



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

Introduction



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

- **New HW landscape:** NVMe devices, architectural heterogeneity, parallelism, distributed facilities. . .
- **RNTuple** (“TTree 2.0”) is optimized to accommodate this leap!



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

- **New HW landscape:** NVMe devices, architectural heterogeneity, parallelism, distributed facilities. . .
- **RNTuple** (“TTree 2.0”) is optimized to accommodate this leap!

RNTuple goals:

- Over 2× better single-core performance
- 1 GB/s/core sustained end-to-end throughput
- . . . and more: check yesterday’s session!
- 25% smaller files
- Robust interfaces and systematic use of exceptions

► talk: RNTuple

► talk: RNTuple in Athena



In a highly-parallel 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)



In a highly-parallel 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)

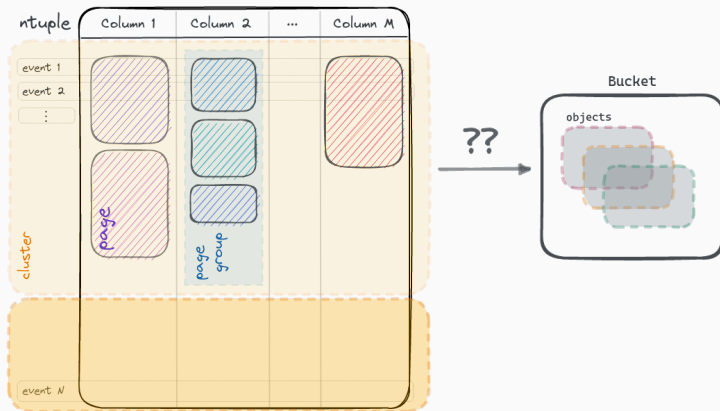
But why invest in native support? **Granular parallel access**

- **Downside:** no storage of other ROOT classes: histograms, etc.

Challenge: Mapping RNTuple → Object Stores



Granularity: **page** $\sim O(\text{KiB})$, **page group**, **cluster**^a $\sim O(\text{MiB})$



Factors to consider

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

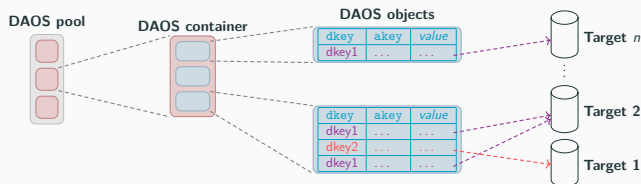
^aFor columnar access, storing clusters wholesale requires byte ranges support

Storing RNTuple data in DAOS



Foundation of the Intel exascale storage stack

- Built for NVMe and persistent memory
- Low-latency, high-bandwidth, high IOPS
- Used in 44% of the top 25 systems in IO500¹ (e.g., **ANL Aurora**)



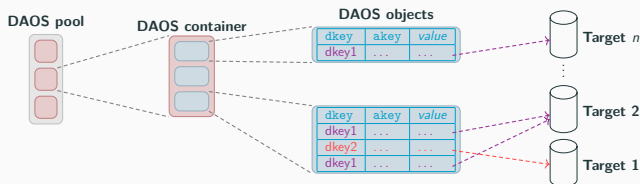
- Locality feature: $\langle object, dkey \rangle \rightarrow target$

¹<https://io500.org/>



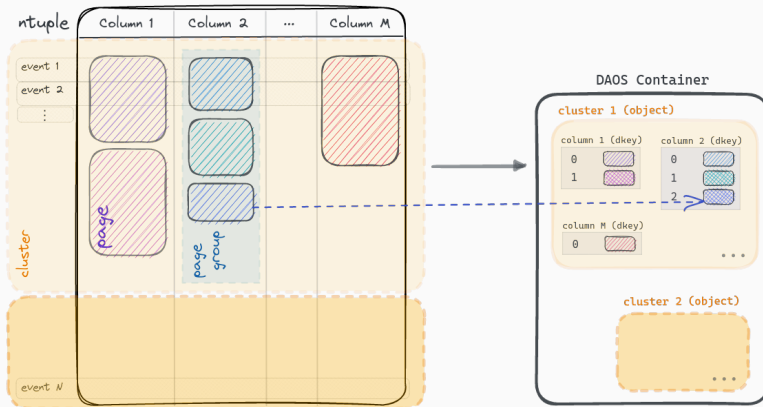
Foundation of the Intel exascale storage stack

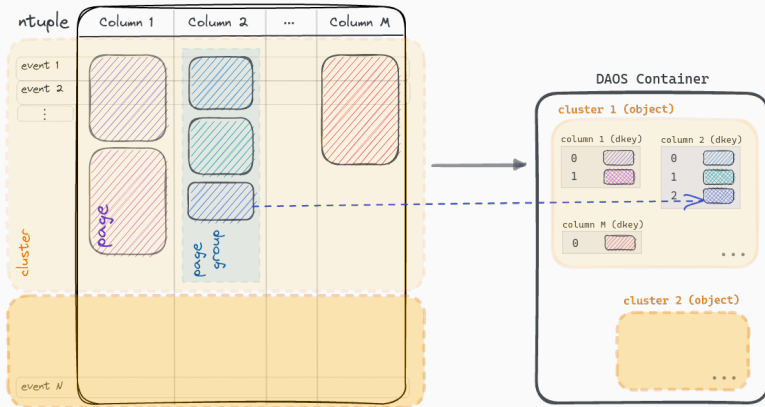
- Built `RDF::Experimental::FromRNTuple("DecayTree", "daos://my-pool/my-container")`
- Low-latency, high-bandwidth, high IOPS
- Used in 44% of the top 25 systems in IO500¹ (e.g., **ANL Aurora**)



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

¹<https://io500.org/>



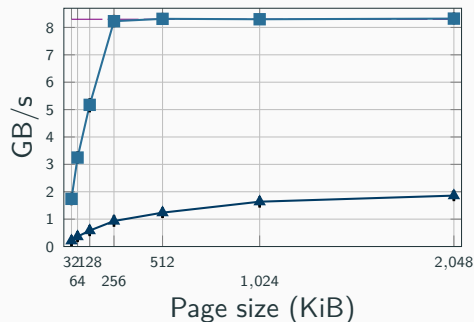


“Page splicing”: empirically optimal with 1 MiB blobs (throughput × granularity)

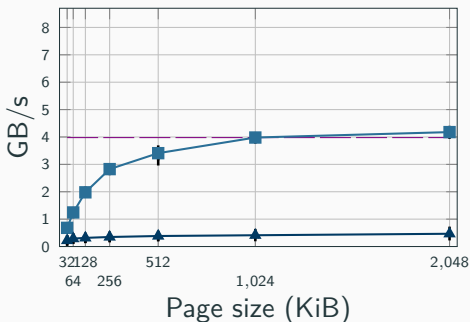


E2E analysis on single DAOS client-server pair (**HPE's Delphi cluster w/ Infiniband interconnect**)

(2.a): write throughput (no compr.)



(2.b): read throughput (no compr.)



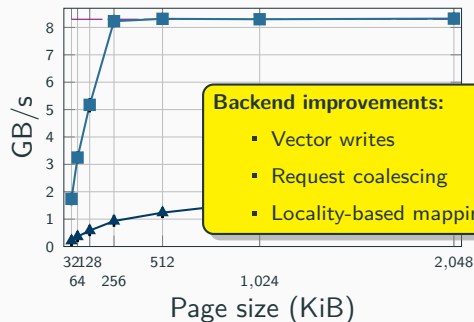
— Current, page-spliced (1 MiB)
—▲— CHEP 2021 [doi]

—■— Current, single page per akey

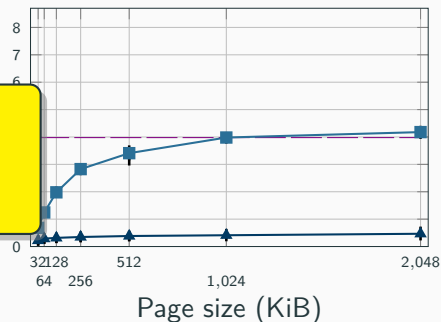


E2E analysis on single DAOS client-server pair (HPE's Delphi cluster w/ **Infiniband** interconnect)

(2.a): write throughput (no compr.)



(2.b): read throughput (no compr.)



Backend improvements:

- Vector writes
- Request coalescing
- Locality-based mapping

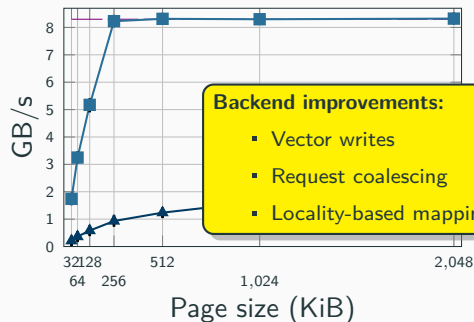
— Current, page-spliced (1 MiB)
— CHEP 2021 [doi]

— Current, single page per akey



E2E analysis on single DAOS client-server pair (HPE's Delphi cluster w/ Infiniband interconnect)

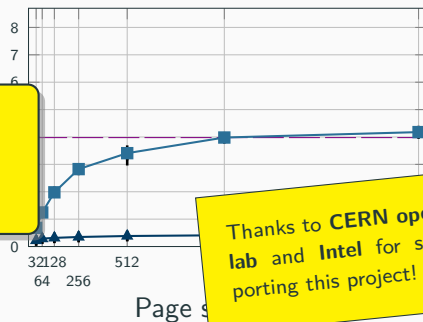
(2.a): write throughput (no compr.)



Backend improvements:

- Vector writes
- Request coalescing
- Locality-based mapping

(2.b): read throughput (no compr.)



Thanks to **CERN open-lab** and **Intel** for supporting this project!

— Current, page-spliced (1 MiB)
— CHEP 2021 [doi]

— Current, single page per akey

Storing RNTuple data in S3



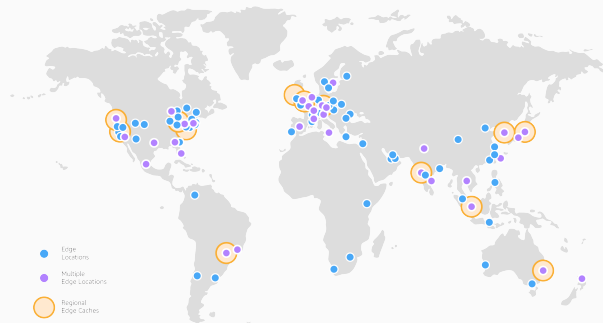
Object storage's de facto standard



Object storage's de facto standard

Main differences w.r.t. DAOS:

- Flat namespace
- Worldwide server distribution across “regions” and “edges”
- Potential interop with research computing infrastructures, e.g., WLCG
- Network latency. . .



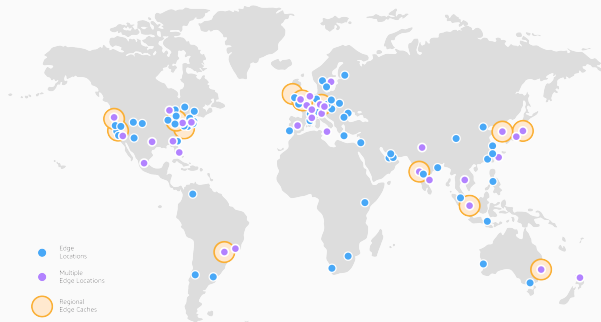
AWS edge locations and regional caches. Amazon.



Object storage's de facto standard

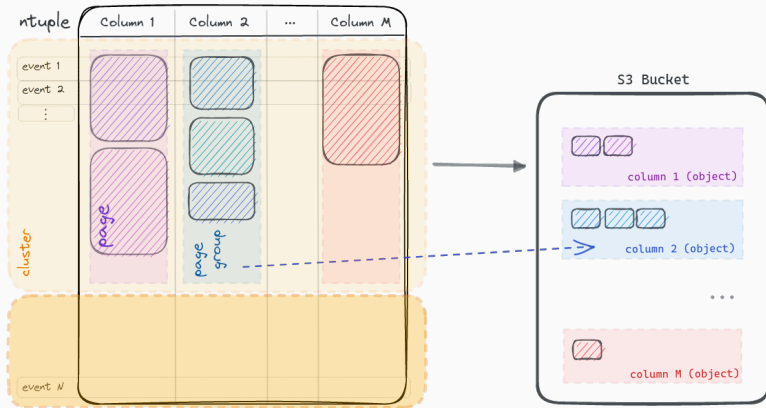
Main differences w.r.t. DAOS:

- Flat namespace
- Worldwide server distribution across “regions” and “edges”
- Potential interop with research computing infrastructures, e.g., WLCG
- Network latency...



AWS edge locations and regional caches. Amazon.

```
RDF::Experimental::FromRNTuple("DecayTree", "s3://aws-region/my-bucket")
```



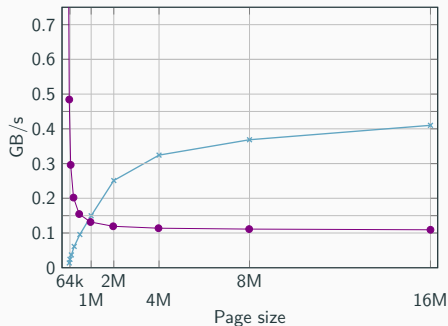
Bigger blobs (**page groups** or **clusters**) may mitigate latency

- Best with scatter-gather I/O, requires Range HTTP requests



E2E analysis on client-server pair (MinIO). Nodes provided by openlab and ROOT.

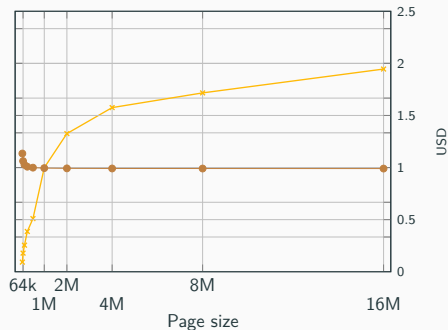
Plot (1.a): write costs, throughput (no compr.)



—x— Current, single page per blob

—●— Projected AWS storage and PUT costs, eu-west-3 region

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



—x— Current, single page per blob

—●— Projected AWS transfer and GET costs, eu-west-3 region

On data migration



RNTuple performance is **very sensitive to its I/O parameters** ► poster: ML + RNTuple

- Even between two object stores!
- Moving data between storage systems requires reshaping for efficient access

Implemented

Fast cloning with field selection (like TTree's CloneTree)

- Avoids page decompression unless necessary

In development

Reshaping ntuples

- Compression algorithm, page size, cluster size, ...

Summary



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



- RNTuple is set to be production-ready in late 2024, leveraging **HPC** and **cloud** object stores
 - **Native, mature DAOS backend** with $8+ \text{ GB/s}$ writes, $4+ \text{ GB/s}$ 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
 - Investigate AWS' C++ SDK, retaining Davix as a compatibility fallback
- Expand RNTuple Migrator



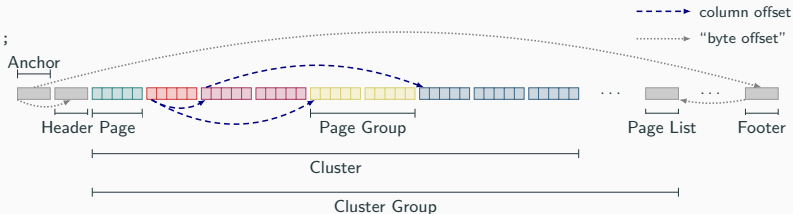
Questions?

glazzari@cern.ch

Acknowledgements: This work benefited from support by the CERN Strategic RD Programme on Technologies for Future Experiments CERN-OPEN-2018-006 and the Intel-CERN openlab collaboration. Access to the hardware for the experimental evaluation was provided by Hewlett-Packard Enterprise.

Breakdown of the RNTuple On-Disk Format

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



Cluster

- Block of consecutive complete events
- Defaults to 50 MB compressed

Page

- Unit of (de-)compression
- Defaults to 64 KiB uncompressed
- Not necessarily aligned on event boundary