



# HEPDATA & IRIS-HEP



### @KyleCranmer

New York University Department of Physics Center for Data Science CILVR Lab

## NSF SUPPORTED HEP SOFTWARE EFFORTS

# First DIANA-HEP (2014-2019)

# Now IRIS-HEP Software Institute (2018-2023)

**dianahep** Team Activities/Products DIANA Fellows Blog





#### **Collaborative Analyses**

Establish infrastructure for a higher-level of collaborative analysis, building on the successful patterns used for the Higgs boson discovery and enabling a deeper communication between the theoretical community and the experimental community



#### **Reproducible Analyses**

Streamline efforts associated to reproducibility, analysis preservation, and data preservation by making these native concepts in the tools





#### Computational and data science research to enable discoveries in fundamental physics

IRIS-HEP is a software institute funded by the National Science Foundation. It aims to develop the state-of-the-art software cyberinfrastructure required for the challenges of data intensive scientific research at the High Luminosity Large Hadron Collider (HL-LHC) at CERN, and other planned HEP experiments of the 2020's. These facilities are discovery machines which aim to understand the fundamental building blocks of nature and their interactions. Full Overview

#### News and Featured Stories:



First USATLAS Bootcamp held in coordination with Software Carpentries and **IRIS-HEP/FIRST-HEP** 

CoDaS-HEP 2019 at Princeton University For the third consecutive summer, high energy physics graduate students, postdocs and

Upcoming Events:	
Jan 15–17, 2020	New York University
ML4Jets2020	
Feb 17–19, 2020	CERN
Analysis Preservation	Bootcamp
Apr 22–24, 2020	Princeton University
Connecting the Dots 2	2020
	View all past events

Upcoming Topical Meetings:
Jan 27, 2020 Allen project
Jan 29, 2020 Primary vertex finding at LHCb
Feb 3, 2020
Modeling computing resource needs/costs
Feb 12, 2020
mkFit and CMS tracking

View all • Indico (recordings) • Vidyo room



### **Analysis Systems**

The goal of the Analysis Systems focus area is to develop sustainable analysis tools to extend the physics reach of the HL-LHC experiments by creating greater functionality, reducing time-to-insight, lowering the barriers for smaller teams, and streamlining analysis preservation, reproducibility, and reuse.



Focus Area Strategies:

- Establish declarative specifications for analysis tasks and workflows that will enable the technical development of analysis systems to be decoupled from the user- facing semantics of physics analysis.
- Leverage and align with developments from industry and the broader scientific software community to enhance sustainability of the analysis systems.
- Develop high-throughput, low-latency systems for analysis for HEP.
- Integrate analysis capture and reuse as first class concepts and capabilities into the analysis systems.

## STATUS 2013

#### You Learn Something New Every Day »

#### **KYLE CRANMER | USLHC | USA**

« Prioritizing the future

#### VIEW BLOG | READ BIO

2013

#### Inspired by the Higgs, a step forward in open access

# The discovery of the Higgs boson is a major step forward in our understanding of nature at the most fundamental levels. In addition to being the last piece of the standard model, it is also at the core of the fine tuning problem — one of the deepest mysteries in particle physics. So it is only natural that our scientific methodology rise to the occasion to provide the most powerful and complete analysis of this breakthrough discovery.

This week the ATLAS collaboration has taken an important step forward by making the likelihood function for three key measurements about the Higgs available to the world digitally. Furthermore, this data is being shared in a way that represents a template for how particle physics operates in the fast-evolving world of open access to data. These steps are a culmination of decades of work, so allow me to elaborate.



Four interactions that can produced a Higgs boson at the LHC

First of all, what are the three key measurements, and why are they important? The three results were presented by ATLAS in this recent paper. Essentially, they are measurements for how often the Higgs is produced at the LHC through different types of interactions (shown above) and how often it decays into three different force carrying particles (photons, W, and Z bosons). In this plot, the black + sign at (1,1) represents the standard model prediction and the three sets of contours represent the measurements performed by ATLAS. These measurements are fundamental tests of the standard model and any deviation could be a sign of new physics like supersymmetry!



### **AAHEP7 Information Provider Summit**

1-3 April 2014 Stony Brook University US/Eastern timezone







Durham

Science & Tecl

~ (Ip3~

The Durham HepData Project

# ATLAS POLICY DOCUMENT

### Level-1. Published results

All scientific output is published in journals, and preliminary results are made available in Conference Notes. All are openly available, without restriction on use by external parties beyond copyright law and the standard conditions agreed by CERN.

Data associated with journal publications are also made available: tables and data from plots (e.g. cross section values, likelihood profiles, selection efficiencies, cross section limits, ...) are stored in appropriate repositories such as HEPDATA[2]. ATLAS also strives to make additional material related to the paper available that allows a reinterpretation of the data in the context of new theoretical models. For example, an extended encapsulation of the analysis is often provided for measurements in the framework of RIVET [3]. For searches information on signal acceptances is also made available to allow reinterpretation of these searches in the context of models developed by theorists after the publication. ATLAS is also exploring how to provide the capability for reinterpretation of searches in the future via a service such as RECAST [4]. RECAST allows theorists to evaluate the sensitivity of a published analysis to a new model they have developed by submitting their model to ATLAS.



Display a menu y from the ATLAS collaboration, with the first open release of full analysis likelihoods





pyhf is a pure-python implementation of the widely-used HistFactory p.d.f. template described in [CERN-OPEN-2012-016]. It also includes interval estimation is based on the asymptotic formulas of "Asymptotic formulae for likelihood-based tests of new physics" [arxiv:1007.1727]. The aim is also to support modern computational graph libraries such as PyTorch and TensorFlow in order to make use of features such as autodifferentiation and GPU acceleration.

## A talk on pyhf by Matthew Feickert GitHub DOI 10.5281/zenodo.1169739 CI/CD passing docker build automated codecov 95% reg code quality: python A+ codefactor A code style black docs master launch binder pypi package 0.4.0 python 3.6 | 3.7 docker stars 2 docker pulls 15k

### Team

- Kyle Cranmer
- Lukas Heinrich
- Matthew Feickert
- Giordon Stark

# DIANA & IRIS-HEP CONTRIBUTING



PERSPECTIVE

https://doi.org/10.1038/s41567-018-0342-2

**Corrected: Publisher Correction** 

### OPEN

### Open is not enough

Xiaoli Chen<sup>1,2</sup>, Sünje Dallmeier-Tiessen<sup>1\*</sup>, Robin Dasler<sup>1,11</sup>, Sebastian Feger<sup>1,3</sup>, Pamfilos Fokianos<sup>1</sup>, Jose Benito Gonzalez<sup>1</sup>, Harri Hirvonsalo<sup>1,4,12</sup>, Dinos Kousidis<sup>1</sup>, Artemis Lavasa<sup>1</sup>, Salvatore Mele<sup>1</sup>, Diego Rodriguez Rodriguez<sup>1</sup>, Tibor Šimko<sup>1\*</sup>, Tim Smith<sup>1</sup>, Ana Trisovic<sup>1,5\*</sup>, Anna Trzcinska<sup>1</sup>, Ioannis Tsanaktsidis<sup>1</sup>, Markus Zimmermann<sup>1</sup>, Kyle Cranmer<sup>6</sup>, Lukas Heinrich<sup>6</sup>, Gordon Watts<sup>7</sup>, Michael Hildreth<sup>8</sup>, Lara Lloret Iglesias<sup>9</sup>, Kati Lassila-Perini<sup>4</sup> and Sebastian Neubert<sup>10</sup>

The solutions adopted by the high-energy physics community to foster reproducible research are examples of best practices that could be embraced more widely. This first experience suggests that reproducibility requires going beyond openness.

# reana

### Reproducible research data analysis platform





#### Schematic view Generate simulated data Werge Fit model to data Physics results b b b c south and b south a



**Fig. 2 | Example of a complex computational workflow on REANA mimicking a beyond the standard model (BSM) analysis**. This figure shows an example where the experimental data is compared to the predictions of the standard model with an additional hypothesized signal component. The example permits one to study the complex computational workflows used in typical particle physics analyses. **a-c**, The computational workflow (**a**) may consist of several tens of thousands of computational steps that are massively parallelizable and run in a cascading 'map-reduce' style of computations on distributed compute clusters. The workflow definition is modelled using the Yadage workflow specification and produces an upper limit on the signal strength of the BSM process. A typical search for BSM physics consists of simulating a hypothetical signal process (**c**), as well as the background processes predicted by the standard model with properties consistent with the hypothetical signal (marked dark green in (**b**)). The background often consists of simulated background estimates (dark blue and light green histograms) and data-driven background estimates (light blue histogram). A statistical model involving both signal (dark green histogram) and background components is built and fit to the observed experimental data (black markers). **b**, Results of the model in its pre-fit configuration at nominal signal strength. We can see the excess of the signal over data, meaning that the nominal setting does not describe the data well. The post-fit distribution would scale down the signal in order to fit the data. This REANA example is publicly available at ref. <sup>35</sup>. For icon credits, see Fig. 1.





**RECAST** is a framework for extending the impact of existing analyses performed by high-energy physics experiments.

- *Request*: Upload alternative signals in the LHE format and request that any given analysis is "recast" for an alternative model. Note: this is a request, there is no obligation for the experiments to respond.
- Subscribe: Anyone can subscribe to an analysis to be informed of activity associated with the analysis
- *Recast*: Experimentalists can accept the request, process these alternative signals with the full simulation, reconstruction, and analysis selection. If authorized by their collaboration, they can respond with an authoritative result for the selection efficiency and cross-section limits for the alternative signal. Note, anyone can provide a non-authoritative result, for instance one based on a phenomenological recasting tool.







Sinclert Pérez joined NYU

- Research Software Engineer
- Previously worked at CERN team on REANA
- 50% IRIS-HEP



Sinclert Pérez New York University

Research Software Engineer

Focus is to improve interoperability of cyberinfrastructure components (eg. RECAST, CAP, HEPData, INSPIRE, ...)

# VISUALIZATION

### https://www.youtube.com/watch?time\_continue=5&v=9egt9ZTm7T0&feature=emb\_title



#### Interactive Exploration of a HistFactory Model

One advantage of a pure-python implementation of Histfactory is the ability to explore the pdf interactively within the setting of a notebook. Try moving the sliders and oberserve the effect on the samples. For example changing the parameter of interest *SigXsecOverSM* (or  $\mu$ ) controls the overall normalization of the (BSM) signal sample ( $\mu$ =0 for background-only and  $\mu$ =1 for the nominal signal-plus-background hypothesis)

interact(plot, order = fixed([1,2,3,0]), ax = fixed(None), \*\*{n[0]: tuple(m) for n,m in zip(sorted(reversed(list(parnam)))









## API: AUTOMATING UPLOADS

### EXAMPLE RECAST → ZENODO

If experiments do adopt something like this, would be nice to have API connection to upload result.

