

PYTHON x

HS³

INFERENCE SUITE

Lorenz Gärtner, Thomas Kuhr, Giordon Stark

PUBLISHING STATISTICAL MODELS

WORKSHOP ON CONFIDENCE LIMITS

CERN, Geneva, Switzerland
17–18 January 2000

CERN 2000–005

Massimo Corradi

Does everybody agree on this statement, to publish likelihoods?

Louis Lyons

Any disagreement ? Carried unanimously. That's actually quite an achievement for this Workshop.

...[Fred James wants to be able to calculate coverage, Don Groom wants to be able to calculate goodness of fit]...

Cousins

I thought the point of unanimity was that publishing the likelihood function was a *necessary* condition, not a sufficient condition.

But a practical problem remained: *How to communicate multi-D likelihood?*

HEP Statistics Serialization Standard - HS³

HS³ provides a **standardized format** for statistical models:

- ❑ Human-readable
 - ❑ Machine-readable
 - ❑ Software-independent
 - ❑ Generic, mathematical definitions
 - ❑ Only from specifications: build your own implementation in language of your choice
- } JSON(-like) format



<https://tinyurl.com/hs3-overleaf>



<https://tinyurl.com/hs3-github>

<https://indico.cern.ch/event/1294577/contributions/5693326/>

HS³ - Overview of supported types and components

Carsten Burgard, Tomas Dado, Jonas Eschle, Matthew Feickert, Cornelius Grunwald, Alexander Held, Jerry Ling, Robin Pelkner, Jonas Rembser, Oliver Schulz

November 6, 2023

About this document



To the extent possible under law, the authors have waived all copyright and related or neighboring rights to this document. For details, please refer to the *Creative Commons* CC0 license [1]. This work is published from: CERN, Geneva.

1 Introduction

This section is non-normative.

With the introduction of `pyhf` [2], a JSON format for likelihood serialization has been put forward. However, an interoperable format that encompasses likelihoods with a scope beyond stacks of binned histograms was sorely lacking. With the release of `ROOT 6.26/00` [3] and the experimental `RoJSONFactoryWSTool` therein, this gap has now been filled.

This document sets out to document the syntax and features of the HEP Statistics Serialization Standard (HS³) for likelihoods and statistical models in general, as to be adopted by any HS³-compatible statistics framework. The examples in this document employ the JSON notation, but are intended to encompass also representations as YAML or TOML.

Please note that this document as well as the HS³ standard are still in development and can still undergo minor and major changes in the future. This document describes the syntax of version 0.2 of the HS³ standard.

CURRENT IMPLEMENTATIONS



ROOT
Data Analysis Framework



BAT.jl
Bayesian Analysis Toolkit



zfit



pyhf
differentiable
Likelihoods

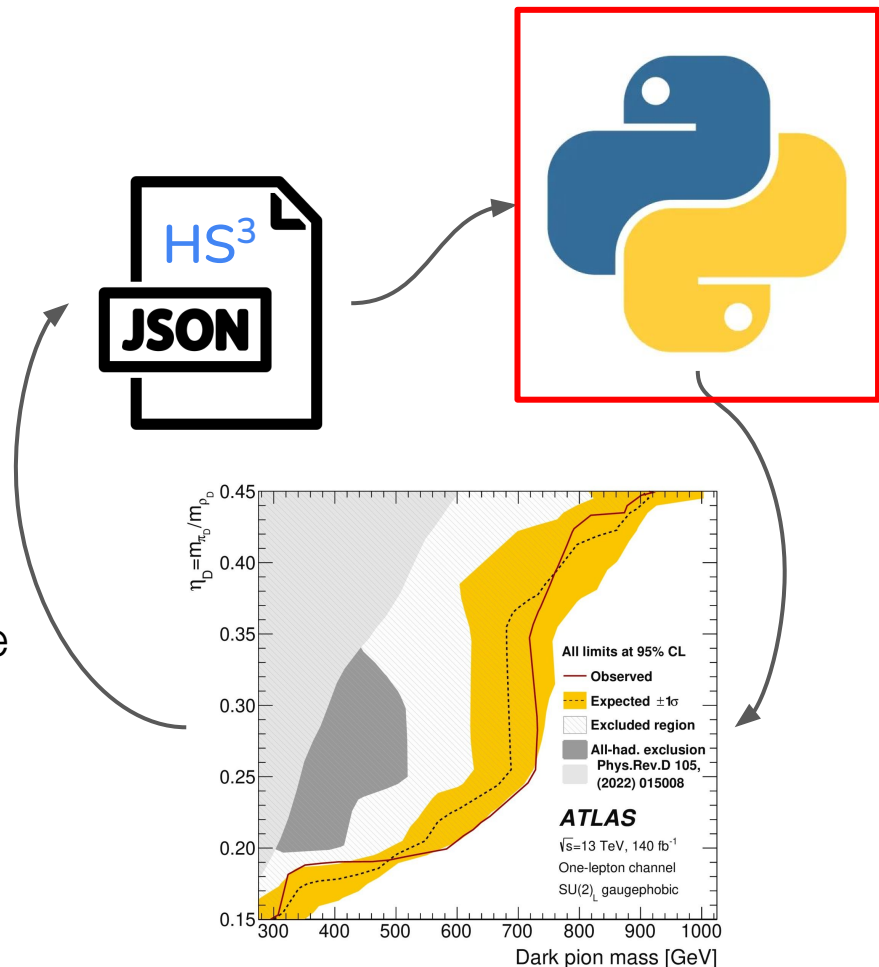


THE MISSING LINK

A pythonic **well-maintained**, **user-friendly** and **efficient implementation** of HS³.

Build a **common core** for probability model serialization and inference.

- ❑ Python-based package that enables statistical inference using HS³-compatible inputs.
- ❑ Interoperable with the HS³ schema for reading and writing.



LEARNING FROM PAST SUCCESS ...

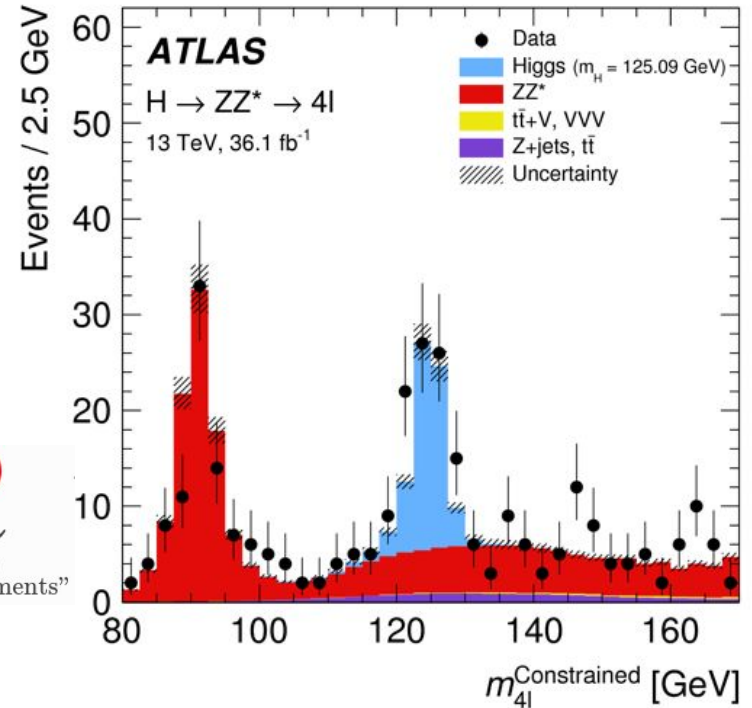
HISTFACTORY

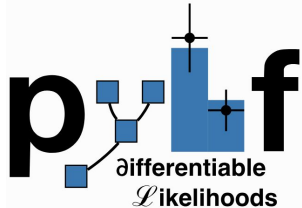
Kyle Cranmer, George Lewis, Lorenzo Moneta, Akira Shibata, Wouter Verkerke

June 20, 2012

- ❑ A **mathematical p.d.f template** specification for the building of statistical models from **binned** distributions and data
- ❑ Widely used by the HEP community for standard model measurements and BSM searches
- ❑ Independent of its implementation in ROOT

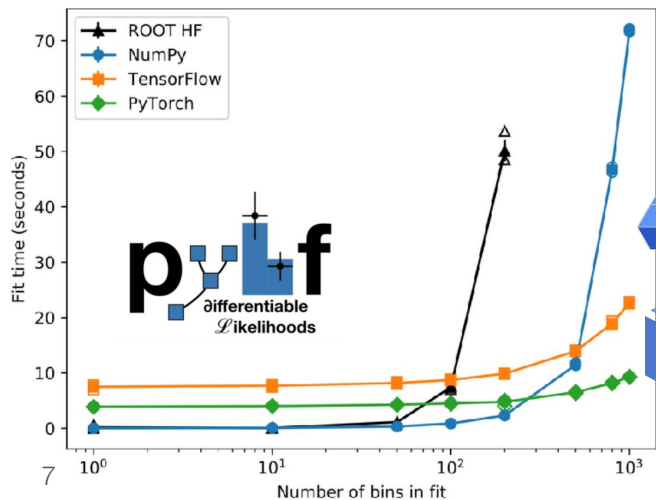
$$f(\mathbf{n}, \mathbf{a} | \boldsymbol{\eta}, \boldsymbol{\chi}) = \underbrace{\prod_{c \in \text{channels}} \prod_{b \in \text{bins}_c} \text{Pois}(n_{cb} | \nu_{cb}(\boldsymbol{\eta}, \boldsymbol{\chi}))}_{\text{Simultaneous measurement of multiple channels}} \underbrace{\prod_{\chi \in \boldsymbol{\chi}} c_{\chi}(a_{\chi} | \boldsymbol{\chi})}_{\text{constraint terms for "auxiliary measurements"}}$$





A python-only implementation of the HistFactory model + inference methods, along with a JSON model serialization.

Orders-of-magnitude faster inference than ROOT



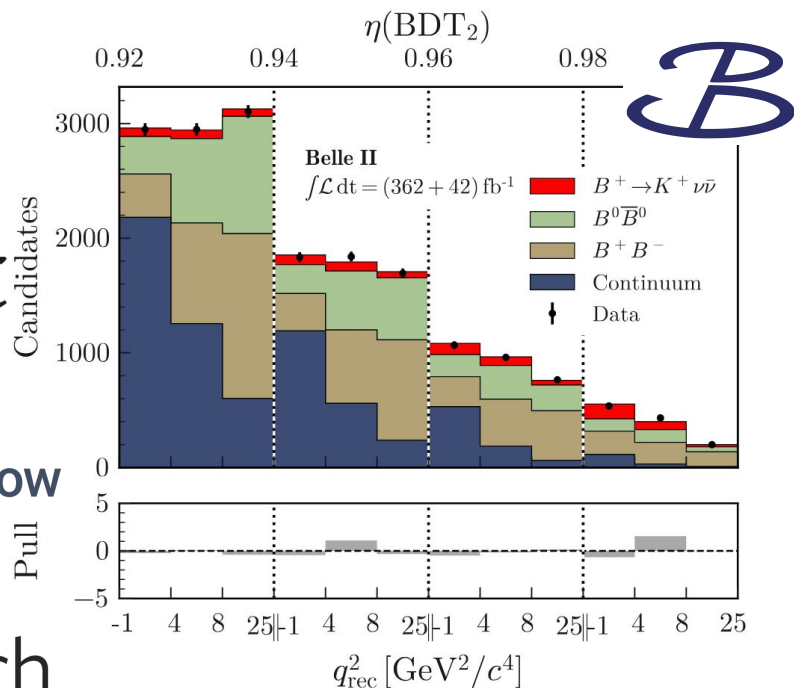
Lukas Heinrich



Matthew Feickert



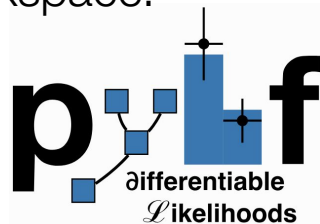
Giordon Stark



PUBLISHING LIKELIHOODS

A single plain-text, human-readable file (JSON) specifies the entire HistFactory workspace!

This enables:



- ❑ Statistical combinations
- ❑ Reinterpretation (w/ signal patching)
- ❑ Propagating updated theory predictions
- ❑ ... **Publishing statistical models: Getting the most out of particle physics experiments**

Kyle Cranmer^{1*}, Sabine Kraml^{2‡}, Harrison B. Prosper^{3§} (editors),
Philip Bechtle⁴, Florian U. Bernlochner⁴, Itay M. Bloch⁵, Enzo Canonero⁶, Marcin
Chrzaszcz⁷, Andrea Coccaro⁸, Jan Conrad⁹, Glen Cowan¹⁰, Matthew Feickert¹¹,
Nahuel Ferreiro Iachellini^{12,13}, Andrew Fowlie¹⁴, Lukas Heinrich¹⁵, Alexander Held¹⁶,
Thomas Kuhr^{13,16}, Anders Kvellestad¹⁷, Maeve Madigan¹⁸, Farvah Mahmoudi^{15,19},
Knut Dundas Morà²⁰, Mark S. Neubauer¹¹, Maurizio Pierini¹⁵, Juan Rojo⁸, Sezen
Sekmen²², Luca Silvestrini²³, Veronica Sanz^{24,25}, **Giordon Stark**²⁶, Riccardo Torre⁸,
Robert Thorne²⁷, Wolfgang Waltenberger²⁸, Nicholas Wardle²⁹, Jonas Wittbrodt³⁰

```
1 {
2   "channels": [
3     {
4       "name": "singlechannel",
5       "samples": [
6         {
7           "name": "signal",
8           "data": [12.0, 11.0],
9           "modifiers": [
10            {
11              "name": "mu",
12              "type": "normfactor",
13              "data": null
14            }
15          ]
16        },
17        {
18          "name": "background",
19          "data": [50.0, 52.0],
20          "modifiers": [
21            {
22              "name": "uncorr_bkguncrt",
23              "type": "shapesys",
24              "data": [3.0, 7.0]
25            }
26          ]
27        }
28      ]
29    }
30 ]
31 }
```


MAIN TAKEAWAYS

- ❑ Simplicity is key for serialization
 - ❑ Simple construction
 - ❑ Simple I/O
- ❑ Backward compatibility essential
 - ❑ pyhf version as metadata
- ❑ Conversion (from other formats)
 - ❑ eg. pyhf <> HistFactory XML (RooWorkspace)
- ❑ 'Closed world' model simplifies serializability and scope of responsibilities

```
>>> with open('workspace.json') as ws:
...     spec = json.load(ws)
...
>>> workspace = pyhf.Workspace(spec)
>>> pdf = workspace.model()
>>> data = workspace.data(model)
>>> pyhf.infer.mle.fit(data, pdf)
array([0.95260667, 0.99635345, 1.02140172])
```

HEPData

analysis:HistFactory Search Reset search Advanced .JSON

Max results Sort by Reverse order Showing 10 of 39 results

Date

Collaboration

Subject_areas

Phrases

Proton-Proton Scattering 10

Cross Section 7

SUSY 6

Supersymmetry 6

Top 5

Next Show 20

Search for flavour-changing neutral-current couplings between the top quark and the Higgs boson in multi-lepton final states in 13 TeV pp collisions with the ATLAS detector

The ATLAS collaboration Aad, Georges; Aakvaag, Erlend; Abbott, Braden Keim; et al.
Eur.Phys.J.C 84 (2024) 757, 2024.
Inspire Record 2773613 DOI 10.117182/hepdata.150988

A search is presented for flavour-changing neutral-current interactions involving the top quark, the Higgs boson and an up-type quark ($q = u, c$) with the ATLAS detector at the Large Hadron Collider. The analysis considers leptonic decays of the top quark along with Higgs boson decays into two W bosons, two Z bosons or a $\tau^+\tau^-$ pair. It focuses on final states...

0 data tables match query

Measurement of ℓ -channel production of single top quarks and antiquarks in pp collisions at 13 TeV using the full ATLAS Run 2 data sample

The ATLAS collaboration Aad, Georges; Abbott, Braden Keim; Abeling, Kira; et al.
JHEP 05 (2024) 305, 2024.
Inspire Record 2764820 DOI 10.117182/hepdata.150693

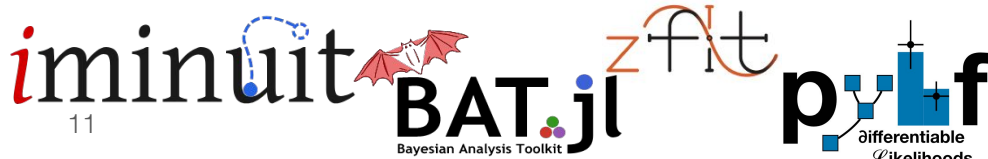
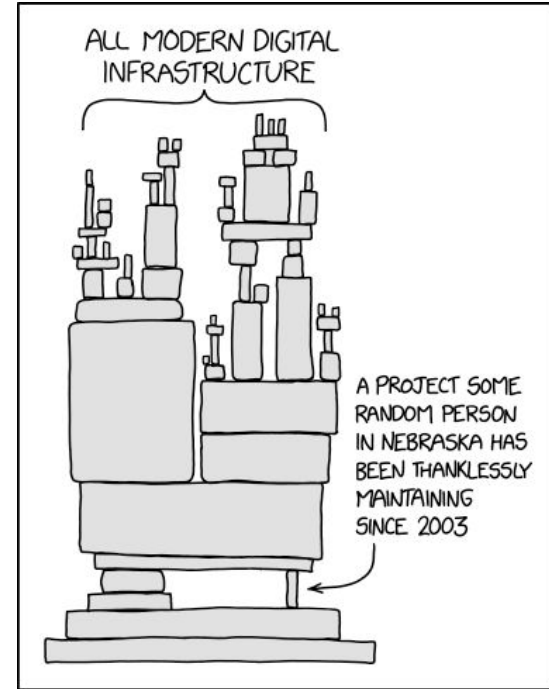
The production of single top quarks and top antiquarks via the ℓ -channel exchange of a virtual W boson is measured in proton-proton collisions at a centre-of-mass energy of 13 TeV at the LHC using 140 fb $^{-1}$ of ATLAS data. The total cross-sections are determined to be $\sigma(\ell q) = 137^{+10}_{-10}$ pb and $\sigma(\ell \bar{q}) = 84^{+6}_{-6}$ pb for top-quark and top-antiquark production,...

39 HistFactory models on HEPData

A NEW COMMON CORE ...

... REQUIRES A SOLID FOUNDATION

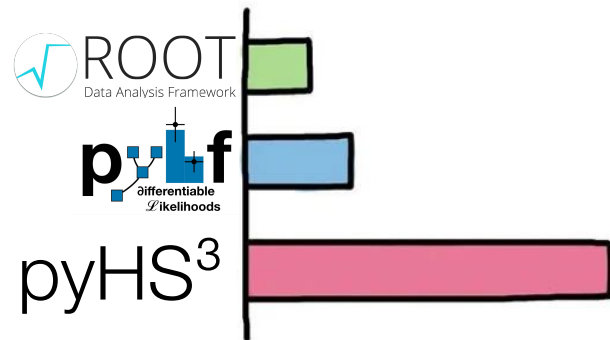
- ❑ Built with HS³ in mind from the start
- ❑ Well maintainable
- ❑ Backwards-compatible
 - ❑ Tools change and so will HS³
 - ❑ Results should remain reproducible
- ❑ Limited scope
 - ❑ Set clear boundaries of the tool
- ❑ Maximized efficiency
 - ❑ Fast evaluation of likelihood and gradients
 - ❑ eg. automatic differentiation
- ❑ Well interfaced with already existing tools



TENTATIVE PLAN

- ❑ Start with representative analyses from different HEP experiments providing
 - ❑ HS³-compatible serialization of their probability models
 - ❑ Public data and statistical inference results for cross-checks
- ❑ Pick the right (modern) tools for the job
 - ❑ TensorFlow-probability, scipy.stats, pytorch, JAX
 - ❑ Minimize dependencies as much as possible
- ❑ Documentation, Tutorials, and Testing
 - ❑ Documentation describes a “**Model Interface**” / “**Inference Interface**”
 - ❑ Translate existing statistics tutorials and lectures to show how to use python-HS³
- ❑ Maintenance
 - ❑ A robust CI/CD to maintain cross-checks with the representative analyses
 - ❑ Enforce code styles and lint with tools like ruff and pylint
 - ❑ Open-source to allow for community contributions

WHAT GIVES PEOPLE FEELINGS OF POWER



OPEN QUESTIONS

SERIALIZE THE OPEN WORLD

- ❑ How open world do we need to be?
- ❑ How to handle custom mathematical functions in a user-friendly way?
- ❑ How to allow users to provide their own custom functions as plugins?

COMPUTATIONAL EFFICIENCY

- ❑ How to build a computational graph that can be exported and used as part of end-to-end differentiation?

```
>>> import numexpr as ne
>>> import numpy as np
>>> a = np.random.rand(3)
>>> b = np.random.rand(3)
>>> ne.evaluate("sin(a) + arcsinh(a/b)")
array([0.24508376, 2.512371 , 4.03901998])
```



ERuM Data project proposal

Lorenz Gärtner, Thomas Kuhr, Giordon Stark

Pythonic HS³ inference

Project summary

In experimental high energy physics, the ability to communicate and publish results in a computationally-accessible format is important for allowing theorists and phenomenologists to (re)use analysis data. In recent years, efforts have made progress on various fronts of the principles of FAIR (Findability, Accessibility, Interoperability, and Reusability) data; including SLHA files for parameter configurations used in Monte Carlo generators to a popular plain-text specification for HistFactory (HiFa) probability models and its corresponding implementation: pyhf. This project proposal is inspired by the latter development and aims to cover a gap that will be missing as

COMMENTS WELCOME

BACKUP

MORE OPEN QUESTIONS

- ❑ Only python? Julia? Rust?
- ❑ What is the best way to serialize probability models following the HS3 schema? JSON? YAML? BSON? Some other plain-text encoding that's more highly-compressible, or a binary encoding?
- ❑ How to deal with plotting?
- ❑ How to deal with debugging?
- ❑ What gets pulled out of pyhf, and how to design in such a way that it is mostly compatible with pyhf's API?

HISTFACTORY - A mathematical prescription

$$f(\mathbf{n}, \mathbf{a} | \boldsymbol{\eta}, \boldsymbol{\chi}) = \underbrace{\prod_{c \in \text{channels}} \prod_{b \in \text{bins}_c} \text{Pois}(n_{cb} | \nu_{cb}(\boldsymbol{\eta}, \boldsymbol{\chi}))}_{\text{Simultaneous measurement of multiple channels}} \underbrace{\prod_{\chi \in \boldsymbol{\chi}} c_{\chi}(a_{\chi} | \boldsymbol{\chi})}_{\text{constraint terms for "auxiliary measurements"}}$$

Multiple, disjoint channels of binned distributions with multiple samples contributing to each with additional (shared) systematics between sample estimates

- ❑ **Poisson p.d.f.** for bins observed in all channels
- ❑ **Constraint p.d.f.** (and data) for auxiliary measurements (systematics: normalization, shape, etc)

$$\nu_{cb}(\boldsymbol{\phi}) = \sum_{s \in \text{samples}} \nu_{scb}(\boldsymbol{\eta}, \boldsymbol{\chi}) = \sum_{s \in \text{samples}} \underbrace{\left(\prod_{\kappa \in \boldsymbol{\kappa}} \kappa_{scb}(\boldsymbol{\eta}, \boldsymbol{\chi}) \right)}_{\text{multiplicative modifiers}} \underbrace{\left(\nu_{scb}^0(\boldsymbol{\eta}, \boldsymbol{\chi}) + \sum_{\Delta \in \boldsymbol{\Delta}} \Delta_{scb}(\boldsymbol{\eta}, \boldsymbol{\chi}) \right)}_{\text{additive modifiers}}$$