



## Schema Evolution and the TTree in HDF5

Tom Eichlersmith

he/him/his

University of Minnesota

eichl008@umn.edu

CHEP 2023

May 8, 2023

- 1 Helpful Aspects of the TTree
- 2 Why re-implement?
- 3 HDF5 and HDTree
- 4 Conclusion

## Interface to Data

Data viewed hierarchically in-memory while still being well organized on-disk.

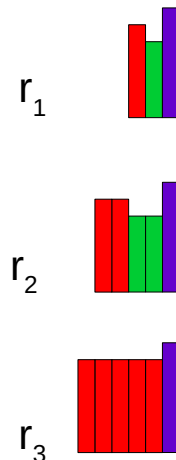
## Schema-Evolution

User-defined classes can evolve with version numbers.

## “Ragged” Data

Directly represent HEP’s common data “awkwardness” where variables change shape from event-to-event

record/row/event



column/branch

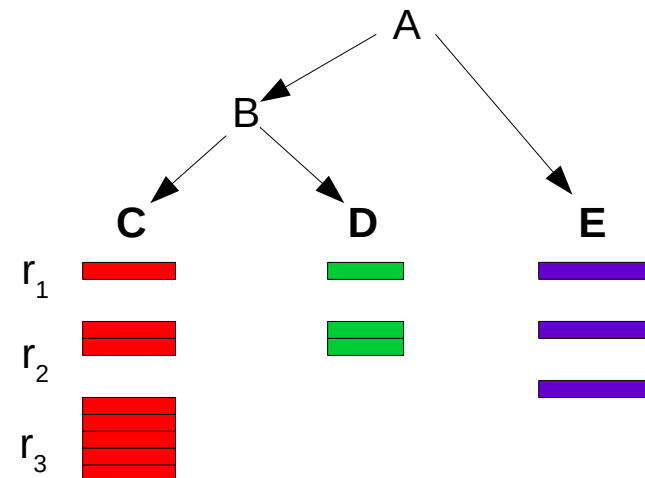


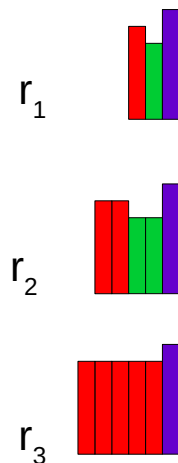
Figure: Credit to J. Pivarski from talk at 2019 CHEP for diagram idea. [▶ Pivarski 2019 CHEP](#) .

# Why re-implement?

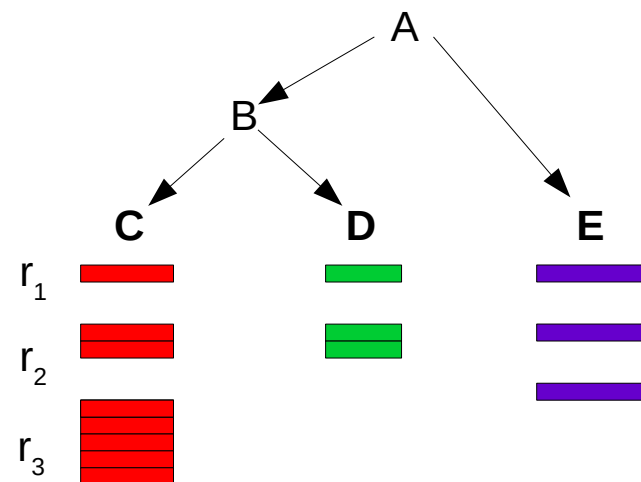


1. Wish to separate the *format* from the *API implementation*
2. Wish to implement *native* APIs in other languages
3. Take lessons learned from using and developing with ROOT – intentionally avoid supporting TTree API to focus on a more modern API

record/row/event



column/branch



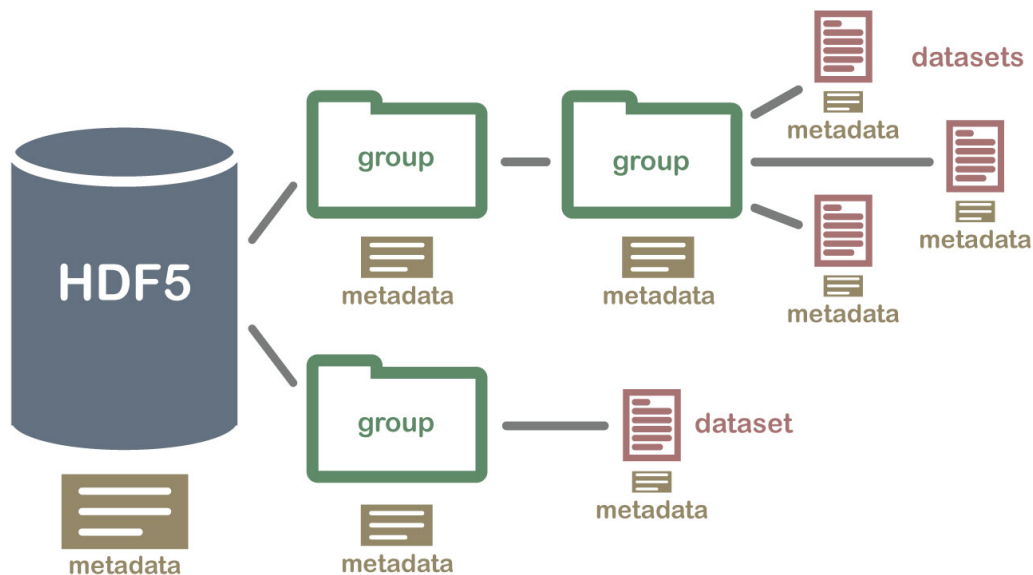


Figure: Diagram of how users can organize data within an HDF5 file. Leah A. Wasser “Hierarchical Data Formats - What is HDF5?” *Neon Science*, Oct 2022. [▶ online](#) .

## Performance

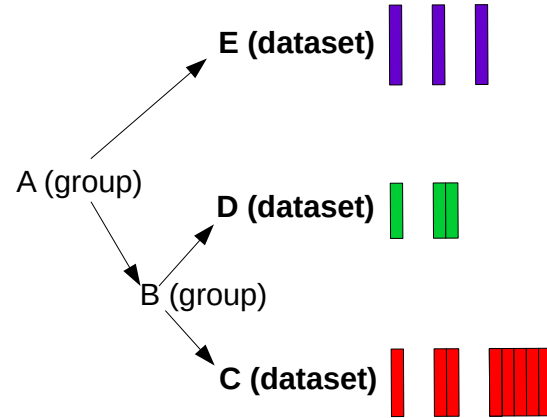
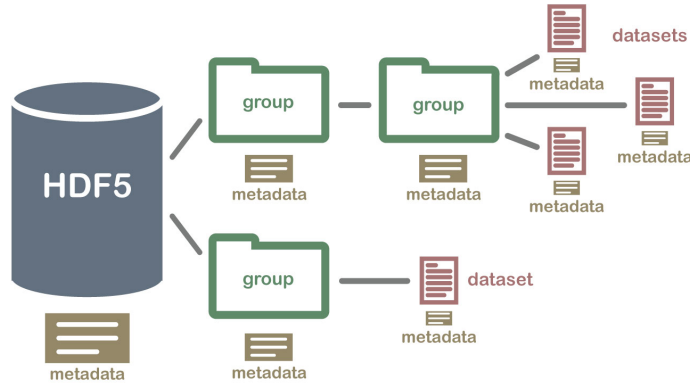
Chunking and compression available to user configuration.

## Usability

C, C++, and Fortran Official APIs - Python, Rust, Julia, and more provided by community

## Support

Industry-supported, popular file format already used within HEP and other research disciplines



1. Atomic data written into HDF5 datasets
2. Description of that data written into HDF5 attributes
3. Datasets organized into groups based on hierarchical in-memory ownership

Specification of HDTree Meta-Format documented online [▶ tomeichlersmith.github.io/hdtree](https://tomeichlersmith.github.io/hdtree)

First API implementation for this meta-format based on C++17 and HighFive<sup>1</sup>

## Capabilities

- Read and write the HDTree meta-format
- Default compression, chunking, and caching behavior similar to ROOT's handling of TTrees
- Schema evolution of user-defined classes

## Distinctions from ROOT

- **No** separate dictionary file required – “binding” user classes to organize data is completely contained in class definition
- **No** user-juggling of memory addresses
- Follows RAII principles – i.e. **No** need to `Close` files or explicitly `Write` trees

---

<sup>1</sup>a header-only, template-focused C++ API for HDF5 – GitHub: [BlueBrain/HighFive](#)



```
1 auto tree = hdtree::Tree::save(  
2     "my-file.hdf5", "/path/to/tree");  
3 auto& i_entry = tree.branch<int>(  
4     "i_entry");  
5 auto& rand_data = tree.branch<  
6     std::vector<float>  
7     >("rand_data");  
8 for (std::size_t i{0}; i < 5; i++)  
9 {  
10     *i_entry = i;  
11     rand_data->resize(i);  
12     for (float& pt : *rand_data) {  
13         pt = i*i;  
14     }  
15     tree.save();  
16 }
```

**Code Task:** Write a tree with two branches: entry index and an array whose size is the index and content is the index squared.

- Create `tree` and new `branches`, accessing branches through handles
- Handles act like smart pointers to underlying data to avoid unnecessary copying
- Final writing and closing done when `tree` is destructed



```
1 class MyData {
2     float x_;
3     friend class hdtree::access;
4     template <typename Branch>
5     void attach(Branch& b) {
6         b.attach("x", x_);
7     }
8     public:
9     hdtree_class_version(2);
10    MyData() = default;
11    MyData(float x)
12        : x_{x} {}
13    void clear() {
14        x_ = 0.;
15    }
16};
```

- Templated `attach` method not necessary but does reduce code boilerplate
- Optional definition of class version – library reports version discrepancies
- `clear` method for resetting class to an “unset” state

Interaction with `Tree` (and underlying `Branch`) the same.

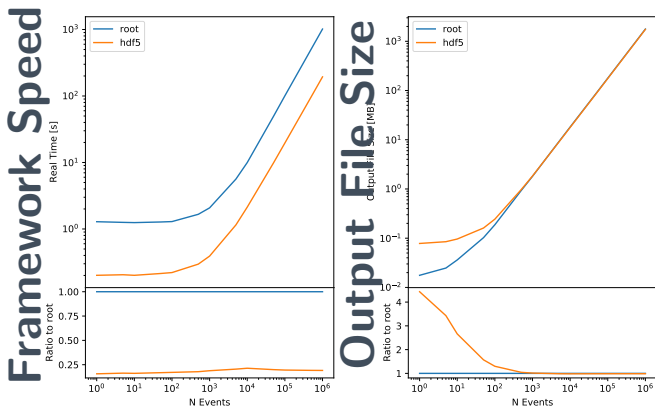
```
1 auto& my_data = tree.branch<std::
    vector<MyData>>("my_data");
```



### Generally – on par with ROOT

- Compared writing and reading of variable-length `std::vector<float>` for a range of entries
- Both close and the faster one depends on computer I run on

Comparison Between HDF5 and ROOT : Production Mode



### Specifics of test highly matters

Embedding this style of serialization into a more generalized event-processing framework led to a vast difference in speed.

HDF5-based Serialization [▶ LDMX-Software/fire](#)

ROOT-based Serialization [▶ LDMX-Software/Framework](#)

## HDTree

- Utilize TTree-like data organization and access patterns while gaining the benefit of an industry-standard file format.
  - C++ API currently available with performance on-par (if not exceeding) ROOT.
- 
- ✓ Implementing specialized APIs for HDF5 files is easier
  - ✓ HDF5 files are supported by other data science libraries (e.g. `pandas` and `pytorch`)
  - ✓ Develop data processing frameworks with less fear of memory issues

## Moving Forward

- More thorough performance testing
- Start Python API based on well-supported and more general `h5py` package