



Storing LHC Data in DAOS and S3 through RNTuple

26th International Conference on Computing in High Energy and Nuclear Physics

Giovanna Lazzari Miotto – Federal University of Rio Grande do Sul (BR)

Javier López-Gómez – CERN (CH)

CHEP 2023, 9 May 2023



- 1 Introduction
- 2 Storing RNTuple data in DAOS
- 3 Storing RNTuple data in S3
- 4 On data migration
- 5 Summary

Introduction



HL-LHC is projected to increase ROOT data **at least tenfold**

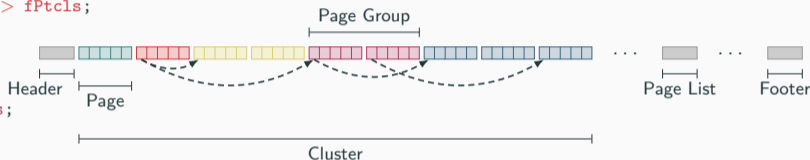
- **HW landscape changing:** NVMe devices, architectural heterogeneity, parallelism, distributed facilities. . .
- **TTree is not designed to accomodate this leap** → enter **RNTuple**

RNTuple goals:

- ×2-5 better single-core performance
- 25% smaller files
- 1 GB/s/core sustained end-to-end throughput
- Robust interfaces and systematic use of exceptions
- . . . and more: check yesterday's RNTuple talks by Jakob [1], Florine [2]



```
struct Event {  
    int fId;  
    vector<Particle> fPtcls;  
};  
struct Particle {  
    float fE;  
    vector<int> fIds;  
};
```



- **Page:** array of values of a fundamental type. Size $O(\text{tens of KiB})$
- **Cluster:** all pages containing data for a specific row range, e.g. 1–1000
- **Page Group:** all pages that contain data for the same column in a given cluster
- **Metadata:** header, page lists, footer



In a highly-parallel WORM¹ setting, object stores align well with our requirements:

- Extremely scalable
- Widely deployed in cloud service providers

Contexts: **HPC** (Intel DAOS) \neq **Cloud** (Amazon S3, Microsoft Azure, Google Cloud)

Downsides

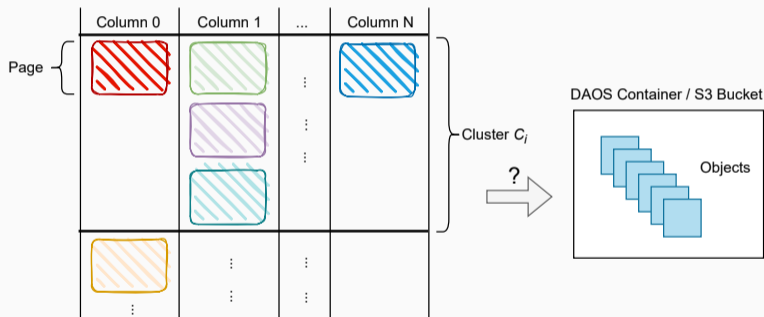
- No hierarchical namespaces
- Not based on POSIX I/O API
- **No storage of other ROOT classes: histograms, etc.**

¹“Write Once, Read Many”



Blow up RNTuple *the right way*

- Data granularity: *page, page group, cluster^a*



Factors to consider

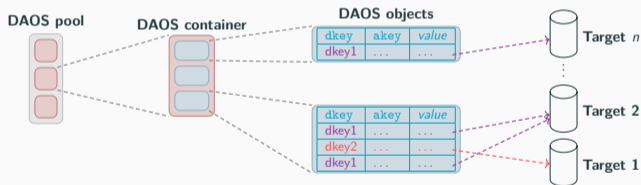
- Analysis pattern
- Throughput, latency
- Cost per request
- Memory consumption

^adependent on byte ranges support

Storing RNTuple data in DAOS



- Fault-tolerant object store for HPC
 - Low-latency, high-bandwidth, high IOPS
- Foundation of the Intel exascale storage stack
- Used in 44% of the top 25 systems in IO500² (e.g., ANL Aurora)



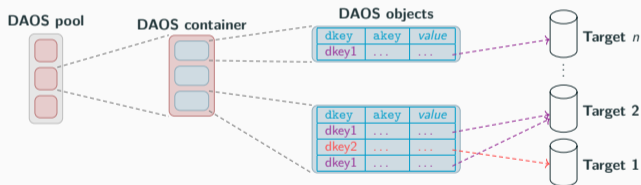
- Locality feature: $\langle \text{object}, \text{dkey} \rangle \rightarrow \text{target}$

```
auto ntuple = RNTupleReader::Open("DecayTree",
    "daos://my-pool/my-container");
```

²<https://io500.org/>



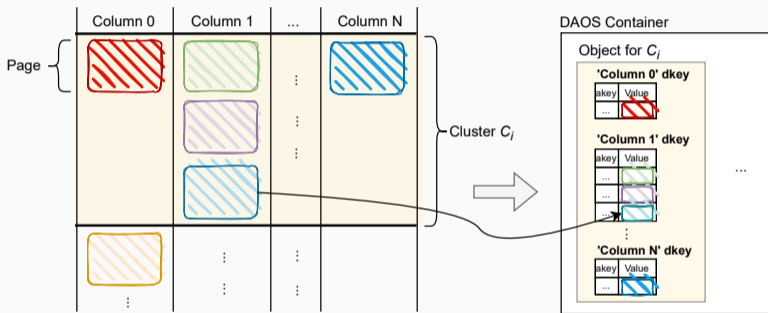
- Fault-tolerant object store for HPC
 - Low-latency, high-bandwidth, high IOPS
- Foundation of the Intel exascale storage stack
- Used in 44% of the top 25 systems in IO500² (e.g., ANL Aurora)



- Locality feature: $\langle \text{object}, \text{dkey} \rangle \rightarrow \text{target}$

```
auto ntuple = RNTupleReader::Open("DecayTree",
  "daos://my-pool/my-container");
```

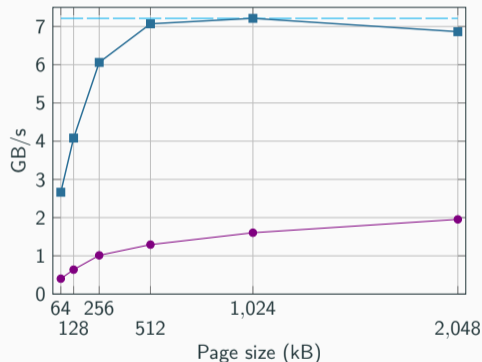
²<https://io500.org/>



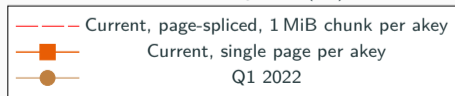
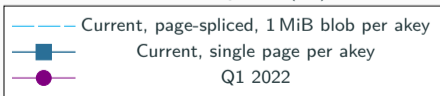
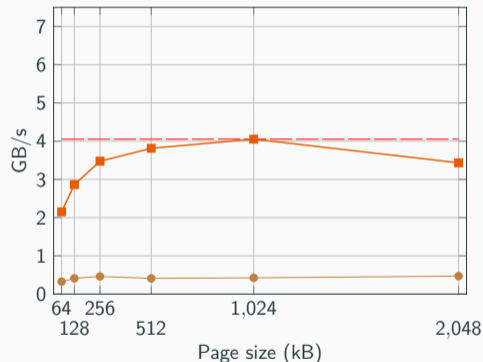
Page splicing: empirically optimal at 1 MiB blobs (throughput × granularity)



Plot (1.a): write throughput (no compr.)



Plot (1.b): read throughput (no compr.)



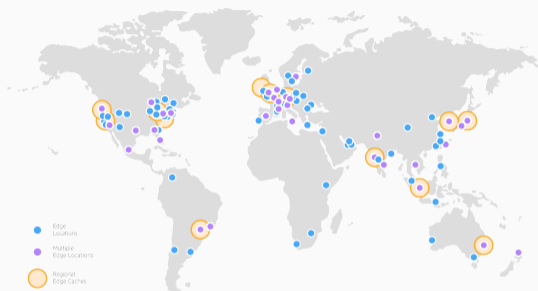
Storing RNTuple data in S3



Industry's de facto standard

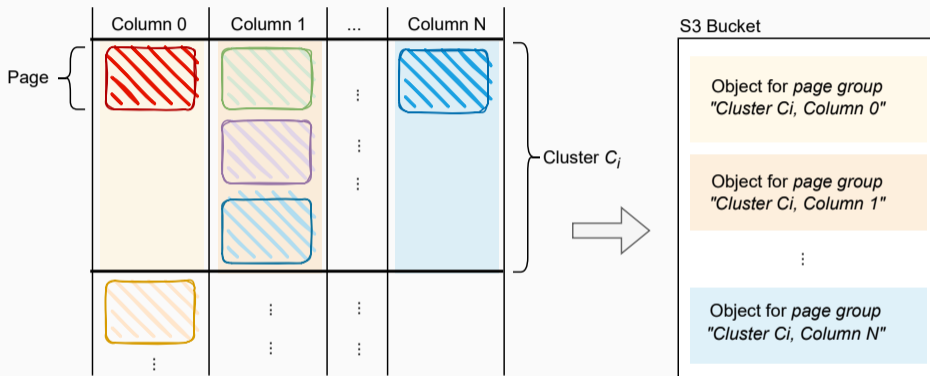
Key differences to DAOS:

- Flat namespace: just buckets with objects
- Worldwide server distribution across “regions” and “edges”
- Potential interop with research computing infrastructures, e.g., WLCG



AWS edge locations and regional caches. Amazon.

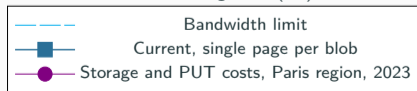
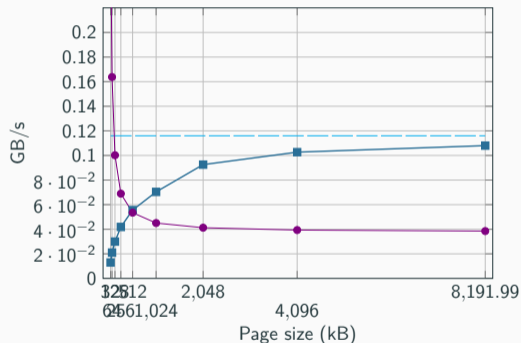
```
auto ntuple = RNTupleReader::Open("DecayTree",  
    "s3://server-region/my-bucket");
```



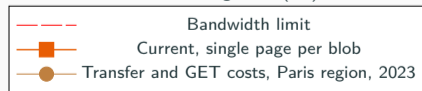
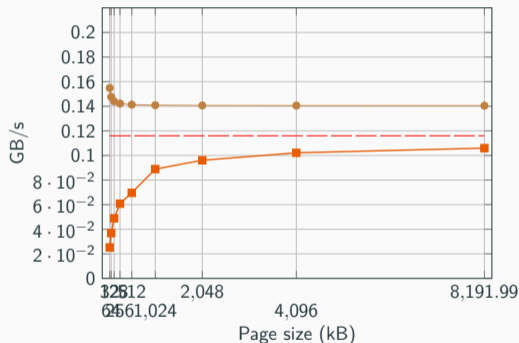
HTTP client (Davix): bigger blobs to counteract latency over network



Plot (1.a): write throughput (no compr.)



Plot (1.b): read throughput (no compr.)



On data migration



RNTuple performance is **very sensitive to its I/O parameters** → see Dante's poster

- As seen, even between two object stores!
- Moving data between storage systems requires reshaping for efficient access

Implemented

Fast cloning (like TTree's CloneTree)

- Avoids page decompression unless necessary

In development

Reshaping ntuples

- Compression algorithm, page size, cluster size, field selection

Summary



- RNTuple is set to be production-ready in 2024, leveraging **HPC** and **cloud object stores** for Run 3 / 4 analyses
 - **Native, mature DAOS backend** with 8+ GB/s writes, 4+ GB/s reads single-node
 - **Native, experimental S3 backend**
- Efficient data migration across storage systems, ntuple reshaping (WIP)

Next steps

- Optimize S3 backend based on first results
 - Leverage AWS C++ SDK with Davix as compatibility fallback
- Expand and improve RNTuple Migrator



CHEP 2023

Thanks!