



DPDK

DATA PLANE DEVELOPMENT KIT

Event Device Drivers

Release 20.05.0

May 26, 2020

1	NXP DPAA Eventdev Driver	2
1.1	Features	2
1.2	Supported DPAA SoCs	2
1.3	Prerequisites	2
1.4	Pre-Installation Configuration	3
1.5	Initialization	3
1.6	Limitations	3
2	NXP DPAA2 Eventdev Driver	4
2.1	Features	4
2.2	Supported DPAA2 SoCs	4
2.3	Prerequisites	4
2.4	Pre-Installation Configuration	5
2.5	Initialization	5
2.6	Enabling logs	5
2.7	Limitations	6
3	Distributed Software Eventdev Poll Mode Driver	7
3.1	Features	7
3.2	Configuration and Options	7
3.3	Limitations	7
4	Software Eventdev Poll Mode Driver	9
4.1	Features	9
4.2	Configuration and Options	9
4.3	Limitations	10
5	OCTEON TX SSOVF Eventdev Driver	12
5.1	Features	12
5.2	Supported OCTEON TX SoCs	12
5.3	Prerequisites	13
5.4	Pre-Installation Configuration	13
5.5	Initialization	13
5.6	Selftest	13
5.7	Enable TIMvf stats	13
5.8	Limitations	14
6	OCTEON TX2 SSO Eventdev Driver	15
6.1	Features	15
6.2	Prerequisites and Compilation procedure	16

6.3	Pre-Installation Configuration	16
7	OPDL Eventdev Poll Mode Driver	18
7.1	Features	18
7.2	Configuration and Options	18
7.3	Limitations	19

The following are a list of event device PMDs, which can be used from an application through the eventdev API.

NXP DPAA EVENTDEV DRIVER

The dpaa eventdev is an implementation of the eventdev API, that provides a wide range of the eventdev features. The eventdev relies on a dpaa based platform to perform event scheduling.

More information can be found at [NXP Official Website](#).

1.1 Features

The DPAA EVENTDEV implements many features in the eventdev API;

- Hardware based event scheduler
- 4 event ports
- 4 event queues
- Parallel flows
- Atomic flows

1.2 Supported DPAA SoCs

- LS1046A/LS1026A
- LS1043A/LS1023A

1.3 Prerequisites

See `../platform/dpaa` for setup information

Currently supported by DPDK:

- NXP SDK **2.0+** or LSDK **18.09+**
- Supported architectures: **arm64 LE**.
- Follow the DPDK Getting Started Guide for Linux to setup the basic DPDK environment.

1.4 Pre-Installation Configuration

1.4.1 Config File Options

The following options can be modified in the `config` file. Please note that enabling debugging options may affect system performance.

- `CONFIG_RTE_LIBRTE_PMD_DPAA_EVENTDEV` (default `y`)
Toggle compilation of the `librte_pmd_dpaa_event` driver.

1.4.2 Driver Compilation

To compile the DPAA EVENTDEV PMD for Linux arm64 gcc target, run the following make command:

```
cd <DPDK-source-directory>
make config T=arm64-dpaa-linux-gcc install
```

1.5 Initialization

The dpaa eventdev is exposed as a vdev device which consists of a set of channels and queues. On EAL initialization, dpaa components will be probed and then vdev device can be created from the application code by

- Invoking `rte_vdev_init("event_dpaa1")` from the application
- Using `--vdev="event_dpaa1"` in the EAL options, which will call `rte_vdev_init()` internally

Example:

```
./your_eventdev_application --vdev="event_dpaa1"
```

- Use dev arg option `disable_intr=1` to disable the interrupt mode

1.6 Limitations

1. DPAA eventdev can not work with DPAA PUSH mode queues configured for ethdev. Please configure `export DPAA_NUM_PUSH_QUEUES=0`

1.6.1 Platform Requirement

DPAA drivers for DPDK can only work on NXP SoCs as listed in the Supported DPAA SoCs.

1.6.2 Port-core Binding

DPAA EVENTDEV driver requires event port 'x' to be used on core 'x'.

NXP DPAA2 EVENTDEV DRIVER

The dpaa2 eventdev is an implementation of the eventdev API, that provides a wide range of the eventdev features. The eventdev relies on a dpaa2 hw to perform event scheduling.

More information can be found at [NXP Official Website](#).

2.1 Features

The DPAA2 EVENTDEV implements many features in the eventdev API;

- Hardware based event scheduler
- 8 event ports
- 8 event queues
- Parallel flows
- Atomic flows

2.2 Supported DPAA2 SoCs

- LX2160A
- LS2084A/LS2044A
- LS2088A/LS2048A
- LS1088A/LS1048A

2.3 Prerequisites

See ../platform/dpaa2 for setup information

Currently supported by DPDK:

- NXP SDK **19.09+**.
- MC Firmware version **10.18.0** and higher.
- Supported architectures: **arm64 LE**.
- Follow the DPDK Getting Started Guide for Linux to setup the basic DPDK environment.

Note: Some part of fslmc bus code (mc flib - object library) routines are dual licensed (BSD & GPLv2).

2.4 Pre-Installation Configuration

2.4.1 Config File Options

The following options can be modified in the `config` file. Please note that enabling debugging options may affect system performance.

- `CONFIG_RTE_LIBRTE_PMD_DPAA2_EVENTDEV` (default `y`)
Toggle compilation of the `lrte_pmd_dpaa2_event` driver.

2.4.2 Driver Compilation

To compile the DPAA2 EVENTDEV PMD for Linux arm64 gcc target, run the following make command:

```
cd <DPDK-source-directory>
make config T=arm64-dpaa-linux-gcc install
```

2.5 Initialization

The dpaa2 eventdev is exposed as a vdev device which consists of a set of dpcon devices and dpci devices. On EAL initialization, dpcon and dpci devices will be probed and then vdev device can be created from the application code by

- Invoking `rte_vdev_init("event_dpaa2")` from the application
- Using `--vdev="event_dpaa2"` in the EAL options, which will call `rte_vdev_init()` internally

Example:

```
./your_eventdev_application --vdev="event_dpaa2"
```

2.6 Enabling logs

For enabling logs, use the following EAL parameter:

```
./your_eventdev_application <EAL args> --log-level=pmd.event.dpaa2,<level>
```

Using `eventdev.dpaa2` as log matching criteria, all Event PMD logs can be enabled which are lower than logging level.

2.7 Limitations

2.7.1 Platform Requirement

DPAA2 drivers for DPDK can only work on NXP SoCs as listed in the `Supported DPAA2 SoCs`.

2.7.2 Port-core binding

DPAA2 EVENTDEV can support only one eventport per core.

DISTRIBUTED SOFTWARE EVENTDEV POLL MODE DRIVER

The distributed software event device is an eventdev driver which distributes the task of scheduling events among all the eventdev ports and the lcore threads using them.

3.1 Features

Queues

- Atomic
- Parallel
- Single-Link

Ports

- Load balanced (for Atomic, Ordered, Parallel queues)
- Single Link (for single-link queues)

3.2 Configuration and Options

The distributed software eventdev is a vdev device, and as such can be created from the application code, or from the EAL command line:

- Call `rte_vdev_init("event_dsw0")` from the application
- Use `--vdev="event_dsw0"` in the EAL options, which will call `rte_vdev_init()` internally

Example:

```
./your_eventdev_application --vdev="event_dsw0"
```

3.3 Limitations

3.3.1 Unattended Ports

The distributed software eventdev uses an internal signaling schema between the ports to achieve load balancing. In order for this to work, the application must perform enqueue and/or dequeue operations on all ports.

Producer-only ports which currently have no events to enqueue should periodically call `rte_event_enqueue_burst()` with a zero-sized burst.

Ports left unattended for longer periods of time will prevent load balancing, and also cause traffic interruptions on the flows which are in the process of being migrated.

3.3.2 Output Buffering

For efficiency reasons, the distributed software eventdev might not send enqueued events immediately to the destination port, but instead store them in an internal buffer in the source port.

In case no more events are enqueued on a port with buffered events, these events will be sent after the application has performed a number of enqueue and/or dequeue operations.

For explicit flushing, an application may call `rte_event_enqueue_burst()` with a zero-sized burst.

3.3.3 Priorities

The distributed software eventdev does not support event priorities.

3.3.4 Ordered Queues

The distributed software eventdev does not support the ordered queue type.

3.3.5 “All Types” Queues

The distributed software eventdev does not support queues of type `RTE_EVENT_QUEUE_CFG_ALL_TYPES`, which allow both atomic, ordered, and parallel events on the same queue.

3.3.6 Dynamic Link/Unlink

The distributed software eventdev does not support calls to `rte_event_port_link()` or `rte_event_port_unlink()` after `rte_event_dev_start()` has been called.

SOFTWARE EVENTDEV POLL MODE DRIVER

The software eventdev is an implementation of the eventdev API, that provides a wide range of the eventdev features. The eventdev relies on a CPU core to perform event scheduling. This PMD can use the service core library to run the scheduling function, allowing an application to utilize the power of service cores to multiplex other work on the same core if required.

4.1 Features

The software eventdev implements many features in the eventdev API;

Queues

- Atomic
- Ordered
- Parallel
- Single-Link

Ports

- Load balanced (for Atomic, Ordered, Parallel queues)
- Single Link (for single-link queues)

Event Priorities

- Each event has a priority, which can be used to provide basic QoS

4.2 Configuration and Options

The software eventdev is a vdev device, and as such can be created from the application code, or from the EAL command line:

- Call `rte_vdev_init("event_sw0")` from the application
- Use `--vdev="event_sw0"` in the EAL options, which will call `rte_vdev_init()` internally

Example:

```
./your_eventdev_application --vdev="event_sw0"
```

4.2.1 Scheduling Quanta

The scheduling quanta sets the number of events that the device attempts to schedule in a single schedule call performed by the service core. Note that is a *hint* only, and that fewer or more events may be scheduled in a given iteration.

The scheduling quanta can be set using a string argument to the vdev create call:

```
--vdev="event_sw0,sched_quanta=64"
```

4.2.2 Credit Quanta

The credit quanta is the number of credits that a port will fetch at a time from the instance's credit pool. Higher numbers will cause less overhead in the atomic credit fetch code, however it also reduces the overall number of credits in the system faster. A balanced number (e.g. 32) ensures that only small numbers of credits are pre-allocated at a time, while also mitigating performance impact of the atomics.

Experimentation with higher values may provide minor performance improvements, at the cost of the whole system having less credits. On the other hand, reducing the quanta may cause measurable performance impact but provide the system with a higher number of credits at all times.

A value of 32 seems a good balance however your specific application may benefit from a higher or reduced quanta size, experimentation is required to verify possible gains.

```
--vdev="event_sw0,credit_quanta=64"
```

4.3 Limitations

The software eventdev implementation has a few limitations. The reason for these limitations is usually that the performance impact of supporting the feature would be significant.

4.3.1 "All Types" Queues

The software eventdev does not support creating queues that handle all types of traffic. An eventdev with this capability allows enqueueing Atomic, Ordered and Parallel traffic to the same queue, but scheduling each of them appropriately.

The reason to not allow Atomic, Ordered and Parallel event types in the same queue is that it causes excessive branching in the code to enqueue packets to the queue, causing a significant performance impact.

The `RTE_EVENT_DEV_CAP_QUEUE_ALL_TYPES` flag is not set in the `event_dev_cap` field of the `rte_event_dev_info` struct for the software eventdev.

4.3.2 Distributed Scheduler

The software eventdev is a centralized scheduler, requiring a service core to perform the required event distribution. This is not really a limitation but rather a design decision.

The `RTE_EVENT_DEV_CAP_DISTRIBUTED_SCHED` flag is not set in the `event_dev_cap` field of the `rte_event_dev_info` struct for the software eventdev.

4.3.3 Dequeue Timeout

The `eventdev` API supports a timeout when dequeuing packets using the `rte_event_dequeue_burst` function. This allows a core to wait for an event to arrive, or until `timeout` number of ticks have passed. Timeout ticks is not supported by the software `eventdev` for performance reasons.

OCTEON TX SSOVF EVENTDEV DRIVER

The OCTEON TX SSOVF PMD (`librte_pmd_octeontx_ssovf`) provides poll mode eventdev driver support for the inbuilt event device found in the **Cavium OCTEON TX** SoC family as well as their virtual functions (VF) in SR-IOV context.

More information can be found at [Cavium, Inc Official Website](#).

5.1 Features

Features of the OCTEON TX SSOVF PMD are:

- 64 Event queues
- 32 Event ports
- HW event scheduler
- Supports 1M flows per event queue
- Flow based event pipelining
- Flow pinning support in flow based event pipelining
- Queue based event pipelining
- Supports ATOMIC, ORDERED, PARALLEL schedule types per flow
- Event scheduling QoS based on event queue priority
- Open system with configurable amount of outstanding events
- HW accelerated dequeue timeout support to enable power management
- SR-IOV VF
- HW managed event timers support through TIMVF, with high precision and time granularity of 1us.
- Up to 64 event timer adapters.

5.2 Supported OCTEON TX SoCs

- CN83xx

5.3 Prerequisites

See `./platform/octeontx` for setup information.

5.4 Pre-Installation Configuration

5.4.1 Config File Options

The following options can be modified in the `config` file. Please note that enabling debugging options may affect system performance.

- `CONFIG_RTE_LIBRTE_PMD_OCTEONTX_SSOVF` (default `y`)
Toggle compilation of the `librte_pmd_octeontx_ssovf` driver.

5.4.2 Driver Compilation

To compile the OCTEON TX SSOVF PMD for Linux arm64 gcc target, run the following make command:

```
cd <DPDK-source-directory>
make config T=arm64-thunderx-linux-gcc install
```

5.5 Initialization

The OCTEON TX eventdev is exposed as a vdev device which consists of a set of SSO group and workslot PCIe VF devices. On EAL initialization, SSO PCIe VF devices will be probed and then the vdev device can be created from the application code, or from the EAL command line based on the number of probed/bound SSO PCIe VF device to DPDK by

- Invoking `rte_vdev_init("event_octeontx")` from the application
- Using `--vdev="event_octeontx"` in the EAL options, which will call `rte_vdev_init()` internally

Example:

```
./your_eventdev_application --vdev="event_octeontx"
```

5.6 Selftest

The functionality of OCTEON TX eventdev can be verified using this option, various unit and functional tests are run to verify the sanity. The tests are run once the vdev creation is successfully complete.

```
--vdev="event_octeontx,selftest=1"
```

5.7 Enable TIMvf stats

TIMvf stats can be enabled by using this option, by default the stats are disabled.


```
--vdev="event_octeontx,timvf_stats=1"
```

5.8 Limitations

5.8.1 Burst mode support

Burst mode is not supported. Dequeue and Enqueue functions accepts only single event at a time.

5.8.2 Rx adapter support

When eth_octeontx is used as Rx adapter event schedule type RTE_SCHED_TYPE_PARALLEL is not supported.

5.8.3 Event timer adapter support

When timvf is used as Event timer adapter the clock source mapping is as follows:

```
RTE_EVENT_TIMER_ADAPTER_CPU_CLK = TIM_CLK_SRC_SCLK  
RTE_EVENT_TIMER_ADAPTER_EXT_CLK0 = TIM_CLK_SRC_GPIO  
RTE_EVENT_TIMER_ADAPTER_EXT_CLK1 = TIM_CLK_SRC_GTI  
RTE_EVENT_TIMER_ADAPTER_EXT_CLK2 = TIM_CLK_SRC_PTP
```

When timvf is used as Event timer adapter event schedule type RTE_SCHED_TYPE_PARALLEL is not supported.

5.8.4 Max mempool size

Max mempool size when using OCTEON TX Eventdev (SSO) should be limited to 128K. When running dpdk-test-eventdev on OCTEON TX the application can limit the number of mbufs by using the option `--pool_sz 131072`

OCTEON TX2 SSO EVENTDEV DRIVER

The OCTEON TX2 SSO PMD (`librte_pmd_octeontx2_event`) provides poll mode eventdev driver support for the inbuilt event device found in the **Marvell OCTEON TX2** SoC family.

More information about OCTEON TX2 SoC can be found at [Marvell Official Website](#).

6.1 Features

Features of the OCTEON TX2 SSO PMD are:

- 256 Event queues
- 26 (dual) and 52 (single) Event ports
- HW event scheduler
- Supports 1M flows per event queue
- Flow based event pipelining
- Flow pinning support in flow based event pipelining
- Queue based event pipelining
- Supports ATOMIC, ORDERED, PARALLEL schedule types per flow
- Event scheduling QoS based on event queue priority
- Open system with configurable amount of outstanding events limited only by DRAM
- HW accelerated dequeue timeout support to enable power management
- HW managed event timers support through TIM, with high precision and time granularity of 2.5us.
- Up to 256 TIM rings aka event timer adapters.
- Up to 8 rings traversed in parallel.
- HW managed packets enqueued from ethdev to eventdev exposed through event eth RX adapter.
- N:1 ethernet device Rx queue to Event queue mapping.
- Lockfree Tx from event eth Tx adapter using `DEV_TX_OFFLOAD_MT_LOCKFREE` capability while maintaining receive packet order.
- Full Rx/Tx offload support defined through ethdev queue config.

6.2 Prerequisites and Compilation procedure

See `../platform/octeontx2` for setup information.

6.3 Pre-Installation Configuration

6.3.1 Compile time Config Options

The following option can be modified in the `config` file.

- `CONFIG_RTE_LIBRTE_PMD_OCTEONTX2_EVENTDEV` (default `y`)
Toggle compilation of the `librte_pmd_octeontx2_event` driver.

6.3.2 Runtime Config Options

- Maximum number of in-flight events (default 8192)

In **Marvell OCTEON TX2** the max number of in-flight events are only limited by DRAM size, the `xae_cnt` devargs parameter is introduced to provide upper limit for in-flight events. For example:

```
-w 0002:0e:00.0,xae_cnt=16384
```

- Force legacy mode

The `single_ws` devargs parameter is introduced to force legacy mode i.e single workslot mode in SSO and disable the default dual workslot mode. For example:

```
-w 0002:0e:00.0,single_ws=1
```

- Event Group QoS support

SSO GGRPs i.e. queue uses DRAM & SRAM buffers to hold in-flight events. By default the buffers are assigned to the SSO GGRPs to satisfy minimum HW requirements. SSO is free to assign the remaining buffers to GGRPs based on a preconfigured threshold. We can control the QoS of SSO GGRP by modifying the above mentioned thresholds. GGRPs that have higher importance can be assigned higher thresholds than the rest. The dictionary format is as follows `[Qx-XAQ-TAQ-IAQ][Qz-XAQ-TAQ-IAQ]` expressed in percentages, 0 represents default. For example:

```
-w 0002:0e:00.0,qos=[1-50-50-50]
```

- Selftest

The functionality of OCTEON TX2 eventdev can be verified using this option, various unit and functional tests are run to verify the sanity. The tests are run once the vdev creation is successfully complete. For example:

```
-w 0002:0e:00.0,selftest=1
```

- TIM disable NPA

By default chunks are allocated from NPA then TIM can automatically free them when traversing the list of chunks. The `tim_disable_npa` devargs parameter disables NPA and uses software mempool to manage chunks For example:

```
-w 0002:0e:00.0,tim_disable_npa=1
```

- TIM modify chunk slots

The `tim_chnk_slots` devargs can be used to modify number of chunk slots. Chunks are used to store event timers, a chunk can be visualised as an array where the last element points to the next chunk and rest of them are used to store events. TIM traverses the list of chunks and enqueues the event timers to SSO. The default value is 255 and the max value is 4095. For example:

```
-w 0002:0e:00.0,tim_chnk_slots=1023
```

- TIM enable arm/cancel statistics

The `tim_stats_ena` devargs can be used to enable arm and cancel stats of event timer adapter. For example:

```
-w 0002:0e:00.0,tim_stats_ena=1
```

- TIM limit max rings reserved

The `tim_rings_lmt` devargs can be used to limit the max number of TIM rings i.e. event timer adapter reserved on probe. Since, TIM rings are HW resources we can avoid starving other applications by not grabbing all the rings. For example:

```
-w 0002:0e:00.0,tim_rings_lmt=5
```

- TIM ring control internal parameters

When using multiple TIM rings the `tim_ring_ctl` devargs can be used to control each TIM rings internal parameters uniquely. The following dict format is expected `[ring-chnk_slots-disable_npa-stats_ena]`. 0 represents default values. For Example:

```
-w 0002:0e:00.0,tim_ring_ctl=[2-1023-1-0]
```

- Lock NPA contexts in NDC

Lock NPA aura and pool contexts in NDC cache. The device args take hexadecimal bitmask where each bit represent the corresponding aura/pool id.

For example:

```
-w 0002:0e:00.0,npa_lock_mask=0xf
```

6.3.3 Debugging Options

Table 6.1: OCTEON TX2 event device debug options

#	Component	EAL log command
1	SSO	<code>-log-level='pmd.event.octeontx2,8'</code>
2	TIM	<code>-log-level='pmd.event.octeontx2.timer,8'</code>

OPDL EVENTDEV POLL MODE DRIVER

The OPDL (Ordered Packet Distribution Library) eventdev is a specific implementation of the eventdev API. It is particularly suited to packet processing workloads that have high throughput and low latency requirements. All packets follow the same path through the device. The order in which packets follow is determined by the order in which queues are set up. Events are left on the ring until they are transmitted. As a result packets do not go out of order.

7.1 Features

The OPDL eventdev implements a subset of features of the eventdev API;

Queues

- Atomic
- Ordered (Parallel is supported as parallel is a subset of Ordered)
- Single-Link

Ports

- Load balanced (for Atomic, Ordered, Parallel queues)
- Single Link (for single-link queues)

7.2 Configuration and Options

The software eventdev is a vdev device, and as such can be created from the application code, or from the EAL command line:

- Call `rte_vdev_init("event_opdl0")` from the application
- Use `--vdev="event_opdl0"` in the EAL options, which will call `rte_vdev_init()` internally

Example:

```
./your_eventdev_application --vdev="event_opdl0"
```

7.2.1 Single Port Queue

It is possible to create a Single Port Queue `RTE_EVENT_QUEUE_CFG_SINGLE_LINK`. Packets dequeued from this queue do not need to be re-enqueued (as is the case with an ordered queue). The

purpose of this queue is to allow for asynchronous handling of packets in the middle of a pipeline. Ordered queues in the middle of a pipeline cannot delete packets.

7.2.2 Queue Dependencies

As stated the order in which packets travel through queues is static in nature. They go through the queues in the order the queues are setup at initialisation `rte_event_queue_setup()`. For example if an application sets up 3 queues, Q0, Q1, Q2 and has 3 associated ports P0, P1, P2 and P3 then packets must be

- Enqueued onto Q0 (typically through P0), then
- Dequeued from Q0 (typically through P1), then
- Enqueued onto Q1 (also through P1), then
- Dequeued from Q2 (typically through P2), then
- Enqueued onto Q3 (also through P2), then
- Dequeued from Q3 (typically through P3) and then transmitted on the relevant eth port

7.3 Limitations

The opdl implementation has a number of limitations. These limitations are due to the static nature of the underlying queues. It is because of this that the implementation can achieve such high throughput and low latency

The following list is a comprehensive outline of the what is supported and the limitations / restrictions imposed by the opdl pmd

- The order in which packets moved between queues is static and fixed (dynamic scheduling is not supported).
- NEW, RELEASE are not explicitly supported. RX (first enqueue) implicitly adds NEW event types, and TX (last dequeue) implicitly does RELEASE event types.
- All packets follow the same path through device queues.
- Flows within queues are NOT supported.
- Event priority is NOT supported.
- Once the device is stopped all inflight events are lost. Applications should clear all inflight events before stopping it.
- Each port can only be associated with one queue.
- Each queue can have multiple ports associated with it.
- Each worker core has to dequeue the maximum burst size for that port.
- For performance, the `rte_event flow_id` should not be updated once packetis enqueued on RX.

7.3.1 Validation & Statistics

Validation can be turned on through a command line parameter

```
--vdev="event_opdl0,do_validation=1,self_test=1"
```

If validation is turned on every packet (as opposed to just the first in each burst), is validated to have come from the right queue. Statistics are also produced in this mode. The statistics are available through the eventdev xstats API. Statistics are per port as follows:

- claim_pkts_requested
- claim_pkts_granted
- claim_non_empty
- claim_empty
- total_cycles