



Design of a request/response buffering application for I/O intensive workloads

Florian Grötschla, Giovanna Lehmann Miotto, Roland Sipos

Overview

1. Detector readout challenges
2. Specification and design principles
3. Components and implementation details
4. Demonstration
5. Outlook

Readout challenges

1. **Wide range of frontend electronics:** TPC electronics, SiPM readout, pixel chips, etc.
2. **Variety of aggregator I/O devices:** COTS servers, PCIe FPGA carrier boards, NICs, etc.
3. **Payload characteristics:** different arrival rate and payload size combinations (fixed/variable)
4. **Quasi real-time performance:** Requirement for high throughput ($\sim 100\text{Gbit/s}$), low latency, and scalability

Frontends

TPC readout



SiPM readout



COTS servers



FPGA boards



different rates and
payload sizes

NICs



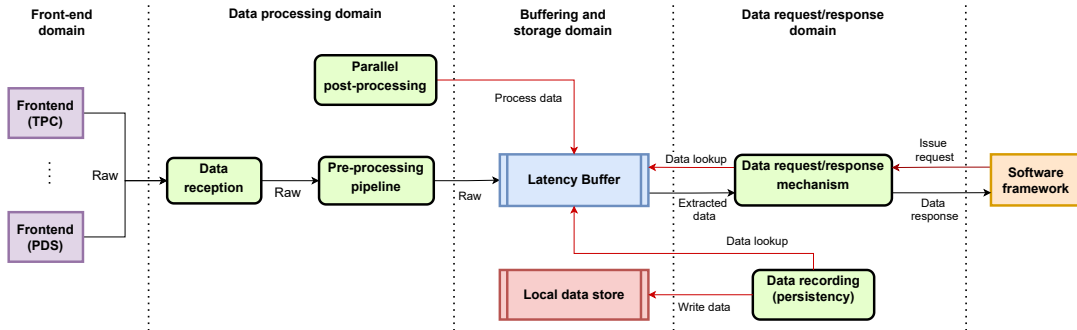
Readout specification

- Support all possible front-end types: be agnostic about data rate and payload size
- Buffer received data for a specified/maximum amount of time
- Respond to data requests (with time-windows)
 - Buffered data can be indexed to maintain search-ability
- In-flight data processing: Error and consistency checks with custom algorithms (e.g.: hit-finding) also supported
- Persist data on command for requested time

Readout components

- **Readout Type:** A wrapper for the raw data that provides functionalities needed by the buffer and other components
- **Latency Buffer:** Provide data structure for data buffering and lookup mechanism of indexed data
- **Frame Processor:** Pre- and post-processing of raw input data
- **Request Handler:** Respond to data requests, handle recording
- **Readout Model:** Contains all other components, calls interface implementations and hands off resources (e.g.: buffer)

Dataflow diagram



Concepts and models

- Diverse set of frontends, having varied characteristics in payload arrival rate, size, and order
- **Main goal:** Avoid reimplementations and code duplication for different frontends
 - Keep the core functionalities fully generic
- **Concepts:** Well defined interfaces
- **Models:** Interchangeable implementations of the concepts

Latency buffers

- Stores data and provides lookup routine based on unique identifier (e.g.: timestamp) in frames
- **Fixed rate queue:** Offset calculation for lookup mechanism
 - Implementation based on [Folly](#) lock-free queue
 - Extension with binary search capability
- **Skiplist:** Uses [Folly](#) concurrent skiplist for non-ordered data
- Memory allocation strategies supported: NUMA aware, aligned

All implementations are interchangeable!

Frame Processor

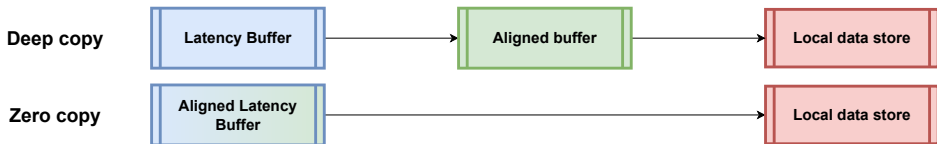
- **Preprocessing pipeline:** Tasks can be registered to a sequence of operations that are executed for every payload
- **Parallel postprocessing:** One thread for every registered function is created that is fed with the payloads
- Preprocessing can change frame, postprocessing works on a constant pointer
- Postprocessing works directly on a pointer to the frame in the latency buffer, no extra copies are involved

Request handling

- Request handler threads deep-copy data from latency buffer
- Read-only, so handle them in parallel via thread pool
- Periodic cleanup thread removes old data from latency buffer
- Synchronize cleanup and requests to avoid interference
- Waiting requests: Postpone requests that are for not yet present data

Data recording

- High bandwidth requirement imposes optimized implementation (exploit kernel features like `O_DIRECT`)
- Use of **boost streams** and support for compression
- 2 recording modes are supported
 - Deep-copy: From latency buffer to boost stream and writer
 - Zero-copy: From aligned latency buffers, directly to storage device

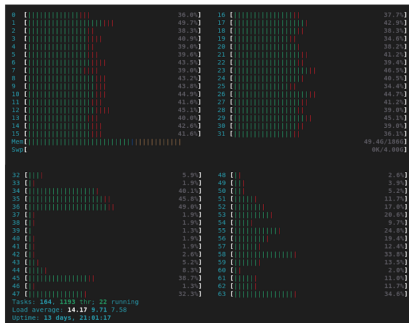


Demonstration

- The library is successfully integrated and used in different scenarios
 - With payload software emulators (fixed rate & size, variable rate & size)
 - With FPGA carrier boards ([FELIX](#) readout)
 - For the TPC readout of DUNE Vertical-Drift ColdBox (WIB frontend)
 - With NICs
 - For the photodetector of ProtoDUNE-SP (SSP frontend)
- Successful tests of the recording implementations using different storage media

Performance snapshot

```
..... Socket 0 ..... Socket 1 .....
----- Memory Channel Monitoring ----- Memory Channel Monitoring -----
-- Mem Ch 0: Reads (MB/s): 4454.53 -- Mem Ch 0: Reads (MB/s): 352.79 --
-- Writes (MB/s): 6823.58 -- Writes (MB/s): 689.62 --
-- PMM Reads (MB/s): 0.00 -- PMM Reads (MB/s): 0.00 --
-- PMM Writes (MB/s): 0.00 -- PMM Writes (MB/s): 0.00 --
-- Mem Ch 1: Reads (MB/s): 4444.70 -- Mem Ch 1: Reads (MB/s): 347.57 --
-- Writes (MB/s): 6884.42 -- Writes (MB/s): 685.34 --
-- PMM Reads (MB/s): 0.00 -- PMM Reads (MB/s): 0.00 --
-- PMM Writes (MB/s): 0.00 -- PMM Writes (MB/s): 0.00 --
-- Mem Ch 2: Reads (MB/s): 4492.81 -- Mem Ch 2: Reads (MB/s): 350.66 --
-- Writes (MB/s): 6893.78 -- Writes (MB/s): 694.25 --
-- PMM Reads (MB/s): 0.00 -- PMM Reads (MB/s): 0.00 --
-- PMM Writes (MB/s): 0.00 -- PMM Writes (MB/s): 0.00 --
-- Mem Ch 3: Reads (MB/s): 4584.88 -- Mem Ch 3: Reads (MB/s): 343.93 --
-- Writes (MB/s): 6858.51 -- Writes (MB/s): 686.36 --
-- PMM Reads (MB/s): 0.00 -- PMM Reads (MB/s): 0.00 --
-- PMM Writes (MB/s): 0.00 -- PMM Writes (MB/s): 0.00 --
-- Mem Ch 4: Reads (MB/s): 4558.34 -- Mem Ch 4: Reads (MB/s): 343.65 --
-- Writes (MB/s): 6139.37 -- Writes (MB/s): 687.56 --
-- PMM Reads (MB/s): 0.00 -- PMM Reads (MB/s): 0.00 --
-- PMM Writes (MB/s): 0.00 -- PMM Writes (MB/s): 0.00 --
-- Mem Ch 5: Reads (MB/s): 4589.56 -- Mem Ch 5: Reads (MB/s): 336.38 --
-- Writes (MB/s): 6839.47 -- Writes (MB/s): 683.22 --
-- PMM Reads (MB/s): 0.00 -- PMM Reads (MB/s): 0.00 --
-- PMM Writes (MB/s): 0.00 -- PMM Writes (MB/s): 0.00 --
-- NODE 0 Mem Read (MB/s): 26955.62 -- NODE 1 Mem Read (MB/s): 2078.98 --
-- NODE 0 Mem Write (MB/s): 36359.13 -- NODE 1 Mem Write (MB/s): 4126.65 --
-- NODE 0 PMM Read (MB/s): 0.00 -- NODE 1 PMM Read (MB/s): 0.00 --
-- NODE 0 PMM Write (MB/s): 0.00 -- NODE 1 PMM Write (MB/s): 0.00 --
-- NODE 0.0 NM read hit rate: 1.00 -- NODE 1.0 NM read hit rate: 1.00 --
-- NODE 0.1 NM read hit rate: 1.00 -- NODE 1.1 NM read hit rate: 1.00 --
-- NODE 0.2 NM read hit rate: 0.00 -- NODE 1.2 NM read hit rate: 0.00 --
-- NODE 0.3 NM read hit rate: 0.00 -- NODE 1.3 NM read hit rate: 0.00 --
-- NODE 0 Memory (MB/s): 63214.15 -- NODE 1 Memory (MB/s): 6286.95 --
-----
-- System DRAM Read Throughput (MB/s): 29033.92 --
-- System DRAM Write Throughput (MB/s): 48482.18 --
-- System Mem Read Throughput (MB/s): 0.00 --
-- System Mem Write Throughput (MB/s): 0.00 --
-- System Read Throughput (MB/s): 29033.92 --
-- System Write Throughput (MB/s): 48482.18 --
-- System Memory Throughput (MB/s): 69521.10 --
-----
```

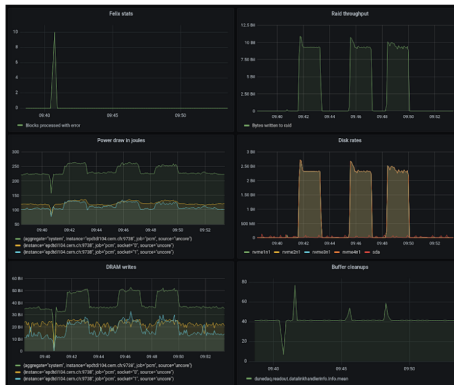
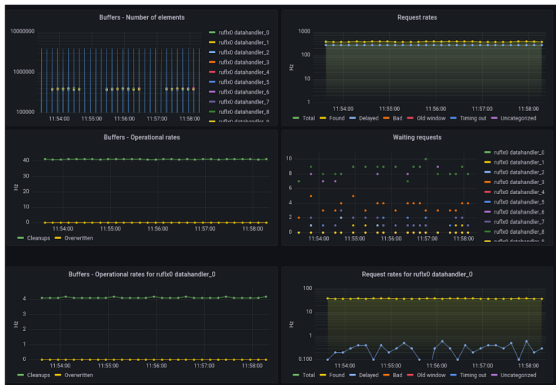


Server: 2S Intel Cascade Lake (Xeon(R) Gold 6242), 64 cores with 192GB RAM

Left: memory I/O: < 60% bandwidth utilization, Right: total CPU load per socket: 50%
Workload: Receiving at 166kHz x 5568 Bytes x 10 links = 8.6 GB/s and software hit finding



Recording performance



Left: same setup as before, Right: 3 consecutive recordings of 100 seconds each
Writing 10 links (8.6GB/s) to a software RAID0 of 4 NVMe SSDs (Samsung 970 Pro 1TB)

Outlook

- Main focus is on **scalability** and **resource locality** studies
 - Optimize CPU and memory bandwidth utilization via dynamic thread affinity balancers
 - Improve performance on multi-socket systems, exploiting NUMA awareness and socket interconnect capabilities
- Integretation of the generic readout library into other DAQ software frameworks (e.g.: [DAQling](#))



Modularity through templating

Readout for WIB

```
auto readout_model = ReadoutModel<
    // Readout Type: WIB Struct
    WIBSS,
    // Request Handler
    DefaultRequestHandlerModel<WIBSS,
        FixedRateQueueModel<WIBSS>>,
    // Latency buffer
    FixedRateQueueModel<WIBSS>,
    // Frame Processor
    WIBFrameProcessor>(run_marker);
readout_model.init(args);
```

Readout for DAPHNE

```
auto readout_model = ReadoutModel<
    // Readout Type: DAPHNE Struct
    DAPHNESS,
    // Request Handler
    DAPHNEListRequestHandler,
    // Latency buffer
    SkipListLatencyBufferModel<DAPHNESS>,
    // Frame Processor
    DAPHNEFrameProcessor>(run_marker);
readout_model.init(args);
```