# hepfile

**Matt Dreyer** (Cornell)
Matt Bellis (Siena College)
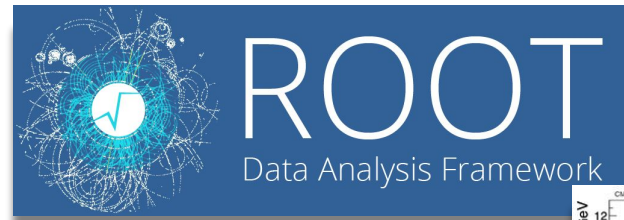
# What is the problem we're trying to solve?

- HEP data is heterogeneous and complicated!
- HEP-native example is particle physics experiment
  - Sensors find muons, electrons, jets, etc.
  - Each particle has specific data attached to it (momentum, charge, etc.)
  - Each **event** (collision) might have different numbers of these particles



- Consider a census of a town, with data gathered per household
  - Each **household** has people, cars, and place of residence
  - Each person has name, gender, and age
  - Each car has age and license plate
  - Each house has # of bedrooms and bathrooms
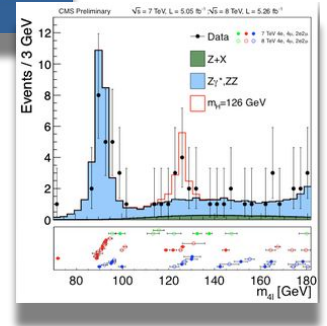
# How to solve it

Current solutions: ROOT!

Problems: Monolithic

PyROOT helps but still tied to a file format that is bound to the analysis toolkit

Makes it difficult to interface with non-HEP people (e.g. broader computing community)



Alternative ROOT-file approaches: uproot/awkward ecosystem

But this is still using the ROOT file format (uproot).

# Early attempt - h5hep

- Package originally called h5hep (*2017*)
- Based on organizing HEP-like data into HDF5 file format
  - Hierarchical Data Format 5
  - Open source
  - Datasets stored in group, stored in groups, etc.
- Used on Particle Physics Playground website for physics pedagogy
  - Data from CMS, BaBar, CLEO, and others
  - File format was originally zipped textfiles





Image credit :
https://www.neonscience.org/resources/learning-hub/tutorials/about-hdf5

# Enter Matt D!

## What we had

- Basics of hepfile
- Could be pip installed
- Basics of ReadThe Docs



## What we wanted to have

- Add some functionality
  - Hold strings
  - Add attributes
- Get documentation better
- Use Henry Schreiner's scikit-cookiecutter to make the package more robust for distribution
- Add functionality to work with awkward
- Make code more robust and fault-tolerant
- Add necessary unit tests
- Get CI working
- Submit to JOSS! *(Journal of Open-Source Software)*

# hepfile is born!

h5hep → hepfile (Summer 2021)

- Define a *schema*
  - How data is organized
  - What metadata needs to be stored to organize the data
- Define minimal *useful* API for flexibility
- *Then*, implement it in
  - Python
  - HDF5
- Hierarchical Data Format 5
  - Open source
  - Groups → datasets
  - Singletons
  - Metadata
- Define structure of two python *dictionaries* to help with packing/unpacking data
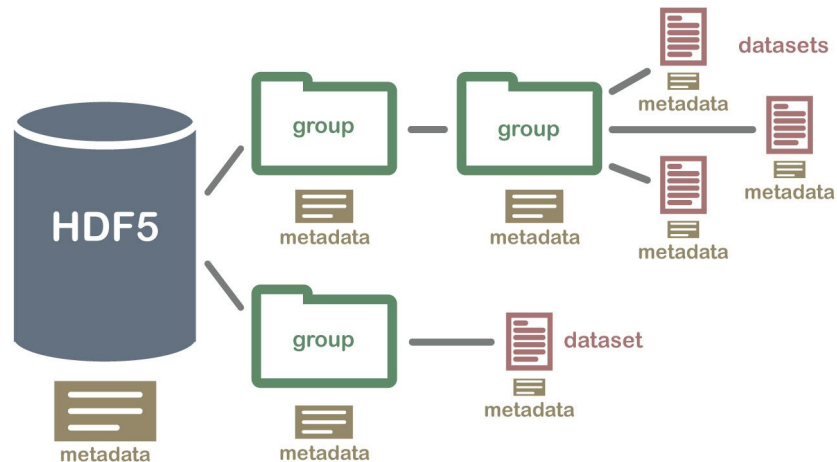  - Analogous to ROOT's TTree/Leaf/Branch



Image credit :
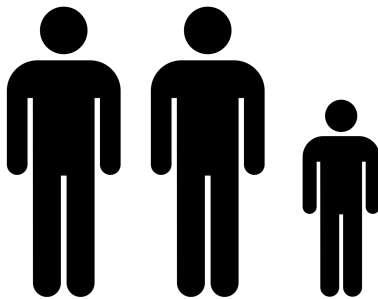https://www.neonscience.org/resources/learning-hub/tutorials/about-hdf5

# What is hepfile?

- Data organized into events/buckets
- hepfile groups data of similar types together in datasets
  - Keep bucket data using 'counter' field in dataset
- Pack takes buckets -> groups and datasets
- Unpack takes groups and datasets -> buckets
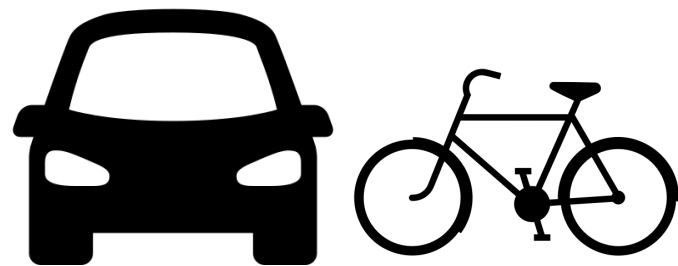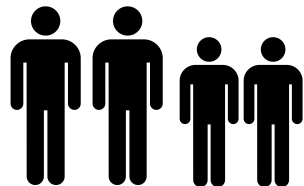  - Extracts specific bucket 'i' from datasets

**Residence**

**People**

**Vehicles**

**Residence:** House, 4, 2.5, 1500, 1955, 250000

**Person 0:** Ollie, Defelice, M, 54, 159, 75000. BS

**Person 1:** Marjorie,Williams,F,52,140,80000,MS

**Person 2:** Tommie,Thoren,NB,18,168,0,12

**Person 3:** David,Haley,F,14,150,0,9
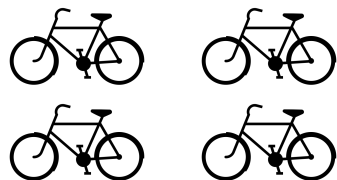
**Vehicle 0:** Car,4,Gas,2005,25000

**Vehicle 1:** Car,5,Electric,2018,40000

**Vehicle 2:** Bike,1,Human,2015,500

**Vehicle 3:** Bike,1,Human,2015,500

**Vehicle 4:** Bike,1,Human,2015,500

**Vehicle 5:** Bike,1,Human,2015,500

**pack**

**Residence:** House, 4, 2.5, 1500, 1955, 250000

**Person 0:** Ollie, Defelice, M, 54, 159, 75000. BS

**Person 1:** Marjorie,Williams,F,52,140,80000,MS

**Person 2:** Tommie,Thoren,NB,18,168,0,12

**Person 3:** David,Haley,F,14,150,0,9

**Vehicle 0:** Car,4,Gas,2005,25000

**Vehicle 1:** Car,5,Electric,2018,40000

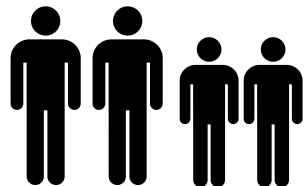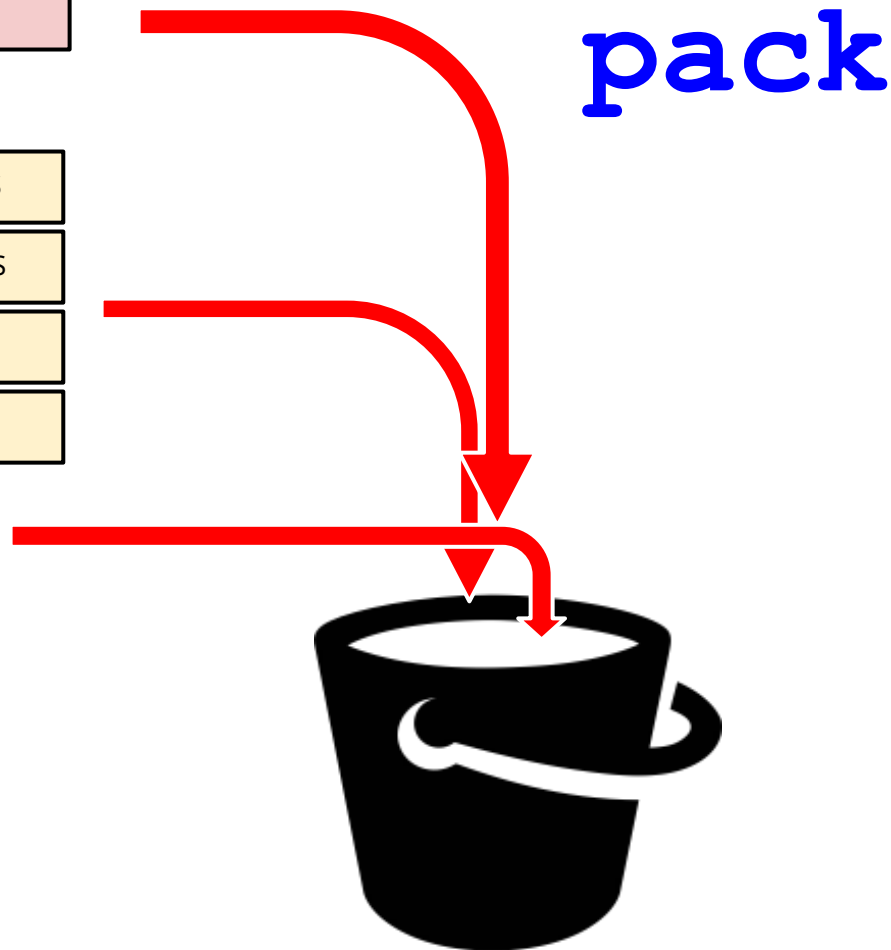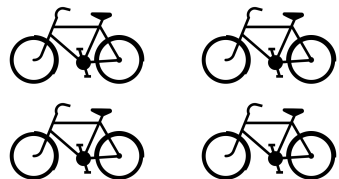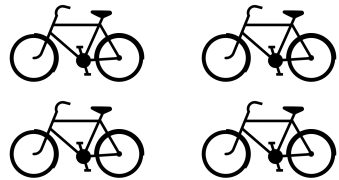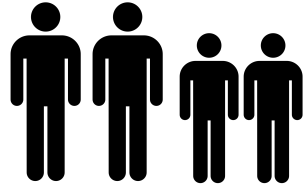**Vehicle 2:** Bike,1,Human,2015,500

**Vehicle 3:** Bike,1,Human,2015,500

**Vehicle 4:** Bike,1,Human,2015,500

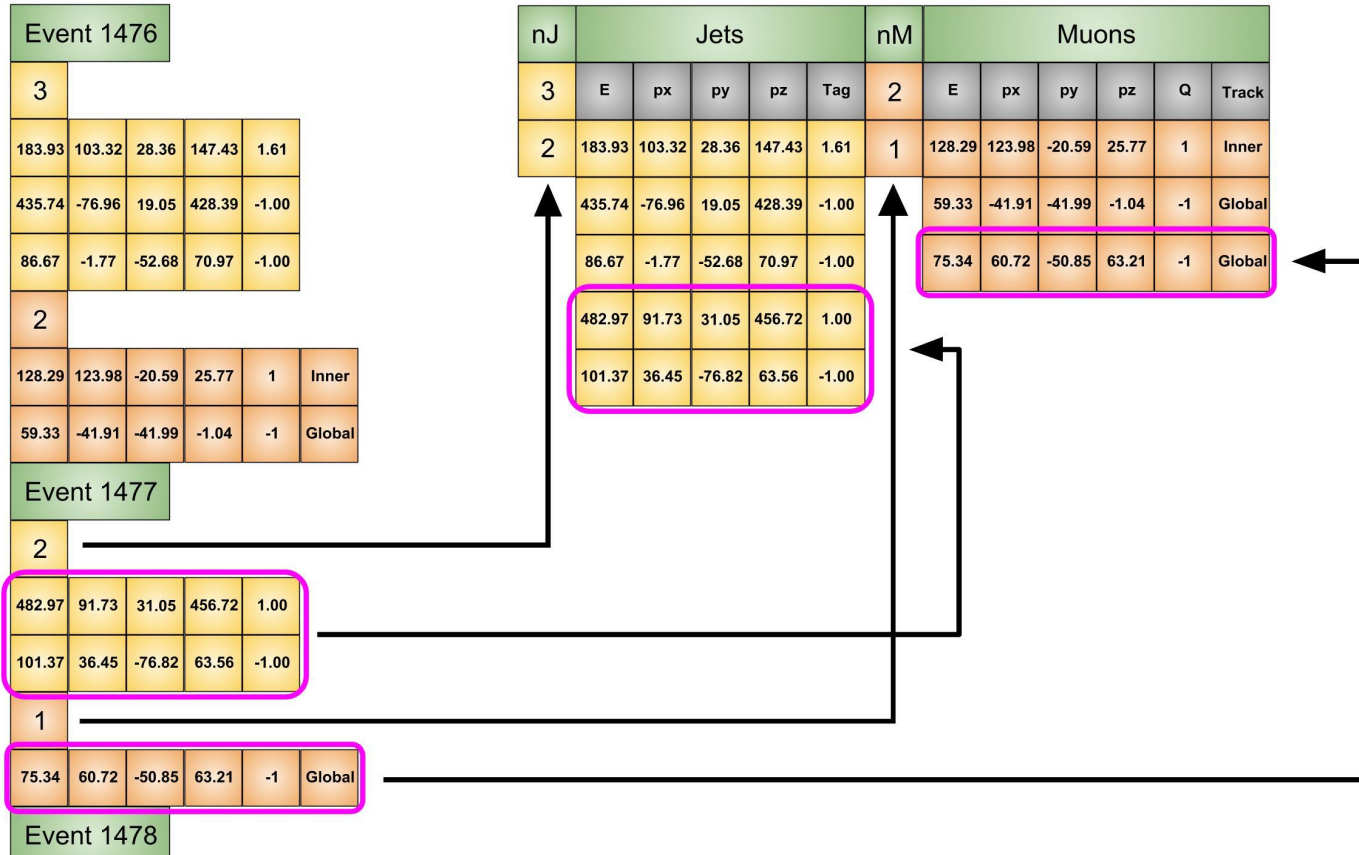**Vehicle 5:** Bike,1,Human,2015,500

pack

# **hepfile** data storage

# hepfile - usage

```python
hepfile.create_group(my_data, 'my_group', counter = 'my_counter')
```

```python
hepfile.create_dataset(my_data, 'my_dataset', group = 'my_group', dtype = str)
```

```python
hepfile.create_dataset(my_data, ['data1', 'data2'] , group = 'my_group')
```

```python
for i in range(5)
    my_bucket['my_group/my_dataset'] = 'yes'
    my_bucket['my_group/data1'] = 1.0
    my_bucket['my_group/data2'] = 2.0

my_bucket['my_unique'] = 3

hepfile.pack(my_data, my_bucket)
```

```python
hepfile.write_to_file('my_file.hdf5', my_data)
```

# hepfile - Read The Docs is getting populated!

## hepfile.read module 🔗

**hepfile.read.calculate_index_from_counters**(*counters*)

**hepfile.read.get_file_metadata**(*filename*)

Get the file metadata and return it as a dictionary

**hepfile.read.get_nbuckets_in_data**(*data*)

Get the number of buckets in the data dictionary.

This is useful in case you've only pulled out subsets of the data

**hepfile.read.get_nbuckets_in_file**(*filename*)

Get the number of buckets in the file.

**hepfile.read.load**(*filename=None, verbose=False, desired_datasets=None, subset=None*)

Reads all, or a subset of the data, from the HDF5 file to fill a data dictionary. Returns an empty dictionary to be filled later with data from individual buckets.

Parameters:
- **filename** (*string*) – Name of the input file
- **verbose** (*boolean*) – True if debug output is required

---

**hepfile.write.pack**(*data, bucket, AUTO_SET_COUNTER=True, EMPTY_OUT_BUCKET=True, STRICT_CHECKING=False, verbose=False*)

Takes the data from an bucket and packs it into the data dictionary, intelligently, so that it can be stored and extracted efficiently. (This is analagous to the ROOT TTree::Fill() member function).

Parameters:
- **data** (*dict*) – Data dictionary to hold the entire dataset EDIT.
- **bucket** (*dict*) – bucket to be packed into data.
- **EMPTY_OUT_BUCKET** (*bool*) – If this is *True* then empty out the *bucket* container in preparation for the next iteration. We used to ask the users to do this "by hand" but now do it automatically by default. We allow the user to not do this, if they are running some sort of debugging.

**hepfile.write.write_file_metadata**(*filename, mydict={}, write_default_values=True, append=True*)

Writes file metadata in the attributes of an HDF5 file

Args: **filename** (string): Name of output file

**mydict** (dictionary): Metadata desired by user

write_default_values (boolean): True if user wants to write/update the

default metadata: date, hepfile version, h5py version, numpy version, and Python version, false if otherwise.

**append** (boolean): True if user wants to keep older metadata, false otherwise.

Returns: **hdoutfile** (HDF5): File with new metadata

# hepfile - Read The Docs is getting populated!

## Writing Data with hepfile

Before anything, we extract the data from the .csv files. (Since *houses* will not be its own group, it is not completely necessary to extract *houses_ID*.)

```python
people = np.loadtxt('sheet1.csv', unpack=True, dtype=str,
                    delimiter=",", comments = '$')[:, 1:]

vehicles = np.loadtxt('sheet2.csv', unpack=True, dtype=str,
                    delimiter=",", comments = '$')[:, 1:]

houses = np.loadtxt('sheet3.csv', unpack=True, dtype=str,
                    delimiter=",", comments = '$')[:, 1:]

people_ID = people[0][1:].astype(np.int32)
vehicles_ID = vehicles[0][1:].astype(np.int32)
houses_ID = houses[0][1:].astype(np.int32)
```

We create the dictionary where we will be storing our data, and then create the groups inside it. For brevity, the **counter** for the buckets will be *ID*. It is fine to repeat the name of the counter because hepfile will store the counter dataset as `f'{groupname}/ID'`.

```python
town = hepfile.initialize()
hepfile.create_group(town, 'people', counter = 'ID')
hepfile.create_group(town, 'vehicles', counter = 'ID')
```

### hepfile
latest

Search docs

Introduction
Fundamentals
Usage
⊟ Examples
  Write to and read from file (generic example)
  Write to file (HEP example)
  Read from file (HEP example)
  Converting .csv files to hepfile
Citation
Contributors

🏠 » Examples       ○ Edit on GitHub

## Examples

### Write to and read from file (generic example)

### Write to file (HEP example)

```python
import numpy as np
import sys
#import hepfile

# For development
sys.path.append('../src/hepfile')
import write as hepfile

data = hepfile.initialize()

hepfile.create_group(data,'jet',counter='njet')
hepfile.create_dataset(data,['e','px','py','pz'],group='jet',dtype=float)
hepfile.create_dataset(data,['algorithm'],group='jet',dtype=int)
hepfile.create_dataset(data,['words'],group='jet',dtype=str)

hepfile.create_group(data,'muons',counter='nmuon')
hepfile.create_dataset(data,['e','px','py','pz'],group='muons',dtype=float)

hepfile.create_dataset(data,['METpx','METpy'],dtype=float)

event = hepfile.create_single_bucket(data)

rando_words = ["hi", "bye", "ciao", "aloha"]

for i in range(0,10000):

    #hepfile.clear_event(event)

    njet = 17
    event['jet/njet'] = njet

    for n in range(njet):
        event['jet/e'].append(np.random.random())
        event['jet/px'].append(np.random.random())
        event['jet/py'].append(np.random.random())
        event['jet/pz'].append(np.random.random())
```
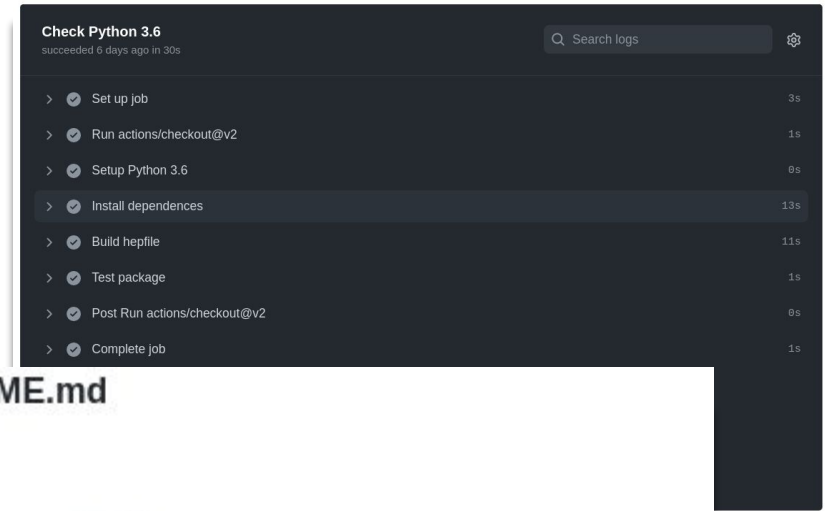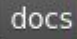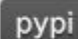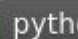
# Unit tests

- *Unit tests* ensure functionality of the program and use cases are considered

- Continuous Integration
  - Evaluates whether every github commit keeps all unit tests working
  - Using **Github Actions** for this

# Installation

*Local install and development*

- Clone from Github
- https://github.com/mattbellis/hepfile

```
git clone https://github.com/mattbellis/hepfile
cd hepfile
flit install
```

*PyPI install*

```
pip install hepfile
```

```
pip install hepfile[awkward]
```

*In development*

# Julia

```julia
using HDF5
using Plots

fname = "output.h5"
fid = h5open(fname, "r")

group_names = keys(fid)

for name in group_names
    println(name)
end

jet = read(fid,"jet")
jet_fields = keys(jet)

for field in jet_fields
  println(field)
end

e = jet["e"]

h = histogram(e,bins=25)

savefig("julia_plot_output.png")

gui()
```



Because we chose a standard underlying file format (HDF5), it makes it easier for other languages to extract data from the file (*assuming there are HDF5 tools written already*)

# JOSS

*"The Journal of Open Source Software (JOSS) is an academic journal (ISSN 2475-9066) with a formal peer review process that is designed to improve the quality of the software submitted. Upon acceptance into JOSS, a Crossref DOI is minted and we list your paper on the JOSS website."*

*"JOSS submissions must:*

- ✔ *Be open source (i.e., have an OSI-approved license).*
- ✔ *Have an obvious research application.*
- ✔ *Be feature-complete (no half-baked solutions) and be designed for maintainable extension (not one-off modifications).*
- ✔ *Minor 'utility' packages, including 'thin' API clients, and single-function packages are not acceptable."*

Apache 2.0

Working on this!

Hope to finish over winter break (Bellis)

# Summary

My perspective on project:

- Learned a lot about package installation
- Learned about Continuous Integration and how useful it is
- Wrote unit tests for full coverage of actual project
- Wrote full documentation for functions
- Clearer perspective on how code should be written for professional use