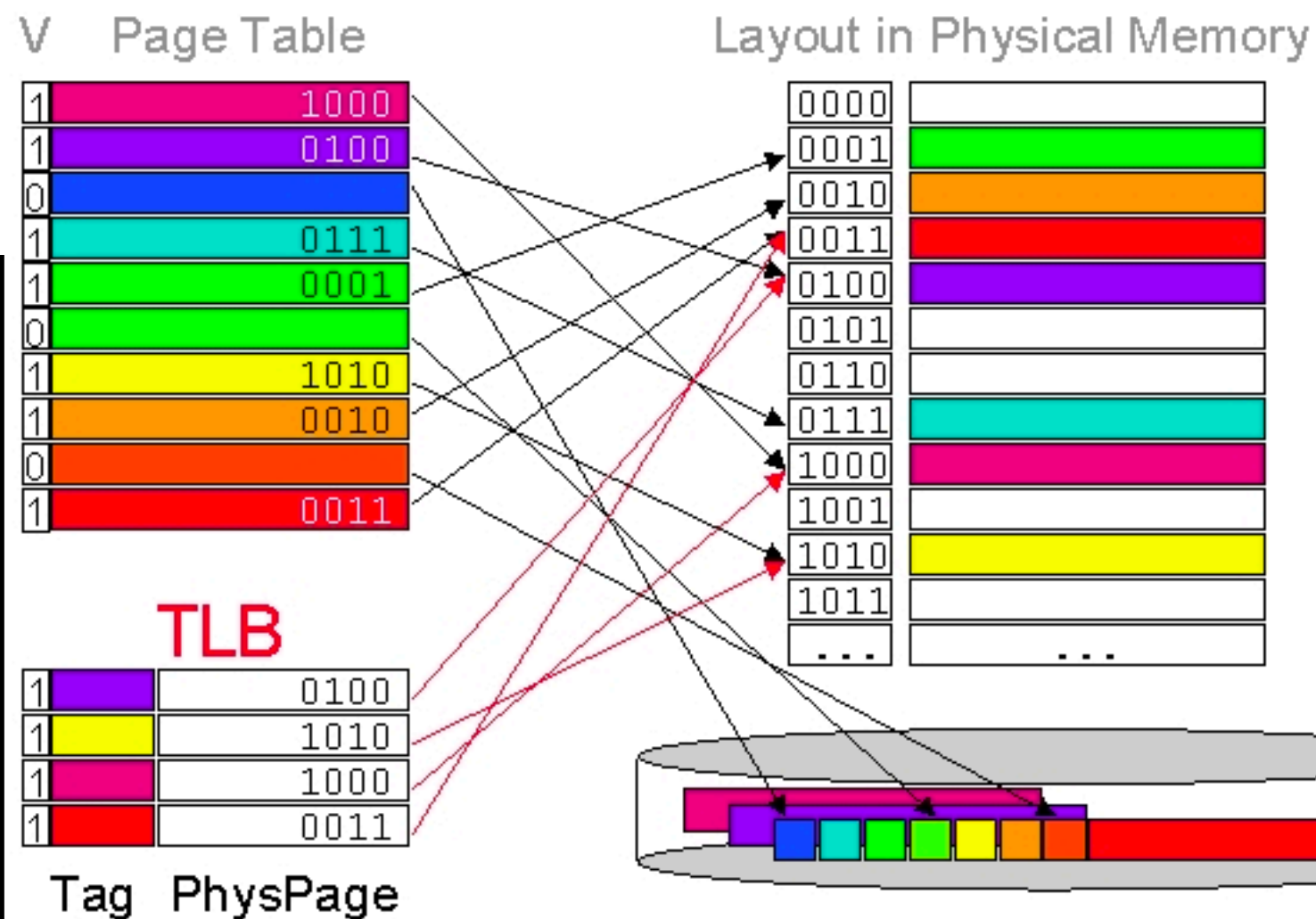
4KB Pages

- 4 packets per page
- 14 million pps
- 3.6 million page table entries every second

Packet \sim 1KB

Translation Lookaside Buffer

How big is the TLB?



Locks



Thread synchronization is expensive

- Tens of nanoseconds to take an uncontested lock
- 10Gbps -> 68ns per packet

Producer/Consumer architecture

- Gather packets from NIC (producer) and ask worker to process them (consumer)

Lock-free communication

- Ring-buffer based message queues

Bulk Operations



PCIe bus uses messaging protocols for CPU to interact with devices (NICs)

Each message incurs some overhead

Better to make larger bulk requests over PCIe

DPDK helps batch requests into bulk operations

- Retrieve a batch (32) of packet descriptors received by NIC
- Enqueue/dequeue beaches of packet descriptors onto rings

Trade-offs?

Using DPDK

Lots of examples! Great docs!

- But still fairly complex...
- <http://dpdk.readthedocs.io/>

Simple hands-on exercise... who will help?

Help me out!

1

I need volunteers to help run DPDK

- Find info at: **<https://github.com/sdnfv/onvm-tutorial>**
- Use nodes 2-9

Connect to your server in a terminal windows

- user=**tutorial** pw=**sigcomm**

```
# become root
sudo -s
# change to ONVM main directory
cd $ONVM_HOME
# configure DPDK to use NICs
./scripts/setup_nics.sh dpdk
```

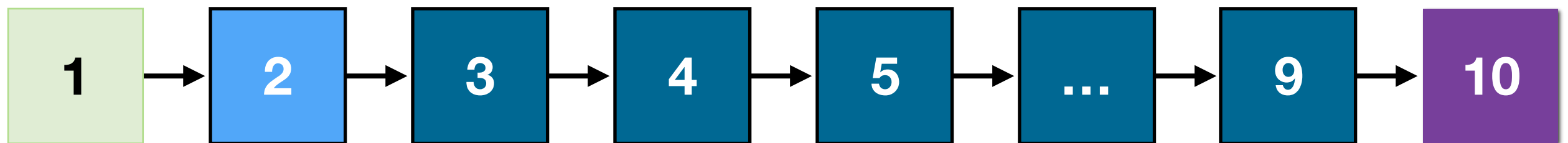
Legal stuff: You may only use these servers for running code related to the tutorial!

Testbed Setup

Servers are courtesy of Cloudlab.us NSF testbed

10 servers connected in a chain

Already have ONVM/DPDK installed and configured



Our goal:

- Send traffic from node 1 to node 10
- Forward all packets through the switches using DPDK

Run the basic forwarder

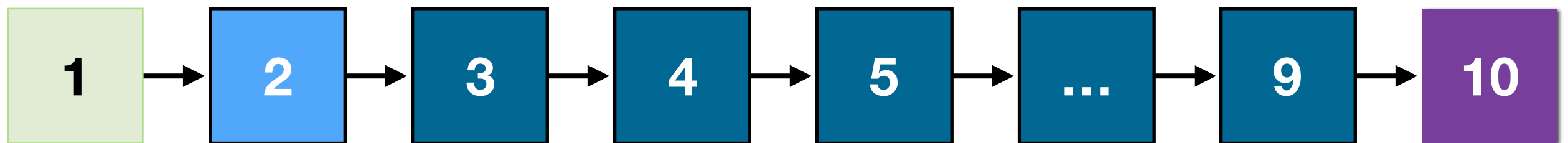
Change to the DPDK forwarding example directory

- **cd \$RTE_SDK/examples/skeleton**

Run the go script to start the program

- **./go.sh**

- this is equivalent to: `./build/basicfwd -l 1 -n 4`



I will run a workload on nodes 1/10

DPDK Basic Forward

Basic forwarding with DPDK

- [dpdk/examples/skeleton/basicfwd.c](#) ([docs](#))
- Create a memory pool to store packets
- Initialize NIC ports and configure with mem pool
- Initialize ring buffer for receiving packets from NIC
- Initialize ring buffer for transmitting packets to NIC
- For each port:
 - Dequeue a batch of packets
 - Transmit each packet out the opposite port

DPDK Limitations

Barebones I/O library

- No protocol processing (e.g., TCP stack)

Focused on running a single NF

- Need to specially design NFs to be able to run multiple functions at once
- NF monopolizes entire NIC port(s)

Extensive use of polling means high CPU usage

No management layer

- addressing, reliability, auto scaling, fault tolerance