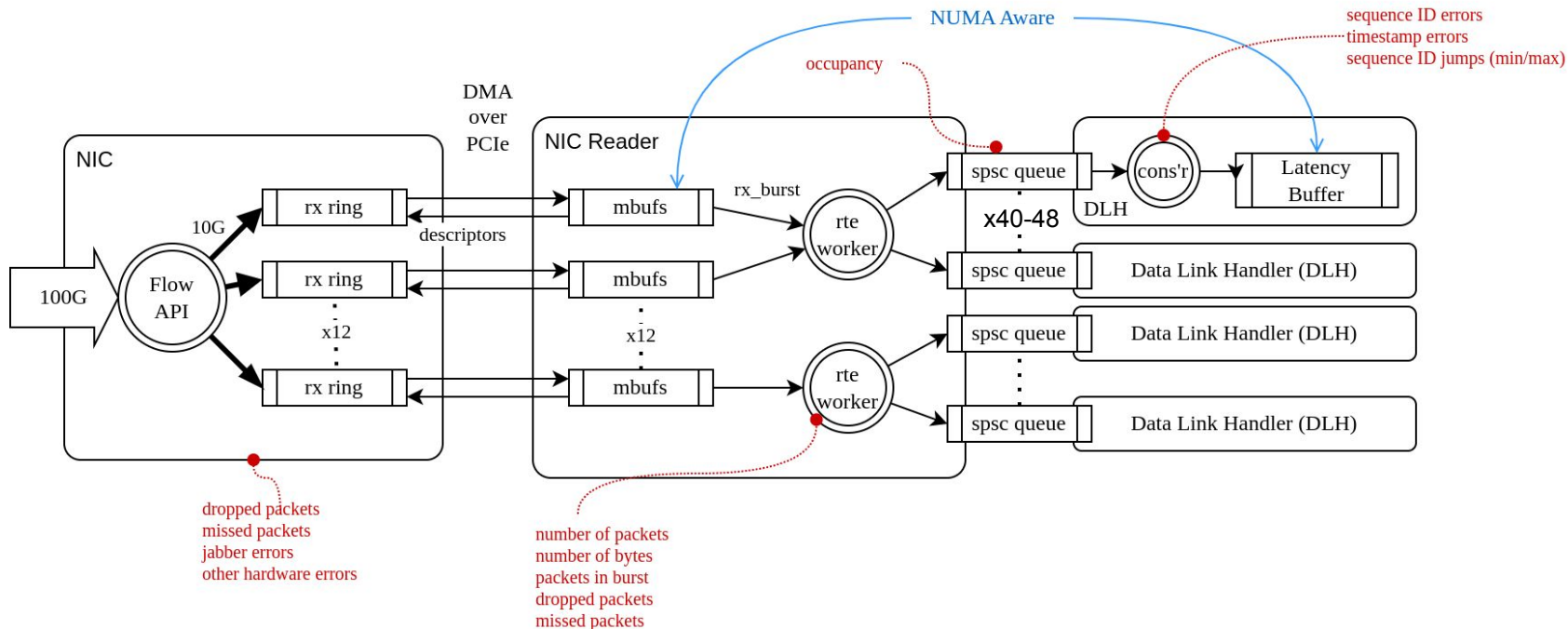# Ethernet readout optimization

Roland Sipos - for the DUNE DAQ
CERN

DUNE DAQ CoreSW Meeting
3rd April 2024

# Data movement - original

# Callback feature in 4.4.0

- Readoutlibs' DataMoveCallbackRegistry for intra-process readout modules
  - Similar API and structure as IOManager, but holds std::functions with specific signature

```
std::shared_ptr<std::function<void(DataType&&)>> m_callback;
```

- DataLinkHandler modules advertise/register a callback function, that moves the data to the preprocessing pipeline, and then writing to the LatencyBuffer.
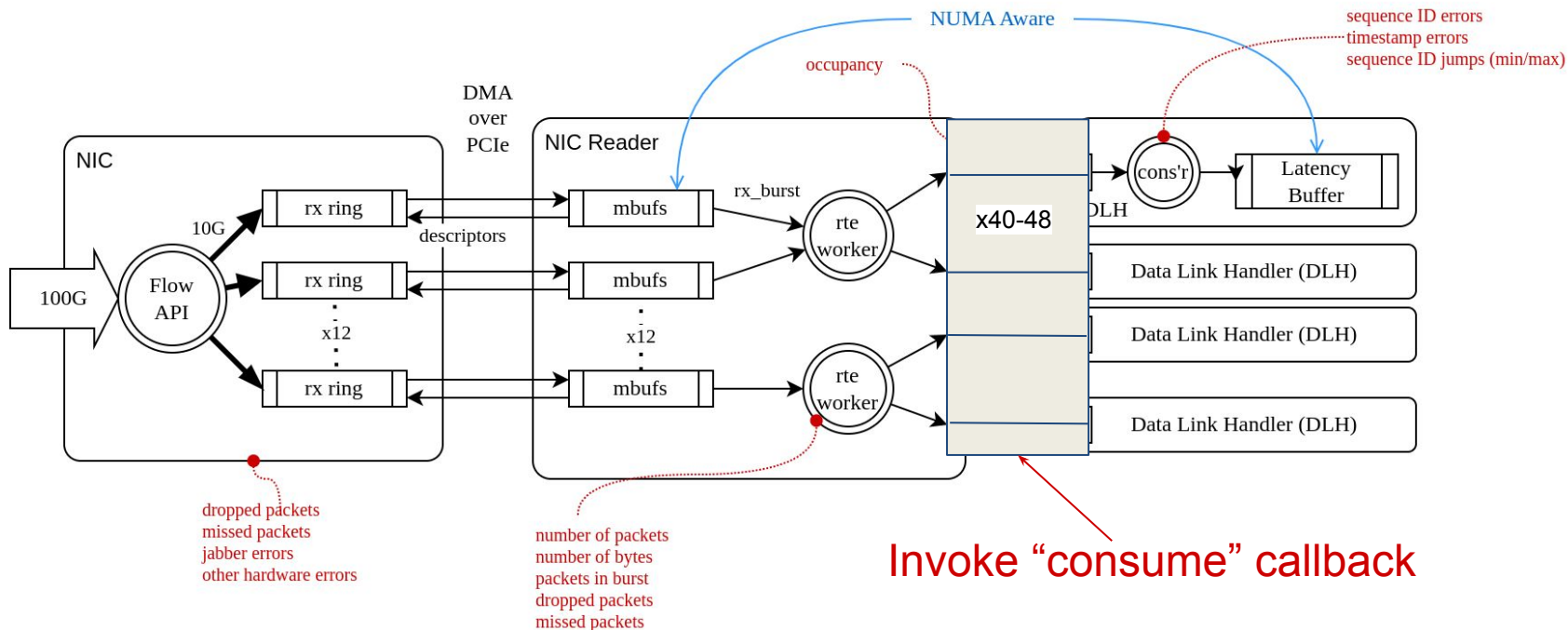  - Registry:

```
template<typename DataType>
void register_callback(const std::string& id, std::function<void(DataType&&)> callback);
```
  - ReadoutModel:

```
// Register callback
auto dmcbr = DataMoveCallbackRegistry::get();
dmcbr->register_callback<RDT>(m_raw_data_receiver_connection_name, m_consume_callback);
```

- I/O card receiver modules (NICReceiver) could look-up which callback functions to invoke to move the data from the producer thread's context directly to the DLHs':

```
// Getting DataMoveCBRegistry
auto dmcbr = readoutlibs::DataMoveCallbackRegistry::get();
m_sink_callback = dmcbr->get_callback<TargetPayloadType>(inherited::m_sink_name);
```

# Data movement - Callbacks
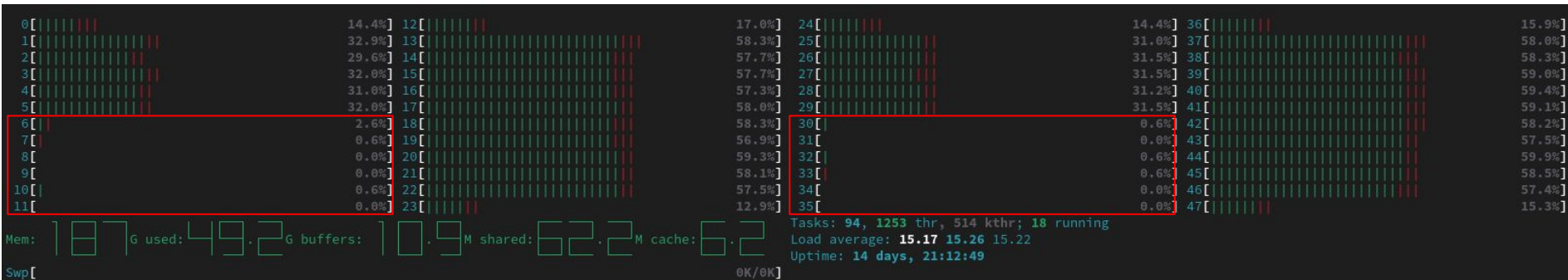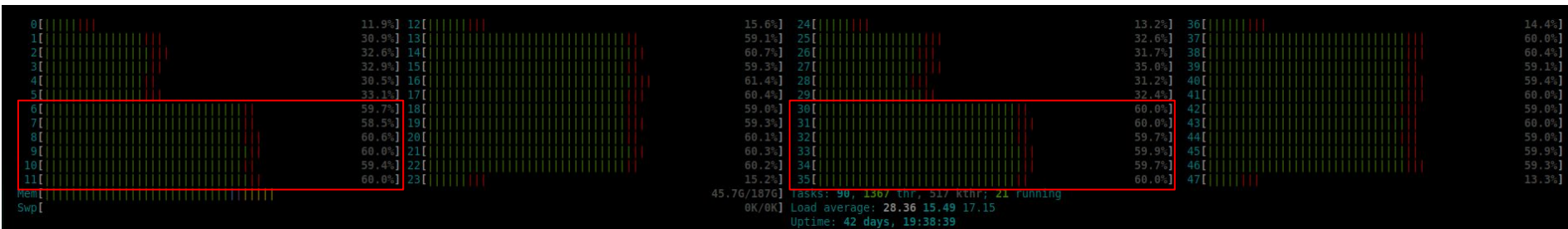
# Callback feature comments

- Optional path, inferred from connection name to be enabled
  - Necessary, as we don't want to break other datatypes (trigger, TPs, etc.)
  - Ugly, "hacky", I know: If the *raw_input* connection has a "*cb_*" prefix

- Would it result with any perf. gain for TPs to TPDataLinkHandlers!?
  - Maybe, but probably callbacks' invocation needs parallel async techniques either on the producer or consume callback owner class. (Tricky, not SPSC.)

- Forward port to "develop"?
  - Feedback for appfwk review first? Or that's too far ahead, seeing the gain?
  - Improvements and common API/utility promotion are possible for sure…
    - Generalize for network: RPC variant?
    - [GenericCallback](#) instead of move signature only functions?

# Resource utilization reduction

- Removes one full copy (intermediate buffering) of the raw data stream:
  - Eliminates consumer threads: no resource utilization due to polling
  - Eliminates context switching between prod/consumers
  - LLC and first level cache miss reduction due to first two points

- Memory bandwidth utilization percentage difference: **~37% less**!
  - Can be utilized as natural headroom for SNB recording!

- CPU utilization reduction: **12 vCores (Phys & their HT) x ~60%**
  - This one is on Intel Xeon Gold 5118 @ 2.3 GHz (Skylake)
  - Heavily depends on CPU!
  - But even on newer generations we see ~10 vCores x ~35% reduction (CPUs with more cores and higher clocks)

# CPU utilization

- **Non negligible on CPUs with lower core counts and clocks!**
- On Intel Xeon Gold 5118 @ 2.3 GHz (Skylake) -> before/after callbacks
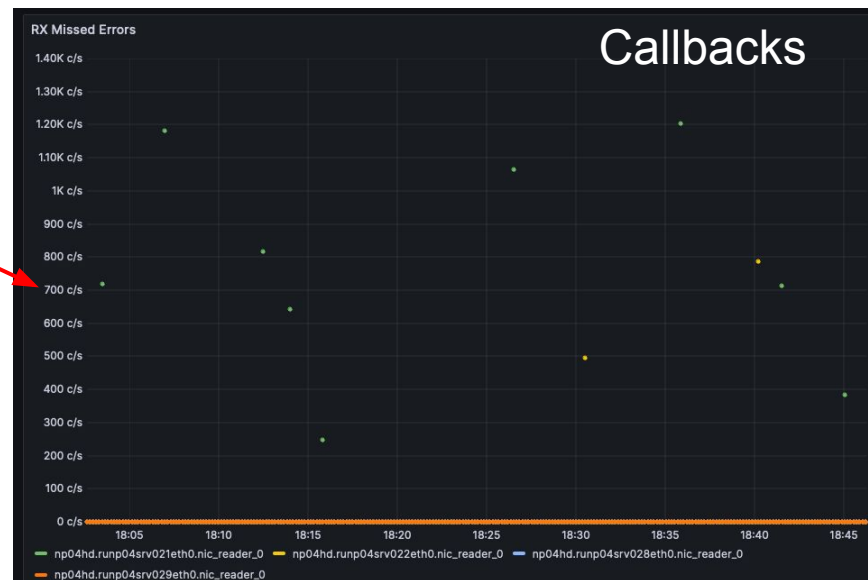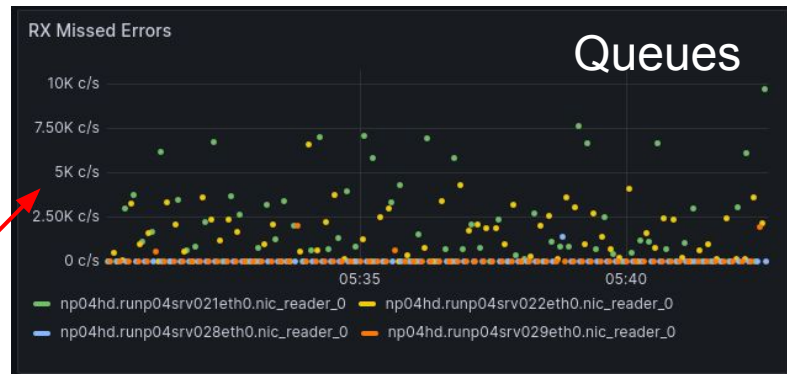
# Missed packets

When the processing (copying out data from the DMA buffers) is lagging. Data is missed!

- Queues: Order of thousands missed packets in every few seconds
- Callbacks: Occasional spikes to a few hundred in every few minutes

**Where are these coming from?**
We knew, that we are sensitive to interrupts, but callbacks only won't help with that, as this comes from the producer thread (in DPDK).

Something is still wrong…

# Alessandro's evaluation and testing

- Alessandro did extensive testing on the AMD server to tweak interface and DMA memory pool configurations in DPDK

- Main find: The interface polling function ("rx burst") returns with **full size** very frequently. Our DPDK/iface configuration needs a revisit.

# DPDK/iface config revisit

New "golden" config

```
"ip_addr": "10.73.139.16",
"mac_addr": "6c:fe:54:59:98:a0",
"parameters": {
    "burst_size": 2048,
    "lcore_sleep_us": 10,
    "mbuf_cache_size": 0,
    "mtu": 9000,
    "num_mbufs": 16384,
    "promiscuous_mode": false,
    "rx_ring_size": 4096,
    "tx_ring_size": 4096,
    "with_flow_control": true
},
"pci_addr": "0000:5e:00.0"
```

- **Burst size**: 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 <u>global loop policies</u>.
  - **Previous value: 256**
- **Lcore_sleep_us**: Microseconds to <u>sleep</u> if there are no full buffers in the burst.
- **Mbuf_cache_size**: Complicated, we didn't understand this properly...
  For details have a look on the mempool <u>API docs</u>
  - **Previous value: 256**
- **Num_mbufs**: The number of mbufs (DPDK DMA buffer segments) to allocate.
  - **Previous value: 8191**
- **RX/TX ring size**: Number of RX/TX descriptors.
  - That is the mechanism you communicate to the NIC hardware - e.g., in the
    case of RX, for passing pointers to empty buffers that you pass to NIC and
    getting pointers for filled buffers that NIC returns.
  - Usually this parameter has hardware limits (max. num. of descriptors).
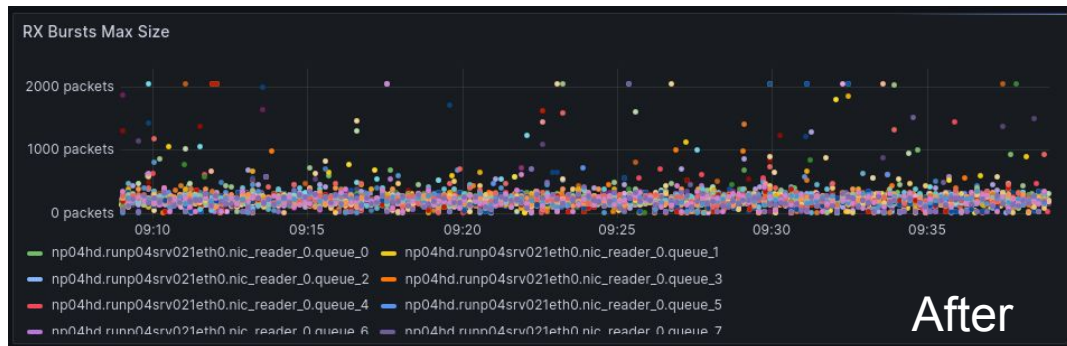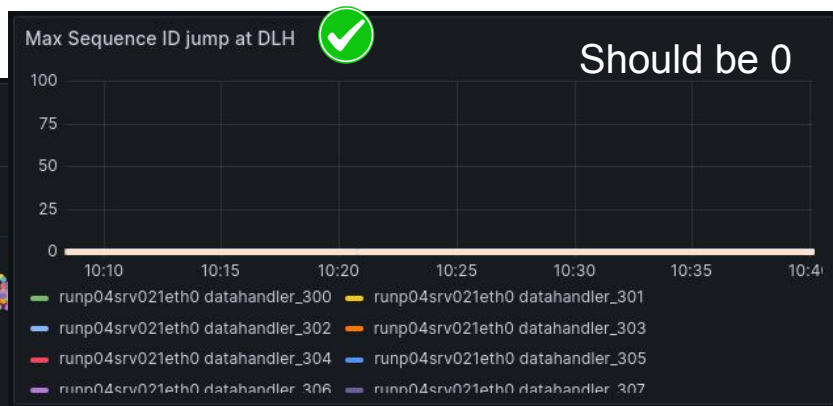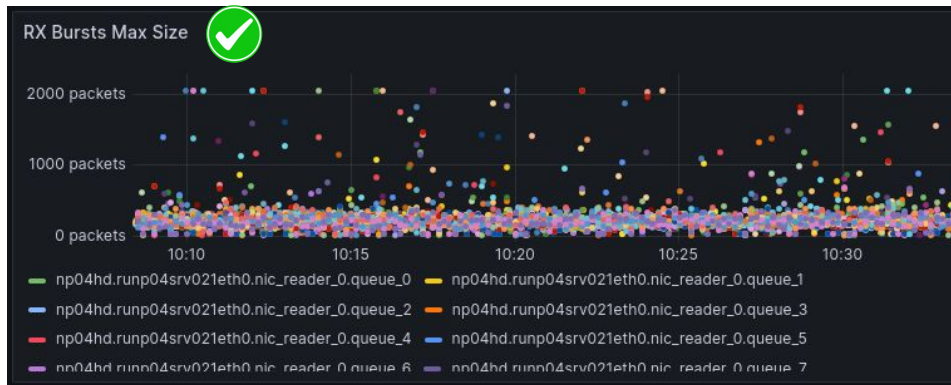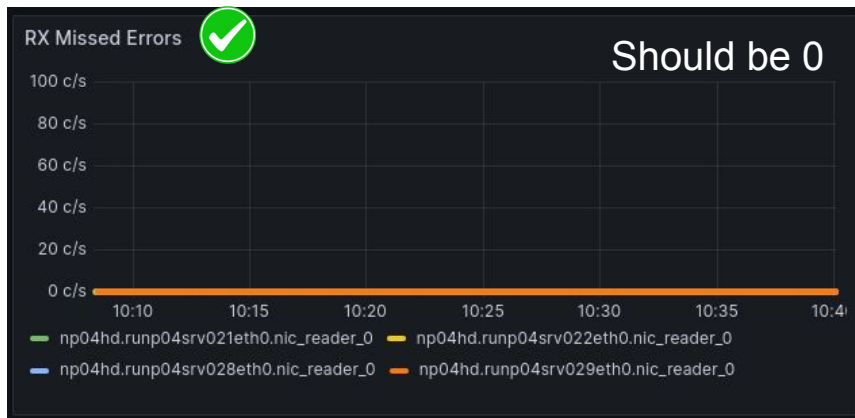  - **Previous value: 1024**

DUNE

# Iface config tuning

- Strategy for our traffic characteristics is to poll (sensitive to interrupts) as many packets as we can in one go **(burst)**. Basically to reduce the impact of interrupts and kernel events.

  - Increase num. descriptors
  - Allocate more mbufs
  - Increase burst size

**Take note of the removed linear roof: Reduced full bursts**

# New performance

# Outlook

- There are still some very rare occurrences of missed packets on 021. (That machine has some kernel parameters under testing, might be related!)

- Kernel isolation should still be pursued, as my suspicion is that will still help the AMD due to the higher I/O die latency on the processor

- This optimization builds up confidence in moving on to scalability aspects
    - 2 NICs in the same server

- Also, SNB recording evaluation and improvements can continue

- Topology evaluation opens up again
    - Symmetric or asymmetric configuration of I/O devices and processing threads

# Summary

- Ethernet readout with TPG, resource headroom available (but no SNB recording) on a dual-socket server with CPUs that are from Q3 of 2017 (mid-range Skylake)

- Great work and sprint from the perf-testing activity team for fddaq-4.4.0!
  - Building up experience and confidence in understanding system characteristics and requirements

- The path ahead is clear what we aim to test in our upcoming sprint
  - Many other things to test and sort out, but we managed to breach through a showstopper problem

- Questions and comments are welcome!