

Intel[®] Data Plane Development Kit (Intel[®] DPDK)

Programmer's Guide

June 2013



INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL PRODUCTS. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO SALE AND/OR USE OF INTEL PRODUCTS INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

A "Mission Critical Application" is any application in which failure of the Intel Product could result, directly or indirectly, in personal injury or death. SHOULD YOU PURCHASE OR USE INTEL'S PRODUCTS FOR ANY SUCH MISSION CRITICAL APPLICATION, YOU SHALL INDEMNIFY AND HOLD INTEL AND ITS SUBSIDIARIES, SUBCONTRACTORS AND AFFILIATES, AND THE DIRECTORS, OFFICERS, AND EMPLOYEES OF EACH, HARMLESS AGAINST ALL CLAIMS COSTS, DAMAGES, AND EXPENSES AND REASONABLE ATTORNEYS' FEES ARISING OUT OF, DIRECTLY OR INDIRECTLY, ANY CLAIM OF PRODUCT LIABILITY, PERSONAL INJURY, OR DEATH ARISING IN ANY WAY OUT OF SUCH MISSION CRITICAL APPLICATION, WHETHER OR NOT INTEL OR ITS SUBCONTRACTOR WAS NEGLIGENT IN THE DESIGN, MANUFACTURE, OR WARNING OF THE INTEL PRODUCT OR ANY OF ITS PARTS.

Intel may make changes to specifications and product descriptions at any time, without notice. Designers must not rely on the absence or characteristics of any features or instructions marked "reserved" or "undefined". Intel reserves these for future definition and shall have no responsibility whatsoever for conflicts or incompatibilities arising from future changes to them. The information here is subject to change without notice. Do not finalize a design with this information.

The products described in this document may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Contact your local Intel sales office or your distributor to obtain the latest specifications and before placing your product order.

Copies of documents which have an order number and are referenced in this document, or other Intel literature, may be obtained by calling 1-800-548-4725, or go to: <http://www.intel.com/design/literature.htm>.

Any software source code reprinted in this document is furnished for informational purposes only and may only be used or copied and no license, express or implied, by estoppel or otherwise, to any of the reprinted source code is granted by this document.

Code Names are only for use by Intel to identify products, platforms, programs, services, etc. ("products") in development by Intel that have not been made commercially available to the public, i.e., announced, launched or shipped. They are never to be used as "commercial" names for products. Also, they are not intended to function as trademarks.

Intel and the Intel logo are trademarks of Intel Corporation in the U.S. and/or other countries.

*Other names and brands may be claimed as the property of others.

Copyright © 2012-2013, Intel Corporation. All rights reserved.



Revision History

Date	Revision	Description
June 2013	-002	Supports public software release 1.3.1
November 2012	-001	Supports public software release 1.2.3 Minor updates to Section 5.5.3.2 and Section 5.5.3.3 since last limited distribution release.



Contents

1.0	Introduction	9
1.1	Documentation Roadmap	9
1.2	Related Publications	9
Part 1: Architecture Overview		10
2.0	Overview	11
2.1	Development Environment	11
2.2	Environment Abstraction Layer	12
2.3	Core Components	12
2.3.1	Memory Manager (librte_malloc)	13
2.3.2	Ring Manager (librte_ring)	13
2.3.3	Memory Pool Manager (librte_mempool)	14
2.3.4	Network Packet Buffer Management (librte_mbuf)	14
2.3.5	Timer Manager (librte_timer)	14
2.4	Ethernet* Poll Mode Driver Architecture	14
2.5	Packet Forwarding Algorithm Support	14
2.6	librte_net	14
3.0	Environment Abstraction Layer	15
3.1	EAL in a Linux-userland Execution Environment	15
3.1.1	Initialization and Core Launching	16
3.1.2	Multi-process Support	17
3.1.3	Memory Mapping Discovery and Memory Reservation	17
3.1.4	PCI Access	17
3.1.5	Per-Core and Shared Variables	17
3.1.6	Logs	17
3.1.6.1	Trace and Debug Functions	17
3.1.7	CPU Feature Identification	17
3.1.8	User Space Interrupt and Alarm Handling	17
3.1.9	Blacklisting	17
3.1.10	Misc Functions	18
3.2	Memory Segments and Memory Zones (memzone)	18
4.0	Malloc Library	19
4.1	Cookies	19
4.2	Alignment and NUMA Constraints	19
4.3	Use Cases	19
5.0	Ring Library	20
5.1	References for Ring Implementation in FreeBSD*	21
5.2	Lockless Ring Buffer in Linux*	21
5.3	Additional Features	21
5.3.1	Name	21
5.3.2	Water Marking	21
5.3.3	Debug	21
5.4	Use Cases	21
5.5	Anatomy of a Ring Buffer	22
5.5.1	Single Producer Enqueue	22
5.5.1.1	Enqueue First Step	22
5.5.1.2	Enqueue Second Step	22
5.5.1.3	Enqueue Last Step	23
5.5.2	Single Consumer Dequeue	23



5.5.2.1	Dequeue First Step	24
5.5.2.2	Dequeue Second Step	24
5.5.2.3	Dequeue Last Step	25
5.5.3	Multiple Producers Enqueue	25
5.5.3.1	MC Enqueue First Step	25
5.5.3.2	MC Enqueue Second Step	26
5.5.3.3	MC Enqueue Third Step	27
5.5.3.4	MC Enqueue Fourth Step	27
5.5.3.5	MC Enqueue Last Step	28
5.5.4	Modulo 32-bit Indexes	28
5.6	References	29
6.0	Mempool Library	30
6.1	Cookies	30
6.2	Stats	30
6.3	Memory Alignment Constraints	30
6.4	Local Cache	31
6.5	Use Cases	32
7.0	Mbuf Library	33
7.1	Design of Packet Buffers	33
7.2	Buffers Stored in Memory Pools	34
7.3	Constructors	35
7.4	Allocating and Freeing mbufs	35
7.5	Manipulating mbufs	35
7.6	Meta Information	35
7.7	Direct and Indirect Buffers	35
7.8	Debug	36
7.9	Use Cases	36
8.0	Poll Mode Driver	37
8.1	Requirements and Assumptions	37
8.2	Design Principles	38
8.3	Logical Cores, Memory and NIC Queues Relationships	39
8.4	Device Identification and Configuration	39
8.4.1	Device Identification	39
8.4.2	Device Configuration	39
8.4.3	On-the-Fly Configuration	39
8.4.4	Configuration of Transmit and Receive Queues	40
8.5	Poll Mode Driver API	41
8.5.1	Generalities	41
8.5.2	Generic Packet Representation	41
8.5.3	Ethernet Device API	41
9.0	Timer Library	42
9.1	Implementation Details	42
9.2	Use Cases	43
9.3	References	43
10.0	Hash Library	44
10.1	Hash API Overview	44
10.2	Implementation Details	45
10.3	Use Case: Flow Classification	45
10.4	References	46
11.0	LPM Library	47
11.1	LPM API Overview	47
11.2	Implementation Details	47



11.3	Use Case: IPv4 Forwarding	47
11.4	References.....	48
12.0	Multi-process Support	49
12.1	Memory Sharing	49
12.2	Deployment Models.....	50
12.2.1	Symmetric/Peer Processes	50
12.2.2	Asymmetric/Non-Peer Processes	50
12.2.3	Running Multiple Independent Intel® DPDK Applications	51
12.2.4	Running Multiple Independent Groups of Intel® DPDK Applications	51
12.3	Multi-process Limitations	51
13.0	IXGBE/IGB Virtual Function Driver	53
13.1	SR-IOV Mode Utilization in an Intel® DPDK Environment	53
13.1.1	Physical and Virtual Function Infrastructure.....	54
13.1.1.1	Intel® 82599 10 Gigabit Ethernet Controller VF Infrastructure	54
13.1.1.2	Intel® 82576 Gigabit Ethernet Controller and Intel® Ethernet Controller I350 Family VF Infrastructure	55
13.1.2	Validated Hypervisors.....	56
13.1.3	Expected Guest Operating System in Virtual Machine.....	56
13.2	Setting Up a KVM Virtual Machine Monitor	56
14.0	Driver for VM Emulated Devices	60
14.1	Validated Hypervisors.....	60
14.2	Expected Guest Operating System in Virtual Machine.....	60
14.3	Setting Up a KVM Virtual Machine	60
14.4	Known Limitations of Emulated Devices	62
15.0	Kernel NIC Interface	64
15.1	The Intel® DPDK KNI Kernel Module.....	65
15.2	KNI Creation and Deletion.....	65
15.3	Intel® DPDK mbuf Flow	65
15.4	Use Case: Ingress.....	66
15.5	Use Case: Egress.....	66
15.6	Ethtool	66
16.0	Thread Safety of Intel® DPDK Functions	67
16.1	Fast-Path APIs.....	67
16.2	Performance Insensitive API	67
16.3	Library Initialization	68
16.4	Interrupt Thread.....	68
Part 2:	Development Environment.....	69
17.0	Source Organization.....	70
17.1	Makefiles and Config	70
17.2	Libraries	70
17.3	Applications	71
18.0	Development Kit Build System	72
18.1	Building the Development Kit Binary.....	72
18.1.1	Build Directory Concept	72
18.2	Building External Applications.....	74
18.3	Makefile Description	74
18.3.1	General Rules For Intel® DPDK Makefiles	74
18.3.2	Makefile Types	75
18.3.2.1	Application	75
18.3.2.2	Library.....	75



18.3.2.3	Install	75
18.3.2.4	Kernel Module	75
18.3.2.5	Objects	75
18.3.2.6	Misc	75
18.3.3	Useful Variables Provided by the Build System	75
18.3.4	Variables that Can be Set/Overridden in a Makefile Only	76
18.3.5	Variables that can be Set/Overridden by the User on the Command Line Only	77
18.3.6	Variables that Can be Set/Overridden by the User in a Makefile or Command Line	77
19.0	Development Kit Root Makefile Help	78
19.1	Configuration Targets	78
19.2	Build Targets	78
19.3	Install Targets	78
19.4	Test Targets	79
19.5	Documentation Targets	79
19.6	Deps Targets	79
19.7	Misc Targets	79
19.8	Other Useful Command-line Variables	79
19.9	Make in a Build Directory	80
19.10	Compiling for Debug	80
20.0	Extending the Intel® DPDK	81
20.1	Example: Adding a New Library libfoo	81
20.1.1	Example: Using libfoo in the Test Application	82
21.0	Building Your Own Application	83
21.1	Compiling a Sample Application in the Development Kit Directory	83
21.2	Build Your Own Application Outside the Development Kit	83
21.3	Customizing Makefiles	83
21.3.1	Application Makefile	83
21.3.2	Library Makefile	84
21.3.3	Customize Makefile Actions	84
22.0	External Application/Library Makefile help	85
22.1	Prerequisites	85
22.2	Build Targets	85
22.3	Help Targets	85
22.4	Other Useful Command-line Variables	85
22.5	Make from Another Directory	86
Part 3:	Performance Optimization	87
23.0	Performance Optimization Guidelines	88
23.1	Introduction	88
24.0	Writing Efficient Code	89
24.1	Memory	89
24.1.1	Memory Copy: Do not Use libc in the Data Plane	89
24.1.2	Memory Allocation	89
24.1.3	Concurrent Access to the Same Memory Area	89
24.1.4	NUMA	90
24.1.5	Distribution Across Memory Channels	90
24.2	Communication Between Icores	90
24.3	PMD Driver	90
24.3.1	Lower Packet Latency	91
24.4	Locks and Atomic Operations	91



24.5	Coding Considerations	91
24.5.1	Inline Functions	91
24.5.2	Branch Prediction	91
24.6	Setting the Target CPU Type	92
25.0	Profile Your Application	93
26.0	Glossary	94

Figures

1	Core Components Architecture	13
2	EAL Initialization in a Linux Application Environment	16
3	Ring Structure	21
4	Two Channels and Quad-ranked DIMM Example	31
5	Three Channels and Two Dual-ranked DIMM Example	31
6	A mempool in Memory with its Associated Ring	32
7	An mbuf with One Segment	34
8	An mbuf with Three Segments	34
9	Timer Operation	43
10	Memory Sharing in the Intel® DPDK Multi-process Sample Application	50
11	Virtualization for a Single Port NIC in SR-IOV Mode	54
12	Performance Benchmark Setup	59
13	Components of a DPDK KNI Application	64
14	Packet Flow via mbufs in the Intel® DPDK KNI	66



1.0 Introduction

This document provides software architecture information, development environment information and optimization guidelines.

For programming examples and for instructions on compiling and running each sample application, see the *Intel® DPDK Sample Application's User Guide* for details.

For general information on compiling and running applications, see the *Intel® DPDK Getting Started Guide*.

1.1 Documentation Roadmap

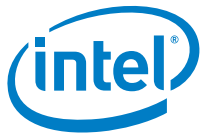
The following is a list of Intel® DPDK documents in the suggested reading order:

- **Release Notes:** Provides release-specific information, including supported features, limitations, fixed issues, known issues and so on. Also, provides the answers to frequently asked questions in FAQ format.
- **Getting Started Guide:** Describes how to install and configure the Intel® DPDK software; designed to get users up and running quickly with the software.
- **Programmer's Guide** (this document): Describes:
 - The software architecture and how to use it (through examples), specifically in a Linux* application (linuxapp) environment
 - The content of the Intel® DPDK, the build system (including the commands that can be used in the root Intel® DPDK Makefile to build the development kit and applications) and guidelines for porting an application
 - Optimizations used in the software and those that should be considered for new developmentA glossary of terms is also provided.
- **API Reference:** Provides detailed information about Intel® DPDK functions, data structures and other programming constructs.
- **Sample Applications User Guide:** Describes a set of sample applications. Each chapter describes a sample application that showcases specific functionality and provides instructions on how to compile, run and use the sample application.

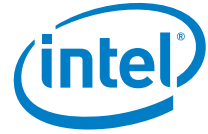
1.2 Related Publications

The follow documents provide information that is relevant to the development of applications using the Intel® DPDK:

- *Intel® 64 and IA-32 Architectures Software Developer's Manual Volume 3A: System Programming Guide*



Part 1: Architecture Overview



2.0 Overview

This section gives a global overview of the architecture of Intel® Data Plane Development Kit (Intel® DPDK).

The main goal of the Intel® DPDK is to provide a simple, complete framework for fast packet processing in data plane applications. Users may use the code to understand some of the techniques employed, to build upon for prototyping or to add their own protocol stacks. Alternative ecosystem options that use the Intel® DPDK are available.

The framework creates a set of libraries for specific environments through the creation of an Environment Abstraction Layer (EAL), which may be specific to a mode of the Intel® architecture (32-bit or 64-bit), Linux* user space compilers or a specific platform. These environments are created through the use of make files and configuration files. Once the EAL library is created, the user may link with the library to create their own applications. Other libraries, outside of EAL, including the Hash, Longest Prefix Match (LPM) and rings libraries are also provided. Sample applications are provided to help show the user how to use various features of the Intel® DPDK.

The Intel® DPDK implements a *run to completion* model for packet processing, where all resources must be allocated prior to calling Data Plane applications, running as execution units on logical processing cores. The model does not support a scheduler and all devices are accessed by polling. The primary reason for not using interrupts is the performance overhead imposed by interrupt processing.

In addition to the run-to-completion model, a pipeline model may also be used by passing packets or messages between cores via the rings. This allows work to be performed in stages and may allow more efficient use of code on cores.

2.1 Development Environment

The Intel® DPDK project installation requires Linux and the associated toolchain, such as one or more compilers, assembler, make utility, editor and various libraries to create the Intel® DPDK components and libraries.

Once these libraries are created for the specific environment and architecture, they may then be used to create the user's data plane application.

When creating applications for the Linux user space, the glibc library is used.

For Intel® DPDK applications, two environmental variables (RTE_SDK and RTE_TARGET) must be configured before compiling the applications. The following are examples of how the variables can be set:

```
export RTE_SDK=/home/user/DPDK
export RTE_TARGET=x86_64-default-linuxapp-gcc
```

See the *Intel® DPDK Getting Started Guide* for information on setting up the development environment.



2.2 Environment Abstraction Layer

The Environment Abstraction Layer (EAL) provides a generic interface that hides the environment specifics from the applications and libraries. The services provided by the EAL are:

- Intel® DPDK loading and launching
- Support for multi-process and multi-thread execution types
- Core affinity/assignment procedures
- System memory allocation/de-allocation
- Atomic/lock operations
- Time reference
- PCI bus access
- Trace and debug functions
- CPU feature identification
- Interrupt handling
- Alarm operations

The EAL is fully described in [Environment Abstraction Layer](#).

2.3 Core Components

The *core components* are a set of libraries that provide all the elements needed for high-performance packet processing applications.

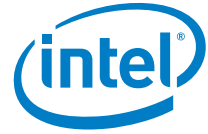
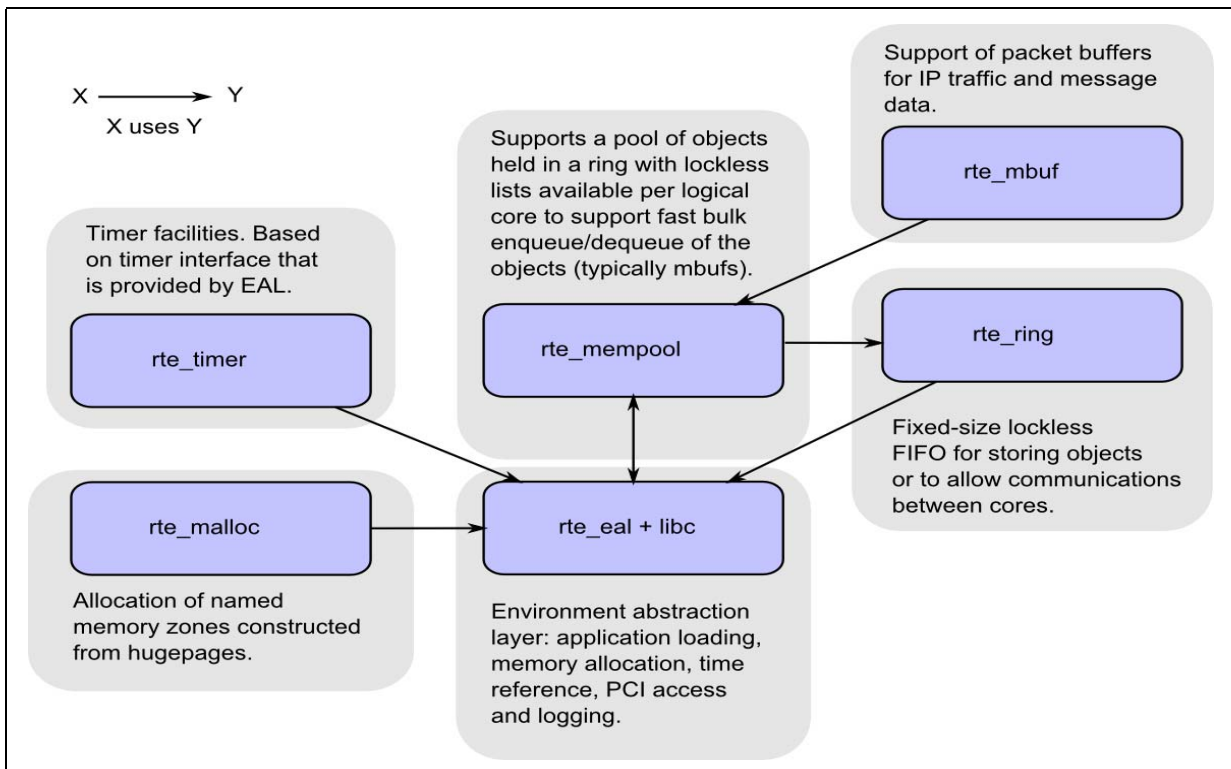


Figure 1. Core Components Architecture



2.3.1 Memory Manager (librte_malloc)

The `librte_malloc` library provides an API to allocate memory from the memzones created from the hugepages instead of the heap. This helps when allocating large numbers of items that may become susceptible to TLB misses when using typical 4k heap pages in the Linux user space environment.

This memory allocator is fully described in [Malloc Library](#).

2.3.2 Ring Manager (librte_ring)

The ring structure provides a lockless multi-producer, multi-consumer FIFO API in a finite size table. It has some advantages over lockless queues; easier to implement, adapted to bulk operations and faster. A ring is used by the [Memory Pool Manager \(librte_mempool\)](#) and may be used as a general communication mechanism between cores and/or execution blocks connected together on a logical core.

This ring buffer and its usage are fully described in [Ring Library](#).



2.3.3 Memory Pool Manager (`librte_mempool`)

The Memory Pool Manager is responsible for allocating pools of objects in memory. A pool is identified by name and uses a ring to store free objects. It provides some other optional services, such as a per-core object cache and an alignment helper to ensure that objects are padded to spread them equally on all RAM channels.

This memory pool allocator is described in [Mempool Library](#).

2.3.4 Network Packet Buffer Management (`librte_mbuf`)

The mbuf library provides the facility to create and destroy buffers that may be used by the Intel® DPDK application to store message buffers. The message buffers are created at startup time and stored in a mempool, using the Intel® DPDK mempool library.

This library provide an API to allocate/free mbufs, manipulate control message buffers (`ctrlmbuf`) which are generic message buffers, and packet buffers (`pkmbuf`) which are used to carry network packets.

Network Packet Buffer Management is described in [Mbuf Library](#).

2.3.5 Timer Manager (`librte_timer`)

This library provides a timer service to Intel® DPDK execution units, providing the ability to execute a function asynchronously. It can be periodic function calls, or just a one-shot call. It uses the timer interface provided by the Environment Abstraction Layer (EAL) to get a precise time reference and can be initiated on a per-core basis as required.

The library documentation is available in [Timer Library](#).

2.4 Ethernet* Poll Mode Driver Architecture

The Intel® DPDK includes Poll Mode Drivers (PMDs) for 1 GbE and 10 GbE Ethernet controllers which are designed to work without asynchronous, interrupt-based signalling mechanisms.

See [Poll Mode Driver](#).

2.5 Packet Forwarding Algorithm Support

The Intel® DPDK includes Hash (`librte_hash`) and Longest Prefix Match (LPM, `librte_lpm`) libraries to support the corresponding packet forwarding algorithms.

See [Hash Library](#) and [LPM Library](#) for more information.

2.6 `librte_net`

The `librte_net` library is a collection of IP protocol definitions and convenience macros. It is based on code from the FreeBSD* IP stack and contains protocol numbers (for use in IP headers), IP-related macros, IPv4/IPv6 header structures and TCP, UDP and SCTP header structures.





3.0 Environment Abstraction Layer

The Environment Abstraction Layer (EAL) is responsible for gaining access to low-level resources such as hardware and memory space. It provides a generic interface that hides the environment specifics from the applications and libraries. It is the responsibility of the initialization routine to decide how to allocate these resources (that is, memory space, PCI devices, timers, consoles, and so on).

Typical services expected from the EAL are:

- Intel® DPDK Loading and Launching: The Intel® DPDK and its application are linked as a single application and must be loaded by some means.
- Core Affinity/Assignment Procedures: The EAL provides mechanisms for assigning execution units to specific cores as well as creating execution instances.
- System Memory Reservation: The EAL facilitates the reservation of different memory zones, for example, physical memory areas for device interactions.
- PCI Address Abstraction: The EAL provides an interface to access PCI address space.
- Trace and Debug Functions: Logs, `dump_stack`, `panic` and so on.
- Utility Functions: Spinlocks and atomic counters that are not provided in `libc`.
- CPU Feature Identification: Determine at runtime if a particular feature, for example, Intel® AVX is supported. Determine if the current CPU supports the feature set that the binary was compiled for.
- Interrupt Handling: Interfaces to register/unregister callbacks to specific interrupt sources.
- Alarm Functions: Interfaces to set/remove callbacks to be run at a specific time.

3.1 EAL in a Linux-userland Execution Environment

In a Linux user space environment, the Intel® DPDK application runs as a user-space application using the `pthread` library. PCI information about devices and address space is discovered through the `/sys` kernel interface and through a module called `igb_uio`. Refer to the *UIO: User-space drivers* documentation in the Linux kernel. This memory is `mmap`'d in the application.

The EAL performs physical memory allocation using `mmap()` in `hugetlbfs` (using huge page sizes to increase performance). This memory is exposed to Intel® DPDK service layers such as the [Mempool Library](#).

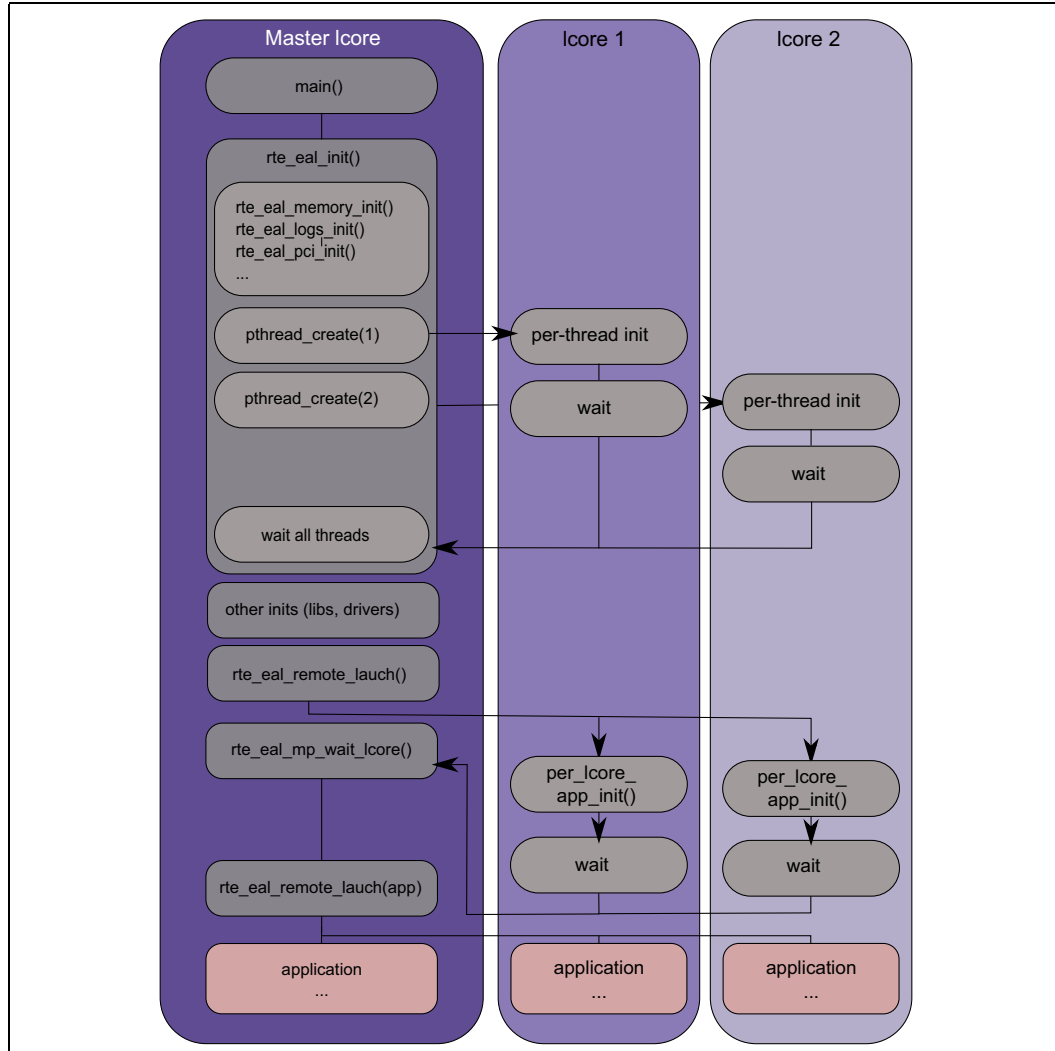
At this point, the Intel® DPDK services layer will be initialized, then through `pthread_setaffinity` calls, each execution unit will be assigned to a specific logical core to run as a user-level thread.

The time reference is provided by the CPU Time-Stamp Counter (TSC) or by the HPET kernel API through a `mmap()` call.

3.1.1 Initialization and Core Launching

Part of the initialization is done by the start function of glibc. A check is also performed at initialization time to ensure that the micro architecture type chosen in the config file is supported by the CPU. Then, the `main()` function is called. The core initialization and launch is done in `rte_eal_init()` (see the API documentation). It consist of calls to the pthread library (more specifically, `pthread_self()`, `pthread_create()`, and `pthread_setaffinity_np()`).

Figure 2. EAL Initialization in a Linux Application Environment



Note: Initialization of objects, such as memory zones, rings, memory pools, lpm tables and hash tables, should be done as part of the overall application initialization on the master lcore. The creation and initialization functions for these objects are not multi-thread safe. However, once initialized, the objects themselves can safely be used in multiple threads simultaneously.



3.1.2 Multi-process Support

The Linuxapp EAL allows a multi-process as well as a multi-threaded (pthread) deployment model. See [Chapter 12.0, “Multi-process Support”](#) for more details.

3.1.3 Memory Mapping Discovery and Memory Reservation

The allocation of large contiguous physical memory is done using the `hugetlbfs` kernel filesystem. The EAL provides an API to reserve named memory zones in this contiguous memory.

3.1.4 PCI Access

The EAL uses the `/sys/bus/pci` utilities provided by the kernel to scan the content on the PCI bus.

To access PCI memory, a kernel module called `igb_uio` provides a `/dev/uioX` device file that can be `mmap`'d to obtain access to PCI address space from the application. It uses the `uio` kernel feature (userland driver).

3.1.5 Per-Core and Shared Variables

Note: *Core* refers to a logical execution unit of the processor, sometimes called a hardware thread.

Shared variables are the default behavior. Per-core variables are implemented using *Thread Local Storage* (TLS) to provide per-thread local storage.

3.1.6 Logs

A logging API is provided by EAL. By default, in a Linux application, logs are sent to `syslog` and also to the console. However, the log function can be overridden by the user to use a different logging mechanism.

3.1.6.1 Trace and Debug Functions

There are some debug functions to dump the stack in `glibc`. The `rte_panic()` function can voluntarily provoke a `SIG_ABORT`, which can trigger the generation of a core file, readable by `gdb`.

3.1.7 CPU Feature Identification

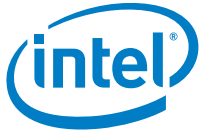
The EAL can query the CPU at runtime (using the `rte_cpu_get_feature()` function) to determine which CPU features are available.

3.1.8 User Space Interrupt and Alarm Handling

The EAL creates a host thread to poll the UIO device file descriptors to detect the interrupts. Callbacks can be registered or unregistered by the EAL functions for a specific interrupt event and are called in the host thread asynchronously. The EAL also allows timed callbacks to be used in the same way as for NIC interrupts.

3.1.9 Blacklisting

The EAL PCI device blacklist functionality can be used to mark certain NIC ports as blacklisted, so they are ignored by the Intel® DPDK. The ports to be blacklisted are identified using the PCIe* description (Domain:Bus:Device.Function).



3.1.10 Misc Functions

Locks and atomic operations are per-architecture (i686 and x86_64).

3.2 Memory Segments and Memory Zones (memzone)

The mapping of physical memory is provided by this feature in the EAL. As physical memory can have gaps, the memory is described in a table of descriptors, and each descriptor (called `rte_memseg`) describes a contiguous portion of memory.

On top of this, the memzone allocator's role is to reserve contiguous portions of physical memory. These zones are identified by a unique name when the memory is reserved.

The `rte_memzone` descriptors are also located in the configuration structure. This structure is accessed using `rte_eal_get_configuration()`. The lookup (by name) of a memory zone returns a descriptor containing the physical address of the memory zone.

Memory zones can be reserved with specific start address alignment by supplying the `align` parameter (by default, they are aligned to cache line size). The alignment value should be a power of two and not less than the cache line size (64 bytes). Memory zones can also be reserved from either 2 MB or 1 GB hugepages, provided that both are available on the system.

§ §



4.0 Malloc Library

The `librte_malloc` library provides an API to allocate any-sized memory.

The objective of this library is to provide malloc-like functions to allow allocation from hugepage memory and to facilitate application porting. The *Intel® DPDK API Reference* manual describes the available functions.

Typically, these kinds of allocations should not be done in data plane processing because they are slower than pool-based allocation and make use of locks within the allocation and free paths. However, they can be used in configuration code.

Refer to the `rte_malloc()` function description in the *Intel® DPDK API Reference* manual for more information.

4.1 Cookies

When `CONFIG_RTE_MALLOC_DEBUG` is enabled, the allocated memory contains overwrite protection fields to help identify buffer overflows.

4.2 Alignment and NUMA Constraints

The `rte_malloc()` takes an `align` argument that can be used to request a memory area that is aligned on a multiple of this value (which must be a power of two).

4.3 Use Cases

This library is needed by an application that requires malloc-like functions.

§ §

5.0 Ring Library

The ring allows the management of queues. Instead of having a linked list of infinite size, the *rte_ring* has the following properties:

- FIFO
- Maximum size is fixed, the pointers are stored in a table
- Lockless implementation
- Multi-consumer or single-consumer dequeue
- Multi-consumer or single-producer enqueue
- Bulk dequeue - Dequeues the specified count of objects if successful; otherwise fails
- Bulk enqueue - Enqueues the specified count of objects if successful; otherwise fails
- Burst dequeue - Dequeue the maximum available objects if the specified count cannot be fulfilled
- Burst enqueue - Enqueue the maximum available objects if the specified count cannot be fulfilled

The advantages of this data structure over a linked list queue are as follows:

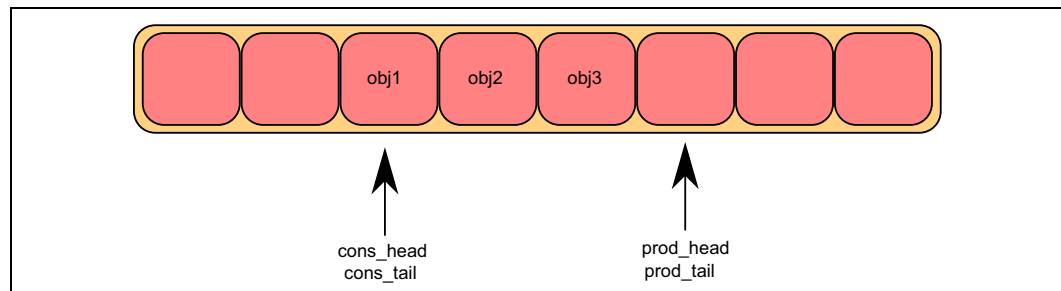
- Faster; only requires a single Compare-And-Swap instruction of `sizeof(void *)` instead of several double-Compare-And-Swap instructions.
- Simpler than a full lockless queue.
- Adapted to bulk enqueue/dequeue operations. As pointers are stored in a table, a dequeue of several objects will not produce as many cache misses as in a linked queue. Also, a bulk dequeue of many objects does not cost more than a dequeue of a simple object.

The disadvantages:

- Size is fixed
- Having many rings costs more in terms of memory than a linked list queue. An empty ring contains at least N pointers.

A simplified representation of a Ring is shown in [Figure 3](#) with consumer and producer head and tail pointers to objects stored in the data structure.

Figure 3. Ring Structure



5.1 References for Ring Implementation in FreeBSD*

The following code was added in FreeBSD 8.0, and is used in some network device drivers (at least in Intel drivers):

- [bufring.c in FreeBSD](#)
- [bufring.h in FreeBSD](#)

5.2 Lockless Ring Buffer in Linux*

The following is a link describing the [Linux Lockless Ring Buffer Design](#).

5.3 Additional Features

5.3.1 Name

A ring is identified by a unique name. It is not possible to create two with the same name (`rte_ring_create()` returns NULL if this is attempted).

5.3.2 Water Marking

The ring can have a high water mark (threshold). Once an enqueue operation reaches the high water mark, the producer is notified, if the water mark is configured.

This mechanism can be used, for example, to exert a back pressure on I/O to inform the LAN to PAUSE.

5.3.3 Debug

When debug is enabled (`CONFIG RTE_LIBRTE_RING_DEBUG` is set), the library stores some per-ring statistic counters about the number of enqueues/dequeues. These statistics are per-core to avoid concurrent accesses or atomic operations.

5.4 Use Cases

Use cases for the Ring library include:

- Communication between applications in the Intel® DPDK
- Used by memory pool allocator

5.5 Anatomy of a Ring Buffer

This section explains how a ring buffer operates. The ring structure is composed of two head and tail couples; one is used by producers and one is used by the consumers. The figures of the following sections refer to them as `prod_head`, `prod_tail`, `cons_head` and `cons_tail`.

Each figure represents a simplified state of the ring, which is a circular buffer. The content of the function local variables is represented on the top of the figure, and the content of ring structure is represented on the bottom of the figure.

5.5.1 Single Producer Enqueue

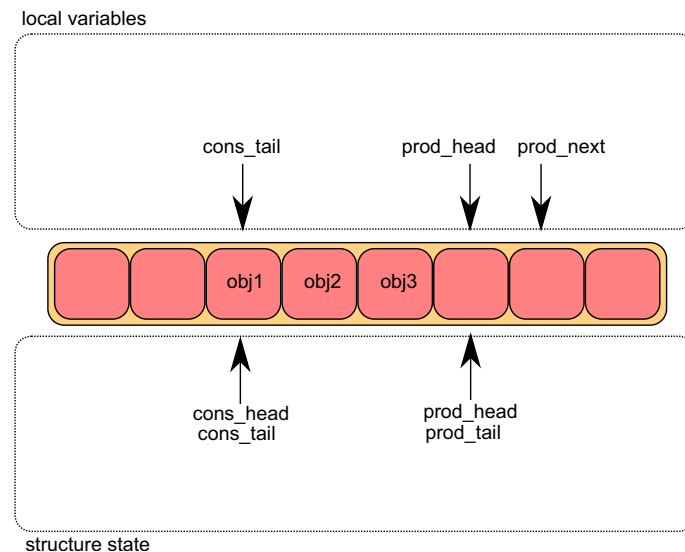
This section explains what occurs when a producer adds an object to the ring. In this example, only the producer head and tail (`prod_head` and `prod_tail`) are modified, and there is only one producer.

The initial state is to have a `prod_head` and `prod_tail` pointing at the same location.

5.5.1.1 Enqueue First Step

First, `ring->prod_head` and `ring->cons_tail` are copied in local variables. The `prod_next` local variable points to the next element of the table, or several elements after in case of bulk enqueue.

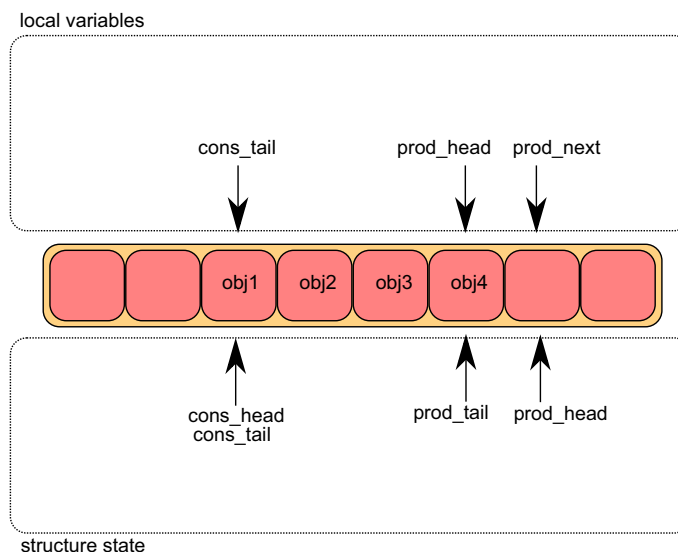
If there is not enough room in the ring (this is detected by checking `cons_tail`), it returns an error.



5.5.1.2 Enqueue Second Step

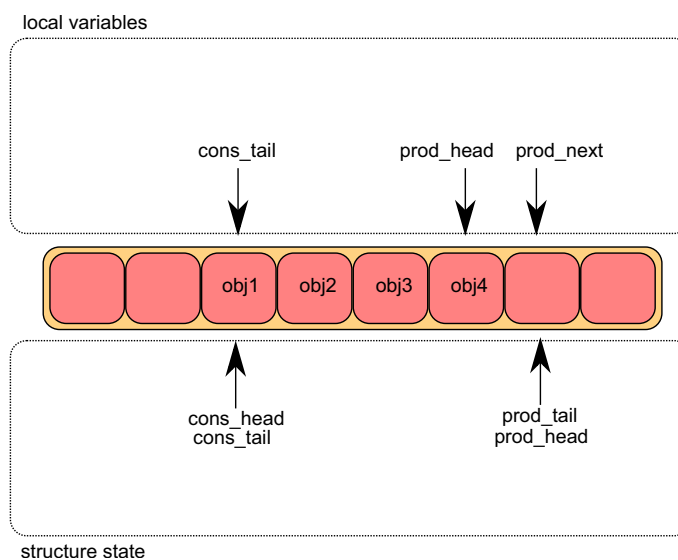
The second step is to modify `ring->prod_head` in ring structure to point to the same location as `prod_next`.

A pointer to the added object is copied in the ring (`obj4`).



5.5.1.3 Enqueue Last Step

Once the object is added in the ring, `ring->prod_tail` in the ring structure is modified to point to the same location as `ring->prod_head`. The enqueue operation is finished.



5.5.2 Single Consumer Dequeue

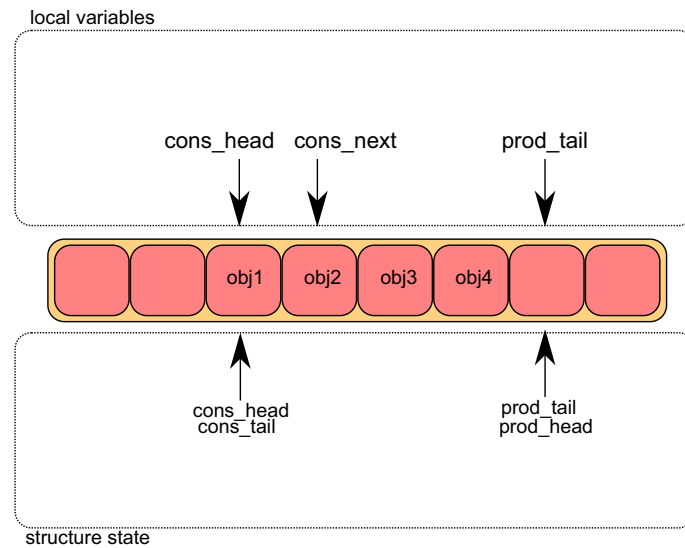
This section explains what occurs when a consumer dequeues an object from the ring. In this example, only the consumer head and tail (`cons_head` and `cons_tail`) are modified and there is only one consumer.

The initial state is to have a `cons_head` and `cons_tail` pointing at the same location.

5.5.2.1 Dequeue First Step

First, `ring->cons_head` and `ring->prod_tail` are copied in local variables. The `cons_next` local variable points to the next element of the table, or several elements after in the case of bulk dequeue.

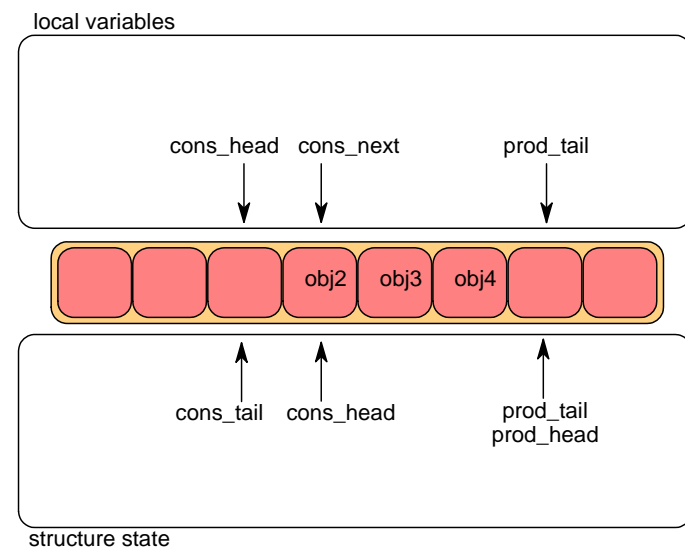
If there are not enough objects in the ring (this is detected by checking `prod_tail`), it returns an error.



5.5.2.2 Dequeue Second Step

The second step is to modify `ring->cons_head` in the ring structure to point to the same location as `cons_next`.

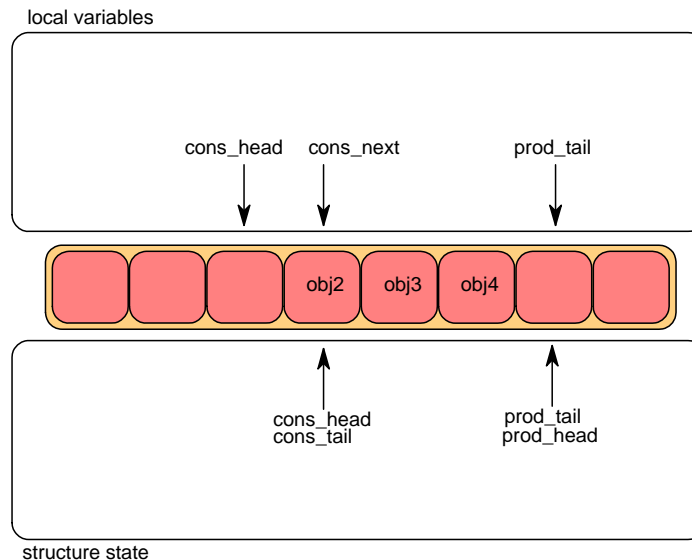
The pointer to the dequeued object (`obj1`) is copied in the pointer given by the user.





5.5.2.3 Dequeue Last Step

Finally, `ring->cons_tail` in the ring structure is modified to point to the same location as `ring->cons_head`. The dequeue operation is finished.



5.5.3 Multiple Producers Enqueue

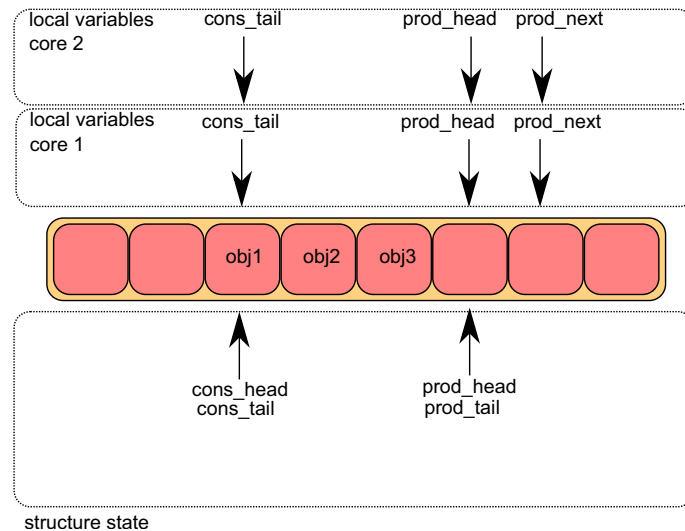
This section explains what occurs when two producers concurrently add an object to the ring. In this example, only the producer head and tail (`prod_head` and `prod_tail`) are modified.

The initial state is to have a `prod_head` and `prod_tail` pointing at the same location.

5.5.3.1 MC Enqueue First Step

On both cores, `ring->prod_head` and `ring->cons_tail` are copied in local variables. The `prod_next` local variable points to the next element of the table, or several elements after in the case of bulk enqueue.

If there are not enough objects in the ring (this is detected by checking `cons_tail`), it returns an error.

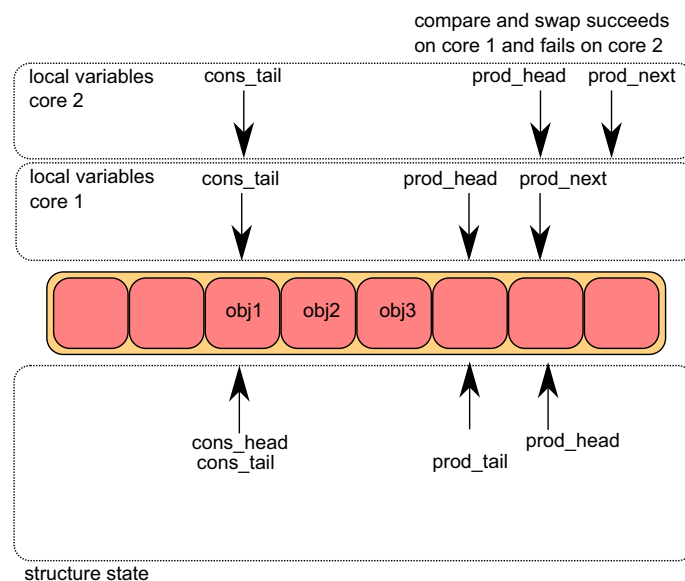


5.5.3.2 MC Enqueue Second Step

The second step is to modify `ring->prod_head` in the ring structure to point to the same location as `prod_next`. This operation is done using a Compare And Swap (CAS) instruction, which does the following operations atomically:

- If `ring->prod_head` is different to local variable `prod_head`, the CAS operation fails, and the code restarts at first step.
- Otherwise, `ring->prod_head` is set to local `prod_next`, the CAS operation is successful, and processing continues.

In the figure, the operation succeeded on core 1, and step one restarted on core 2.

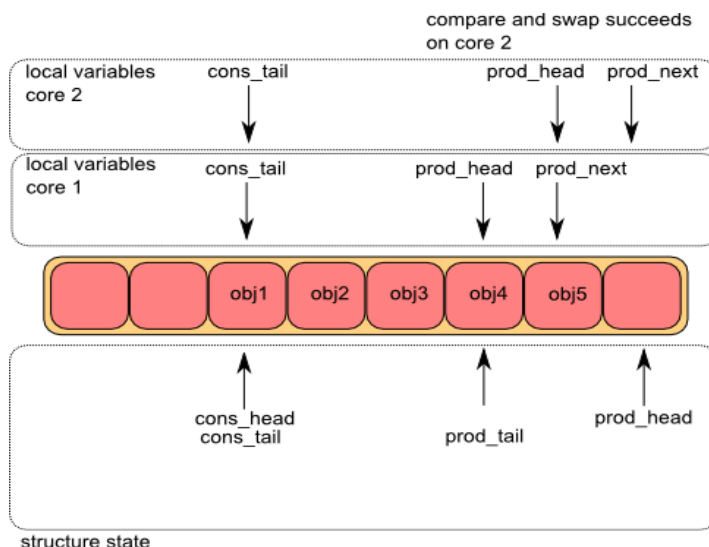




5.5.3.3 MC Enqueue Third Step

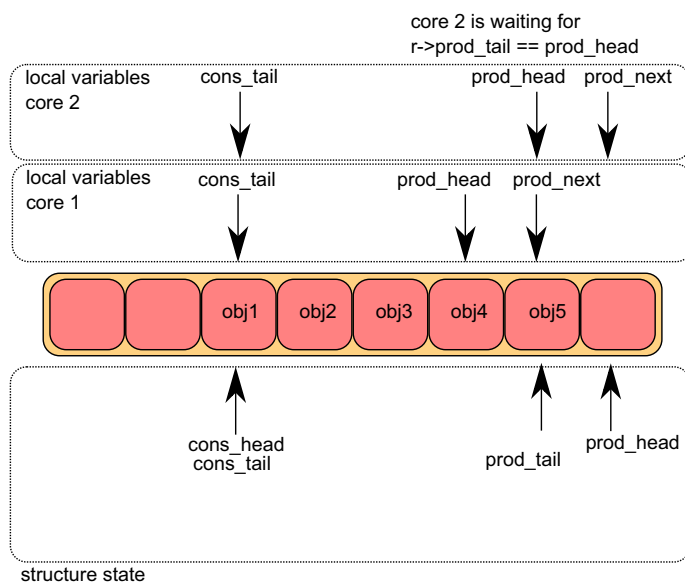
The CAS operation is retried on core 2 with success.

The core 1 updates one element of the ring (obj4), and the core 2 updates another one (obj5).



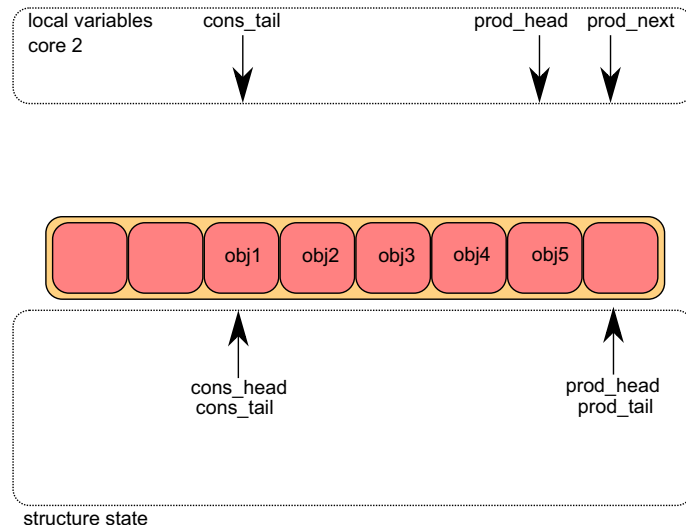
5.5.3.4 MC Enqueue Fourth Step

Each core now wants to update `ring->prod_tail`. A core can only update it if `ring->prod_tail` is equal to the `prod_head` local variable. This is only true on core 1. The operation is finished on core 1.



5.5.3.5 MC Enqueue Last Step

Once `ring->prod_tail` is updated by core 1, core 2 is allowed to update it too. The operation is also finished on core 2.



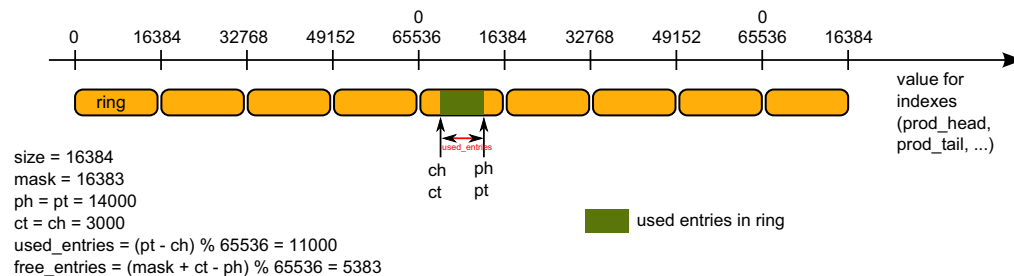
5.5.4 Modulo 32-bit Indexes

In the preceding figures, the `prod_head`, `prod_tail`, `cons_head` and `cons_tail` indexes are represented by arrows. In the actual implementation, these values are not between 0 and `size(ring) - 1` as would be assumed. The indexes are between 0 and $2^{32}-1$, and we mask their value when we access the pointer table (the ring itself). 32-bit modulo also implies that operations on indexes (such as, add/subtract) will automatically do 2^{32} modulo if the result overflows the 32-bit number range.

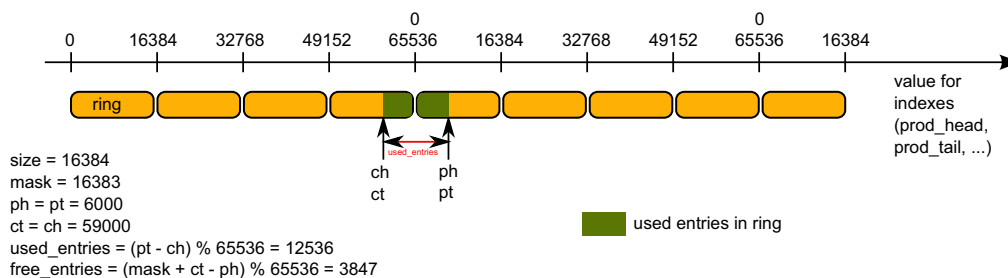
The following are two examples that help to explain how indexes are used in a ring.

Note:

To simplify the explanation, operations with modulo 16-bit are used instead of modulo 32-bit. In addition, the four indexes are defined as unsigned 16-bit integers, as opposed to unsigned 32-bit integers in the more realistic case.



This ring contains 11000 entries.



This ring contains 12536 entries.

Note: For ease of understanding, we use modulo 65536 operations in the above examples. In real execution cases, this is redundant for low efficiency, but is done automatically when the result overflows.

The code always maintains a distance between producer and consumer between 0 and `size(ring) - 1`. Thanks to this property, we can do subtractions between 2 index values in a modulo-32bit base: that's why the overflow of the indexes is not a problem.

At any time, `entries` and `free_entries` are between 0 and `size(ring) - 1`, even if only the first term of subtraction has overflowed:

```

uint32_t entries = (prod_tail - cons_head);
uint32_t free_entries = (mask + cons_tail - prod_head);

```

5.6 References

- [bufring.c in FreeBSD \(version 8\)](#)
- [bufring.h in FreeBSD \(version 8\)](#)
- [Linux Lockless Ring Buffer Design](#)

6.0 Mempool Library

A memory pool is an allocator of a fixed-sized object. In the Intel® DPDK, it is identified by name and uses a ring to store free objects. It provides some other optional services such as a per-core object cache and an alignment helper to ensure that objects are padded to spread them equally on all DRAM or DDR3 channels.

This library is used by the [Mbuf Library](#) and the [Environment Abstraction Layer](#) (for logging history).

6.1 Cookies

In debug mode (`CONFIG_RTE_LIBRTE_MEMPOOL_DEBUG` is enabled), cookies are added at the beginning and end of allocated blocks. The allocated objects then contain overwrite protection fields to help debugging buffer overflows.

6.2 Stats

In debug mode (`CONFIG_RTE_LIBRTE_MEMPOOL_DEBUG` is enabled), statistics about get from/put in the pool are stored in the mempool structure. Statistics are per-lcore to avoid concurrent access to statistics counters.

6.3 Memory Alignment Constraints

Depending on hardware memory configuration, performance can be greatly improved by adding a specific padding between objects. The objective is to ensure that the beginning of each object starts on a different channel and rank in memory so that all channels are equally loaded.

This is particularly true for packet buffers when doing L3 forwarding or flow classification. Only the first 64 bytes are accessed, so performance can be increased by spreading the start addresses of objects among the different channels.

The number of ranks on any DIMM is the number of independent sets of DRAMs that can be accessed for the full data bit-width of the DIMM. The ranks cannot be accessed simultaneously since they share the same data path. The physical layout of the DRAM chips on the DIMM itself does not necessarily relate to the number of ranks.

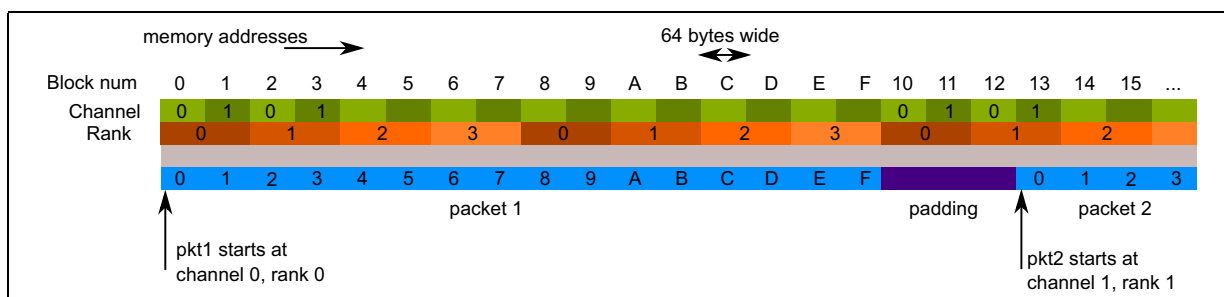
When running an application, the EAL command line options provide the ability to add the number of memory channels and ranks.

Note: The command line must always have the number of memory channels specified for the processor.

Examples of alignment for different DIMM architectures are shown in [Figure 4](#) and [Figure 5](#).



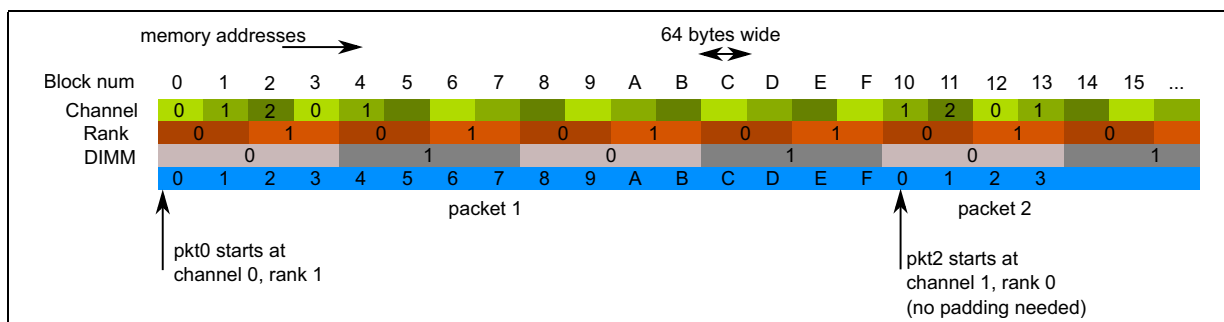
Figure 4. Two Channels and Quad-ranked DIMM Example



In this case, the assumption is that a packet is 16 blocks of 64 bytes, which is not true.

The Intel® 5520 chipset has three channels, so in most cases, no padding is required between objects (except for objects whose size are $n \times 3 \times 64$ bytes blocks).

Figure 5. Three Channels and Two Dual-ranked DIMM Example



When creating a new pool, the user can specify to use this feature or not.

6.4 Local Cache

In terms of CPU usage, the cost of multiple cores accessing a memory pool's ring of free buffers may be high since each access requires a compare-and-set (CAS) operation. To avoid having too many access requests to the memory pool's ring, the memory pool allocator can maintain a per-core cache and do bulk requests to the memory pool's ring, via the cache with many fewer locks on the actual memory pool structure. In this way, each core has full access to its own cache (with locks) of free objects and only when the cache fills does the core need to shuffle some of the free objects back to the pools ring or obtain more objects when the cache is empty.

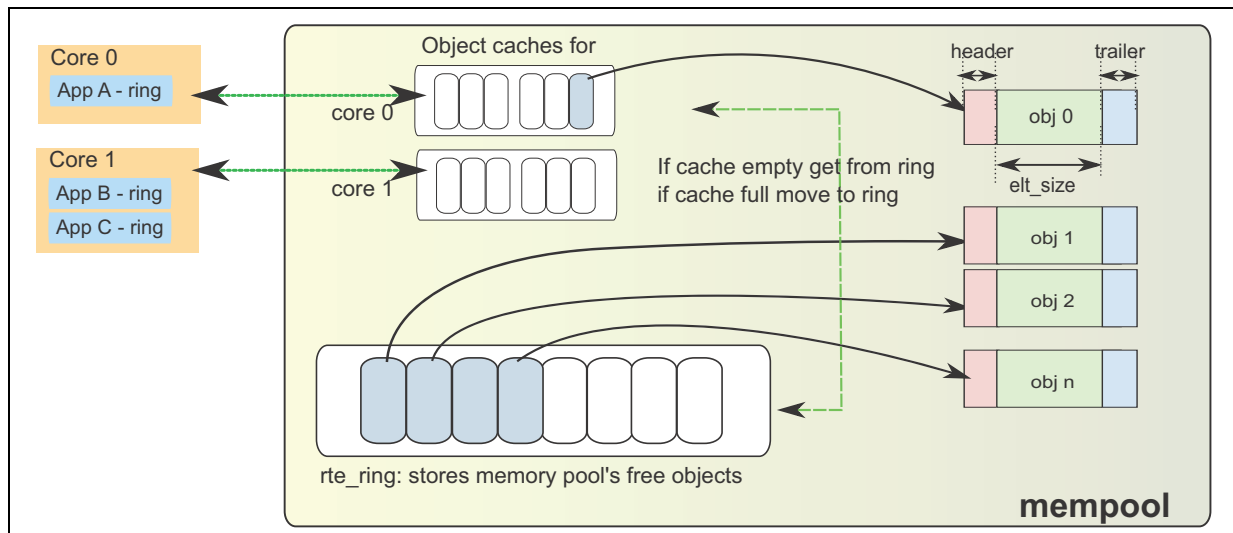
While this may mean a number of buffers may sit idle on some core's cache, the speed at which a core can access its own cache for a specific memory pool without locks provides performance gains.

The cache is composed of a small, per-core table of pointers and its length (used as a stack). This cache can be enabled or disabled at creation of the pool.

The maximum size of the cache is static and is defined at compilation time (CONFIG RTE_MEMPOOL_CACHE_MAX_SIZE).

Figure 6 shows a cache in operation.

Figure 6. A mempool in Memory with its Associated Ring



6.5 Use Cases

All allocations that require a high level of performance should use a pool-based memory allocator. Below are some examples:

- [Mbuf Library](#)
- [Environment Abstraction Layer](#), for logging service
- Any application that needs to allocate fixed-sized objects in the data plane and that will be continuously utilized by the system.

§ §



7.0 Mbuf Library

The mbuf library provides the ability to allocate and free buffers (*mbufs*) that may be used by the Intel® DPDK application to store message buffers. The message buffers are stored in a mempool, using the [Mempool Library](#).

A `rte_mbuf` struct can carry network packet buffers (type is `RTE_MBUF_PKT`) or generic control buffers (type is `RTE_MBUF_CTRL`). This can be extended to other types. The `rte_mbuf` is kept as small as possible (one cache line if possible).

7.1 Design of Packet Buffers

For the storage of the packet data (including protocol headers), two approaches were considered:

1. Embed metadata within a single memory buffer the structure followed by a fixed size area for the packet data.
2. Use separate memory buffers for the metadata structure and for the packet data.

The advantage of the first method is that it only needs one operation to allocate/free the whole memory representation of a packet. On the other hand, the second method is more flexible and allows the complete separation of the allocation of metadata structures from the allocation of packet data buffers.

The first method was chosen for the Intel® DPDK. The metadata contains control information such as message type, length, pointer to the start of the data and a pointer for additional mbuf structures allowing buffer chaining.

Message buffers that are used to carry network packets can handle buffer chaining where multiple buffers are required to hold the complete packet. This is the case for jumbo frames that are composed of many mbufs linked together through their `pkt.next` field.

For a newly allocated mbuf, the area at which the data begins in the message buffer is `RTE_PKTMBUF_HEADROOM` bytes after the beginning of the buffer, which is cache aligned. Message buffers may be used to carry control information, packets, events, and so on between different entities in the system. Message buffers may also use their data pointers to point to other message buffer data sections or other structures.

[Figure 7](#) and [Figure 8](#) show some of these scenarios.

Figure 7. An mbuf with One Segment

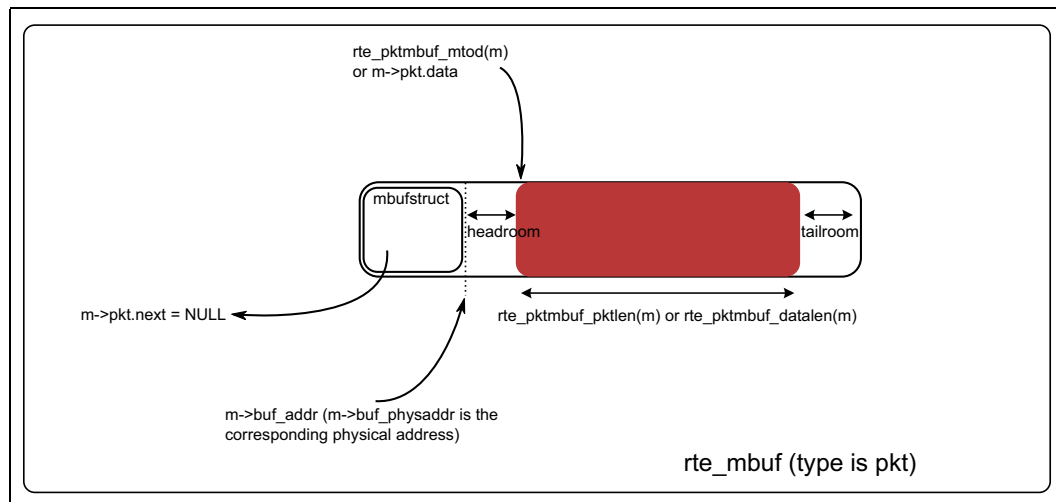
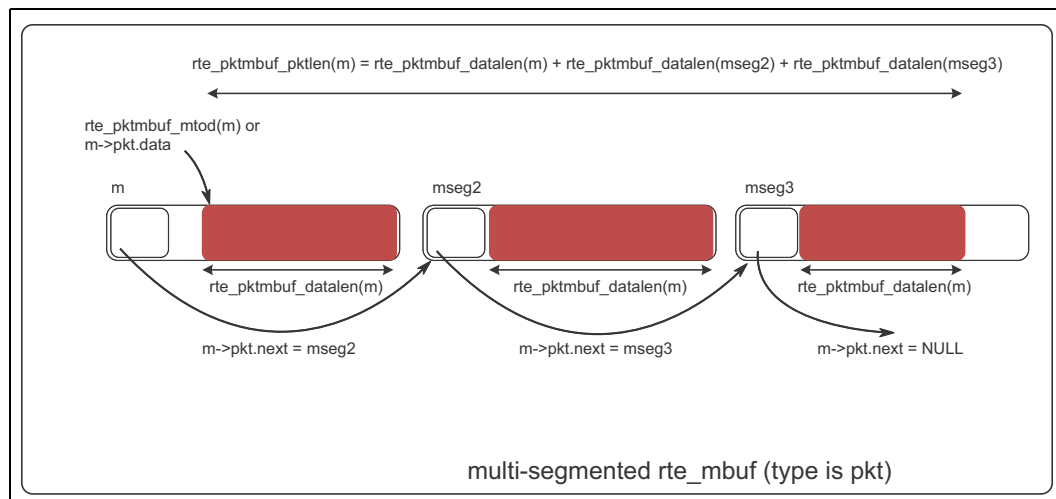


Figure 8. An mbuf with Three Segments



The Buffer Manager implements a fairly standard set of buffer access functions to manipulate network packets.

7.2 Buffers Stored in Memory Pools

The Buffer Manager uses the [Mempool Library](#) to allocate buffers. Therefore, it ensures that the packet header is interleaved optimally across the channels and ranks for L3 processing. An mbuf contains a field indicating the pool that it originated from. When calling `rte_ctrlmbuf_free(m)` or `rte_pktmbuf_free(m)`, the mbuf returns to its original pool.



7.3 Constructors

Packet and control mbuf constructors are provided by the API. The `rte_pktmbuf_init()` and `rte_ctrlmbuf_init()` functions initialize some fields in the mbuf structure that are not modified by the user once created (mbuf type, origin pool, buffer start address, and so on). This function is given as a callback function to the `rte_mempool_create()` function at pool creation time.

7.4 Allocating and Freeing mbufs

Allocating a new mbuf requires the user to specify the mempool from which the mbuf should be taken. For a packet mbuf, it contains one segment, with a length of 0. The pointer to data is initialized to have some bytes of headroom in the buffer (`RTE_PKTMBUF_HEADROOM`). For a control mbuf, it is initialized with data pointing to the beginning of the buffer and a length of zero.

Freeing a mbuf means returning it into its original mempool. The content of an mbuf is not modified when it is stored in a pool (as a free mbuf). Fields initialized by the constructor do not need to be re-initialized at mbuf allocation.

When freeing a packet mbuf that contains several segments, all of them are freed and returned to their original mempool.

7.5 Manipulating mbufs

This library provides some functions for manipulating the data in a packet mbuf. For instance:

- Get data length
- Get a pointer to the start of data
- Prepend data before data
- Append data after data
- Remove data at the beginning of the buffer (`rte_pktmbuf_adj()`)
- Remove data at the end of the buffer (`rte_pktmbuf_trim()`)

Refer to the *Intel® DPDK API Reference* for details.

7.6 Meta Information

Some information is retrieved by the network driver and stored in an mbuf to make processing easier. For instance, the VLAN, the RSS hash result (see [Poll Mode Driver](#)) and a flag indicating that the checksum was computed by hardware.

An mbuf also contains the input port (where it comes from), and the number of segment mbufs in the chain.

For chained buffers, only the first mbuf of the chain stores this meta information.

7.7 Direct and Indirect Buffers

A direct buffer is a buffer that is completely separate and self-contained. An indirect buffer behaves like a direct buffer but for the fact that the data pointer it contains points to data in another direct buffer. This is useful in situations where packets need to be duplicated or fragmented, since indirect buffers provide the means to reuse the same packet data across multiple buffers.

A buffer becomes indirect when it is “attached” to a direct buffer using the `rte_pktmbuf_attach()` function. Each buffer has a reference counter field and whenever an indirect buffer is attached to the direct buffer, the reference counter on the direct buffer is incremented. Similarly, whenever the indirect buffer is detached, the reference counter on the direct buffer is decremented. If the resulting reference counter is equal to 0, the direct buffer is freed since it is no longer in use.

There are a few things to remember when dealing with indirect buffers. First of all, it is not possible to attach an indirect buffer to another indirect buffer. Secondly, for a buffer to become indirect, its reference counter must be equal to 1, that is, it must not be already referenced by another indirect buffer. Finally, it is not possible to reattach an indirect buffer to the direct buffer (unless it is detached first).

While the attach/detach operations can be invoked directly using the recommended `rte_pktmbuf_attach()` and `rte_pktmbuf_detach()` functions, it is suggested to use the higher-level `rte_pktmbuf_clone()` function, which takes care of the correct initialization of an indirect buffer and can clone buffers with multiple segments.

Since indirect buffers are not supposed to actually hold any data, the memory pool for indirect buffers should be configured to indicate the reduced memory consumption. Examples of the initialization of a memory pool for indirect buffers (as well as use case examples for indirect buffers) can be found in several of the sample applications, for example, the IPv4 Multicast sample application.

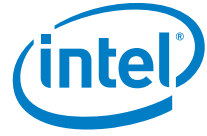
7.8 Debug

In debug mode (`CONFIG_RTE_MBUF_DEBUG` is enabled), the functions of the mbuf library perform sanity checks before any operation (such as, buffer corruption, bad type, and so on).

7.9 Use Cases

All networking application should use mbufs to transport network packets.

§ §



8.0 Poll Mode Driver

The Intel® DPDK includes 1 Gigabit and 10 Gigabit Poll Mode Drivers.

A Poll Mode Driver (PMD) consists of APIs, provided through the BSD driver running in user space, to configure the devices and their respective queues. In addition, a PMD accesses the RX and TX descriptors directly without any interrupts (with the exception of Link Status Change interrupts) to quickly receive, process and deliver packets in the user's application. This section describes the requirements of the PMDs, their global design principles and proposes a high-level architecture and a generic external API for the Ethernet PMDs.

8.1 Requirements and Assumptions

The Intel® DPDK environment for packet processing applications allows for two models, run-to-completion and pipe-line:

- In the *run-to-completion* model, a specific port's RX descriptor ring is polled for packets through an API. Packets are then processed on the same core and placed on a port's TX descriptor ring through an API for transmission.
- In the *pipe-line* model, one core polls one or more port's RX descriptor ring through an API. Packets are received and passed to another core via a ring. The other core continues to process the packet which then may be placed on a port's TX descriptor ring through an API for transmission.

In a synchronous run-to-completion model, each logical core assigned to the Intel® DPDK executes a packet processing loop that includes the following steps:

- Retrieve input packets through the PMD receive API
- Process each received packet one at a time, up to its forwarding
- Send pending output packets through the PMD transmit API

Conversely, in an asynchronous pipe-line model, some logical cores may be dedicated to the retrieval of received packets and other logical cores to the processing of previously received packets. Received packets are exchanged between logical cores through rings. The loop for packet retrieval includes the following steps:

- Retrieve input packets through the PMD receive API
- Provide received packets to processing lcores through packet queues

The loop for packet processing includes the following steps:

- Retrieve the received packet from the packet queue
- Process the received packet, up to its retransmission if forwarded

To avoid any unnecessary interrupt processing overhead, the execution environment must not use any asynchronous notification mechanisms. Whenever needed and appropriate, asynchronous communication should be introduced as much as possible through the use of rings.



Avoiding lock contention is a key issue in a multi-core environment. To address this issue, PMDs are designed to work with per-core private resources as much as possible. For example, a PMD maintains a separate transmit queue per-core, per-port. In the same way, every receive queue of a port is assigned to and polled by a single logical core (lcore).

To comply with Non-Uniform Memory Access (NUMA), memory management is designed to assign to each logical core a private buffer pool in local memory to minimize remote memory access. The configuration of packet buffer pools should take into account the underlying physical memory architecture in terms of DIMMS, channels and ranks. The application must ensure that appropriate parameters are given at memory pool creation time. See [Mempool Library](#).

8.2 Design Principles

The API and architecture of the Ethernet* PMDs are designed with the following guidelines in mind.

PMDs must help global policy-oriented decisions to be enforced at the upper application level. Conversely, NIC PMD functions should not impede the benefits expected by upper-level global policies, or worse prevent such policies from being applied.

For instance, both the receive and transmit functions of a PMD have a maximum number of packets/descriptors to poll. This allows a run-to-completion processing stack to statically fix or to dynamically adapt its overall behavior through different global loop policies, such as:

- Receive, process immediately and transmit packets one at a time in a piecemeal fashion.
- Receive as many packets as possible, then process all received packets, transmitting them immediately.
- Receive a given maximum number of packets, process the received packets, accumulate them and finally send all accumulated packets to transmit.

To achieve optimal performance, overall software design choices and pure software optimization techniques must be considered and balanced against available low-level hardware-based optimization features (CPU cache properties, bus speed, NIC PCI bandwidth, and so on). The case of packet transmission is an example of this software/hardware tradeoff issue when optimizing burst-oriented network packet processing engines. In the initial case, the PMD could export only an `rte_eth_tx_one` function to transmit one packet at a time on a given queue. On top of that, one can easily build an `rte_eth_tx_burst` function that loops invoking the `rte_eth_tx_one` function to transmit several packets at a time. However, an `rte_eth_tx_burst` function is effectively implemented by the PMD to minimize the driver-level transmit cost per packet through the following optimizations:

- Share among multiple packets the un-amortized cost of invoking the `rte_eth_tx_one` function.
- Enable the `rte_eth_tx_burst` function to take advantage of burst-oriented hardware features (prefetch data in cache, use of NIC head/tail registers) to minimize the number of CPU cycles per packet, for example by avoiding unnecessary read memory accesses to ring transmit descriptors, or by systematically using arrays of pointers that exactly fit cache line boundaries and sizes.
- Apply burst-oriented software optimization techniques to remove operations that would otherwise be unavoidable, such as ring index wrap back management.



Burst-oriented functions are also introduced via the API for services that are intensively used by the PMD. This applies in particular to buffer allocators used to populate NIC rings, which provide functions to allocate/free several buffers at a time. For example, an `mbuf_multiple_alloc` function returning an array of pointers to `rte_mbuf` buffers which speeds up the receive poll function of the PMD when replenishing multiple descriptors of the receive ring.

8.3 Logical Cores, Memory and NIC Queues Relationships

The Intel® DPDK supports NUMA allowing for better performance when a processor's logical cores and interfaces utilize its local memory. Therefore, mbuf allocation associated with local PCIe* interfaces should be allocated from memory pools created in the local memory. The buffers should, if possible, remain on the local processor to obtain the best performance results and RX and TX buffer descriptors should be populated with mbufs allocated from a mempool allocated from local memory.

The run-to-completion model also performs better if packet or data manipulation is in local memory instead of a remote processors memory. This is also true for the pipe-line model provided all logical cores used are located on the same processor.

Multiple logical cores should never share receive or transmit queues for interfaces since this would require global locks and hinder performance.

8.4 Device Identification and Configuration

8.4.1 Device Identification

Each NIC port is uniquely designated by its (bus/bridge, device, function) PCI identifiers assigned by the PCI probing/enumeration function executed at Intel® DPDK initialization. Based on their PCI identifier, NIC ports are assigned two other identifiers:

- A port index used to designate the NIC port in all functions exported by the PMD API.
- A port name used to designate the port in console messages, for administration or debugging purposes. For ease of use, the port name includes the port index.

8.4.2 Device Configuration

The configuration of each NIC port includes the following operations:

- Allocate PCI resources
- Reset the hardware (issue a Global Reset) to a well-known default state
- Set up the PHY and the link
- Initialize statistics counters

The PMD API must also export functions to start/stop the all-multicast feature of a port and functions to set/unset the port in promiscuous mode.

Some hardware offload features must be individually configured at port initialization through specific configuration parameters. This is the case for the Receive Side Scaling (RSS) and Data Center Bridging (DCB) features for example.

8.4.3 On-the-Fly Configuration

All device features that can be started or stopped “on the fly” (that is, without stopping the device) do not require the PMD API to export dedicated functions for this purpose.



All that is required is the mapping address of the device PCI registers to implement the configuration of these features in specific functions outside of the drivers.

For this purpose, the PMD API exports a function that provides all the information associated with a device that can be used to set up a given device feature outside of the driver. This includes the PCI vendor identifier, the PCI device identifier, the mapping address of the PCI device registers, and the name of the driver.

The main advantage of this approach is that it gives complete freedom on the choice of the API used to configure, to start, and to stop such features.

As an example, refer to the configuration of the IEEE1588 feature for the Intel® 82576 Gigabit Ethernet Controller and the Intel® 82599 10 Gigabit Ethernet Controller controllers in the `testpmd` application.

Other features such as the L3/L4 5-Tuple packet filtering feature of a port can be configured in the same way. Ethernet* flow control (pause frame) can be configured on the individual port. Refer to the `testpmd` source code for details. Also, L4 (UDP/TCP/SCTP) checksum offload by the NIC can be enabled for an individual packet as long as the packet mbuf is set up correctly. Refer to the `testpmd` source code (specifically the `csumonly.c` file) for details.

That being said, the support of some offload features implies the addition of dedicated status bit(s) and value field(s) into the `rte_mbuf` data structure, along with their appropriate handling by the receive/transmit functions exported by each PMD.

For instance, this is the case for the IEEE1588 packet timestamp mechanism, the VLAN tagging and the IP checksum computation, as described in the [Section 7.6, “Meta Information” on page 35](#).

8.4.4 Configuration of Transmit and Receive Queues

Each transmit queue is independently configured with the following information:

- The number of descriptors of the transmit ring
- The socket identifier used to identify the appropriate DMA memory zone from which to allocate the transmit ring in NUMA architectures
- The values of the Prefetch, Host and Write-Back threshold registers of the transmit queue
- The *minimum transmit packets to free* threshold (`tx_free_thresh`). When the number of descriptors used to transmit packets exceeds this threshold, the network adaptor should be checked to see if it has written back descriptors. A value of 0 can be passed during the TX queue configuration to indicate the default value should be used. The default value for `tx_free_thresh` is 32. This ensures that the PMD does not search for completed descriptors until at least 32 have been processed by the NIC for this queue.
- The *minimum RS bit* threshold. The minimum number of transmit descriptors to use before setting the Report Status (RS) bit in the transmit descriptor. Note that this parameter may only be valid for Intel 10 GbE network adapters. The RS bit is set on the last descriptor used to transmit a packet if the number of descriptors used since the last RS bit setting, up to the first descriptor used to transmit the packet, exceeds the transmit RS bit threshold (`tx_rs_thresh`). In short, this parameter controls which transmit descriptors are written back to host memory by the network adapter. A value of 0 can be passed during the TX queue configuration to indicate that the default value should be used. The default value for `tx_rs_thresh` is 32. This ensures that at least 32 descriptors are used before the network adapter writes back the most recently used descriptor. This saves upstream PCIe* bandwidth resulting from TX descriptor write-backs. It is important to note that the TX Write-back threshold (`TX_wthresh`) should be set to 0 when



`tx_rs_thresh` is greater than 1. Refer to the [Intel® 82599 10 Gigabit Ethernet Controller Datasheet](#) for more details.

The following constraints must be satisfied for `tx_free_thresh` and `tx_rs_thresh`:

- `tx_rs_thresh` must be greater than 0.
- `tx_rs_thresh` must be less than the size of the ring minus 2.
- `tx_rs_thresh` must be less than or equal to `tx_free_thresh`.
- `tx_free_thresh` must be greater than 0.
- `tx_free_thresh` must be less than the size of the ring minus 3.
- For optimal performance, TX `wthresh` should be set to 0 when `tx_rs_thresh` is greater than 1.

One descriptor in the TX ring is used as a sentinel to avoid a hardware race condition, hence the maximum threshold constraints.

Note: When configuring for DCB operation, at port initialization, both the number of transmit queues and the number of receive queues must be set to 128.

8.5 Poll Mode Driver API

8.5.1 Generalities

By default, all functions exported by a PMD are lock-free functions that are assumed not to be invoked in parallel on different logical cores to work on the same target object. For instance, a PMD receive function cannot be invoked in parallel on two logical cores to poll the same RX queue of the same port. Of course, this function can be invoked in parallel by different logical cores on different RX queues. It is the responsibility of the upper-level application to enforce this rule.

If needed, parallel accesses by multiple logical cores to shared queues can be explicitly protected by dedicated inline lock-aware functions built on top of their corresponding lock-free functions of the PMD API.

8.5.2 Generic Packet Representation

A packet is represented by an `rte_mbuf` structure, which is a generic metadata structure containing all necessary housekeeping information. This includes fields and status bits corresponding to offload hardware features, such as checksum computation of IP headers or VLAN tags.

The `rte_mbuf` data structure includes specific fields to represent, in a generic way, the offload features provided by network controllers. For an input packet, most fields of the `rte_mbuf` structure are filled in by the PMD receive function with the information contained in the receive descriptor. Conversely, for output packets, most fields of `rte_mbuf` structures are used by the PMD transmit function to initialize transmit descriptors.

The mbuf structure is fully described in the [Mbuf Library](#) chapter.

8.5.3 Ethernet Device API

The Ethernet device API exported by the Ethernet PMDs is described in the [Intel® DPDK API Reference](#).

§ §

9.0 Timer Library

The Timer library provides a timer service to Intel® DPDK execution units to enable execution of callback functions asynchronously. Features of the library are:

- Timers can be periodic (multi-shot) or single (one-shot).
- Timers can be loaded from one core and executed on another. It has to be specified in the call to `rte_timer_reset()`.
- Timers provide high precision (depends on the call frequency to `rte_timer_manage()` that checks timer expiration for the local core).
- If not required in the application, timers can be disabled at compilation time by not calling the `rte_timer_manage()` to increase performance.

The timer library uses the `rte_get_hpet_cycles()` function that uses the High Precision Event Timer (HPET) to provide a reliable time reference.

This library provides an interface to add, delete and restart a timer. The API is based on BSD `callout()` with a few differences. Refer to the [callout manual](#).

9.1 Implementation Details

Each lcore contains three lists: a list of pending timers, a list of expired timers, and a list of done timers. The roles of these lists are explained in the following figure.

These lists are accessed when resetting or stopping a timer and in the `rte_timer_manage()` function, which runs the expired timers. To read or update these lists, the core must take a per-lcore lock to avoid concurrent modifications.

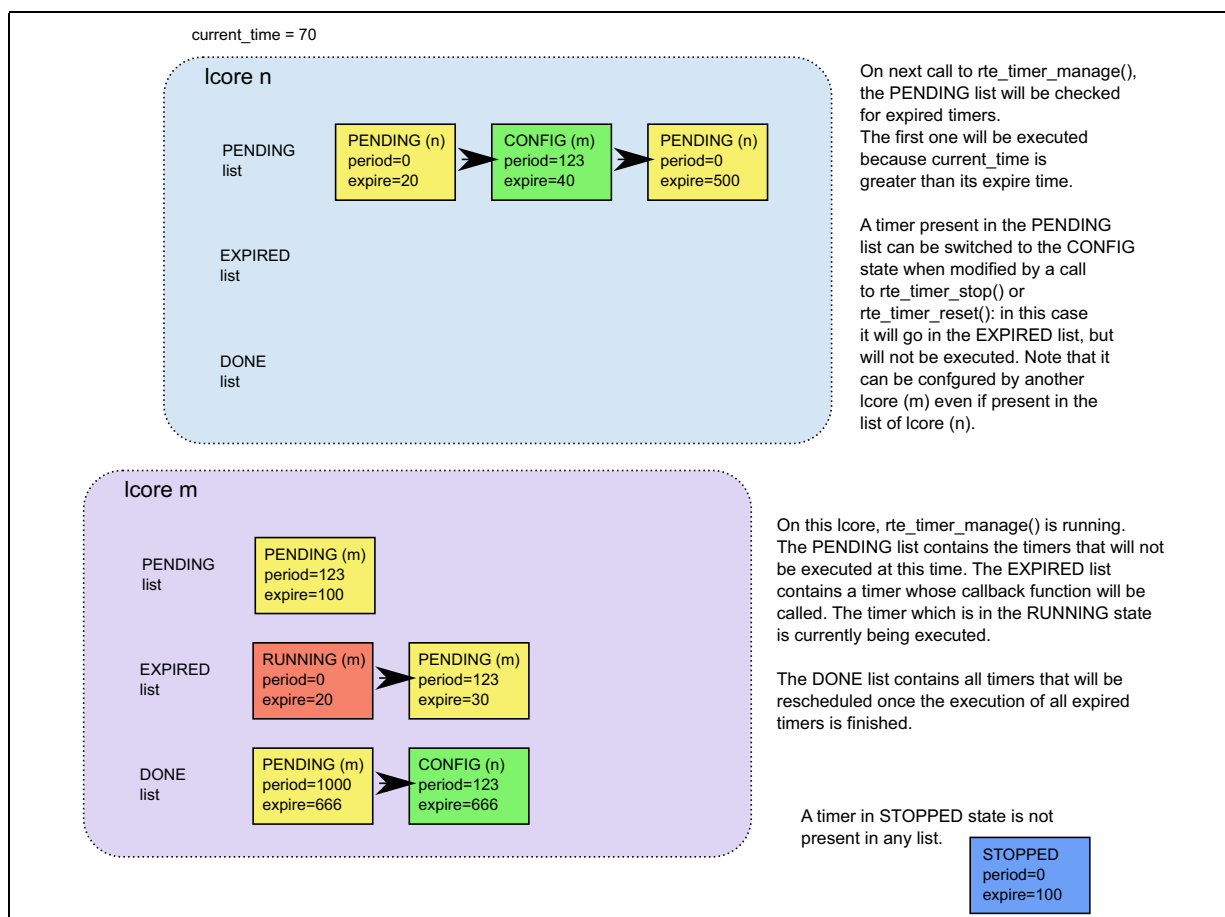
A timer structure contains a special field called `status`, which is a union of a timer state (stopped, pending, running, config) and an owner (lcore id). Depending on the timer state, we know if a timer is present in a list or not:

- STOPPED: no owner, not in a list
- CONFIG: owned by a core, must not be modified by another core, maybe in a list or not, depending on previous state
- PENDING: owned by a core, present in a list (pending, expired, or done list)
- RUNNING: owned by a core, must not be modified by another core, present in a list (pending, expired, or done list)

Resetting or stopping a timer while it is in a CONFIG or RUNNING state is not allowed. When modifying the state of a timer, a Compare and Swap instruction should be used to guarantee that the status (state+owner) is modified atomically.

[Figure 9](#) shows an example of a timer operation.

Figure 9. Timer Operation



9.2 Use Cases

The timer library is used for periodic calls, such as garbage collectors, or some state machines (ARP, bridging, and so on).

9.3 References

- [callout manual](#) - The callout facility that provides timers with a mechanism to execute a function at a given time.
- [HPET](#) - Information about the High Precision Event Timer (HPET).

§ §

10.0 Hash Library

The Intel® DPDK provides a Hash Library for creating hash table for fast lookup. The hash table is a data structure optimized for searching through a set of entries that are each identified by a unique key. For increased performance the Intel® DPDK Hash requires that all the keys have the same number of bytes which is set at the hash creation time.

10.1 Hash API Overview

The main configuration parameters for the hash are:

- Total number of hash entries
- Size of the key in bytes

The hash also allows the configuration of some low-level implementation related parameters such as:

- Hash function to translate the key into a bucket index
- Number of entries per bucket

The main methods exported by the hash are:

- Add entry with key: The key is provided as input. If a new entry is successfully added to the hash for the specified key, or there is already an entry in the hash for the specified key, then the position of the entry is returned. If the operation was not successful, for example due to lack of free entries in the hash, then a negative value is returned;
- Delete entry with key: The key is provided as input. If an entry with the specified key is found in the hash, then the entry is removed from the hash and the position where the entry was found in the hash is returned. If no entry with the specified key exists in the hash, then a negative value is returned
- Lookup for entry with key: The key is provided as input. If an entry with the specified key is found in the hash (lookup hit), then the position of the entry is returned, otherwise (lookup miss) a negative value is returned.

The current hash implementation handles the key management only. The actual data associated with each key has to be managed by the user using a separate table that mirrors the hash in terms of number of entries and position of each entry, as shown in the Flow Classification use case describes in the following sections.

The example hash tables in the L2/L3 Forwarding sample applications defines which port to forward a packet to based on a packet flow identified by the five-tuple lookup. However, this table could also be used for more sophisticated features and provide many other functions and actions that could be performed on the packets and flows.



10.2 Implementation Details

The hash table is implemented as an array of entries which is further divided into buckets, with the same number of consecutive array entries in each bucket. For any input key, there is always a single bucket where that key can be stored in the hash, therefore only the entries within that bucket need to be examined when the key is looked up. The lookup speed is achieved by reducing the number of entries to be scanned from the total number of hash entries down to the number of entries in a hash bucket, as opposed to the basic method of linearly scanning all the entries in the array. The hash uses a hash function (configurable) to translate the input key into a 4-byte key signature. The bucket index is the key signature modulo the number of hash buckets. Once the bucket is identified, the scope of the hash add, delete and lookup operations is reduced to the entries in that bucket.

To speed up the search logic within the bucket, each hash entry stores the 4-byte key signature together with the full key for each hash entry. For large key sizes, comparing the input key against a key from the bucket can take significantly more time than comparing the 4-byte signature of the input key against the signature of a key from the bucket. Therefore, the signature comparison is done first and the full key comparison done only when the signatures matches. The full key comparison is still necessary, as two input keys from the same bucket can still potentially have the same 4-byte hash signature, although this event is relatively rare for hash functions providing good uniform distributions for the set of input keys.

10.3 Use Case: Flow Classification

Flow classification is used to map each input packet to the connection/flow it belongs to. This operation is necessary as the processing of each input packet is usually done in the context of their connection, so the same set of operations is applied to all the packets from the same flow.

Applications using flow classification typically have a flow table to manage, with each separate flow having an entry associated with it in this table. The size of the flow table entry is application specific, with typical values of 4, 16, 32 or 64 bytes.

Each application using flow classification typically has a mechanism defined to uniquely identify a flow based on a number of fields read from the input packet that make up the flow key. One example is to use the DiffServ 5-tuple made up of the following fields of the IP and transport layer packet headers: Source IP Address, Destination IP Address, Protocol, Source Port, Destination Port.

The Intel® DPDK hash provides a generic method to implement an application specific flow classification mechanism. Given a flow table implemented as an array, the application should create a hash object with the same number of entries as the flow table and with the hash key size set to the number of bytes in the selected flow key.

The flow table operations on the application side are described below:

- **Add flow:** Add the flow key to hash. If the returned position is valid, use it to access the flow entry in the flow table for adding a new flow or updating the information associated with an existing flow. Otherwise, the flow addition failed, for example due to lack of free entries for storing new flows.
- **Delete flow:** Delete the flow key from the hash. If the returned position is valid, use it to access the flow entry in the flow table to invalidate the information associated with the flow.
- **Lookup flow:** Lookup for the flow key in the hash. If the returned position is valid (flow lookup hit), use the returned position to access the flow entry in the flow table. Otherwise (flow lookup miss) there is no flow registered for the current packet.



10.4 References

- Donald E. Knuth, *The Art of Computer Programming, Volume 3: Sorting and Searching (2nd Edition)*, 1998, Addison-Wesley Professional

§ §



11.0 LPM Library

The Intel® DPDK LPM library component implements the Longest Prefix Match (LPM) table search method for 32-bit keys that is typically used to find the best route match in IP forwarding applications.

11.1 LPM API Overview

The main configuration parameter for LPM component instances is the maximum number of rules to support. An LPM prefix is represented by a pair of parameters (32-bit key, depth), with depth in the range of 1 to 32. An LPM rule is represented by an LPM prefix and some user data associated with the prefix. The prefix serves as the unique identifier of the LPM rule. In this implementation, the user data is 1-byte long and is called next hop, in correlation with its main use of storing the ID of the next hop in a routing table entry.

The main methods exported by the LPM component are:

- Add LPM rule: The LPM rule is provided as input. If there is no rule with the same prefix present in the table, then the new rule is added to the LPM table. If a rule with the same prefix is already present in the table, the next hop of the rule is updated. An error is returned when there is no available rule space left.
- Delete LPM rule: The prefix of the LPM rule is provided as input. If a rule with the specified prefix is present in the LPM table, then it is removed.
- Lookup LPM key: The 32-bit key is provided as input. The algorithm selects the rule that represents the best match for the given key and returns the next hop of that rule. In the case that there are multiple rules present in the LPM table that have the same 32-bit key, the algorithm picks the rule with the highest depth as the best match rule, which means that the rule has the highest number of most significant bits matching between the input key and the rule key.

11.2 Implementation Details

The current implementation uses a variation of the DIR-24-8 algorithm that trades memory usage for improved LPM lookup speed. The algorithm allows the lookup operation to be performed with typically a single memory read access. In the statistically rare case when the best match rule has a depth bigger than 24, the lookup operation requires two memory read accesses. Therefore, the performance of the LPM lookup operation is greatly influenced by whether the specific memory location is present in the processor cache or not.

11.3 Use Case: IPv4 Forwarding

The LPM algorithm is used to implement Classless Inter-Domain Routing (CIDR) strategy used by routers implementing IPv4 forwarding.



11.4 References

- *RFC1519 Classless Inter-Domain Routing (CIDR): an Address Assignment and Aggregation Strategy*, <http://www.ietf.org/rfc/rfc1519>
- Pankaj Gupta, *Algorithms for Routing Lookups and Packet Classification*, PhD Thesis, Stanford University, 2000 (http://klamath.stanford.edu/~pankaj/thesis/thesis_1sided.pdf)

§ §

§ §



12.0 Multi-process Support

In the Intel® DPDK, multi-process support is designed to allow a group of Intel® DPDK processes to work together in a simple transparent manner to perform packet processing, or other workloads, on Intel® architecture hardware. To support this functionality, a number of additions have been made to the core Intel® DPDK Environment Abstraction Layer (EAL).

The EAL has been modified to allow different types of Intel® DPDK processes to be spawned, each with different permissions on the hugepage memory used by the applications. For now, there are two types of process specified:

- primary processes, which can initialize and which have full permissions on shared memory
- secondary processes, which cannot initialize shared memory, but can attach to pre-initialized shared memory and create objects in it.

Standalone Intel® DPDK processes are primary processes, while secondary processes can only run alongside a primary process or after a primary process has already configured the hugepage shared memory for them.

To support these two process types, and other multi-process setups described later, two additional command-line parameters are available to the EAL:

- `--proc-type`: for specifying a given process instance as the primary or secondary Intel® DPDK instance
- `--file-prefix`: to allow processes that do not want to co-operate to have different memory regions

A number of example applications are provided that demonstrate how multiple Intel® DPDK processes can be used together. These are more fully documented in the “Multi-process Sample Application” chapter in the *Intel® DPDK Sample Application's User Guide*.

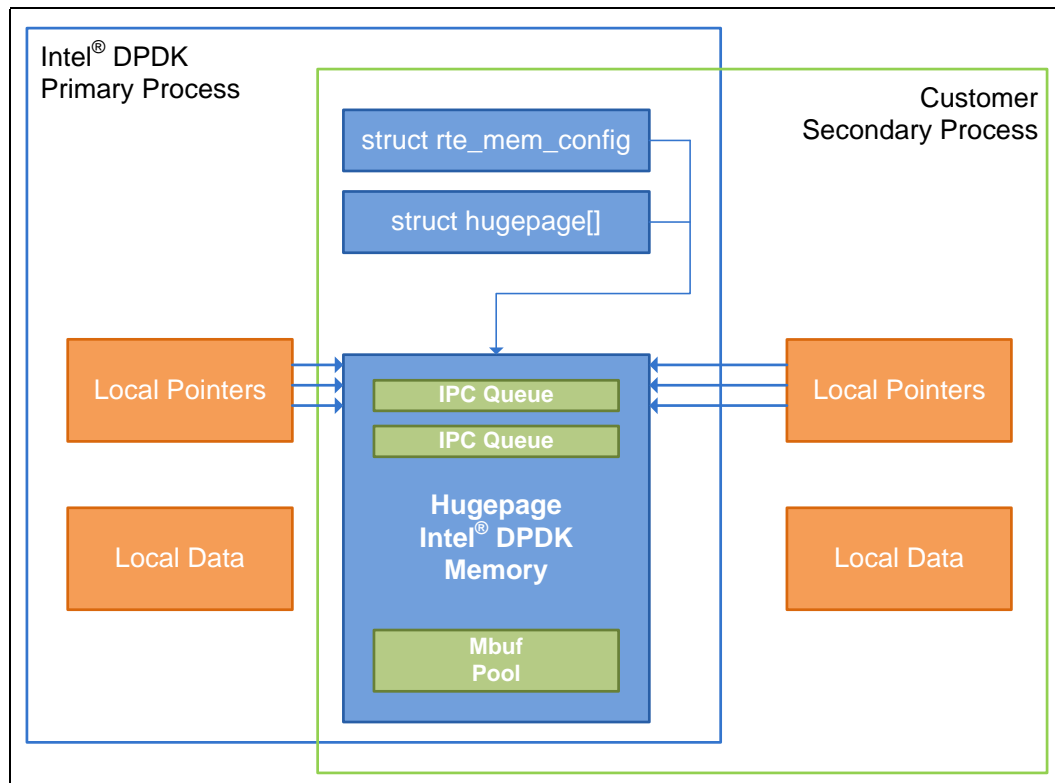
12.1 Memory Sharing

The key element in getting a multi-process application working using the Intel® DPDK is to ensure that memory resources are properly shared among the processes making up the multi-process application. Once there are blocks of shared memory available that can be accessed by multiple processes, then issues such as inter-process communication (IPC) becomes much simpler.

On application start-up in a primary or standalone process, the Intel DPDK records to memory-mapped files the details of the memory configuration it is using - hugepages in use, the virtual addresses they are mapped at, the number of memory channels present, etc. When a secondary process is started, these files are read and the EAL recreates the same memory configuration in the secondary process so that all memory zones are shared between processes and all pointers to that memory are valid, and point to the same objects, in both processes.

Note: Refer to the [Section 12.3, “Multi-process Limitations” on page 51](#) for details of how Linux kernel Address-Space Layout Randomization (ASLR) can affect memory sharing.

Figure 10. Memory Sharing in the Intel® DPDK Multi-process Sample Application



The EAL also supports an auto-detection mode (set by EAL `--proc-type=auto` flag), whereby an Intel® DPDK process is started as a secondary instance if a primary instance is already running.

12.2 Deployment Models

12.2.1 Symmetric/Peer Processes

Intel® DPDK multi-process support can be used to create a set of peer processes where each process performs the same workload. This model is equivalent to having multiple threads each running the same main-loop function, as is done in most of the supplied Intel® DPDK sample applications. In this model, the first of the processes spawned should be spawned using the `--proc-type=primary` EAL flag, while all subsequent instances should be spawned using the `--proc-type=secondary` flag.

The `simple_mp` and `symmetric_mp` sample applications demonstrate this usage model. They are described in the “Multi-process Sample Application” chapter in the *Intel® DPDK Sample Application's User Guide*.

12.2.2 Asymmetric/Non-Peer Processes

An alternative deployment model that can be used for multi-process applications is to have a single primary process instance that acts as a load-balancer or server distributing received packets among worker or client threads, which are run as secondary processes. In this case, extensive use of `rte_ring` objects is made, which are located in shared hugepage memory.



The `client_server_mp` sample application shows this usage model. It is described in the “Multi-process Sample Application” chapter in the *Intel® DPDK Sample Application's User Guide*.

12.2.3 Running Multiple Independent Intel® DPDK Applications

In addition to the above scenarios involving multiple Intel® DPDK processes working together, it is possible to run multiple Intel® DPDK processes side-by-side, where those processes are all working independently. Support for this usage scenario is provided using the `--file-prefix` parameter to the EAL.

By default, the EAL creates hugepage files on each `hugetlbfs` filesystem using the `rtemap_X` filename, where `X` is in the range 0 to the maximum number of hugepages -1. Similarly, it creates shared configuration files, memory mapped in each process, using the `/var/run/.rte_config` filename, when run as root (or `$HOME/.rte_config` when run as a non-root user; if filesystem and device permissions are set up to allow this). The `rte` part of the filenames of each of the above is configurable using the `file-prefix` parameter.

In addition to specifying the `file-prefix` parameter, any Intel® DPDK applications that are to be run side-by-side must explicitly limit their memory use. This is done by passing the `-m` flag to each process to specify how much hugepage memory, in megabytes, each process can use (or passing `--socket-mem` to specify how much hugepage memory on each socket each process can use).

Note: Independent Intel® DPDK instances running side-by-side on a single machine cannot share any network ports. Any network ports being used by one process should be blacklisted in every other process.

12.2.4 Running Multiple Independent Groups of Intel® DPDK Applications

In the same way that it is possible to run independent Intel® DPDK applications side-by-side on a single system, this can be trivially extended to multi-process groups of Intel® DPDK applications running side-by-side. In this case, the secondary processes must use the same `--file-prefix` parameter as the primary process whose shared memory they are connecting to.

Note: All restrictions and issues with multiple independent Intel® DPDK processes running side-by-side apply in this usage scenario also.

12.3 Multi-process Limitations

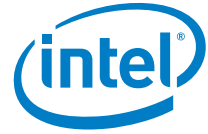
There are a number of limitations to what can be done when running Intel® DPDK multi-process applications. Some of these are documented below:

- The multi-process feature requires that the exact same hugepage memory mappings be present in all applications. The Linux security feature - Address-Space Layout Randomization (ASLR) can interfere with this mapping, so it may be necessary to disable this feature in order to reliably run multi-process applications.
Note: Disabling Address-Space Layout Randomization (ASLR) may have security implications, so it is recommended that it be disabled only when absolutely necessary, and only when the implications of this change have been understood.
- All Intel® DPDK processes running as a single application and using shared memory must have distinct `coremask` arguments. It is not possible to have a primary and secondary instance, or two secondary instances, using any of the same logical cores. Attempting to do so can cause corruption of memory pool caches, among other issues.



- Since the Intel® DPDK is now thread-safe, memory (memzone) allocation/reservation from hugepage memory is no longer limited to the primary process as in earlier releases of the Intel® DPDK. In the primary process, all that is necessary is the initialization of hugepage memory early during EAL startup.
- The delivery of interrupts, such as Ethernet* device link status interrupts, do not work in secondary processes. All interrupts are triggered inside the primary process only. Any application needing interrupt notification in multiple processes should provide its own mechanism to transfer the interrupt information from the primary process to any secondary process that needs the information.
- The use of function pointers between multiple processes running based on different compiled binaries is not supported, since the location of a given function in one process may be different to its location in a second. This prevents the `librte_hash` library from behaving properly as in a multi-threaded instance, since it uses a pointer to the hash function internally.
- Depending upon the hardware in use, and the number of Intel® DPDK processes used, it may not be possible to have HPET timers available in each Intel® DPDK instance. The minimum number of HPET comparators available to Linux* userspace can be just a single comparator, which means that only the first, primary Intel® DPDK process instance can open and `mmap /dev/hpet`. All other Intel® DPDK process instances use the chip TSC as a fallback if the HPET cannot be used.

§ §



13.0 IXGBE/IGB Virtual Function Driver

Supported Intel® Ethernet Controllers (see the *Intel® DPDK Release Notes* for details) support the following modes of operation in a virtualized environment:

- **SR-IOV mode:** Involves direct assignment of part of the port resources to different guest operating systems using the PCI-SIG Single Root I/O Virtualization (SR IOV) standard, also known as “native mode” or “pass-through” mode. In this chapter, this mode is referred to as IOV mode.
- **VMDq mode:** Involves central management of the networking resources by an IO Virtual Machine (IOVM) or a Virtual Machine Monitor (VMM), also known as “software switch acceleration” mode. In this chapter, this mode is referred to as the Next Generation VMDq mode.

13.1 SR-IOV Mode Utilization in an Intel® DPDK Environment

The Intel® DPDK uses the SR-IOV feature for hardware-based I/O sharing in IOV mode. Therefore, it is possible to partition SR-IOV capability on Ethernet controller NIC resources logically and expose them to a virtual machine as a separate PCI function called a “Virtual Function”. Refer to [Figure 11](#).

Therefore, a NIC is logically distributed among multiple virtual machines (as shown in [Figure 11](#)), while still having global data in common to share with the Physical Function and other Virtual Functions. The Intel® DPDK igbvf or ixgbev as a Poll Mode Driver (PMD) serves for the Intel® 82576 Gigabit Ethernet Controller, Intel® Ethernet Controller I350 family, or Intel® 82599 10 Gigabit Ethernet Controller NIC’s virtual PCI function. Meanwhile the Intel® DPDK Poll Mode Driver (PMD) also supports “Physical Function” of such NIC’s on the host.

The Intel® DPDK PF/VF Poll Mode Driver (PMD) supports the Layer 2 switch on Intel® 82576 Gigabit Ethernet Controller, Intel® Ethernet Controller I350 family, and Intel® 82599 10 Gigabit Ethernet Controller NICs so that guest can choose it for inter virtual machine traffic in SR-IOV mode.

For more detail on SR-IOV, please refer to the following documents:

- [SR-IOV provides hardware based I/O sharing](#)
- [PCI-SIG-Single Root I/O Virtualization Support on IA](#)
- [Scalable I/O Virtualized Servers](#)

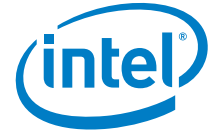
The diagram illustrates the SR-IOV architecture. At the top, three VMs (Virtual Machines) are shown in yellow, green, and red, each connected to a VF Driver (Virtual Function Driver). These VF Drivers are connected to the Virtual Machine Monitor (VMM). The VMM is connected to the Intel® VT-d (Virtualization Technology for Directed I/O) layer. The VT-d layer is connected to the PCI Express* interface. The PCI Express* interface is connected to the Virtual Functions (VFs) and the Physical Function (PF). The VFs are connected to the Virtual Ethernet Bridge and Classifier. The PF is connected to the Physical Ethernet Network. The Virtual Ethernet Bridge and Classifier is connected to the Physical Ethernet Network. The diagram shows the flow of traffic from the VMs through the VMM, VT-d, PCI Express, and the Virtual Functions to the Virtual Ethernet Bridge and Classifier, which then connects to the Physical Ethernet Network.

The following describes the Physical Function and Virtual Functions infrastructure for the supported Ethernet Controller NICs.

Virtual Functions operate under the respective Physical Function on the same NIC Port and therefore have no access to the global NIC resources that are shared between other functions for the same NIC port.

A Virtual Function has basic access to the queue resources and control structures of the queues assigned to it. For global resource access, a Virtual Function has to send a request to the Physical Function for that port, and the Physical Function operates on the global resources on behalf of the Virtual Function. For this out-of-band communication, an SR-IOV enabled NIC provides a memory buffer for each Virtual Function, which is called a "Mailbox".

The programmer can enable a maximum of *63 Virtual Functions* and there must be *one Physical Function* per Intel® 82599 10 Gigabit Ethernet Controller NIC port. The reason for this is that the device allows for a maximum of 128 queues per port and a



virtual/physical function has to have at least one queue pair (RX/TX). The current implementation of the Intel® DPDK ixgbevf driver supports a single queue pair (RX/TX) per Virtual Function and the host must have a Physical Function configured by the Linux* ixgbe driver (in the case of the Linux Kernel-based Virtual Machine [KVM]).

For example,

```
rmmod ixgbe (To remove the ixgbe module)
insmod ixgbe max_vfs=2,2 (To enable two Virtual Functions per port)
```

Virtual Function enumeration is performed in the following sequence by the Linux* pci driver for a dual-port NIC. When you enable the four Virtual Functions with the above command, the four enabled functions have a Function# represented by (Bus#, Device#, Function#) in sequence starting from 0 to 3. However:

- Virtual Functions 0 and 2 belong to Physical Function 0
- Virtual Functions 1 and 3 belong to Physical Function 1

Note: The above is an important consideration to take into account when targeting specific packets to a selected port.

13.1.1.2 Intel® 82576 Gigabit Ethernet Controller and Intel® Ethernet Controller I350 Family VF Infrastructure

In a virtualized environment, an Intel® 82576 Gigabit Ethernet Controller serves up to eight virtual machines (VMs). The controller has 16 TX and 16 RX queues. They are generally referred to (or thought of) as queue pairs (one TX and one RX queue). This gives the controller 16 queue pairs.

A pool is a group of queue pairs for assignment to the same VF, used for transmit and receive operations. The controller has eight pools, with each pool containing two queue pairs, that is, two TX and two RX queues assigned to each VF.

In a virtualized environment, an Intel® Ethernet Controller I350 family device serves up to eight virtual machines (VMs) per port. The eight queues can be accessed by eight different VMs if configured correctly, that means, one Transmit and one Receive queues assigned to each VF.

For example,

```
rmmod igb (To remove the igb module)
insmod igb max_vfs=2,2 (To enable two Virtual Functions per port)
```

Virtual Function enumeration is performed in the following sequence by the Linux* pci driver for a four-port NIC. When you enable the four Virtual Functions with the above command, the four enabled functions have a Function# represented by (Bus#, Device#, Function#) in sequence, starting from 0 to 7. However:

- Virtual Functions 0 and 4 belong to Physical Function 0
- Virtual Functions 1 and 5 belong to Physical Function 1
- Virtual Functions 2 and 6 belong to Physical Function 2
- Virtual Functions 3 and 7 belong to Physical Function 3

Note: The above is an important consideration to take into account when targeting specific packets to a selected port.



13.1.2 Validated Hypervisors

The validated hypervisor is:

- KVM (Kernel Virtual Machine) with Qemu, version 0.14.0

However, the hypervisor is bypassed to configure the Virtual Function devices using the Mailbox interface, the solution is hypervisor-agnostic. Xen* and VMware* (when SR-IOV is supported) will also be able to support the Intel® DPDK with Virtual Function driver support.

13.1.3 Expected Guest Operating System in Virtual Machine

The expected guest operating systems in a virtualized environment are:

- Fedora* 14 (64-bit)
- Ubuntu* 10.04 (64-bit)

For supported kernel versions, refer to the *Intel® DPDK Release Notes*.

13.2 Setting Up a KVM Virtual Machine Monitor

The following describes a target environment:

- Host Operating System: Fedora 14
- Hypervisor: KVM (Kernel Virtual Machine) with Qemu version 0.14.0
- Guest Operating System: Fedora 14
- Linux Kernel Version: Refer to the *Intel® DPDK Getting Started Guide*
- Target Applications: l2fwd-vf, l3fwd-vf

The setup procedure is as follows:

1. Before booting the Host OS, open **BIOS setup** and **enable Intel® VT features**.
2. While booting the Host OS kernel, pass the `intel_iommu=on` kernel command line argument using GRUB.
3. Download `qemu-kvm-0.14.0` from <http://sourceforge.net/projects/kvm/files/qemu-kvm/> and install it in the Host OS using the following steps:

When using a recent kernel (2.6.25+) with `kvm` modules included:

```
tar xzf qemu-kvm-release.tar.gz
cd qemu-kvm-release
./configure --prefix=/usr/local/kvm
make
sudo make install
sudo /sbin/modprobe kvm-intel
```

When using an older kernel, or a kernel from a distribution without the `kvm` modules, you must download (from the same link), compile and install the modules yourself:

```
tar xjf kvm-kmod-release.tar.bz2
cd kvm-kmod-release
./configure
make
sudo make install
sudo /sbin/modprobe kvm-intel
```




gemu-kvm installs in the /usr/local/bin directory.

For more details about KVM configuration and usage, please refer to:
<http://www.linux-kvm.org/page/HOWTO1>.

4. Create a Virtual Machine and install Fedora 14 on the Virtual Machine. This is referred to as the Guest Operating System (Guest OS).
5. Download and install the latest ixgbe driver from: http://downloadcenter.intel.com/Detail_Desc.aspx?agr=Y&DwnldID=14687
6. In the Host OS, unload the Linux ixgbe driver and reload it with the max_vfs=2,2 argument:

```
rmmod ixgbe
"modprobe ixgbe max_vfs=2,2"
```

Note that you need to explicitly specify number of vfs for each port, for example, in the command above, it creates two vfs for the first two ixgbe ports.

Let say we have a machine with four physical ixgbe ports:

```
0000:02:00.0
0000:02:00.1
0000:0e:00.0
0000:0e:00.1
```

The command above creates two vfs for device 0000:02:00.0:

```
ls -alrt /sys/bus/pci/devices/0000\:02\:00.0/virt*

lrwxrwxrwx. 1 root root 0 Apr 13 05:40 /sys/bus/pci/devices/0000:02:00.0/
virtfn1 -> ../0000:02:10.2
lrwxrwxrwx. 1 root root 0 Apr 13 05:40 /sys/bus/pci/devices/0000:02:00.0/
virtfn0 -> ../0000:02:10.0
```

It also creates two vfs for device 0000:02:00.1:

```
ls -alrt /sys/bus/pci/devices/0000\:02\:00.1/virt*
lrwxrwxrwx. 1 root root 0 Apr 13 05:51 /sys/bus/pci/devices/0000:02:00.1/
virtfn1 -> ../0000:02:10.3
lrwxrwxrwx. 1 root root 0 Apr 13 05:51 /sys/bus/pci/devices/0000:02:00.1/
virtfn0 -> ../0000:02:10.1
```

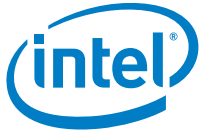
7. List the PCI devices connected and notice that the Host OS shows two Physical Functions (traditional ports) and four Virtual Functions (two for each port). This is the result of the previous step.
8. Insert the pci_stub module to hold the PCI devices that are freed from the default driver using the following command:

```
sudo /sbin/modprobe pci-stub
```

Unbind the default driver from the PCI devices representing the Virtual Functions. A script to perform this action is as follows:

```
echo "8086 10ed" > /sys/bus/pci/drivers/pci-stub/new_id
echo 0000:08:10.0 > /sys/bus/pci/devices/0000:08:10.0/driver/unbind
echo 0000:08:10.0 > /sys/bus/pci/drivers/pci-stub/bind
```

where, 0000:08:10.0 belongs to the Virtual Function visible in the Host OS.



9. Now, start the Virtual Machine by running the following command:

```
/usr/local/kvm/bin/qemu-system-x86_64 -m 4096 -smp 4 -boot c -hda  
lucid.qcow2 -device pci-assign,host=08:10.0
```

where:

- -m = memory to assign
- -smp = number of smp cores
- -boot = boot option
- -hda = virtual disk image
- -device = device to attach

Notes:

- The `pci-assign,host=08:10.0` value indicates that you want to attach a PCI device to a Virtual Machine and the respective (Bus:Device.Function) numbers should be passed for the Virtual Function to be attached.
- `qemu-kvm-0.14.0` allows a maximum of four PCI devices assigned to a VM, but this is `qemu-kvm` version dependent since `qemu-kvm-0.14.1` allows a maximum of five PCI devices.
- `qemu-system-x86_64` also has a `-cpu` command line option that is used to select the `cpu_model` to emulate in a Virtual Machine. Therefore, it can be used as:

```
/usr/local/kvm/bin/qemu-system-x86_64 -cpu ?  
(to list all available cpu_models)
```

```
/usr/local/kvm/bin/qemu-system-x86_64 -m 4096 -cpu host -smp 4 -boot c -hda  
lucid.qcow2 -device pci-assign,host=08:10.0  
(to use the same cpu_model equivalent to the host cpu)
```

For more information, please refer to: <http://wiki.qemu.org/Features/CPUModels>

10. Finally, access the Guest OS using `vncviewer` with the `localhost:5900` port and check the `lspci` command output in the Guest OS. The virtual functions will be listed as available for use.
11. Configure and install the Intel® DPDK with an `x86_64-default-linuxapp-gcc` configuration on the Guest OS as normal, that is, there is no change to the normal installation procedure.

```
make config T=x86_64-default-linuxapp-gcc O=x86_64-default-linuxapp-gcc  
cd x86_64-default-linuxapp-gcc  
make
```

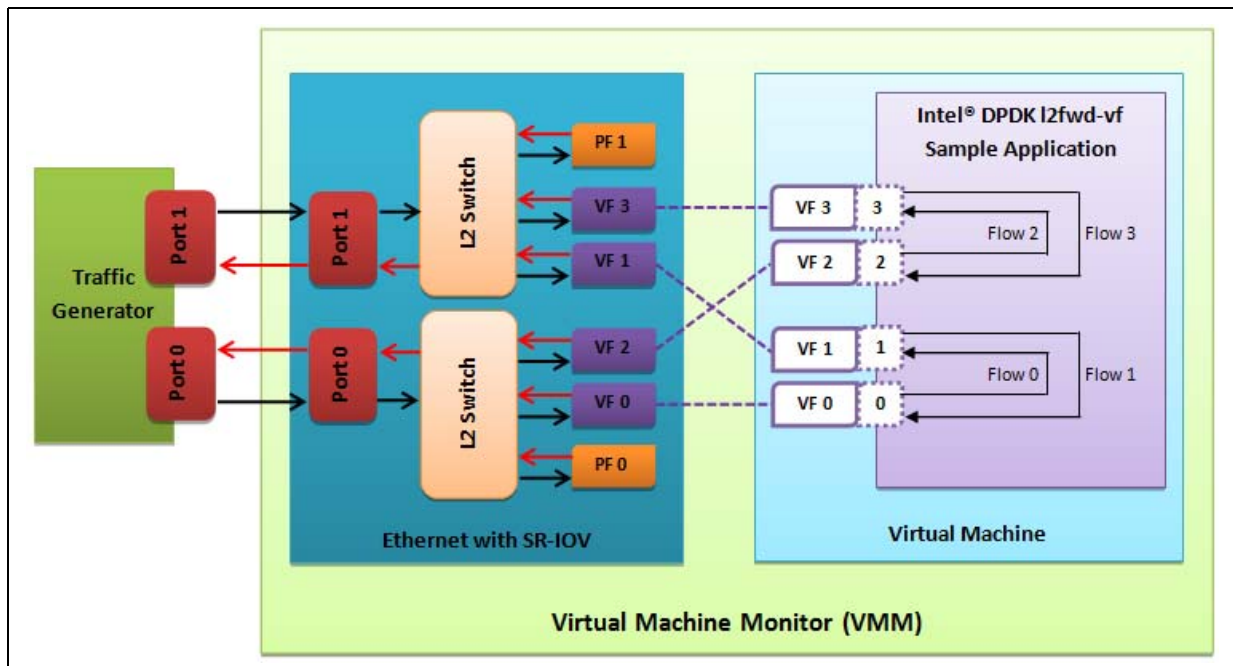
Note: If you are unable to compile the Intel® DPDK and you are getting “error: CPU you selected does not support x86-64 instruction set”, power off the Guest OS and start the virtual machine with the correct `-cpu` option in the `qemu-system-x86_64` command as shown in step 9. You must select the best `x86_64 cpu_model` to emulate or you can select `host` option if available.

12. Run the Intel® DPDK `l2fwd-vf` sample application in the Guest OS with Hugepages enabled. For the expected benchmark performance, you must pin the cores from the Guest OS to the Host OS (`taskset` can be used to do this) and you must also look at the PCI Bus layout on the board to ensure you are not running the traffic over the QPI Interface.

Notes:

- The Virtual Machine Manager (the Fedora package name is `virt-manager`) is a utility for virtual machine management that can also be used to create, start, stop and delete virtual machines. If this option is used, step 2 and 6 in the instructions provided will be different.
- `virsh`, a command line utility for virtual machine management, can also be used to bind and unbind devices to a virtual machine in Ubuntu. If this option is used, step 6 in the instructions provided will be different.
- The Virtual Machine Monitor (see [Figure 12](#)) is equivalent to a Host OS with KVM installed as described in the instructions.

Figure 12. Performance Benchmark Setup





14.0 Driver for VM Emulated Devices

The Intel® DPDK EM poll mode driver supports the following emulated devices:

- qemu-kvm emulated Intel® 82540EM Gigabit Ethernet Controller (qemu_e1000 device)
- VMware* emulated Intel® 82545EM Gigabit Ethernet Controller
- VMware emulated Intel® 8274L Gigabit Ethernet Controller.

14.1 Validated Hypervisors

The validated hypervisors are:

- KVM (Kernel Virtual Machine) with Qemu, version 0.14.0
- KVM (Kernel Virtual Machine) with Qemu, version 0.15.1
- VMware ESXi 5.0, Update 1

14.2 Expected Guest Operating System in Virtual Machine

The expected guest operating system in a virtualized environment is:

- Fedora* 14 (64-bit)

For supported kernel versions, refer to the *Intel® DPDK Release Notes*.

14.3 Setting Up a KVM Virtual Machine

The following describes a target environment:

- Host Operating System: Fedora 14
- Hypervisor: KVM (Kernel Virtual Machine) with Qemu version, 0.14.0
- Guest Operating System: Fedora 14
- Linux Kernel Version: Refer to the *Intel® DPDK Getting Started Guide*
- Target Applications: testpmd

The setup procedure is as follows:

1. Download qemu-kvm-0.14.0 from <http://sourceforge.net/projects/kvm/files/qemu-kvm/> and install it in the Host OS using the following steps:

When using a recent kernel (2.6.25+) with kvm modules included:

```
tar xzf qemu-kvm-release.tar.gz
cd qemu-kvm-release
./configure --prefix=/usr/local/kvm
```



```
make
sudo make install
sudo /sbin/modprobe kvm-intel
```

When using an older kernel or a kernel from a distribution without the `kvm` modules, you must download (from the same link), compile and install the modules yourself:

```
tar xjf kvm-kmod-release.tar.bz2
cd kvm-kmod-release
./configure
make
sudo make install
sudo /sbin/modprobe kvm-intel
```

Note that `qemu-kvm` installs in the `/usr/local/bin` directory.

For more details about KVM configuration and usage, please refer to:

<http://www.linux-kvm.org/page/HOWTO1>.

2. Create a Virtual Machine and install Fedora 14 on the Virtual Machine. This is referred to as the Guest Operating System (Guest OS).
3. Start the Virtual Machine with at least one emulated `e1000` device.

Note that Qemu provides several choices for the emulated network device backend. Most commonly used is a TAP networking backend that uses a TAP networking device in the host. For more information about Qemu supported networking backends and different options for configuring networking at Qemu, please refer to:

- <http://www.linux-kvm.org/page/Networking>
- <http://wiki.qemu.org/Documentation/Networking>
- <http://qemu.weilnetz.de/qemu-doc.html>

For example, to start a VM with two emulated `e1000` devices, issue the following command:

```
/usr/local/kvm/bin/qemu-system-x86_64 -cpu host -smp 4 -hda qemu1.raw -m 1024
-net nic,model=e1000,vlan=1,macaddr=DE:AD:1E:00:00:01
-net tap,vlan=1,ifname=tapvm01,script=no,downscript=no
-net nic,model=e1000,vlan=2,macaddr=DE:AD:1E:00:00:02
-net tap,vlan=2,ifname=tapvm02,script=no,downscript=no
```

where:

- `-m` = memory to assign
- `-smp` = number of smp cores
- `-hda` = virtual disk image

This command starts a new virtual machine with two emulated `82540EM` devices, backed up with two TAP networking host interfaces, `tapvm01` and `tapvm02`.

```
# ip tuntap show
tapvm01: tap
tapvm02: tap
```

4. Configure your TAP networking interfaces using `ip/ifconfig` tools.



5. Log in to the guest OS and check that the expected emulated devices exist:

```
# lspci -d 8086:100e

00:04.0 Ethernet controller: Intel Corporation 82540EM Gigabit Ethernet
Controller (rev 03)

00:05.0 Ethernet controller: Intel Corporation 82540EM Gigabit Ethernet
Controller (rev 03)
```

6. Install the Intel® DPDK and run testpmd.

14.4 Known Limitations of Emulated Devices

The following are know limitations:

1. The Qemu e1000 RX path does not support multiple descriptors/buffers per packet. Therefore, `rte_mbuf` should be big enough to hold the whole packet. For example, to allow testpmd to receive jumbo frames, use the following:

```
testpmd [options] -- --mbuf-size=<your-max-packet-size>
```

2. Qemu e1000 does not validate the checksum of incoming packets.



§ §

15.0 Kernel NIC Interface

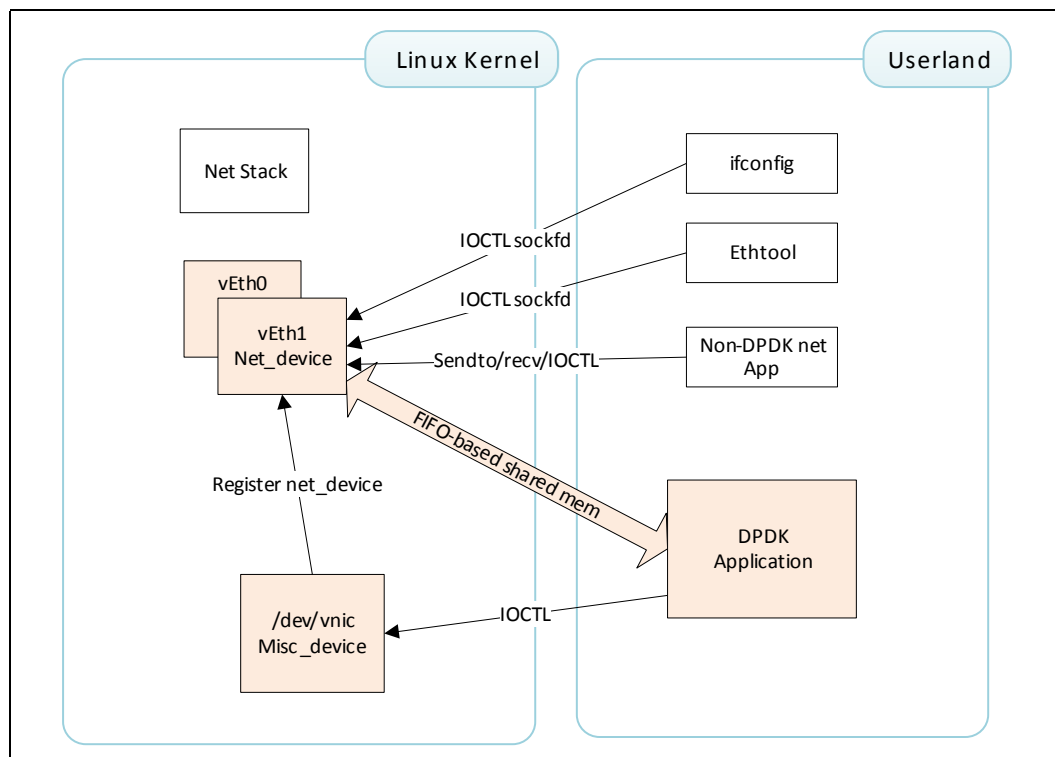
The Intel® DPDK Kernel NIC Interface (KNI) allows userspace applications access to the Linux* control plane.

The benefits of using the Intel® DPDK KNI are:

- Faster than existing Linux TUN/TAP interfaces (by eliminating system calls and `copy_to_user()/copy_from_user()` operations).
- Allows management of Intel® DPDK ports using standard Linux net tools such as `ethtool`, `ifconfig` and `tcpdump`.
- Allows an interface with the kernel network stack.

The components of an application using the Intel® DPDK Kernel NIC Interface are shown in Figure 13.

Figure 13. Components of a DPDK KNI Application





15.1 The Intel® DPDK KNI Kernel Module

The KNI kernel loadable module provides support for two types of devices:

- A Miscellaneous device (/dev/kni) that:
 - Creates net devices (via `ioctl` calls).
 - Maintains a kernel thread context shared by all KNI instances (simulating the RX side of the net driver).
- Net device:
 - Net functionality provided by implementing several operations such as `netdev_ops`, `header_ops`, `ethtool_ops` that are defined by `struct net_device`, including support for Intel® DPDK mbufs and FIFOs.
 - The interface name is provided from userspace.
 - The MAC address is the real NIC MAC address.

15.2 KNI Creation and Deletion

The KNI interfaces are created by an Intel® DPDK application dynamically. The interface name and FIFO details are provided by the application through an `ioctl` call using the `rte_kni_device_info` struct which contains:

- The interface name.
- Physical addresses of the corresponding memzones for the relevant FIFOs.
- Mbuf mempool details, both physical and virtual (to calculate the offset for mbuf pointers).
- PCI information.

Refer to `rte_kni_common.h` in the Intel® DPDK source code for more details.

The physical addresses will be re-mapped into the kernel address space and stored in separate KNI contexts.

Once KNI interfaces are created, the KNI context information can be queried by calling the `rte_kni_info_get()` function.

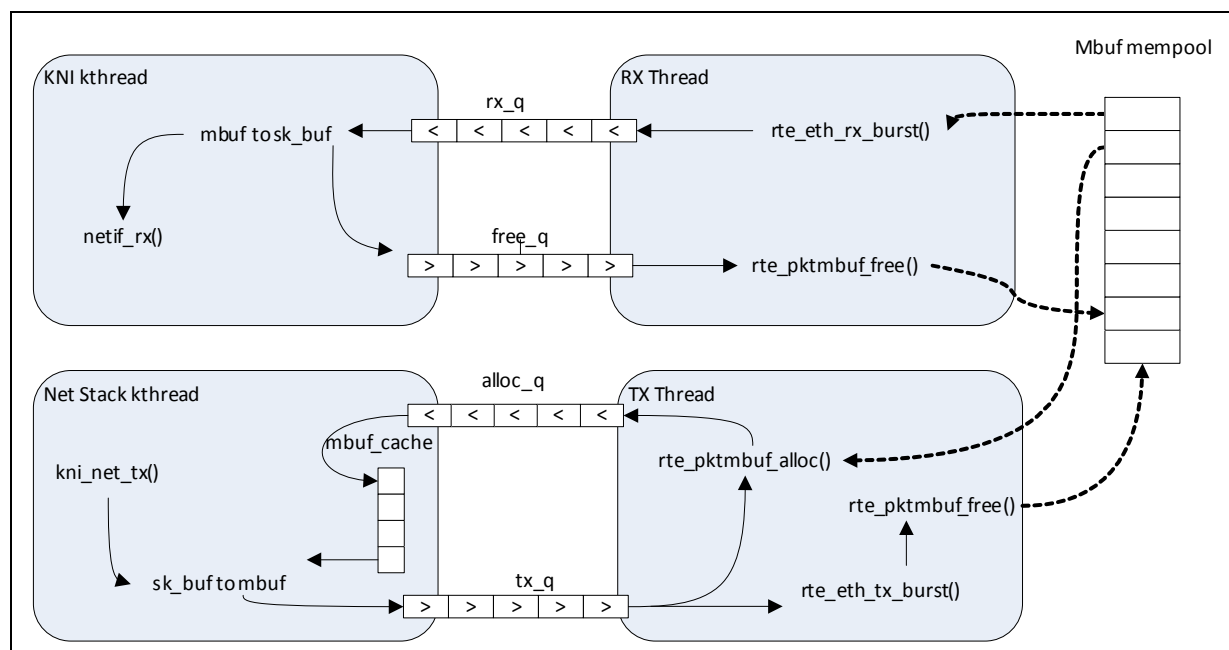
The KNI interfaces can be deleted by an Intel® DPDK application dynamically after being created. Furthermore, all those KNI interfaces not deleted will be deleted on the release operation of the miscellaneous device (when the Intel® DPDK application is closed).

15.3 Intel® DPDK mbuf Flow

To minimize the amount of Intel® DPDK code running in kernel space, the mbuf mempool is managed in userspace only. The kernel module will be aware of mbufs, but all mbuf allocation and free operations will be handled by the Intel® DPDK application only.

Figure 14 shows a typical scenario with packets sent in both directions.

Figure 14. Packet Flow via mbufs in the Intel® DPDK KNI



15.4 Use Case: Ingress

On the Intel® DPDK RX side, the mbuf is allocated by the PMD in the RX thread context. This thread will enqueue the mbuf in the rx_q FIFO. The KNI thread will poll all KNI active devices for the rx_q. If an mbuf is dequeued, it will be converted to a sk_buff and sent to the net stack via netif_rx(). The dequeued mbuf must be freed, so the same pointer is sent back in the free_q FIFO.

The RX thread, in the same main loop, polls this FIFO and frees the mbuf after dequeuing it.

15.5 Use Case: Egress

For packet egress the Intel® DPDK application must first enqueue several mbufs to create an mbuf cache on the kernel side.

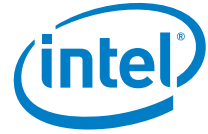
The packet is received from the Linux net stack, by calling the kni_net_tx() callback. The mbuf is dequeued (without waiting due the cache) and filled with data from sk_buff. The sk_buff is then freed and the mbuf sent in the tx_q FIFO.

The Intel® DPDK TX thread dequeues the mbuf and sends it to the PMD (via rte_eth_tx_burst()). It then puts the mbuf back in the cache.

15.6 Ethtool

Ethtool is a Linux-specific tool with corresponding support in the kernel where each net device must register its own callbacks for the supported operations. The current implementation uses the igb/ixgbe modified Linux drivers for ethtool support. Ethtool is not supported in VMs (VF or EM devices).





16.0 Thread Safety of Intel® DPDK Functions

The Intel® DPDK is comprised of several libraries. Some of the functions in these libraries can be safely called from multiple threads simultaneously, while others cannot. This section allows the developer to take these issues into account when building their own application.

The run-time environment of the Intel® DPDK is typically a single thread per logical core. In some cases, it is not only multi-threaded, but multi-process. Typically, it is best to avoid sharing data structures between threads and/or processes where possible. Where this is not possible, then the execution blocks must access the data in a thread-safe manner. Mechanisms such as atomics or locking can be used that will allow execution blocks to operate serially. However, this can have an effect on the performance of the application.

16.1 Fast-Path APIs

Applications operating in the data plane are performance sensitive but certain functions within those libraries may not be safe to call from multiple threads simultaneously. The hash, LPM and mempool libraries and RX/TX in the PMD are examples of this.

The hash and LPM libraries are, by design, thread unsafe in order to maintain performance. However, if required the developer can add layers on top of these libraries to provide thread safety. Locking is not needed in all situations, and in both the hash and LPM libraries, lookups of values can be performed in parallel in multiple threads. Adding, removing or modifying values, however, cannot be done in multiple threads without using locking.

The RX and TX of the PMD are the most critical aspects of an Intel® DPDK application and it is recommended that no locking be used as it will impact performance. Note, however, that these functions can safely be used from multiple threads when each thread is performing I/O on a different NIC queue. If multiple threads are to use the same hardware queue on the same NIC port, then locking, or some other form of mutual exclusion, is necessary.

The ring library is based on a lockless ring-buffer algorithm that maintains its original design for thread safety. Moreover, it provides high performance for either multi- or single-consumer/producer enqueue/dequeue operations. The mempool library is based on the Intel® DPDK lockless ring library and therefore is also multi-thread safe.

16.2 Performance Insensitive API

Outside of the performance sensitive areas described in [Section 16.1](#), the Intel® DPDK provides a thread-safe API for most other libraries. For example, `malloc(librte_malloc)` and `memzone` functions are safe for use in multi-threaded and multi-process environments.

The setup and configuration of the PMD is not performance sensitive, but is not thread safe either. It is possible that the multiple read/writes during PMD setup and configuration could be corrupted in a multi-thread environment. Since this is not



performance sensitive, the developer can choose to add their own layer to provide thread-safe setup and configuration. It is expected that, in most applications, the initial configuration of the network ports would be done by a single thread at startup.

16.3 Library Initialization

It is recommended that Intel® DPDK libraries are initialized in the main thread at application startup rather than subsequently in the forwarding threads. However, the Intel® DPDK performs checks to ensure that libraries are only initialized once. If initialization is attempted more than once, an error is returned.

In the multi-process case, the configuration information of shared memory will only be initialized by the master process. Thereafter, both master and secondary processes can allocate/release any objects of memory that finally rely on `rte_malloc` or `memzones`.

16.4 Interrupt Thread

The Intel® DPDK works almost entirely in Linux user space in polling mode. For certain infrequent operations, such as receiving a PMD link status change notification, callbacks may be called in an additional thread outside the main Intel® DPDK processing threads. These function callbacks should avoid manipulating Intel® DPDK objects that are also managed by the normal Intel® DPDK threads, and if they need to do so, it is up to the application to provide the appropriate locking or mutual exclusion restrictions around those objects.

§ §



Part 2: Development Environment



17.0 Source Organization

This section describes the organization of sources in the Intel® DPDK framework.

17.1 Makefiles and Config

Note: In the following descriptions, `RTE_SDK` is the environment variable that points to the base directory into which the tarball was extracted. See [Useful Variables Provided by the Build System](#) for descriptions of other variables.

Makefiles that are provided by the Intel® DPDK libraries and applications are located in `$(RTE_SDK)/mk`.

Config templates are located in `$(RTE_SDK)/config`. The templates describe the options that are enabled for each target. The config file also contains items that can be enabled and disabled for many of the Intel® DPDK libraries, including debug options. The user should look at the config file and become familiar with the options. The config file is also used to create a header file, which will be located in the new build directory.

17.2 Libraries

Libraries are located in subdirectories of `$(RTE_SDK)/lib`. By convention, we call a *library* any code that provides an API to an application. Typically, it generates an archive file (`.a`), but a kernel module should also go in the same directory.

The `lib` directory contains:

```
lib
+++ librte_cmdline      # command line interface helper
+++ librte_eal          # environment abstraction layer
+++ librte_ether        # generic interface to poll mode driver
+++ librte_hash         # hash library
+++ librte_lpm          # longest prefix match library
+++ librte_malloc       # malloc-like functions
+++ librte_mbuf         # packet and control mbuf manipulation library
+++ librte_mempool      # memory pool manager (fixedsize objects)
+++ librte_net          # various IP-related headers
+++ librte_pmd_e1000    # 1GbE poll mode drivers (igb and em)
+++ librte_pmd_ixgbe    # 10GbE poll mode driver
+++ librte_ring         # software rings (act as lockless FIFOs)
```



17.3 Applications

Applications are sources that contain a `main()` function. They are located in the `$(RTE_SDK)/app` and `$(RTE_SDK)/examples` directories.

The `app` directory contains sample applications that are used to test the Intel® DPDK (autotests). The `examples` directory contains sample applications that show how libraries can be used.

```

app
+-- chkincs                # test prog to check include depends
+-- test                   # autotests, to validate DPDK features
+-- test-pmd               # test and bench poll mode driver examples
+-- cmdline                # Example of using cmdline library
+-- dpdk_gat               # Example showing integration with Intel QuickAssist
+-- exception_path         # Sending packets to and from Linux ethernet device (TAP)
+-- helloworld             # Helloworld basic example
+-- ipv4_frag              # Example showing IPv4 Fragmentation
+-- ipv4_multicast         # Example showing IPv4 Multicast
+-- l2fwd                  # L2 Forwarding example with and without SR-IOV
+-- l3fwd                  # L3 Forwarding example
+-- l3fwd-vf               # L3 Forwarding example with SR-IOV
+-- link_status_interrupt  # Link status change interrupt example
+-- load_balancer          # Load balancing across multiple cores/sockets
+-- multi_process          # Example applications with multiple DPDK processes
+-- timer                  # Example of using librte_timer library
+-- vmdq_dcb               # Intel 82599 Ethernet Controller VMDQ and DCB receiving

```

Note: The actual `examples` directory may contain additional sample applications to those shown above. Check the latest Intel® DPDK source files for details.

§ §



18.0 Development Kit Build System

The Intel® DPDK requires a build system for compilation activities and so on. This section describes the constraints and the mechanisms used in the Intel® DPDK framework.

There are two use-cases for the framework:

- Compilation of the Intel® DPDK libraries and sample applications; the framework generates specific binary libraries, include files and sample applications
- Compilation of an external application or library, using an installed binary Intel® DPDK

18.1 Building the Development Kit Binary

The following provides details on how to build the Intel® DPDK binary.

18.1.1 Build Directory Concept

After installation, a build directory structure is created. Each build directory contains include files, libraries, and applications:

```
~/DPDK$ ls
app                               MAINTAINERS
config                           Makefile
COPYRIGHT                        mk
doc                              scripts
examples                         lib
tools                           x86_64-default-linuxapp-gcc
x86_64-default-linuxapp-icc      i686-default-linuxapp-gcc
i686-default-linuxapp-icc
...
~/DEV/DPDK$ ls i686-default-linuxapp-gcc
app build hostapp include kmod lib Makefile

~/DEV/DPDK$ ls i686-default-linuxapp-gcc/app/
chkinsc      dump_cfg      test      testpmd
chkinsc.map  dump_cfg.map test.map  testpmd.map

~/DEV/DPDK$ ls i686-default-linuxapp-gcc/lib/
libethdev.a      librte_hash.a      librte_mbuf.a      librte_pmd_ixgbe.a
librte_cmdline.a librte_lpm.a        librte_mempool.a   librte_ring.a
librte_eal.a     librte_malloc.a    librte_pmd_e1000.a librte_timer.a

~/DEV/DPDK$ ls i686-default-linuxapp-gcc/include/
arch
cmdline_cirbuf.h      rte_cpuflags.h      rte_memcpy.h
cmdline_parse.h       rte_cycles.h         rte_memory.h
cmdline_parse_ethaddr.h  rte_debug.h         rte_mempool.h
cmdline_parse_ipaddr.h  rte_eal.h            rte_memzone.h
cmdline_parse_num.h     rte_errno.h          rte_pci_dev_ids.h
cmdline_parse_portlist.h  rte_ethdev.h         rte_pci.h
cmdline_parse_string.h   rte_ether.h           rte_per_lcore.h
                        rte_fbk_hash.h       rte_prefetch.h
                        rte_hash_crc.h      rte_random.h
```




```

cmdline_rdtline.h      rte_hash.h             rte_ring.h
cmdline_socket.h      rte_interrupts.h    rte_rwlock.h
cmdline_vt100.h       rte_ip.h             rte_sctp.h
exec-env              rte_jhash.h          rte_spinlock.h
rte_alarm.h           rte_launch.h         rte_string_fns.h
rte_atomic.h          rte_lcore.h          rte_tailq.h
rte_branch_prediction.h  rte_log.h          rte_tcp.h
rte_byteorder.h       rte_lpm.h           rte_timer.h
rte_common.h          rte_malloc.h        rte_udp.h
rte_config.h          rte_mbuf.h

```

A build directory is specific to a configuration that includes architecture + execution environment + toolchain. It is possible to have several build directories sharing the same sources with different configurations.

For instance, to create a new build directory called `my_sdk_build_dir` using the default configuration template `config/defconfig_x86_64-linuxapp`, we use:

```

cd ${RTE_SDK}
make config T=x86_64-default-linuxapp-gcc O=my_sdk_build_dir

```

This creates a new `my_sdk_build_dir` directory. After that, we can compile by doing:

```

cd my_sdk_build_dir
make

```

which is equivalent to:

```

make O=my_sdk_build_dir

```

The content of the `my_sdk_build_dir` is then:

```

-- .config                # used configuration

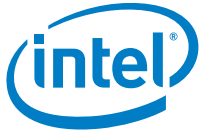
-- Makefile               # wrapper that calls head Makefile
                        # with $PWD as build directory

-- build                  # All temporary files used during build
+-- app                   # process, including .o, .d, and .cmd files.
|   |-- test              # For libraries, we have the .a file.
|   |   +-- test.o        # For applications, we have the elf file.
|   |   |-- ...
|   |-- lib
|       +-- librte_eal
|       |   |-- ...
|       +-- librte_mempool
|           +-- mempool-file1.o
|           +-- .mempool-file1.o.cmd
|           +-- .mempool-file1.o.d
|           +-- mempool-file2.o
|           +-- .mempool-file2.o.cmd
|           +-- .mempool-file2.o.d
|           |-- mempool.a
|       |-- ...
-- include                # All include files installed by libraries
+-- librte_mempool.h      # and applications are located in this
+-- rte_eal.h             # directory. The installed files can depend
+-- rte_spinlock.h        # on configuration if needed (environment,
+-- rte_atomic.h          # architecture, ...)
|-- *.h ...

-- lib                    # all compiled libraries are copied in this
+-- librte_eal.a          # directory
+-- librte_mempool.a
|-- *.a ...

-- app                    # All compiled applications are installed
+-- test                  # here. It includes the binary in elf format

```



Refer to [Development Kit Root Makefile Help](#) for details about make commands that can be used from the root of Intel® DPDK.

18.2 Building External Applications

Since Intel® DPDK is in essence a development kit, the first objective of end users will be to create an application using this SDK. To compile an application, the user must set the RTE_SDK and RTE_TARGET environment variables.

```
export RTE_SDK=/opt/DPDK
export RTE_TARGET=x86_64-default-linuxapp-gcc
cd /path/to/my_app
```

For a new application, the user must create their own Makefile that includes some .mk files, such as \${RTE_SDK}/mk/DPDK.vars.mk, and \${RTE_SDK}/mk/DPDK.app.mk. This is described in [Building Your Own Application](#).

Depending on the chosen target (architecture, machine, executive environment, toolchain) defined in the Makefile or as an environment variable, the applications and libraries will compile using the appropriate .h files and will link with the appropriate .a files. These files are located in \${RTE_SDK}/arch-machine-execenv-toolchain, which is referenced internally by \${RTE_BIN_SDK}.

To compile their application, the user just has to call make. The compilation result will be located in /path/to/my_app/build directory.

Sample applications are provided in the examples directory.

18.3 Makefile Description

18.3.1 General Rules For Intel® DPDK Makefiles

In the Intel® DPDK, Makefiles always follow the same scheme:

1. Include \$(RTE_SDK)/mk/DPDK.vars.mk at the beginning.
2. Define specific variables for RTE build system.
3. Include a specific \$(RTE_SDK)/mk/DPDK.XYZ.mk, where XYZ can be app, lib, extapp, extlib, obj, gnuconfigure, and so on, depending on what kind of object you want to build. See [Makefile Types](#) below.
4. Include user-defined rules and variables.

The following is a very simple example of an external application Makefile:

```
include $(RTE_SDK)/mk/DPDK.vars.mk

# binary name
APP = helloworld

# all source are stored in SRCS-y
SRCS-y := main.c

CFLAGS += -O3
CFLAGS += $(WERROR_FLAGS)

include $(RTE_SDK)/mk/DPDK.extapp.mk
```



18.3.2 Makefile Types

Depending on the .mk file which is included at the end of the user Makefile, the Makefile will have a different role. Note that it is not possible to build a library and an application in the same Makefile. For that, the user must create two separate Makefiles, possibly in two different directories.

In any case, the `rte.vars.mk` file must be included in the user Makefile as soon as possible.

18.3.2.1 Application

These Makefiles generate a binary application.

- `rte.app.mk`: Application in the development kit framework
- `rte.extapp.mk`: External application
- `rte.hostapp.mk`: Host application in the development kit framework

18.3.2.2 Library

Generate a .a library.

- `rte.lib.mk`: Library in the development kit framework
- `rte.extlib.mk`: external library
- `rte.hostlib.mk`: host library in the development kit framework

18.3.2.3 Install

- `rte.install.mk`: Does not build anything, it is only used to create links or copy files to the installation directory. This is useful for including files in the development kit framework.

18.3.2.4 Kernel Module

- `rte.module.mk`: Build a kernel module in the development kit framework.

18.3.2.5 Objects

- `rte.obj.mk`: Object aggregation (merge several .o in one) in the development kit framework.
- `rte.extobj.mk`: Object aggregation (merge several .o in one) outside the development kit framework.

18.3.2.6 Misc

- `rte.doc.mk`: Documentation in the development kit framework
- `rte.gnuconfigure.mk`: Build an application that is configure-based (used to compile *newlib*).
- `rte.subdir.mk`: Build several directories in the development kit framework.

18.3.3 Useful Variables Provided by the Build System

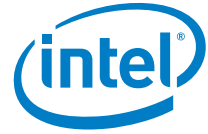
- `RTE_SDK`: The absolute path to the Intel® DPDK sources. When compiling the development kit, this variable is automatically set by the framework. It has to be defined by the user as an environment variable if compiling an external application.



- **RTE_SRCDIR:** The path to the root of the sources. When compiling the development kit, `RTE_SRCDIR = RTE_SDK`. When compiling an external application, the variable points to the root of external application sources.
- **RTE_OUTPUT:** The path to which output files are written. Typically, it is `$(RTE_SRCDIR)/build`, but it can be overridden by the `O=` option in the make command line.
- **RTE_TARGET:** A string identifying the target for which we are building. The format is `arch-machine-execenv-toolchain`. When compiling the SDK, the target is deduced by the build system from the configuration (`.config`). When building an external application, it must be specified by the user in the Makefile or as an environment variable.
- **RTE_SDK_BIN:** References `$(RTE_SDK)/$(RTE_TARGET)`.
- **RTE_ARCH:** Defines the architecture (`i686`, `x86_64`). It is the same value as `CONFIG_RTE_ARCH` but without the double-quotes around the string.
- **RTE_MACHINE:** Defines the machine. It is the same value as `CONFIG_RTE_MACHINE` but without the double-quotes around the string.
- **RTE_TOOLCHAIN:** Defines the toolchain (`gcc`, `icc`). It is the same value as `CONFIG_RTE_TOOLCHAIN` but without the double-quotes around the string.
- **RTE_EXEC_ENV:** Defines the executive environment (`linuxapp`). It is the same value as `CONFIG_RTE_EXEC_ENV` but without the double-quotes around the string.
- **RTE_KERNELDIR:** This variable contains the absolute path to the kernel sources that will be used to compile the kernel modules. The kernel headers must be the same as the ones that will be used on the target machine (the machine that will run the application). By default, the variable is set to `/lib/modules/$(shell uname -r)/build`, which is correct when the target machine is also the build machine.

18.3.4 Variables that Can be Set/Overridden in a Makefile Only

- **VPATH:** The path list that the build system will search for sources. By default, `RTE_SRCDIR` will be included in `VPATH`.
- **CFLAGS:** Flags to use for C compilation. The user should use `+=` to append data in this variable.
- **LDFLAGS:** Flags to use for linking. The user should use `+=` to append data in this variable.
- **ASFLAGS:** Flags to use for assembly. The user should use `+=` to append data in this variable.
- **CPPFLAGS:** Flags to use to give flags to C preprocessor (only useful when assembling `.S` files). The user should use `+=` to append data in this variable.
- **LDLIBS:** In an application, the list of libraries to link with (for example, `-L /path/to/libfoo -lfoo`). The user should use `+=` to append data in this variable.
- **SRC-y:** A list of source files (`.c`, `.S`, or `.o` if the source is a binary) in case of application, library or object Makefiles. The sources must be available from `VPATH`.
- **INSTALL-y-\$(INSTPATH):** A list of files to be installed in `$(INSTPATH)`. The files must be available from `VPATH` and will be copied in `$(RTE_OUTPUT)/$(INSTPATH)`. Can be used in almost any RTE Makefile.
- **SYMLINK-y-\$(INSTPATH):** A list of files to be installed in `$(INSTPATH)`. The files must be available from `VPATH` and will be linked (symbolically) in `$(RTE_OUTPUT)/$(INSTPATH)`. This variable can be used in almost any Intel® DPDK Makefile.



- **PREBUILD:** A list of prerequisite actions to be taken before building. The user should use += to append data in this variable.
- **POSTBUILD:** A list of actions to be taken after the main build. The user should use += to append data in this variable.
- **PREINSTALL:** A list of prerequisite actions to be taken before installing. The user should use += to append data in this variable.
- **POSTINSTALL:** A list of actions to be taken after installing. The user should use += to append data in this variable.
- **PRECLEAN:** A list of prerequisite actions to be taken before cleaning. The user should use += to append data in this variable.
- **POSTCLEAN:** A list of actions to be taken after cleaning. The user should use += to append data in this variable.
- **DEPDIR-y:** Only used in the development kit framework to specify if the build of the current directory depends on build of another one. This is needed to support parallel builds correctly.

18.3.5 Variables that can be Set/Overridden by the User on the Command Line Only

Some variables can be used to configure the build system behavior. They are documented in [Development Kit Root Makefile Help](#) and [External Application/Library Makefile help](#).

- **WERROR_CFLAGS:** By default, this is set to a specific value that depends on the compiler. Users are encouraged to use this variable as follows:

```
CFLAGS += $(WERROR_CFLAGS)
```

This avoids the use of different cases depending on the compiler (icc or gcc). Also, this variable can be overridden from the command line, which allows bypassing of the flags for testing purposes.

18.3.6 Variables that Can be Set/Overridden by the User in a Makefile or Command Line

- **CFLAGS_my_file.o:** Specific flags to add for C compilation of my_file.c.
- **LDFLAGS_my_app:** Specific flags to add when linking my_app.
- **NO_AUTOLIBS:** If set, the libraries provided by the framework will not be included in the LDLIBS variable automatically.
- **EXTRA_CFLAGS:** The content of this variable is appended after CFLAGS when compiling.
- **EXTRA_LDFLAGS:** The content of this variable is appended after LDFLAGS when linking.
- **EXTRA_ASFLAGS:** The content of this variable is appended after ASFLAGS when assembling.
- **EXTRA_CPPFLAGS:** The content of this variable is appended after CPPFLAGS when using a C preprocessor on assembly files.

§ §



19.0 Development Kit Root Makefile Help

The Intel® DPDK provides a root level Makefile with targets for configuration, building, cleaning, testing, installation and others. These targets are explained in the following sections.

19.1 Configuration Targets

The configuration target requires the name of the target, which is specified using `T=mytarget` and it is mandatory. The list of available targets are in `$(RTE_SDK) / config` (remove the `defconfig_` prefix).

Configuration targets also support the specification of the name of the output directory, using `O=mybuilddir`. This is an optional parameter, the default output directory is `build`.

- `config`
This will create a build directory, and generates a configuration from a template. A Makefile is also created in the new build directory.

Example: `make config O=mybuild T=x86_64-default-linuxapp-gcc`

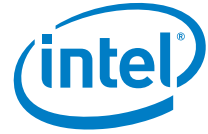
19.2 Build Targets

Build targets support the optional specification of the name of the output directory, using `O=mybuilddir`. The default output directory is `build`.

- `all`, `build` or just `make`
Build the Intel® DPDK in the output directory previously created by a `make config`.
Example: `make O=mybuild`
- `clean`
Clean all objects created using `make build`.
Example: `make clean O=mybuild`
- `%_sub`
Build a subdirectory only, without managing dependencies on other directories.
Example: `make lib/librte_eal_sub O=mybuild`
- `%_clean`
Clean a subdirectory only.
Example: `make lib/librte_eal_clean O=mybuild`

19.3 Install Targets

- `install`
Build the Intel® DPDK binary. Actually, this builds each supported target in a separate directory. The name of each directory is the name of the target.



The name of the targets to install can be optionally specified using `T=mytarget`. The target name can contain wildcard `*` characters. The list of available targets are in `$(RTE_SDK)/config` (remove the `defconfig_` prefix).

Example: `make install T=x86_64-*`

- `uninstall`
Remove installed target directories.

19.4 Test Targets

- `test`
Launch automatic tests for a build directory specified using `O=mybuilddir`. It is optional, the default output directory is `build`.
Example: `make test O=mybuild`
- `testall`
Launch automatic tests for all installed target directories (after a `make install`). The name of the targets to test can be optionally specified using `T=mytarget`. The target name can contain wildcard (`*`) characters. The list of available targets are in `$(RTE_SDK)/config` (remove the `defconfig_` prefix).
Examples: `make testall`, `make testall T=x86_64-*`

19.5 Documentation Targets

- `doxydoc`
Generate the Doxygen documentation (pdf only).

19.6 Deps Targets

- `depdirs`
This target is implicitly called by `make config`. Typically, there is no need for a user to call it, except if `DEPDIRS-y` variables have been updated in Makefiles. It will generate the file `$(RTE_OUTPUT)/.depdirs`.
Example: `make depdirs O=mybuild`
- `depgraph`
This command generates a dot graph of dependencies. It can be displayed to debug circular dependency issues, or just to understand the dependencies.
Example: `make depgraph O=mybuild > /tmp/graph.dot && dotty /tmp/graph.dot`

19.7 Misc Targets

- `help`
Show this help.

19.8 Other Useful Command-line Variables

The following variables can be specified on the command line:

- `V=`
Enable verbose build (show full compilation command line, and some intermediate commands).
- `D=`



Enable dependency debugging. This provides some useful information about why a target is built or not.

- EXTRA_CFLAGS=, EXTRA_LDFLAGS=, EXTRA_ASFLAGS=, EXTRA_CPPFLAGS=
Append specific compilation, link or asm flags.
- CROSS=
Specify a cross toolchain header that will prefix all gcc/binutils applications. This only works when using gcc.

19.9 Make in a Build Directory

All targets described above are called from the SDK root \$(RTE_SDK). It is possible to run the same Makefile targets inside the build directory. For instance, the following command:

```
cd $(RTE_SDK)
make config O=mybuild T=x86_64-default-linuxapp-gcc
make O=mybuild
```

is equivalent to:

```
cd $(RTE_SDK)
make config O=mybuild T=x86_64-default-linuxapp-gcc
cd mybuild
# no need to specify O= now
make
```

19.10 Compiling for Debug

To compile the Intel® DPDK and sample applications with debugging information included and the optimization level set to 0, the EXTRA_CFLAGS environment variable should be set before compiling as follows:

```
export EXTRA_CFLAGS='-O0 -g'
```

The Intel® DPDK and any user or sample applications can then be compiled in the usual way. For example:

```
make install T=x86_64-default-linuxapp-gcc
make -C examples/<theapp>
```

§ §



20.0 Extending the Intel® DPDK

This chapter describes how a developer can extend the Intel® DPDK to provide a new library, a new target, or support a new target.

20.1 Example: Adding a New Library libfoo

To add a new library to the Intel® DPDK, proceed as follows:

1. Add a new configuration option:

```
for f in config/*; do \
    echo CONFIG_RTE_LIBFOO=y >> $f; done
```

2. Create a new directory with sources:

```
mkdir ${RTE_SDK}/lib/libfoo
touch ${RTE_SDK}/lib/libfoo/foo.c
touch ${RTE_SDK}/lib/libfoo/foo.h
```

3. Add a `foo()` function in `libfoo`.

Definition is in `foo.c`:

```
void foo(void)
{
}
```

Declaration is in `foo.h`:

```
extern void foo(void);
```

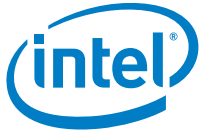
4. Update `lib/Makefile`:

```
vi ${RTE_SDK}/lib/Makefile
# add:
# DIRS-$(CONFIG_RTE_LIBFOO) += libfoo
```

5. Create a new Makefile for this library, for example, derived from `mempool` Makefile:

```
cp ${RTE_SDK}/lib/librte_mempool/Makefile ${RTE_SDK}/lib/libfoo/
vi ${RTE_SDK}/lib/libfoo/Makefile
# replace:
# librte_mempool -> libfoo
# rte_mempool -> foo
```

6. Update `mk/DPDK.app.mk`, and add `-lfoo` in `LDLIBS` variable when the option is enabled. This will automatically add this flag when linking an Intel® DPDK application.



7. Build the Intel® DPDK with the new library (we only show a specific target here):

```
cd ${RTE_SDK}
make config T=x86_64-default-linuxapp-gcc
make
```

8. Check that the library is installed:

```
ls build/lib
ls build/include
```

20.1.1 Example: Using libfoo in the Test Application

The test application is used to validate all functionality of the Intel® DPDK. Once you have added a library, a new test case should be added in the test application.

- A new `test_foo.c` file should be added, that includes `foo.h` and calls the `foo()` function from `test_foo()`. When the test passes, the `test_foo()` function should return 0.
- `Makefile`, `test.h` and `commands.c` must be updated also, to handle the new test case.
- Test report generation: `autotest.py` is a script that is used to generate the test report that is available in the `${RTE_SDK}/doc/rst/test_report/autotests` directory. This script must be updated also. If `libfoo` is in a new test family, the links in `${RTE_SDK}/doc/rst/test_report/test_report.rst` must be updated.
- Build the Intel® DPDK with the updated test application (we only show a specific target here):

```
cd ${RTE_SDK}
make config T=x86_64-default-linuxapp-gcc
make
```





21.0 Building Your Own Application

21.1 Compiling a Sample Application in the Development Kit Directory

When compiling a sample application (for example, hello world), the following variables must be exported: RTE_SDK and RTE_TARGET.

```
~/DPDK$ cd examples/helloworld/
~/DPDK/examples/helloworld$ export RTE_SDK=/home/user/DPDK
~/DPDK/examples/helloworld$ export RTE_TARGET=x86_64-default-linuxapp-gcc
~/DPDK/examples/helloworld$ make
CC main.o
CC commands.o
CC parse_obj_list.o
LD test-helloworld
INSTALL-APP helloworld
INSTALL-MAP helloworld.map
```

The binary is generated in the build directory by default:

```
~/DPDK/examples/helloworld$ ls build/app
test-helloworld test-helloworld.map
```

21.2 Build Your Own Application Outside the Development Kit

The sample application (Hello World) can be duplicated in a new directory as a starting point for your development:

```
~$ cp -r DPDK/examples/helloworld my_rte_app
~$ cd my_rte_app/
~/DPDK/examples/helloworld$ export RTE_SDK=/home/user/DPDK
~/DPDK/examples/helloworld$ export RTE_TARGET=x86_64-default-linuxapp-gcc
~/my_rte_app$ make
CC main.o
CC commands.o
CC parse_obj_list.o
LD test-helloworld
INSTALL-APP test-helloworld
INSTALL-MAP test-helloworld.map
```

21.3 Customizing Makefiles

21.3.1 Application Makefile

The default makefile provided with the Hello World sample application is a good starting point. It includes:

- \$(RTE_SDK)/mk/DPDK.vars.mk at the beginning
- \$(RTE_SDK)/mk/DPDK.extapp.mk at the end



The user must define several variables:

- APP: Contains the name of the application.
- SRCS-y: List of source files (*.c, *.S).

21.3.2 Library Makefile

It is also possible to build a library in the same way:

- Include \$(RTE_SDK)/mk/DPDK.vars.mk at the beginning.
- Include \$(RTE_SDK)/mk/DPDK.extlib.mk at the end.

The only difference is that APP should be replaced by LIB, which contains the name of the library. For example, libfoo.a.

21.3.3 Customize Makefile Actions

Some variables can be defined to customize Makefile actions. The most common are listed below. Refer to [Makefile Description](#) section in [Development Kit Build System](#) chapter for details.

- VPATH: The path list where the build system will search for sources. By default, RTE_SRCDIR will be included in VPATH.
- CFLAGS_my_file.o: The specific flags to add for C compilation of my_file.c.
- CFLAGS: The flags to use for C compilation.
- LDFLAGS: The flags to use for linking.
- CPPFLAGS: The flags to use to provide flags to the C preprocessor (only useful when assembling .S files)
- LDLIBS: A list of libraries to link with (for example, -L /path/to/libfoo -lfoo)
- NO_AUTOLIBS: If set, the libraries provided by the framework will not be included in the LDLIBS variable automatically.

§ §



22.0 External Application/Library Makefile help

External applications or libraries should include specific Makefiles from RTE_SDK, located in mk directory. These Makefiles are:

- `${RTE_SDK}/mk/DPDK.extapp.mk`: Build an application
- `${RTE_SDK}/mk/DPDK.extlib.mk`: Build a static library
- `${RTE_SDK}/mk/DPDK.extobj.mk`: Build objects (.o)

22.1 Prerequisites

The following variables must be defined:

- `${RTE_SDK}`: Points to the root directory of the Intel® DPDK.
- `${RTE_TARGET}`: Reference the target to be used for compilation (for example, `x86_64-default-linuxapp-gcc`).

22.2 Build Targets

Build targets support the specification of the name of the output directory, using `O=mybuilddir`. This is optional; the default output directory is `build`.

- `all`, “nothing” (meaning just `make`)
Build the application or the library in the specified output directory.
Example: `make O=mybuild`
- `clean`
Clean all objects created using `make build`.
Example: `make clean O=mybuild`

22.3 Help Targets

- `help`
Show this help.

22.4 Other Useful Command-line Variables

The following variables can be specified at the command line:

- `S=`
Specify the directory in which the sources are located. By default, it is the current directory.
- `M=`
Specify the Makefile to call once the output directory is created. By default, it uses `$(S)/Makefile`.



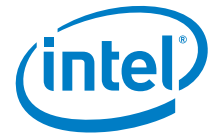
- V=
Enable verbose build (show full compilation command line and some intermediate commands).
- D=
Enable dependency debugging. This provides some useful information about why a target must be rebuilt or not.
- EXTRA_CFLAGS=, EXTRA_LDFLAGS=, EXTRA_ASFLAGS=, EXTRA_CPPFLAGS=
Append specific compilation, link or asm flags.
- CROSS=
Specify a cross-toolchain header that will prefix all gcc/binutils applications. This only works when using gcc.

22.5 Make from Another Directory

It is possible to run the Makefile from another directory, by specifying the output and the source dir. For example:

```
export RTE_SDK=/path/to/DPDK
export RTE_TARGET=x86_64-default-linuxapp-icc
make -f /path/to/my_app/Makefile S=/path/to/my_app O=/path/to/build_dir
```

§ §



Part 3: Performance Optimization



23.0 Performance Optimization Guidelines

23.1 Introduction

The following sections describe optimizations used in the Intel® DPDK and optimizations that should be considered for a new applications.

They also highlight the performance-impacting coding techniques that should, and should not be, used when developing an application using the Intel® DPDK.

And finally, they give an introduction to application profiling using a Performance Analyzer from Intel to optimize the software.

§ §



24.0 Writing Efficient Code

This chapter provides some tips for developing efficient code using the Intel® DPDK. For additional and more general information, please refer to the [Intel® 64 and IA-32 Architectures Optimization Reference Manual](#) which is a valuable reference to writing efficient code.

24.1 Memory

This section describes some key memory considerations when developing applications in the Intel® DPDK environment.

24.1.1 Memory Copy: Do not Use libc in the Data Plane

Many `libc` functions are available in the Intel® DPDK, via the Linux* application environment. This can ease the porting of applications and the development of the configuration plane. However, many of these functions are not designed for performance. Functions such as `memcpy()` or `strcpy()` should not be used in the data plane. To copy small structures, the preference is for a simpler technique that can be optimized by the compiler. Refer to the *VTune™ Performance Analyzer Essentials* publication from Intel Press for recommendations.

For specific functions that are called often, it is also a good idea to provide a self-made optimized function, which should be declared as `static inline`.

The Intel® DPDK API provides an optimized `rte_memcpy()` function.

24.1.2 Memory Allocation

Other functions of `libc`, such as `malloc()`, provide a flexible way to allocate and free memory. In some cases, using dynamic allocation is necessary, but it is really not advised to use `malloc`-like functions in the data plane because managing a fragmented heap can be costly and the allocator may not be optimized for parallel allocation.

If you really need dynamic allocation in the data plane, it is better to use a memory pool of fixed-size objects. This API is provided by `librte_mempool`. This data structure provides several services that increase performance, such as memory alignment of objects, lockless access to objects, NUMA awareness, bulk get/put and per-lcore cache. The `rte_malloc()` function uses a similar concept to mempools.

24.1.3 Concurrent Access to the Same Memory Area

Read-Write (RW) access operations by several lcores to the same memory area can generate a lot of data cache misses, which are very costly. It is often possible to use per-lcore variables, for example, in the case of statistics. There are at least two solutions for this:

- Use `RTE_PER_LCORE` variables. Note that in this case, data on lcore X is not available to lcore Y.

- Use a table of structures (one per lcore). In this case, each structure must be cache-aligned.

Read-mostly variables can be shared among lcores without performance losses if there are no RW variables in the same cache line.

24.1.4 NUMA

On a NUMA system, it is preferable to access local memory since remote memory access is slower. In the Intel® DPDK, the memzone, ring, rte_malloc and mempool APIs provide a way to create a pool on a specific socket.

Sometimes, it can be a good idea to duplicate data to optimize speed. For read-mostly variables that are often accessed, it should not be a problem to keep them in one socket only, since data will be present in cache.

24.1.5 Distribution Across Memory Channels

Modern memory controllers have several memory channels that can load or store data in parallel. Depending on the memory controller and its configuration, the number of channels and the way the memory is distributed across the channels varies. Each channel has a bandwidth limit, meaning that if all memory access operations are done on the first channel only, there is a potential bottleneck.

By default, the [Mempool Library](#) spreads the addresses of objects among memory channels.

24.2 Communication Between lcores

To provide a message-based communication between lcores, it is advised to use the Intel® DPDK ring API, which provides a lockless ring implementation.

The ring supports *bulk* and *burst* access, meaning that it is possible to read several elements from the ring with only one costly atomic operation (see [Chapter 5.0, “Ring Library”](#)). Performance is greatly improved when using bulk access operations.

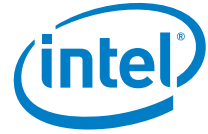
The code algorithm that dequeues messages may be something similar to the following:

```
#define MAX_BULK 32
while (1) {
    /* Process as many elements as can be dequeued. */
    count = rte_ring_dequeue_burst(ring, obj_table, MAX_BULK);
    if (unlikely(count == 0))
        continue;
    my_process_bulk(obj_table, count);
}
```

24.3 PMD Driver

The Intel® DPDK Poll Mode Driver (PMD) is also able to work in bulk/burst mode, allowing the factorization of some code for each call in the send or receive function.

Avoid partial writes. When PCI devices write to system memory through DMA, it costs less if the write operation is on a full cache line as opposed to part of it. In the PMD code, actions have been taken to avoid partial writes as much as possible.



24.3.1 Lower Packet Latency

Traditionally, there is a trade-off between throughput and latency. An application can be tuned to achieve a high throughput, but the end-to-end latency of an average packet will typically increase as a result. Similarly, the application can be tuned to have, on average, a low end-to-end latency, at the cost of lower throughput.

In order to achieve higher throughput, the Intel® DPDK attempts to aggregate the cost of processing each packet individually by processing packets in bursts.

Using the `testpmd` application as an example, the burst size can be set on the command line to a value of 16 (also the default value). This allows the application to request 16 packets at a time from the PMD. The `testpmd` application then immediately attempts to transmit all the packets that were received, in this case, all 16 packets.

The packets are not transmitted until the tail pointer is updated on the corresponding TX queue of the network port. This behavior is desirable when tuning for high throughput because the cost of tail pointer updates to both the RX and TX queues can be spread across 16 packets, effectively hiding the relatively slow MMIO cost of writing to the PCIe* device. However, this is not very desirable when tuning for low latency because the first packet that was received must also wait for another 15 packets to be received. It cannot be transmitted until the other 15 packets have also been processed because the NIC will not know to transmit the packets until the TX tail pointer has been updated, which is not done until all 16 packets have been processed for transmission.

To consistently achieve low latency, even under heavy system load, the application developer should avoid processing packets in bunches. The `testpmd` application can be configured from the command line to use a burst value of 1. This will allow a single packet to be processed at a time, providing lower latency, but with the added cost of lower throughput.

24.4 Locks and Atomic Operations

Atomic operations imply a `lock` prefix before the instruction, causing the processor's `LOCK#` signal to be asserted during execution of the following instruction. This has a big impact on performance in a multicore environment.

Performance can be improved by avoiding lock mechanisms in the data plane. It can often be replaced by other solutions like per-core variables. Also, some locking techniques are more efficient than others. For instance, the Read-Copy-Update (RCU) algorithm can frequently replace simple `rwlocks`.

24.5 Coding Considerations

24.5.1 Inline Functions

Small functions can be declared as `static inline` in the header file. This avoids the cost of a `call` instruction (and the associated context saving). However, this technique is not always efficient; it depends on many factors including the compiler.

24.5.2 Branch Prediction

The Intel® C/C++ Compiler (`icc`)/`gcc` built-in helper functions `likely()` and `unlikely()` allow the developer to indicate if a code branch is likely to be taken or not. For instance:

```
if (likely(x > 1))
    do_stuff();
```



24.6 Setting the Target CPU Type

The Intel® DPDK supports CPU microarchitecture-specific optimizations by means of `CONFIG_RTE_MACHINE` option in the Intel® DPDK configuration file. The degree of optimization depends on the compiler's ability to optimize for a specific microarchitecture, therefore it is preferable to use the latest compiler versions whenever possible.

If the compiler version does not support the specific feature set (for example, the Intel® AVX instruction set), the build process gracefully degrades to whatever latest feature set is supported by the compiler.

Since the build and runtime targets may not be the same, the resulting binary also contains a platform check that runs before the `main()` function and checks if the current machine is suitable for running the binary.

Along with compiler optimizations, a set of preprocessor defines are automatically added to the build process (regardless of the compiler version). These defines correspond to the instruction sets that the target CPU should be able to support. For example, a binary compiled for any SSE4.2-capable processor will have `RTE_MACHINE_CPUFLAG_SSE4_2` defined, thus enabling compile-time code path selection for different platforms.

§ §



25.0 Profile Your Application

Intel processors provide performance counters to monitor events. Some tools provided by Intel can be used to profile and benchmark an application. See the *VTune™ Performance Analyzer Essentials* publication from Intel Press for more information.

For an Intel® DPDK application, this can be done in a Linux* application environment only.

The main situations that should be monitored through event counters are:

- Cache misses
- Branch mis-predicts
- DTLB misses
- Long latency instructions and exceptions

Refer to the [Intel Performance Analysis Guide](#) for details about application profiling.

26.0 Glossary

API	Application Programming Interface
ASLR	Linux* kernel Address-Space Layout Randomization
BSD	Berkeley Software Distribution
CIDR	Classless Inter-Domain Routing
Control Plane	The control plane is concerned with the routing of packets and with providing a start or end point.
Core	A core may include several <i>lcores</i> or <i>threads</i> if the processor supports hyperthreading.
Core Components	A set of libraries provided by the Intel® DPDK, including eal, ring, mempool, mbuf, timers, and so on.
CPU	Central Processing Unit
CRC	Cyclic Redundancy Check
ctrlmbuf	An <i>mbuf</i> carrying control data.
Data Plane	In contrast to the control plane, the data plane in a network architecture are the layers involved when forwarding packets. These layers must be highly optimized to achieve good performance.
DIMM	Dual In-line Memory Module
Doxygen	A documentation generator used in the Intel® DPDK to generate the API reference.
DPDK	Data Plane Development Kit
DRAM	Dynamic Random Access Memory
EAL	The Environment Abstraction Layer (EAL) provides a generic interface that hides the environment specifics from the applications and libraries. The services expected from the EAL are: development kit loading and launching, core affinity/assignment procedures, system memory allocation/description, PCI bus access, inter-partition communication.
FIFO	First In First Out
GbE	Gigabit Ethernet
HPET	High Precision Event Timer; a hardware timer that provides a precise time reference on x86 platforms.
I/O	Input/Output
lcore	A logical execution unit of the processor, sometimes called a <i>hardware thread</i> .
LAN	Local Area Network



LPM	Longest Prefix Match
master lcore	The execution unit that executes the <code>main()</code> function and that launches other <i>lcores</i> .
mbuf	An mbuf is a data structure used internally to carry messages (mainly network packets). The name is derived from BSD stacks. To understand the concepts of packet buffers or mbuf, refer to <i>TCP/IP Illustrated, Volume 2: The Implementation</i> .
MESI	Modified Exclusive Shared Invalid (CPU cache coherency protocol)
NIC	Network Interface Card
NUMA	Non-uniform Memory Access
PCI	Peripheral Connect Interface
PHY	An abbreviation for the physical layer of the OSI model.
pktmbuf	An <i>mbuf</i> carrying a network packet.
PMD	Poll Mode Driver
RCU	Read-Copy-Update algorithm, an alternative to simple rwlocks.
RSS	Receive Side Scaling
RTE	Run Time Environment. Provides a fast and simple framework for fast packet processing, in a lightweight environment as a Linux* application and using Poll Mode Drivers (PMDs) to increase speed.
Rx	Reception
Slave lcore	Any <i>lcore</i> that is not the <i>master lcore</i> .
Socket	A physical CPU, that includes several <i>cores</i> .
SLA	Service Level Agreement
Target	In the Intel® DPDK, the target is a combination of architecture, machine, executive environment and toolchain. For example: <code>i686-default-linuxapp-gcc</code> .
TLB	Translation Lookaside Buffer
TLS	Thread Local Storage
TSC	Time Stamp Counter
Tx	Transmission
TUN/TAP	TUN and TAP are virtual network kernel devices.
VLAN	Virtual Local Area Network