Sorry if this is a little disorganized, COVID has complicated work life

# ACKNOWLEDGEMENTS

EP-IT Data science seminars

## HEP in the Cloud Computing and Open Science Era

by Lukas Alexander Heinrich (CERN)
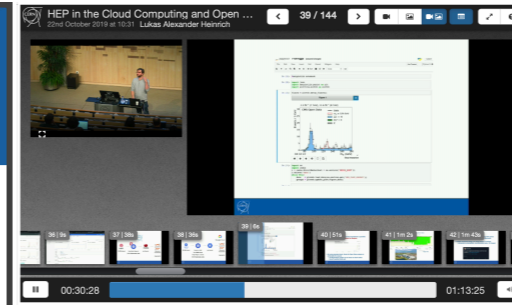
📅 Wednesday 23 Oct 2019, 11:00 → 12:00 Europe/Zurich

📍 500/1-001 - Main Auditorium (CERN)  https://indico.cern.ch/event/840837/

Description  As the LHC readies for Run-3 and its second decade of data-taking, the world around us is changing rapidly. Since the discovery of the Higgs boson in 2012 cloud computing has fundamentally changed the style and access of distributed computing, Deep Learning and Data Science have entered the public vocabulary and Open Science and Reproducibility has been grown in importance. The LHC experiment, with their vast amounts of data, unique dataset is necessarily find themselves at the forefront of these developments. In this talk, I will discuss about how these trends enable new research avenues and data anlaysis capabilities, such as a systematic reinterpretation program for Beyond the Standard Model search using cloud-native workflows and RECAST, "rediscovering" the Higgs boson in CERN Open Data within a few minutes on the cloud, to enabling third-party research through open access to high-fidelity data products of LHC searches the wider HEP community.

📎 DataScience23.10...  📎 data_science_semi...  🔗 Recording

Organized by  M. Girone, M. Elsing, L. Moneta, M. Pierini.......... Coffee will be served at 10h30

Webcast  📹 There is a live webcast for this event  [Watch]

## Scalable cyberinfrastructure applications

Team: J. Brehmer[1,2], K. Cranmer[1,2], **Irina Espejo**[1], S. Macaluso[1,2] and H. Müller[1]
Institutions: [1] Center for Data Science, New York University
[2] Department of Physics, New York University  The SCAILFIN Project
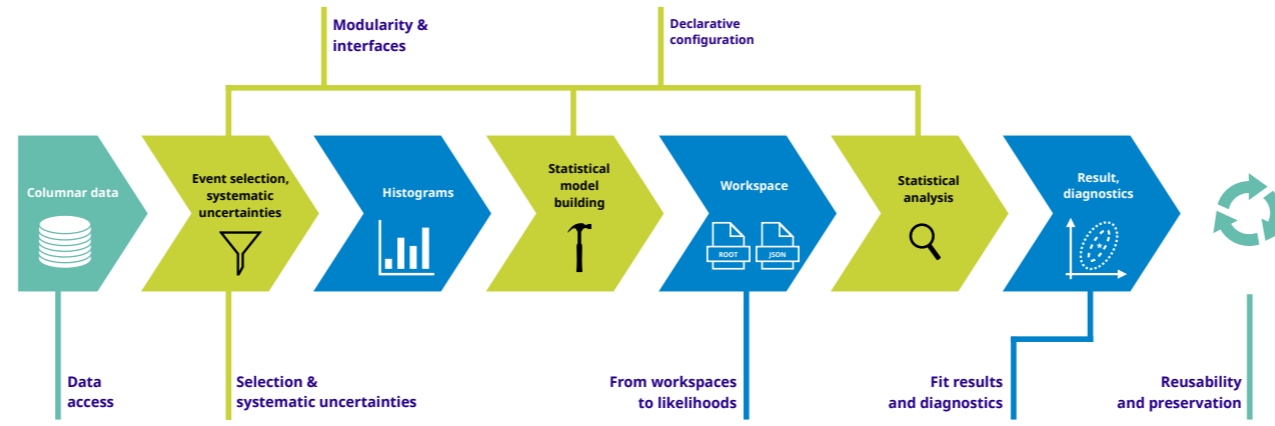
**Matthew Feickert**
@HEPfeickert

Thanks to everyone at #SciPy2019 who came and asked me great questions about pyhf!

pyhf is a pure Python statistical fitting library that uses tensors and autograd to speed up physics analysis at the LHC

4

# TOPICS

Accelerating analysis design

- more powerful observables

- end-to-end optimization

- benchmarking of algorithms

Accelerating fitting

- pyhf and a fitting service

More efficient simulation

- excursion
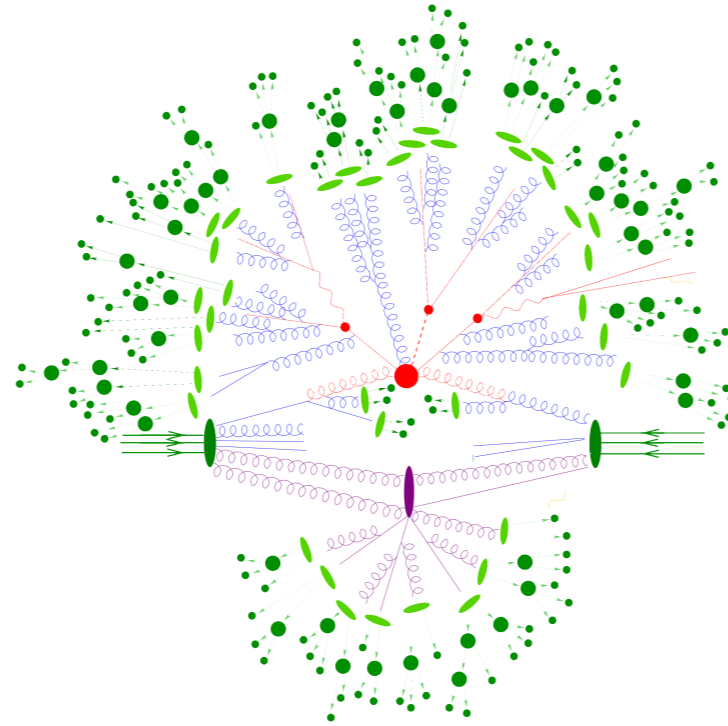
- Probabilistic programming

Extending impact of results

- RECAST

Core technologies:

- automatic differentiation

- GPUs & TPUs

- Cloud-native : docker, kubernetes

- Workflows & REANA

- Functions as a service Accelerating analysis design
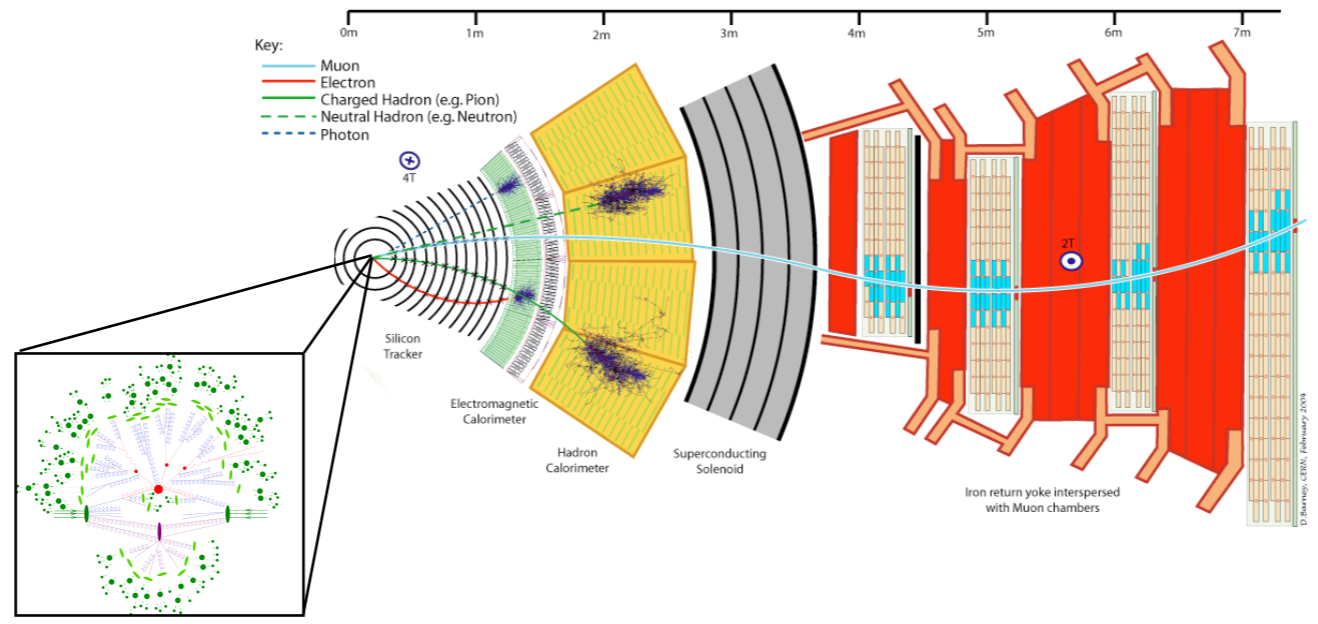
$$\mathcal{L}_{SM} = \underbrace{\frac{1}{4}\mathbf{W}_{\mu\nu}\cdot\mathbf{W}^{\mu\nu} - \frac{1}{4}B_{\mu\nu}B^{\mu\nu} - \frac{1}{4}G^a_{\mu\nu}G^{\mu\nu}_a}_{\text{kinetic energies and self-interactions of the gauge bosons}}$$

$$+ \underbrace{\bar{L}\gamma^\mu(i\partial_\mu - \frac{1}{2}g\tau\cdot\mathbf{W}_\mu - \frac{1}{2}g'YB_\mu)L + \bar{R}\gamma^\mu(i\partial_\mu - \frac{1}{2}g'YB_\mu)R}_{\text{kinetic energies and electroweak interactions of fermions}}$$

$$+ \underbrace{\frac{1}{2}\left|(i\partial_\mu - \frac{1}{2}g\tau\cdot\mathbf{W}_\mu - \frac{1}{2}g'YB_\mu)\phi\right|^2 - V(\phi)}_{W^\pm, Z, \gamma, \text{and Higgs masses and couplings}}$$

$$+ \underbrace{g''(\bar{q}\gamma^\mu T_a q)G^a_\mu}_{\text{interactions between quarks and gluons}} \quad + \quad \underbrace{(G_1\bar{L}\phi R + G_2\bar{L}\phi_c R + h.c.)}_{\text{fermion masses and couplings to Higgs}}$$

## DETECTOR SIMULATION

**Conceptually:** Prob(detector response | particles )

**Implementation:** Monte Carlo integration over micro-physics

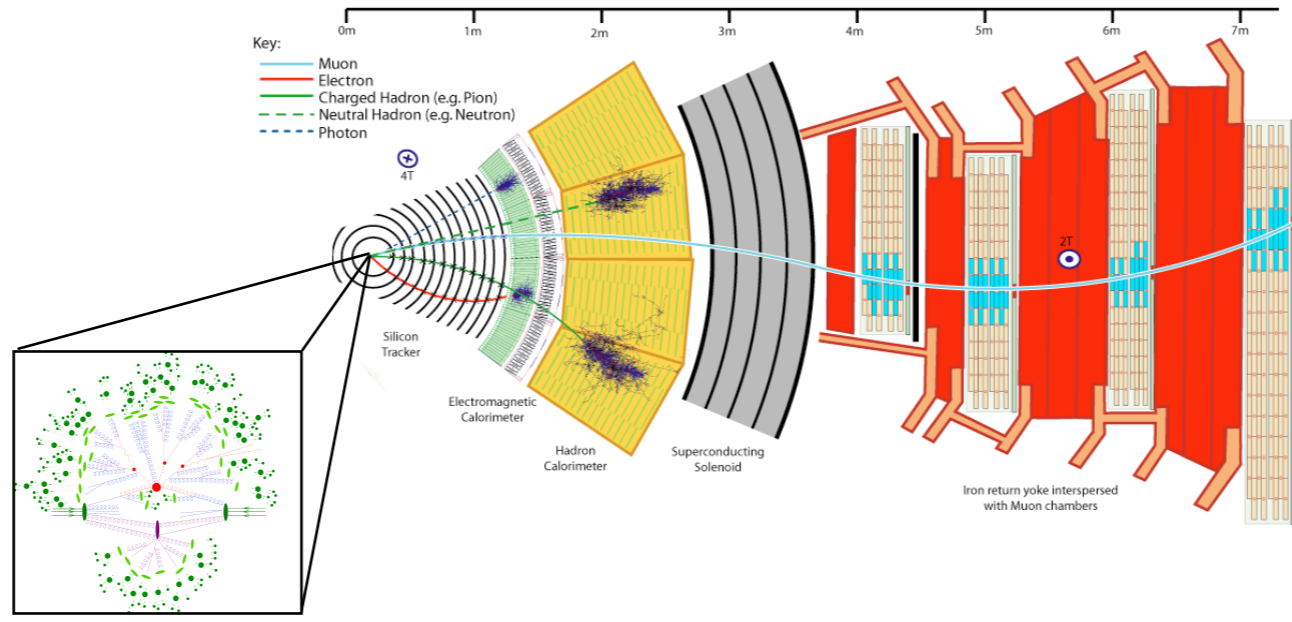**Consequence:** evaluation of the likelihood is intractable

# THE CAUSAL, GENERATIVE MODEL

| Observables | Detector interactions | Shower splittings | Parton-level momenta | Theory parameters |
|---|---|---|---|---|

$$x \longleftarrow z_d \longleftarrow z_s \longleftarrow z_p \longleftarrow \theta$$

$$p(x|\theta) = \int \mathrm{d}z_d \int \mathrm{d}z_s \int \mathrm{d}z_p \; p(x|z_d) \qquad p(z_d|z_s) \qquad p(z_s|z_p) \qquad p(z_p|\theta)$$

# $10^8$ SENSORS → 1 REAL-VALUED QUANTITY

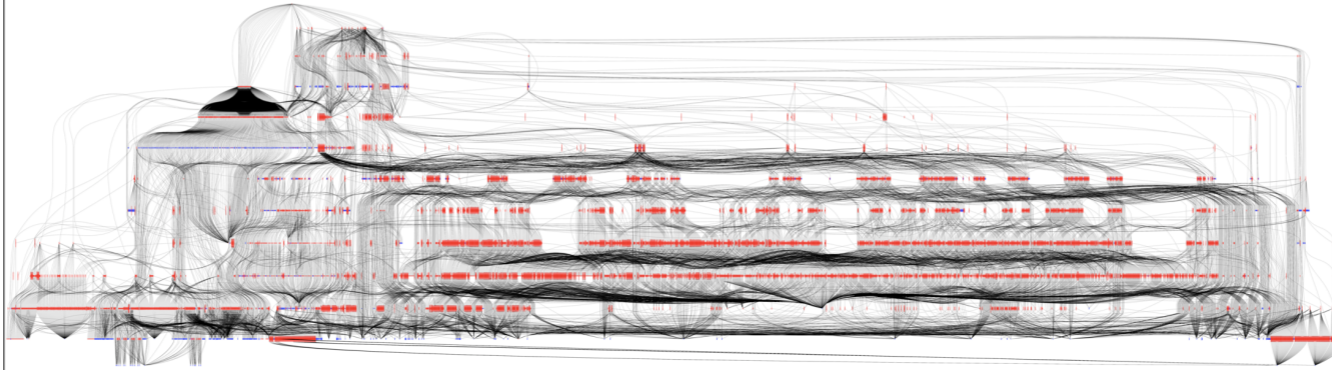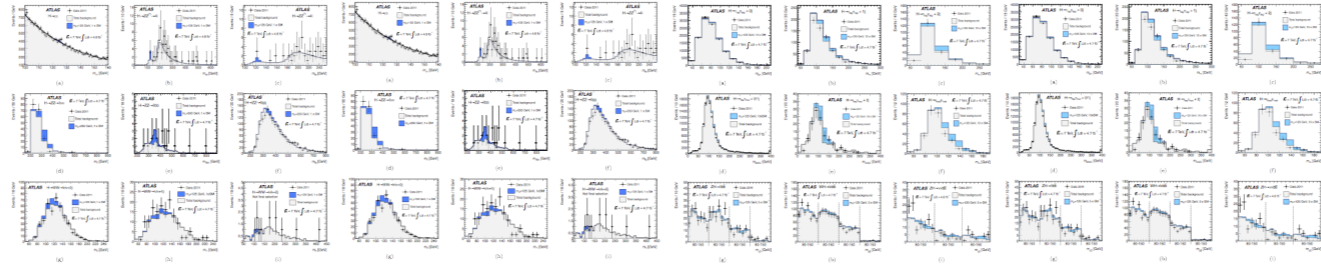Most measurements and searches for new particles at the LHC are based on the distribution of a single summary statistic

- choosing a good summary statistic (feature engineering) is a task for a skilled physicist and tailored to the goal of measurement or new particle search

- likelihood $p(x|\theta)$ **approximated** using histograms (univariate density estimation)
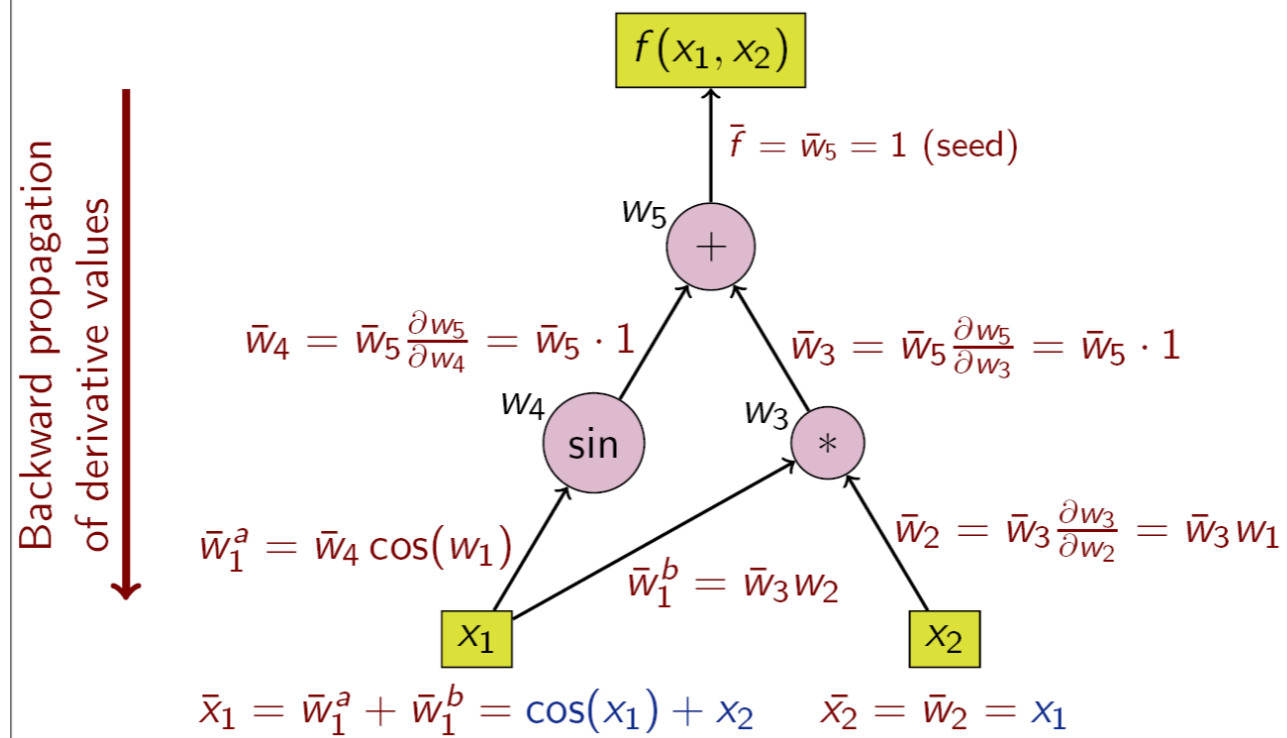


**This doesn't scale if x is high dimensional!**

$$\mathbf{f}_{\text{tot}}(\mathcal{D}_{\text{sim}}, \mathcal{G}|\boldsymbol{\alpha}) = \prod_{c \in \text{channels}} \left[ \text{Pois}(n_c|\nu_c(\boldsymbol{\alpha})) \prod_{e=1}^{n_c} f_c(x_{ce}|\boldsymbol{\alpha}) \right] \cdot \prod_{p \in \mathbb{S}} f_p(a_p|\alpha_p)$$

2. **reverse accumulation** computes the recursive relation: $\dfrac{dy}{dw_i} = \dfrac{dy}{dw_{i+1}} \dfrac{dw_{i+1}}{dw_i}$ with $w_0 = x$.

$f(x_1, x_2)$

$\bar{f} = \bar{w}_5 = 1$ (seed)

$w_5$ $+$

$\bar{w}_4 = \bar{w}_5 \dfrac{\partial w_5}{\partial w_4} = \bar{w}_5 \cdot 1$ $\qquad$ $\bar{w}_3 = \bar{w}_5 \dfrac{\partial w_5}{\partial w_3} = \bar{w}_5 \cdot 1$

$w_4$ $\sin$ $\qquad$ $w_3$ $*$

$\bar{w}_1^a = \bar{w}_4 \cos(w_1)$ $\qquad$ $\bar{w}_2 = \bar{w}_3 \dfrac{\partial w_3}{\partial w_2} = \bar{w}_3 w_1$

$\bar{w}_1^b = \bar{w}_3 w_2$

$x_1$ $\qquad$ $x_2$

$\bar{x}_1 = \bar{w}_1^a + \bar{w}_1^b = \cos(x_1) + x_2$ $\qquad$ $\bar{x}_2 = \bar{w}_2 = x_1$

Backward propagation of derivative values

12

# pyhf: auto-differentiable binned HEP likelihoods

Kyle Cranmer (NYU) Matthew Feickert (SMU)
Lukas Heinrich (CERN), Giordon Stark (UCSC)

implementation of HistFactory likelihood (1) as a **computational graph of multi-dimensional array operations**

Use of array ("tensor") operations through a common API layer around high-performance tensor libraries: e.g.

NumPy · TensorFlow · PyTorch

**Installation:**

```
$> pip install pyhf
```

Example: simple number-counting experiment

$$\mathcal{P}(n_c, x_e, a_p \,|\, \phi_p, \alpha_p, \gamma_b) = \prod_{c\in\text{channels}} \left[ \text{Pois}(n_c|\nu_c) \prod_{e=1}^{n_c} f_c(x_e|\boldsymbol{\alpha}) \right] \cdot G(L_0|\lambda, \Delta_L) \cdot \prod_{p\in\mathbb{S}+\Gamma} f_p(a_p|\alpha_p)$$

## Auto Differentiation:

Tensor libraries from ML communty provide **exact gradients** for use in minimization. $\frac{\partial\mathcal{L}}{\partial\mu}, \frac{\partial\mathcal{L}}{\partial\theta_i}$

## Optimizer

pyhf likeliho... ...n functions. Can use multiple minimization algorithms, such as `scipy.minimize` or MINUIT

$$\frac{\partial\mathcal{L}}{\partial\mu}, \frac{\partial\mathcal{L}}{\partial\theta_i}$$

Standard Model · SUSY · Exotics

```
{"op": "replace", "path": "/channels/0/samples/0/data", "value": [5., 6.]}
```

```
$> curl http://url-to-json/workspace.json|pyhf cls
```

13

# OTHER ADVANTAGES OF USING TENSOR BACKENDS

- All numerical operations implemented in **tensor backends** through an API of $n$-dimensional array operations

- Using deep learning frameworks as computational backends allows for **exploitation of auto differentiation (autograd) and GPU acceleration**

- As huge buy in from industry we benefit for free as these frameworks are **continually improved** by professional software engineers



- Preliminary results

- Show hardware acceleration giving **order of magnitude speedup** for some models!

- Hardware acceleration benchmarking planned

- Improvements over traditional
  - 10 hrs to 30 min; 20 min to 10 sec

# MAKING IT STANDARD

10 years later: community embraces publishing likelihoods as a standard

- Moved to JSON schema



CERN Council appoints Fabiola Gianotti for second term of office as CERN Director General

Press release | 6 November, 2019

## LATEST NEWS

Particle physicists formulate the f...

Knowledge sharing | News

LHCb explores the beauty of lepton universali...

Physics | News | 15 January, 2020

New open release allows theorists to explore ...

Knowledge sharing | News
9 January, 2020

View all news ›

## New open release allows theorists to explore LHC data in a new way

The ATLAS collaboration releases full analysis likelihoods, a first for an LHC experiment

9 JANUARY, 2020 | By Katarina Anthony



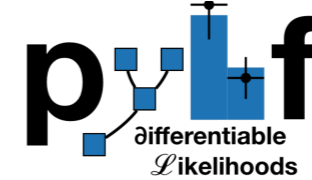Explore ATLAS open likelihoods on the HEPData platform (Image: CERN)

What if you could test a new theory against LHC data? Better yet, what if the expert knowledge needed to do this was captured in a convenient format? This tall order is now on offer from the ATLAS collaboration, with the first open release of full analysis likelihoods

### Related Articles

Open Data: pushing back the frontiers togethe...

Knowledge sharing | News | 22 October, 2018

The invisible structure providing Open Acces...

Knowledge sharing | News | 30 July, 2018

Sit down for coffee with the Standard Model

At CERN | News
7 April, 2017

View all news ›

# FITTING SERVICE

With JSON format, it is much easier to stream necessary data to a fitting service.

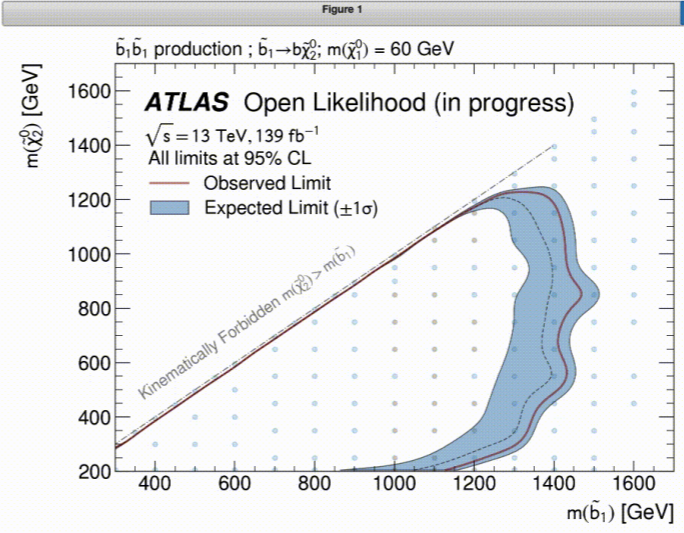- Lukas prototyped this using functions-as-a-service

Ideally have a machine with a big GPU or TPU for this

- From 10 min to a few seconds!

Accessing Fitting Service

```
def func(data):
    filename = data['filename']
    region = data['region']
    m = re.compile("sbottom_(\d+)_(\d+)_(\d+)").search(filename).group(0)
    outname = 'results/region().result.{}.json'.format(region,m)
    for i in range(10):
        try:
            d = requests.post(
                '{}/region{}'.format(FITTING_SVC,region),
                data = open(filename), headers = {'Content-Type': 'application/json'}
            ).json()
            json.dump(d,open(outname,'w'))
            break
        except:
            pass#retry
```

```
[4]: %matplotlib notebook
fig,ax = plt.subplots(1,1)
fig.set_size_inches(9.33,7)
apply_decorations(ax,label = 'Open Likelihood (in progress)')
```

Figure 1

$\tilde{b}_1\tilde{b}_1$ production ; $\tilde{b}_1 \to b\tilde{\chi}_2^0$; $m(\tilde{\chi}_1^0)$ = 60 GeV

**ATLAS** Open Likelihood (in progress)

$\sqrt{s} = 13$ TeV, 139 fb$^{-1}$
All limits at 95% CL

Observed Limit
Expected Limit ($\pm 1\sigma$)

Kinematically Forbidden $m(\tilde{\chi}_2^0) > m(\tilde{b}_1)$

$m(\tilde{\chi}_2^0)$ [GeV]

$m(\tilde{b}_1)$ [GeV]

No handles with labels found to put in legend.

```
[*]: for x in glob.glob('results/*.json'):
    os.unlink(x)
cA = [{'region': 'A', 'filename': f} for f in glob.glob('RegionA/patch*_60.json')]
cC = [{'region': 'C', 'filename': f} for f in glob.glob('RegionC/patch*_60.json')]
configs = cA[:] + cC[:]
# np.random.shuffle(configs)

import time
import concurrent.futures
fig.canvas.draw()

with concurrent.futures.ThreadPoolExecutor(max_workers=MAX_WORKERS) as executor:
    for i,_ in enumerate(tqdm(executor.map(func, configs),total = len(configs))):
        if i > 5 and i % 5 == 0:
            make_plot(ax,label = 'Open Likelihood (in progress)', color = 'steelblue', showPoints = True)
            fig.canvas.draw()
        time.sleep(.005)
    make_plot(ax, label = 'Open Likelihood', color = 'gold', showPoints = False)
```
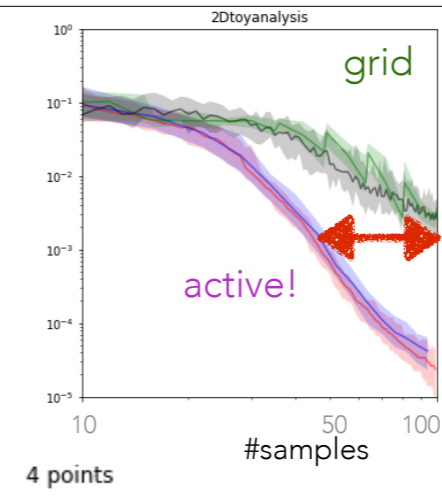
35%                91/259 [01:55<01:54, 1.46it/s]

```
/home/jovyan/interpolate.py:365: UserWarning: No contour levels were found within the data range.
  c = ax.contour(xi,yi,zi, [level])
```
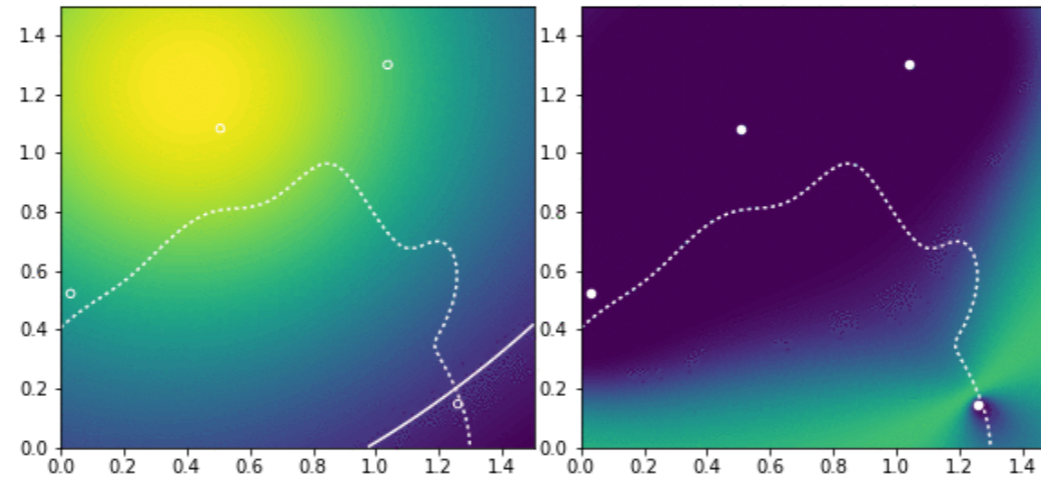
# ACTIVE LEARNING

Instead of generating Monte Carlo a priori, generate it on demand where it is relevant!

↓ An algorithm for finding exclusion contours

Drastically more efficient use of computing resources →

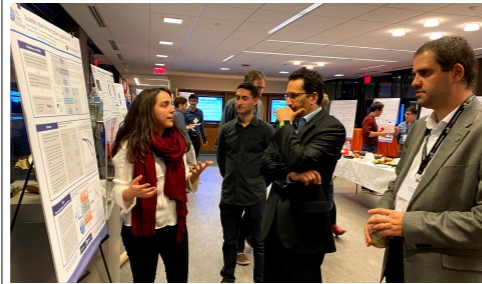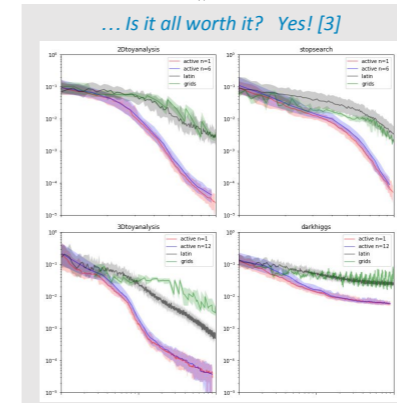K.C., Lukas Heinrich, Gilles Louppe, ACAT2019

# Scalable cyberinfrastructure applications

Team: J. Brehmer[1,2], K. Cranmer[1,2], **Irina Espejo**[1], S. Macaluso[1,2] and H. Müller[1]

Institutions: [1] Center for Data Science, New York University
[2] Department of Physics, New York University

**iris hep** — Institute for Research & Innovation in Software for High Energy Physics

**The SCAILFIN Project**

## Excursion

in collaboration with G. Louppe[2] and L. Heinrich[3]
[2] University of Liège

 diana-hep/excursion
 irinaespejo/excursion

Physics simulations

- Goal is to find *level sets of black-box functions* that are expensive to evaluate. Examples: test statistics from complex simulations.

- Evaluate the black blox function at *interesting points only* instead of evaluating at whole regular grid. We use a *Gaussian process* to: interpolate between samples and model uncertainty in the knowledge of the black box function.
  The *acquisition function* regulates the exploration vs exploitation tradeoff. Select one that *minimizes global uncertainty* of the location of the excursion set.

- Future: efforts will focus on *scaling* the dimensionality of the function domain. Example, likelihood ratio as function of mass, charge, spin,...

configure

generate simulations
physics + detector

combine

sample

train + evaluate

docker — dockerhub/madminertool/
docker-madminer-physics

**reana**

Machine Learning Inference

docker — dockerhub/madminertool/
docker-madminer-ml

... Is it all worth it?  Yes! [3]

[3] L. Heinrich, G. Louppe, K. Cranmer, *Excursion Set Estimation using Sequential Entropy Reduction for Efficient Searches for New Physics at the LHC*, ACAT 2019

## README.md

https://indico.cern.ch/event/708041/contributions/3269754/

## `excursion` — Efficient Excursion Set Estimation

DOI 10.5281/zenodo.1634427  launch binder  build passing

This package implements a Bayesian Optimization procedure based on Gaussian Processes to efficiently determine excursion sets (or equivalently iso-surfaces) of one or many expensive black-box functions.
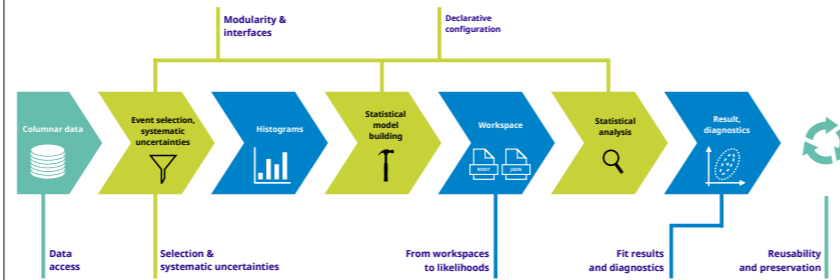
### Installation and Example

Install via `pip install excursion==0.0.1a0` .

**Front-End**

**flow≶erv**

**Back-End**

dianahep  NYU Center for Data Science

MSDSE  DeepMind  FULBRIGHT

NSF

# DIFFERENTIABLE PROGRAMING

Automatic differentiation is not just for Machine Learning!

- Differentiable Programming

- **Attitude**: we can auto-diff through analysis and reconstruction
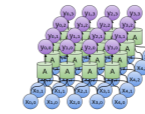
- End-to-end optimization

Backpropagate: dLimit /dSelection including full statistical treatment with systematics.



https://github.com/pyhf/neos

# Effective field theory

When looking for deviations from the standard model Higgs, we would like to look at all sorts of kinematic correlations

- thus each observation **x** is high-dimensional

# LEARNING THE LIKELIHOOD RATIO

parameter $\theta$

latent $z$

observable $x$

$r(x, z|\theta)$

$t(x, z|\theta)$

augmented data

$\arg\min_{g} L[g]$

approximate likelihood ratio

$\hat{r}(x|\theta)$

$\theta_j$

$\theta_i$

**Simulation**　　　　**Machine Learning**　　　　**Inference**



**Amortized likelihood**

**Amortized posterior**

**Amortized likelihood ratio**

**Amortized surrogates trained with augmented data**

e)　　f)　　g)　　h)

23

# LIKELIHOOD-FREE INFERENCE

parameter $\theta$

latent $z$

observable $x$

$r(x, z|\theta)$

$t(x, z|\theta)$

$\arg\min_{g} L[g]$

approximate likelihood ratio

$\hat{r}(x|\theta)$

$\theta_j$

$\theta_i$

Machine Learning

Inference

We can use **augmented data** to dramatically improve training

with augmented data

without augmented data

Technically similar to PDF weights

2D histogram
CARL
ROLR
SALLY
CASCAL
RASCAL

Estimation error

Training sample size

(based on a 42-Dim observation **x**)



J Brehmer, J Pavez, G Louppe, K.C. PRL & PRD 2018 [arXiv:1805.00013 & arXiv:1805.00020]
"Better Higgs Measurements Through Information Geometry" [arXiv:1612.05261] & CARL [arxiv:1506.02169]

[JB, S. Dawson, S. Homiller, F. Kling, T. Plehn 1908.06980]

- Simplified Template Cross-Sections (STXS) define observable bins that are supposed to capture as much information on NP as possible
  [N. Berger et al. 1906.02754; HXSWG YR4]



- Let's check! How much information on

$$\tilde{\mathcal{O}}_{HD} = \mathcal{O}_{H\square} - \frac{\mathcal{O}_{HD}}{4} = (\phi^\dagger\phi)\square(\phi^\dagger\phi) - \frac{1}{4}(\phi^\dagger D^\mu\phi)^*(\phi^\dagger D_\mu\phi)$$
$$\mathcal{O}_{HW} = \phi^\dagger\phi W_{\mu\nu}^a W^{\mu\nu a}$$
$$\mathcal{O}_{Hq}^{(3)} = (\phi^\dagger i \overset{\leftrightarrow}{D}{}_\mu^a \phi)(\overline{Q}_L\sigma^a\gamma^\mu Q_L) ,$$

can we extract from $pp \to WH \to \ell\nu \, b\bar{b}$ ?

- Results: STXS are indeed sensitive to operators, adding a few more bins improve them, but a multivariate analysis is still stronger

We accomplished a lot!

From scratch:

- Generate simulated data for EFT with MadGraph

- Fast detector simulation

- Trained neural network to learn likelihood ratio

- Trained neural network to learn Score (Optimal Observable)

- Calculated expected limit for both approaches and compared to simple 1-d histogram approach

- Calculated Fisher information matrix

This is workflow for several published papers

- To speed this up, working to streamline MadMiner with REANA

https://cranmer.github.io/madminer-tutorial/



**MadMiner Tutorial**

Introduction

**MadMiner Tutorial**

Preliminaries

Overview

Define process to study *

   Morphing

   Interactive Morphing Demo

Create training data

   Set MadGraph Directory

   Parton Level *

   With Delphes

Train model

   Likelihood Ratio *

   Score *

   Likelihood

Statistical Analysis

   Limits on EFT parameters *

   Fisher Information

   Information Geometry

Congratulations

## Introduction

**MadMiner tutorial**

This is a tutorial on MadMiner developed by Johann Brehmer, Felix Kling, Irina Espejo, and Kyle Cranmer. It is built using Jupyter Book.

**Introduction to MadMiner**

Particle physics processes are usually modelled with complex Monte-Carlo simulations of the hard process, parton shower, and detector interactions. These simulators typically do not admit a tractable likelihood function: given a (pot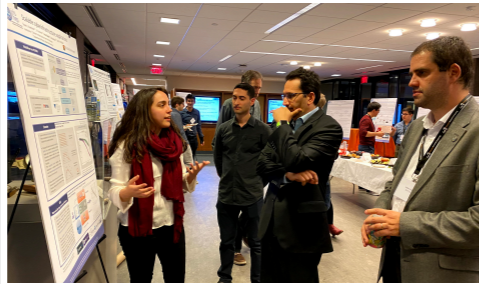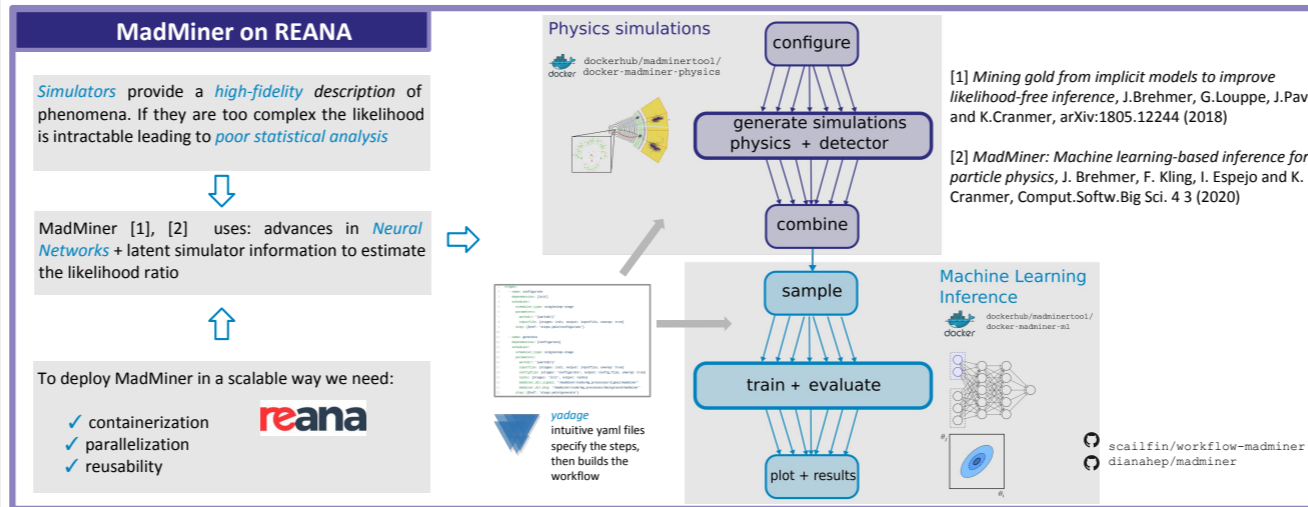entially high-dimensional) set of observables, it is usually not possible to calculate the probability of these observables for some model parameters. Particle physicists usually tackle this problem of "likelihood-free inference" by hand-picking a few "good" observables or summary statistics and filling histograms of them. But this conventional approach discards the information in all other observables and often does not scale well to high-dimensional problems.

In the three publications "Constraining Effective Field Theories With Machine Learning", "A Guide to Constraining Effective Field Theories With Machine Learning", and "Mining gold from implicit models to improve likelihood-free inference", a new approach has been developed. In a nut shell, additional information is extracted from the simulations that is closely related to the matrix elements that determine the hard process. This "augmented data" can be used to train neural networks to efficiently approximate arbitrary likelihood ratios. We playfully call this process "mining gold" from the simulator, since this information may be hard to get, but turns out to be very valuable for inference.

# Scalable cyberinfrastructure applications

Team: J. Brehmer[1 2], K. Cranmer[1 2], **Irina Espejo**[1], S. Macaluso[1 2] and H. Müller[1]

Institutions: [1] Center for Data Science, New York University
[2] Department of Physics, New York University

The SCAILFIN Project

## MadMiner on REANA

*Simulators* provide a *high-fidelity* *description* of phenomena. If they are too complex the likelihood is intractable leading to *poor statistical analysis*

MadMiner [1], [2] uses: advances in *Neural Networks* + latent simulator information to estimate the likelihood ratio

To deploy MadMiner in a scalable way

- ✓ containerization
- ✓ parallelization
- ✓ reusability

Physics simulations

configure

generate simulations
physics + detector

combine

sample

Machine Learning Inference

train + evaluate

[1] *Mining gold from implicit models to improve likelihood-free inference*, J.Brehmer, G.Louppe, J.Pavez and K.Cranmer, arXiv:1805.12244 (2018)

[2] *MadMiner: Machine learning-based inference for particle physics*, J. Brehmer, F. Kling, I. Espejo and K. Cranmer, Comput.Softw.Big Sci. 4 3 (2020)

scailfin/workflow-madminer
dianahep/madminer

Front-End

End

Back-End

- **HEP analysis on Cloud w/ standard ingredients is possible**
  - **Terabit per second throughput for analysis.**



- **All using s** ‌ ‌ ‌ ‌ **is achievable for anyone**
  - **O(10k)/analysis nodes not an issue on short notice**
  - **haven't pushed limits yet.**

# SHIFTING FROM REPRODUCIBILITY TO REUSE

## Open is not enough

Xiaoli Chen[1,2], Sünje Dallmeier-Tiessen[1]*, Robin Dasler[1,11], Sebastian Feger[1,3], Pamfilos Fokianos[1], Jose Benito Gonzalez[1], Harri Hirvonsalo[1,4,12], Dinos Kousidis[1], Artemis Lavasa[1], Salvatore Mele[1], Diego Rodriguez Rodriguez[1], Tibor Šimko[1]*, Tim Smith[1], Ana Trisovic[1,5]*, Anna Trzcinska[1], Ioannis Tsanaktsidis[1], Markus Zimmermann[1], Kyle Cranmer[6], Lukas Heinrich[6], Gordon Watts[7], Michael Hildreth[8], Lara Lloret Iglesias[9], Kati Lassila-Perini[4] and Sebastian Neubert[10]
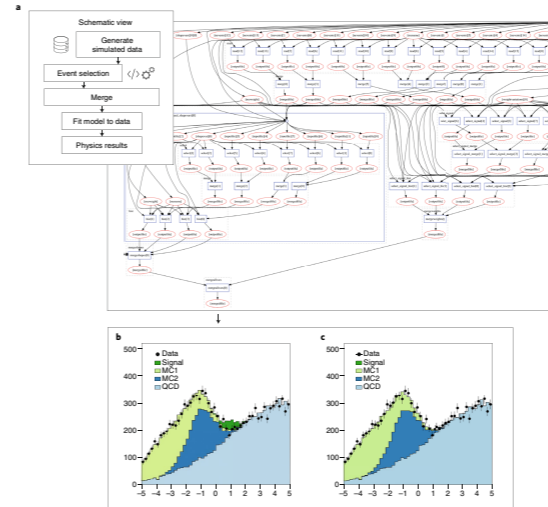
The solutions adopted by the high-energy physics community to foster reproducible research are examples of best practices that could be embraced more widely. This first experience suggests that reproducibility requires going beyond openness.

# reana

Reproducible research data analysis platform

| Flexible | Scalable | Reusable | Free |
|---|---|---|---|
| Run many computational workflow engines. | Support for remote compute clouds. | Containerise once, reuse elsewhere. Cloud-native. | Free Software. MIT licence. Made with ♥ at CERN. |

COMMON WORKFLOW LANGUAGE

kubernetes

CERN

http://reanahub.io

---

**2 | Example of a complex computational workflow on REANA mimicking a beyond the standard model (BSM) analysis.** This figure shows an ...ple where the experimental data is compared to the predictions of the standard model with an additional hypothesized signal component. The ...ple permits one to study the complex computational workflows used in typical particle physics analyses. **a–c.** The computational workflow (**a**) may ...ist of several tens of thousands of computational steps that are massively parallelizable and run in a cascading 'map-reduce' style of computations ...stributed compute clusters. The workflow definition is modelled using the Yadage workflow specification and produces an upper limit on the ...l strength of the BSM process. A typical search for BSM physics consists of simulating a hypothetical signal process (**c**), as well as the background ...sses predicted by the standard model with properties consistent with the hypothetical signal (marked dark green in (**b**)). The background often ...ists of simulated background estimates (dark blue and light green histograms) and data-driven background estimates (light blue histogram). ...tistical model involving both signal (dark green histogram) and background components is built and fit to the observed experimental data (black ...ers). **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 ...nal 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 ...cly available at ref. [25]. For icon credits, see Fig. 1.

# BUILD IT AND THEY WILL COME

In 2010 we <u>identified</u> a use-case with high scientific value for community
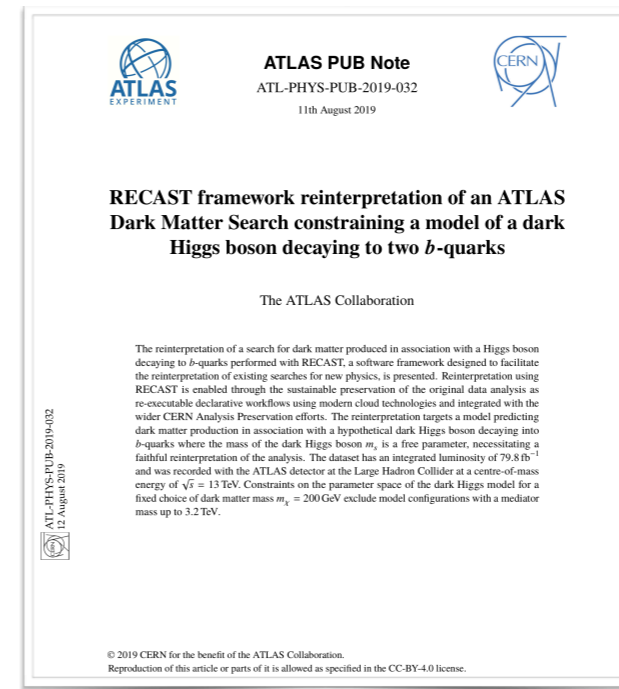
- Conservative narrative compared to "open data"

- Not conservative enough for many. Lots of resistance

- People said it couldn't be done, our workflows are too complicated

- Hard to get effort to work on it.

Got lucky with an amazing student that took a risk and just built it.

- Containers & Cloud technology

- 9 years later …

recast

Lukas Heinrich

**ATLAS PUB Note**
ATL-PHYS-PUB-2019-032
11th August 2019

**RECAST framework reinterpretation of an ATLAS Dark Matter Search constraining a model of a dark Higgs boson decaying to two $b$-quarks**

The ATLAS Collaboration

The reinterpretation of a search for dark matter produced in association with a Higgs boson decaying to $b$-quarks performed with RECAST, a software framework designed to facilitate the reinterpretation of existing searches for new physics, is presented. Reinterpretation using RECAST is enabled through the sustainable preservation of the original data analysis as re-executable declarative workflows using modern cloud technologies and integrated with the wider CERN Analysis Preservation efforts. The reinterpretation targets a model predicting dark matter production in association with a hypothetical dark Higgs boson decaying into $b$-quarks where the mass of the dark Higgs boson $m_s$ is a free parameter, necessitating a faithful reinterpretation of the analysis. The dataset has an integrated luminosity of $79.8\,\text{fb}^{-1}$ and was recorded with the ATLAS detector at the Large Hadron Collider at a centre-of-mass energy of $\sqrt{s} = 13$ TeV. Constraints on the parameter space of the dark Higgs model for a fixed choice of dark matter mass $m_\chi = 200$ GeV exclude model configurations with a mediator mass up to 3.2 TeV.

ATL-PHYS-PUB-2019-032
12 August 2019

## Analysis preservation bootcamp



Participants in Analysis Preservation Bootcamp showing off their ability to reproduce an LHC analysis. Photo Credit: Samuel Meehan
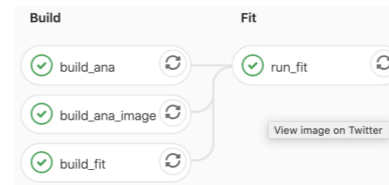


**Josh McFayden**
@JoshMcFayden

Thoroughly enjoying myself at an @iris_hep/@diana_hep analysis preservation bootcamp @CERN today!
indico.cern.ch/event/854880/o…

**Josh McFayden** @JoshMcFayden · Feb 18, 2020
Replying to @JoshMcFayden and 3 others
PROGRESS ✅

**Josh McFayden**
@JoshMcFayden

Today: REANA ✅

34

## JSON Patch for signal model (reinterpretation)

```
$ pyhf cls example.json | jq .CLs_obs
0.053994246621274014

$ cat new_signal.json
[{
    "op": "replace",
    "path": "/channels/0/samples/0/data",
    "value": [10.0, 6.0]
}]

$ pyhf cls example.json --patch new_signal.json | jq .CLs_obs
0.3536906623262466
```

Original analysis (model A)

Recast analysis (model B)

active learning / sequential design / black box optimization

**Active Sciencing**

reusable workflows

simulation-based /
likelihood-free
inference engines

# Scalable cyberinfrastructure applications

Team: J. Brehmer[1,2], K. Cranmer[1,2], **Irina Espejo**[1], S. Macaluso[1,2] and H. Müller[1]

Institutions: [1] Center for Data Science, New York University
[2] Department of Physics, New York University

The SCAILFIN Project

## ROB
### Reproducible Open Benchmark Platform

in collaboration with S. Hsu[4], A. Maritz[4], A. Rawat[4] and C. Suaysom[4]
[4]University of W

- ROB is an experimental prototype for enabling *community benchmarks* of data analysis algorithms. The goal of ROB is to allow user communities to evaluate the performance of their different data analysis algorithms in a *reproducible competition-style* format.

- The *workflow template* and input data are defined by a coordinator. The template contains placeholders for workflow steps that are implemented by the participants (e.g., with Docker containers). The backend submission workflows. The user participants to *submit new runs* results.

Physics simulations

Leaderboard

configure

generate simulations physics + detector

Front-End

flowServ
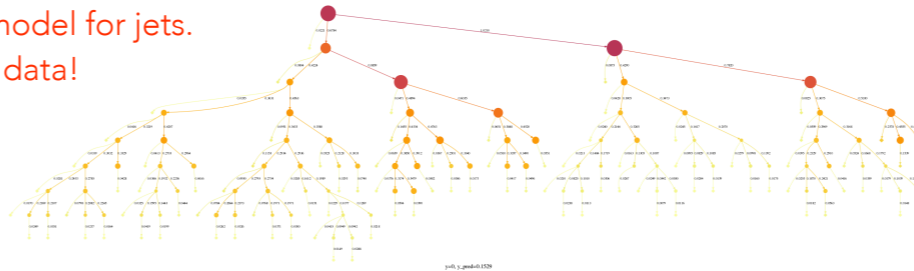
combine

sample

Machine Learning Inference

train + evaluate

Scientist

ailfin/flowserv-core
ailfin/rob-webapi-flask
ailfin/rob-ui
astianMacaluso/TopTagComparison

We can use the sa... ...ology to streamline compar... of up-stream tools like ML-based jet taggers.
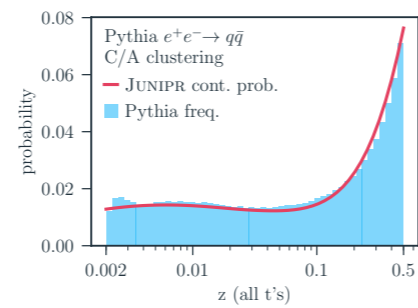
- Building ROB & F... ...on top of

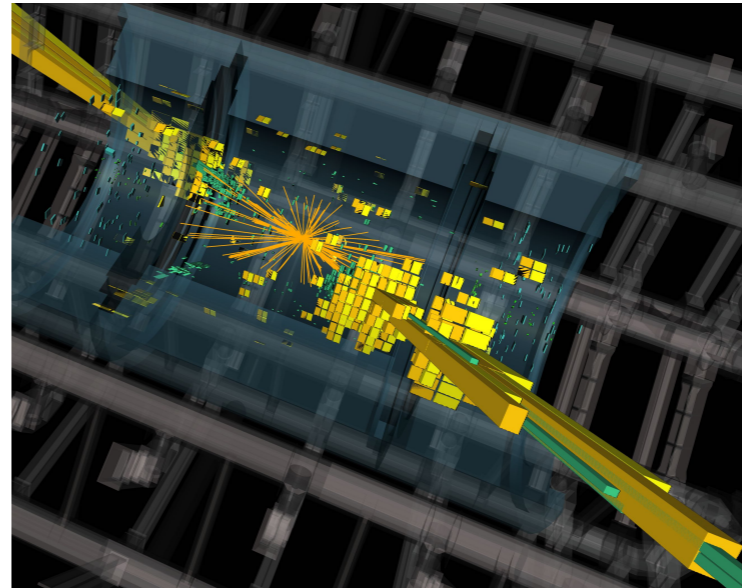JUNIPR is a generative model for jets.
Can train on real data!



tractable likelihood

$$P_{\text{jet}}(\{p_1, \ldots p_n\}) = \left[ \prod_{t=1}^{n-1} P_t\big(k_1^{(t+1)}, \ldots, k_{t+1}^{(t+1)} \big| k_1^{(t)}, \ldots, k_t^{(t)}\big) \right] \times P_n\big(\text{end} \big| k_1^{(n)}, \ldots, k_n^{(n)}\big).$$

… and it is interpretable
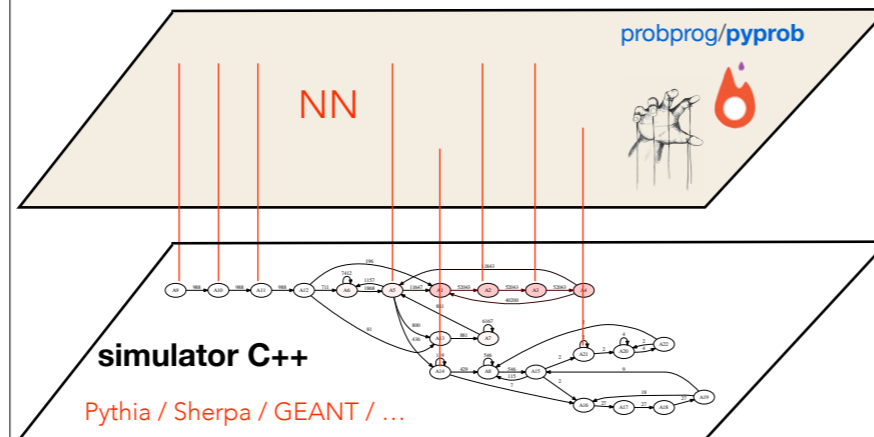




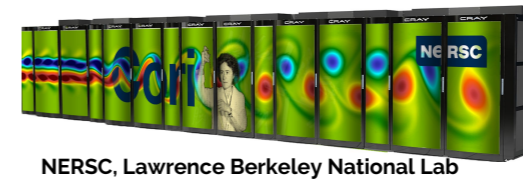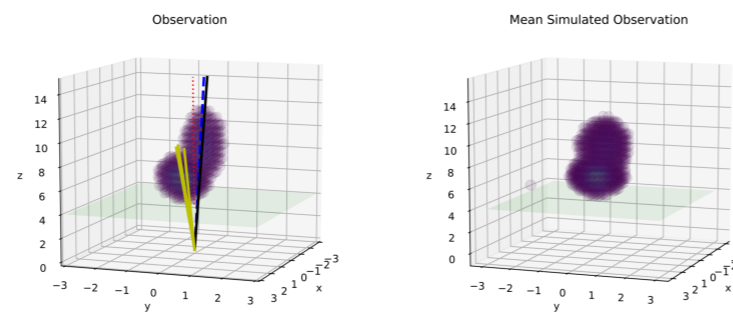Andreassen, Feige, Frye, Schwartz  arXiv:1804.09720

**Idea:** hijack the random number generators and use Neural Network to perform a *very* fancy type of importance sampling



probprog/**pyprob**

NN

**simulator C++**

Pythia / Sherpa / GEANT / …

- Neural Network powered inference engine (python)

- real-world scientific simulator (C++)

Observation

Mean Simulated Observation

**NERSC, Lawrence Berkeley National Lab**

arXiv:1807.07706

39

Highlight

https://arxiv.org/abs/1907.03382
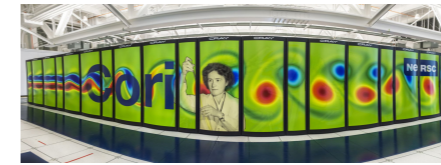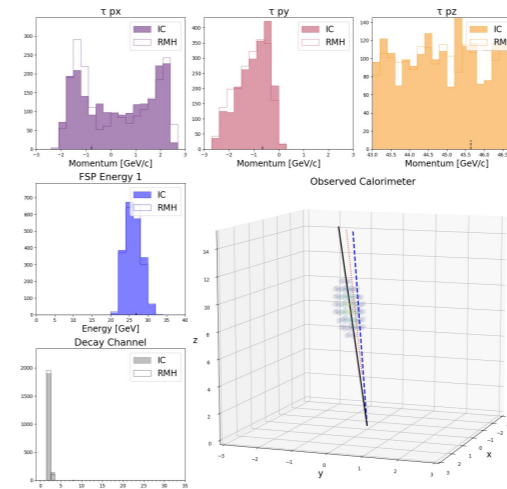
Finalist for best paper award at SC19 (Super Computing)

- Largest scale Bayesian inference ever using in a universal probabilistic programming language
  - *Applied to complex LHC Physics use case: Sherpa code base of ~1M lines of code in C++*

- 230x speedup for synchronous data parallel training of a 3DCNN-LSTM neural network
  - *1,024 nodes (32,768 CPU cores)*
  - *128k minibatch size, largest for this NN architecture*
  - *One of the largest-scale use of PyTorch built-in MPI*

- Novel protocol (PPX) to execute & control existing, large-scale, scientific simulator code bases

# SUMMARY

Accelerating analysis design

- more powerful observables

- end-to-end optimization

- benchmarking of algorithms

Accelerating fitting

- pyhf and a fitting service

More efficient simulation

- excursion

- Probabilistic programming

Extending impact of results

- RECAST

Core technologies:

- automatic differentiation

- GPUs & TPUs

- Cloud-native : docker, kubernetes

- Workflows & REANA

- Functions as a service Accelerating analysis design

Sorry if this is a little disorganized, COVID has complicated work life