

**NAME**

**iflib** - Network Interface Driver Framework

**SYNOPSIS**

**device pci**  
**device iflib**

**DESCRIPTION**

**iflib** is a framework for network interface drivers for FreeBSD. It is designed to remove a large amount of the boilerplate that is often needed for modern network interface devices, allowing driver authors to focus on the specific code needed for their hardware. This allows for a shared set of sysctl(8) names, rather than each driver naming them individually.

**SYSCTL VARIABLES**

These variables must be set before loading the driver, either via loader.conf(5) or through the use of kenv(1). They are all prefixed by *dev.X.Y.iflib*, where *X* is the driver name, and *Y* is the instance number.

*override\_nrxds*

Override the number of RX descriptors for each queue. The value is a comma separated list of positive integers. Some drivers only use a single value, but others may use more. These numbers must be powers of two, and zero means to use the default. Individual drivers may have additional restrictions on allowable values. Defaults to all zeros.

*override\_ntxds*

Override the number of TX descriptors for each queue. The value is a comma separated list of positive integers. Some drivers only use a single value, but others may use more. These numbers must be powers of two, and zero means to use the default. Individual drivers may have additional restrictions on allowable values. Defaults to all zeros.

*override\_qs\_enable*

When set, allows the number of transmit and receive queues to be different. If not set, the lower of the number of TX or RX queues will be used for both.

*override\_nrxqs*

Set the number of RX queues. If zero, the number of RX queues is derived from the number of cores on the socket connected to the controller. Defaults to 0.

*override\_ntxqs*

Set the number of TX queues. If zero, the number of TX queues is derived from the number of

cores on the socket connected to the controller.

*disable\_msix*

Disables MSI-X interrupts for the device.

*core\_offset*

Specifies a starting core offset to assign queues to. If the value is unspecified or 65535, cores are assigned sequentially across controllers.

*separate\_txrx*

Requests that RX and TX queues not be paired on the same core. If this is zero or not set, an RX and TX queue pair will be assigned to each core. When set to a non-zero value, TX queues are assigned to cores following the last RX queue.

These sysctl(8) variables can be changed at any time:

*tx\_abdicate*

Controls how the transmit ring is serviced. If set to zero, when a frame is submitted to the transmission ring, the same task that is submitting it will service the ring unless there's already a task servicing the TX ring. This ensures that whenever there is a pending transmission, the transmit ring is being serviced. This results in higher transmit throughput. If set to a non-zero value, task returns immediately and the transmit ring is serviced by a different task. This returns control to the caller faster and under high receive load, may result in fewer dropped RX frames.

*rx\_budget*

Sets the maximum number of frames to be received at a time. Zero (the default) indicates the default (currently 16) should be used.

There are also some global sysctls which can change behaviour for all drivers, and may be changed at any time.

*net.iflib.min\_tx\_latency*

If this is set to a non-zero value, iflib will avoid any attempt to combine multiple transmits, and notify the hardware as quickly as possible of new descriptors. This will lower the maximum throughput, but will also lower transmit latency.

*net.iflib.no\_tx\_batch*

Some NICs allow processing completed transmit descriptors in batches. Doing so usually increases the transmit throughput by reducing the number of transmit interrupts. Setting this to a non-zero value will disable the use of this feature.

These sysctl(8) variables are read-only:

*driver\_version*

A string indicating the internal version of the driver.

There are a number of queue state sysctl(8) variables as well:

*txqZ* The following are repeated for each transmit queue, where Z is the transmit queue instance number:

*r\_abdications*

Number of consumer abdications in the MP ring for this queue. An abdication occurs on every ring submission when tx\_abdicate is true.

*r\_restarts*

Number of consumer restarts in the MP ring for this queue. A restart occurs when an attempt to drain a non-empty ring fails, and the ring is already in the STALLED state.

*r\_stalls*

Number of consumer stalls in the MP ring for this queue. A stall occurs when an attempt to drain a non-empty ring fails.

*r\_starts*

Number of normal consumer starts in the MP ring for this queue. A start occurs when the MP ring transitions from IDLE to BUSY.

*r\_drops*

Number of drops in the MP ring for this queue. A drop occurs when there is an attempt to add an entry to an MP ring with no available space.

*r\_enqueues*

Number of entries which have been enqueued to the MP ring for this queue.

*ring\_state*

MP (soft) ring state. This provides a snapshot of the current MP ring state, including the producer head and tail indexes, the consumer index, and the state. The state is one of "IDLE", "BUSY", "STALLED", or "ABDICATED".

*txq\_cleaned*

The number of transmit descriptors which have been reclaimed. Total cleaned.

*txq\_processed*

The number of transmit descriptors which have been processed, but may not yet have been reclaimed.

*txq\_in\_use*

Descriptors which have been added to the transmit queue, but have not yet been cleaned. This value will include both untransmitted descriptors as well as descriptors which have been processed.

*txq\_cidx\_processed*

The transmit queue consumer index of the next descriptor to process.

*txq\_cidx*

The transmit queue consumer index of the oldest descriptor to reclaim.

*txq\_pidx*

The transmit queue producer index where the next descriptor to transmit will be inserted.

*no\_tx\_dma\_setup*

Number of times DMA mapping a transmit mbuf failed for reasons other than EFBIG.

*txd\_encap\_efbig*

Number of times DMA mapping a transmit mbuf failed due to requiring too many segments.

*tx\_map\_failed*

Number of times DMA mapping a transmit mbuf failed for any reason (sum of *no\_tx\_dma\_setup* and *txd\_encap\_efbig*)

*no\_desc\_avail*

Number of times a descriptor couldn't be added to the transmit ring because the transmit ring was full.

*mbuf\_defrag\_failed*

Number of times both *m\_collapse(9)* and *m\_defrag(9)* failed after an EFBIG error result from DMA mapping a transmit mbuf.

*m\_pullups*

Number of times *m\_pullup(9)* was called attempting to parse a header.

*mbuf\_defrag*

Number of times `m_defrag(9)` was called.

*rxqZ* The following are repeated for each receive queue, where *Z* is the receive queue instance number:

*rxq\_flo.credits*

Credits currently available in the receive ring.

*rxq\_flo.cidx*

Current receive ring consumer index.

*rxq\_flo.pidx*

Current receive ring producer index.

Additional OIDs useful for driver and iflib development are exposed when the `INVARIANTS` and/or `WITNESS` options are enabled in the kernel.

## SEE ALSO

`iflib(9)`

## HISTORY

This framework was introduced in FreeBSD 11.0.