# The Scikit-HEP project overview and future

**Eduardo Rodrigues, for the Scikit-HEP project admins**
**University of Liverpool**
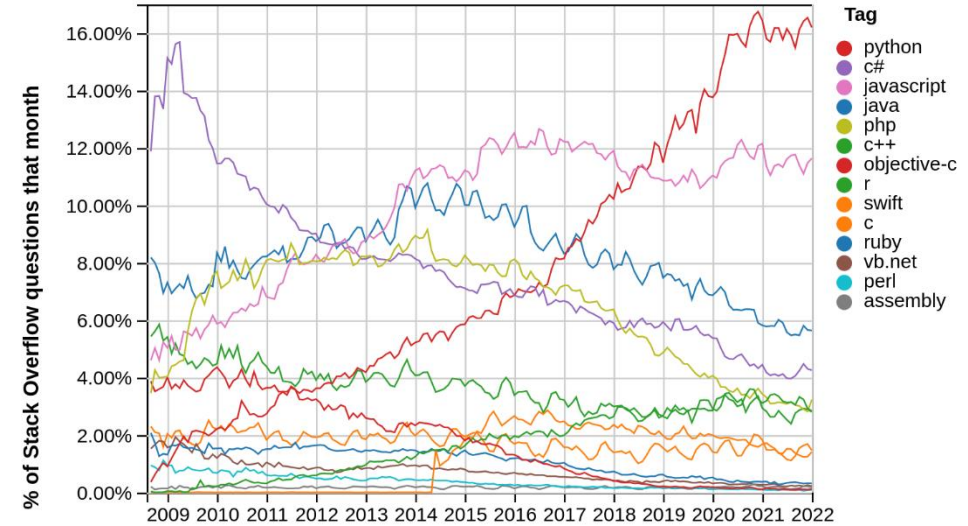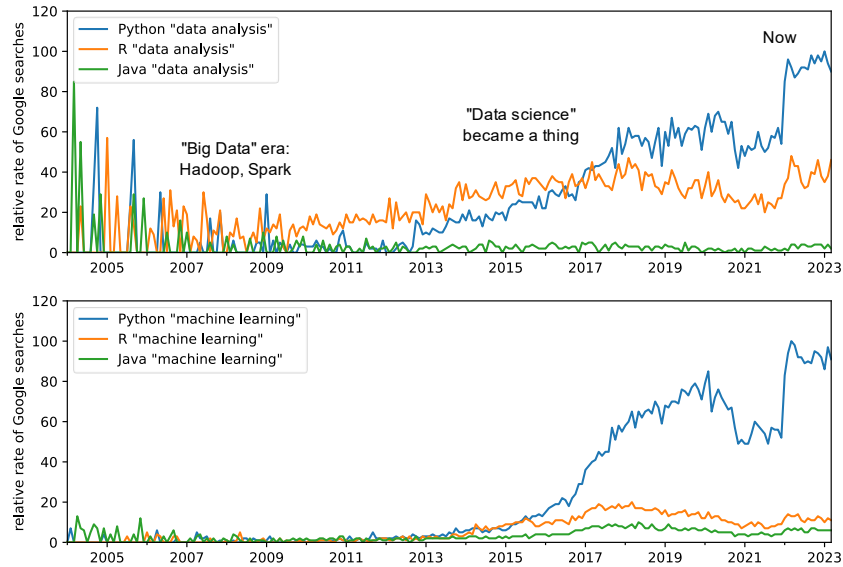
# Historical perspective and project grand idea

- **Python and Pythonic physics analysis in the 2010s**

- **The very early days of Scikit-HEP**

# Recap – some motivation as to why many looked seriously at Python

❑ **Very well-known these days – Python for data analysis was the thing back in the days !**



❑ **Python continues to dominate a decade down the line !**

   - **Example from the TIOBE index:**

| May 2024 | May 2023 | Change | Programming Language | Ratings | Change |
|---|---|---|---|---|---|
| 1 | 1 | | Python | 16.33% | +2.88% |
| 2 | 2 | | C | 9.98% | -3.37% |
| 3 | 4 | ⌃ | C++ | 9.53% | -2.43% |
| 4 | 3 | ⌄ | Java | 8.69% | -3.53% |
| 5 | 5 | | C# | 6.49% | -0.94% |
| 6 | 7 | ⌃ | JavaScript | 3.01% | +0.57% |

# Python and Pythonic physics analysis in the 2000-2010s

- ❑ **In HEP, Python was being used**
  **by some early and enthusiastic adopters**

- ❑ **"Python emerging as standard glue"**

- ❑ **PyROOT official bindings coming to life around 2004 …**

- ❑ **Though the scientific Python ecosystem around NumPy**
  **was already reasonably mature,**
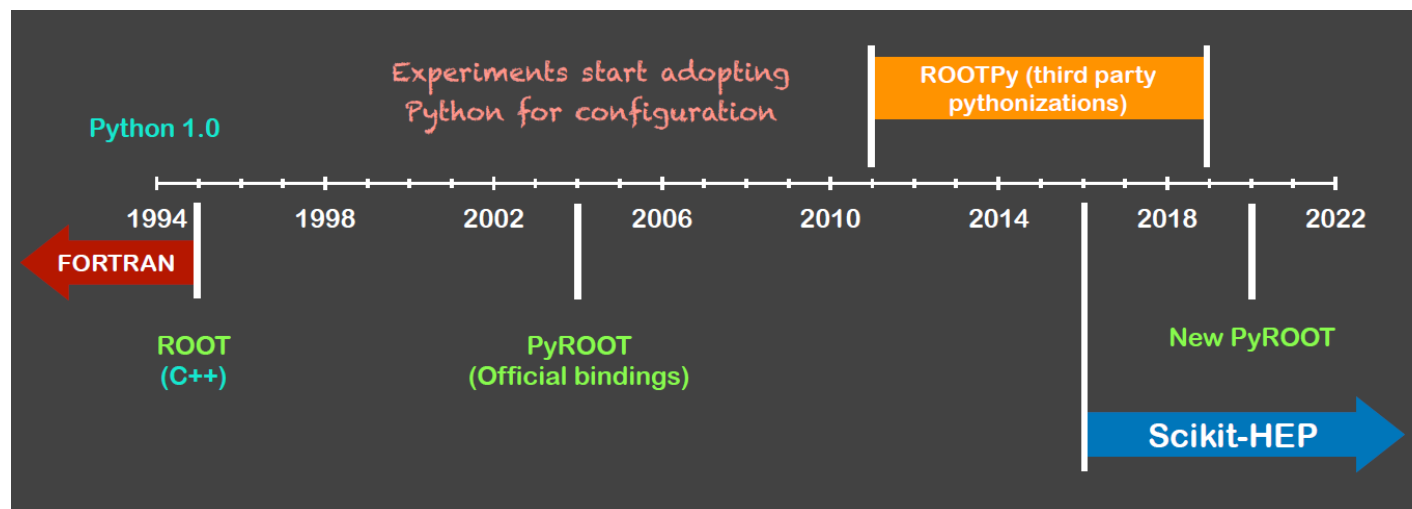  **not much of that was popular with HEP analysts**



**Emerging Standard ?**
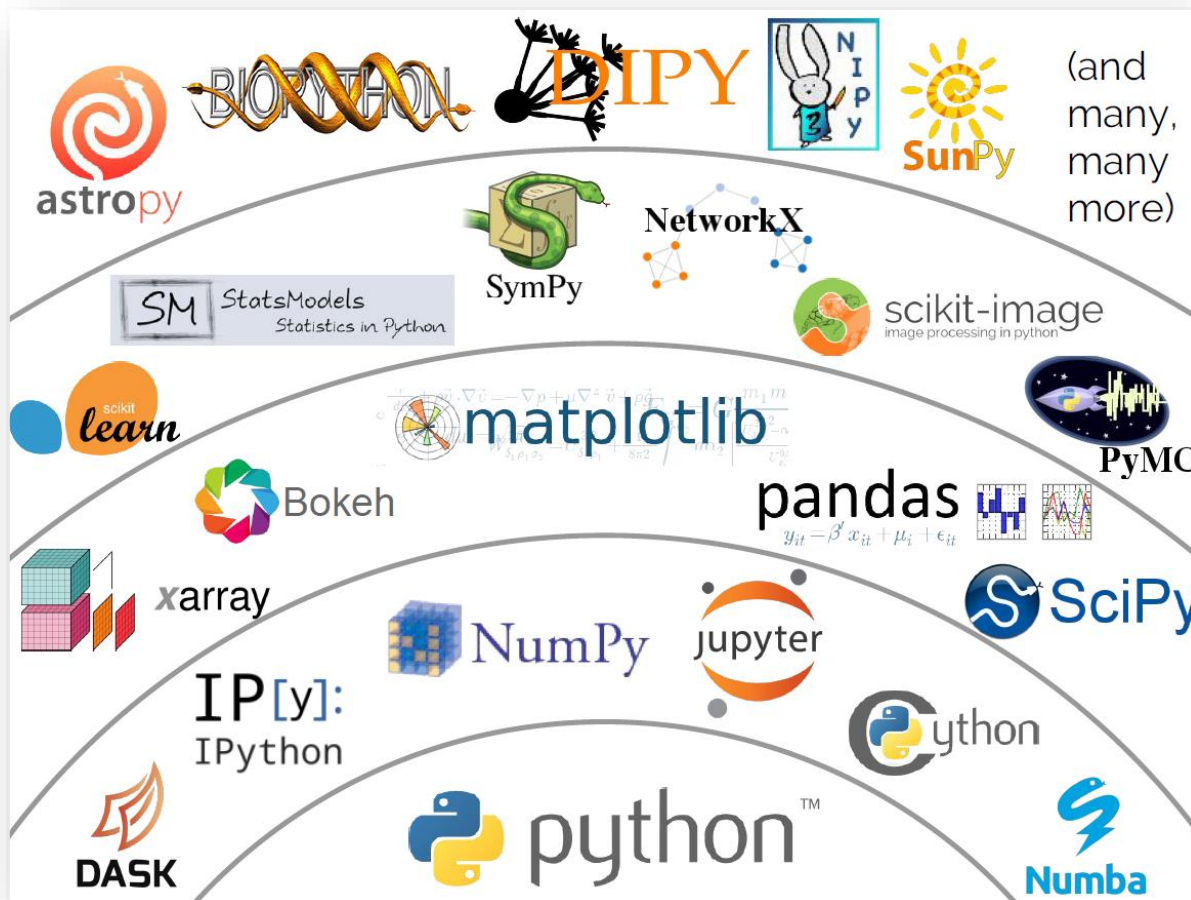**Python as "Software Glue"**

**CHEP 2001 !**

- ■ Clear trend towards Python
  - ❖ Used by: ATLAS (Athena),CMS, D0, LHCb (Gaudi), SND,...
  - ❖ Used by: Lizard/Anaphe, HippoDraw, JAS (Jython)...
  - ❖ Architecturally, scripting is "just another service"
  - ❖ ROOT is the exception to the "Python rule"
    - ➢ CINT interpreter plays a central role
    - ➢ Developers and users seem happy
- ■ Python is popular with developers...
  - ❖ Rapid prototyping;   gluing together code
  - ❖ (Almost) auto-generation of wrappers (SWIG)
- ■ ...but acceptance by users not yet proven
  - ❖ Another language to learn, syntax,...

*"Summary of Track 2: Data Analysis and Visualis. Lucas Taylor, Northeastern U., CHEP 01, Beijing, 3-7 S*



Experiments start adopting
Python for configuration

ROOTPy (third party pythonizations)

Python 1.0

1994   1998   2002   2006   2010   2014   2018   2022

FORTRAN

ROOT
(C++)

PyROOT
(Official bindings)

New PyROOT

Scikit-HEP

❑ **Famous "shell diagram" of the Python scientific ecosystem outside HEP, back in the days:**

❑ **What about HEP? Why was HEP not present back in 2017?**

❑ **Sort of a philosophical question at this point. But there was a niche …**

*Domain-specific*

*Python's*

*Scientific*

*stack*

Jake VanderPlas,
*The Unexpected Effectiveness
of Python in Science,*
PyCon 2017

# The very early days of Scikit-HEP

❑ **Scikit-HEP started with the goal to:**

- Improve the interoperability between HEP tools and the scientific Python ecosystem

- Build a community of developers and users – be community-driven and community-oriented

- Improve the discoverability of relevant tools

- …

❑ **Project started around concept of pillars: datasets, data aggregation, modelling, visualisation and simulation (and "utilities")**
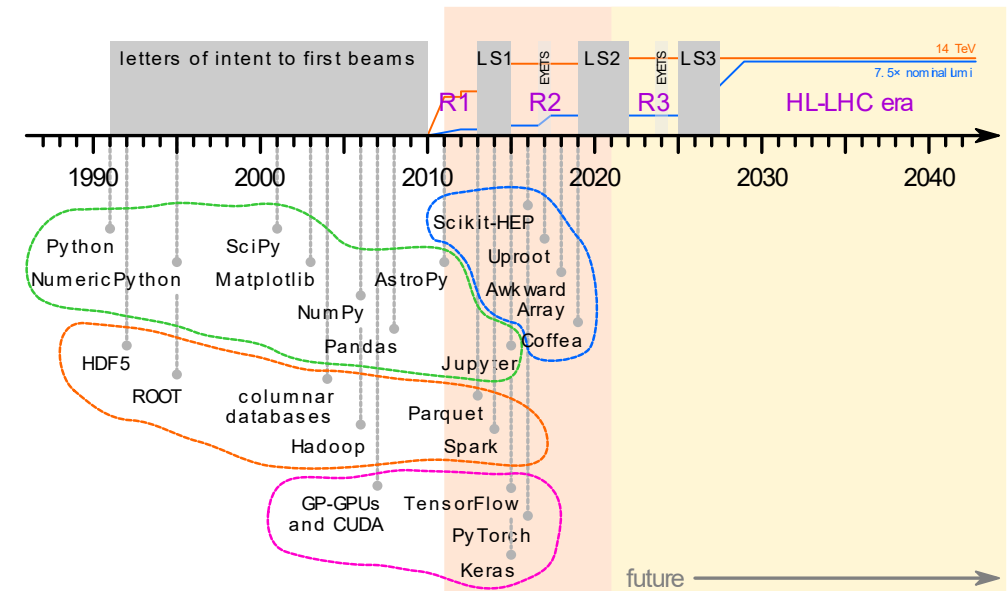
❑ **Initial approach quickly evolved towards a toolset, not a toolkit:**



**Henry Schreiner et al., "Awkward Packaging: building Scikit-HEP" ,SciPy 2022**



**From first public talk on Scikit-HEP, in 2017**

# Some biased aspects of the timeline of Scikit-HEP

❑ **Much of the initial challenge / work was (1) attract people and (2) attract relevant packages**

❑ **Started off with inspiration from LHCb's Ostap package and Noel Dawe's rootpy,**
   **and putting together in a GitHub organisation packages root_numpy, then numpythia, pyjet, iminuit, etc.**

❑ **Uproot came out of conversations with Noel about root_numpy's installation difficulties**

❑ **uproot got showcased at CHEP 2018 and was a game-changer – for Python analysis in HEP and Scikit-HEP also**
   - Loading/saving HEP data was *the* technical bottleneck that required a full ROOT installation. uproot made the barrier the lightest possible

❑ **So was awkward array, as iminuit a few years before (used a lot in Astro)**

❑ **By CHEP 2018 (July 2018), the GitHub org was 8 people and 12 repositories**

❑ **Early 2019 – key developments on histogramming (boost-histogram)**

❑ **End of 2019 Belle II and CMS were both customers**
   **Also Flavio, a Flavour Physics phenomenology software project**

❑ **Community projects Coffea and zfit had been created in the meantime**

❑ **Many packages got archived along the way**
   - No worries, they got superseded by better products !

❑ **Fast-forward to 2024 – around 30 packages maintained**

# Scikit-HEP project overview

- **Not a description of all packages**

- **Not attempting to be comprehensive in any aspect**

- **The following catches your attention?
  Feel encouraged to look further at the repos … ☺**

# Topics and overview of main packages

| Basics |
|---|
| Awkward(array) |
| Vector |
| hepunits |

| Data manipulation & interoperability |
|---|
| uproot |
| uproot-browser |
| hepconvert |

| Histogramming |
|---|
| boost-histogram |
| Hist |
| histoprint |
| UHI |

| HEP specific libraries and interfaces to HEP libraries |
|---|
| Particle |
| DecayLanguage |
| fastjet |
| pylhe |
| pyhepmc |

| Fitting and Statistics |
|---|
| iminuit |
| pyhf |
| cabinetry |
| resample |
| hepstats |

| Visualisation |
|---|
| mplhep |

# Basics, data manipulation, interoperability

| Basics | Data manipulation & interoperability |
|---|---|
| Awkward(Array) | uproot |
| Vector | uproot-browser |
| hepunits | hepconvert |

❑ **Data manipulation is key**

❑ **So is interoperability to leverage efficiently both HEP specific tools and tools from the wider scientific ecosystem (based on NumPy)**

❑ **Uproot then awkward "filled the gap" and acted as an enabler / catalyst - uproot got showcased at CHEP 2018**

❑ **uproot:**
**library for reading and writing ROOT files in pure Python and NumPy**

❑ **Awkward:**
**library for nested, variable-sized data, including arbitrary-length lists, records, mixed types, and missing data, using NumPy-like idioms**

❑ **Vector:**
**arrays of 2D, 3D, and Lorentz vectors**
**(very recently: added a SymPy backend atop "standard backends")**

# Histogramming

| Histogramming |
| :---: |
| boost-histogram |
| Hist |
| histoprint |
| UHI |

❑ **No HEP analysis without histogramming ;-) !**

❑ **boost-histogram – Python bindings to C++ Boost.Histogram triggered a lot of activity …**

❑ **Came:**

    **Hist**, an analyst-friendly front-end for boost-histogram

    **UHI**, a protocol package with interoperability in mind (see next slide)

❑ **The "histogramming suite" has been one of the major successes**

    - **Illustration of the level of developer activity:**

```
import numpy as np

from hist import Hist

Hist.new.Reg(50, 1, 2).Double().fill(np.random.normal(1.5, 0.3, 10_000))
```

Regular(50, 1, 2, label='Axis 0')

Double() Σ=9039.0 *(10000.0 with flow)*

# Histogramming

| Histogramming |
| --- |
| boost-histogram |
| Hist |
| histoprint |
| UHI |

❑ **UHI – Universal Histogram Interface:**
**Primarily for documenting histogram indexing and the PlottableHistogram Protocol and any future cross-library standards**
**⇒ packages can collaborate thanks to agreeing on a shared protocol !**

❑ **Protocol followed by other org packages:**
**boost-histogram, Hist, mplhep, histoprint, uproot**

### PlottableHistogram

Defines expected behaviors for producers

Consumers expect these behaviors

# Fitting and Statistics

| Fitting and Statistics |
| --- |
| iminuit |
| pyhf |
| cabinetry |
| resample |
| hepstats |

## *iminuit*

❑ **Builds on the Minuit2 C++ code in ROOT**

❑ **Extensively used by other packages such as zfit, and also outside HEP, in Astro in particular**

❑ **Thorough unit testing has been revealing bugs in the C++ Minuit2,**
   **for which fixes upstream to ROOT have been submitted – iminuit also a gain for ROOT**

## *resample*

❑ **Originally developed by David Saxton who works in the private sector as a statistician, "taken over" by us**

❑ **One of few libraries that implements the "extended bootstrap" technique**

## *Some future directions*

❑ **Working to adopt and contribute to the HS$^3$ statistical representation to try to really unite the whole field**
   **on a common open statistical format**

❑ **Excellent to see a collaboration from**

HS$^3$
High Energy Physics
Statistics Serialization Standard

# Fitting and Statistics

***pyhf, cabinetry***

❑ **pyhf has been a big hit ! The more that 100 citations attest this**

❑ **Also, widespread adoption of its use across multiple subfields (1) given it being Python and (2) given it provides the ability for open statistical model publication**

**Example of a full HistFactory DAG visualisation:**

# Fitting and Statistics

***Future of pyhf as as "building block"***

❑ **Working to become a tool that integrates well with others and/or can be leveraged**
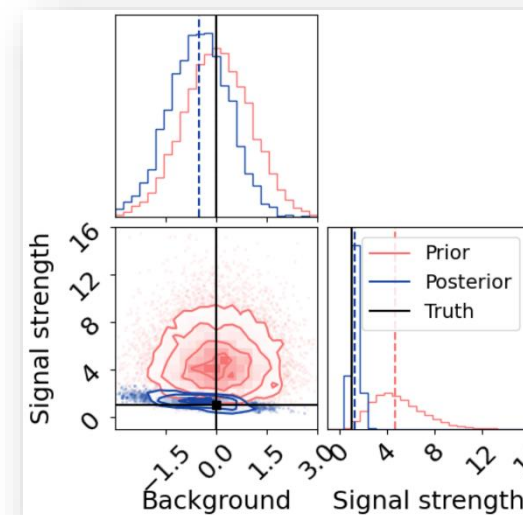
❑ **Specific examples:**

- <u>Paper</u> **"Bayesian Methodologies with pyhf":**

   *"<u>bayesian_pyhf</u> is a Python package that allows for the parallel Bayesian and frequentist evaluation of multi-channel binned statistical models. The Python library pyhf is used to build such models according to the HistFactory framework ..."*

- <u>spey-pyhf</u> **plugin for inference in phenomenology**

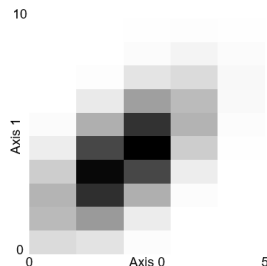**Smooth inference for reinterpretation studies**



(a) Correlation for signal and background distribution. Indicated in black is the underlying truth.

Taken from Paper "Bayesian Methodologies with pyhf"

# Visualisation

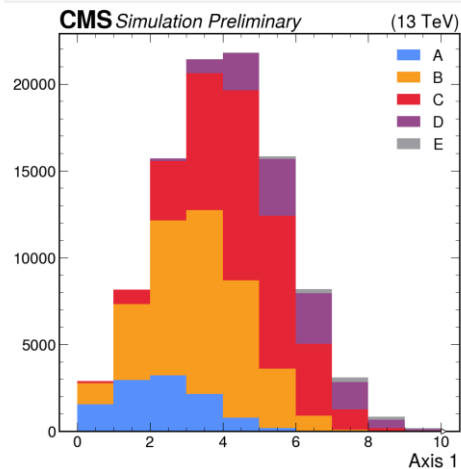❏ **mplhep is (only!) about easy matplotlib histogramming / styling for HEP plots**

- **Convenient user interface for histogram plotting**

- **Interface to matplotlib, fully integrated with Hist**

- **Can also ingest ROOT TH1/TH2**



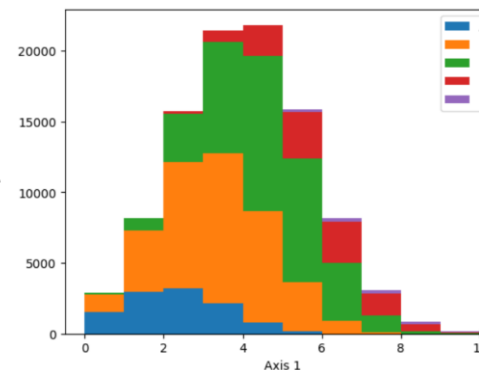Publication quality plots conforming to collaboration requirements (ATLAS, CMS, LHCb, ALICE) can be produced with 2 lines of code

For CMS - CVD friendly color scheme (M. Petroff colors)

# HEP domain-specific libraries

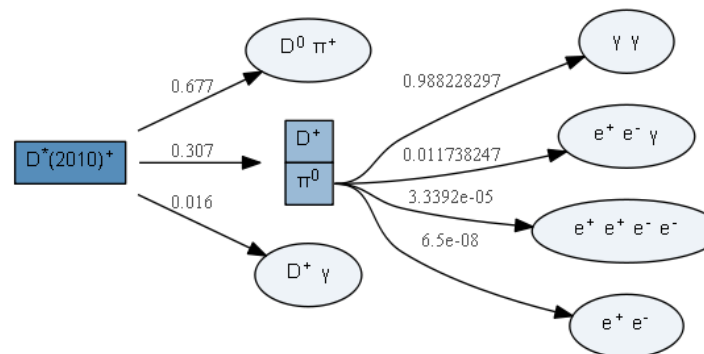| HEP specific libraries and interfaces to HEP libraries |
| :---: |
| Particle |
| DecayLanguage |
| fastjet |
| pylhe |
| pyhepmc |

### pylhe

```python
# Use an example LHE file from package scikit-hep-testdata
from skhep_testdata import data_path

lhe_file = data_path("pylhe-drell-yan-ll-lhe.gz")

arr = pylhe.to_awkward(pylhe.read_lhe_with_attributes(lhe_file))

arr
```

```
[{eventinfo: {nparticles: 4, pid: 1, ...}, particles: [...]},
 {eventinfo: {nparticles: 5, pid: 1, ...}, particles: [...]},
```
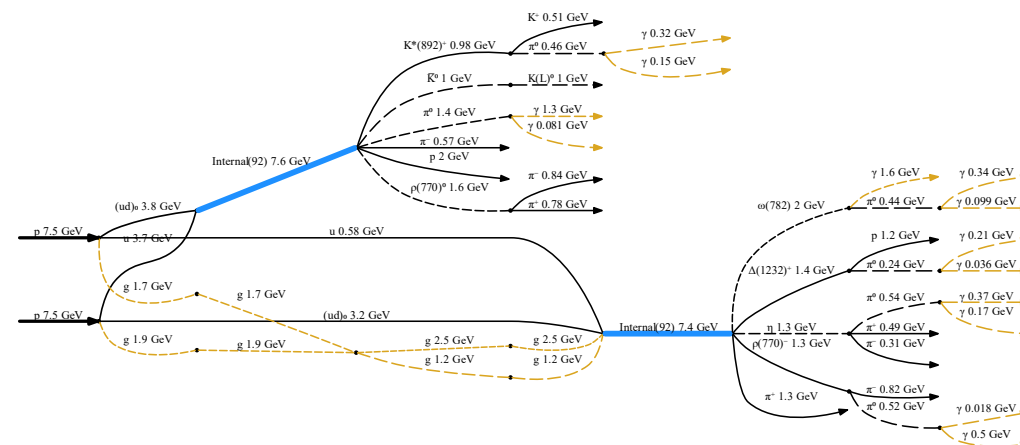
### DecayLanguage



### fastjet

Python bindings for the C++ FastJet library.
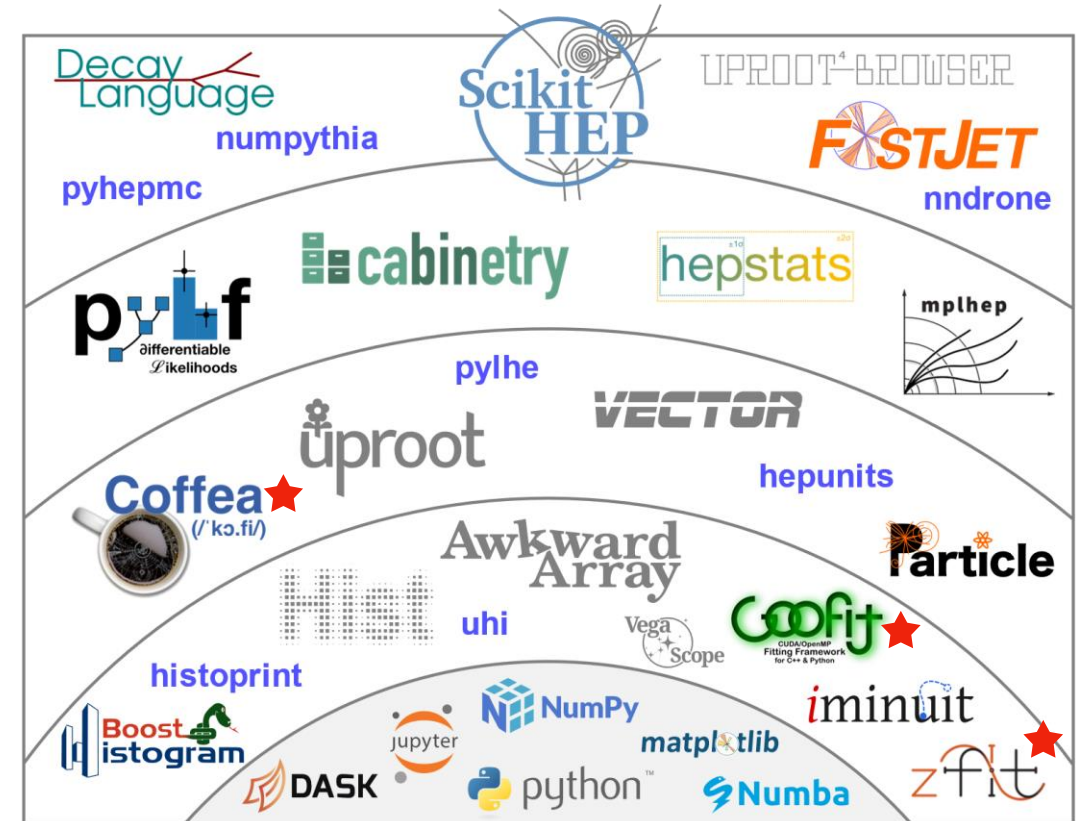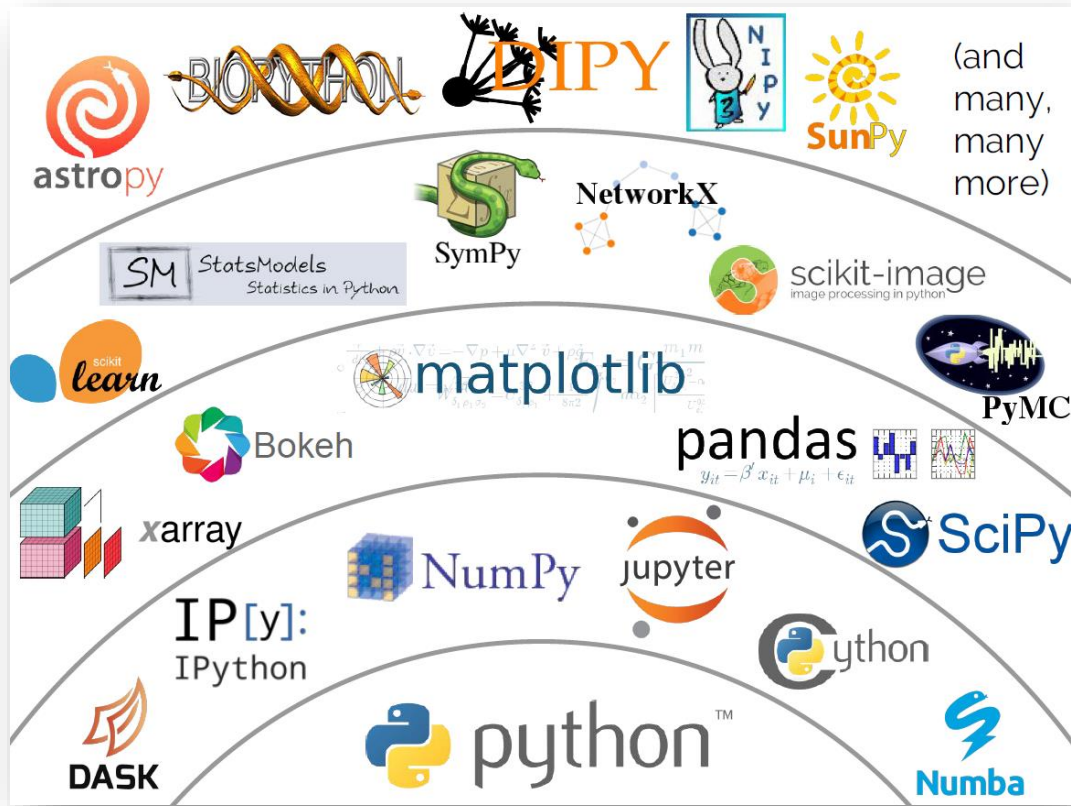Contains 2 interfaces – the Classic and the Awkward(array)

**Find all pseudoscalar charm mesons known to the PDG:**

```python
Particle.findall(lambda p: p.pdgid.is_meson and p.pdgid.has_charm and p.spin_type==SpinType.PseudoScalar)
```

```
[<Particle: name="D+", pdgid=411, mass=1869.65 ± 0.05 MeV>,
 <Particle: name="D-", pdgid=-411, mass=1869.65 ± 0.05 MeV>,
 <Particle: name="D0", pdgid=421, mass=1864.83 ± 0.05 MeV>,
 <Particle: name="D~0", pdgid=-421, mass=1864.83 ± 0.05 MeV>,
 <Particle: name="D(s)+", pdgid=431, mass=1968.34 ± 0.07 MeV>,
 <Particle: name="D(s)-", pdgid=-431, mass=1968.34 ± 0.07 MeV>,
 <Particle: name="eta(c)(1S)", pdgid=441, mass=2983.9 ± 0.5 MeV>,
 <Particle: name="B(c)+", pdgid=541, mass=6274.9 ± 0.8 MeV>,
 <Particle: name="B(c)-", pdgid=-541, mass=6274.9 ± 0.8 MeV>,
 <Particle: name="eta(c)(2S)", pdgid=100441, mass=3637.5 ± 1.1 MeV>]
```
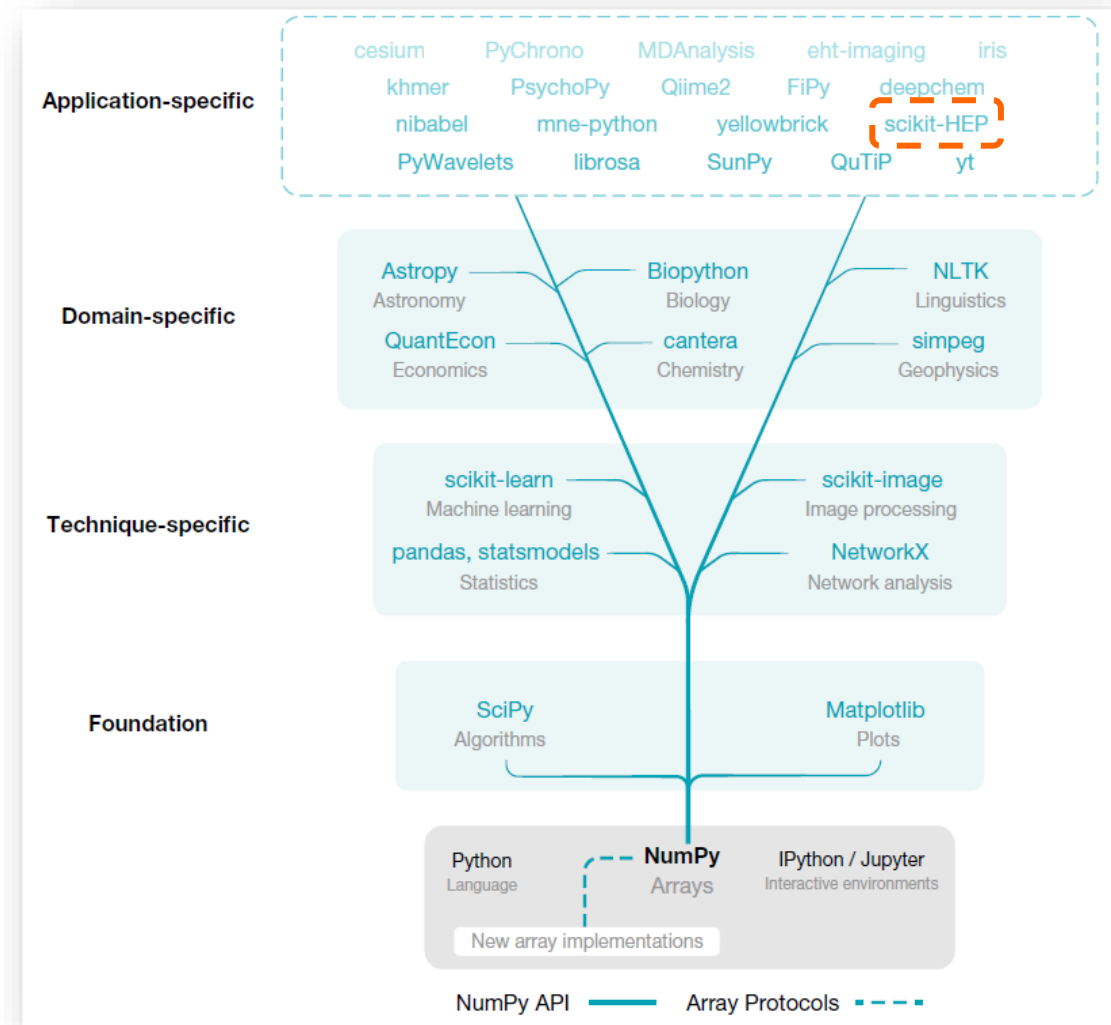
### Particle

### pyhepmc

# Our contribution to the HEP domain-specific Python analysis ecosystem

❑ **After 7-ish years we can proudly splash our "NumPy fueled" shell ecosystem side-by-side ☺ !**

    - **The full HEP ecosystem is of course wider, ROOT being prominent, in particular**

# Our contribution to the HEP domain-specific Python analysis ecosystem

☐ **It was a good surprise to see Scikit-HEP in the "Array programming with NumPy" Nature paper in 2020!**
**- Recognition that "HEP is out there as well"**

**(It should be under "Domain-specific", but fine ☺.)**



**Array programming with NumPy, Nature 585, 357 (2020)**

# HEP community engagement

- **Recognition of contributions and contributors**

- **Training material**

- **Community events**

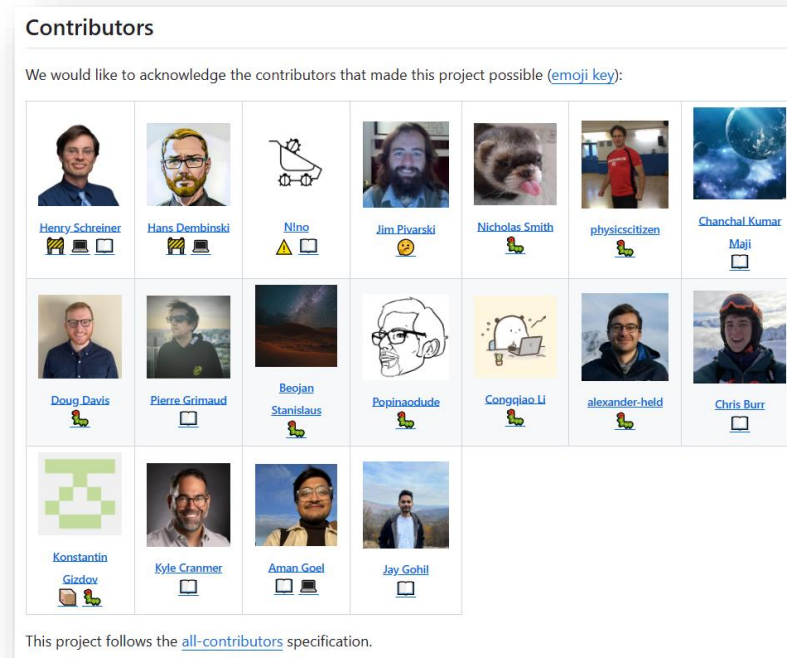# Engagement with community also means recognition of contributors !



- ❑ **That seems totally natural to be "bijective" in matters of recognition and citation**

- ❑ **Seems effectively compulsory for community projects !**

- ❑ **Contributions to software are easily found in GitHub and GitLab repositories**
  **- Example in https://github.com/scikit-hep/boost-histogram**

- ❑ **But there are means to highlight contributions even more**

- ❑ **Our guidelines – use the all-contributors specification**



  **- Caveat: only works on Markdown (README) files, not ReST**

# Engagement via training and training material

❑ **Training material on Scikit-HEP packages contributed to the HSF (Training WG) courses**

❑ **More recently, engagement also in training in general with an HSF Training WG convenership**

❑ **Next event "HSF/IRIS-HEP Python for Analysis Training" to take place on June 5th - https://indico.cern.ch/event/1408846/**

   **- Worth emphasising fact – it's community colleagues who take time to deliver these courses on ROOT, Scikit-HEP, etc. !**



## Scikit-HEP Tutorial

Welcome!

This tutorial aims to demonstrate how to quickly get started with Scikit-HEP, a collection of packages for particle physics analysis in Python.

The tutorial was written by Jim Pivarski and was first taught during a Software Carpentry Workshop on December 15, 2021.

☀ Prerequisites
- Basic Python knowledge, e.g. through the Software Carpentry Programming with Python lesson

📌 HSF Software Training

This training module is part of the **HSF Software Training Center**, a series of training modules that serves HEP newcomers the software skills needed as they enter the field, and in parallel, instill best practices for writing software.

**https://github.com/hsf-training/hsf-training-scikit-hep-webpage**

# Engagement via workshops / conferences

**General events**

❑ **Scikit-HEP gets represented regularly**
- We engage proactively
- Also in non-HEP events such as the SciPy Conferences

**PyHEP workshops**

❑ **Crucial series of workshops to exchange on "Python in HEP"**

❑ **In 2023 a <u>first PyHEP.dev "developers workshop"</u> took place**

❑ **Preparations have been done via a dedicated repository <u>https://github.com/HSF/PyHEP.dev-workshops</u>**

❑ **List of community libraries / projects presented is very large!**

## Library maintainers who will be present

Do you want to connect with the maintainers of a particular library, but don't know who they are? The following people and the packages they maintain are listed below. If the following is incorrect or incomplete, please help us fix it with a pull request. Thank you!

- scikit-build, improved build system generator for CPython C, C++, Cython and Fortran extensions: @henryiii
- cookie, Scientific Python library development guide and cookie-cutter: @henryiii
- SkyHookDM (website), tabular data management in object storage, now part of the Apache Arrow project: @JayjeetAtGithub
- ServiceX (website), a data delivery service pilot for IRIS-HEP DOMA: @talvandaalen, @masonproffitt, @gordonwatts
- hepfile, general file schema defined to handle heterogeneous data but implemented in Python and HDF5: @mattbellis
- Uproot, ROOT I/O in pure Python and NumPy: @ioanaif, @Moelf, @jpivarski
- pylhe, lightweight Python interface to read Les Houches Event (LHE) files: @lukasheinrich, @matthewfeickert
- ATLAS PHYSLITE, a new reduced common data format in ROOT files for ATLAS: @nikoladze, @gordonwatts
- Pythia (website), a program for generating high-energy physics collision events, high energy collisions between electrons, protons, photons, and heavy nuclei: philten
- fastjet, wrapper for FastJet in the Scikit-HEP ecosystem: @rkansal47
- GNN Tracking, charged particle tracking with graph neural networks: @klieret
- Coffea, basic tools and wrappers for enabling not-too-alien syntax when running columnar collider HEP analysis: @lgray, @nsmith-
- hist, histogramming for analysis powered by boost-histogram: @amangoel185, @henryiii
- uproot-browser, a TUI viewer for ROOT files: @amangoel185, @henryiii
- JetNet, for developing and reproducing ML + HEP projects: @rkansal47
- Awkward Array, manipulate JSON-like data with NumPy-like idioms: @agoose77, @ianna, @jpivarski
- Cling, a C++ interpreter based on LLVM: @vgvassilev
- cppyy, fully automatic, dynamic Python-C++ bindings by leveraging the Cling C++ interpreter and LLVM: @sudo-panda
- FCCAnalyses, common analysis framework for the Future Circular Collider (FCC): @kjvbrt
- func_adl, constructs hierarchical data queries using SQL-like concepts in Python: @masonproffitt, @gordonwatts
- Common Partial Wave Analysis, makes amplitude analysis accessible through transparent and interactive documentation, modern software development tools, and collaboration-independent frameworks: @redeboer
- Rio+, a C++ package for amplitude analysis: @JMolinaHN
- GooFit, a massively-parallel framework for maximum-likelihood fits, implemented in CUDA/OpenMP: @mdsokoloff
- zfit, model manipulation and fitting library based on TensorFlow and optimised for simple and direct manipulation of probability density functions, with a focus on scalability, parallelisation and user friendly experience: @jonas-eschle
- pyhf, pure-Python HistFactory implementation with tensors and autodiff: @matthewfeickert, @lukasheinrich
- cabinetry, for designing and steering profile likelihood fits: @alexander-held
- correctionlib, a generic correction library: @nsmith-
- hepstats, statistical tools and utilities: @jonas-eschle
- coffea-casa, prototype Analysis Facility: @oshadura
- Work Queue (website), a framework for building large distributed applications that span thousands of machines drawn from clusters, clouds, and grids: @btovar
- dask-awkward, native Dask collection for Awkward Arrays: @agoose77, @lgray
- hepdata_lib, a library for getting your data into HEPData: @clelange
- Analysis Grand Challenge, a project performing performing large scale end-to-end data analysis for HEP use cases: @oshadura, @alexander-held

# PyHEP workshops – aside advert ☺

❑ **Crucial series of workshops to exchange on "Python in HEP"**

❑ **2023 was the first year with 2 workshops !**
   **- It became relevant to complement the "Users workshop" held online since 2020**
    **with an in-person, informal workshop for developers of Python software in HEP**
    **to plan a coherent roadmap and make priorities for the upcoming year**

❑ **Upcoming events:**
   **PyHEP 2024 "Users workshop", online, July 1-4**
   **PyHEP.dev 2024 "Developers workshop", Aachen, August 26-30**
   **- Registration is open, as well as the call for abstracts !**

Number of PyHEP registrations by year, showing the huge jump when we went online.



Snapshots from **PyHEP.dev 2023** sessions

- **Scikit-HEP packages cited in a wide range of publications, within HEP and far beyond !**

- **Interestingly enough, HEP phenomenologists seem to be doing a much better job at citing the software they use !**

## Scientific publications using Scikit-HEP

▼ TABLE OF CONTENTS                                    (CLOSE)

### Peer reviewed journal papers

**Experimental Particle Physics**

- Scientific Data 9, 31 (2022) (arXiv:2108.02214 [hep-ex]) - cites `awkward` and `uproot`.
- LHCb Collaboration, JHEP 01 (2022) 166 ((arXiv:2107.10090 [hep-ex])), first LHCb analysis performed fully in Python with standard scientific packages and Scikit-HEP packages - cites `boost-histogram`, `iminuit`, `Particle`, `uproot`.
- Phys. Rev. D 104, 112011 (2021) (arXiv:2109.01685 [hep-ex]) - cites `iminuit`.
- JHEP 12 (2021) 133 (arXiv:2107.2105.13330 [hep-ex]) - cites `awkward`.
- Comput. Softw. Big Sci. 5, 22 (2021) (arxiv:2103.00659 [hep-ex]) - cites `uproot`.

**Particle Physics Phenomenology**

- Phys. Rev. D 109, 014040 (2024) (arXiv:2306.01578 [hep-ph]) - cites the Scikit-HEP project.
- Eur. Phys. J. C 84, 366 (2024) (arXiv:2312.04646 [hep-ph]) - cites `pylhe`.
- Phys. Rev. D 108, 035001 (2023) (arXiv:2302.08281 [hep-ph]) - cites `pylhe`.
- SciPost Phys. 14, 009 (2023) (arXiv:2206.14870 [hep-ph]) - cites `pyhf`.
- Eur. Phys. J. C 82, 326 (2022) (arXiv:2104.08921 [hep-ph]) - cites `iminuit`.
- JHEP 10 (2022) 092 (arxiv:2208.10406 [hep-ph]) - cites `pylhe`.
- Int. J. of Modern Physics A (arXiv:2208.09538 [hep-ph]) - cites the Scikit-HEP project.
- Phys. Rev. D 106, 052011 (2022) (arXiv:2206.09932 [hep-ph]) - cites the Scikit-HEP project.
- JHEP 05 (2022) 085 (arXiv:2110.02174 [hep-ph]) - cites `iminuit`.
- JHEP 04 (2022) 119 (arXiv:2201.01788 [hep-ph]) - cites `iminuit`.
- SciPost Phys. 12, 037 (2022) (arXiv:2109.04981 [hep-ph]) - cites `cabinetry`, `pyhf`.
- Eur. Phys. J. C 82, 194 (2022) (arXiv:2107.07537 [hep-ph]) - cites `iminuit`.
- Eur. Phys. J. C 82, 46 (2022) (arXiv:2109.08159 [hep-ph]) - cites `pyjet`.
- Rep. Prog. Phys. 85, 052201 (2022) (arXiv:2012.09874 [hep-ph]) - cites `iminuit`.
- Phys. Rev. D 105, 094033 (2022) (arXiv:2108.05383 [hep-ph]) - cites `iminuit`.
- Phys. Rev. D 105, 075001 (2021) (arXiv:2111.10343 [hep-ph]) - cites the Scikit-HEP project.
- Eur. Phys. J. C 81, 952 (2021) (arXiv:2103.13370 [hep-ph]) - cites `iminuit`.
- JHEP 11 (2021) 066 (arXiv:2107.06625 [hep-ph]) - cites `pyjet`.
- JHEP 05 (2021) 151 (arXiv:2103.03270 [hep-ph]) - cites `iminuit`.
- Rep. Prog. Phys. 84, 124201 (2021) (arXiv:2101.08320 [hep-ph]) - cites `pyjet`.
- Phys. Rev. D 104, 113008 (2021) (arXiv:2105.08270 [hep-ph]) - cites `uproot`, Scikit-HEP project.
- Phys. Rev. D 104, 035012 (2021) (arXiv:2105.07077 [hep-ph]) - cites `uproot`, Scikit-HEP project.
- Phys. Rev. D 104, 035014 (2021) (arXiv:2010.07941 [hep-ph]) - cites the Scikit-HEP project.
- Phys. Rev. D 103, 056014 (2021) (arXiv:2012.10500 [hep-ph]) - cites `uproot`, Scikit-HEP project.
- SciPost Phys. 10, 139 (2021) (arXiv:2008.06545 [hep-ph]) - cites `iminuit`, `probfit`.
- Phys. Rev. D 102, 015032 (2020) (arXiv:2006.10630 [hep-ph]) - cites `pylhe`, Scikit-HEP project.
- Phys. Rev. D 102, 015007 (2020) (arXiv:2005.03594 [hep-ph]) - cites `pylhe`, Scikit-HEP project.
- Comput. Softw. Big Sci. 4, 3 (2020) (arxiv:1907.10621 [hep-ph]) - cites `pylhe`, `scikit-hep`, `uproot`, Scikit-HEP project.
- JHEP 03 (2020) 094 (arXiv:2001.11041 [hep-ph]) - cites `iminuit`.
- JHEP 03 (2020) 076 (arXiv:1912.09760 [hep-ph]) - cites `iminuit`.
- JHEP 11 (2019) 034 (arXiv:1908.06980 [hep-ph]) - cites `pylhe`, `uproot`, Scikit-HEP project.
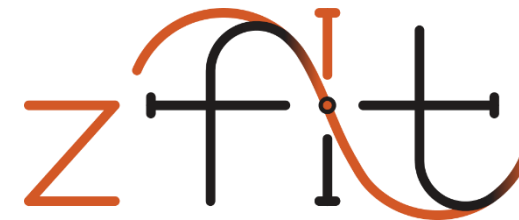- Phys. Rev. Lett. 123 (2019) 14, 141801 (arxiv:1811.06465 [hep-ph]) - cites `pylhe`.

Flavio

Coffea

# Engagement with
# broader scientific Python community

- **Contributions "upstream" to Python scientific ecosystem**

- **Engagement and impact**

# Engagement with the broader scientific Python community

*Python ecosystem in general*

❑ **We made contributions "upstream" to matplotlib and scipy**

❑ **Long ago we built our own wheel building tooling (for Azure DevOps),**
**then merged the features we needed into cibuildwheel, and have been a major contributor ever since**

❑ **We contributed to pybind11 and scikit-build and support scikit-build-core, the modern CMake-based build backend**

❑ **HEP specific – we also did a fair share of the work to make ROOT available on Conda**

**NUMFOCUS**
**OPEN CODE = BETTER SCIENCE**

*NumFOCUS*

❑ **2 packages (pyhf, then awkward) are <u>NumFOCUS Affiliated Projects</u>**

❑ **Affiliation brings visibility and various other benefits**

❑ **Working with the broader Scientific Python ecosystem**
**to make sure that we adopt better sustainability practices**

## Project news

**February 2024:** Packages `hepconvert` and `ragged` joined Scikit-HEP.

**June 2023:** Developer pages, `cookie`, and `repo-review` generalized and accepted into scientific-python.

**May 2023:** `pyhf` project updates published in the May 2023 NumFOCUS newsletter.

**March 2023:** `awkward` accepted as an affiliated project of NumFOCUS.

**January 2023:** Package `numpythia` deprecated and archived. Package `pyjet` archived.

**December 2022:** `pyhf` accepted as an affiliated project of NumFOCUS.

**December 2022:** Major releases `awkward` 2.0.0 and `uproot` 5.0.0.

**August 2022:** Major metapackage release `scikit-hep` 4.0.0.

**July 2022:** Scikit-HEP packages being taught in the "Python for Analysis" training events of the HEP Software Foundation (HSF) Training Working Group.

**May 2022:** Package `repo-review` created in Scikit-HEP.

# Engagement with the broader scientific Python community

**Scientific Python** 

*scientific-python.org*

❑ **Such engagement is paramount and a win-win !**

❑ **Our developer guide, cookiecutter, and repo-review,**
   **were adopted into scientific-python,**
   **and are now used by other projects, including AstroPy**
   **(see next page …)**

Graph shows how people who contribute to Scikit-HEP repos
are connected to people who contribute to core Scientific Python repos
- We're still a bit farther than many others …



**Taken from <u>Blog - Developer Summit 1 (scientific-python.org)</u>**

# Engagement with the broader scientific Python community

❑ **Check out**

**https://learn.scientific-python.org/**

# Engagement with the broader scientific Python community

☐ **We gave/give but also get in return**

☐ **Concrete and useful tool we have access to – statistics from Plausible Analystics !**



☐ **Unique visitors of scikit-hep.org site and pages most viewed**

   - Interesting to see that Top Pages are largely about info for developers ☺

# Lessons learned and future

- **Future directions and developments**

- **Matters of sustainability**

# Lessons learned and future

*Early days and avoiding the "quantum fluctuation drama"*

❑ **A vision is not enough without a dedicated team *and* a community**

❑ **It is not necessary to find many people in early days to bring a project over threshold.**
   **2-3 very motivated and skilful colleagues is enough**

❑ **Scikit-HEP seems to have come to the arena at the right moment and constant engagement with the community was key**


*Future and sustainability*

❑ **It is the ecosystem surrounding a piece of software, built not only by original authors but also by community members, that forms a long-term hence sustainable solution in a continuously evolving space**


❑ **Small package granularity – a set of "building blocks" and well-scoped generic or domain-specific libraries providing tools that address a well-defined class of problems at one level of abstraction has been a success story**
   **- Increasing number of projects/libraries build on top of these**


❑ **Last but not least, without adequate support, career path and funding, nothing of this is sustainable !**

# Going forward – family of small and scoped packages

❑ **Packages that surround Awkward and Uproot: new developments are mostly going to be in new packages now, rather than turning the base packages into monoliths**

❑ **Tiny packages are great – the easiest to try something new. Did not work / was a success? No worries, try another approach (with no significant loss for the ecosystem)!**

❑ **Areas under active development / investigation:**
- **Autodifferentiation**
- **Fitting with JAX as backend (zfit2, evermore)**
- **Leverage Dask but not only**
- **Even better interoperability**
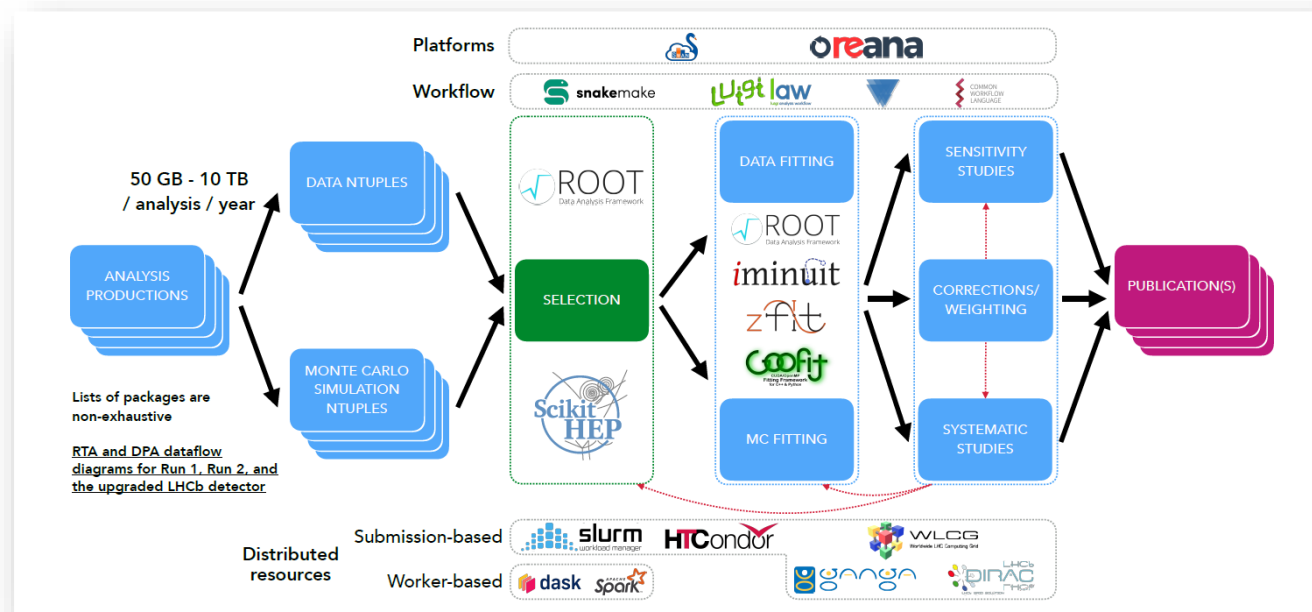- **Etc.**



Jim Pivarski, ACAT 2024

# Side-comment wrapper: Heterogeneous hardware *and* software, as it were

*diverse in character or content

❑ **The recent decade has been very rich as far as Computing & Software is concerned**

❑ **Not only the hardware we use is ever more heterogeneous\*,**
  **the same is true for the ecosystem of software libraries !**

  - **Biased overview example from LHCb for offline analysis:**



LHCb-FIGURE-2024-002 (CERN CDS link)

# Thank you for listening

**And thank you for using/contributing to/… Scikit-HEP ☺ !**