# DPDK ARCHITECTURE AND ROADMAP DISCUSSION

KANNAN BABU RAMIA, INTEL
DEEPAK KUMAR JAIN, INTEL

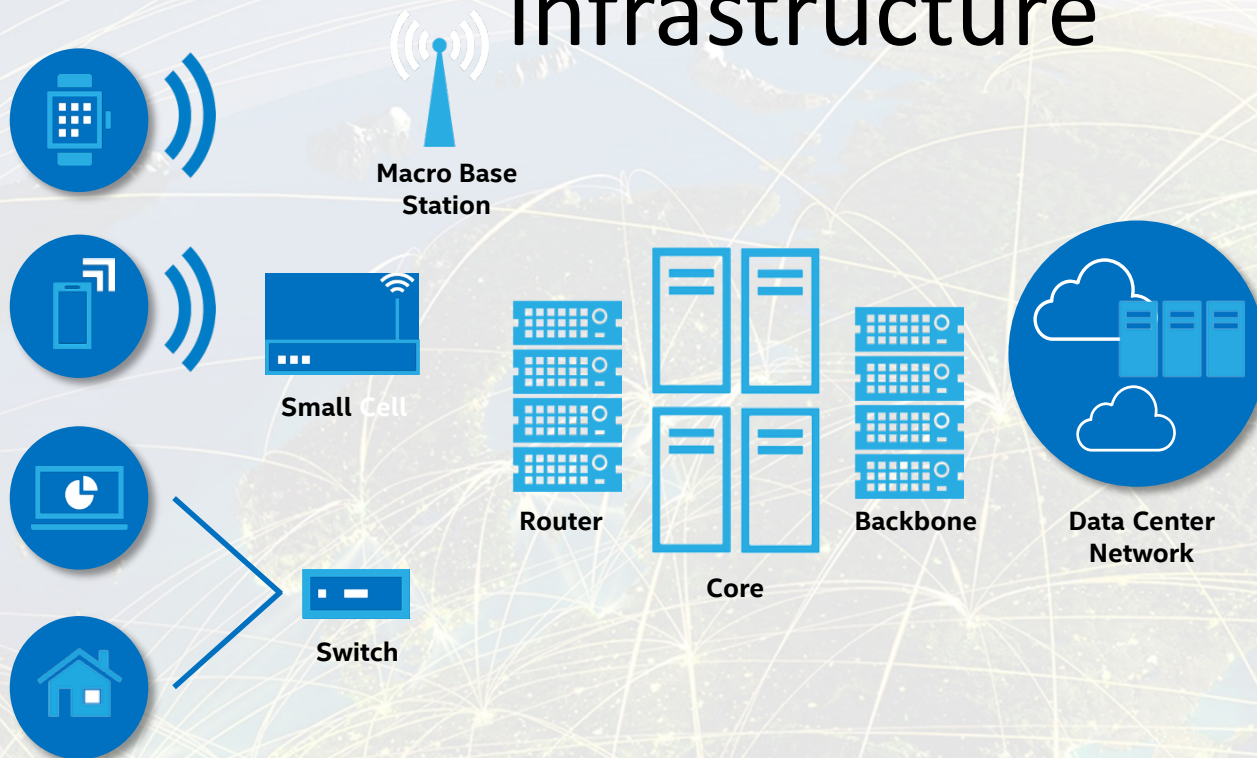# LEGAL DISCLAIMER

# Agenda

- ➤ Key trends in network transformation
- ➤ DPDK role
- ➤ DPDK Architecture
- ➤ Multi Architecture/ Multi vendor support
- ➤ Open source projects using DPDK
- ➤ DPDK Roadmap
- ➤ Open Questions

# End-to-end network infrastructure

**Macro Base Station**

**Small Cell**

**Switch**

**Router**

**Core**

**Backbone**

**Data Center Network**

**WORKLOAD CONVERGENCE**

**NETWORK VIRTUALIZATION**

**END-TO-END TRANSFORMATION**

# DPDK Generational Performance Gains

**L3Fwd Performance (MPPS)**

**IPV4 L3 Forwarding Performance of 64Byte Packets**

| Year | Value (MPPS) | Gbps |
|------|------|------|
| 2010 | 55 | 37 Gbps |
| 2011 | 80.1 | 53.8 Gbps |
| 2012 | 164.9 | 110.8 Gbps |
| 2014 | 255 | 171.4 Gbps |
| 2015 | 279.9 | 187.2 Gbps |
| 2016 | *346.7 | 233 Gbps |

**Year**

| *System Configuration | |
|------|------|
| **Hardware** | |
| Platform | SuperMicro* - X10DRX |
| CPU | Intel® Xeon® E5-2658 v4 Processor |
| Chipset | Intel® C612 chipset |
| Sockets | 2 |
| Cores per Socket | 14 (28 threads) |
| LL CACHE | 30 MB |
| QPI/DMI | 9.6GT/s |
| PCIe | Gen3x8 |
| MEMORY | DDR4 2400 MHz, 1Rx4 8GB (total 64GB), 4 Channel per Socket |
| NIC | 10 x Intel® Ethernet CNA XL710-QDA2PCI-Express Gen3 x8 Dual Port 40 GbE Ethernet NIC (1x40G/card) |
| NIC Mbps | 40,000 |
| BIOS | BIOS version: 1.0c (02/12/2015) |
| **Software** | |
| OS | Debian* 8.0 |
| Kernel version | 3.18.2 |
| Other | DPDK 2.2.0 |

# NFV – Life of a Packet



VEB – Virtual Embedded Bridge
TEP – Tunnel End Point
ACL – Access Control List
GFT – Generic Flow Table

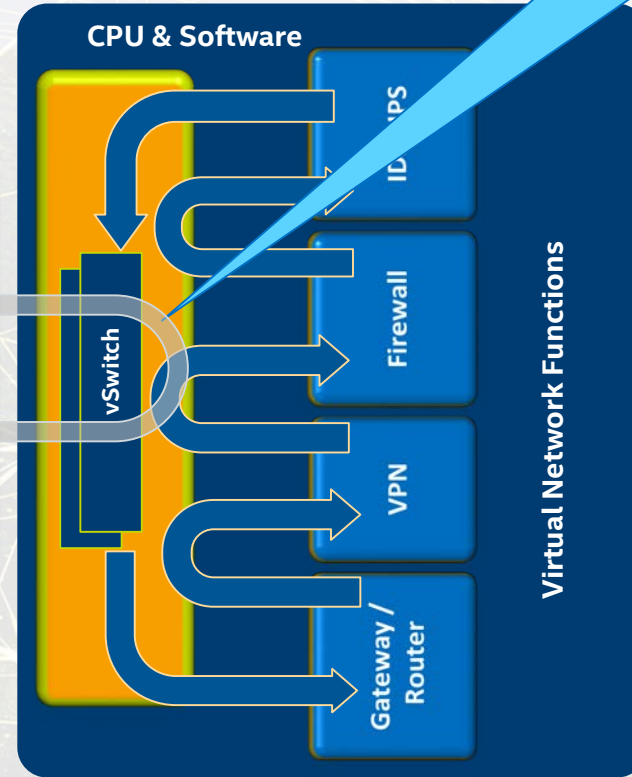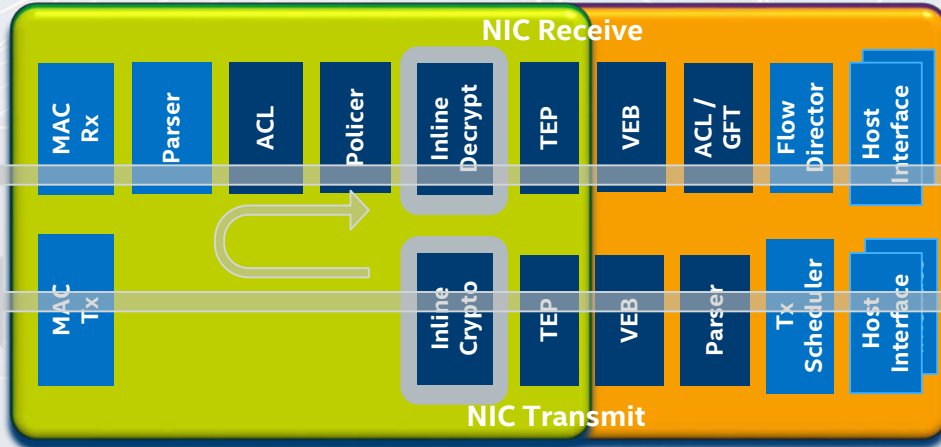| Outer Header (Underlay) | | | | | Inner Payload (Overlay) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Outer Ethernet Header | Outer IP Header | IPSec ESP transport | Outer UDP Header | Overlay Header (e.g. VXLAN) | Inner Ethernet Header | Inner IP Header | Inner L4 Header | Inner Data | IPSec ESP trailer | Outer CRC |

# NFV – Service Function Chaining



Intel® QuickAssist Technology **Crypto Accelerator**

VM to VM communication between Virtual Network Functions, needs multi-100s of Gbps throughput

**CPU & Software**

vSwitch

**Virtual Network Functions**

IDS / IPS
Firewall
VPN
Gateway / Router

**NIC Receive**

MAC Rx | Parser | ACL | Policer | Inline Decrypt | TEP | VEB | ACL / GFT | Flow Director | Host Interface

MAC Tx | Inline Crypto | TEP | VEB | Parser | Tx Scheduler | Host Interface
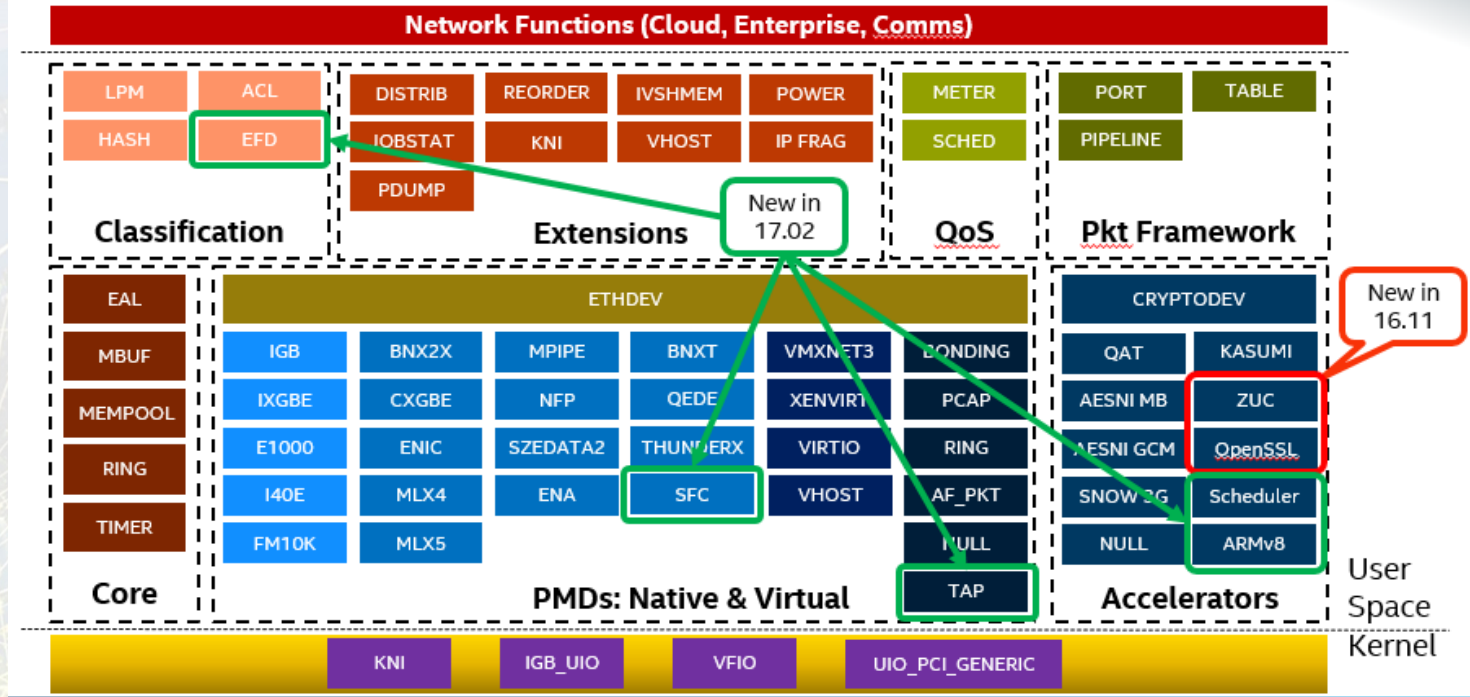
**NIC Transmit**

VEB – Virtual Embedded Bridge
TEP – Tunnel End Point
ACL – Access Control List
GFT – Generic Flow Table
VNF – Virtual Network Functions

# DPDK Architecture

## DPDK Fundamentals

- Implements run-to-completion and pipeline models
- No scheduler - all devices accessed by polling
- Supports 32-bit and 64-bit OSs with and without NUMA
- Scales from Intel® Atom™ to Intel® Xeon® processors
- Number of cores and processors is not limited
- Optimal packet allocation across DRAM channels
- Use of 2M & 1G hugepages and cache aligned structures
- Uses bulk concepts - processing 'n' packets simultaneously

# DPDK CONSUMPTION

## vSwitches

VPP

Open vSwitch

BESS

Lagopus

## Packet Generators

TRex

Pktgen

MoonGen

Ostinato

## DPDK in OS Distros

redhat.
Version 7.1 +

freeBSD
Version 10.1 +

ubuntu
Version 15.10 +

WIND
Version 6 +

CentOS
Version 7.1 +

fedora
Version 22 +

## Storage

SPDK

+ Many more

## vRouters

OPENCONTRAIL

CloudRouter

VPP

## TCP/IP Stacks

mTCP

TLDK & VPP

Seastar

LWIP DPDK

# DPDK Roadmap

*Q2'17 (v 17.05)*

- **Added Eventdev PMD.**
- **Added event driven programming model library (rte_eventdev).**
- **Added bit-rate calculation, latency stats and information metric library.**
- **Kept consistent PMD batching behaviour.**
- **Added VFIO hotplug and vmxnet3 version 3 support.**
- **Added MTU feature support to Virtio and Vhost.**
- **Added interrupt mode support for virtio-user.**

# DPDK Roadmap

**Q2'17 (v 17.08)**

- Generic QoS API
- Cryptodev Multi-Core SW Scheduler
- Generic Receive Offload
- Generic Flow Enhancements
- VF Port Reset for IXGBE
- API to Configure Queue Regions for RSS
- Support for IPFIX

# OPEN QUESTIONS?

▶ What is missing from DPDK?

▶ What are the major pain-points in using DPDK?

▶ What can be improved in DPDK? Build process? Logging?

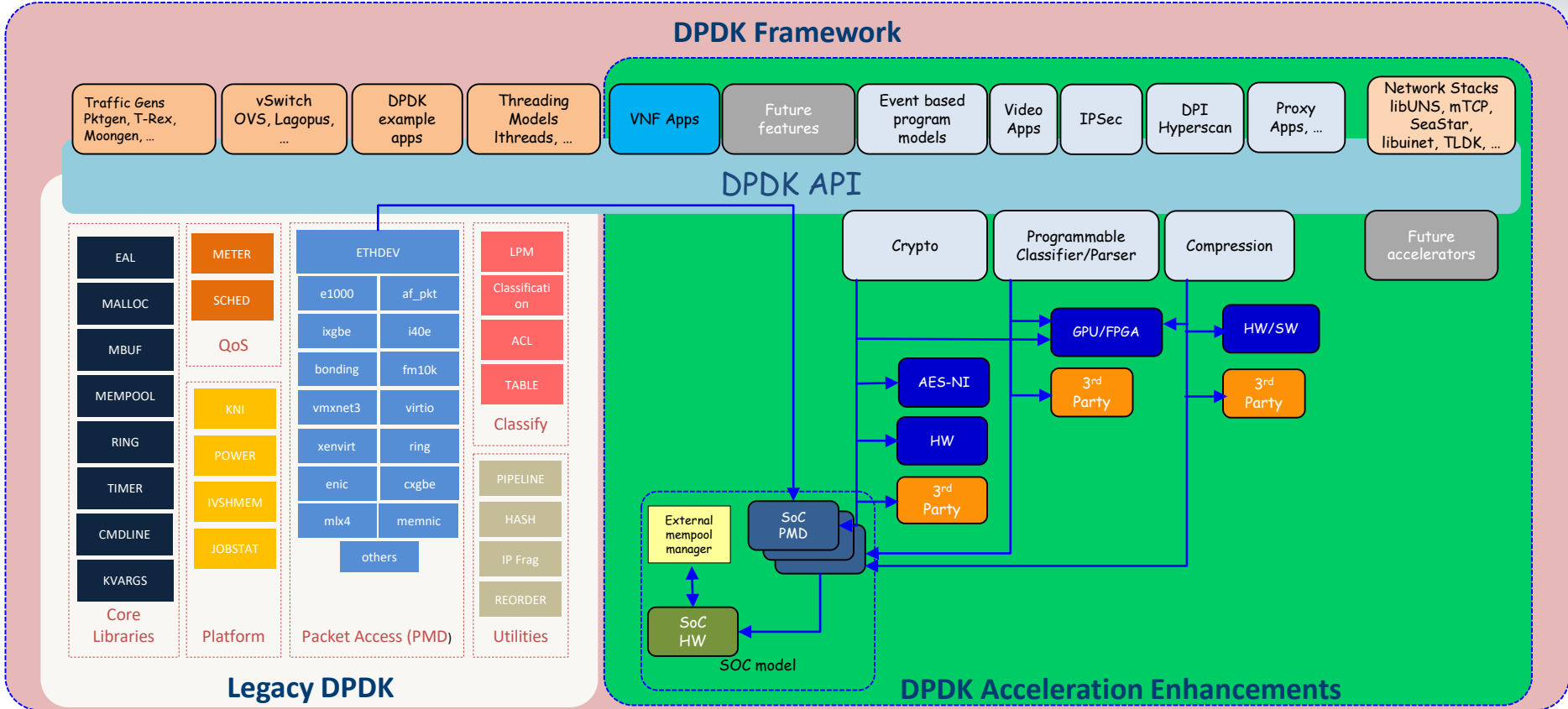▶ What are the big performance bottlenecks?
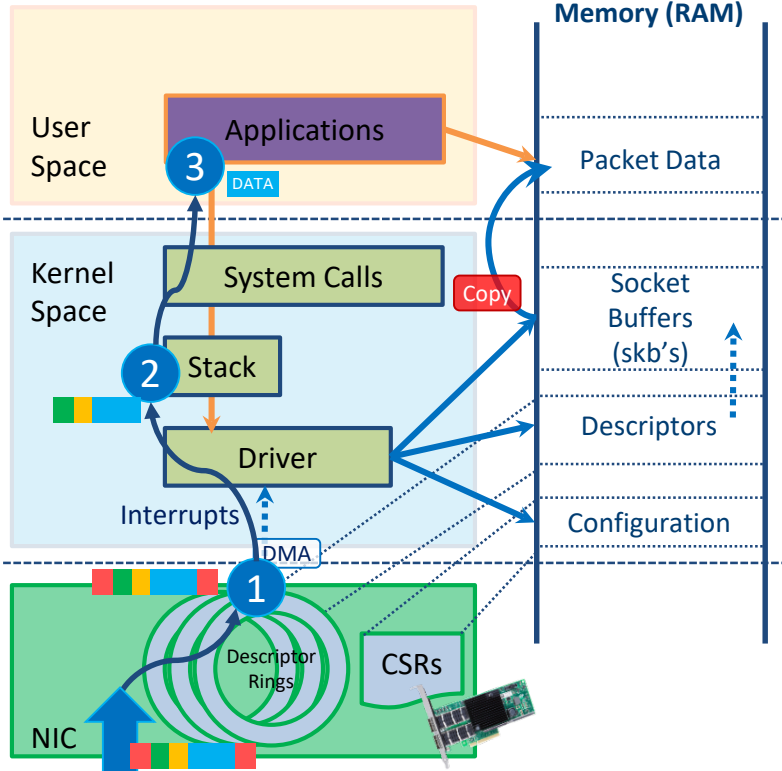
▶ Working with Kernel?

# DPDK Sample Apps

# DPDK Acceleration Enhancements

# Packet Processing Kernel vs. User Space
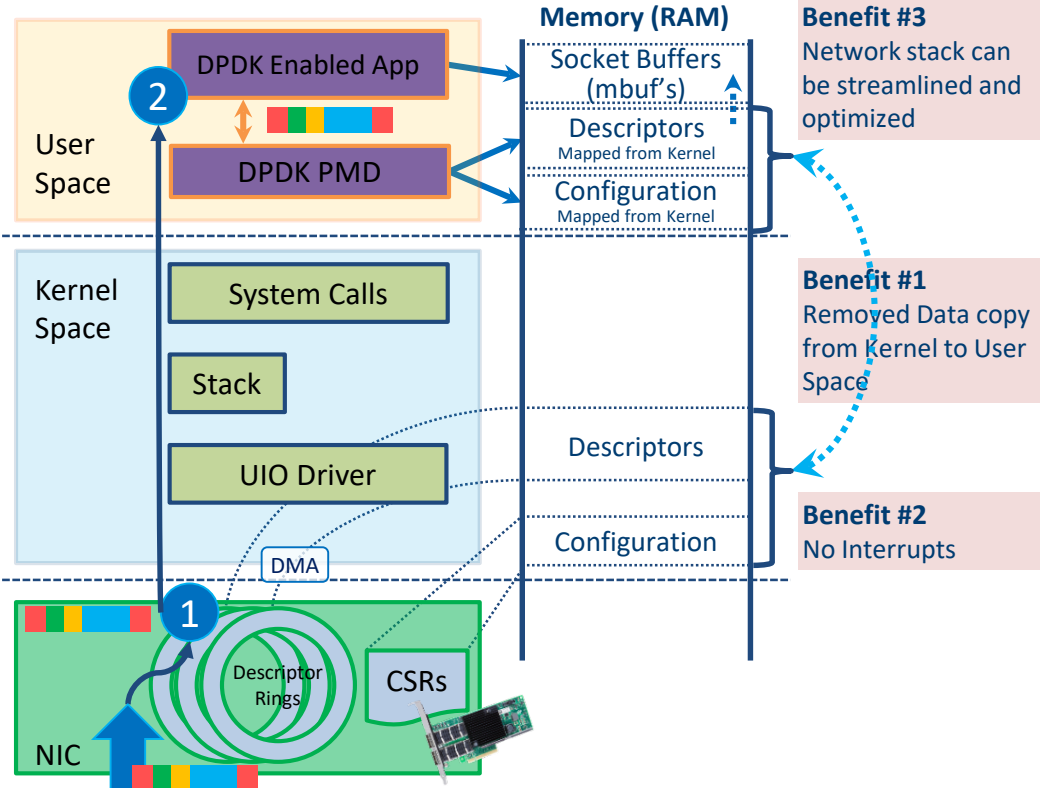


**Kernel Space Driver**

**User Space Driver with Zero Copy**

**Benefit #3**
Network stack can be streamlined and optimized

**Benefit #1**
Removed Data copy from Kernel to User Space

**Benefit #2**
No Interrupts

# DPDK IN-DEPTH

# PCIe* Connectivity and Core Usage

Using run-to-completion or pipeline software models
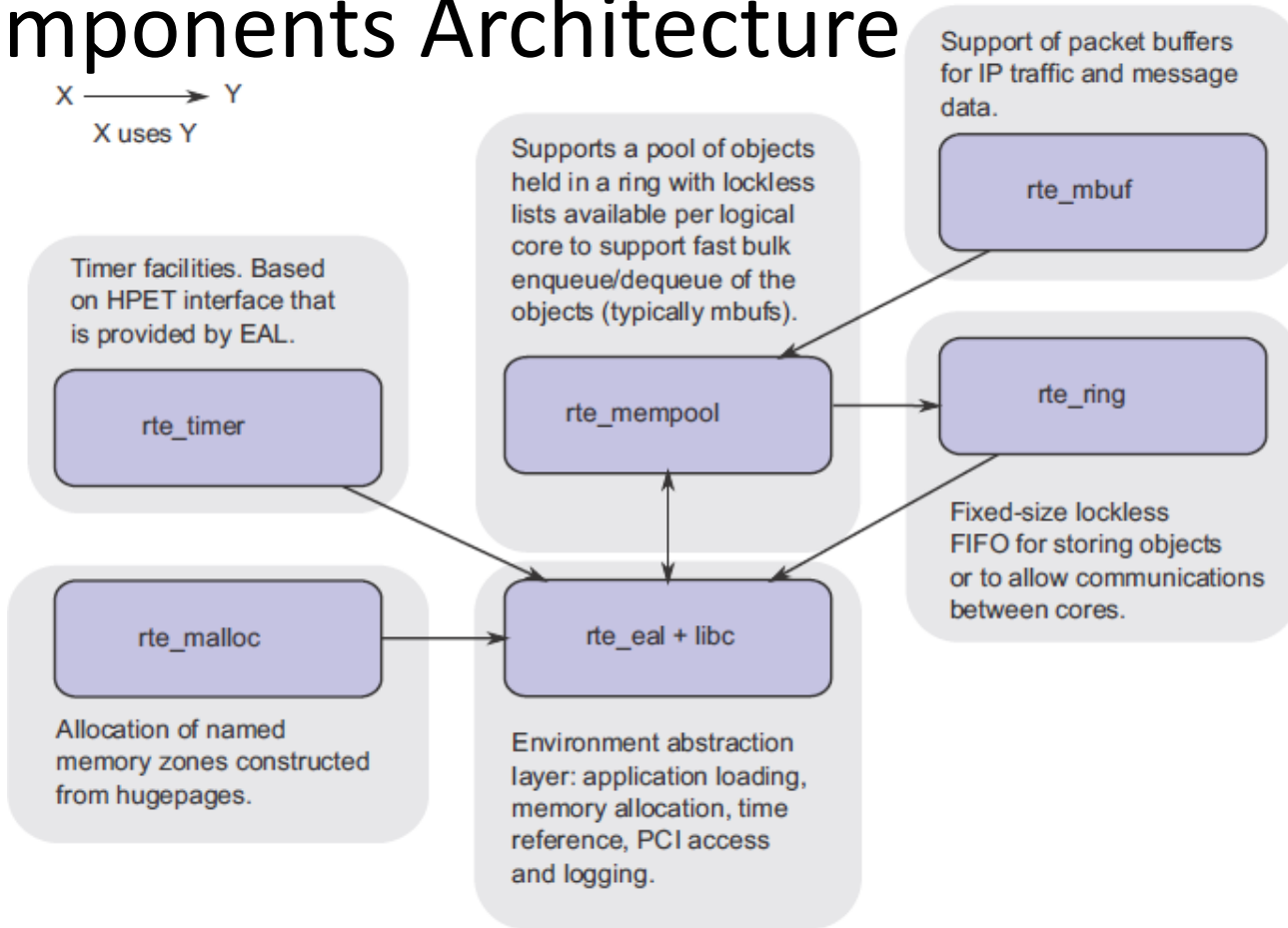


**Run to Completion Model**
- I/O and Application workload can be handled on a single core
- I/O can be scaled over multiple cores

**Pipeline Model**
- I/O application disperses packets to other cores
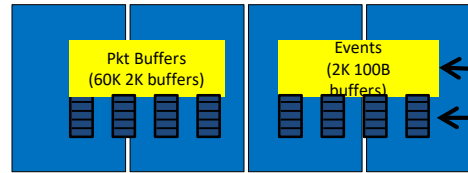- Application work performed on other cores

Can handle more I/O on fewer cores with vectorization

# Core Components Architecture

# DPDK model

# High Performance Components of DPDK

- Environment Abstraction Layer
  - Abstracts huge-page file system, provides multi-thread and multi-process support, etc.
- Memory Manager
  - Responsible for allocating pools of objects in memory. A pool is created in huge page memory space and uses a ring to store free objects. It also provides an alignment helper to ensure that objects are padded to spread them equally on all DRAM channels.
- Buffer Manager
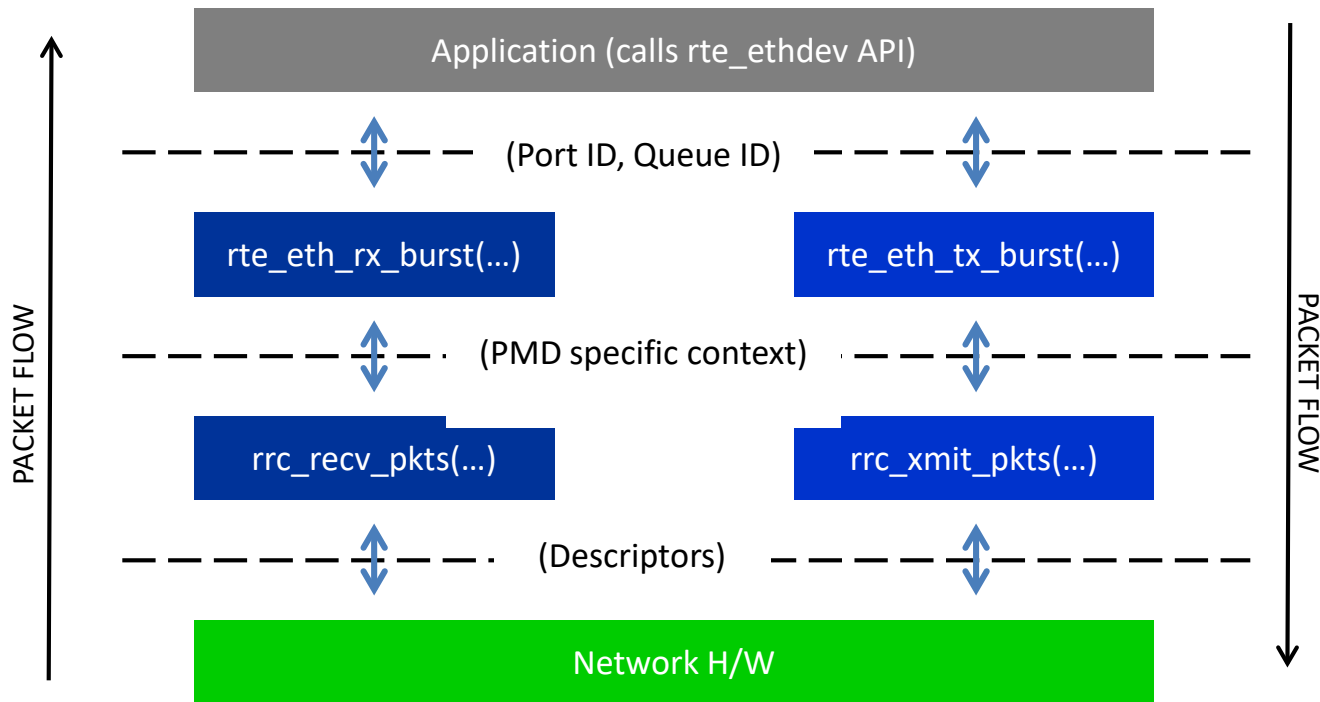  - Reduces by a significant amount the time the operating system spends allocating and de-allocating buffers. The Intel® DPDK pre-allocates fixed size buffers which are stored in memory pools.
- Queue Manager
  - Implements safe lockless queues, instead of using spinlocks, that allow different software components to process packets, while avoiding unnecessary wait times.
- Flow Classification
  - Provides an efficient mechanism which incorporates Intel® Streaming SIMD Extensions (Intel® SSE) to produce a hash based on tuple information so that packets may be placed into flows quickly for processing, thus greatly improving throughput.
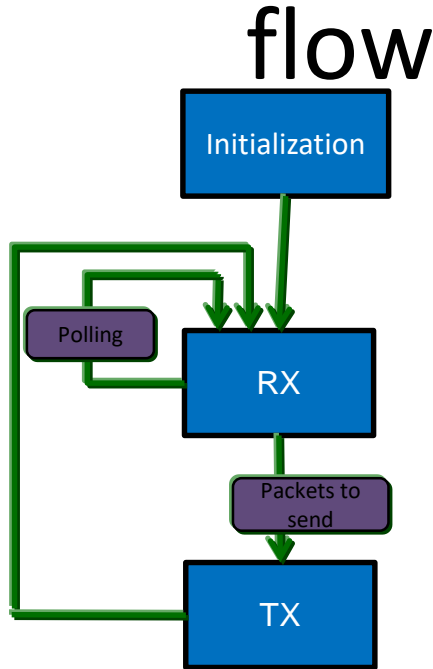
# EAL Initialization in a Linux Environment

# Ethernet Device Framework

# 30,000 ft overview of packet flow



1. Initialization
   o Init Memory Zones and Pools
   o Init Devices and Device Queues
   o Start Packet Forwarding Application
2. Packet Reception (RX)
   o Poll Devices' RX queues and receive packets in bursts
   o Allocate new RX buffers from per queue memory pools to stuff into descriptors
3. Packet Transmission (TX)
   o Transmit the received packets from RX
   o Free the buffers that we used to store the packets