



Parallel Writing of Nested Data in Columnar Formats

Jonas Hahnfeld^{1,2} Jakob Blomer¹ Thorsten Kollegger²
jonas.hahnfeld@cern.ch

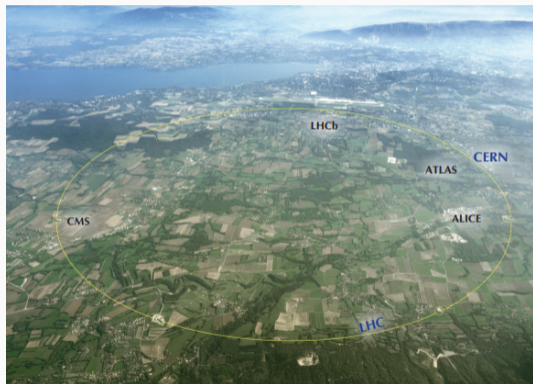
¹ CERN, Geneva, Switzerland

² Goethe University Frankfurt, Institute of Computer Science, Frankfurt, Germany

Euro-Par 2024 – August 29, 2024

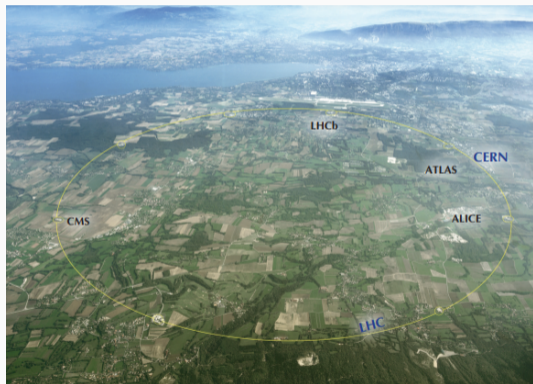


- Large Hadron Collider at CERN
 - More than 2 exabytes since 2010
 - Stored in binary columnar format



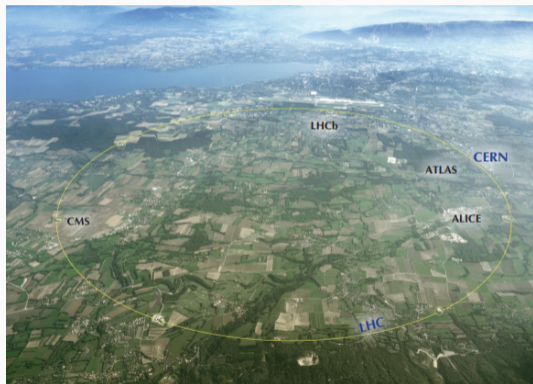


- Large Hadron Collider at CERN
 - More than 2 exabytes since 2010
 - Stored in binary columnar format
- High-Luminosity LHC (from 2029)
 - Further increase of data rate





- Large Hadron Collider at CERN
 - More than 2 exabytes since 2010
 - Stored in binary columnar format
- High-Luminosity LHC (from 2029)
 - Further increase of data rate
- RNTuple: evolution of the currently used TTree columnar format
 - Opportunity to include parallel writing from the start





RNTuple Overview

Concepts for Parallel Writing of Columnar Data

Evaluation of Parallel RNTuple Writing

Conclusions and Future Work

RNTuple Overview



- Serialize acyclic C++ data structures into a columnar format
 - Important to support nested collections for HEP use cases

```
struct Event {  
    int fld;  
    std::vector<Track> fTracks;  
};
```

```
struct Track {  
    float fEnergy;  
    std::vector<int> flds;  
};
```

Figure 1: Simplified example of nested data structures. Real-world HEP data models often have thousands of fields.



- Recursively decompose data structures into *fields*
 - *Columns* at leafs of field tree: primitive, fixed-size types



- Recursively decompose data structures into *fields*
 - *Columns* at leafs of field tree: primitive, fixed-size types
- Columns partitioned into *pages*
 - Transparently compressed (default for RNTuple: Zstandard)



- Recursively decompose data structures into *fields*
 - *Columns* at leafs of field tree: primitive, fixed-size types
- Columns partitioned into *pages*
 - Transparently compressed (default for RNTuple: Zstandard)
- *Cluster*: all pages of a consecutive range of rows, or *entries*

**Table 1:** Example of columnar representation for the nested data structure shown in Figure 1.

fId	fTracks	fTracks._0.fEnergy	fTracks._0.fIds	fTracks._0.fIds._0
6873	2	25.4f	2	42 27
		32.8f	3	16
6874	3	14.7f	5	21 8
⋮	⋮	⋮	⋮	⋮



Table 1: Example of columnar representation for the nested data structure shown in Figure 1.

fId	fTracks	fTracks._0.fEnergy	fTracks._0.fIds	fTracks._0.fIds._0
6873	2	25.4f	2	42
		32.8f	3	27
6874	3	14.7f	5	16
				21
				8
⋮	⋮	⋮	⋮	⋮

Diagram annotations:

- A green box highlights the column containing `fTracks._0.fEnergy` values (25.4f, 32.8f, 14.7f, and vertical ellipsis). A green label `column` is positioned below this box.
- A red box highlights the row containing `fId` 6874 and its associated `fTracks` data (3, 14.7f, 5, 21, 8). A red label `entry` is positioned below this row.



Table 1: Example of columnar representation for the nested data structure shown in Figure 1.

fId	fTracks	fTracks._0.fEnergy	fTracks._0.fIds	fTracks._0.fIds._0
6873	2	25.4f	2	42
		32.8f	3	27
		14.7f	5	16
6874	3			21
				8
⋮	⋮	⋮	⋮	⋮

Diagram annotations:

- A green box labeled "column" encloses the `fTracks._0.fEnergy` column.
- A red box labeled "entry" encloses the row for `fId` 6874.
- A blue box labeled "page" encloses the `fTracks._0.fIds._0` column.



Table 1: Example of columnar representation for the nested data structure shown in Figure 1.

fId	fTracks	fTracks._0.fEnergy	fTracks._0.fIds	fTracks._0.fIds._0
6873	2	25.4f	2	42
		32.8f	3	27
6874	3	14.7f	5	16
				21
				8
⋮	⋮	⋮	⋮	⋮

Diagram annotations:

- A green box highlights the `fTracks._0.fEnergy` column, labeled "column".
- A red box highlights the row for `fId` 6874, labeled "entry".
- A blue box highlights the `fTracks._0.fIds._0` column, labeled "page".
- A brown box highlights the entire table structure, labeled "cluster".

Concepts for Parallel Writing of Columnar Data



- Support for nested data implies variable row sizes
 - Transparent compression leads to variable page sizes
- Data cannot be organized in a regular grid



- Support for nested data implies variable row sizes
 - Transparent compression leads to variable page sizes
 - Data cannot be organized in a regular grid

- Clusters are *relocatable*
 - Serialization and compression without synchronization
 - Then: final size is known, can be written into binary file format



- Support for nested data implies variable row sizes
 - Transparent compression leads to variable page sizes
 - Data cannot be organized in a regular grid
- Clusters are *relocatable*
 - Serialization and compression without synchronization
 - Then: final size is known, can be written into binary file format
- Metadata updated in critical section

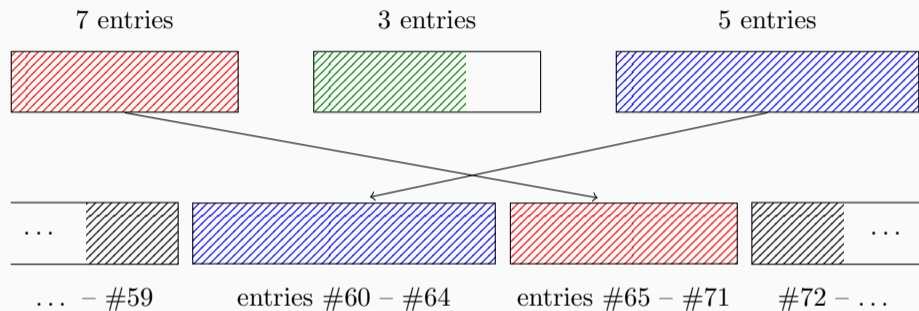


Figure 2: Illustration of filling three RNTuple clusters in parallel. After buffering in memory (top), entries are appended into a sequential file (bottom).

Evaluation of Parallel RNTuple Writing



- Evaluate weak scaling behavior for a synthetic benchmark
 - Fixed number of 20 million entries per thread
 - Two top-level fields: “event ID” and a vector of floats



- Evaluate weak scaling behavior for a synthetic benchmark
 - Fixed number of 20 million entries per thread
 - Two top-level fields: “event ID” and a vector of floats
- Write data on Samsung PM1733 NVMe SSD formatted with ext4
 - Measured bandwidth with [Flexible I/O Tester \(fio\)](#): 768 MB/s



- Evaluate weak scaling behavior for a synthetic benchmark
 - Fixed number of 20 million entries per thread
 - Two top-level fields: “event ID” and a vector of floats
- Write data on Samsung PM1733 NVMe SSD formatted with ext4
 - Measured bandwidth with [Flexible I/O Tester \(fio\)](#): 768 MB/s
 - Possible optimization: pre-allocate file with `fallocate`
 - Increases bandwidth to 1062 MB/s

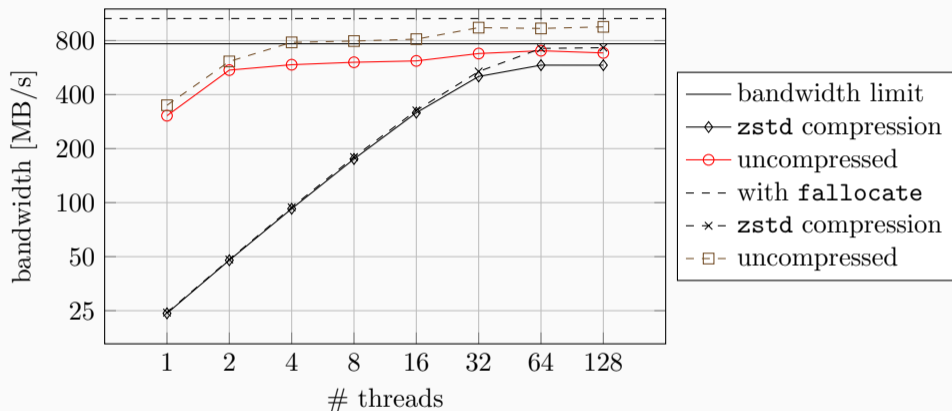


Figure 4: Bandwidth measured with the synthetic benchmark writing on a server SSD.¹

¹Figure adapted from the paper, see Backup Slide 3 for the original plot.



- **Analysis Grand Challenge:** test workflows at scales required for HL-LHC
 - “ttbar” analysis, input dataset derived from 2015 Open Data of the CMS experiment
 - Conversion to RNTuple: 969 GB across 787 files, divided into nine partitions



- **Analysis Grand Challenge:** test workflows at scales required for HL-LHC
 - “ttbar” analysis, input dataset derived from 2015 Open Data of the CMS experiment
 - Conversion to RNTuple: 969 GB across 787 files, divided into nine partitions
- Evaluation benchmark: reduce size of the dataset by *skimming*
 - Retain only fields used by analysis, filter events (entries) based on coarse cuts
 - Output: nine files, total size of 19 GB

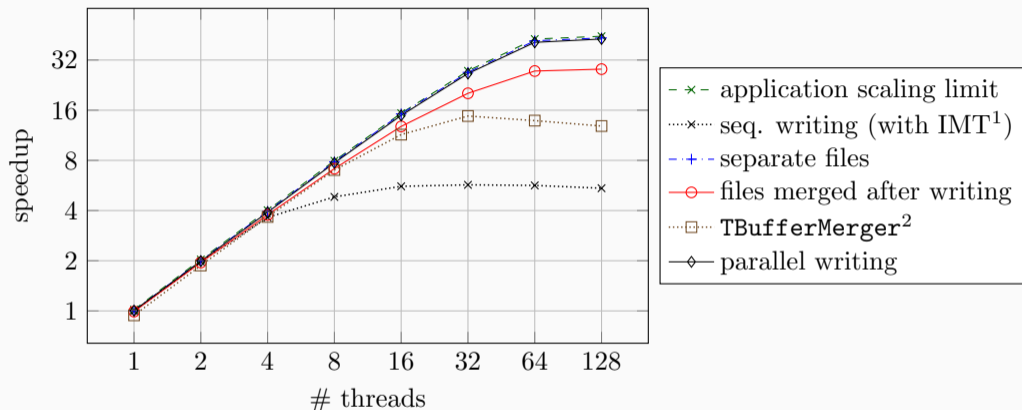


Figure 6: Speedup of the AGC dataset skimming benchmark compared to a full sequential run.

¹ROOT's *implicit multithreading* (IMT) used to parallelize compression of pages.

²TBufferMerger merges files in memory. It is the current way of parallel writing in ROOT.

Conclusions and Future Work



- Presented concept for parallel writing of nested data in columnar formats
- Implemented for RNTuple: <https://github.com/root-project/root/>
- Evaluation with synthetic benchmark and application of dataset skimming



- Presented concept for parallel writing of nested data in columnar formats
- Implemented for RNTuple: <https://github.com/root-project/root/>
- Evaluation with synthetic benchmark and application of dataset skimming

- Synthetic benchmark scales up to storage bandwidth limit
 - Plan to investigate usage of Direct I/O
 - [Presentation at CHEP 2024](#) in October

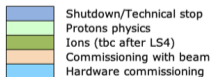
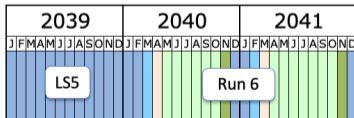
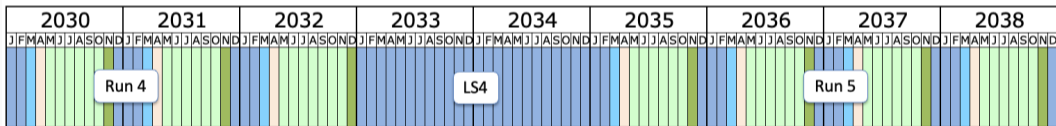
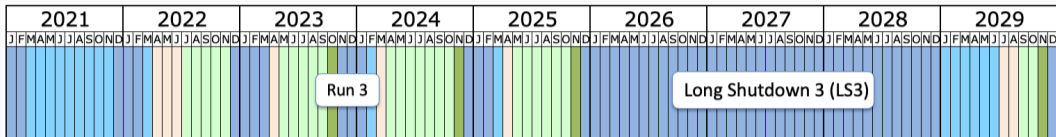


- Presented concept for parallel writing of nested data in columnar formats
- Implemented for RNTuple: <https://github.com/root-project/root/>
- Evaluation with synthetic benchmark and application of dataset skimming

- Synthetic benchmark scales up to storage bandwidth limit
 - Plan to investigate usage of Direct I/O
 - [Presentation at CHEP 2024](#) in October

- Extend parallel writing to process-level parallelism
- Distributed parallelism with MPI and writing to cluster filesystems

Backup Slides



Last update: June 24

