# Upstream DAQ Readout software integration to appfwk

Roland Sipos
CERN

DUNE DAQ/SC Meeting
8th October 2020

# Outline of the talk

- Upstream DAQ Readout system
- ProtoDUNE versions
- OnHost BR
- Modularization
- Integration to appfwk
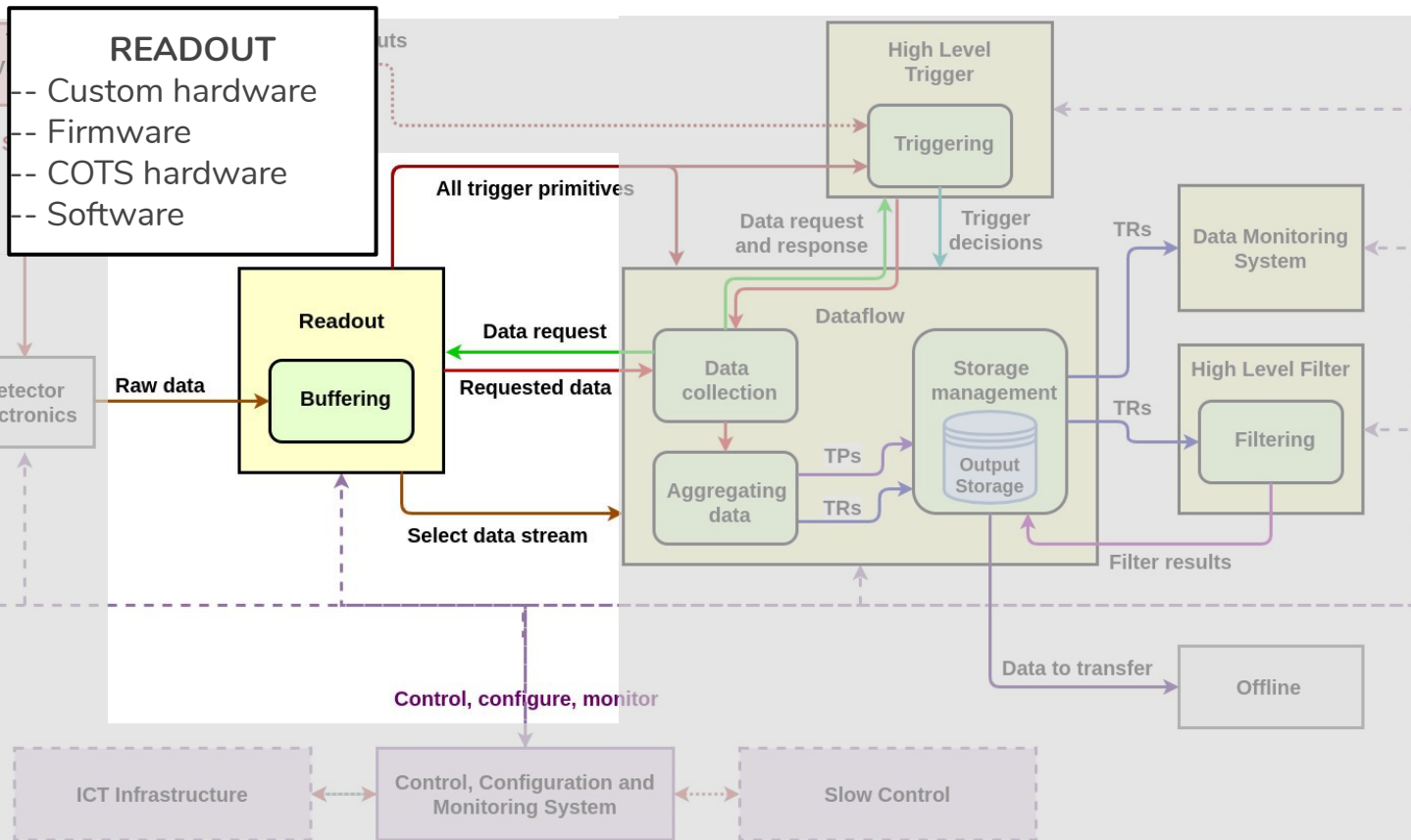- Outlook and Summary

# The goal of this talk

Based on the experience with the FELIX readout in ProtoDUNE-SP, I would like to make some remarks on the existing system's problems and what lead to those.

Then I will present the new approach in order to try solving those.

Finally, the first experiences of the integration of the readout DAQModules to the DUNE DAQ appfwk.
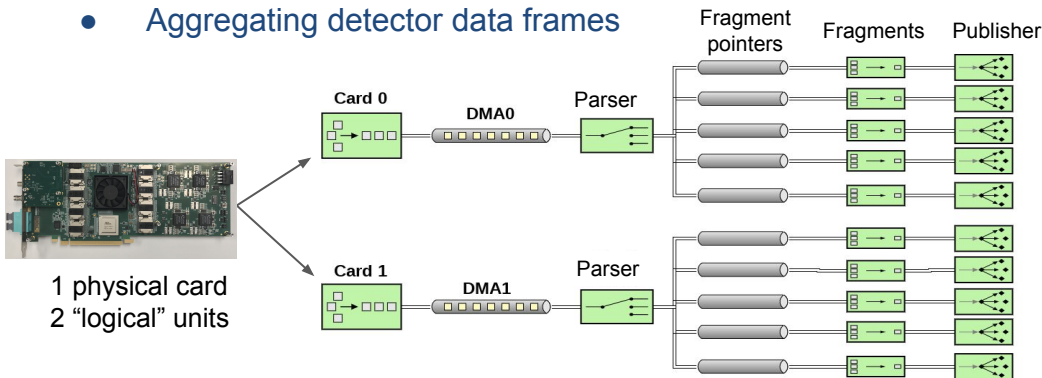
# Upstream DAQ
## Readout System functions



**READOUT**
- Custom hardware
- Firmware
- COTS hardware
- Software

# FELIX in ProtoDUNE



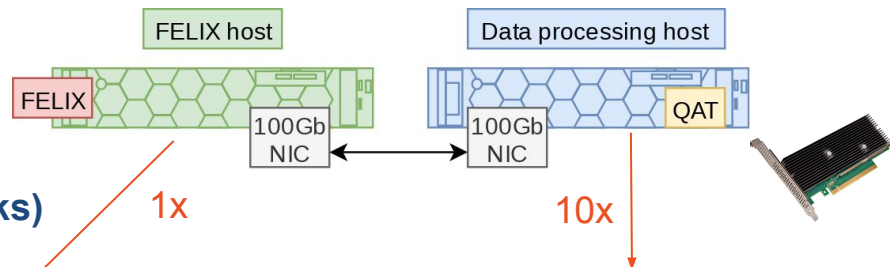**1 FELIX card handles a single APA (10x ~1GB/s links)**

1x          10x

Modest modifications on FPGA gateware
for lower memory I/O rate:
- Increasing DMA payload size
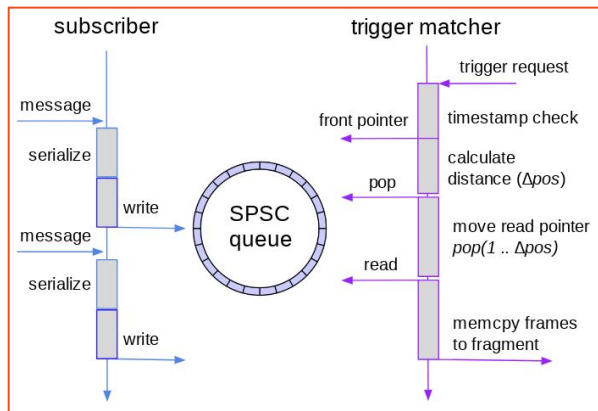- Aggregating detector data frames

1 physical card
2 "logical" units

Data routing software is customized:
- Scatter-gather: collects pointers for detector data fragments
- Single copy pipeline: serialization to user buffer
- Data published on Infiniband over Ethernet
- **Felix-star uses a similar principle!**

BoardReader process
- Subscribes to single link & buffers data
- Extracts data fragments for triggers
- Reorders data (AVX2 & 512)
- Hardware accelerated compression
  - Intel® QuickAssist (QAT)

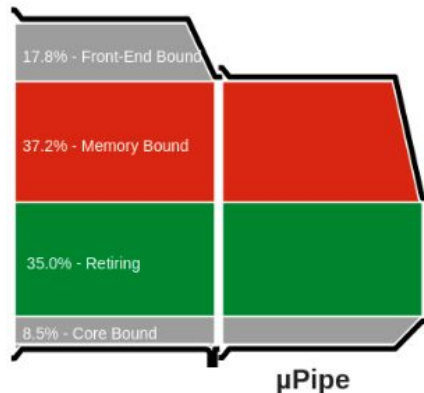# OnHost BoardReader



FELIX & processing host

FELIX    AVX    QAT

**Merged FELIX data routing software, with data selection (trigger-matching) algorithm**

- Eliminated 100Gb P2P connection
- Stuffed every existing RO feature to a single ArtDAQ BoardReader
- Extensive server evaluation
  - PCIe riser configurations,
  - BIOS settings
- Optimized for memory throughput
  - Performance profiling
- Heavy NUMA balancing between processing threads and allocated memory
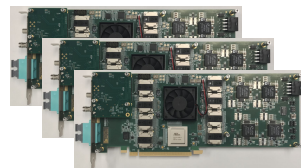- Interrupt moderation of 10Gb NIC

2 x OnHost process

1 x DMA parser
5 x link parser    →    High memory I/O
5 x hit finder
5 x trigger matcher    →    High CPU util.



17.8% - Front-End Bound
37.2% - Memory Bound
35.0% - Retiring
8.5% - Core Bound

µPipe

# OnHost BoardReader 2



FELIX & processing host

FELIX    AVX    QAT

**6 x OnHost process**

1 x DMA parser
5 x link parser
5 x trigger matcher

High memory I/O

**Locality and performance characteristics understood**

- 2 x APAs - demonstrated within ProtoDUNE run environment
- One full socket available for post-processing

Other studies also done:
- 3 x APAs - emulated data with fake triggering demonstrated on the same system

**The solution scales on a well understood manner**

The talk does not focus on performance indicators, but purely about implementation shortcomings in the FELIX OnHost Boardreader.

# OnHost BoardReader - Problems

**Extremely monolithic approach**
**The result became the anti-pattern: "big ball of mud"**

- Backward compatibility attempt and bloated feature set
  - P2P mode, frame reordering, QAT compression
  - Linked against unnecessary and heavy dependencies (libqat, ibverbs, etc.)
  - Forces platform dependent build (avx intrinsics, QAT)
  - No interfaces straightforward interfaces for TP handling
- Hardwired logic based on configuration combinations with override
  - E.g.: If (mode == "OnHost" && num_sources == 2)  // this must be 10 links! <- bad
  - Hardcoded constants and cryptic identifiers
- Tree hierarchy in resource management
  - Hierarchy between different resources and processing element are not clear
  - Central "glue" holds every member and prepare resources for them
- Custom and unconventional ways in some components
  - Avoided framework features (tracing, logging, statistics handling)
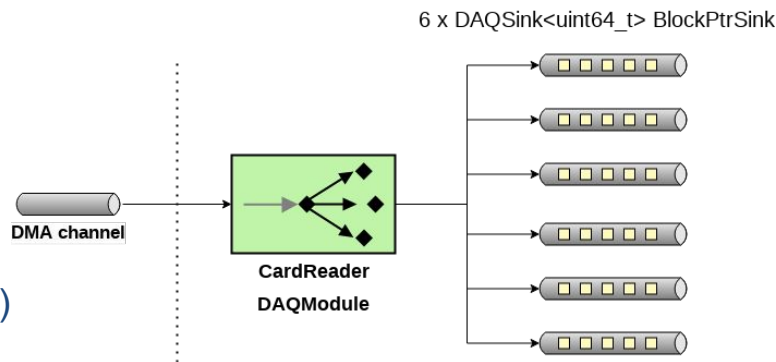- Complete lack of documentation

# OnHost BoardReader - What to do?

**Cleanup or structural redesign and reimplementation?**

- Fearing from the cleanup resulting in "yet another OhBR version", I decided to make an attempt to re-design the Readout software, following the principles of the appfwk

1. Identify resources
   a. Physical FELIX Card, Memory that holds WIB frames, Sockets that publish Hits
2. Identify processing elements
   a. DMA Parser, Block/Link Parser, TriggerMatcher, RequestReceiver
3. Identify connections between the resources and processing elements
   a. DMA Parser reads card, produces pointers
   b. Parsers consume pointers and produce WIB frames
   c. TriggerMatchers use WIB frame containers and publish on their custom network layer
4. Group the components into a processing pipeline

# CardReaderModule

**Reads the FELIX card's DMA channel and based on the Block header, it queues in its pointer to the corresponding DAQ sink**

- Resources
  - FLX DMA buffer (populated by every enabled logical link's output)
- Processing stage
  - Communicates with the card's driver to advance the read pointers, indicating that we have "seen" the new block
  - The pointer to the "seen" block is stored in the link's specific queue
- Sinks
  - Queues for pointers (for links that are enabled!)
- Remarks:
  - Only this module needs to be linked against libFlxCard.so from the FELIX dependencies

6 x DAQSink<uint64_t> BlockPtrSink

DMA channel

CardReader
DAQModule

# LinkReaderModule

**Parses FELIX SuperChunks from a single BlockPtrSource resource**

- <u>Resources</u>
  - BlockPtrSource of link that the Module need to process
- <u>Processing stage</u>
  - Merge the fragments of the chunks that span across multiple blocks
  - Report problems on data corruption
  - Push chunks to WIB frame buffer
- <u>Sinks</u>
  - Latency buffers that hold the actual raw data



DAQSource<uint64_t> BlockPtrSource
(elink_tag: 0 to 6x64)

**LinkReader
DAQModule**

DAQSink<SUPERCHUNK> ChunkSink
(elink_tag: 0 to 6x64)

**Should it be available to read from multiple BlockPtrSource and spawn processing units in the same Module?**
**Sometimes an extra scaling capability may compromise clarity…**

# TriggerMatcherModule

**Finds requested data in the latency buffers**

- Resources
    - Circular buffer of raw WIB binary data
- Processing stage
    - Extracts data from buffers based on the requests
- Sinks
    - Requested data sinks

# Other features/components

- SNB Buffer needs some experimentation how to connect in the the topology
  - E.g.: specific DataRequest, streamed by the TriggerMatcher to dedicated, special sinks
- DataRequest handling need to be extended with specific RequestTypes
  - Might need to drive the TriggerMatching mechanism based on that
- TriggerMatcher redesign
  - Need to support overlapping triggers and variable request window sizes
  - Direction is clear, easier solution exist than expected
  - Non-intrusive to existing sink/source types
- Fake mode
  - I want to extend the LinkReader with a "fake mode" that is able to replay raw data files over-and-over and re-create O(10s) latency buffer with content in memory.
  - Add fake-trigger to TriggerMatcher: extract last available data with random window (ms-s)
  - Add fake-trigger to request mode
- Stretch goals
  - Compression DAQmodule using QuickAssist (Compress data from inputs to outputs)
  - Post-processor DAQModule (Async callback wrapper to do "whatever with extracted data)

# Integration

**Missing modules are coming at a steady pace...**
**https://github.com/roland-sipos/udaq-readout.git**

# Summary & Outlook

Summary:
- All in all, really good experience so far
  - Fought with some build issues on the beginning. Needed to get familiar with the framework a bit
  - Great support and guidance from the developers!
    Thanks a lot for it.

Outlook:
- Ex. CardReader I want every Module to operate in emulation mode, providing a full readout layer for testing
- Provide fake-mode chain ASAP for Dataflow development (HDF5 Module)
- Provide full-chain (with CardReader) ASAP for MiniDAQ

# End

- Feedback is more than welcome!

# Backup - Readout software performance