

Writing files with uproot

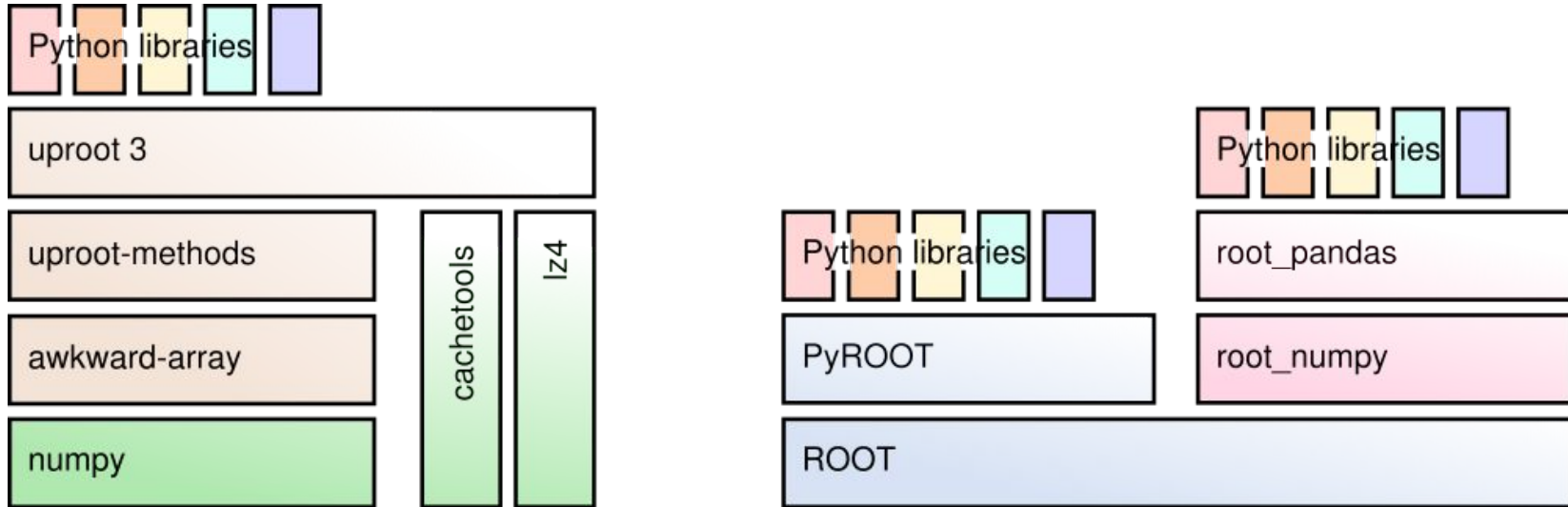
Pratyush Das, Jim Pivarski

- ROOT Users' Workshop - Pratyush Das (10th September)
- LPC Computing Discussion - Jim Pivarski (28th September)

- Python reimplementations of ROOT I/O (until recently, only “input”).
- An array-centric view of ROOT TTree data:
 - branches of simple types are simple arrays
 - branches of complex types are “jagged arrays”
- High performance for large baskets, despite Python’s slowness (because all per-entry operations are performed in Numpy).
- Although originally intended as a temporary replacement for BulkIO, physicists like it for its simplicity: minimal installation, set-up, and affinity with machine learning interfaces.

Layers of dependencies

Each box is a separately-installable package that depends on the ones below it:

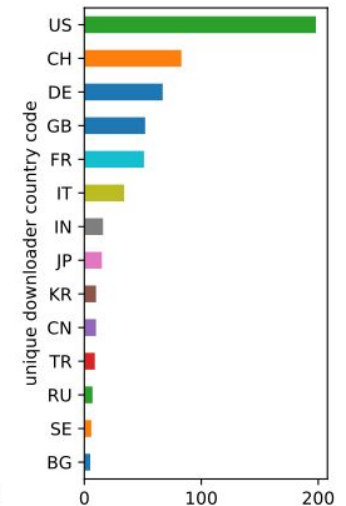
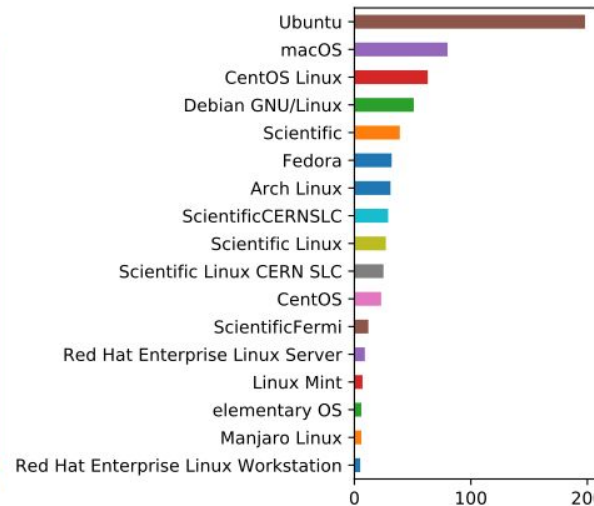
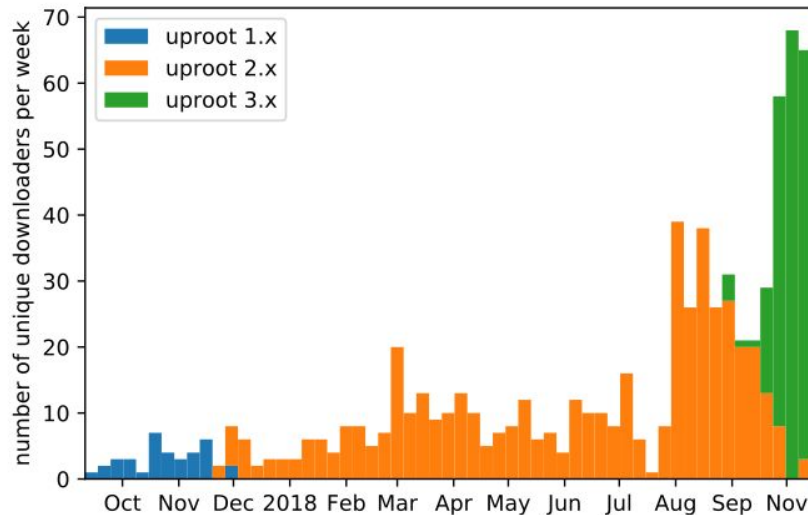
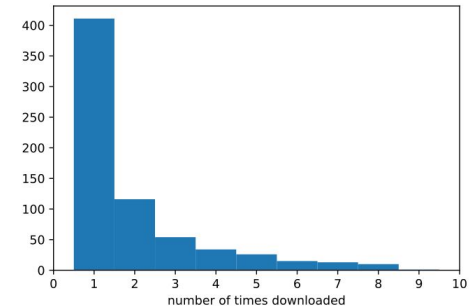


Download statistics

About **732** unique downloaders by pip since version 1.0 was released on September 14, 2017.

(Uniqueness [estimated](#) from country code, OS version, distribution, and release number.)

Primarily U.S. and Linux, but some Europe and MacOS.



Overturn the notion that uproot is to “escape” the ROOT ecosystem, make it a tool for interoperability in both directions.

However, still focus on the end-user physicist:

- Object types useful for analysis: histograms, simple TTrees, TLorentzVectors, and std::vectors for jaggedness.

Not complex types like AODs.

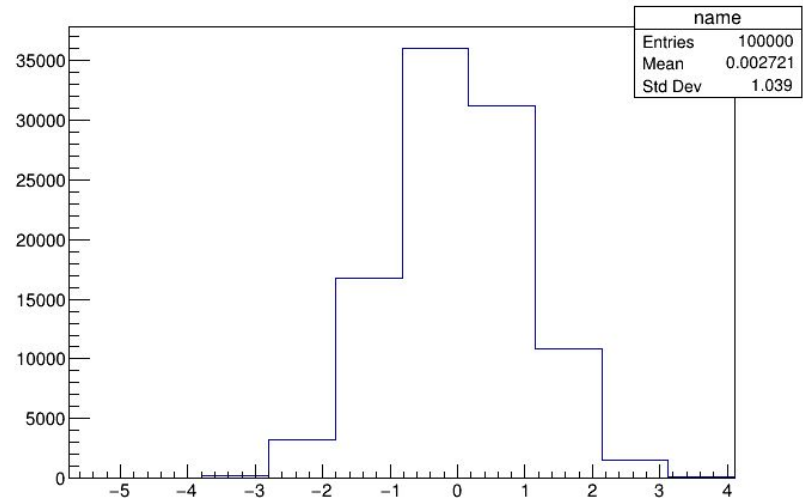
- Optimize for “write once,” not “file as a database” or “partial write recovery.” Single-threaded writing.
- Simple dict-like interfaces wherever possible.

Write a Numpy histogram to a ROOT file (also Pandas, Physt, HEPData YAML...):

```
>>> import uproot, numpy
>>> f = uproot.recreate("tmp.root")
>>> f["name"] = numpy.histogram(numpy.random.normal(0, 1, 100000))
```

Read it back in ROOT:

```
>>> import ROOT
>>> f = ROOT.TFile("tmp.root")
>>> h = f.Get("name")
>>> h.Draw()
```



- Using uproot's read feature to find how data is stored in ROOT files.
- Replacing parts of ROOT files with bytes written by uproot.
- Figuring out which byte strings are relevant for uproot.
- Comparing bytes of file created by uproot and ROOT.

- Module structure
 - Reading and Writing are kept separate
- Standard file handle instead of memmap
- Streamer handling
 - Adding all streamers
- Growing buffer reallocation

Different classes and different modules for reading and writing.

3 layers of abstraction -

1. Objects - Recognised by their methods(Duck typing) and written by a streamer for a single version of the class (corresponding to a recent ROOT release)
2. Primitives - Strings and numbers of particular types (Handled by the Cursor class)
3. Bytes - Interface with the physical layer, such as file handle and in the future xrootd (Handled by the Sink class - writing, Source class - reading)

Considered two different ways to write:

- memmap
 - Easy manipulation of data due to array like structure.
 - Takes a long time to resize - ROOT file size may continuously grow with the addition of new objects.
 - Work around - Grow file exponentially, starting at a high initial size. Disadvantage - Huge ROOT files.
- Standard file handle
 - Need separate pointers at each variable value location.
 - File resizing is ~100 times faster than when using memmap. (Tested)

uproot uses memmap to read files (no resizing, allows parallel reads) and a standard file handle for writing.

Streamer interface for writing is different from reading.

Writer - Static view of streamers

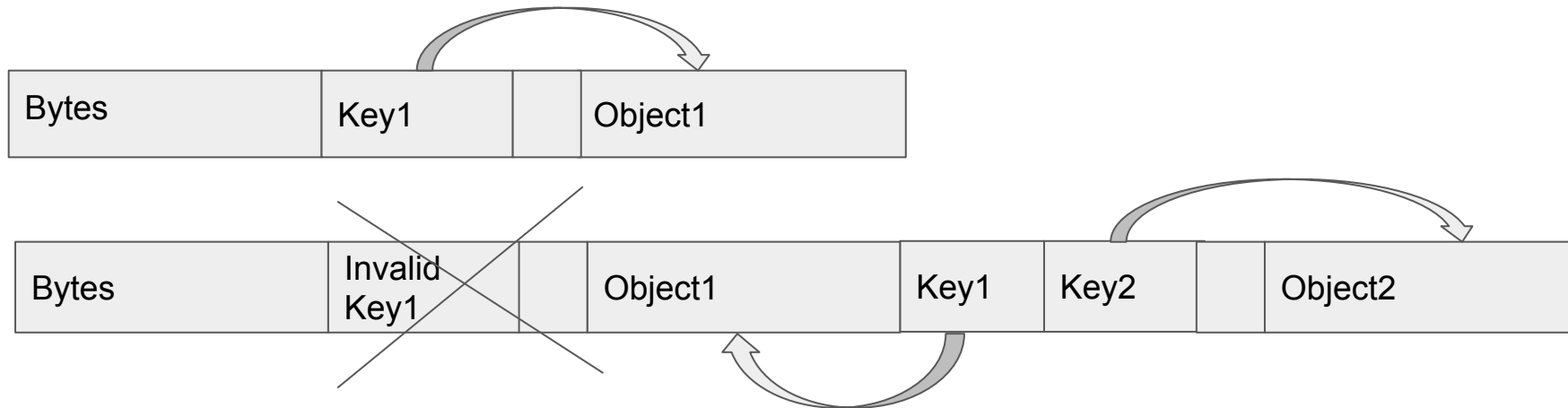
- Design - Writing each streamer as an uninterpreted byte string
 - Only supports a subset of ROOT's own classes
 - Does not generate any new streamers
- Currently, writes a set of commonly used streamers as one big byte string
 - For simplicity
 - Classes supported - TObjString, TH*, TProfile*, simple TTree, TLorentzVector, TVector2, STL collections

Growing buffer reallocation

Common problem when objects in a fixed byte stream grow: what to do if they grow past their allocated bounds?

We create a new (larger) allocation, point the header to the new one, and consider the old one invalid. Unused space grows only logarithmically with #keys.

Example - TKeys in a TDirectory as objects are added:



The following objects can be written on a testing branch of uproot:

- TObjString - Basic object, used as a first test case
- TH1* - Useful for physics; can extend to TH* and TProfile*

Writing TTrees will be supported in 6–12 months.

First-stage TTrees will be flat (numerical types); following that, jagged arrays and object arrays like TLorentzVector may be supported.

- Status
 - Generalising histograms
- Future Plans
 - Compressed objects
 - Support all histogram types
 - Support simple TTrees
 - Generalize streamer writing