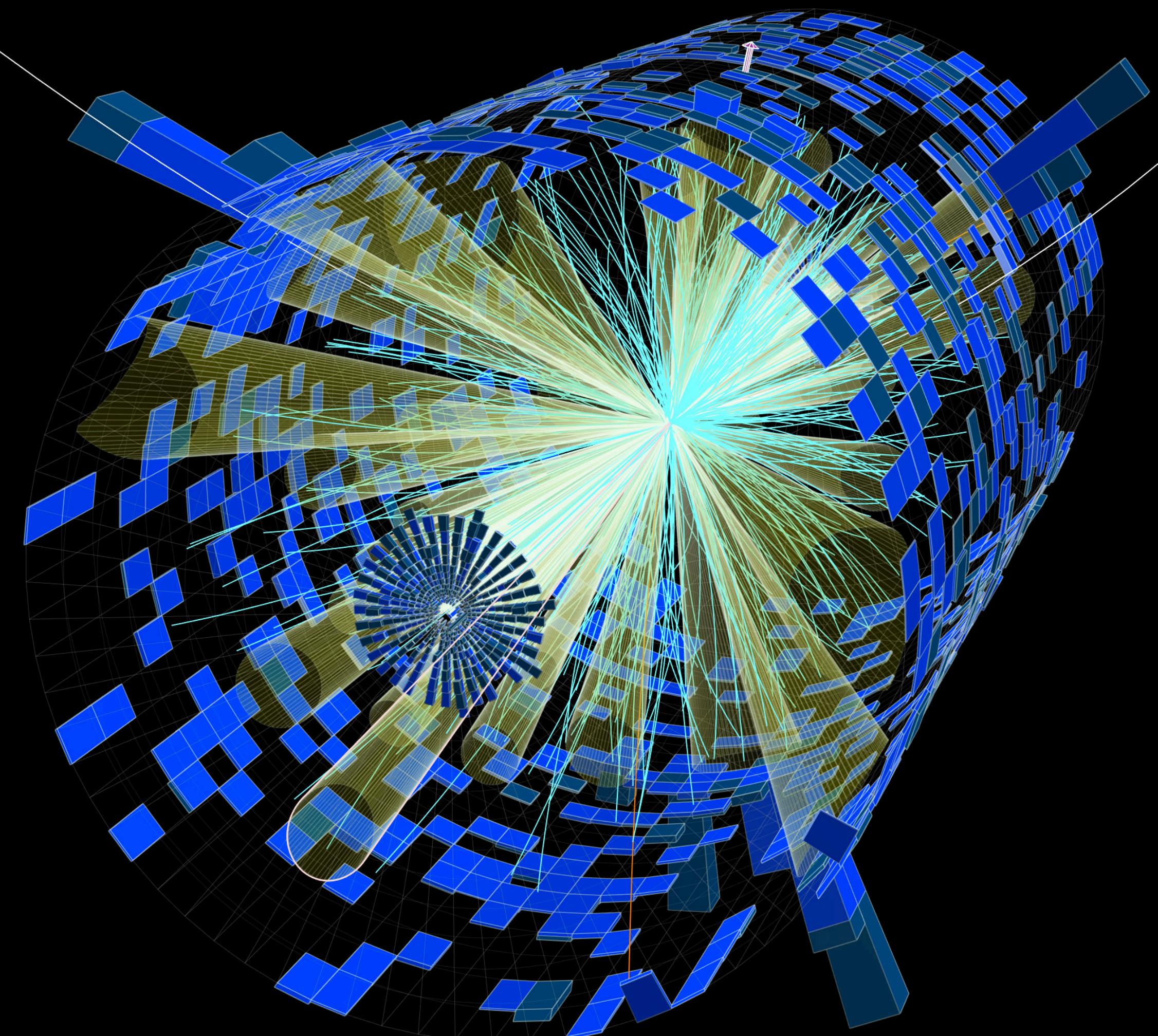




SciPy 2024

Particles, People, and Pull Requests



@KyleCranmer

University of Wisconsin-Madison

Data Science Institute, Physics, Statistics, Computer Science

Acknowledgements



GORDON AND BETTY
MOORE
FOUNDATION



Alfred P. Sloan
FOUNDATION

DASPOS

dianahep

iris
hep

FAIROS-HEP
SCAILFIN



Matthew Feickert



Lukas Heinrich



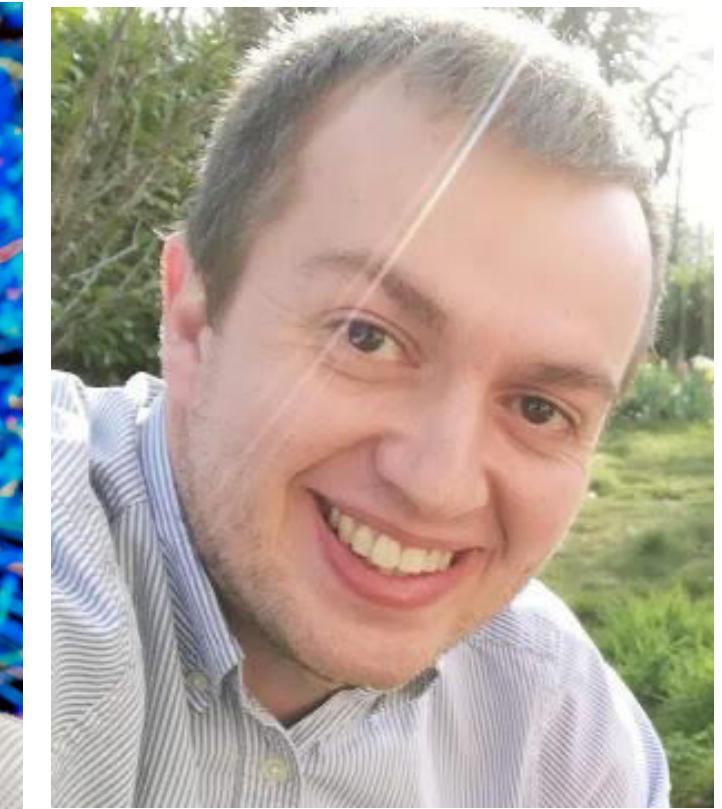
Giordon Stark



Alex Held



Oksana Shadura



Gilles Louppe



Johann Brehmer



Jim Pivarski



Ianna Osborne



Henry Schreiner



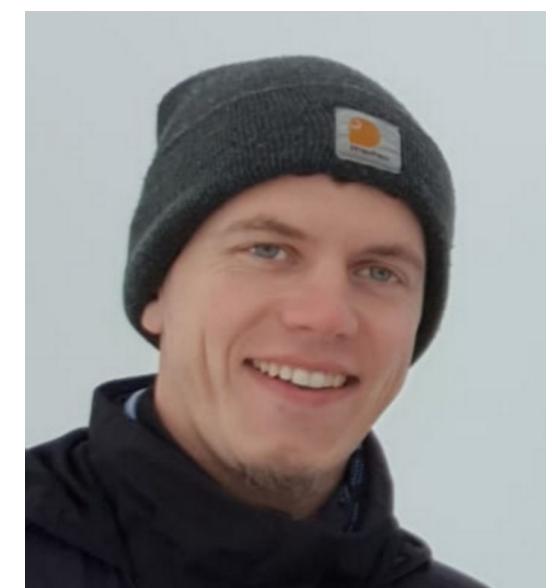
Jonas Eschle



Gordon Watts



Vangelis Kourlitis



Patrick Rieck



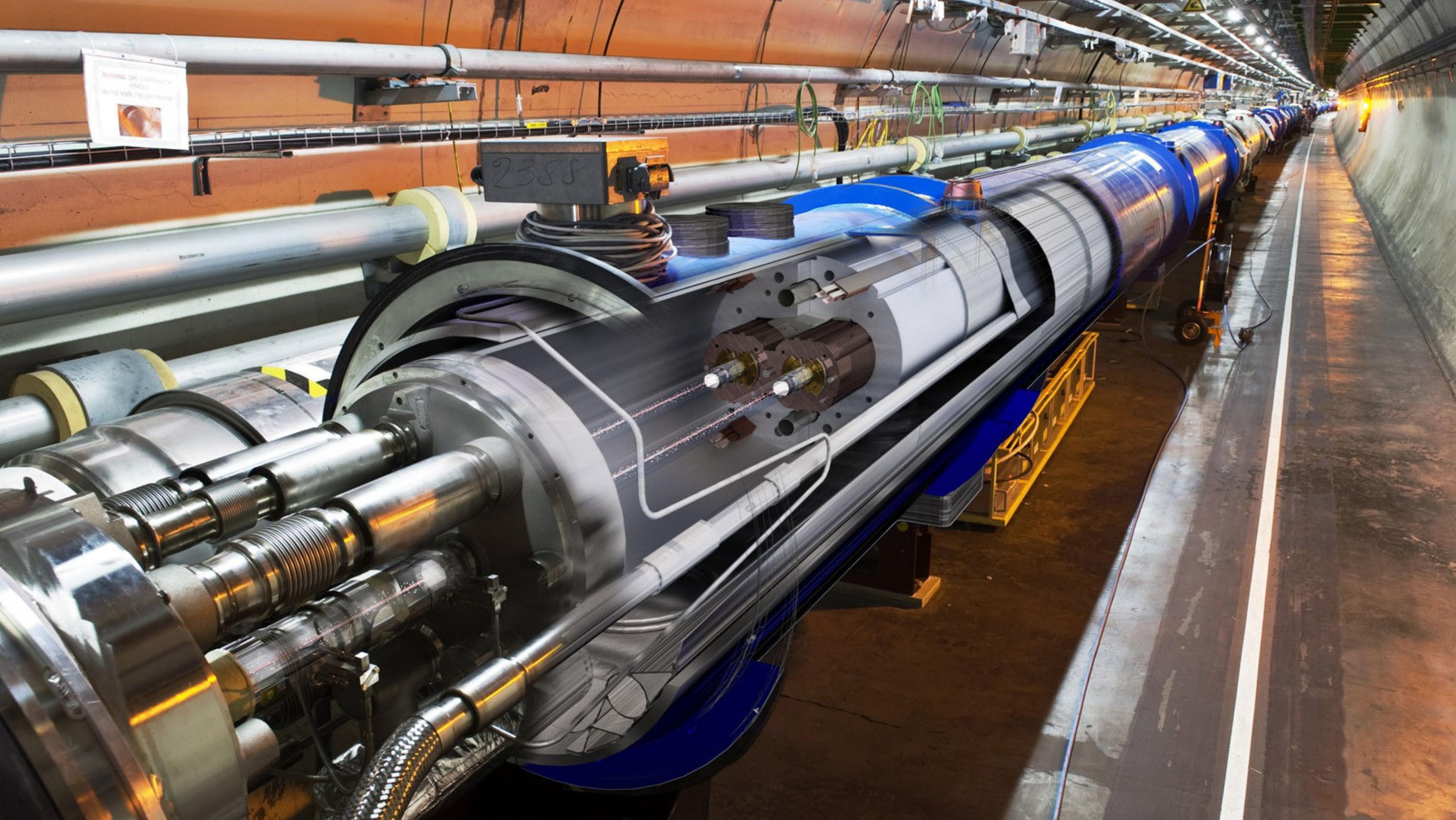
Zubair Bhatti

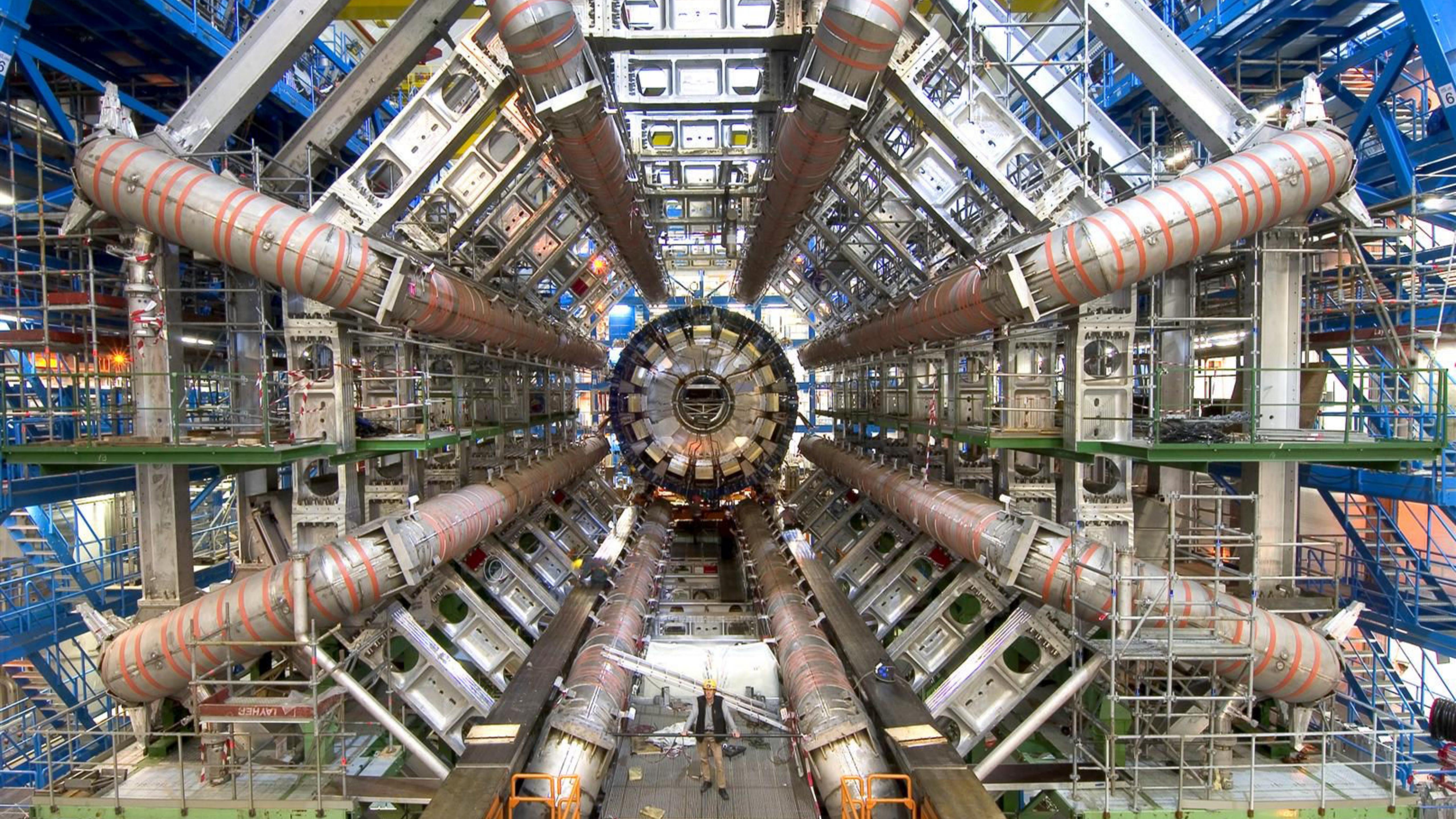


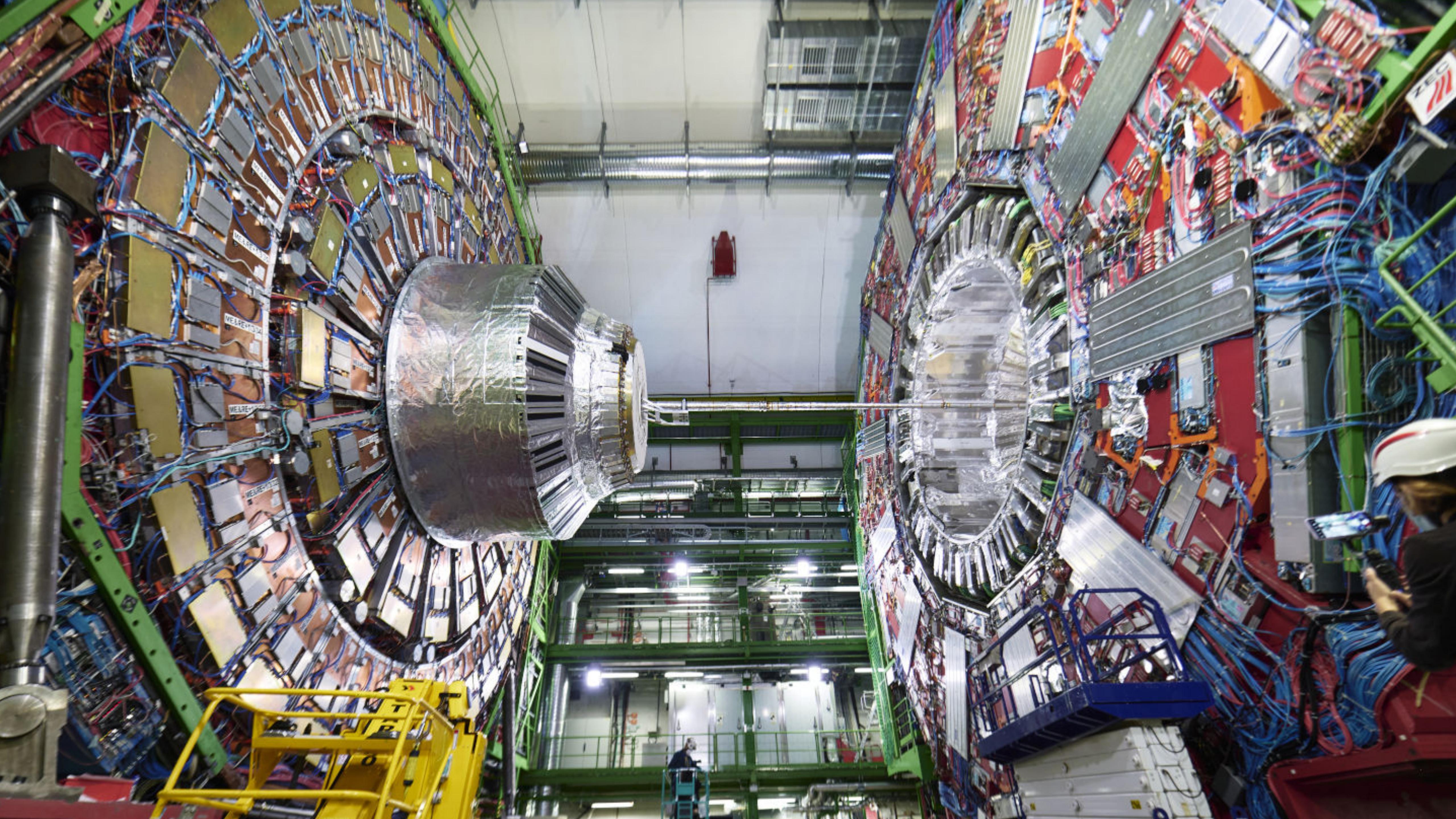
Irina Espejo

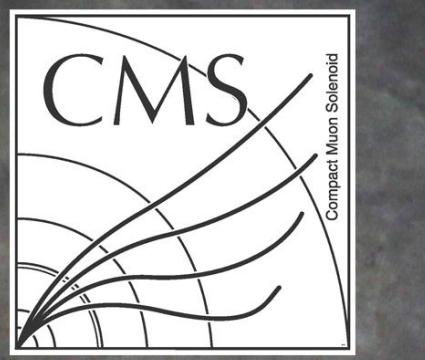
High Energy Physics
&
The Large Hadron Collider











1992

2017



Big Numbers

40 quadrillion collisions

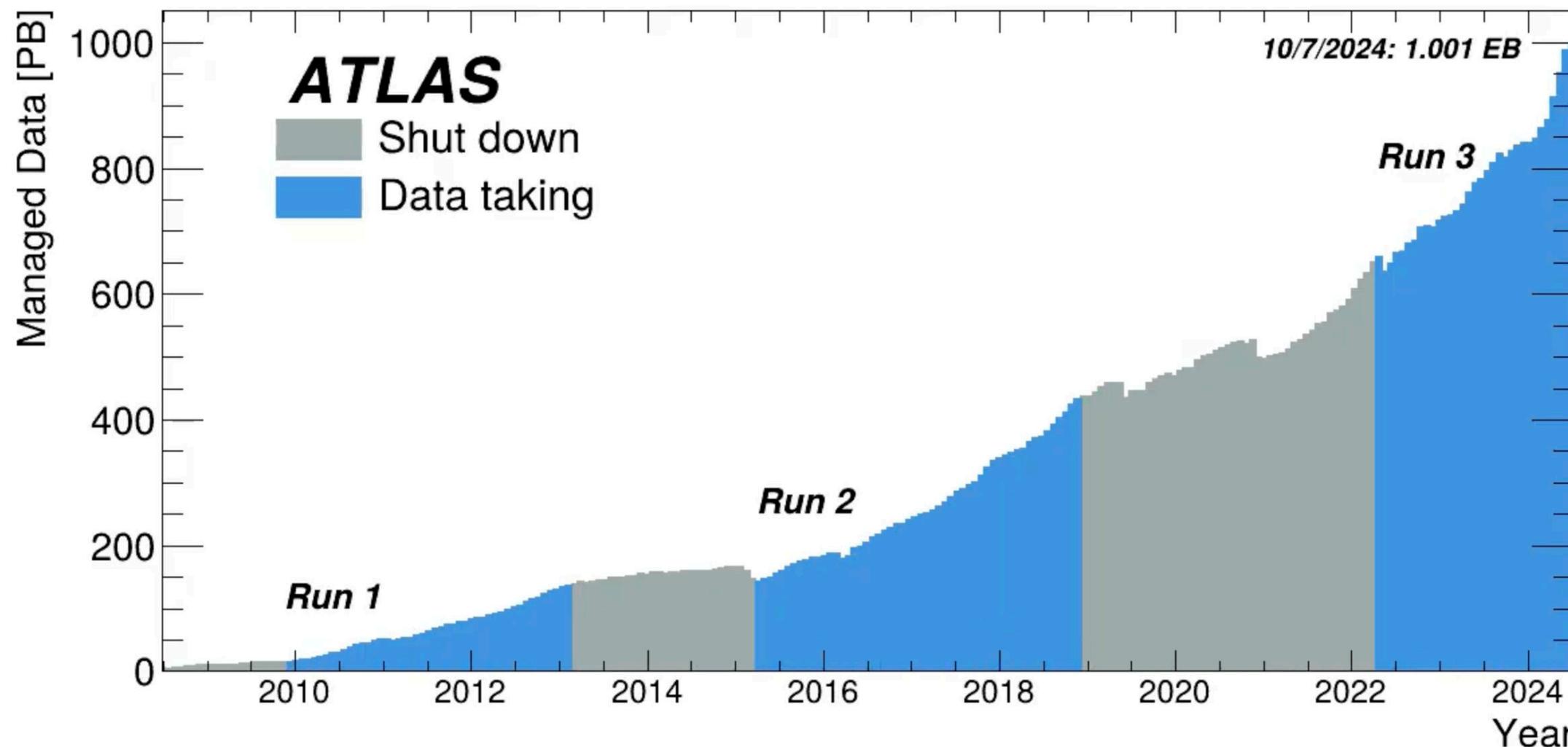
The number of collisions recorded by the four main experiments at the LHC is close to 40 quadrillion, or, as physicists say it, 395 “inverse femtobarns.” (Each inverse femtobarn corresponds to about 100 trillion collisions.)

7.5 billion

Worldwide LHC Computing Grid requests

Physicists need a huge amount of computing power to do their research—much more than a standard laptop can support. Every day several thousand physicists submit a total of about 2 million “jobs” to the WLCG. Each “job” is an important brick in the growing body of scientific work.

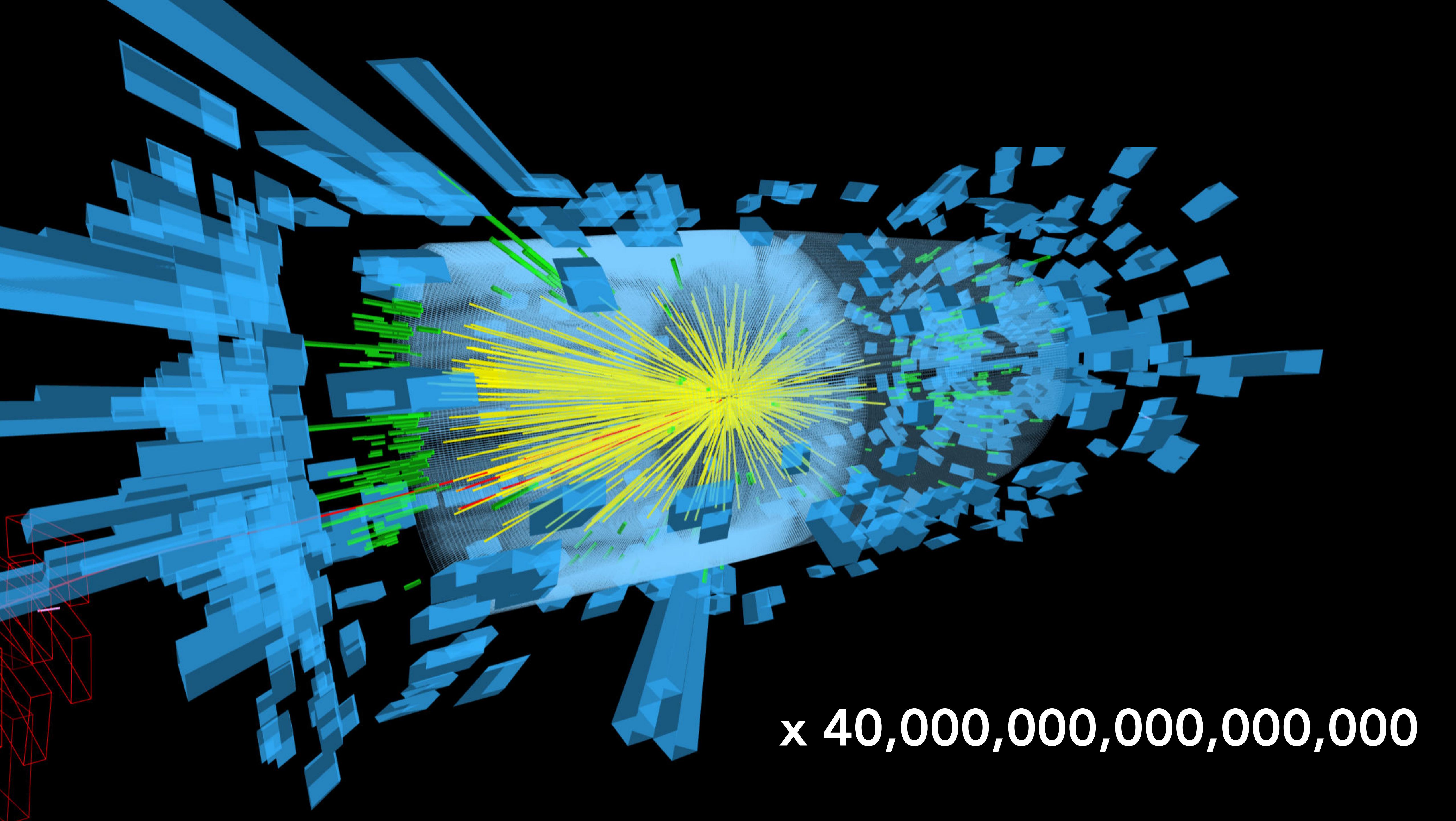
1 exabyte



2,725

scientific papers

Every week the number of scientific papers that LHC scientists have published steadily increases as they comb through the data to study rare phenomena and search for new physics. This includes work by thousands of graduate students on their way to earning their PhDs.



$\times 40,000,000,000,000$

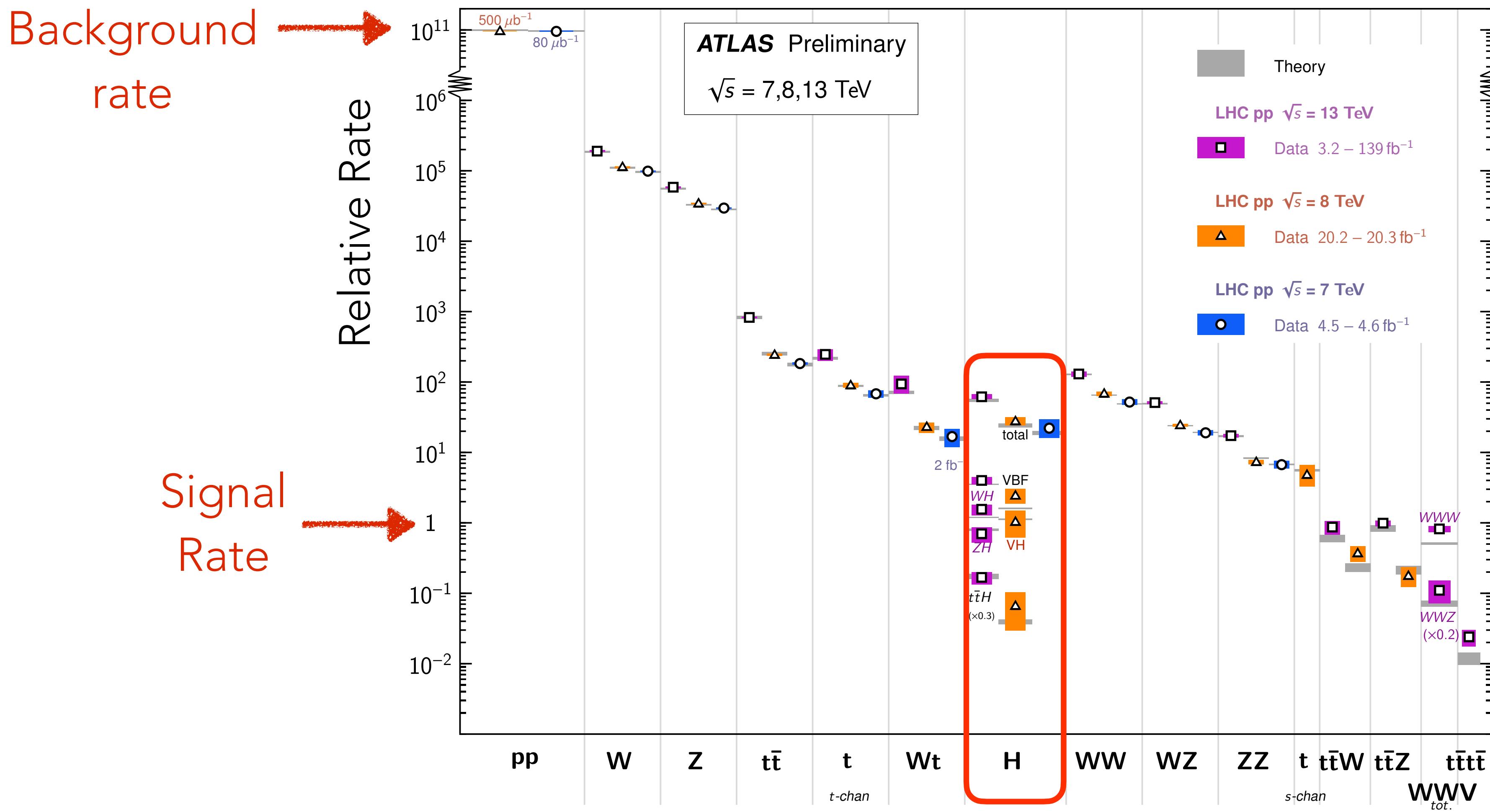
40,000,000,000,000,000



Needle in haystack

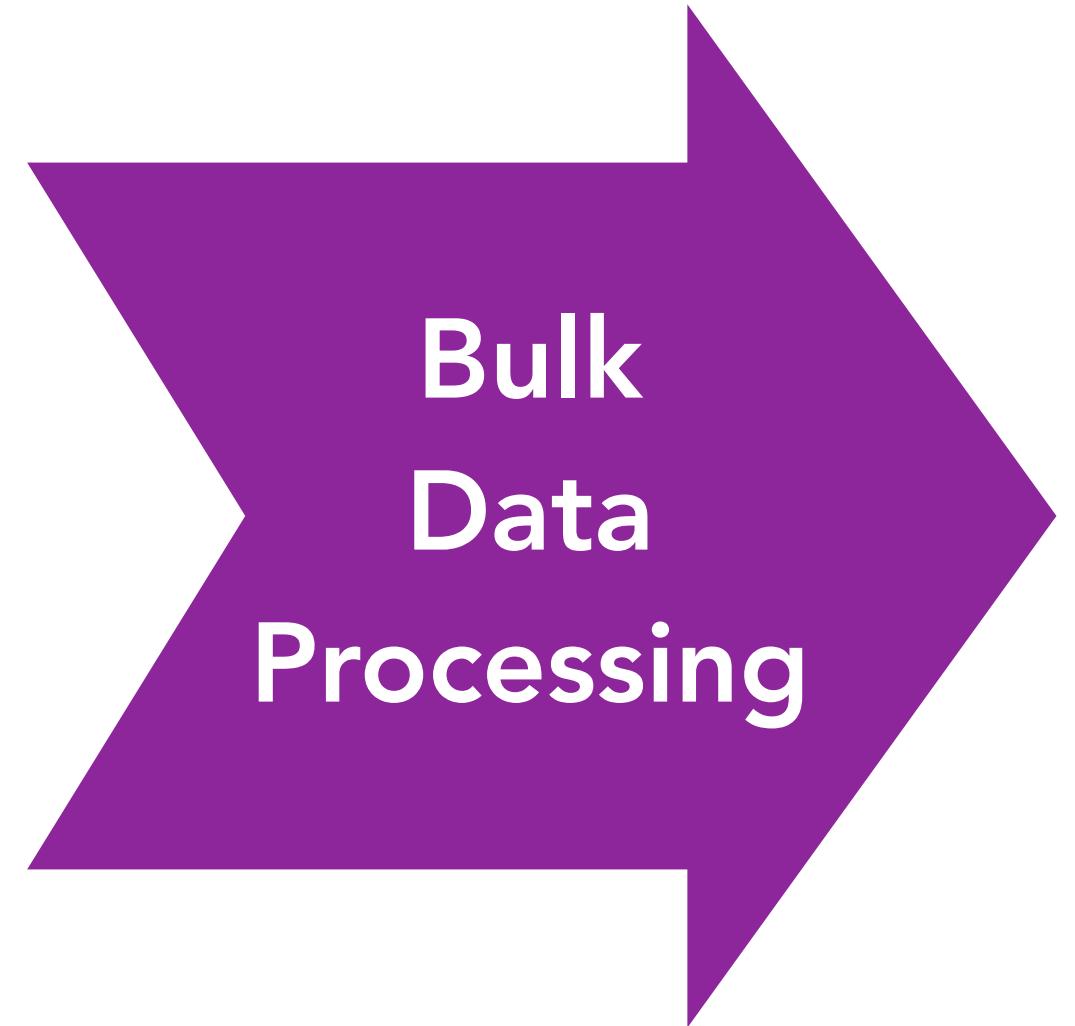
LHC collides bunches of protons at 40 MHz. Detector has $\sim 100M$ sensors

- Signals of interest are rare, with S/B $\sim 10^{-10}$

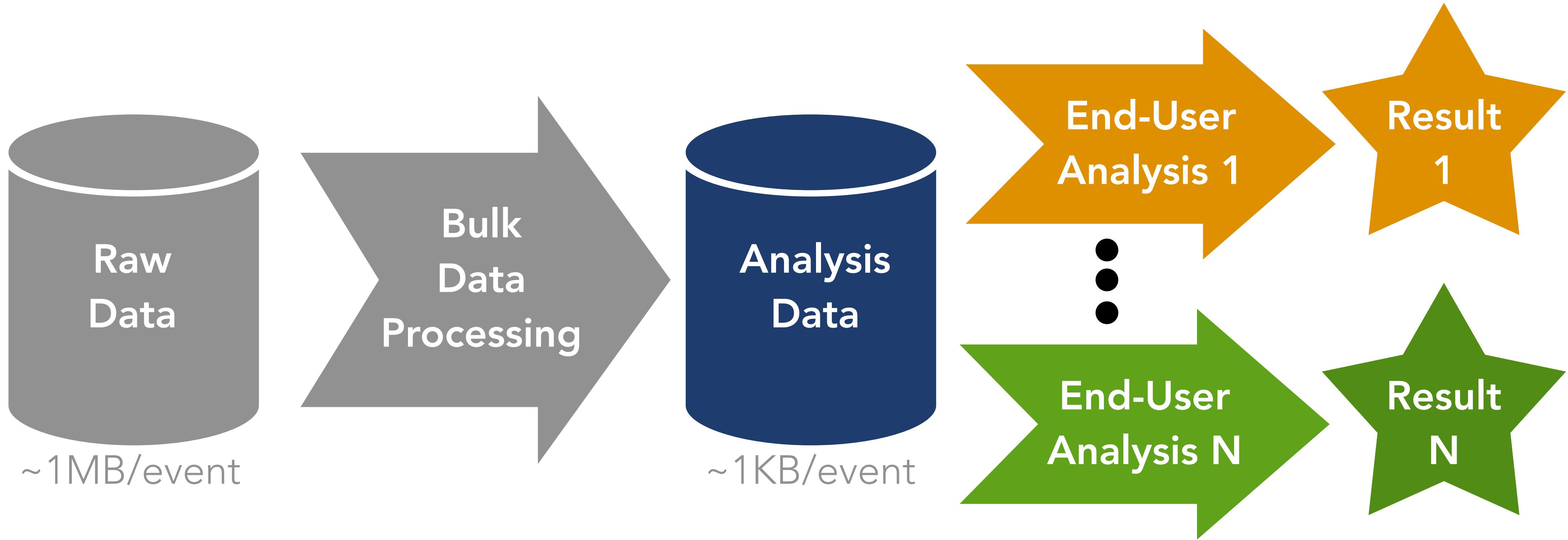


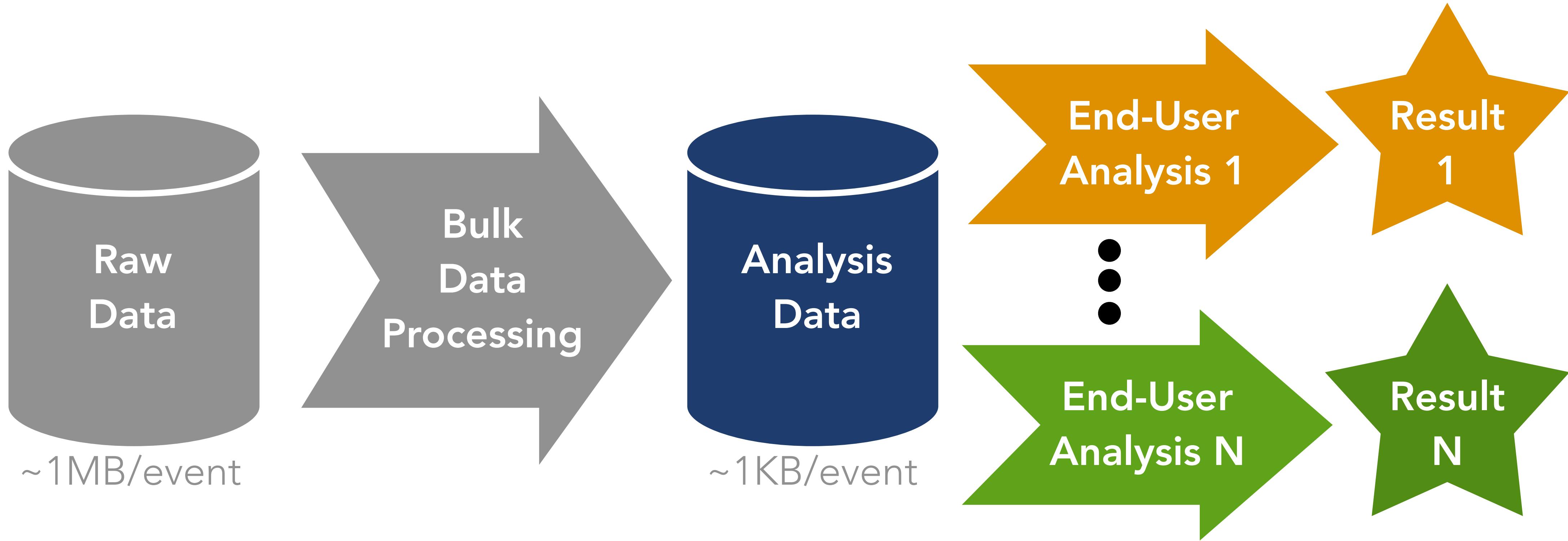


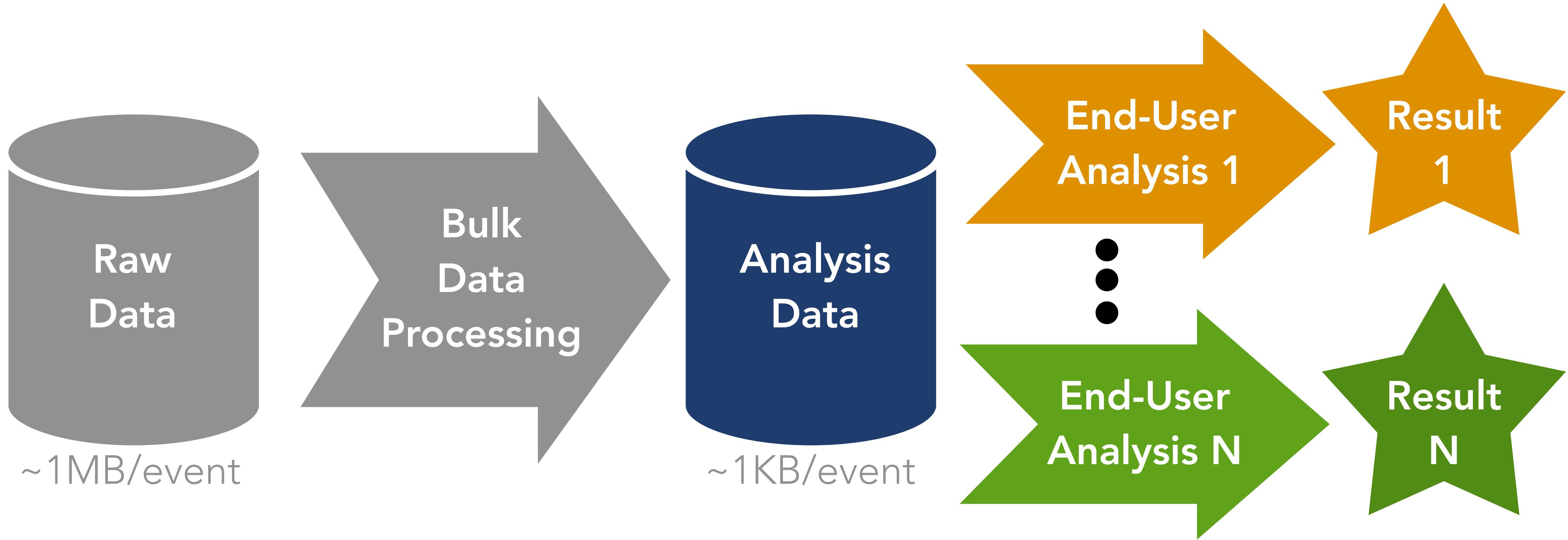
~1MB/event

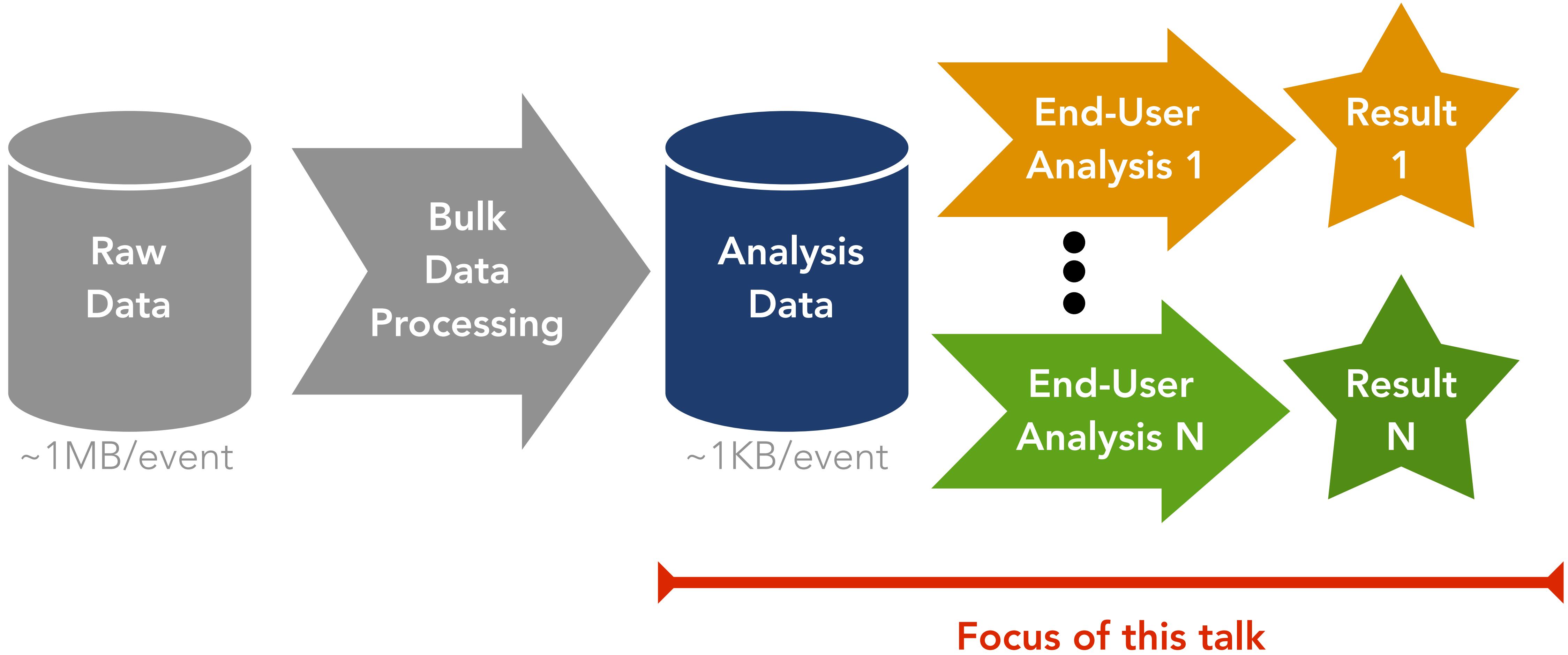


~1KB/event







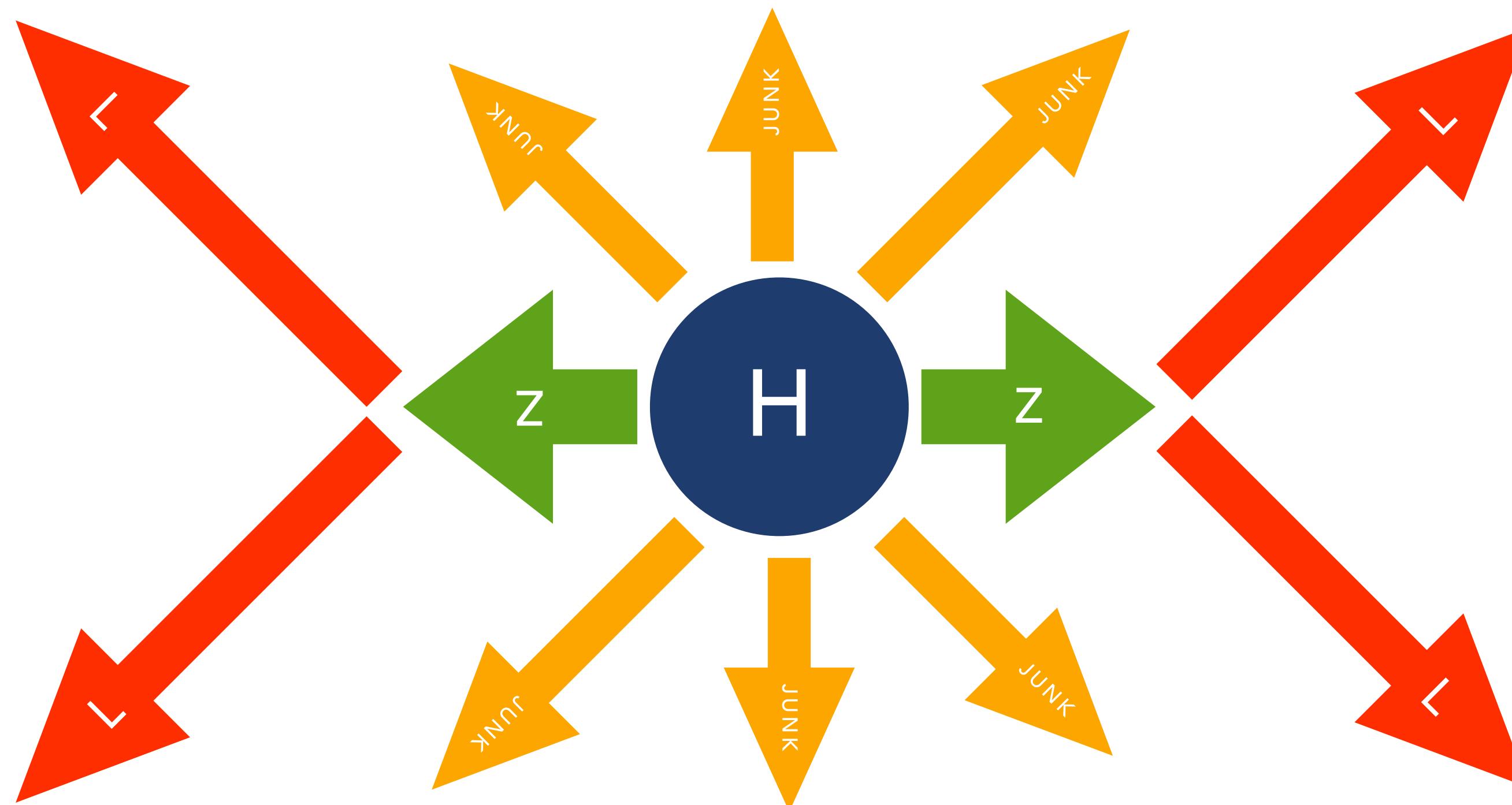


Example hypothesis

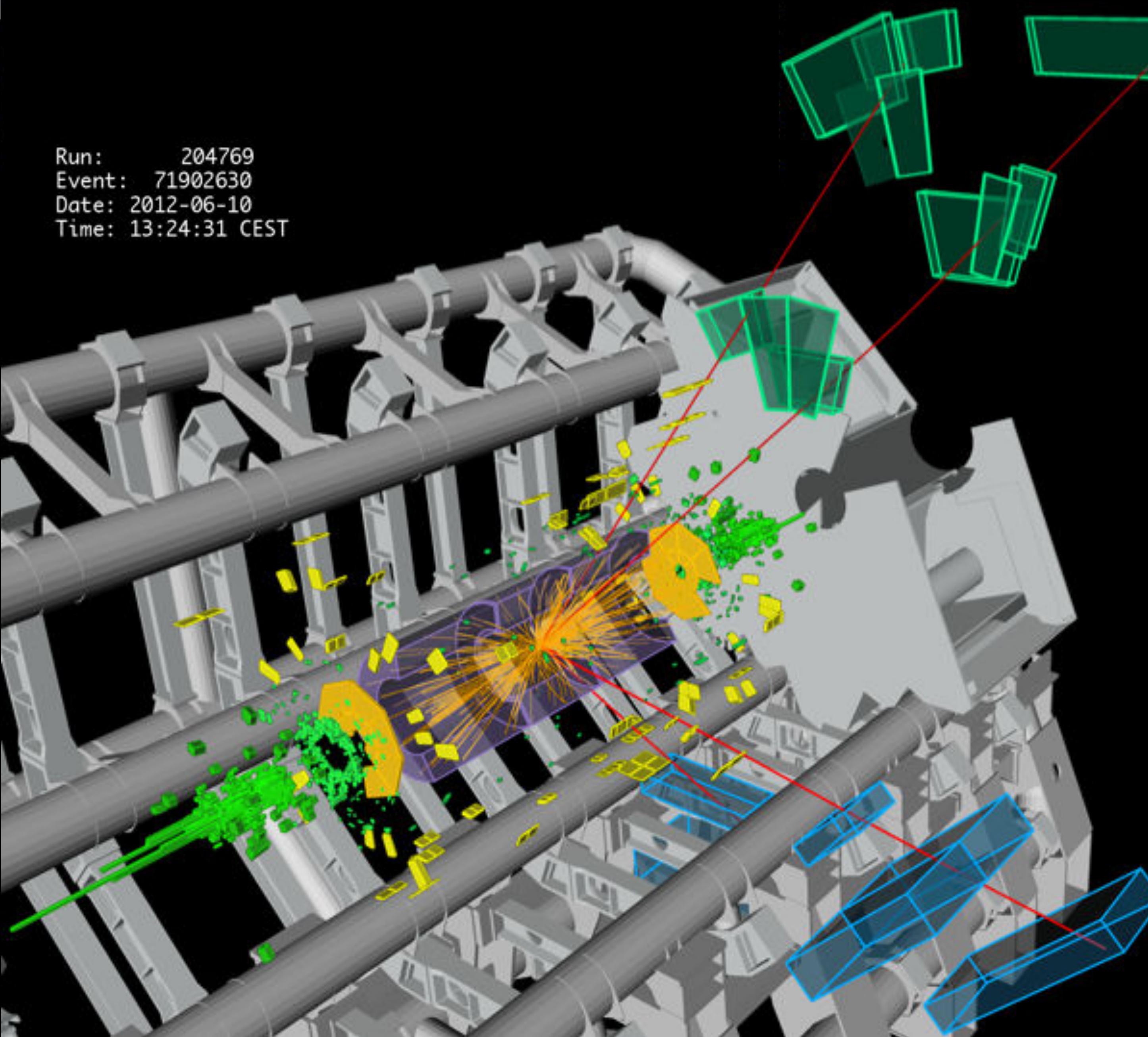
What if we produce (along with some other **junk**) an unstable particle **H**

- that decays into two other unstable particles **Z**,
- that each decay into two particles **L** that we detect in our sensors

What would be a good strategy to search for **H**?



Analysis basics



This “end-user analysis” is highly customized for a particular goal

- First **filter the data** to find collisions that are useful for testing a particular hypothesis
- Experts design *features / observables / summary statistics* that distinguish signal from background
- Perform statistical analysis

Expert feature engineering

Don't believe the media:

$$E \neq mc^2$$

What Einstein really said:

$$E^2 = (mc^2)^2 + (|\vec{p}|c)^2$$

Every physics student knows energy and momentum are conserved

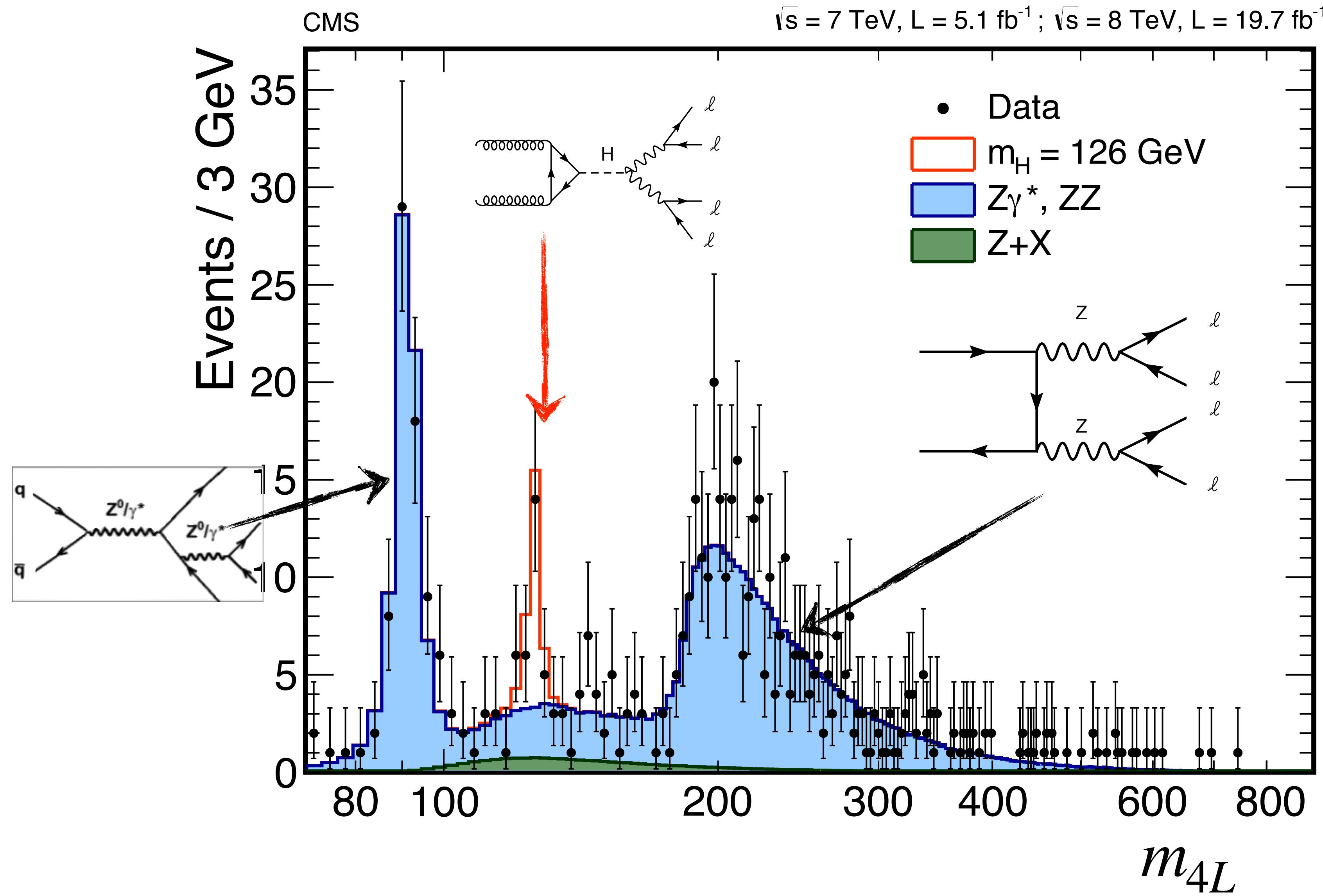
$$E_{\text{Higgs}} = E_{\text{before}} = E_{\text{after}} = \sum_i E_i$$

$$\vec{p}_{\text{Higgs}} = \vec{p}_{\text{before}} = \vec{p}_{\text{after}} = \sum_i \vec{p}_i$$

Thus, for our hypothesis, we expect a peak in this feature / observable:

$$m_{4L} = \sqrt{E_{\text{after}}^2/c^4 - |\vec{p}_{\text{after}}|^2/c^2}$$

A signal emerges

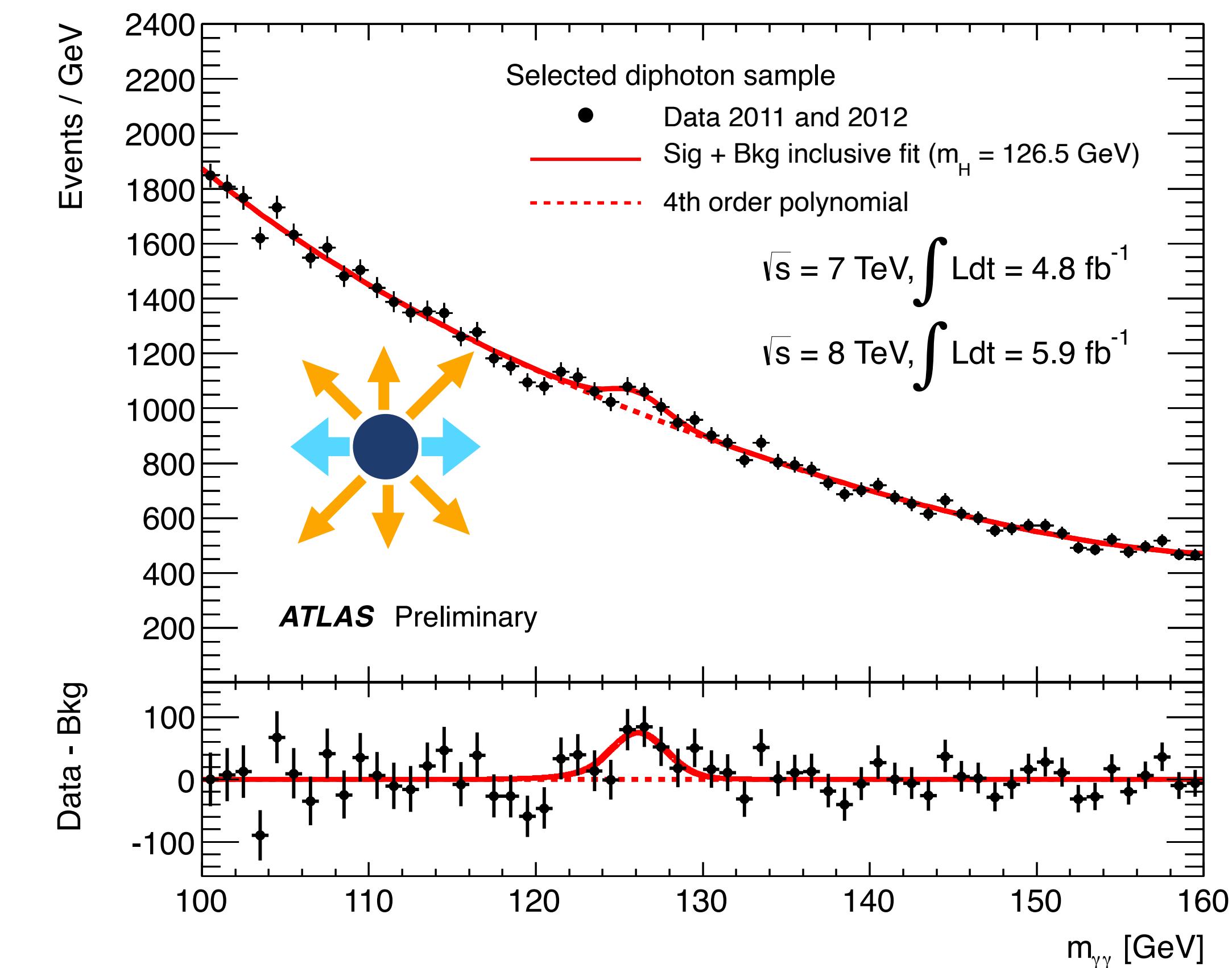
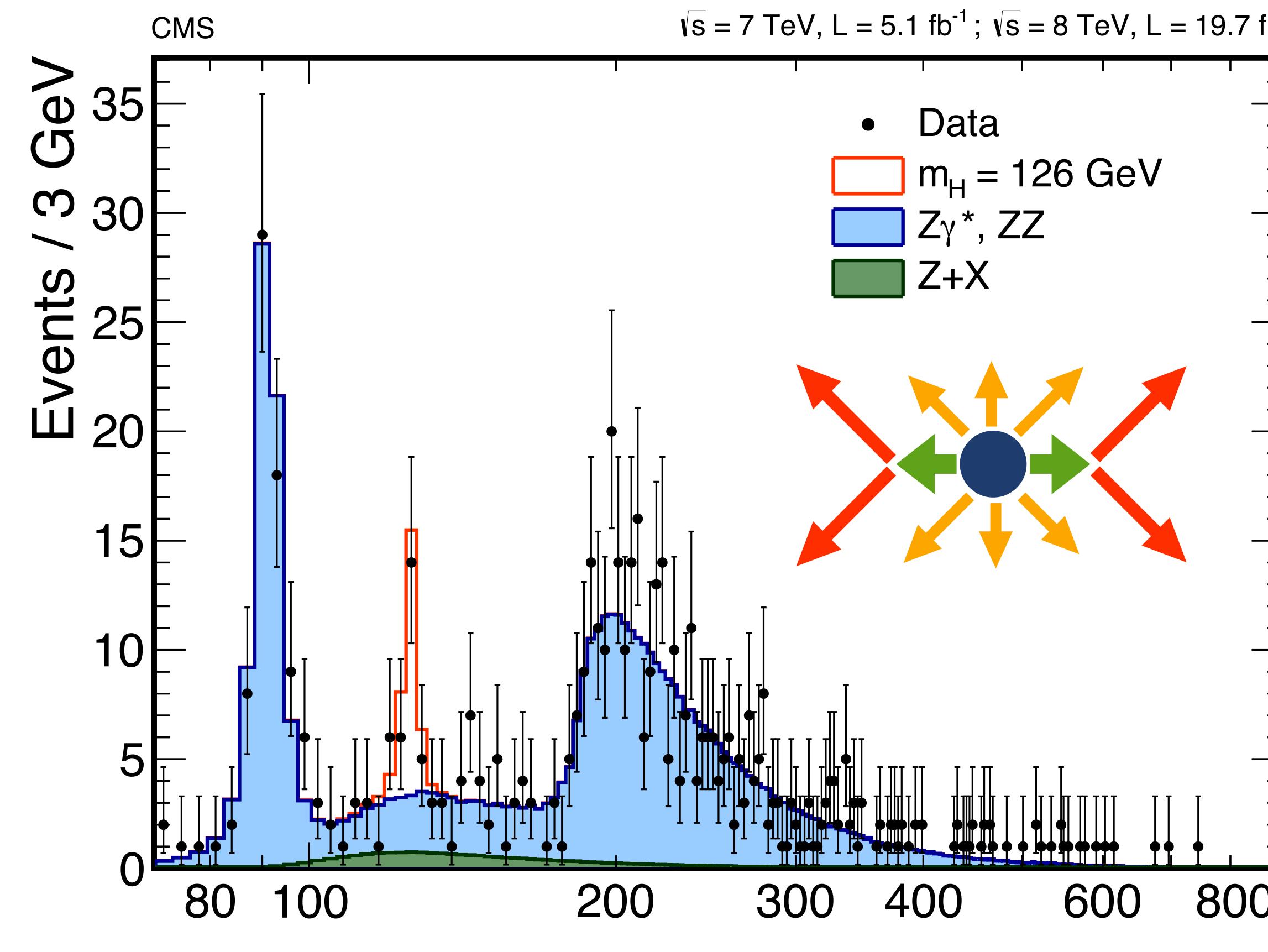


Collaborative Statistical Modeling

The power of an underlying theory

Hypotheses are connected to an underlying theory, which allows us to relate their predictions across multiple signatures in the data

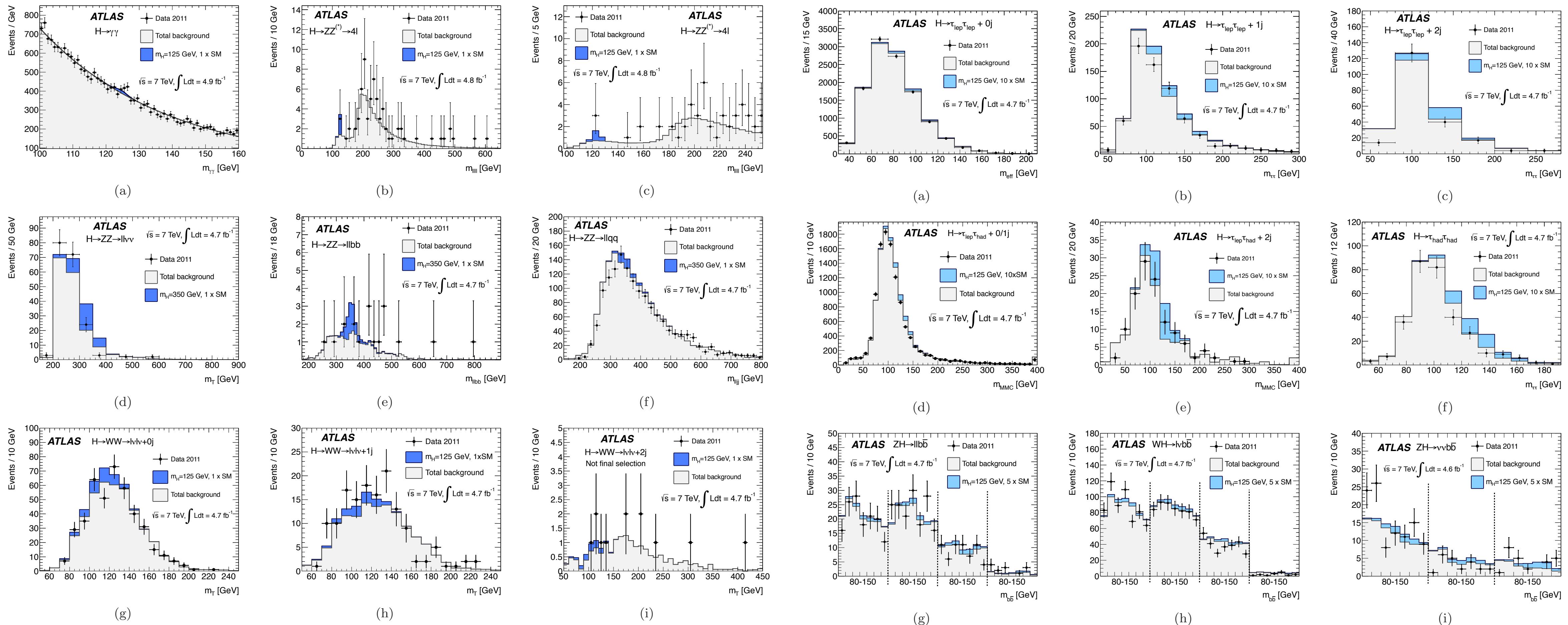
- Unstable particle **H** can decay to **Z** or **G**



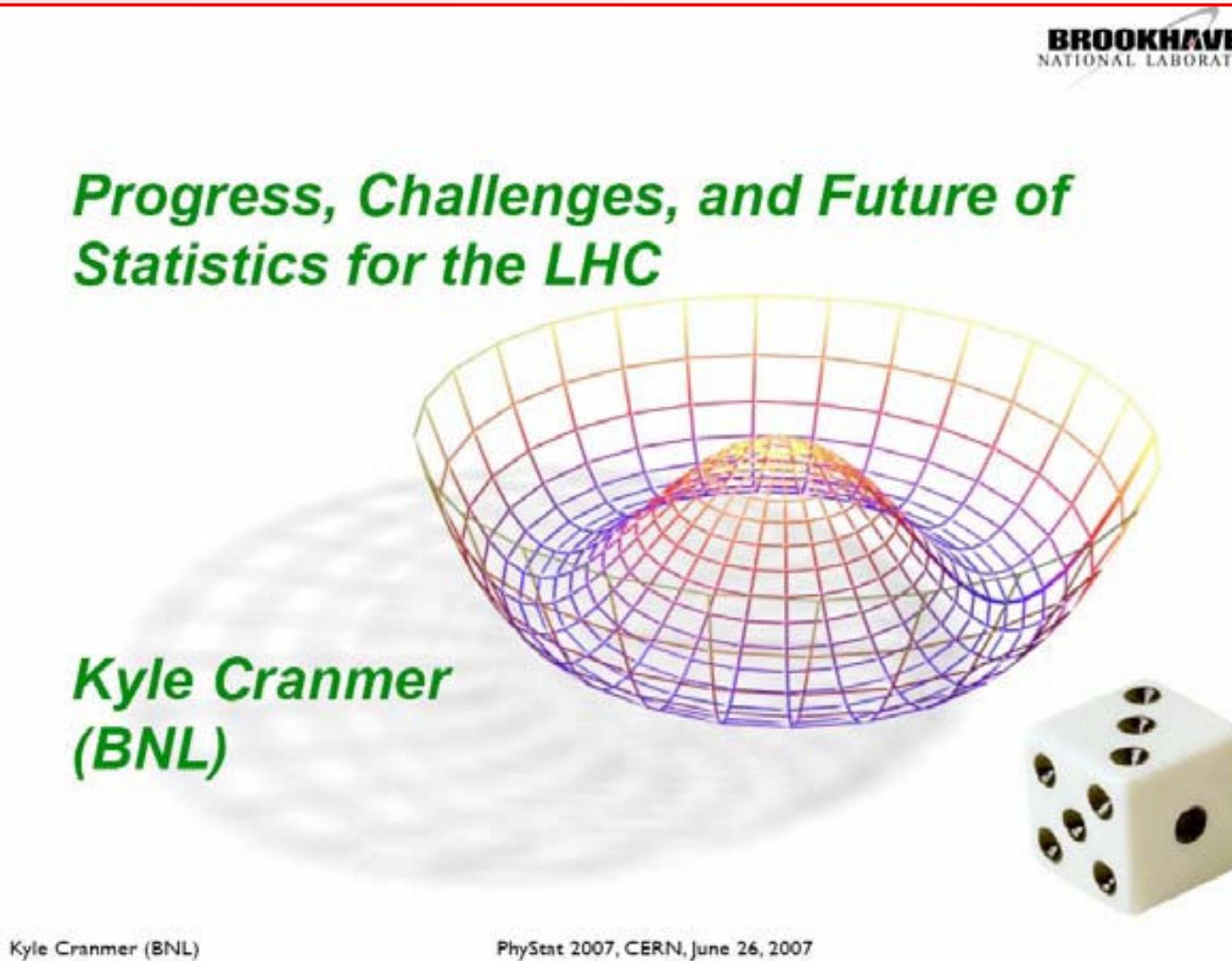
We want to combine the evidence statistically

The traditional approach was very top-down and centralized

- Tightly coupled, required a lot of coordination, moved slowly, very rigid.



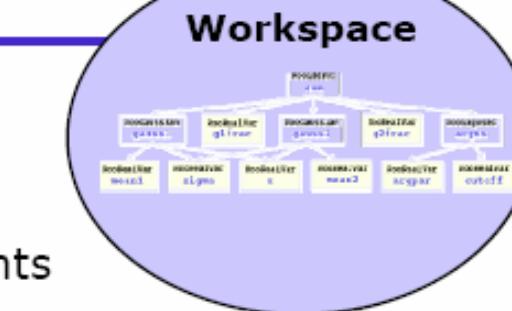
Preparing for the LHC in 2007



Statistics software for the LHC

The Workspace as publication

- Now have functional `RooWorkspace` class that can contain
 - Probability density functions and its components
 - (Multiple) Datasets
 - Supporting interpretation information (`RooModelConfig`)
 - Can be stored in file with regular ROOT persistence



Ultimate publication of analysis...

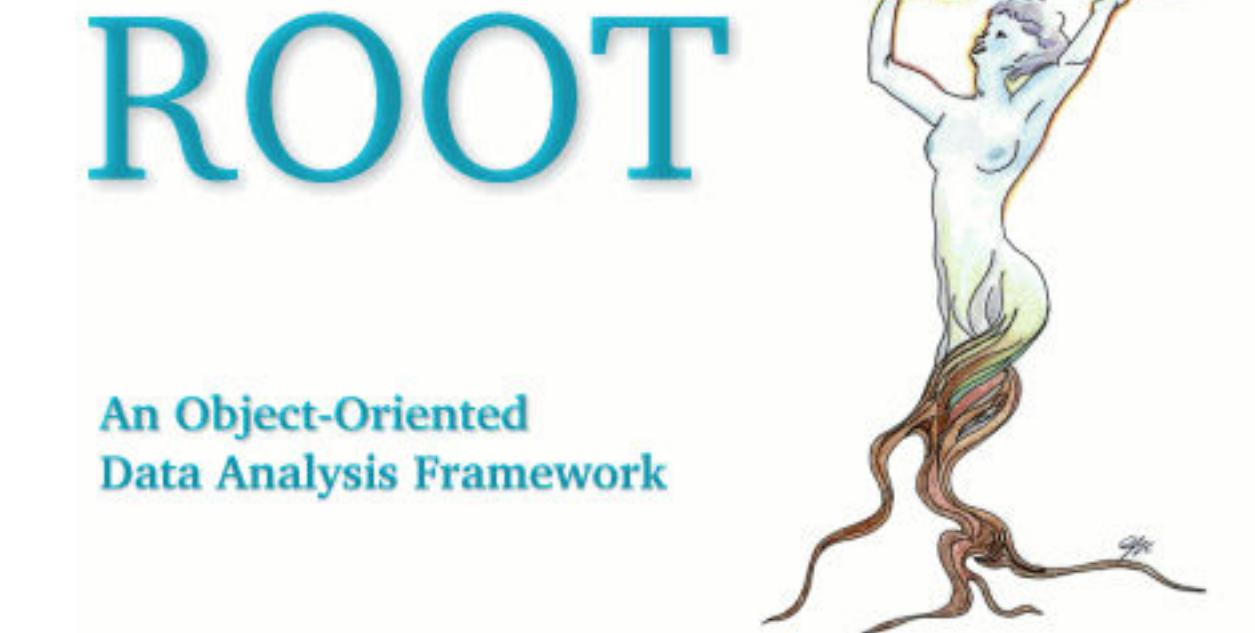
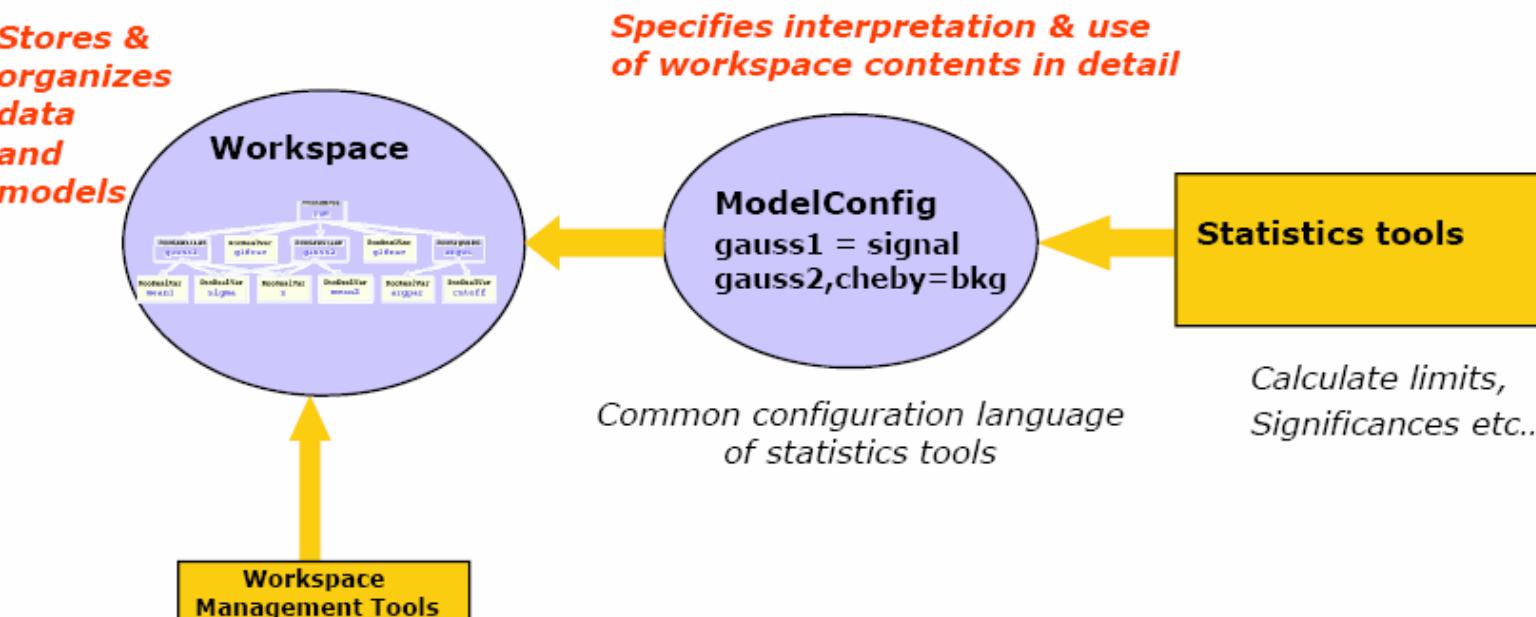
- Full likelihood available for Bayesian analysis
- Probability density function available for Frequentist analysis
- Information can be easily extracted, combined etc...
- Common format for sharing, combining of various physics results

Wouter Verkerke



Framework design & RooFit adaptations

- Have had more meetings last 3 months to review RooFit lessons from BaBar
 - Kyle, Amir Farbin (ex-Babar), Frank Wrinkmeyer (ex-Babar), WV
 - Design for `WorkSpace` and `ModelConfig` concept in RooFit to interface with statistics tools



ROOT

An Object-Oriented
Data Analysis Framework

ROOT Team Meeting

Friday 18 Apr 2008, 11:00 → 13:00 Europe/Zurich
32/1-A24 (CERN)

11:00 → 11:05 News

- New schema Evolution document will be ready early next week (Lukasz)
- Bertrand poster
- Document by Axel/Philippe on new multi-threaded CINT C++ design
- Progress report next meeting
- David/Lorenzo,
- Axel
- Olivier: progress with TGraph restructure

11:05 → 11:30 SVN restructuring status

- still to be done
 - /doc and ReleaNotes.html in each top level dir
 - \$ROOTSYS/doc pointing to the top chapters and release notes
- misc/table
 - /minicern (still to be ported on Windows and Solaris)
 - removal of references to cernlib/shift from configure

11:30 → 11:50 5.19/04 dev release scheduled for May 7

- expect new rootfit/rootstats package from Kyle Cranmer and Wouter Verkerke
- material in branches must be moved before middle of next week
- developments by Ilka/Roj will be introduced after the release
- fixes to get "make static" and the code checker working again
- THtml must be modified to support the new dir structure (urgent)

Object Serialization in C++ with ROOT

Obvious idea now, but an important shift was moving from **centralized model**:

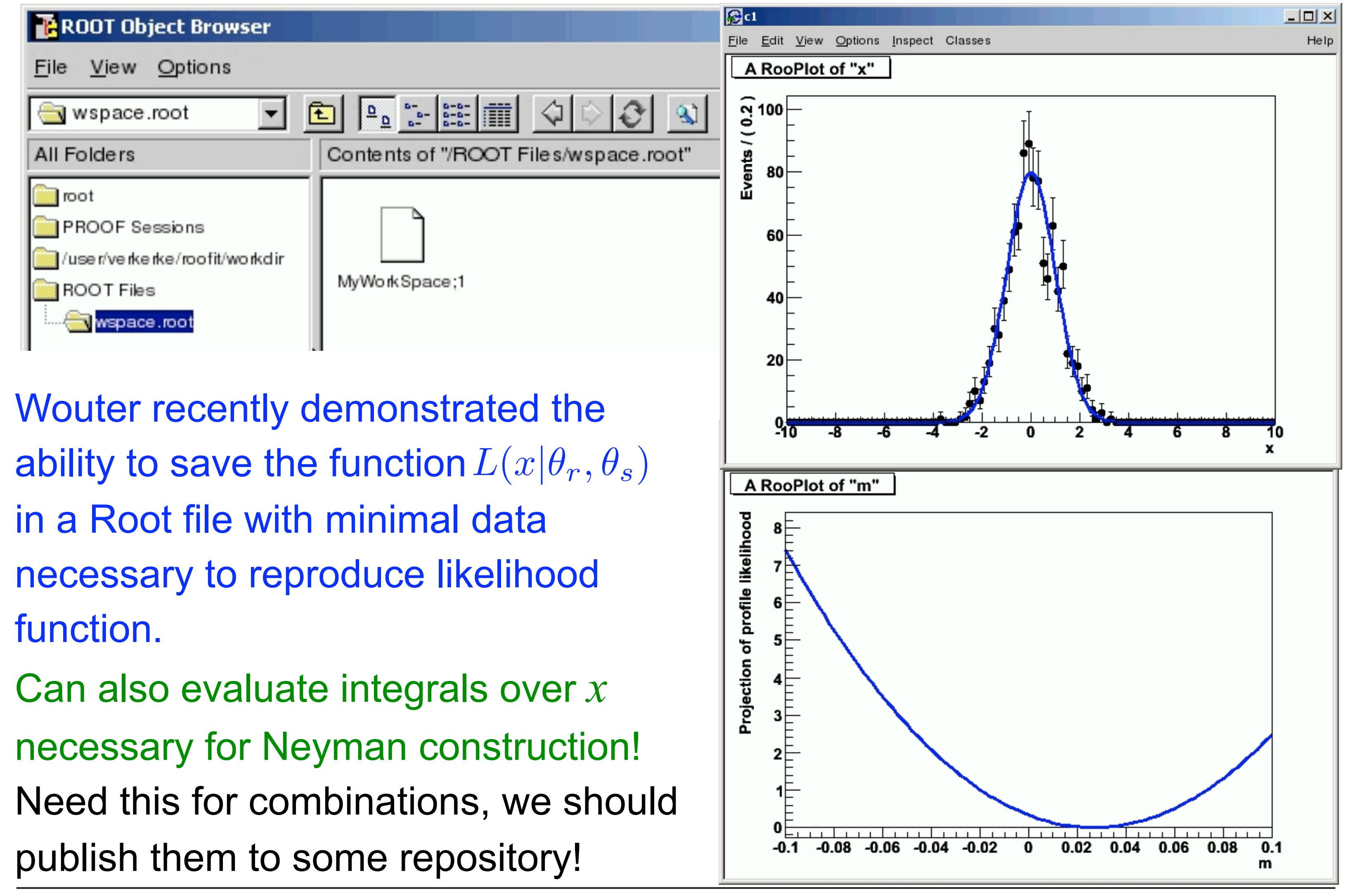
- scripts that create everything from scratch at run-time

to **distributed model**:

- Different groups **sharing digital artifacts**

At the time serializing complicated data structures composed of C++ objects was non-trivial

Example of Digital Publishing



Wouter recently demonstrated the ability to save the function $L(x|\theta_r, \theta_s)$ in a Root file with minimal data necessary to reproduce likelihood function.

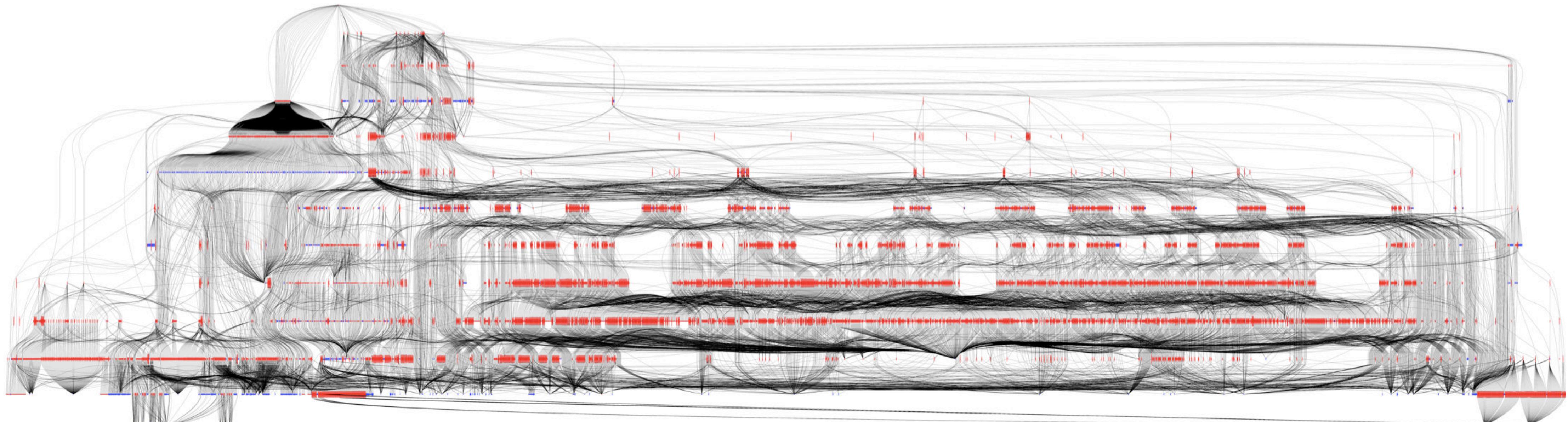
Can also evaluate integrals over x necessary for Neyman construction! Need this for combinations, we should publish them to some repository!

Combining evidence from multiple datasets

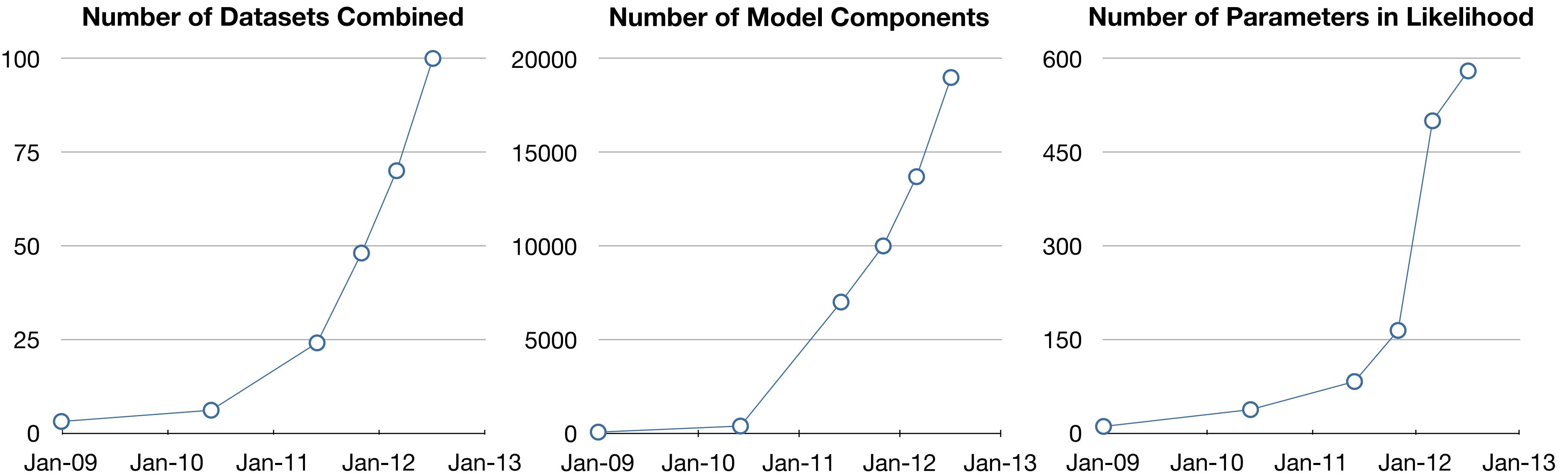
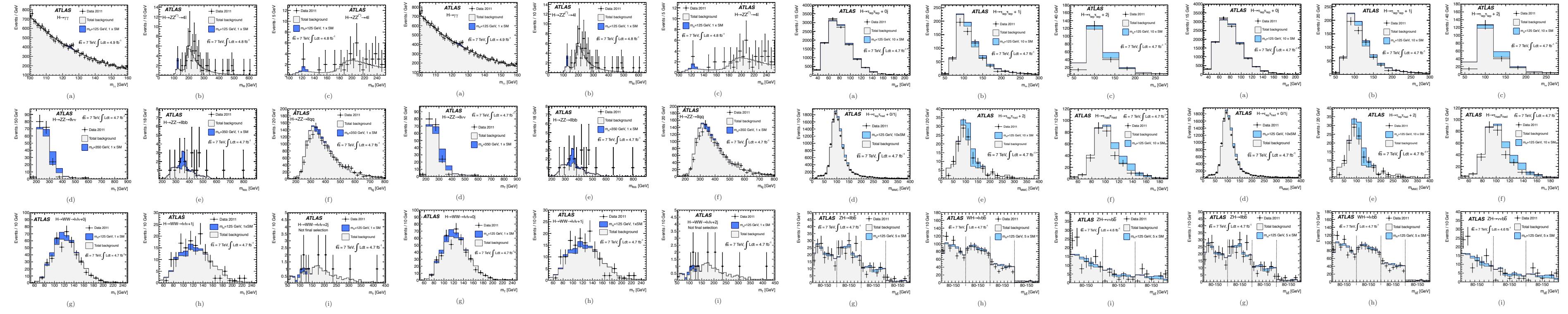
We developed a formalism, **a specification**, and a software implementation that allowed us **to combine** (statistically) **evidence** from **multiple datasets**

- Provides enormous flexibility to **remix** the evidence (different map theories to same signature)

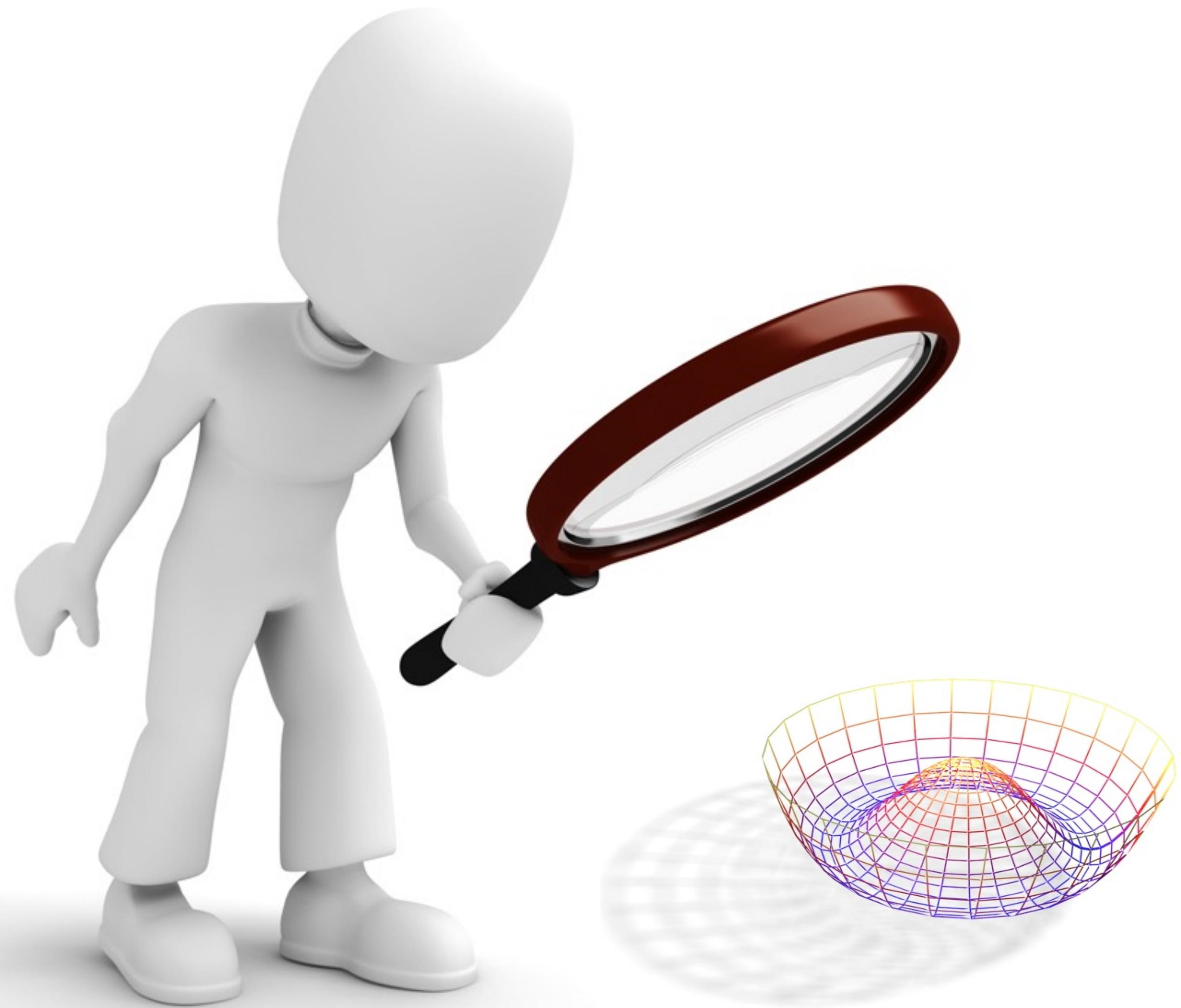
$$f_{\text{tot}}(\mathcal{D}_{\text{sim}}, \mathcal{G} | \alpha) = \prod_{c \in \text{channels}} \left[\text{Pois}(n_c | \nu_c(\alpha)) \prod_{e=1}^{n_c} f_c(x_{ce} | \alpha) \right] \cdot \prod_{p \in \mathbb{S}} f_p(a_p | \alpha_p)$$



Collaborative Statistical Modeling

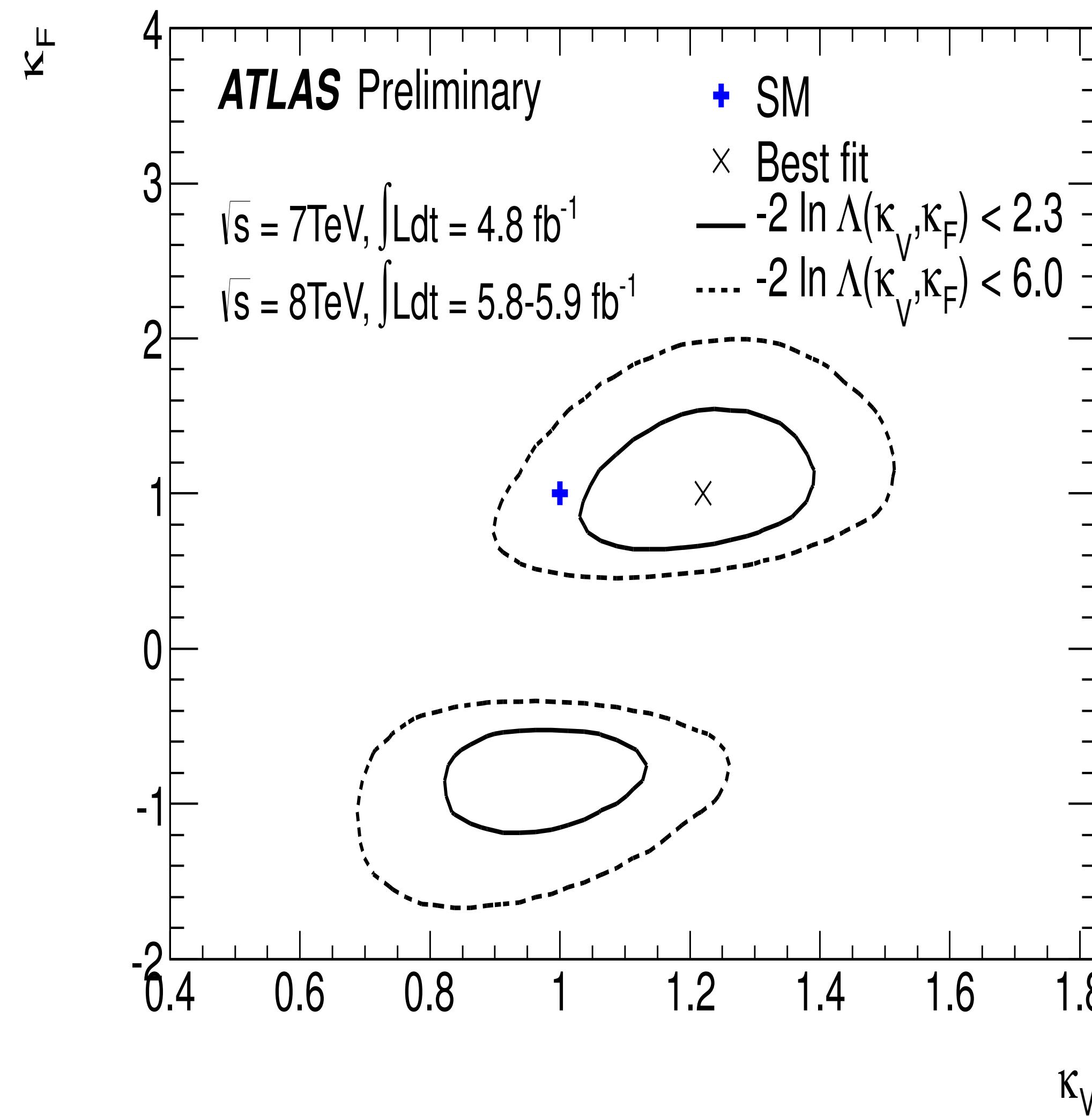






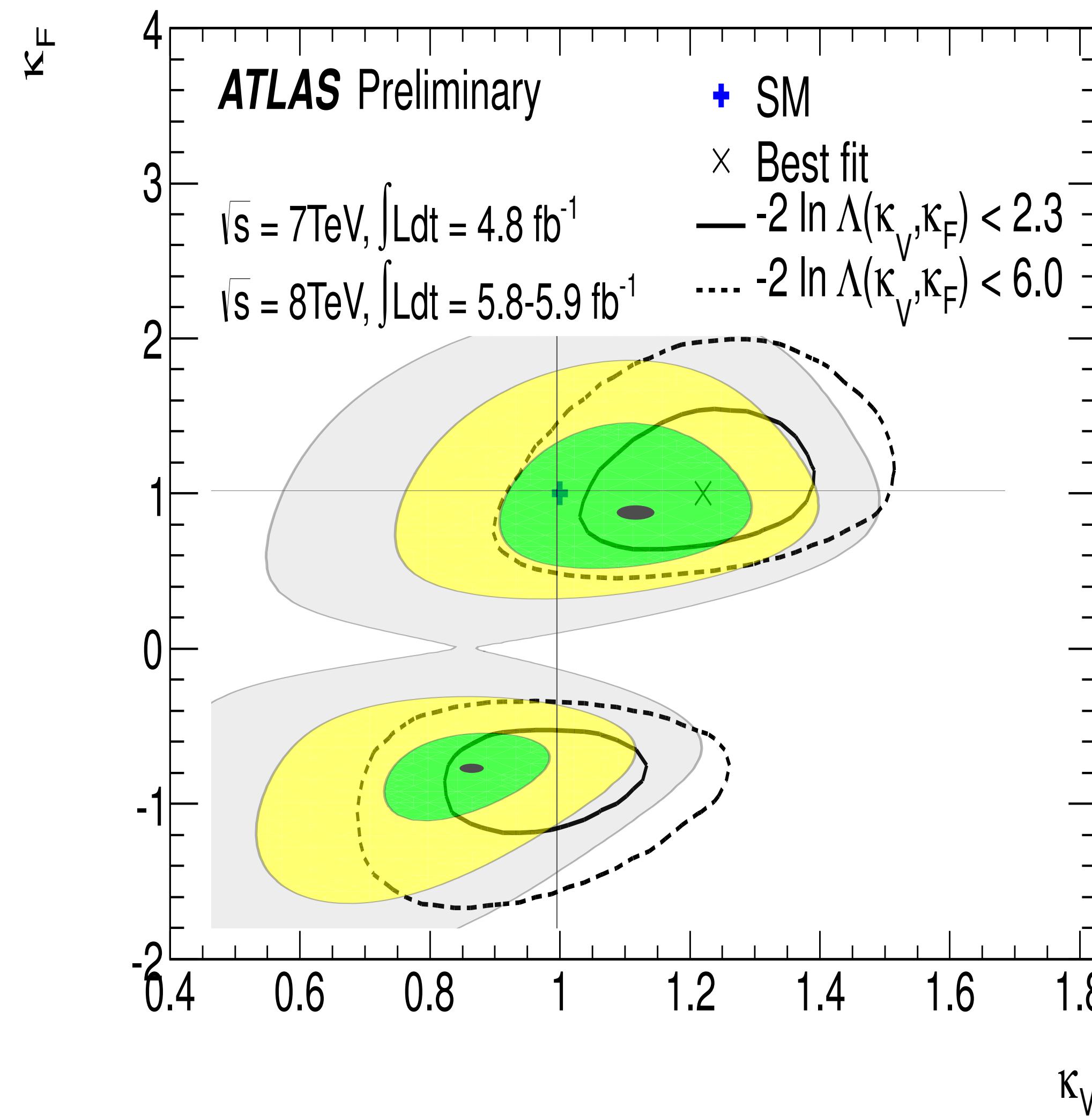
Reproducibility problem

Not possible for others to reproduce results from paper.



Reproducibility problem

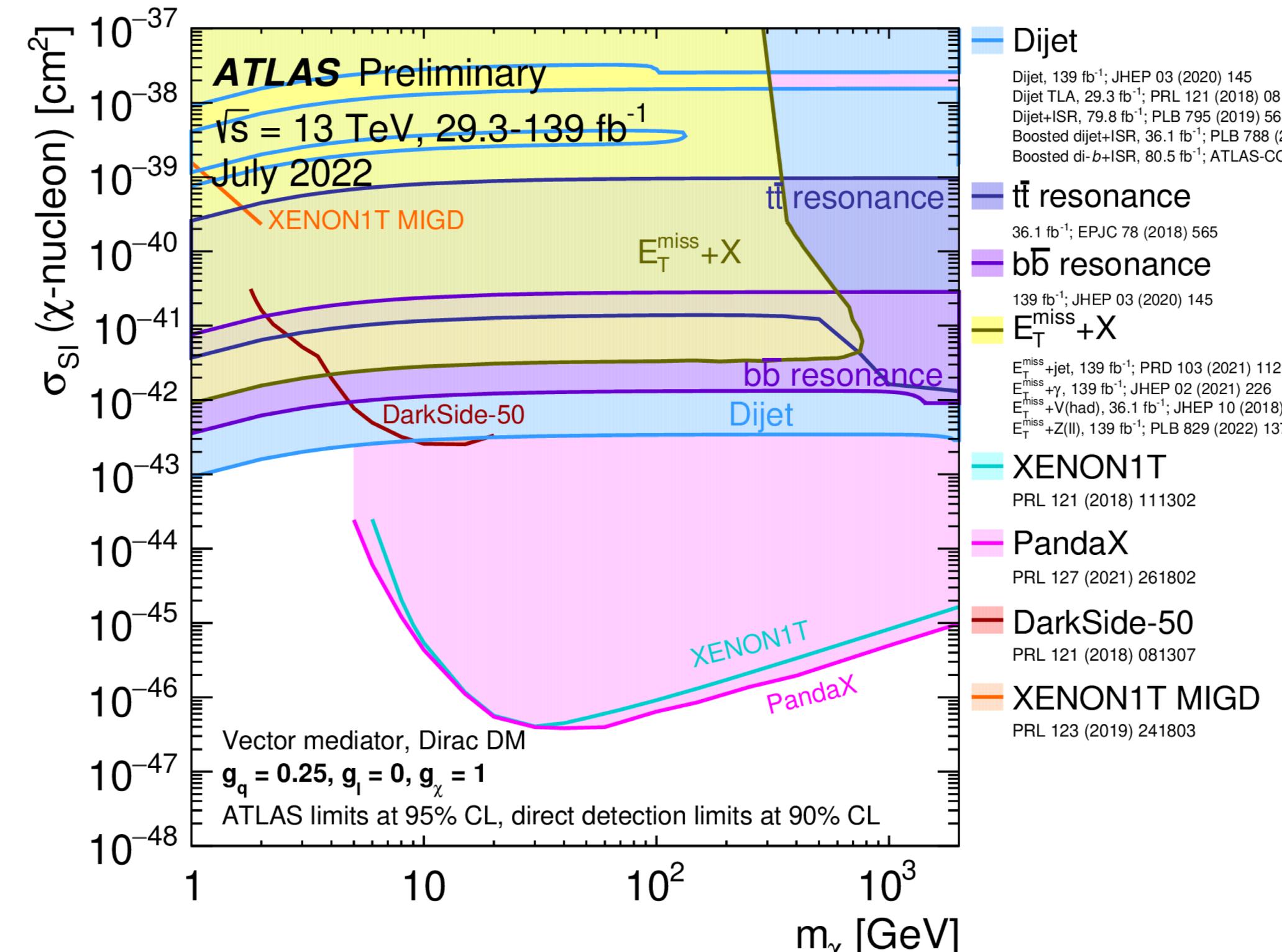
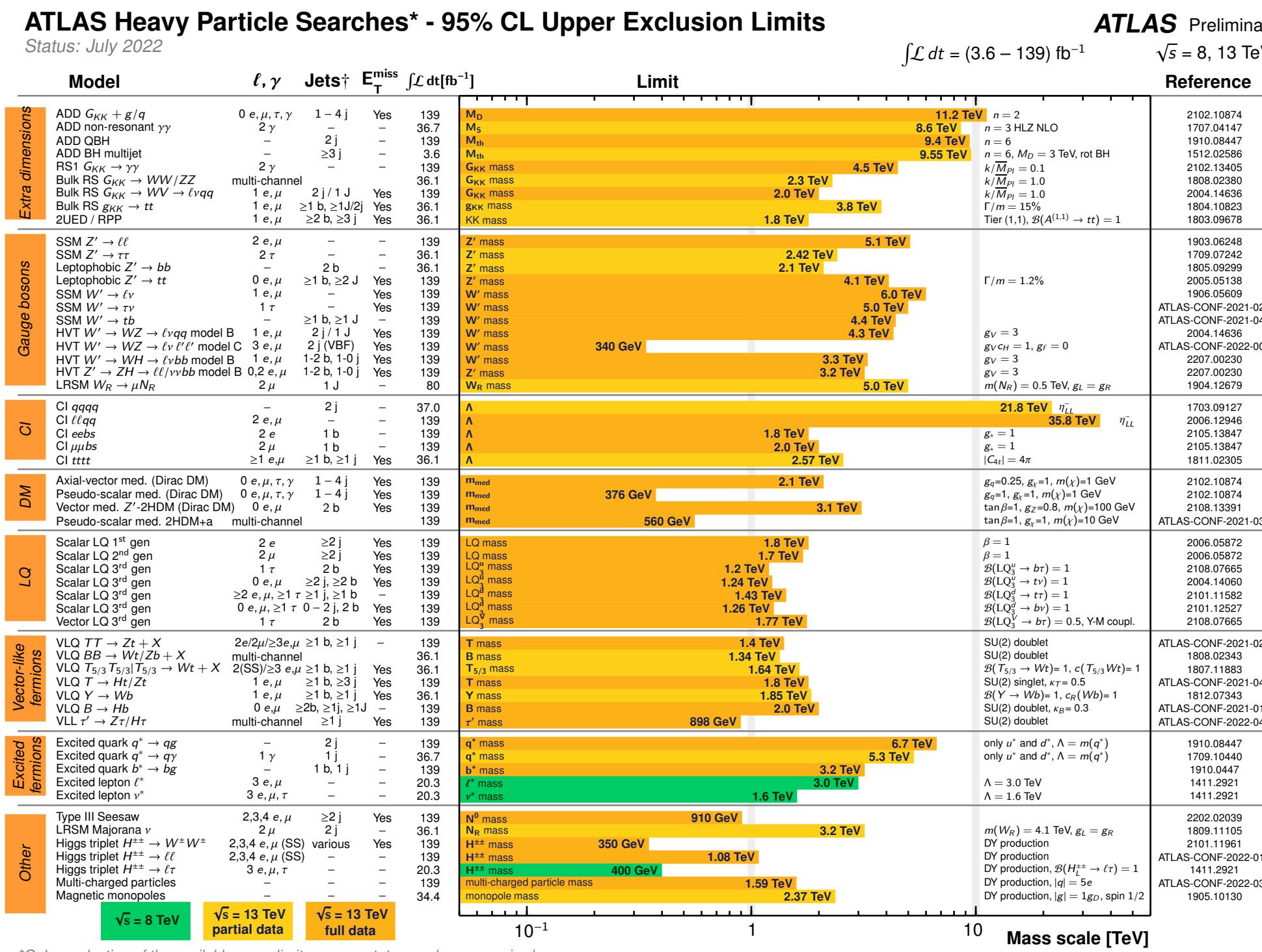
Not possible for others to reproduce results from paper.



Null results and reinterpretation

Usually, we **do not see** the signature that would be evidence for a hypothesis

- Is that because we weren't sensitive to it, or because it just wasn't there?
- Huge unmet need for reinterpretation — many questions went unanswered



Reinterpretation as a Service

$\mathcal{L}_{SM} =$

$$\frac{1}{4}\mathbf{W}_{\mu\nu} \cdot \mathbf{W}^{\mu\nu} - \frac{1}{4}B_{\mu\nu}B^{\mu\nu} - \frac{1}{4}G_a^a G_a^{\mu\nu}$$

kinetic energies and self-interactions of the gauge bosons

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

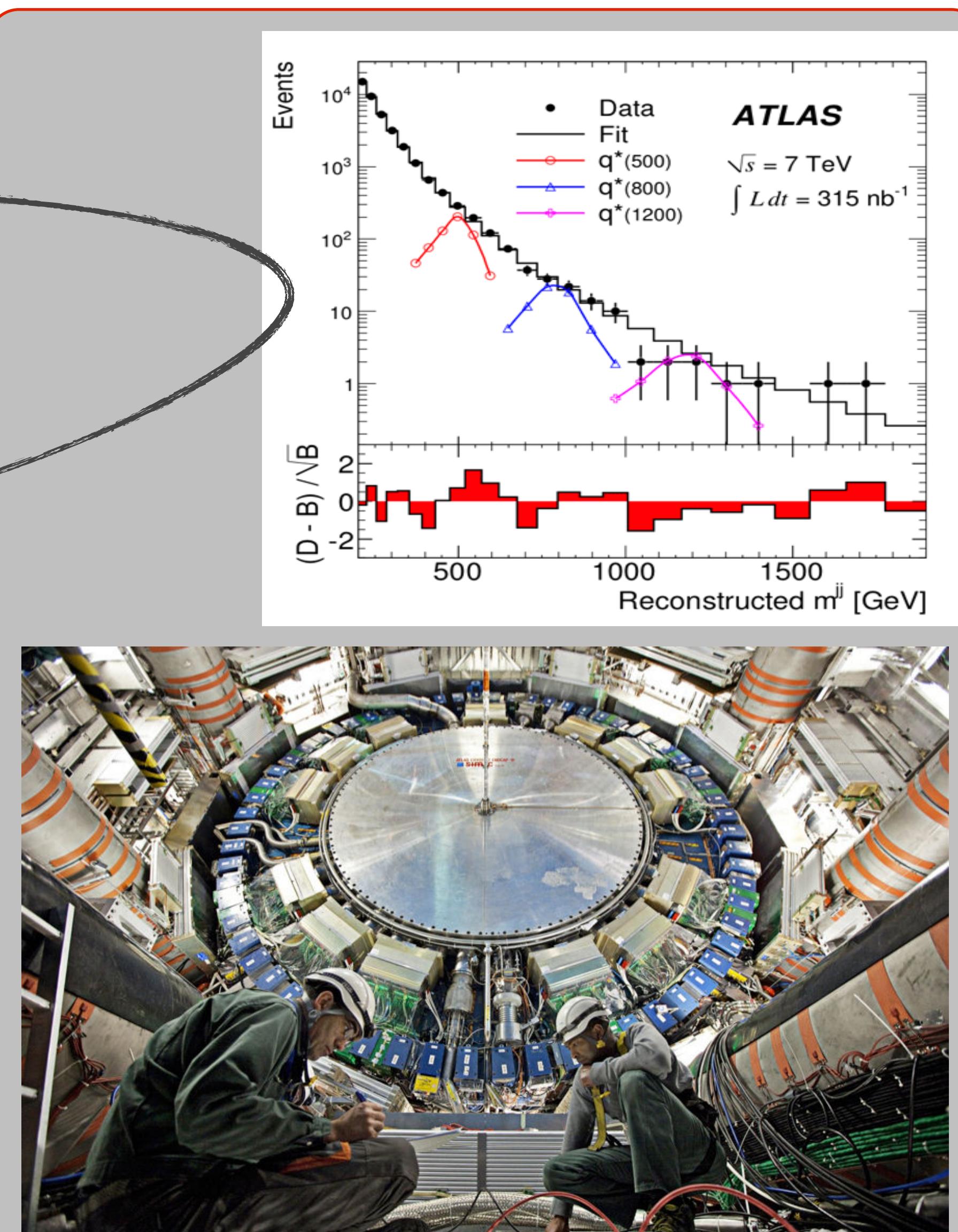
$$+ \underbrace{\frac{1}{2}|(i\partial_\mu - \frac{1}{2}g\tau \cdot \mathbf{W}_\mu - \frac{1}{2}g'YB_\mu)\phi|^2 - V(\phi)}_{W^\pm, Z, \gamma, \text{and Higgs masses and couplings}}$$

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

Q



A



ATLAS

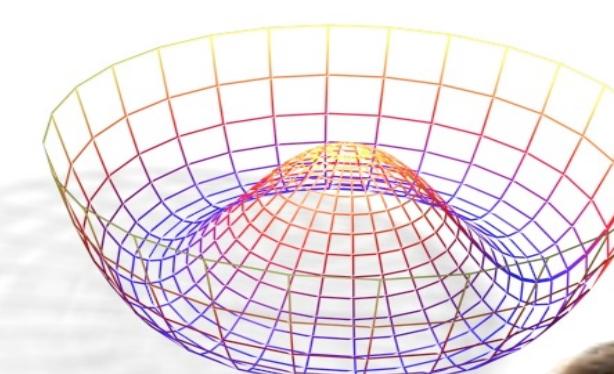
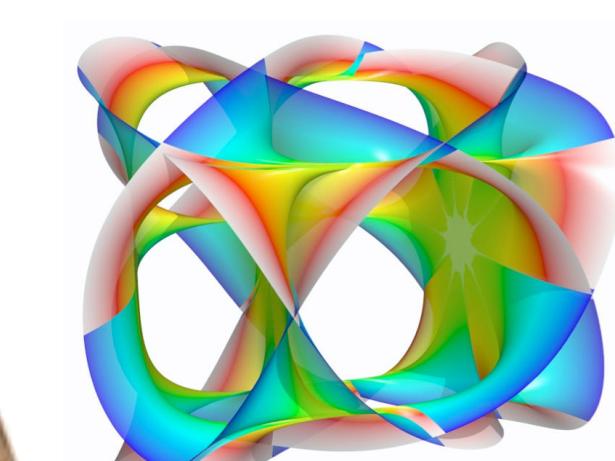
Events

• Data
— Fit
—○— $q^*(500)$
—△— $q^*(800)$
—◆— $q^*(1200)$

$\sqrt{s} = 7 \text{ TeV}$
 $\int L dt = 315 \text{ nb}^{-1}$

(D - B) / \sqrt{B}

Reconstructed m^{jj} [GeV]



<https://doi.org/10.1007/JHEP04%282011%29038>

Where were we 10 years ago?

Incentives

- **Bulk data processing** was respected as essential, given resources and incentives, well managed with dedicated professionals
- **End-user analysis**, in contrast, was considered ‘the fun part’ and not given resources or incentives
 - Neglected the need for high-quality, specialized tools used by many
 - Wild West. Lack of harmonization led to inefficient process
 - Analysis code rarely preserved, results were difficult to reproduce or reinterpret

ROOT (a general-purpose C++ based framework for analysis) was used by almost everyone

- Open source, but had a small development team with few contributions from community

Overall our toolkit was not aligned with industry, data science, or other areas of science

- Better alignment ⇒ better leverage of tools in the data science ecosystem, better competitiveness for PhD’s entering job market, better appreciation of HEP’s value to society

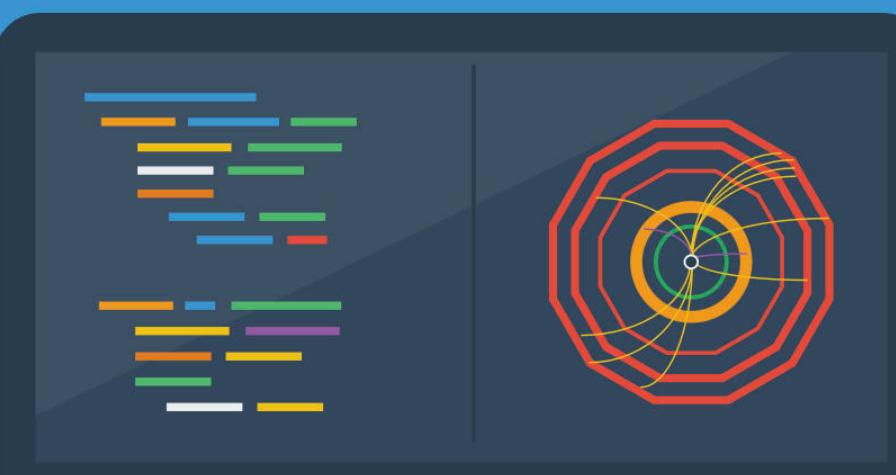
A new approach

In 2014 we were awarded NSF funding for a new software effort aimed at an upgrade of the LHC with an approach that was more aligned with python & data science

 dianahep Team Activities/Products DIANA Fellows Blog

Advanced software plays a fundamental role in large scientific projects.

The primary goal of DIANA/HEP is to develop state-of-the-art software tools for experiments which acquire, reduce, and analyze petabytes of data. Improving performance, interoperability, and collaborative tools through modifications and additions to ROOT and other software packages broadly used by the community will allow users to more fully exploit the data being acquired at CERN's Large Hadron Collider (LHC) and other facilities. As part of the NSF's Software Infrastructure for Sustained Innovation (SI2) program, DIANA is concerned with the overarching goal of transforming innovations in research and education into sustained software resources that are an integral part of the cyberinfrastructure.



Bogdan Mihaila



Bogdan Mihaila

 iris hep About Connect Activities Fellows Jobs

Institute for Research and Innovation in Software for High Energy Physics (IRIS-HEP)

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:



Upcoming Events:

Jul 8–14, 2024	Tacoma, Washington
Scientific Computing with Python (SciPy) 2024	
Jul 18–19, 2024	University of Washington
USATLAS/IRIS-HEP Software Training	
Jul 21–26, 2024	Fermilab
Coding Camp - Fermilab	
Jul 22–26, 2024	Princeton University
CoDaS-HEP 2024 - Computational and Data Science Training for High Energy Physics	
Aug 26–30, 2024	Aachen, Germany
PyHEP.dev 2024 - "Python in HEP" Developer's Workshop	
Sep 4–6, 2024	University of Washington
IRIS-HEP Institute Retreat	
Sep 23–25, 2024	Valencia (Spain)
Fourth MODE Workshop on Differentiable Programming for Experiment Design	

[View all past events](#)

IRIS-HEP Receives \$25M Funding for Another Five Years of Research

"IRIS-HEP received funding from the Office of

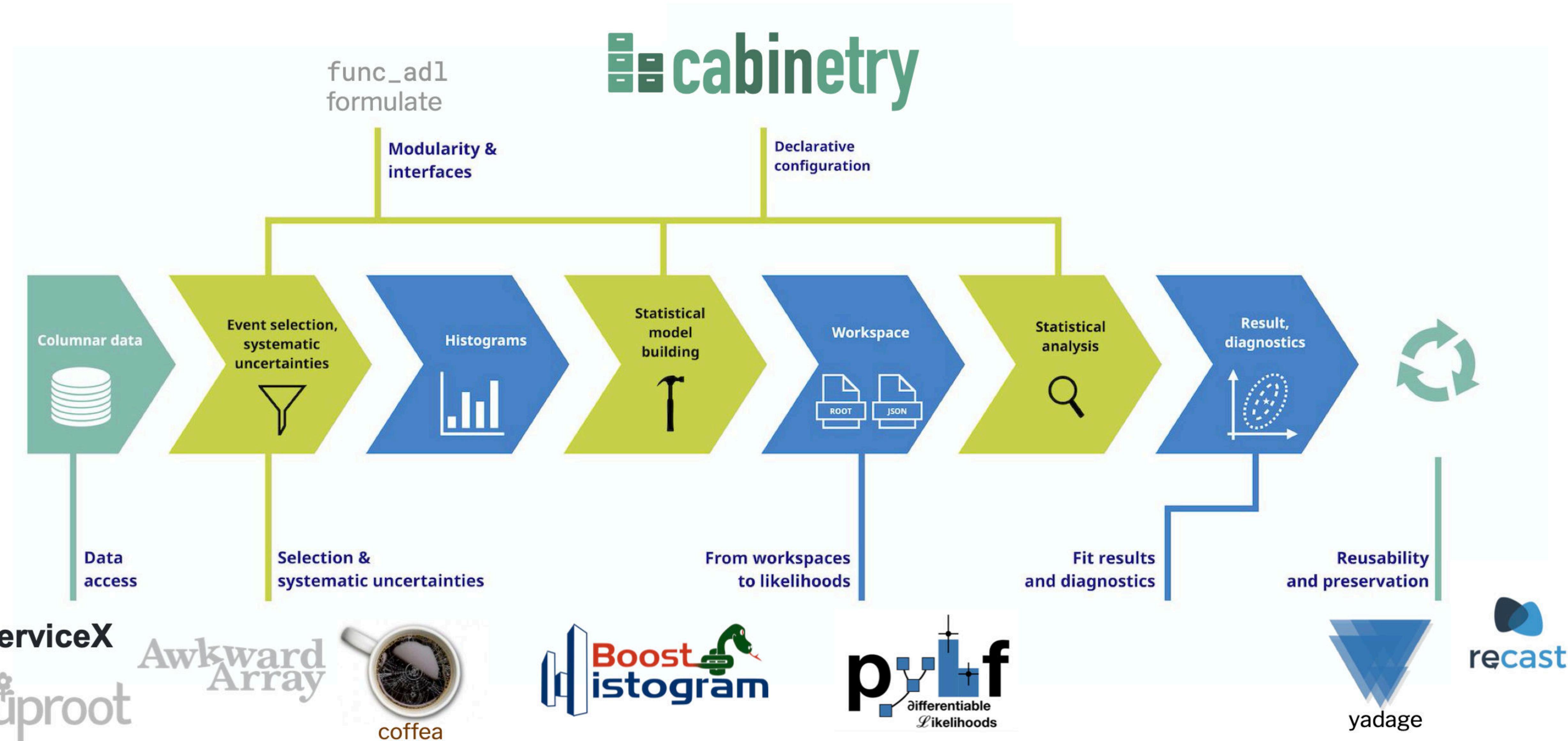
Out of harm's way: Physics research program supports Ukrainian students displaced by war

"Ukrainia students escape the war and pursue

2014 - 2018

2018 - 2029

An ecosystem of analysis tools



A thriving open source community



Search Scikit-HEP

Scikit-HEP on GitHub

Packages

You can see our [affiliated projects](#), get links to documentation, see the [guidelines](#) for a new Scikit-HEP package to follow, and see our [statement on Python version support](#) here too.

Basics:

Awkward Array Manipulate JSON-like data with NumPy-like idioms.
conda-forge v2.6.5 pypi v2.6.6 wheel yes docs available

hepunits Units and constants in the HEP system of units.
conda-forge v2.3.4 pypi v2.3.4 wheel yes

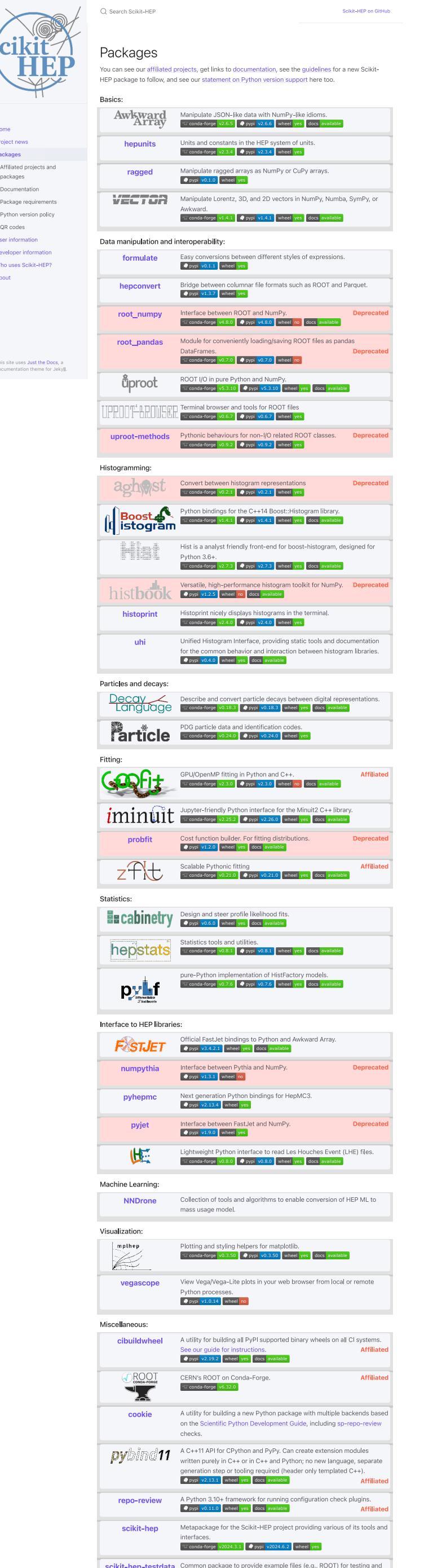
ragged Manipulate ragged arrays as NumPy or CuPy arrays.
pypi v0.1.0 wheel yes

VECTOR Manipulate Lorentz, 3D, and 2D vectors in NumPy, Numba, SymPy, or Awkward.
conda-forge v1.4.1 pypi v1.4.1 wheel yes docs available

Data manipulation and interoperability:

formulate Easy conversions between different styles of expressions.
pypi v0.1.1 wheel yes

hepconvert Bridge between columnar file formats such as ROOT and Parquet.
pypi v1.3.7 wheel yes



Packages

You can see our affiliated projects, get links to documentation, see the guidelines for a new Scikit-HEP package to follow, and see our statement on Python version support here too.

Basics:

- Awkward Array** Manipulate JSON-like data with NumPy-like idioms.
conda-forge v2.6.5 pypi v2.6.6 wheel yes docs available
- hepunits** Units and constants in the HEP system of units.
conda-forge v2.3.4 pypi v2.3.4 wheel yes docs available
- ragged** Manipulate ragged arrays as NumPy or CuPy arrays.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available
- VECTOR** Manipulate Lorentz, 3D, and 2D vectors in NumPy, Numba, SymPy, or Awkward.
conda-forge v1.4.1 pypi v1.4.1 wheel yes docs available

Data manipulation and interoperability:

- formulate** Easy conversions between different styles of expressions.
conda-forge v0.1.1 pypi v0.1.1 wheel yes docs available
- hepconvert** Bridge between columnar file formats such as ROOT and Parquet.
conda-forge v1.3.7 pypi v1.3.7 wheel yes docs available
- root_numpy** Interface between ROOT and NumPy.
conda-forge v0.5.0 pypi v0.5.0 wheel yes docs available
- root_pandas** Module for conveniently loading/saving ROOT files as pandas DataFrames.
conda-forge v0.7.0 pypi v0.7.0 wheel yes docs available
- uproot** ROOT I/O in pure Python and NumPy.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available
- uproot-avro** Terminal browser and tools for ROOT files.
conda-forge v0.0.1 pypi v0.0.1 wheel yes docs available
- uproot-methods** Pythonic behaviours for non-I/O related ROOT classes.
conda-forge v0.0.2 pypi v0.0.2 wheel yes docs available

Histogramming:

- aghost** Convert between histogram representations.
conda-forge v0.2.0 pypi v0.2.0 wheel yes docs available
- BoostHistogram** Python bindings for the C++14 Boost::Histogram library.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available
- Hist** Hist is an analyst friendly front-end for boost-histogram, designed for Python 3.6+.
conda-forge v0.2.0 pypi v0.2.0 wheel yes docs available
- histbook** Versatile, high-performance histogram toolkit for NumPy.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available
- histprint** Histogram nicely displays histograms in the terminal.
conda-forge v0.0.1 pypi v0.0.1 wheel yes docs available
- uhf** Unified Histogram Interface, providing static tools and documentation for the common behavior and interaction between histogram libraries.
conda-forge v0.0.1 pypi v0.0.1 wheel yes docs available

Particles and decays:

- DecayLanguage** Describe and convert particle decays between digital representations.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available
- Particle** PDG particle data and identification codes.
conda-forge v0.2.0 pypi v0.2.0 wheel yes docs available

Fitting:

- coffit** GPU/OpenMP fitting in Python and C++.
conda-forge v2.0.0 pypi v2.0.0 wheel yes docs available
- iminuit** Jupyter-friendly Python interface for the Minuit2 C++ library.
conda-forge v2.2.2 pypi v2.2.2 wheel yes docs available
- problif** Cost function builder. For fitting distributions.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available
- zfit** Scalable Pythonic fitting.
conda-forge v0.2.0 pypi v0.2.0 wheel yes docs available

Statistics:

- cabinetry** Design and steer profile likelihood fits.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available
- hepstats** Statistics tools and utilities.
conda-forge v0.0.1 pypi v0.0.1 wheel yes docs available
- pystl** pure-Python implementation of HistFactory models.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available

Interface to HEP libraries:

- FastJet** Official FastJet bindings to Python and Awkward Array.
conda-forge v0.4.2 pypi v0.4.2 wheel yes docs available
- numhepmath** Interface between Python and NumPy.
conda-forge v1.3.3 pypi v1.3.3 wheel yes docs available
- pyhepmc** Next generation Python bindings for HepMC3.
conda-forge v1.2.1 pypi v1.2.1 wheel yes docs available
- pyjet** Interface between FastJet and NumPy.
conda-forge v1.9.0 pypi v1.9.0 wheel yes docs available
- HL** Lightweight Python interface to read Les Houches Event (LHE) files.
conda-forge v0.0.2 pypi v0.0.2 wheel yes docs available

Machine Learning:

- NNDrone** Collection of tools and algorithms to enable conversion of HEP ML to mass usage model.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available

Visualization:

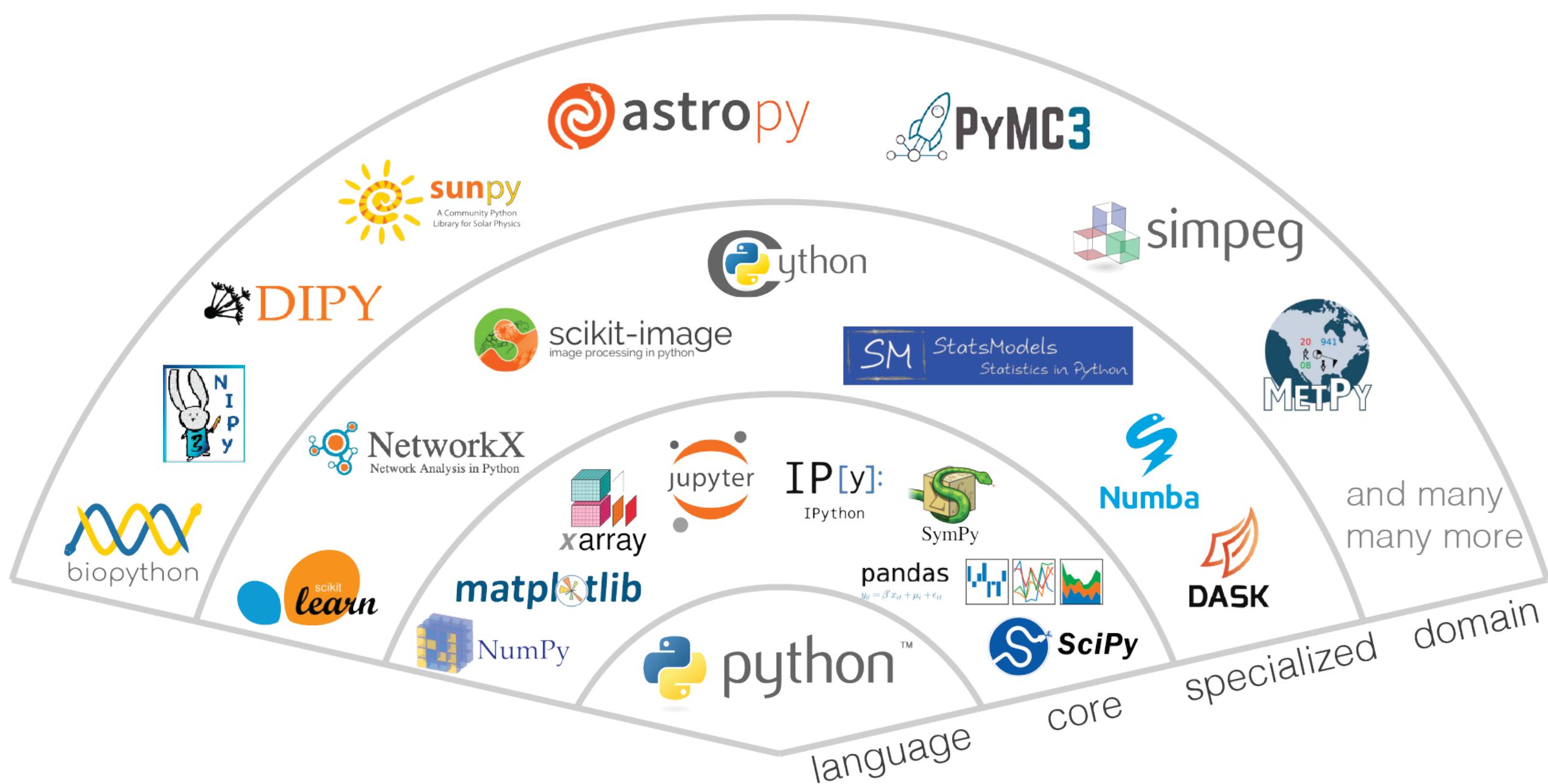
- matplotlib** Plotting and styling helpers for matplotlib.
conda-forge v3.0.2 pypi v3.0.2 wheel yes docs available
- vegascope** View Vega/Vega-Lite plots in your web browser from local or remote Python processes.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available

Miscellaneous:

- cibuildwheel** A utility for building all PyPI supported binary wheels on all CI systems. See our guide for instructions.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available
- ROOT** CERN's ROOT on Conda-Forge.
conda-forge v5.32.0 pypi v5.32.0 wheel yes docs available
- cookie** A utility for building a new Python package with multiple backends based on the Scientific Python Development Guide, including so-repo-review checks.
conda-forge v0.1.0 pypi v0.1.0 wheel yes docs available
- pybind11** A C++11 API for CPython and PyPy. Can create extension modules written purely in C++ or in C++ and Python; no new language, separate generation step or tooling required (header only templated C++).
conda-forge v2020.5 pypi v2020.5 wheel yes docs available
- repo-review** A Python 3.10+ framework for running configuration check plugins.
conda-forge v0.1.1 pypi v0.1.1 wheel yes docs available
- scikit-hep** Metapackage for the Scikit-HEP project, providing various of its tools and interfaces.
conda-forge v0.0.1 pypi v0.0.1 wheel yes docs available
- scikit-hep-testdata** Common package to provide example files (e.g., ROOT) for testing and developing packages against.
conda-forge v0.0.5 pypi v0.0.5 wheel yes docs available

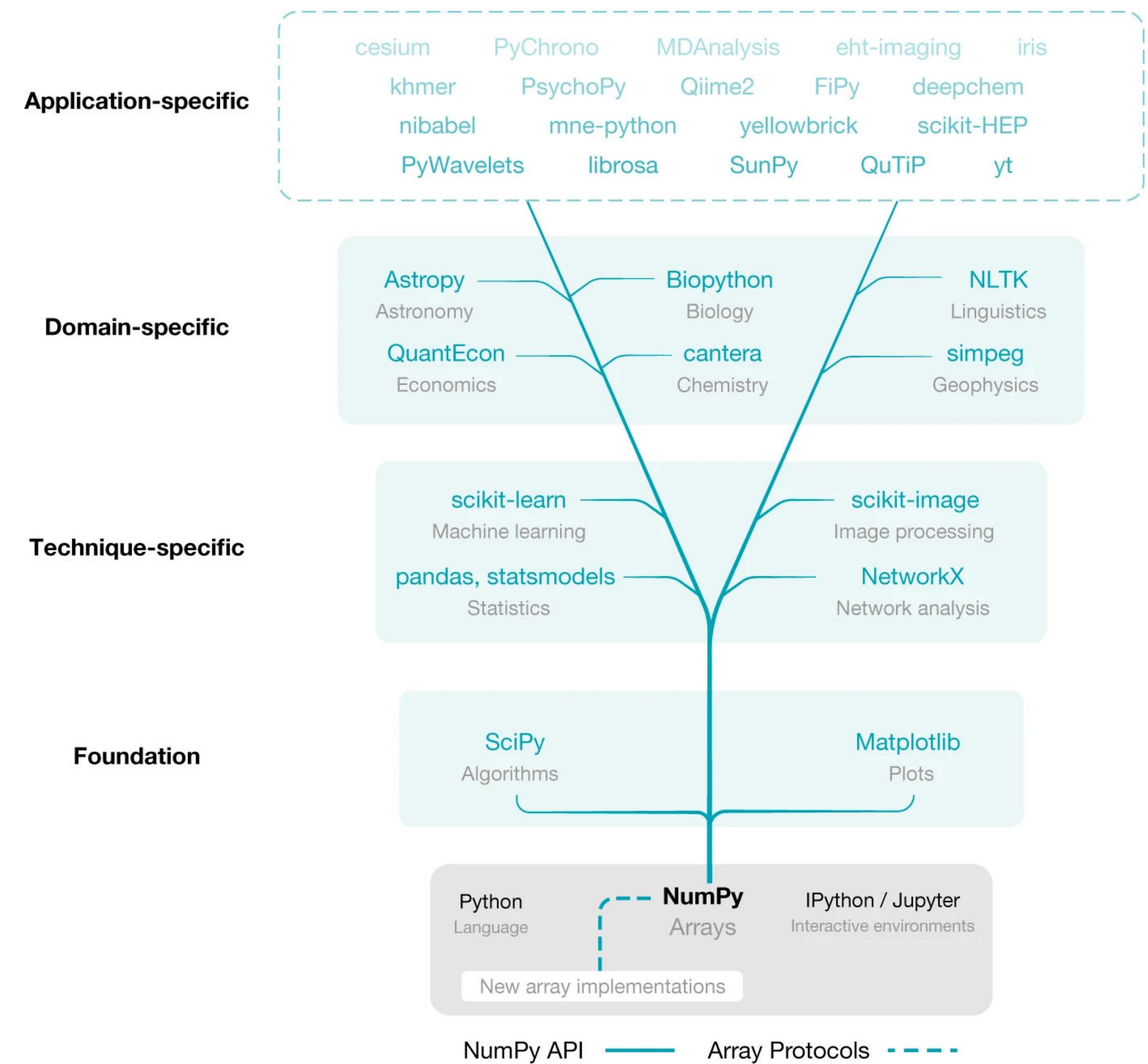
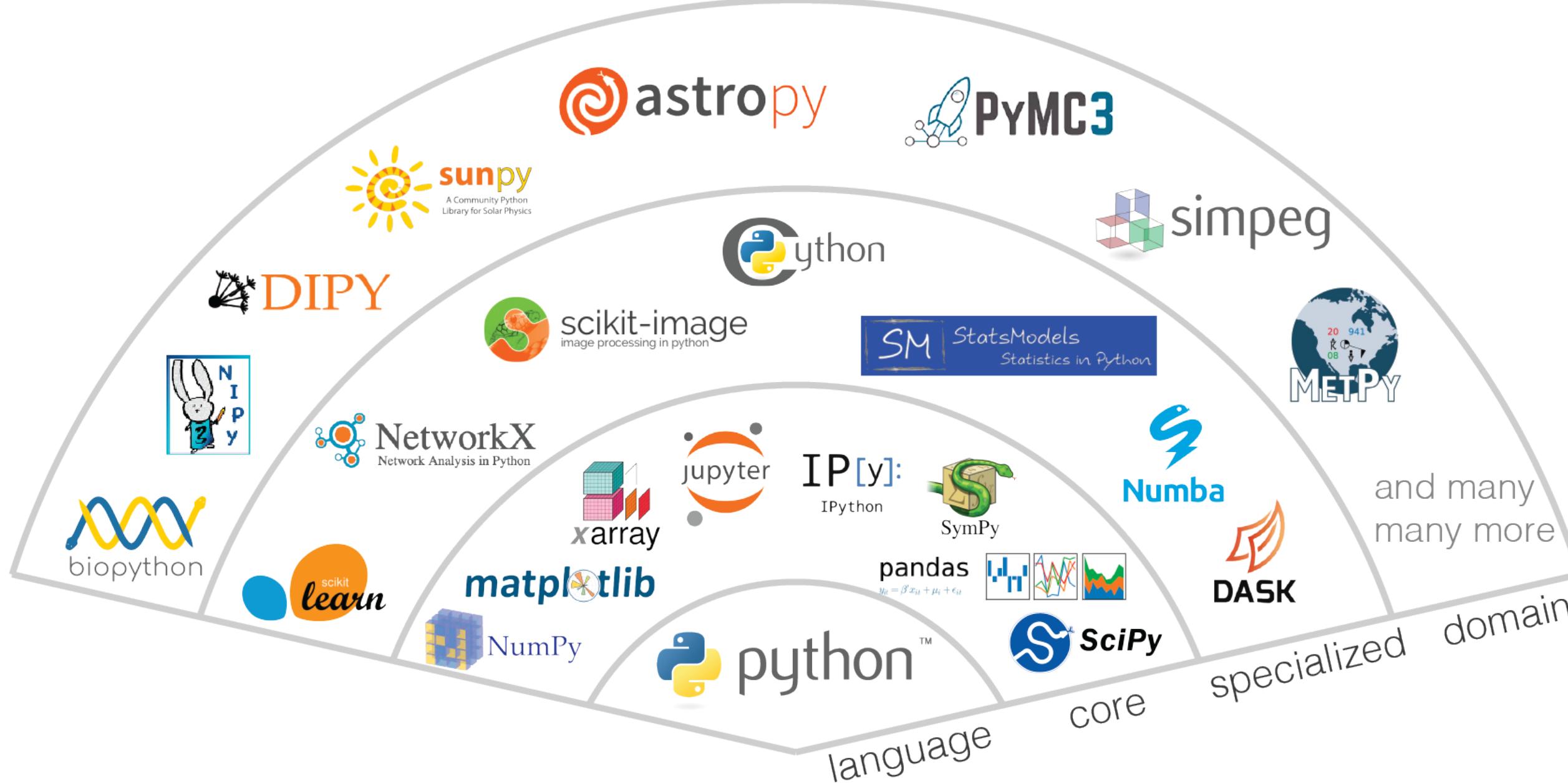
Recognition within broader scientific community

Until recently, HEP was still not well integrated or represented in broader scientific python/software community



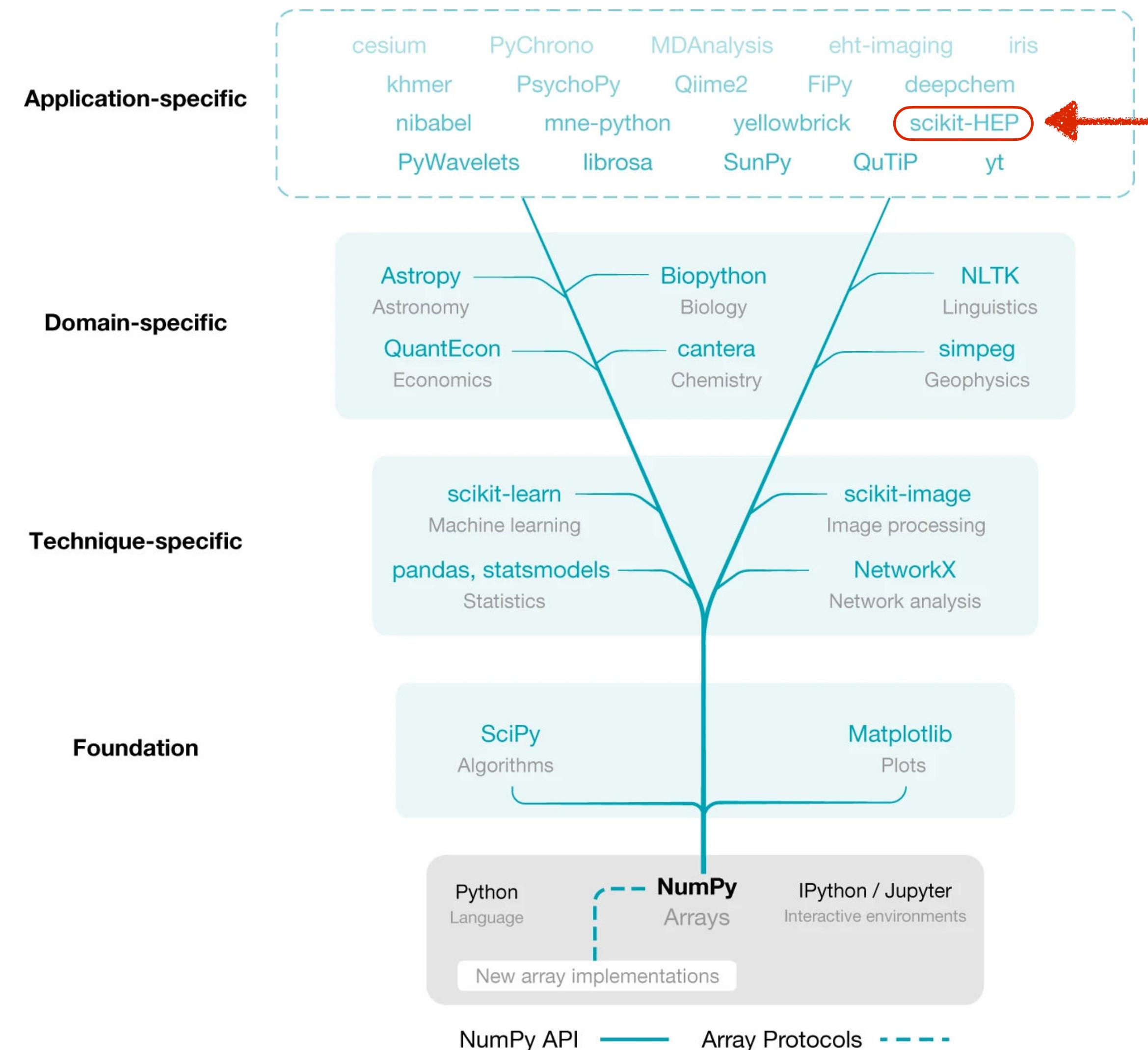
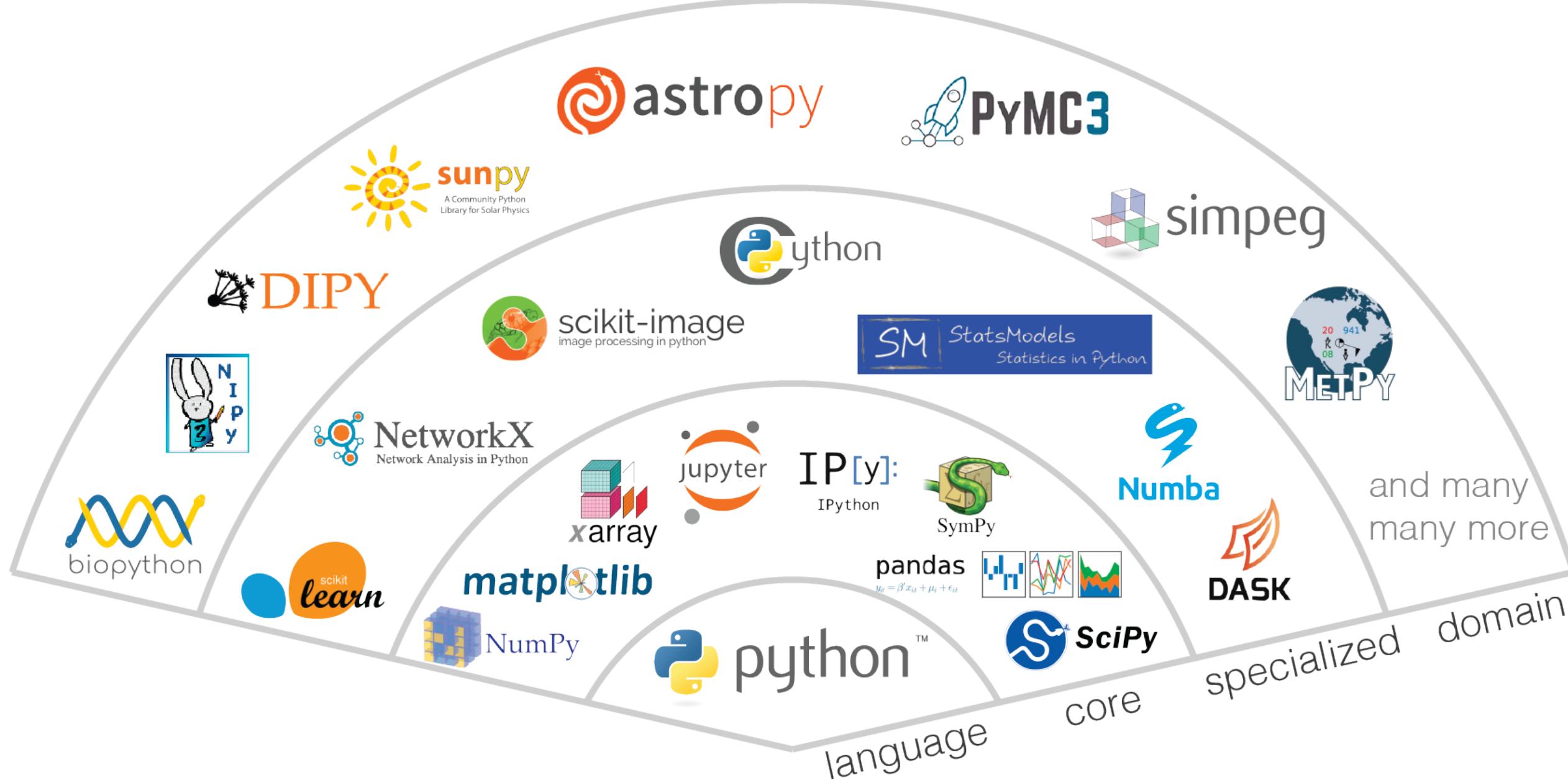
Recognition within broader scientific community

Until recently, HEP was still not well integrated or represented in broader scientific python/software community



Recognition within broader scientific community

Until recently, HEP was still not well integrated or represented in broader scientific python/software community



Some Highlights

The power of a good specification

The original C++ version of the HistFactory tool distributed in ROOT used a configuration file with clear semantics

- The statistical model was very well defined mathematically

Later recognized this as embracing a **declarative specification**.

- Facilitated a pythonic implementation

HistFactory

HistFactory tool that ships with ROOT targeting binned analyses

- XML files organize the histograms
- conventions define model exactly
 - [CERN-OPEN-2012-016](#)
- command line tool creates likelihood

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

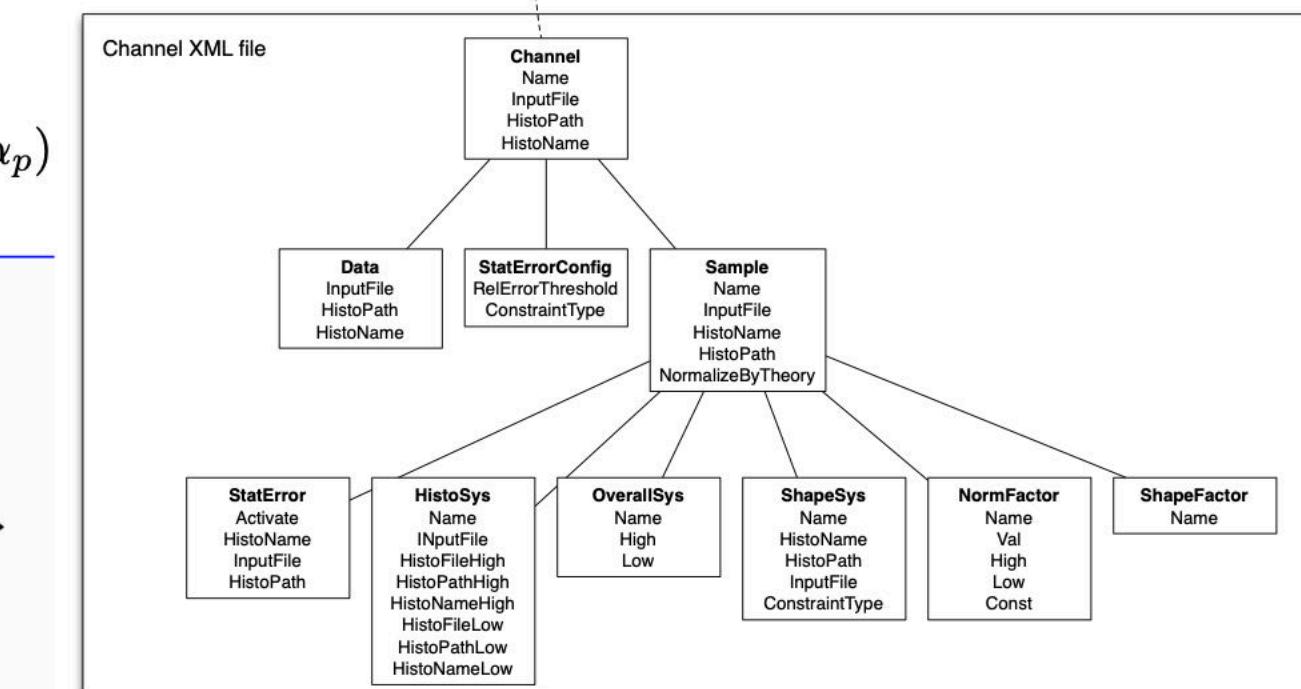
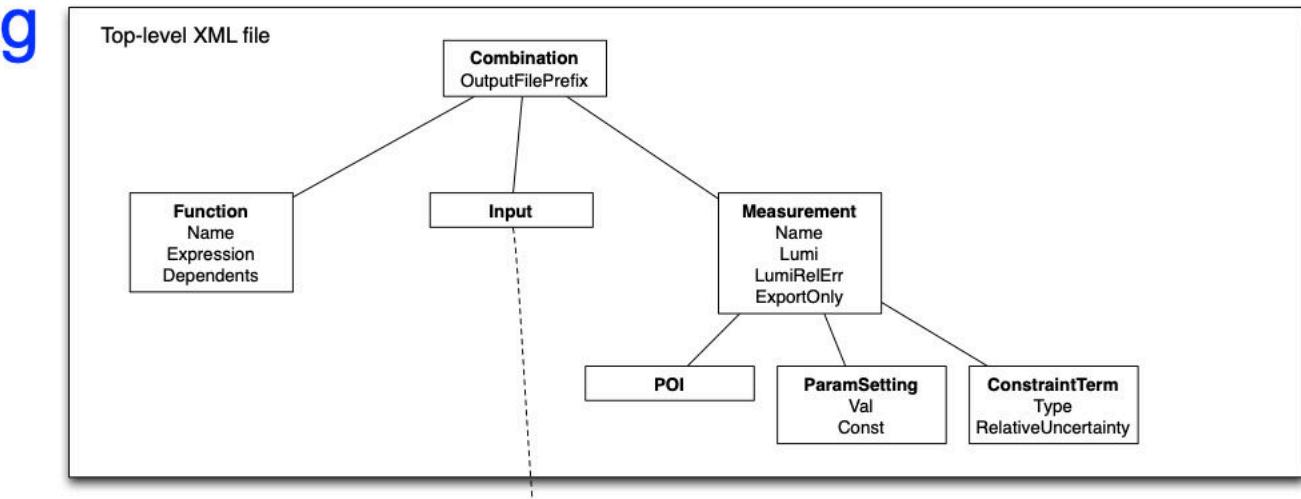
```
<!DOCTYPE Channel SYSTEM 'HistFactorySchema.dtd'>
<Channel Name="A" InputFile="./data/ABCD.root" >
  <Data HistoName="A_data" HistoPath="" />

  <!-- This is the signal (eg. mu)-->
  <Sample Name="A_signal" HistoPath="" HistoName="unit_histogram">
    <!-- now mu is number of events-->
    <NormFactor Name="mu" Val="1" Low="0" High="200" />
    <OverallSys Name="syst1" High="1.01" Low="0.99" />
  </Sample>

  <!-- This bkg is estimated from MC (eg. mu_A^K) -->
  <Sample Name="A_backgroundMC" HistoPath="" NormalizeByTheory="True" HistoName="unit_histogram" >
    <NormFactor Name="mu_K_A" Val="100" Low="0" High="200" />
  </Sample>

  <!-- Background 2 is completely Data-Driven -->
  <Sample Name="A_backgroundDD" HistoPath="" NormalizeByTheory="False" HistoName="unit_histogram" >
    <NormFactor Name="mu_D_U" Val="100" Low="24500" High="26000" />
    <NormFactor Name="etaB" Val="1" Low="0." High="0.02" Const="False" />
    <NormFactor Name="etaC" Val="1" Low="0." High="0.3" Const="False" />
    <!-- NormFactor and ShapeFactor same for a 1-bin histogram. But we can name NormFactor-->
  </Sample>
</Channel>
```

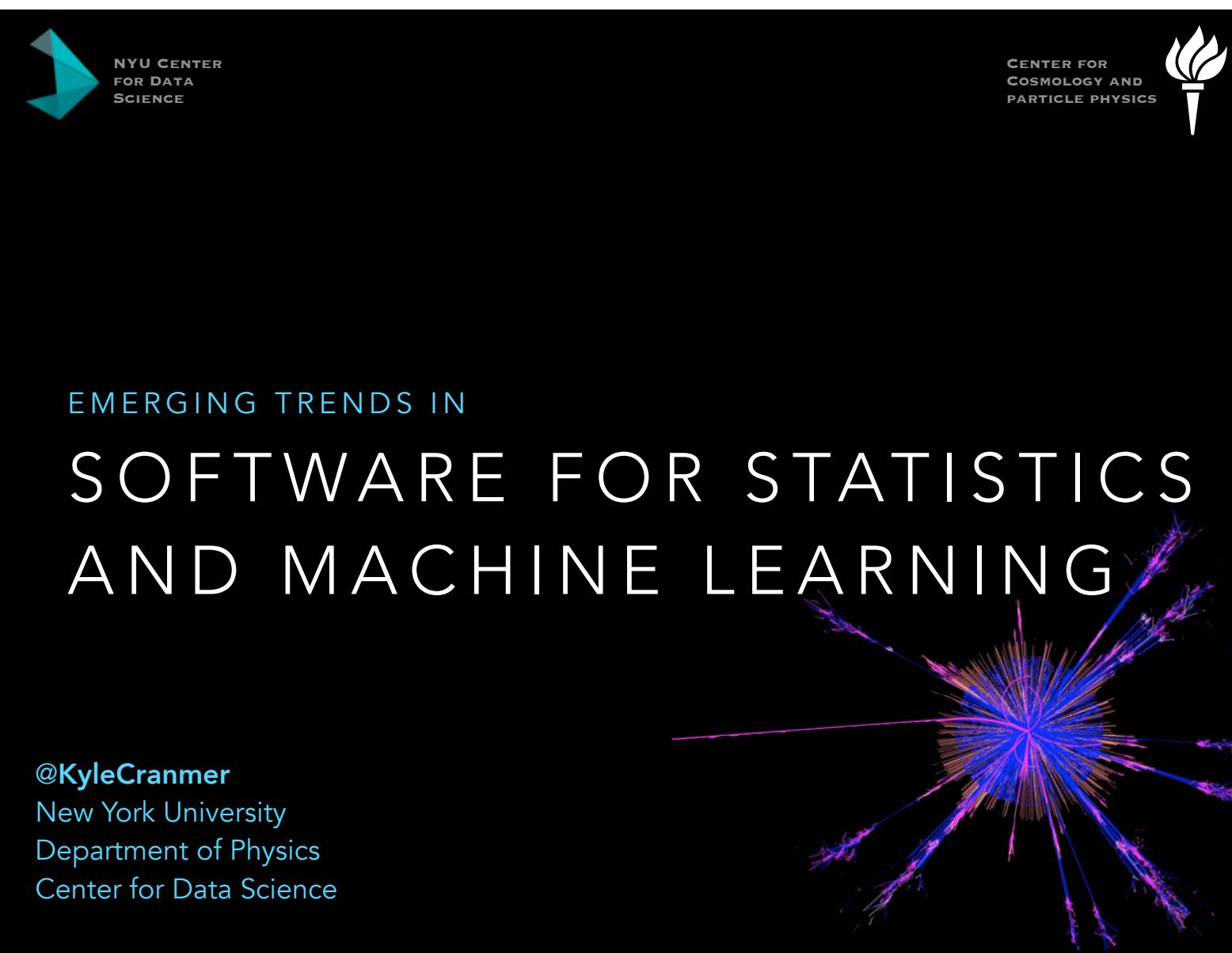
CENTER FOR
COSMOLOGY AND
PARTICLE PHYSICS



The impact of AI/ML frameworks

In 2014, I was investigating how to leverage advances in AI/ML in statistics (more on that later)

- In 2015, TensorFlow was announced, and we immediately saw advantages of reimplementing statistical tools in a AI/ML framework with automatic differentiation



2nd S2I2 HEP/CS Workshop

1 May 2017, 07:00 → 3 May 2017, 13:35 US/Eastern

Princeton University

<https://indico.cern.ch/event/622920/timetable/#sc-5-2-lightning-talk-emerging>

Probabilistic programming frameworks

RooFit serves us well, but shows limits in terms of **scalability**.

Using a data flow graph framework, RooFit would be **distributed**, GPU-enabled and automatically **differentiable**.

Feasibility? Certainly **within reach!** As illustrated by our tentative proof-of-concepts `carl.distributions` [[Gilles Louppe](#)] and `tensorprob` [[Igor Babuschkin](#), now at DeepMind]. See also `Edward`.

Edward
A library for probabilistic modeling, inference, and criticism.

Dustin Tran
Ph.D. Student
Columbia University
<http://dustintran.com>

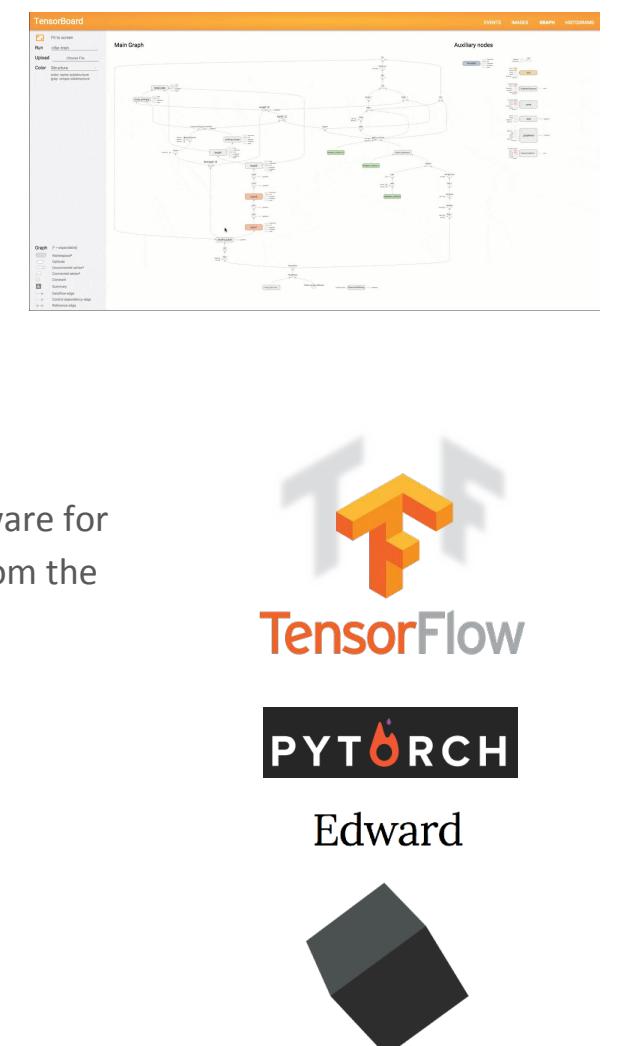
Matthew Feickert
High Energy Physics Ph.D. Candidate
Southern Methodist University
matthew.feickert@cern.ch or mfieckert@smu.edu
GitHub: [@HFFeickert](https://github.com/HFFeickert)

The modern AI/ML software stack

Recent switch to

- Numerical computations with data flow graphs
 - TensorFlow, Theano, MXNet, etc
 - Support for CPUs and GPUs out of the box.
 - Automatic differentiation
 - Enable new ways of thinking (model composition, learning to learn, etc)
- Probabilistic programming languages
 - Stan, Anglican, Edward, etc

Recommendation. The next generation of physics software for high-level analysis should take notice and inspiration from the AI/ML community.



pyhf matures

Thanks to the amazing developer team, pyhf quickly matured along with scikit-hep &PyHEP community

- Now it is a NumFOCUS affiliated project



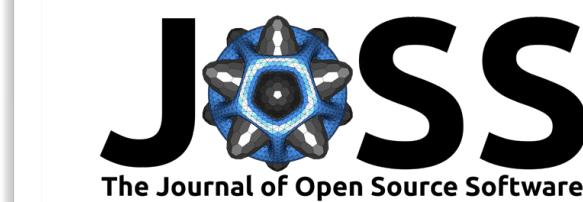
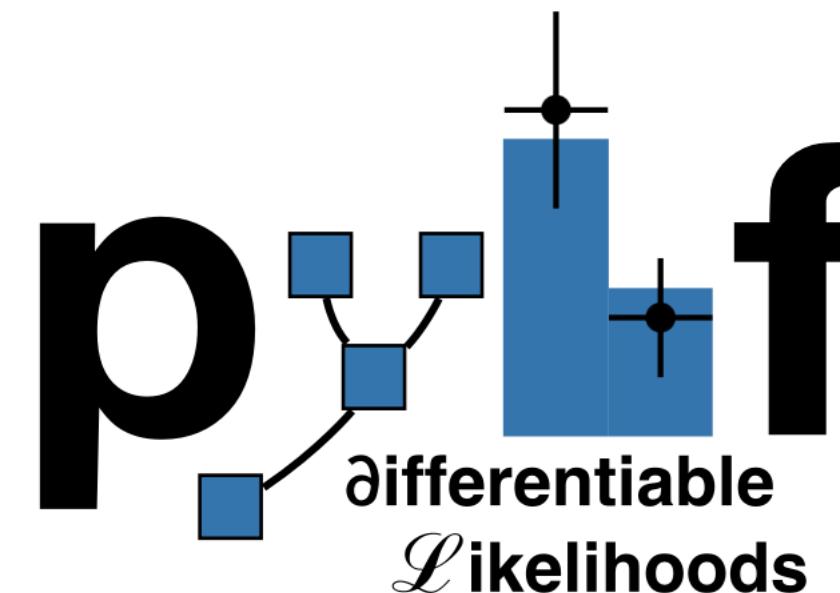
Matthew Feickert



Giordon Stark



Lukas Heinrich



The Journal of Open Source Software

pyhf: pure-Python implementation of HistFactory statistical models

Lukas Heinrich¹, Matthew Feickert^{*2}, Giordon Stark³, and Kyle Cranmer⁴

¹ CERN ² University of Illinois at Urbana-Champaign ³ SCIPP, University of California, Santa Cruz
⁴ New York University

DOI: [10.21105/joss.02823](https://doi.org/10.21105/joss.02823)

Software

- [Review ↗](#)
- [Repository ↗](#)
- [Archive ↗](#)

Editor: Eloisa Bentivegna ↗

Reviewers:

- @suchitakulkarni
- @bradkav

Submitted: 07 October 2020
Published: 04 February 2021

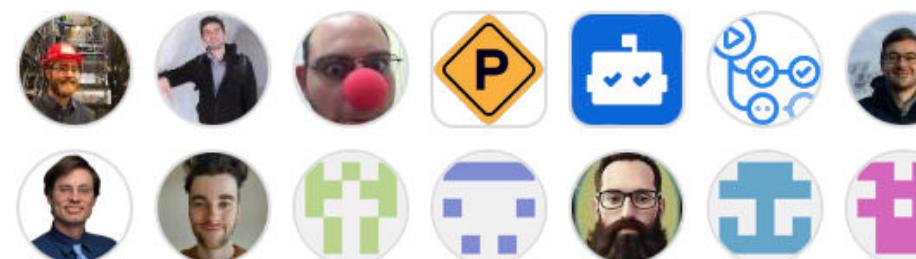
License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

Statistical analysis of High Energy Physics (HEP) data relies on quantifying the compatibility of observed collision events with theoretical predictions. The relationship between them is often formalised in a statistical model $f(\mathbf{x}|\phi)$ describing the probability of data \mathbf{x} given model parameters ϕ . Given observed data, the likelihood $\mathcal{L}(\phi)$ then serves as the basis for inference on the parameters ϕ . For measurements based on binned data (histograms), the HistFactory family of statistical models (Cranmer et al., 2012) has been widely used in both Standard Model measurements (ATLAS Collaboration, 2013) as well as searches for new physics (ATLAS Collaboration, 2018). pyhf is a pure-Python implementation of the HistFactory model specification and implements a declarative, plain-text format for describing HistFactory-based likelihoods that is targeted for reinterpretation and long-term preservation in analysis data repositories such as HEPData (Maguire et al., 2017). The source code for pyhf has been archived on Zenodo with the linked DOI: (Heinrich, Lukas and Feickert, Matthew and Stark, Giordon, 2020). At the time of writing this paper, the most recent release of pyhf is v0.5.4.

Contributors 36

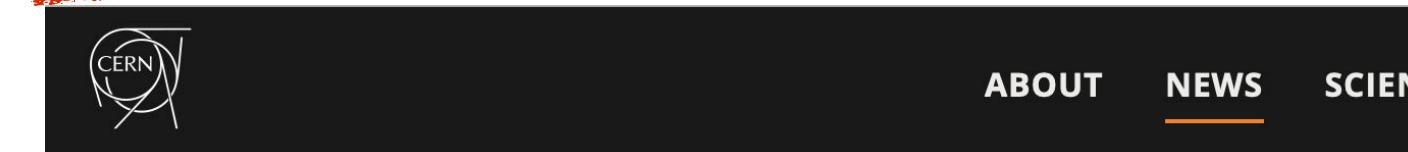


+ 22 contributors



Realizing a 20 year old dream

The pyhf JSON serialization format
facilitated publishing statistical models,
realizing a 20 year old dream for the field.

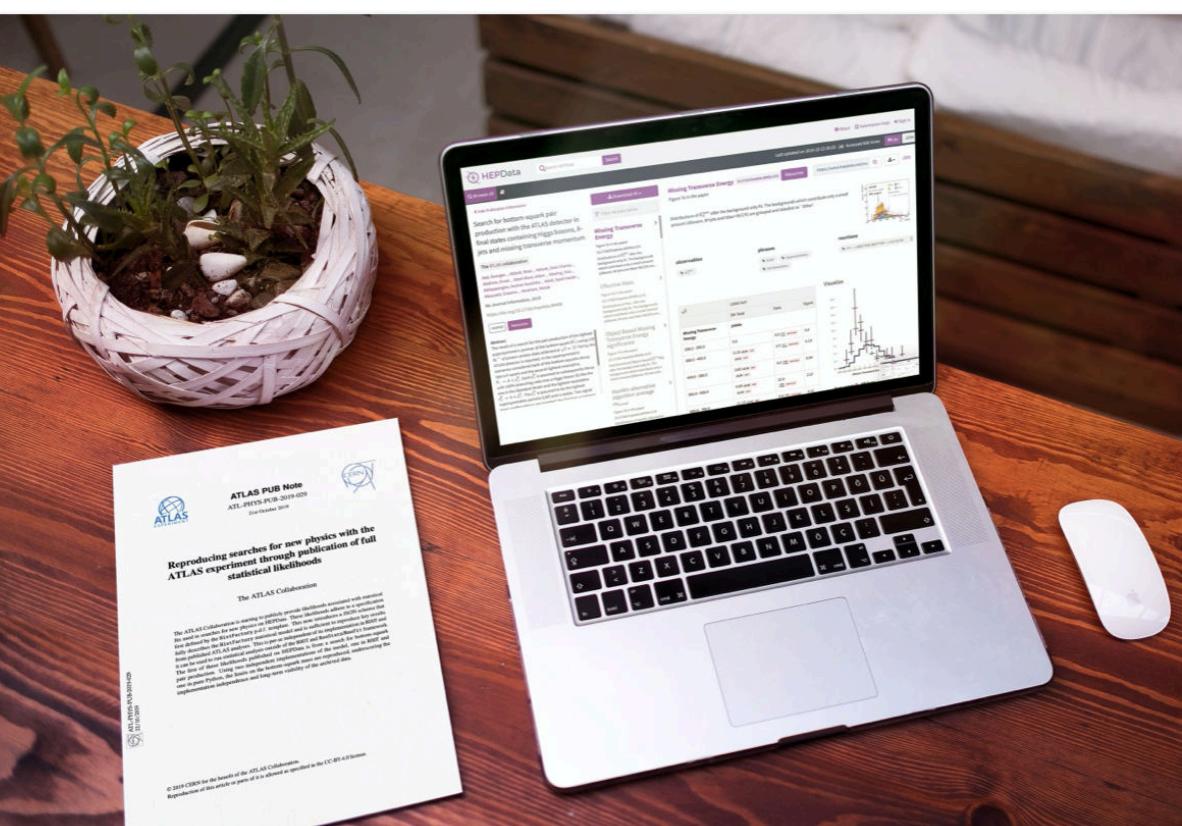


ABOUT NEWS SCIENCE RESOURCES SEARCH | EN

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

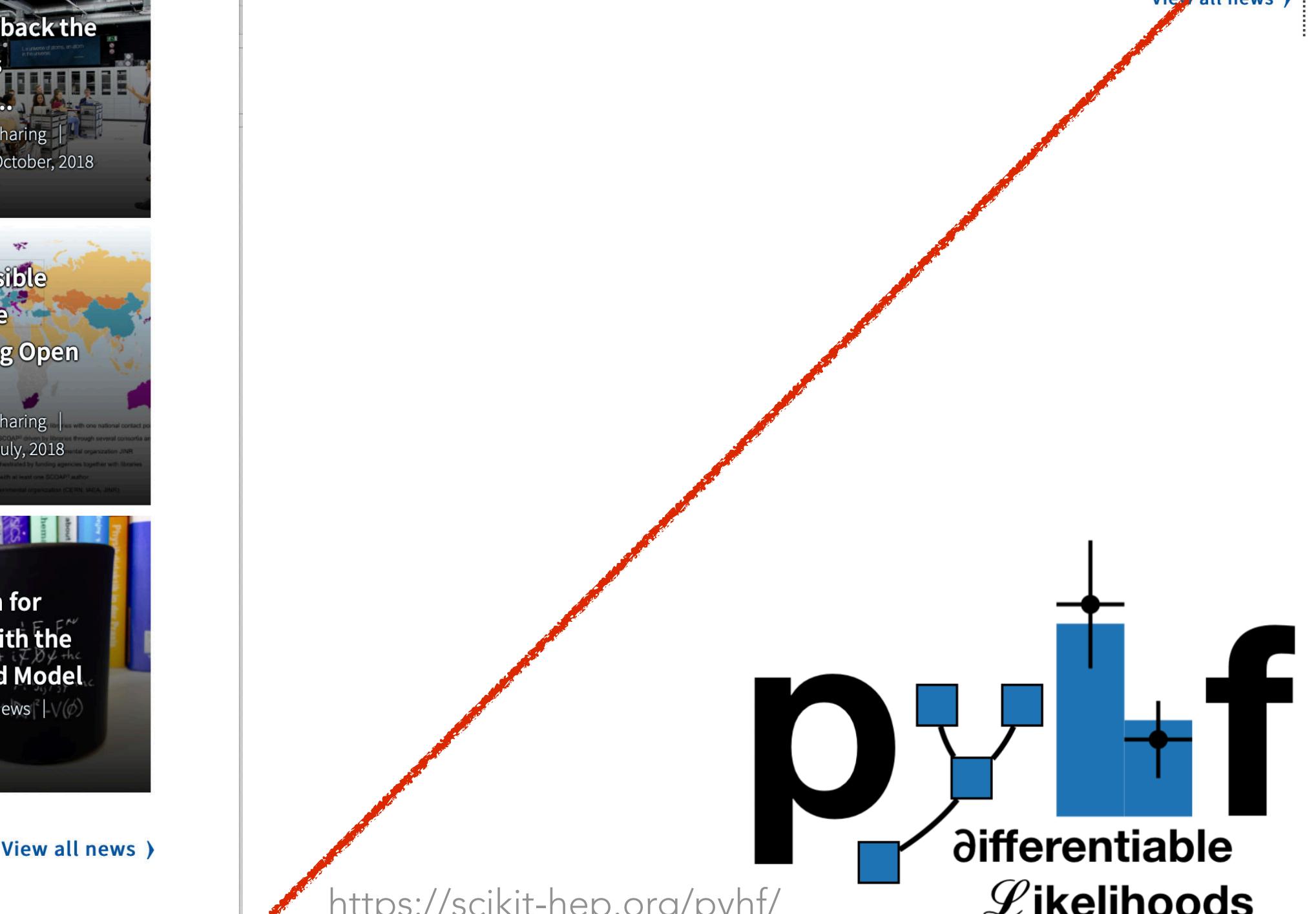
9 JANUARY, 2020 | By Katarina Anthony



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

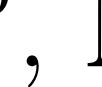
What if you could test a new theory against LHC data? Better yet, what if the expert knowledge needed to do this was captured in a convenient format? This tall order is now on display a menu from the ATLAS collaboration, with the first open release of full analysis likelihoods

A screenshot of the CERN website. At the top right, there's a navigation bar with links for 'ABOUT', 'NEWS', 'SCIENCE', 'RESOURCES', 'SEARCH | EN'. Below the navigation, there's a banner image of a large yellow dome at night. To the left of the banner, a news article is shown: 'CERN Council appoints Fabiola Gianotti for second term of office as CERN Director General' (Press release | 6 November, 2019). To the right of the banner, there's a section titled 'LATEST NEWS' with three smaller news cards: 'Particle physicists formulate the f...' (Physics | News | 20 January, 2020), 'LHCb explores the beauty of lepton universali...' (Physics | News | 15 January, 2020), and 'New open release allows theorists to explore ...' (Knowledge sharing | News | 9 January, 2020). On the far right, there's a vertical sidebar with links for 'Welcome', 'News', 'Explore CERN', 'Events', and 'Key Achievements'. A red line starts from the bottom right corner of the image and curves upwards towards the top left, highlighting the news section.



<https://scikit-hep.org/pyhf/>

Publishing statistical models: Getting the most out of particle physics experiments

Kyle Cranmer ^{1*}, Sabine Kraml ^{2†}, Harrison B. Prosper ^{3§} (editors), Philip Bechtle ⁴, Florian U. Bernlochner ⁴, Itay M. Bloch ⁵, Enzo Canonero ⁶, Marcin Chrzaszcz ⁷, Andrea Coccato ⁸, Jan Conrad ⁹, Glen Cowan ¹⁰, Matthew Feickert ¹¹, Nahuel Ferreiro Iachellini ^{12,13}, Andrew Fowlie ¹⁴, Lukas Heinrich ¹⁵, Alexander Held ¹, Thomas Kuhr ^{13,16}, Anders Kvellestad ¹⁷, Maeve Madigan ¹⁸, Farvah Mahmoudi ^{15,19}, Knut Dundas Morå ²⁰, Mark S. Neubauer ¹¹, Maurizio Pierini ¹⁵, Juan Rojo ⁸, Sezen Sekmen ²², Luca Silvestrini ²³, Veronica Sanz ^{24,25}, Giordon Stark ²⁶, Riccardo Torre ⁸, Robert Thorne ²⁷, Wolfgang Waltenberger ²⁸, Nicholas Wardle ²⁹, Jonas Wittbrodt ³⁰

It's a reality

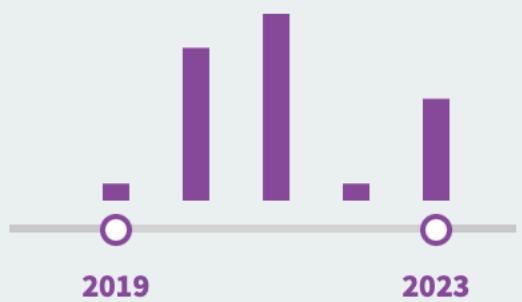


Search analysis:HistFactory

Max results ▾ Sort by ▾ Reverse order

Showing 10 of 28 results

Date



Collaboration

ATLAS 28

Subject_areas

hep-ex 28

Phrases

Proton-Proton Scattering 3

Cross Section 2

SUSY 2

Supersymmetry 2

Top 2

[Next 5](#) [Show All](#)

Reactions

P P → CHARGINO+ CHARGINO- 1

P P → CHARGINO+ NEUTRALINO 1

P P → CHARGINO+ NEUTRALINO 1

<https://www.hepdata.net/search/?q=analysis:HistFactory>

Find all papers which include specific types of analysis.

analysis:rivet (Rivet analysis)

analysis:MadAnalysis (MadAnalysis 5 analysis)

analysis:HistFactory (likelihoods in HistFactory format)

« < 1 2 > »

HistFactory HistFactory HistFactory HistFactory Search for flavour-changing neutral-current couplings between the top quark and the photon with the ATLAS detector at $\sqrt{s} = 13$ TeV

The ATLAS collaboration Aad, Georges ; Abbott, Braden Keim ; Abbott, Dale ; et al.

Phys.Lett.B 842 (2023) 137379, 2023.

Inspire Record 2077557 DOI 10.17182/hepdata.129959

This letter documents a search for flavour-changing neutral currents (FCNCs), which are strongly suppressed in the Standard Model, in events with a photon and a top quark with the ATLAS detector. The analysis uses data collected in pp collisions at $\sqrt{s} = 13$ TeV during Run 2 of the LHC, corresponding to an integrated luminosity of 139 fb^{-1} . Both FCNC top-quark production and decay are considered. The final state consists of a charged lepton, missing transverse momentum, a b -tagged jet, on...

0 data tables match query

HistFactory HistFactory Measurement of the $t\bar{t}t\bar{t}$ production cross section in $\mu\mu$ collisions at $\sqrt{s}=13$ TeV with the ATLAS detector

The ATLAS collaboration Aad, Georges ; Abbott, Braden Keim ; Abbott, Dale ; et al.

JHEP 11 (2021) 118, 2021.

Inspire Record 1869695 DOI 10.17182/hepdata.105039

A measurement of four-top-quark production using proton-proton collision data at a centre-of-mass energy of 13 TeV collected by the ATLAS detector at the Large Hadron Collider corresponding to an integrated luminosity of 139 fb^{-1} is presented. Events are selected if they contain a single lepton (electron or muon) or an opposite-sign lepton pair, in association with multiple jets. The events are categorised according to the number of jets and how likely these are to contain b -hadrons. A...

0 data tables match query

HistFactory HistFactory Observation of single-top-quark production in association with a photon using the ATLAS detector

It's a reality

INSPIRE HEP

literature

Help Submit Login

Literature Authors Jobs Seminars Conferences More...

pyhf: pure-Python implementation of HistFactory statistical models

Lukas Heinrich (CERN), Matthew Feickert (Illinois U., Urbana), Giordon Stark (UC, Santa Cruz, Inst. Part. Phys.), Kyle Cranmer (New York U.)
Feb 4, 2021
2 pages
Published in: *J.Open Source Softw.* 6 (2021) 58, 2823
Published: Feb 4, 2021
DOI: [10.21105/joss.02823](https://doi.org/10.21105/joss.02823)
View in: [CERN Document Server](#)

pdf cite claim reference search 128 citations

Citations per year

Year	Citations
2020	0
2021	15
2022	38
2023	45
2024	25

Use and Citations

Warning: This is a development version and should not be cited. To find the specific version to cite, please go to [ReadTheDocs](#).

Citation

The preferred BibTeX entry for citation of `pyhf` includes both the [Zenodo](#) archive and the [JOSS](#) paper:

```
@software{pyhf,
    author = {Lukas Heinrich and Matthew Feickert and Giordon Stark},
    title = "{pyhf: v0.7.6}",
    version = {0.7.6},
    doi = {10.5281/zenodo.1169739},
    url = {https://doi.org/10.5281/zenodo.1169739},
    note = {https://github.com/scikit-hep/pyhf/releases/tag/v0.7.6}
}

@article{pyhf_joss,
    doi = {10.21105/joss.02823},
    url = {https://doi.org/10.21105/joss.02823},
    year = {2021},
    publisher = {The Open Journal},
    volume = {6},
    number = {58},
    pages = {2823},
    author = {Lukas Heinrich and Matthew Feickert and Giordon Stark and Kyle Cranmer},
    title = {pyhf: pure-Python implementation of HistFactory statistical models},
    journal = {Journal of Open Source Software}
}
```

HEPData

analysis:HistFactory

Search Reset search Advanced JSON

Max results Sort by Reverse order Showing 10 of 28 results

Date: 2019 - 2023

Collaboration: ATLAS (28)

Subject_areas: hep-ex (28)

Phrases: Proton-Proton Scattering (3), Cross Section (2), SUSY (2), Supersymmetry (2), Top (2)

Reactions: P P -> CHARGINO+ CHARGINO- (1), P P -> CHARGINO+ NEUTRALINO (1), P P -> CHARGINO- NEUTRALINO (1)

Search results:

- Search for flavour-changing neutral-current couplings between the top quark and the photon with the ATLAS detector at $\sqrt{s} = 13$ TeV**
The ATLAS collaboration Aad, Georges; Abbott, Braden Keim; Abbott, Dale; et al.
Phys.Lett.B 842 (2023) 137379, 2023.
Inspire Record 2077557 % DOI 10.17182/hepdata.129959
- Measurement of the $t\bar{t}t\bar{t}$ production cross section in $p\mu$ collisions at $\sqrt{s}=13$ TeV with the ATLAS detector**
The ATLAS collaboration Aad, Georges; Abbott, Braden Keim; Abbott, Dale; et al.
JHEP 11 (2021) 118, 2021.
Inspire Record 1869695 % DOI 10.17182/hepdata.105039
- Observation of single-top-quark production in association with a photon using the ATLAS detector**

Use in Publications

The following is an updating list of citations and use cases of `pyhf`. There is an incomplete but automatically updated [list of citations on INSPIRE](#) as well.

Use Citations

- Alexander Held, Elliott Kauffman, Oksana Shadura, and Andrew Wightman. Physics analysis for the HL-LHC: concepts and pipelines in practice with the Analysis Grand Challenge. 1 2024. [arXiv:2401.02766](#).
- Priyotosh Bandyopadhyay, Snehashis Parashar, Chandrima Sen, and Jeonghyeon Song. Probing Inert Triplet Model at a multi-TeV muon collider via vector boson fusion with forward muon tagging. 1 2024. [arXiv:2401.02697](#).
- Elliott Kauffman, Alexander Held, and Oksana Shadura. Machine Learning for Columnar High Energy Physics Analysis. 1 2024. [arXiv:2401.01802](#).
- Mohammad Mahdi Altakach, Sabine Kraml, Andre Lessa, Sahana Narasimha, Timothée Pascal, Théo Reymermier, and Wolfgang Waltenberger. Global LHC constraints on electroweak-inos with SModels v2.3. 12 2023. [arXiv:2312.16635](#).
- MicroBooNE Collaboration. First search for dark-trident processes using the MicroBooNE detector. 12 2023. [arXiv:2312.13945](#).
- Leandro Da Rold, Manuel Epele, Anibal D. Medina, Nicolás I. Mileo, and Alejandro Szynkman. Double Higgs production at the HL-LHC: probing a loop-enhanced model with kinematical distributions. 12 2023. [arXiv:2312.13149](#).

It's a reality

INSPIRE HEP

literature

Help Submit Login

Literature Authors Jobs Seminars Conferences More...

pyhf: pure-Python implementation of HistFactory statistical models

Lukas Heinrich (CERN), Matthew Feickert (Illinois U., Urbana), Giordon Stark (UC, Santa Cruz, Inst. Part. Phys.), Kyle Cranmer (New York U.)
Feb 4, 2021
2 pages
Published in: *J.Open Source Softw.* 6 (2021) 58, 2823
Published: Feb 4, 2021
DOI: [10.21105/joss.02823](https://doi.org/10.21105/joss.02823)
View in: [CERN Document Server](#)

128 citations

Citations per year

Year	Citations
2020	2
2021	15
2022	38
2023	45
2024	26

11. Measurements of the inclusive and differential production cross sections of a top-quark-antiquark pair in association with a $Z\bar{Z}$ boson at $\sqrt{s} = 13$ TeV with the ATLAS detector. 2021. URL: <https://doi.org/10.17182/hepdata.100351>, doi:10.17182/hepdata.100351.
 12. Search for pair production of third-generation scalar leptoquarks decaying into a top quark and a τ -lepton in pp collisions at $\sqrt{s} = 13$ TeV with the ATLAS detector. 2021. URL: <https://doi.org/10.17182/hepdata.100174>, doi:10.17182/hepdata.100174.
 13. Search for squarks and gluinos in final states with one isolated lepton, jets, and missing transverse momentum at $\sqrt{s}=13$ with the ATLAS detector. 2021. URL: <https://doi.org/10.17182/hepdata.97041>, doi:10.17182/hepdata.97041.
 14. Search for trilepton resonances from chargino and neutralino pair production in $\sqrt{s}=13$ TeV pp collisions with the ATLAS detector. 2020. URL: <https://doi.org/10.17182/hepdata.99806>, doi:10.17182/hepdata.99806.
 15. Search for Displaced Leptons in $\sqrt{s} = 13$ TeV pp Collisions with the ATLAS Detector. 2020. URL: <https://doi.org/10.17182/hepdata.98796>, doi:10.17182/hepdata.98796.
 16. Search for squarks and gluinos in final states with jets and missing transverse momentum using 139 fb^{-1} of $\sqrt{s}=13$ TeV pp collision data with the ATLAS detector. 2020. URL: <https://doi.org/10.17182/hepdata.95664>, doi:10.17182/hepdata.95664.
 17. Evidence for $t\bar{t}\bar{b}t\bar{b}$ production in the multilepton final state in proton-proton collisions at $\sqrt{s}=13$ TeV with the ATLAS detector. 2020. URL: <https://doi.org/10.17182/hepdata.100170>, doi:10.17182/hepdata.100170.
 18. Measurement of the $t\bar{t}$ production cross-section in the lepton+jets channel at $\sqrt{s}=13$ TeV with the ATLAS experiment. 2020. URL: <https://doi.org/10.17182/hepdata.95748>, doi:10.17182/hepdata.95748.
 19. Search for long-lived, massive particles in events with a displaced vertex and a muon with large impact parameter in pp collisions at $\sqrt{s}=13$ TeV with the ATLAS detector. 2020. URL: <https://doi.org/10.17182/hepdata.91760>, doi:10.17182/hepdata.91760.
 20. Search for chargino-neutralino production with mass splittings near the electroweak scale in three-lepton final states in $\sqrt{s}=13$ TeV pp collisions with the ATLAS detector. 2019. URL: <https://doi.org/10.17182/hepdata.91127>, doi:10.17182/hepdata.91127.
 21. Searches for electroweak production of supersymmetric particles with compressed mass spectra in $\sqrt{s}=13$ TeV pp collisions with the ATLAS detector. 2019. URL: <https://doi.org/10.17182/hepdata.91374>, doi:10.17182/hepdata.91374.
 22. Search for direct stau production in events with two hadronic τ -leptons in $\sqrt{s}=13$ TeV pp collisions with the ATLAS detector. 2019. URL: <https://doi.org/10.17182/hepdata.92006>, doi:10.17182/hepdata.92006.
 23. Search for direct production of electroweakinos in final states with one lepton, missing transverse momentum and a Higgs boson decaying into two b -jets in (pp) collisions at $\sqrt{s}=13$ TeV with the ATLAS detector. 2019. URL: <https://doi.org/10.17182/hepdata.90607.v2>, doi:10.17182/hepdata.90607.v2.
 24. Search for squarks and gluinos in final states with same-sign leptons and jets using 139 fb^{-1} of data collected with the ATLAS detector. 2019. URL: <https://doi.org/10.17182/hepdata.91214.v3>, doi:10.17182/hepdata.91214.v3.
 25. Search for electroweak production of charginos and sleptons decaying into final states with two leptons and missing transverse momentum in $\sqrt{s}=13$ TeV pp collisions using the ATLAS detector. 2019. URL: <https://doi.org/10.17182/hepdata.89413>, doi:10.17182/hepdata.89413.
 26. Search for bottom-squark pair production with the ATLAS detector in final states containing Higgs bosons, b -jets and missing transverse momentum. 2019. URL: <https://doi.org/10.17182/hepdata.89408>, doi:10.17182/hepdata.89408.
- Note**

ATLAS maintains a public listing of all published statistical models on the [ATLAS public results page](#) which can be found by filtering all public results by the "Likelihood available" analysis characteristics keyword.

HEPData

analysis:HistFactory

Search Reset search Advanced JSON

Max results Sort by Reverse order Showing 10 of 28 results

Date: 2019 - 2023

Collaboration: ATLAS (28)

Subject_areas: hep-ex (28)

Phrases: Proton-Proton Scattering (3), Cross Section (2), SUSY (2), Supersymmetry (2), Top (2)

Reactions: P P --> CHARGINO+ CHARGINO- (1), P P --> CHARGINO+ NEUTRALINO (1), P P --> CHARGINO- NEUTRALINO (1)

Search for flavour-changing neutral-current couplings between the top quark and the photon with the ATLAS detector at $\sqrt{s} = 13$ TeV

The ATLAS collaboration Aad, Georges; Abbott, Braden Keim; Abbott, Dale; et al.
Phys.Lett.B 842 (2023) 137379, 2023.

Inspire Record 2077557 % DOI 10.17182/hepdata.129959

This letter documents a search for flavour-changing neutral currents (FCNCs), which are strongly suppressed in the Standard Model, in events with a photon and a top quark with the ATLAS detector. The analysis uses data collected in pp collisions at $\sqrt{s} = 13$ TeV during Run 2 of the LHC, corresponding to an integrated luminosity of 139 fb^{-1} . Both FCNC top-quark production and decay are considered. The final state consists of a charged lepton, missing transverse momentum, a b -tagged jet, on...

Measurement of the $t\bar{t}\bar{t}\bar{t}$ production cross section in pp collisions at $\sqrt{s}=13$ TeV with the ATLAS detector

The ATLAS collaboration Aad, Georges; Abbott, Braden Keim; Abbott, Dale; et al.
JHEP 11 (2021) 118, 2021.

Inspire Record 1869695 % DOI 10.17182/hepdata.105039

A measurement of four-top-quark production using proton-proton collision data at a centre-of-mass energy of 13 TeV collected by the ATLAS detector at the Large Hadron Collider corresponding to an integrated luminosity of 139 fb^{-1} is presented. Events are selected if they contain a single lepton (electron or muon) or an opposite-sign lepton pair, in association with multiple jets. The events are categorised according to the number of jets and how likely these are to contain b -hadrons. A...

Observation of single-top-quark production in association with a photon using the ATLAS detector

Note

ATLAS maintains a public listing of all published statistical models on the [ATLAS public results page](#) which can be found by filtering all public results by the "Likelihood available" analysis characteristics keyword.

Using published likelihoods



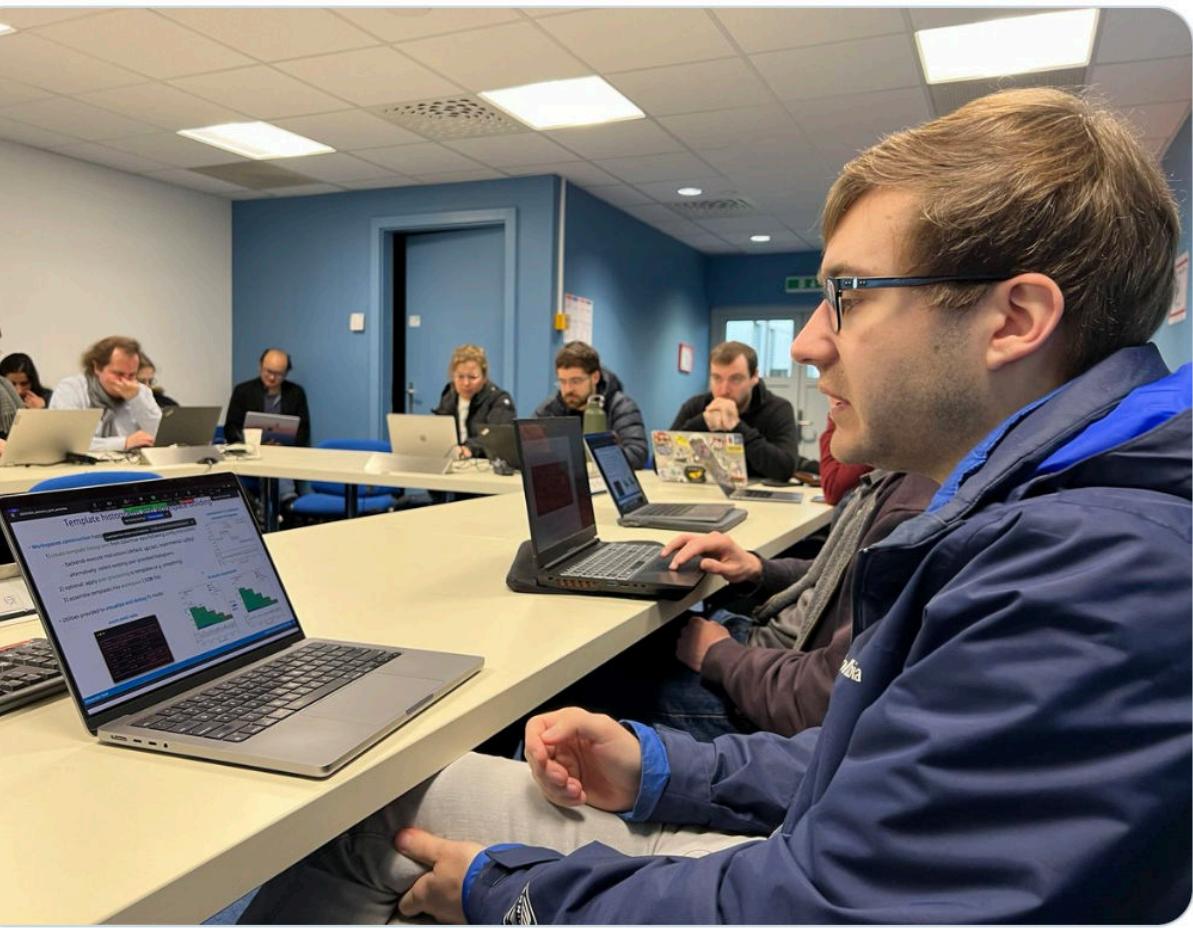
<https://github.com/scikit-hep/cabinetry>

```
1 import json
2 import cabinetry
3 import pyhf
4 from cabinetry.model_utils import prediction
5 from pyhf.contrib.utils import download
6
7 # download the ATLAS bottom-squarks analysis probability models from HEPData
8 download("https://www.hepdata.net/record/resource/1935437?view=true", "bottom-squarks")
9
10 # construct a workspace from a background-only model and a signal hypothesis
11 bkg_only_workspace = pyhf.Workspace(json.load(open("bottom-squarks/RegionC/BkgOnly.json")))
12 patchset = pyhf.PatchSet(json.load(open("bottom-squarks/RegionC/patchset.json")))
13 workspace = patchset.apply(bkg_only_workspace, "sbottom_600_280_150")
14
15 # construct the probability model and observations
16 model, data = cabinetry.model_utils.model_and_data(workspace)
17
18 # produce visualizations of the pre-fit model and observed data
19 prefit_model = prediction(model)
20 cabinetry.visualize.data_mc(prefit_model, data)
21
22 # fit the model to the observed data
23 fit_results = cabinetry.fit.fit(model, data)
24
25 # produce visualizations of the post-fit model and observed data
26 postfit_model = prediction(model, fit_results=fit_results)
27 cabinetry.visualize.data_mc(postfit_model, data)
```



pyhf - python HistFactory @pyhf_ · 17h

Now we're moving on to an overview of cabinetry from [@alheld!](#) [indico.cern.ch/event/1294577/...](https://indico.cern.ch/event/1294577/)



Alexander Held

Using published likelihoods

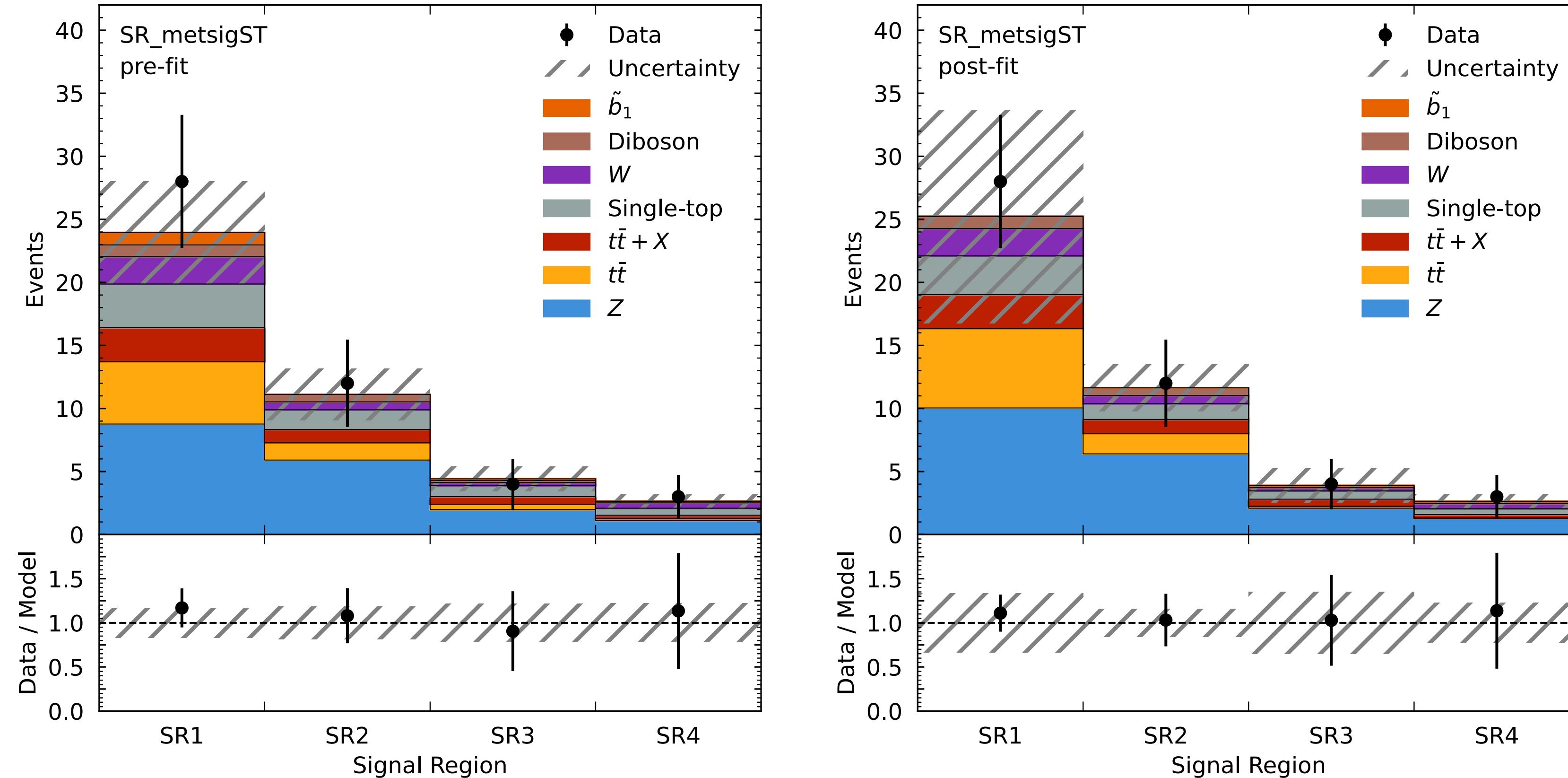


Figure 3: Pre-fit (left) and post-fit (right) visualizations of a selected signal hypothesis for four signal regions of the ATLAS search [41] of a bottom-squark of mass 600 GeV with a second-lightest neutralino of mass 280 GeV and lightest supersymmetric particle of mass 150 GeV generated from the full statistical models published in Ref. [20] using code from Ref. [40].

Declarative specifications FTW!

We are now seeing implementations in **julia** and an exciting initiative to expand this strategy to other statistical modeling tools



J. Ling 🐣 @l_II_llI · May 3, 2022

...
happy to report that after a few days of wrestling with channel-sample-bin structure tunneling/masking, as well as trying not to blow up compilation time, **LiteHF.jl** is slightly faster than even Jax+pyhf. It's AD friendly out of the box (trade-off being increased allocation).



JuliaPackages @JuliaPackages · Dec 16, 2022

...
Automated
New package: LiteHF v0.1.0 announced **#JuliaLang**

LiteHF: Light-weight **HistFactory** in pure **Julia**, attempts to be compatible with `pyhf` json format

Registration: [github.com/JuliaRegistries...](https://github.com/JuliaRegistries/registries.jl/blob/main/registries/JuliaHEP/registration.toml)

Repository: [github.com/JuliaHEP/LiteH...](https://github.com/JuliaHEP/LiteHF.jl)

HS³ - A serialization standard for statistical models in high energy physics

Carsten Burgard¹, Robin Pelkner¹

Many people involved: Matthew Feickert, Lukas Heinrich, Alexander Held, Cornelius Grunwald, Oliver Schulz, Mikhail Mikhasenko, Jerry Ling, Wouter Verkerke, Jonas Eschle, Lorenzo Moneta, Louis Moureaux, Tomas Dado and many others

pyhf Workshop 2023 - 04.12.2023

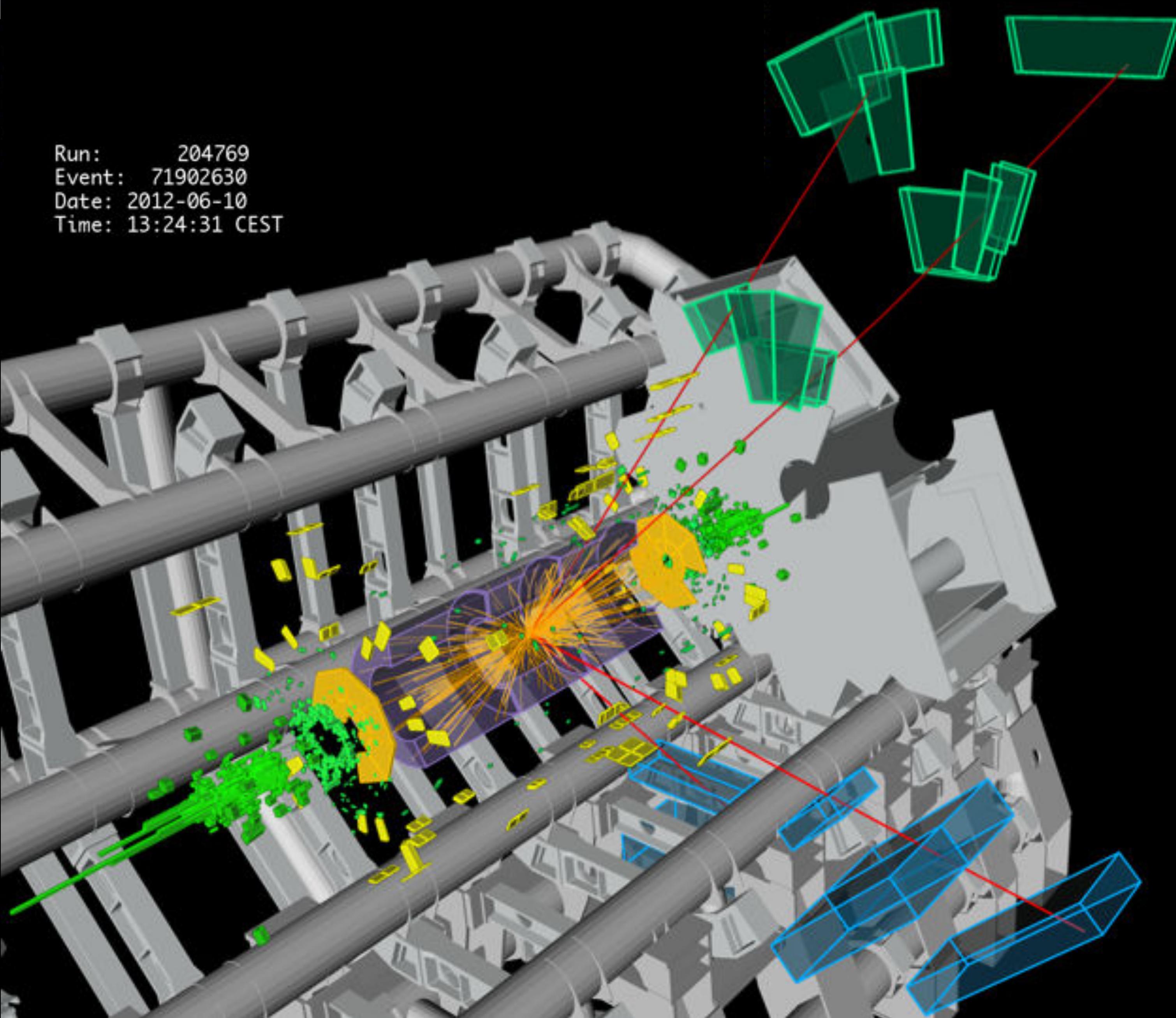
¹ TU Dortmund University

Thinking in Arrays:
the move to columnar data & array-oriented programming



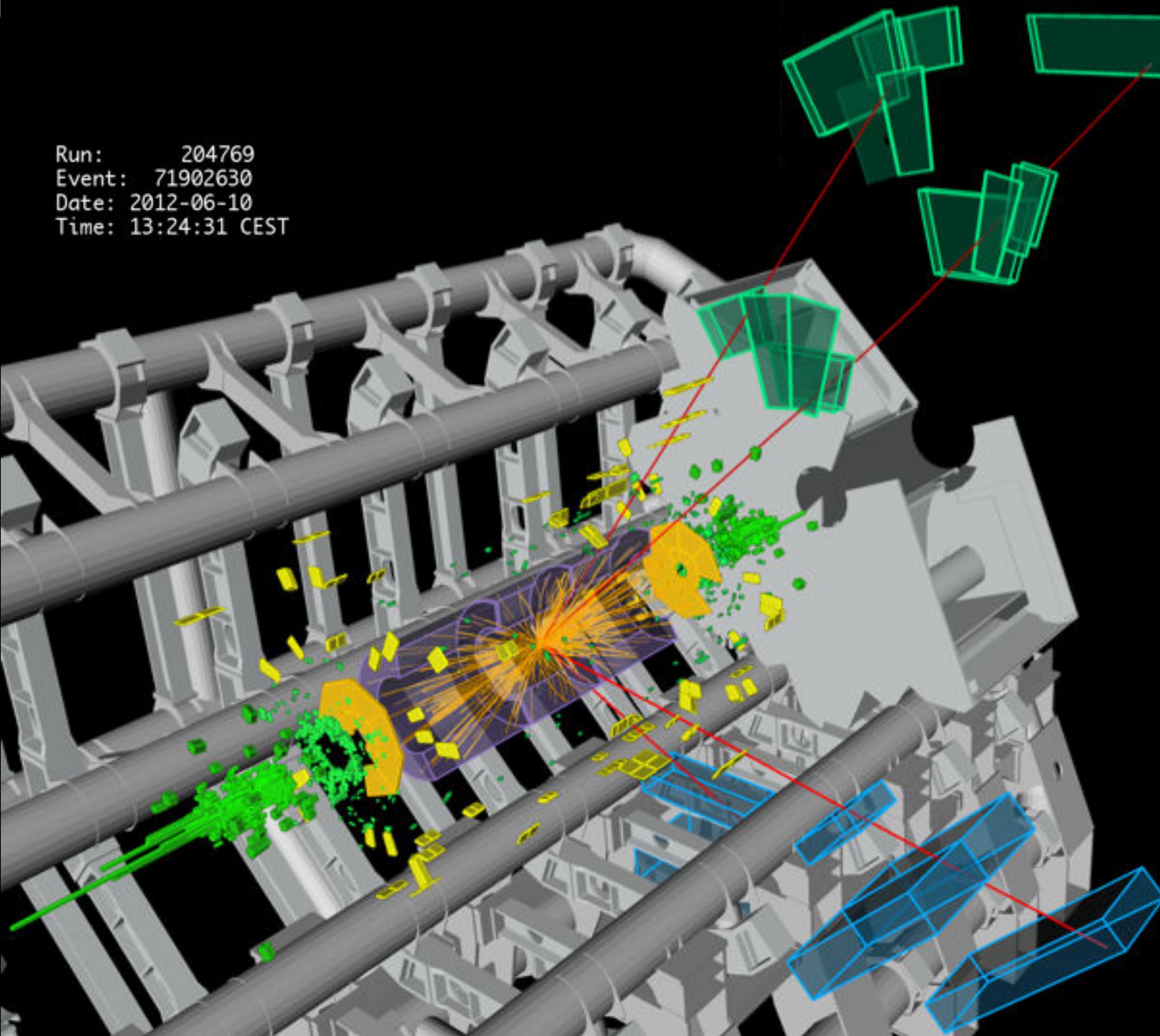
Bulk data analysis pipeline

Run: 204769
Event: 71902630
Date: 2012-06-10
Time: 13:24:31 CEST



Bulk data analysis pipeline

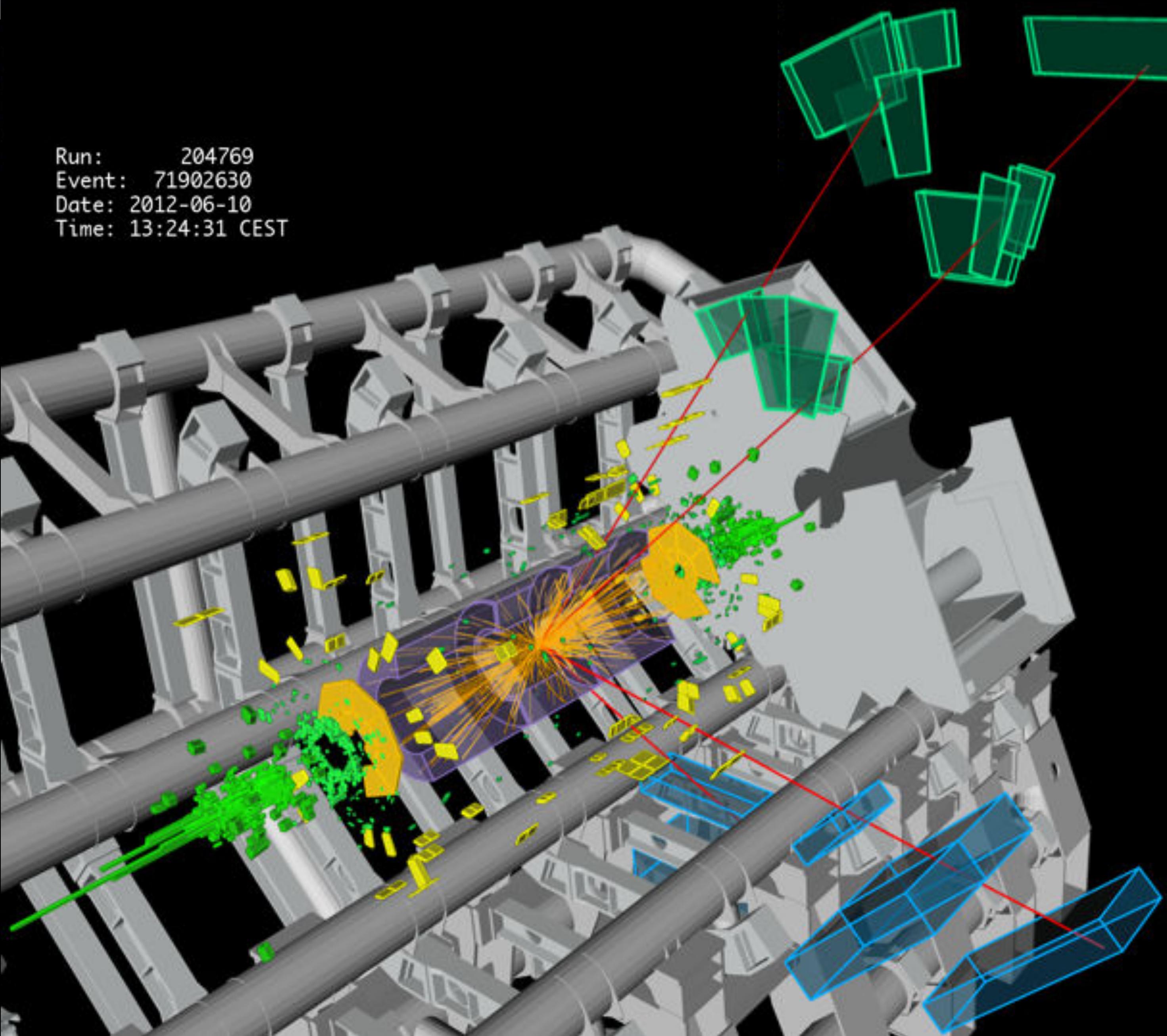
Run: 204769
Event: 71902630
Date: 2012-06-10
Time: 13:24:31 CEST



The raw data associated to one collision is complex & not convenient to work with

- Roughly 100M sensors. Sparse.

Bulk data analysis pipeline



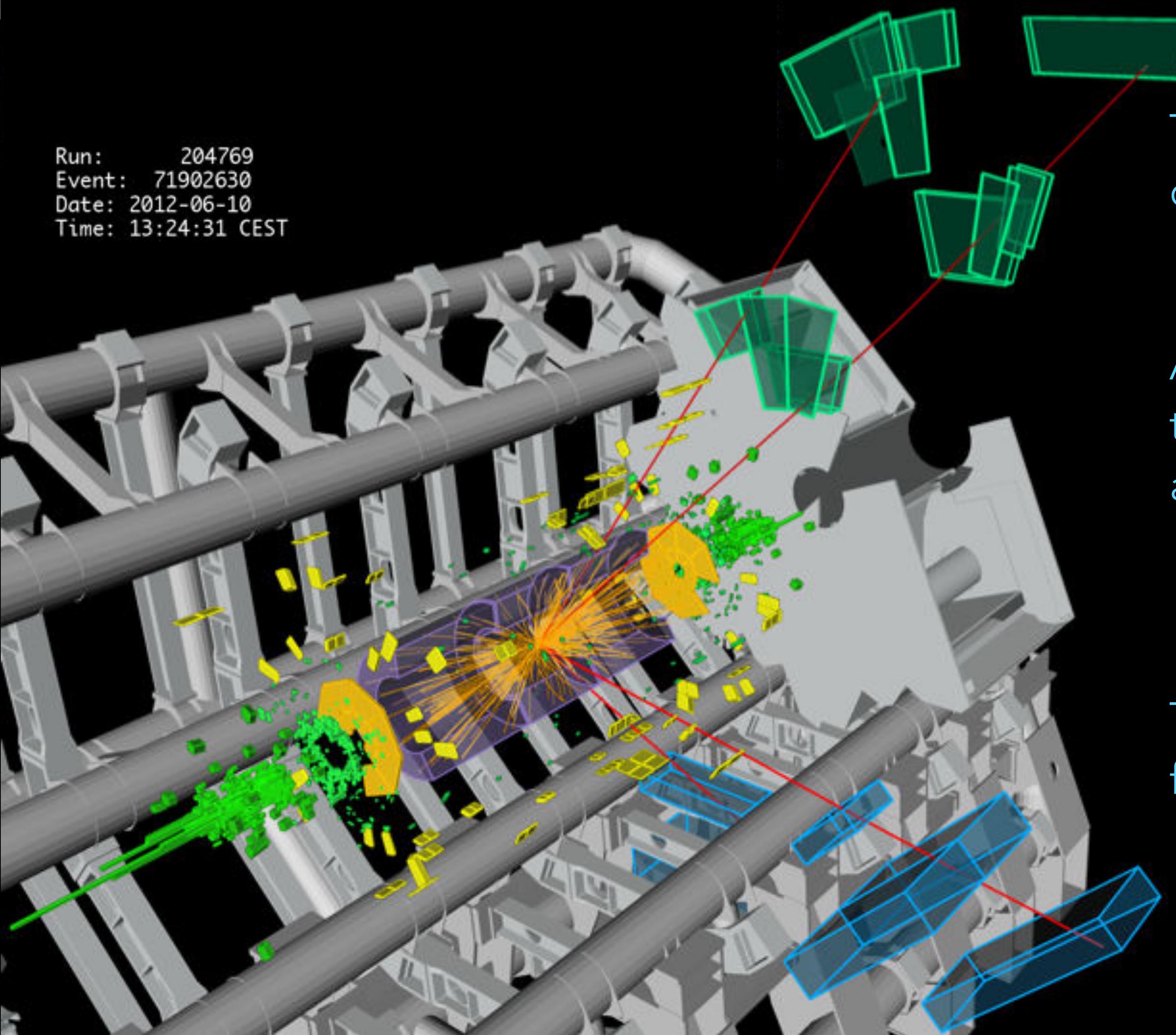
The raw data associated to one collision is complex & not convenient to work with

- Roughly 100M sensors. Sparse.

A **data analysis pipeline** runs algorithms that identify and characterize particles by analyzing energy deposits in detector

- e.g. the red lines & yellow lines are identified as particles of different types

Bulk data analysis pipeline



The raw data associated to one collision is complex & not convenient to work with

- Roughly 100M sensors. Sparse.

A **data analysis pipeline** runs algorithms that identify and characterize particles by analyzing energy deposits in detector

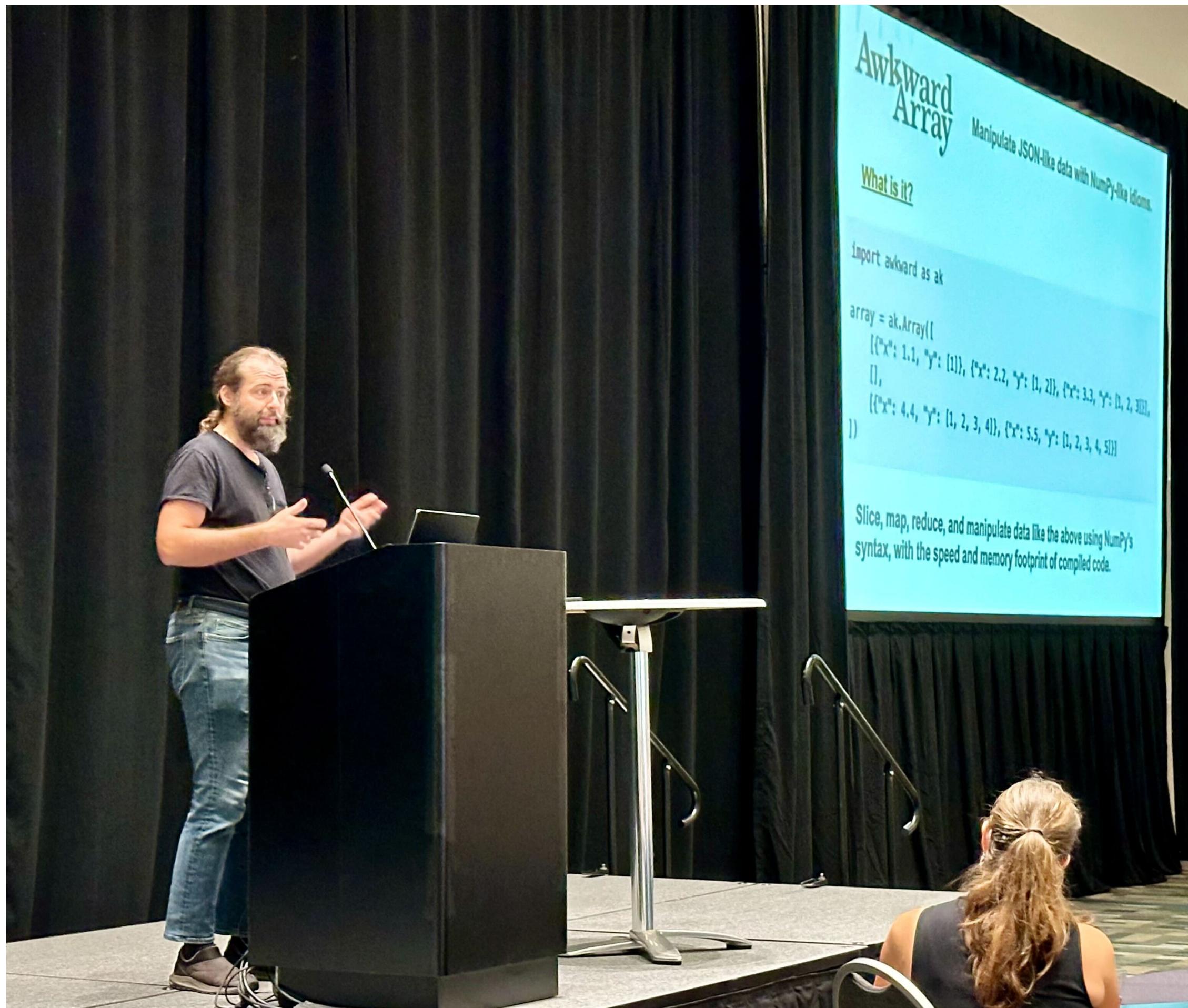
- e.g. the red lines & yellow lines are identified as particles of different types

The number of identified particles varies from collision to collision

- Output is **not tabular data**, but a nested, variable-length data structure

Awkward Array

Awkward Array is a library for nested, variable-sized data, including arbitrary-length lists, records, mixed types, and missing data, **using NumPy-like idioms.**



Awkward Array

Awkward Array is a library for nested, variable-sized data, including arbitrary-length lists, records, mixed types, and missing data, **using NumPy-like idioms.**

Awkward Array

Awkward Array is a library for nested, variable-sized data, including arbitrary-length lists, records, mixed types, and missing data, **using NumPy-like idioms.**

```
import awkward as ak

array = ak.Array([
    [{"x": 1.1, "y": [1]}, {"x": 2.2, "y": [1, 2]}, {"x": 3.3, "y": [1, 2, 3]}],
    [],
    [{"x": 4.4, "y": [1, 2, 3, 4]}, {"x": 5.5, "y": [1, 2, 3, 4, 5]}]
])
```



Awkward Array

```
array = ak.Array([
    [{"x": 1.1, "y": [1]}, {"x": 2.2, "y": [1, 2]}, {"x": 3.3, "y": [1, 2, 3]}],
    [],
    [{"x": 4.4, "y": [1, 2, 3, 4]}, {"x": 5.5, "y": [1, 2, 3, 4, 5]}]
])
```

Awkward Array

```
array = ak.Array([
    [{"x": 1.1, "y": [1]}, {"x": 2.2, "y": [1, 2]}, {"x": 3.3, "y": [1, 2, 3]}],
    [],
    [{"x": 4.4, "y": [1, 2, 3, 4]}, {"x": 5.5, "y": [1, 2, 3, 4, 5]}]
])
```

the following slices out the `y` values, drops the first element from each inner list, and runs NumPy's `np.square` function on everything that is left:

```
output = np.square(array["y", ..., 1:])
```



Awkward Array

```
array = ak.Array([
    [{"x": 1.1, "y": [1]}, {"x": 2.2, "y": [1, 2]}, {"x": 3.3, "y": [1, 2, 3]}],
    [],
    [{"x": 4.4, "y": [1, 2, 3, 4]}, {"x": 5.5, "y": [1, 2, 3, 4, 5]}]
])
```

the following slices out the `y` values, drops the first element from each inner list, and runs NumPy's `np.square` function on everything that is left:

```
output = np.square(array["y", ..., 1:])
```

The result is

```
[[], [4], [4, 9]],
[],
[[4, 9, 16], [4, 9, 16, 25]]
```

Awkward Array

The equivalent loop-based approach in pure Python is:

```
output = []
for sublist in array:
    tmp1 = []
    for record in sublist:
        tmp2 = []
        for number in record["y"][1:]:
            tmp2.append(np.square(number))
        tmp1.append(tmp2)
    output.append(tmp1)
```



- the Awkward Array one-liner takes 1.5 seconds to run and uses 2.1 GB of memory,
- the equivalent using Python lists and dicts takes 140 seconds to run and uses 22 GB of memory.

Awkward Array

The equivalent loop-based approach in pure Python is:

```
output = []
for sublist in array:
    tmp1 = []
    for record in sublist:
        tmp2 = []
        for number in record["y"][1:]:
            tmp2.append(np.square(number))
        tmp1.append(tmp2)
    output.append(tmp1)
```

Note: we would traditionally write
loop-based analysis code in C++

Move to array-oriented paradigm is a
big cultural shift.

- the Awkward Array one-liner takes 1.5 seconds to run and uses 2.1 GB of memory,
- the equivalent using Python lists and dicts takes 140 seconds to run and uses 22 GB of memory.

Playing nice

arXiv > cs > arXiv:2404.18170

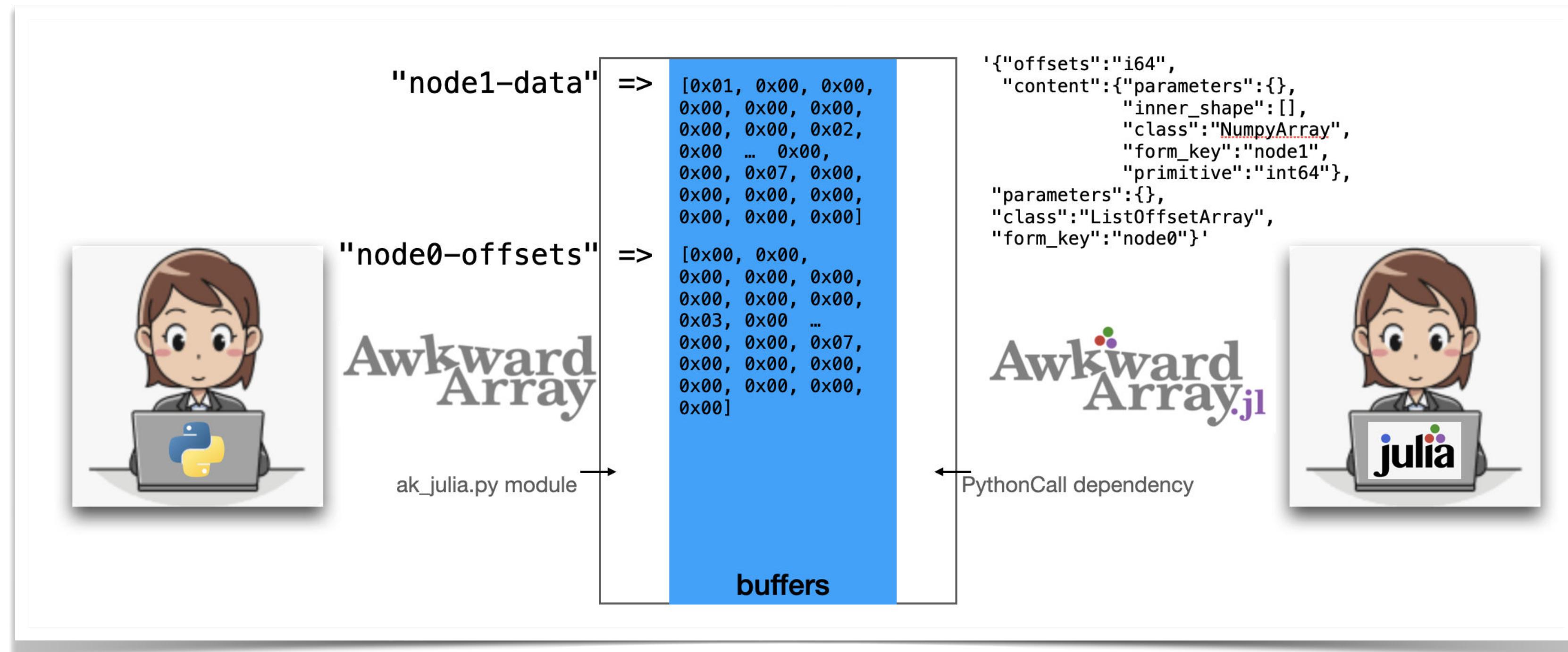
Computer Science > Programming Languages

[Submitted on 28 Apr 2024]

Bridging Worlds: Achieving Language Interoperability between Julia and Python in Scientific Computing

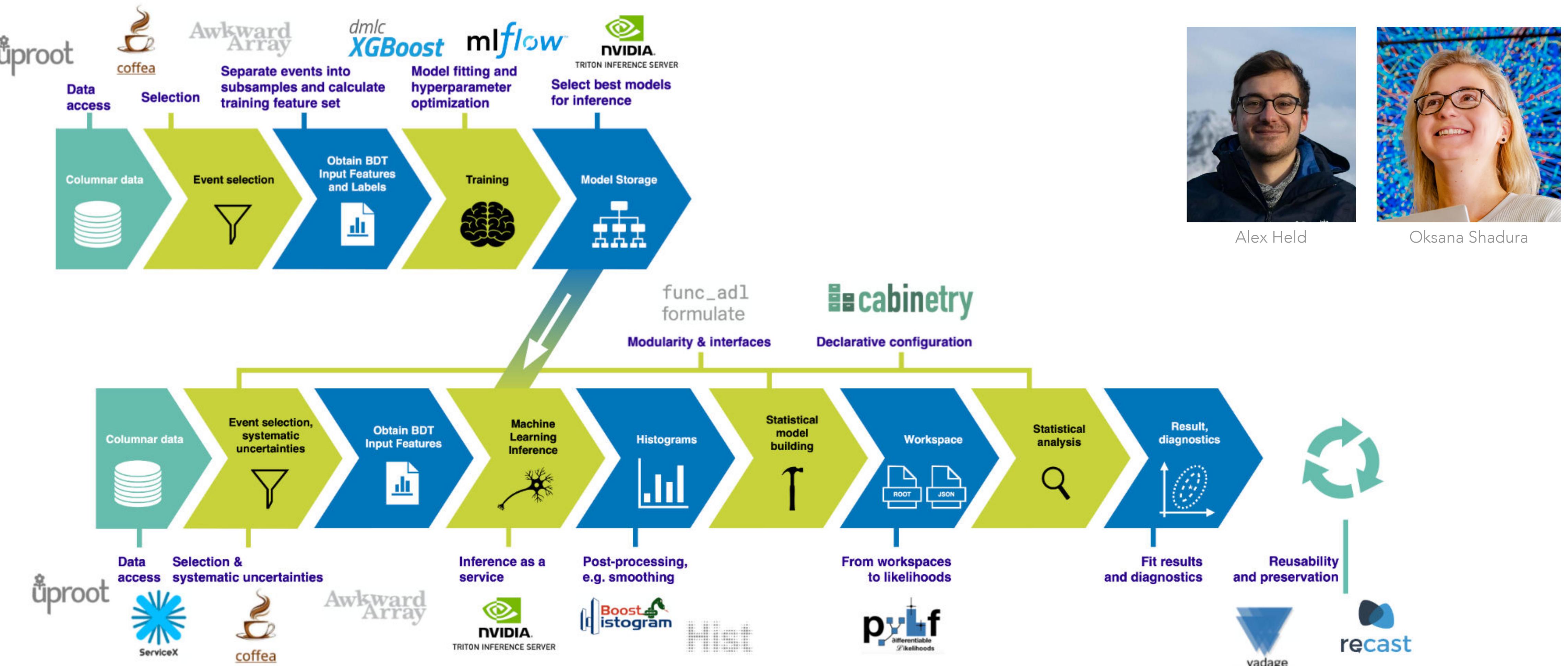
Ianna Osborne, Jim Pivarski, Jerry Ling

In the realm of scientific computing, both Julia and Python have established themselves as powerful tools. Within the context of High Energy Physics (HEP) data analysis, Python has been traditionally favored, yet there exists a compelling case for migrating legacy software to Julia. This article focuses on language interoperability, specifically exploring how Awkward Array data structures can bridge the gap between Julia and Python. The talk offers insights into key considerations such as memory management, data buffer copies, and dependency handling. It delves into the performance enhancements achieved by invoking Julia from Python and vice versa, particularly for intensive array-oriented calculations involving large-scale, though not excessively dimensional, arrays of HEP data. The advantages and challenges inherent in achieving interoperability between Julia and Python in the domain of scientific computing are discussed.



Grand Challenges

We have embraced the use of Grand Challenges as a mechanism to improve integration of the ecosystem of tools & test performance on analysis facilities.



Grand Challenges

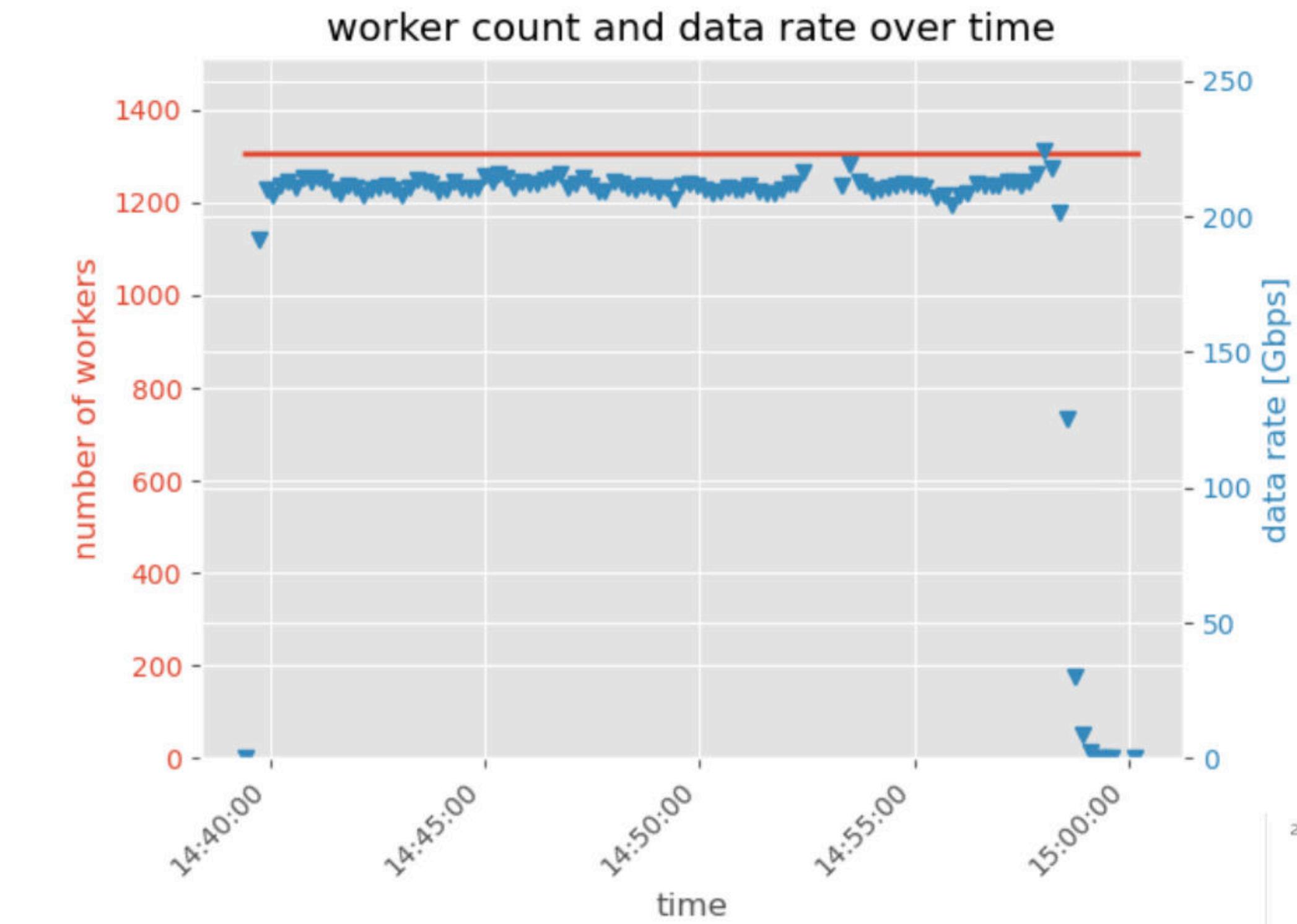
We have embraced the use of Grand Challenges as a mechanism to improve integration of the ecosystem of tools & test performance on analysis facilities.



Process 100B events in 30 min
(55MHz event processing rate)

2KB/event \Rightarrow 200TB (200 Gbps)

Analysis running on 2,200 cores

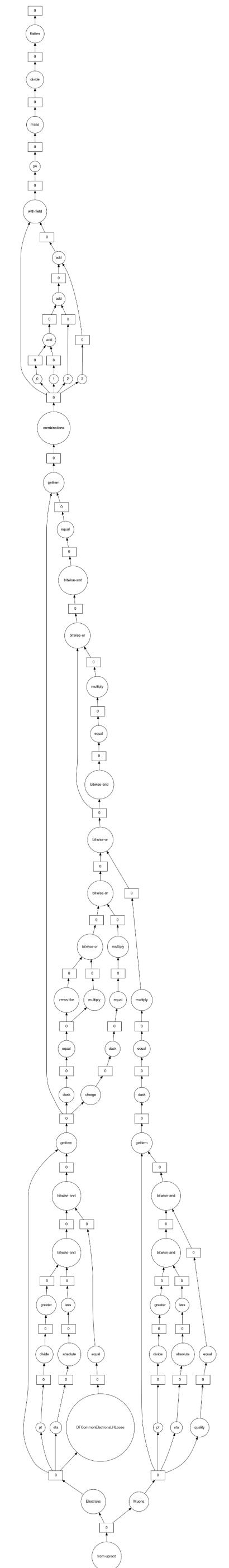
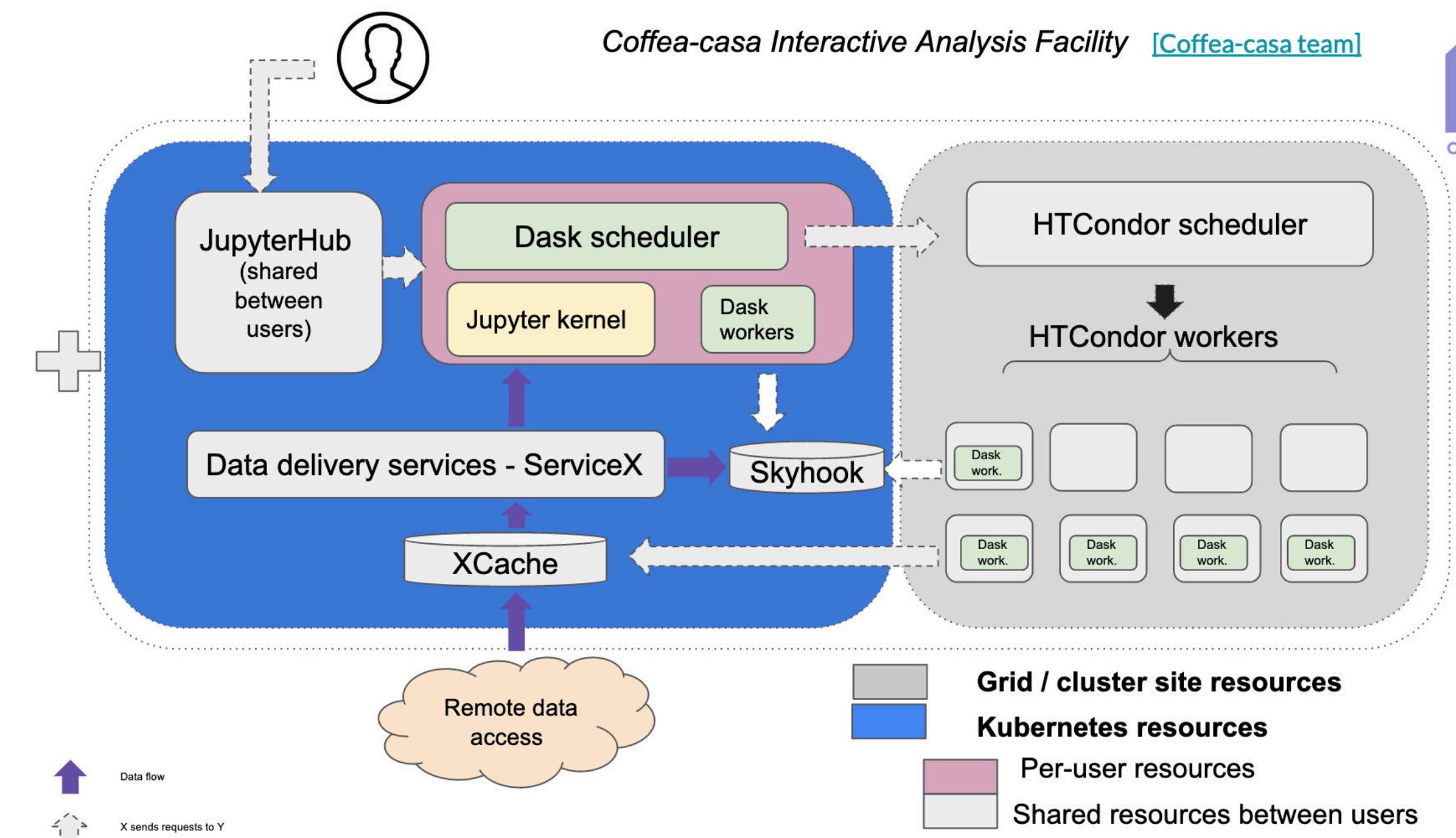


TL;DR We learned a lot about Dask

Quick shoutout



Elements featured at SciPy this week by Vangelis Kourlitis in
How the Scientific Python ecosystem helps answering fundamental questions of the Universe



Analysis Preservation & Reuse

Reinterpretation as a Service

$\mathcal{L}_{SM} =$

$$\frac{1}{4}\mathbf{W}_{\mu\nu} \cdot \mathbf{W}^{\mu\nu} - \frac{1}{4}B_{\mu\nu}B^{\mu\nu} - \frac{1}{4}G_a^a G_a^{\mu\nu}$$

kinetic energies and self-interactions of the gauge bosons

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

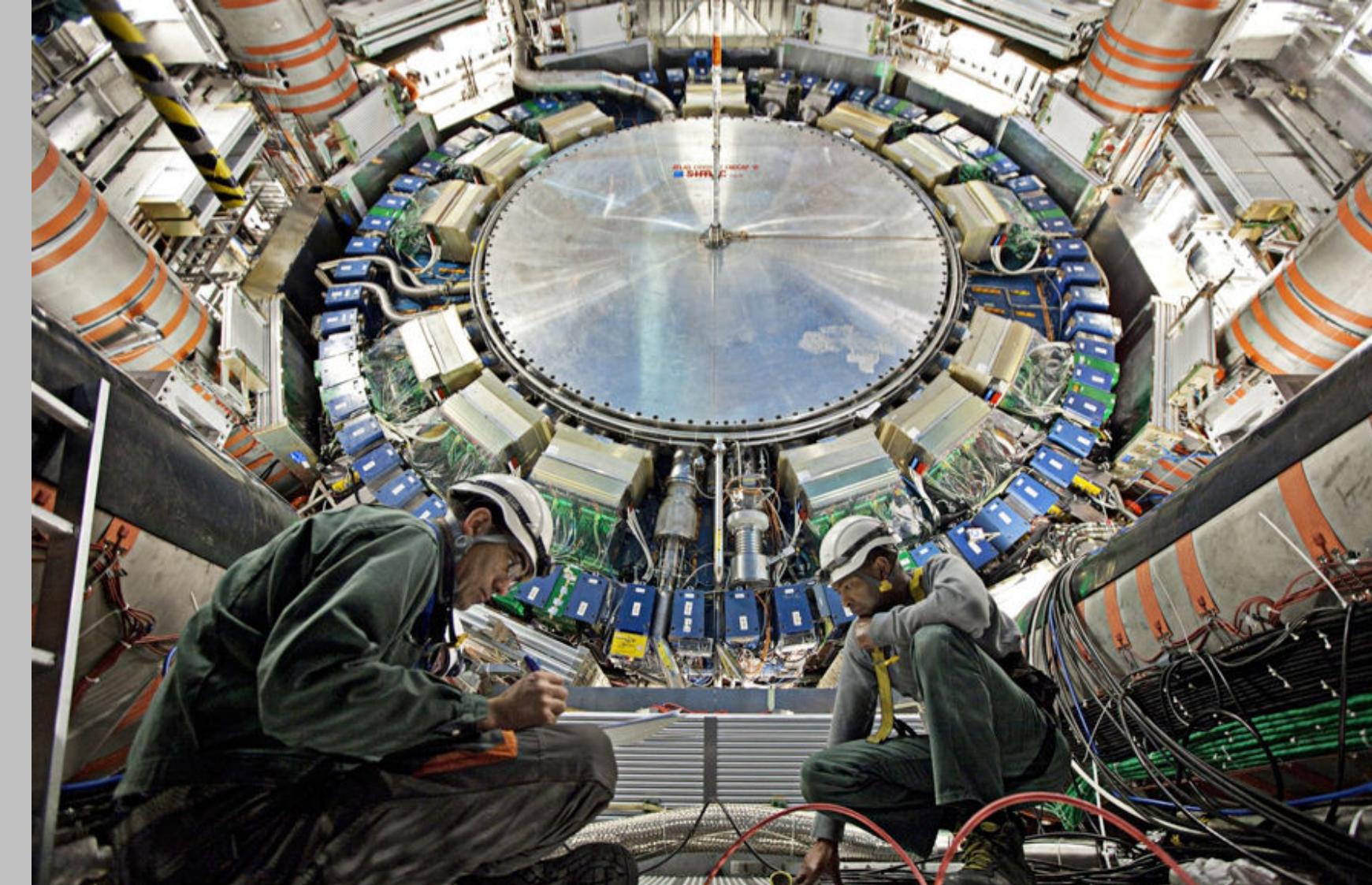
$$+ \underbrace{\frac{1}{2}|(i\partial_\mu - \frac{1}{2}g\tau \cdot \mathbf{W}_\mu - \frac{1}{2}g'YB_\mu)\phi|^2 - V(\phi)}_{W^\pm, Z, \gamma, \text{and Higgs masses and couplings}}$$

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

Q



A



ATLAS

Events

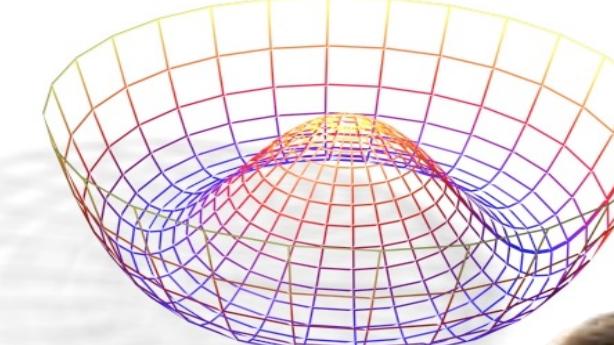
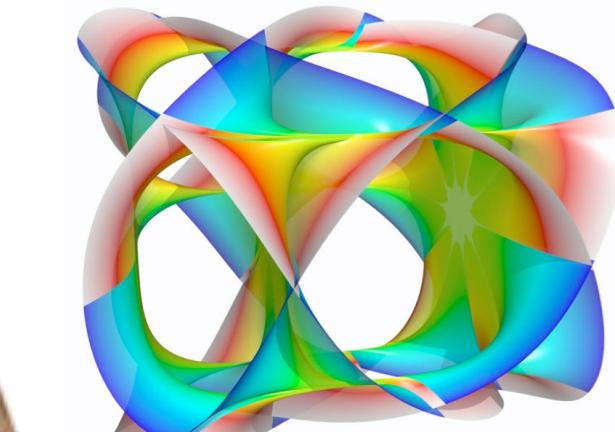
• Data
— Fit
—○— $q^*(500)$
—△— $q^*(800)$
—◆— $q^*(1200)$

$\sqrt{s} = 7 \text{ TeV}$
 $\int L dt = 315 \text{ nb}^{-1}$

(D - B) / \sqrt{B}

500 1000 1500

Reconstructed m^{jj} [GeV]

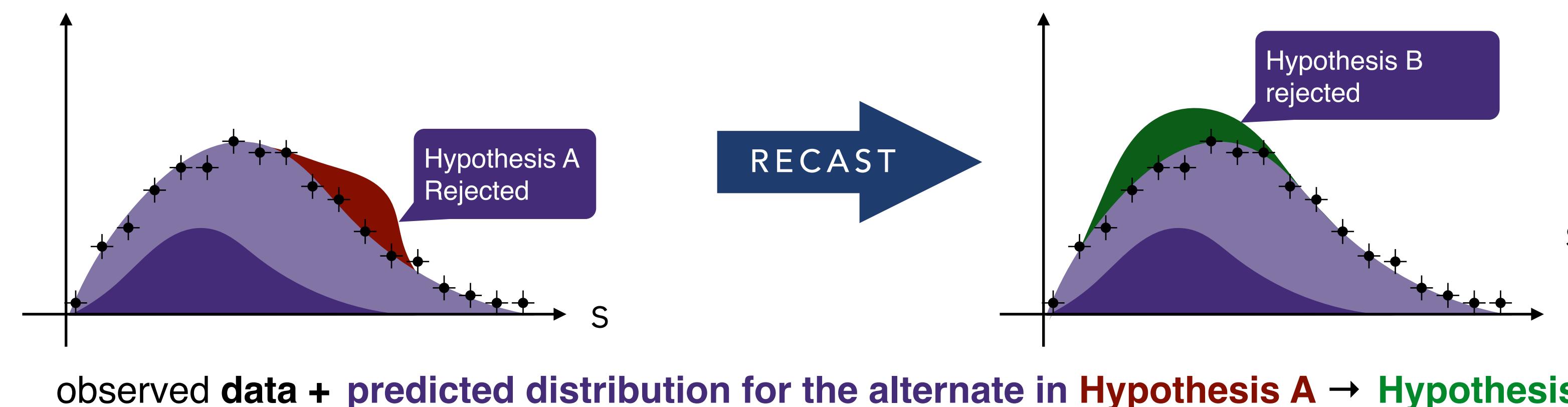




Reinterpretation infrastructure

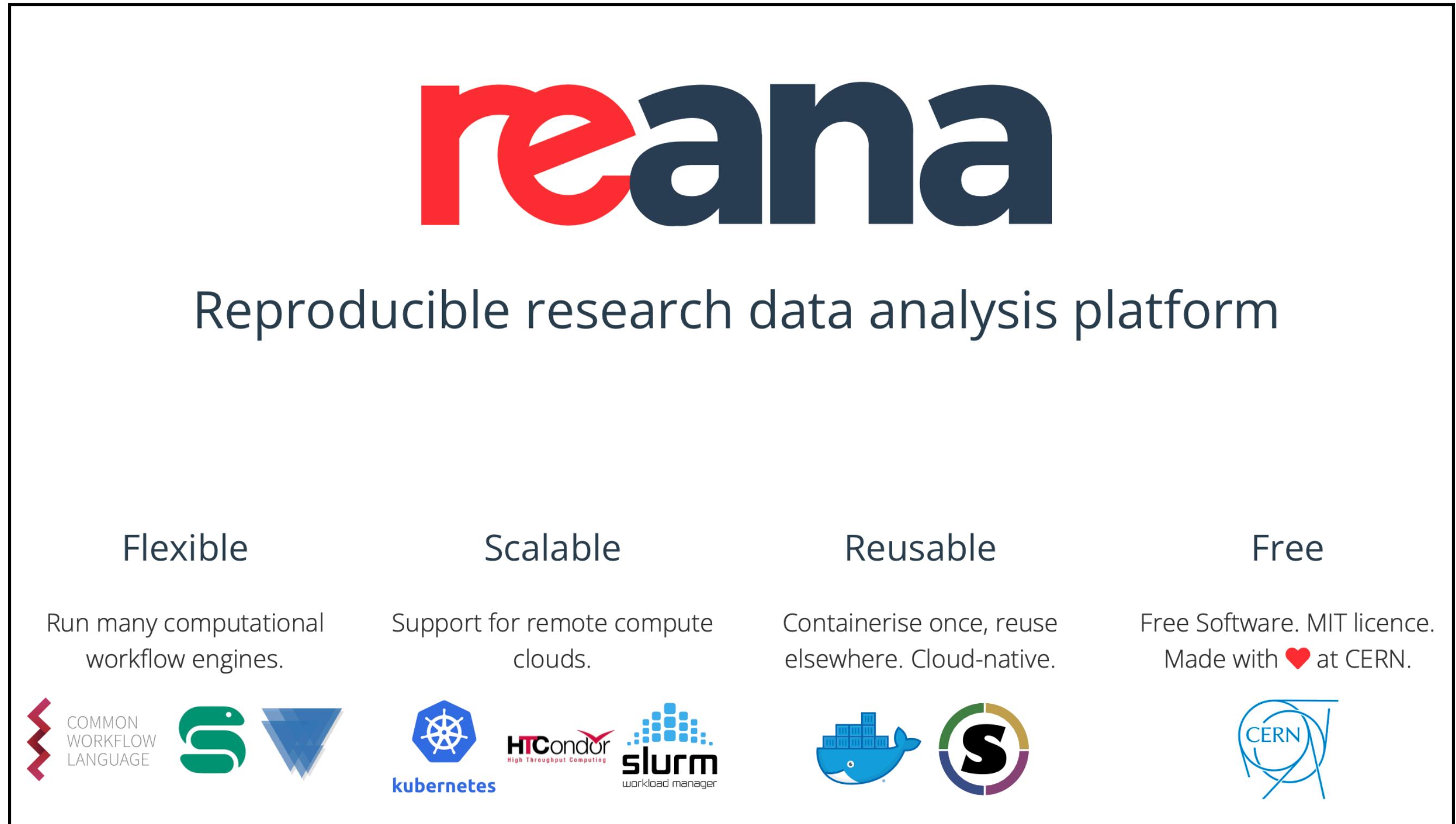
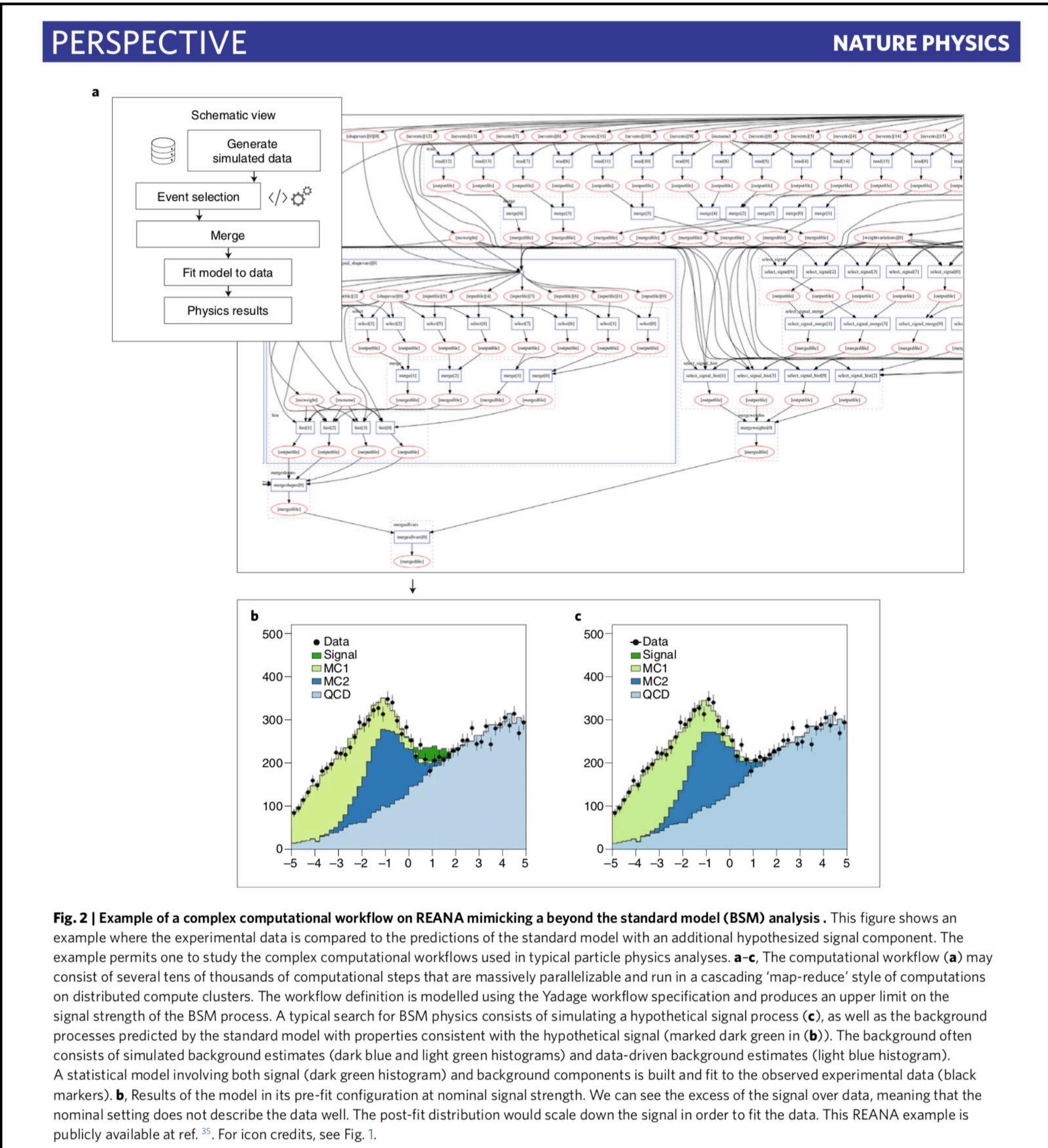
The data analysis pipeline developed to test **Hypothesis A** can be reused

- If **Hypothesis B** predicts a similar signature in the data as **Hypothesis A**, then that same pipeline may be good enough to draw conclusions for **Hypothesis B**
- Not optimal, but not wrong
- This depends on the details of pipeline and hypotheses. Need the code!

We have invested effort to preserve analysis pipeline so that it can be reused



Open is not enough!



Training

Shifting the emphasis from reproducibility to reuse was key

- Train community to build containers & create reusable workflows



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

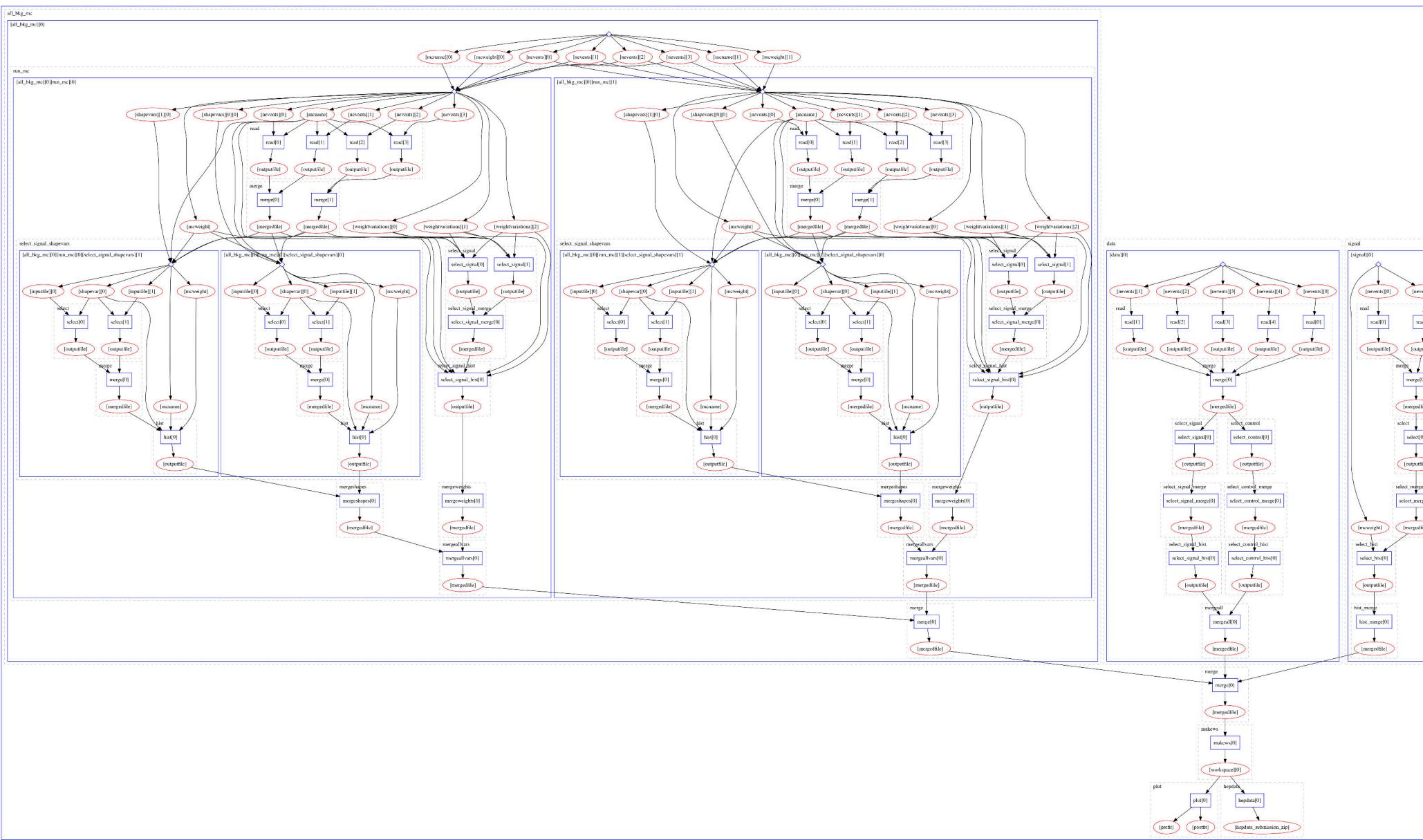


Instructors Danika MacDonnel and Giordon Stark working with participants. Photo Credit: Samuel Meehan.

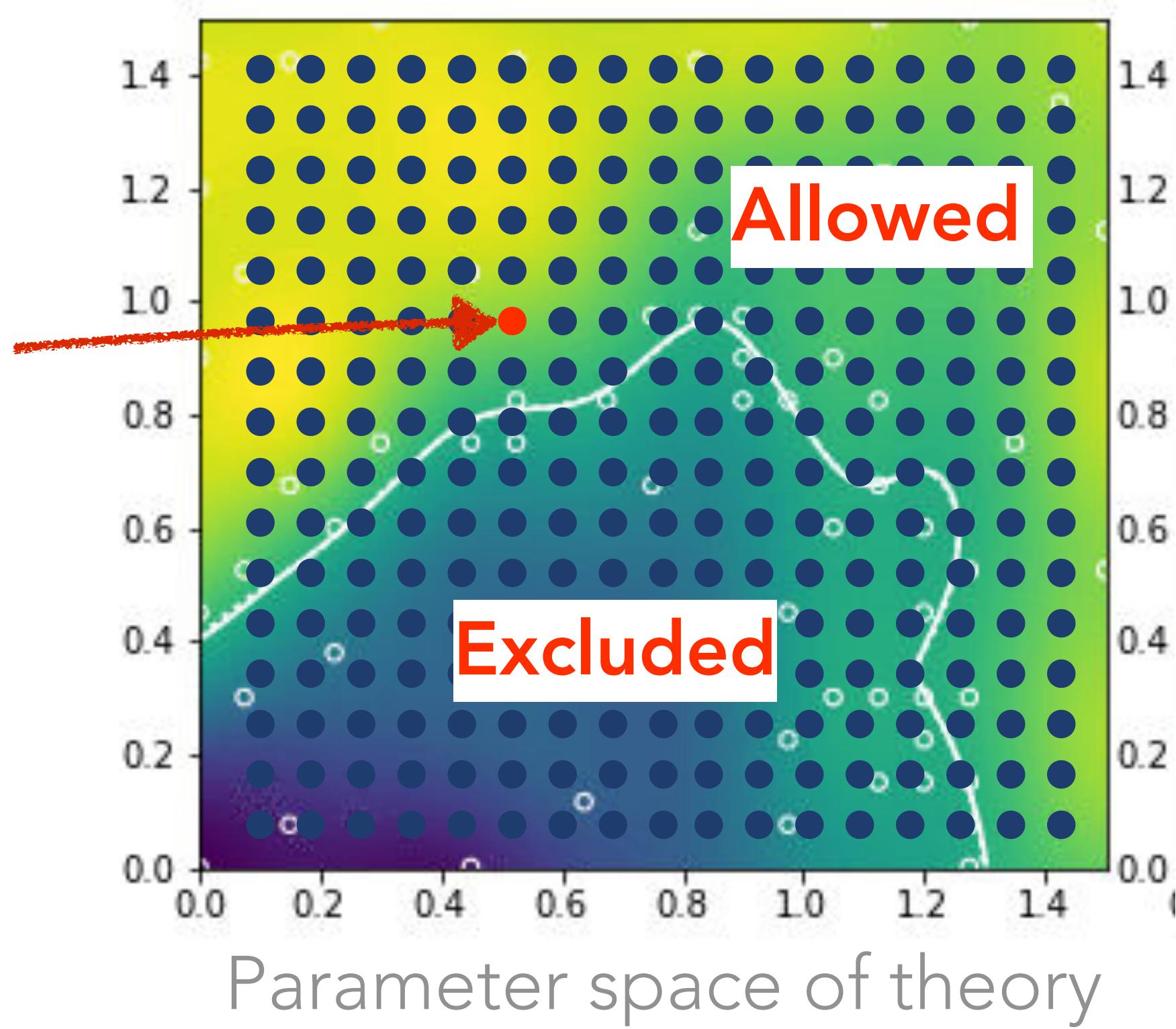
Reinterpretation with RECAST

Basic approach:

- scan parameter space of theory, simulate signal for each point
- execute complex workflow that implements analysis for each parameter point
- determine which regions of parameter space are excluded



Complex computational workflow

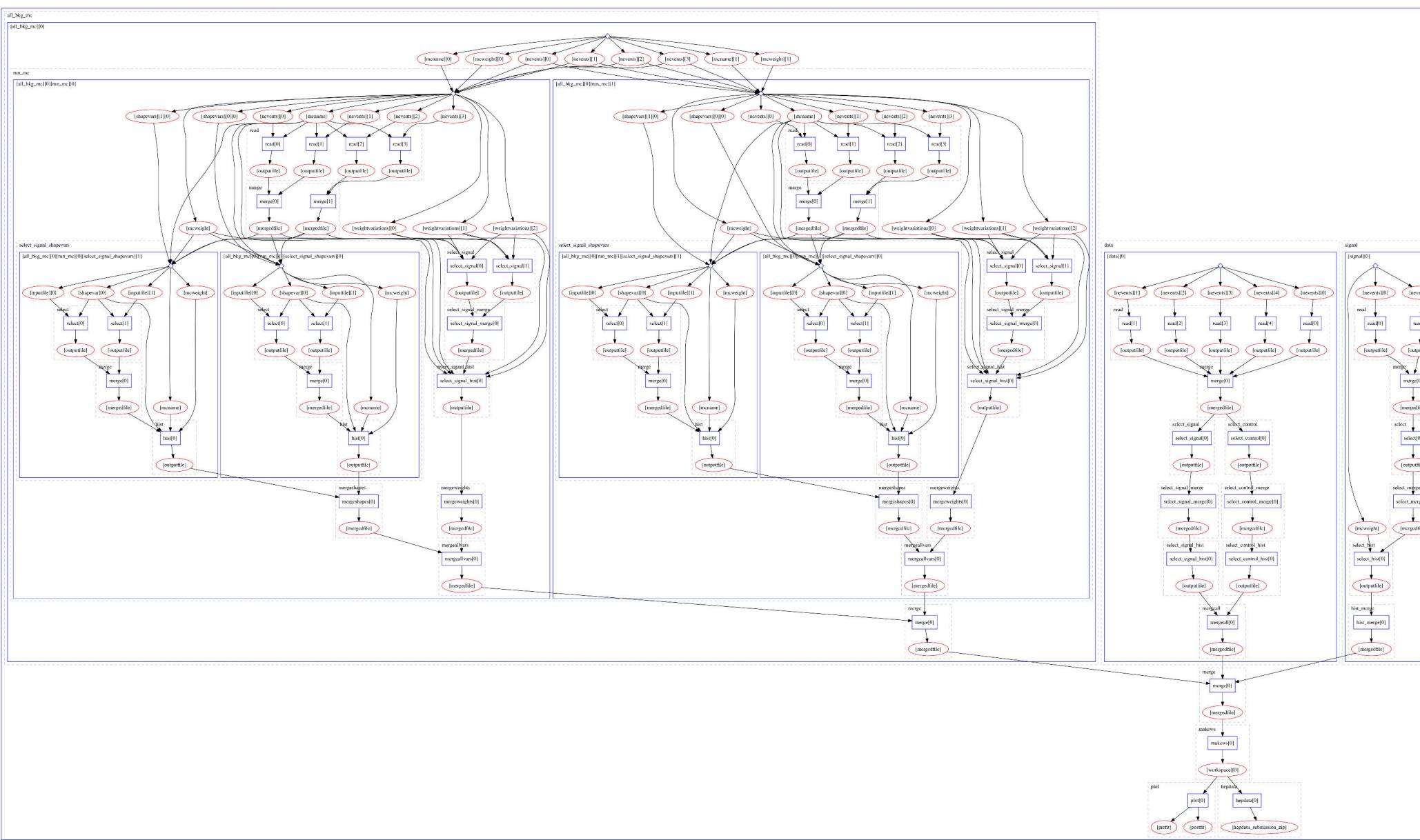


Parameter space of theory

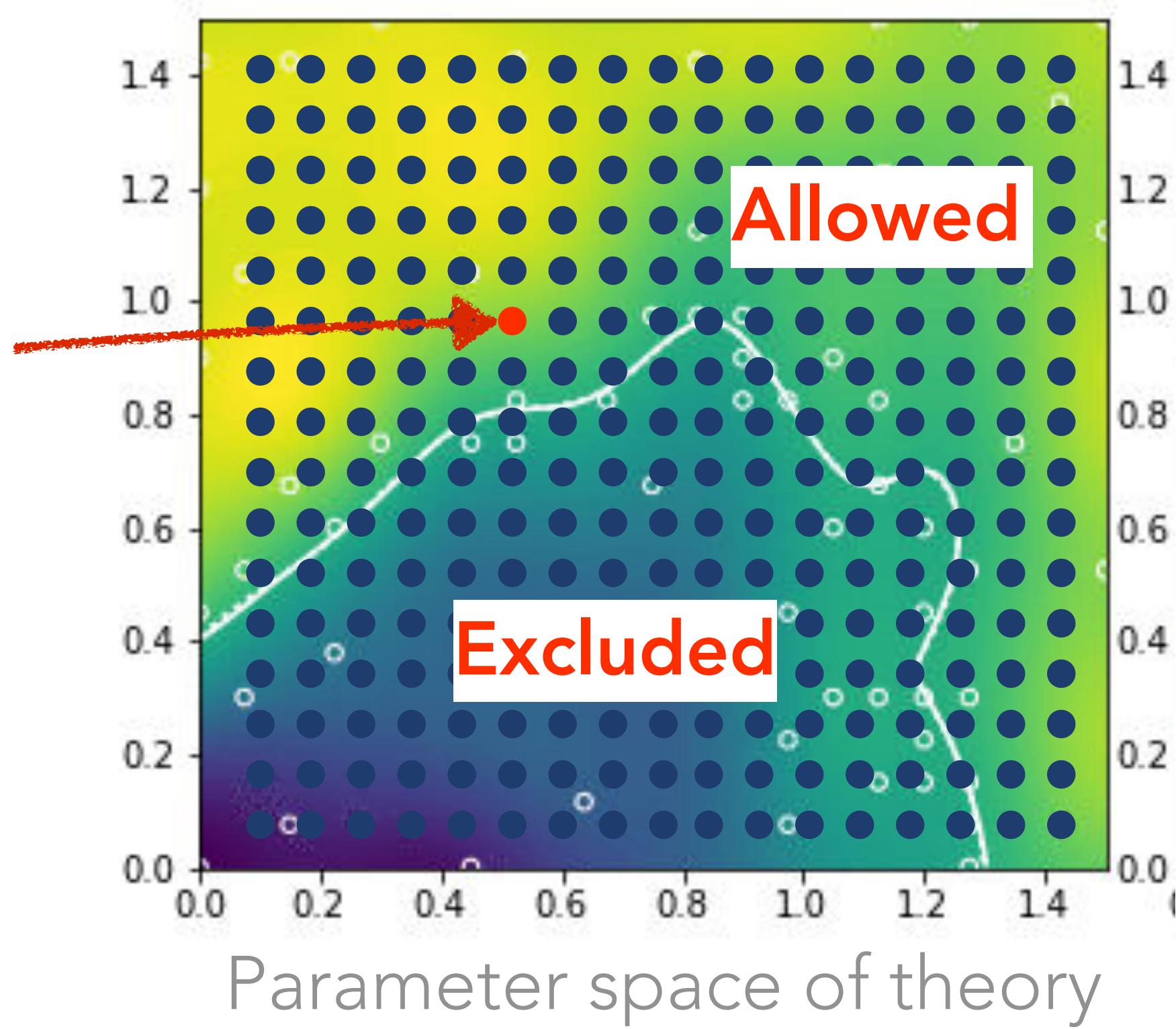
Reinterpretation with RECAST

Basic approach:

- scan parameter space of theory, simulate signal for each point
- execute complex workflow that implements analysis for each parameter point
- determine which regions of parameter space are excluded



Complex computational workflow



Parameter space of theory

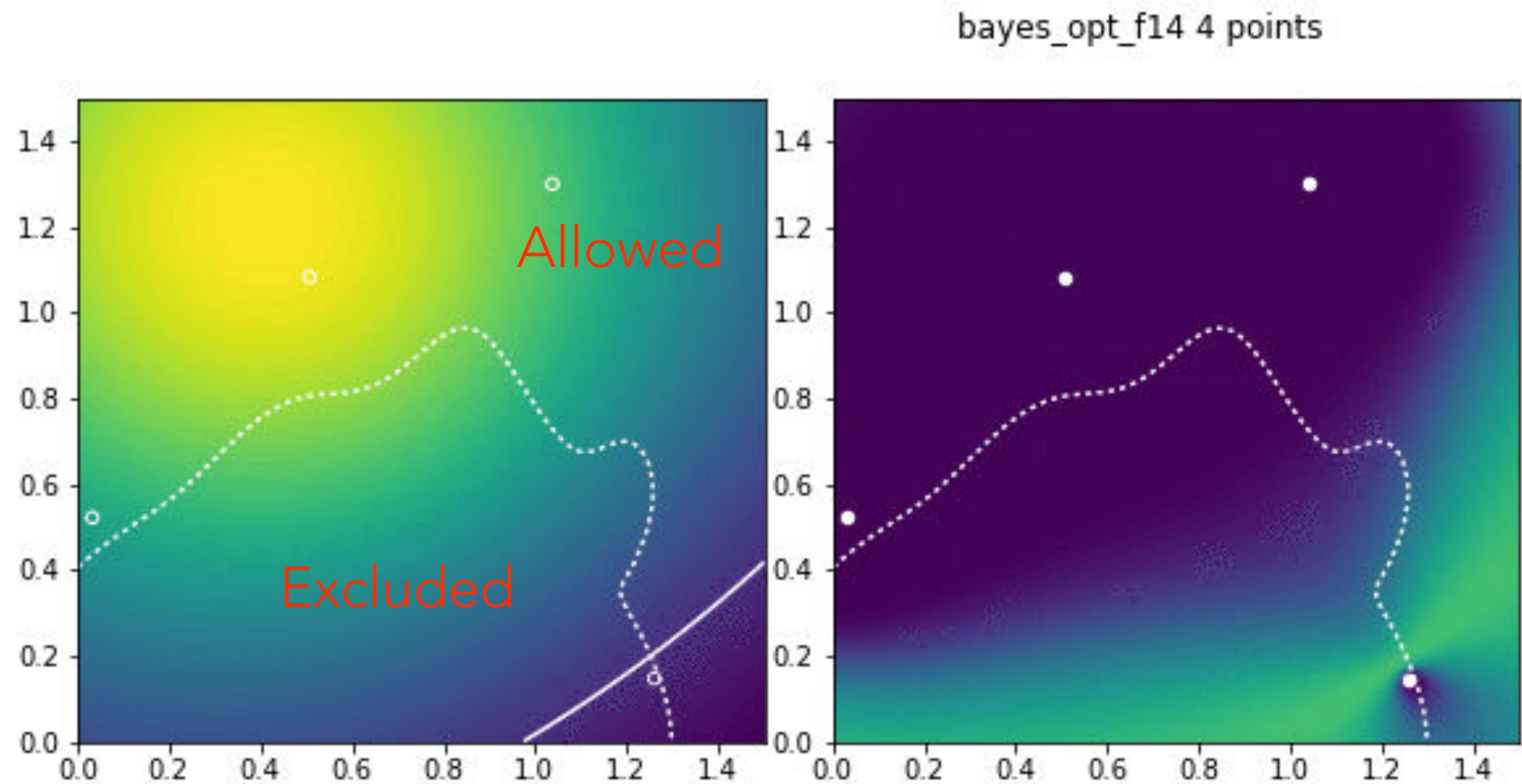
Active Learning for Reinterpretation

Instead of simulating on a grid a priori, simulate on demand where it is relevant!

Drastically more efficient use of computing resources

Changes traditional relationship between production system & analysis

Active learning algorithm steers enormous computational effort



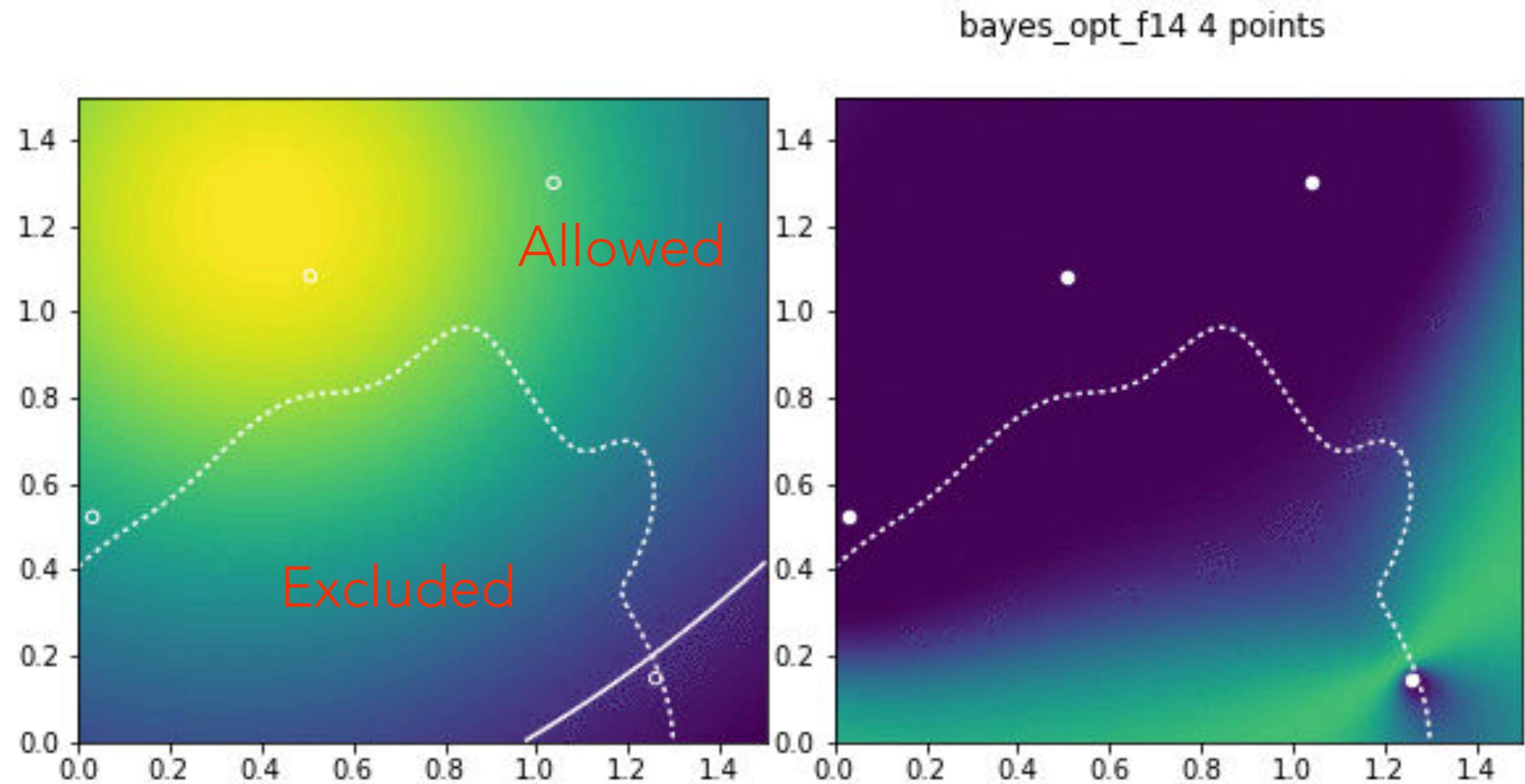
Active Learning for Reinterpretation

Instead of simulating on a grid a priori, simulate on demand where it is relevant!

Drastically more efficient use of computing resources

Changes traditional relationship between production system & analysis

Active learning algorithm steers enormous computational effort

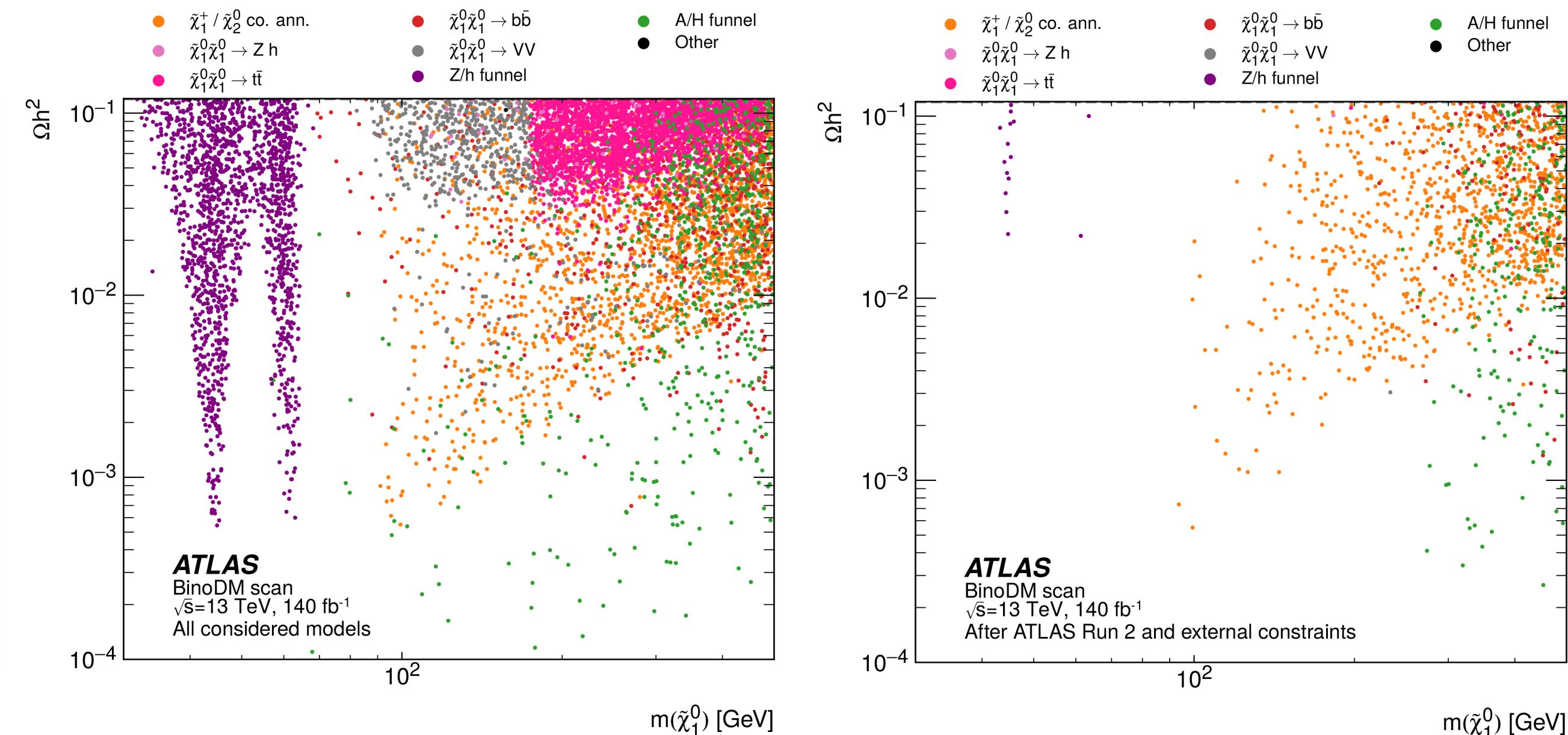
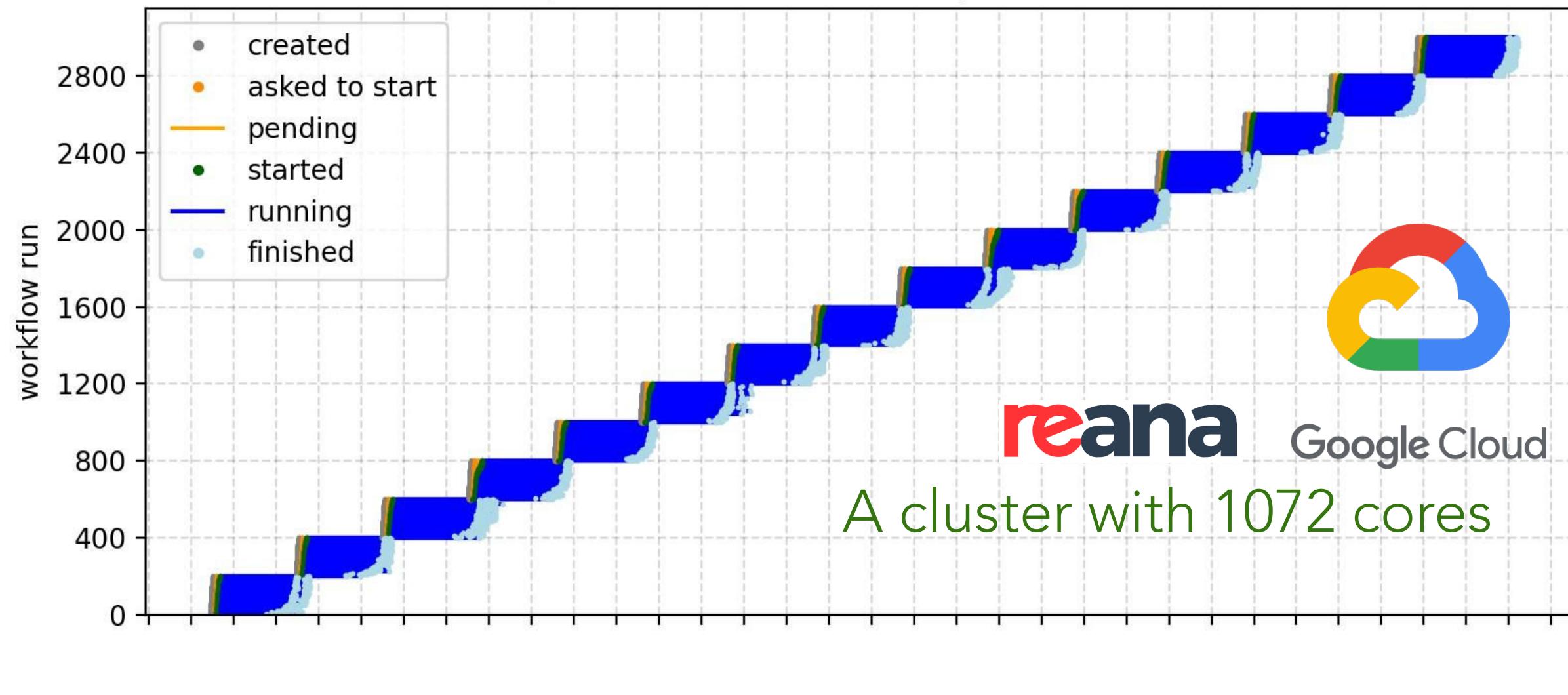


Reuse enables new science

Scalable ATLAS pMSSM computational workflows using containerised REANA reusable analysis platform

Marco Donadoni¹, Matthew Feickert², Lukas Heinrich³, Yang Liu⁴, Audrius Mečionis¹, Vladyslav Moisieienkov¹, Tibor Šimko^{1,*}, Giordon Stark⁵, and Marco Vidal García¹

susy201816-submit200-sleep1140-total3000



(a) All considered models

(d) Models not excluded by ATLAS or external constraints

"This result was a tour de force from ATLAS's supersymmetry physics group, pulling together results from eight separate ATLAS searches using data collected during Run 2 of the LHC and evaluating tens of thousands of supersymmetric models in the 19-dimensional parameter space to set new, more restrictive constraints on models with supersymmetric dark matter particles. " — Matthew Feickert

Do we really need another workflow language?

Probably not...

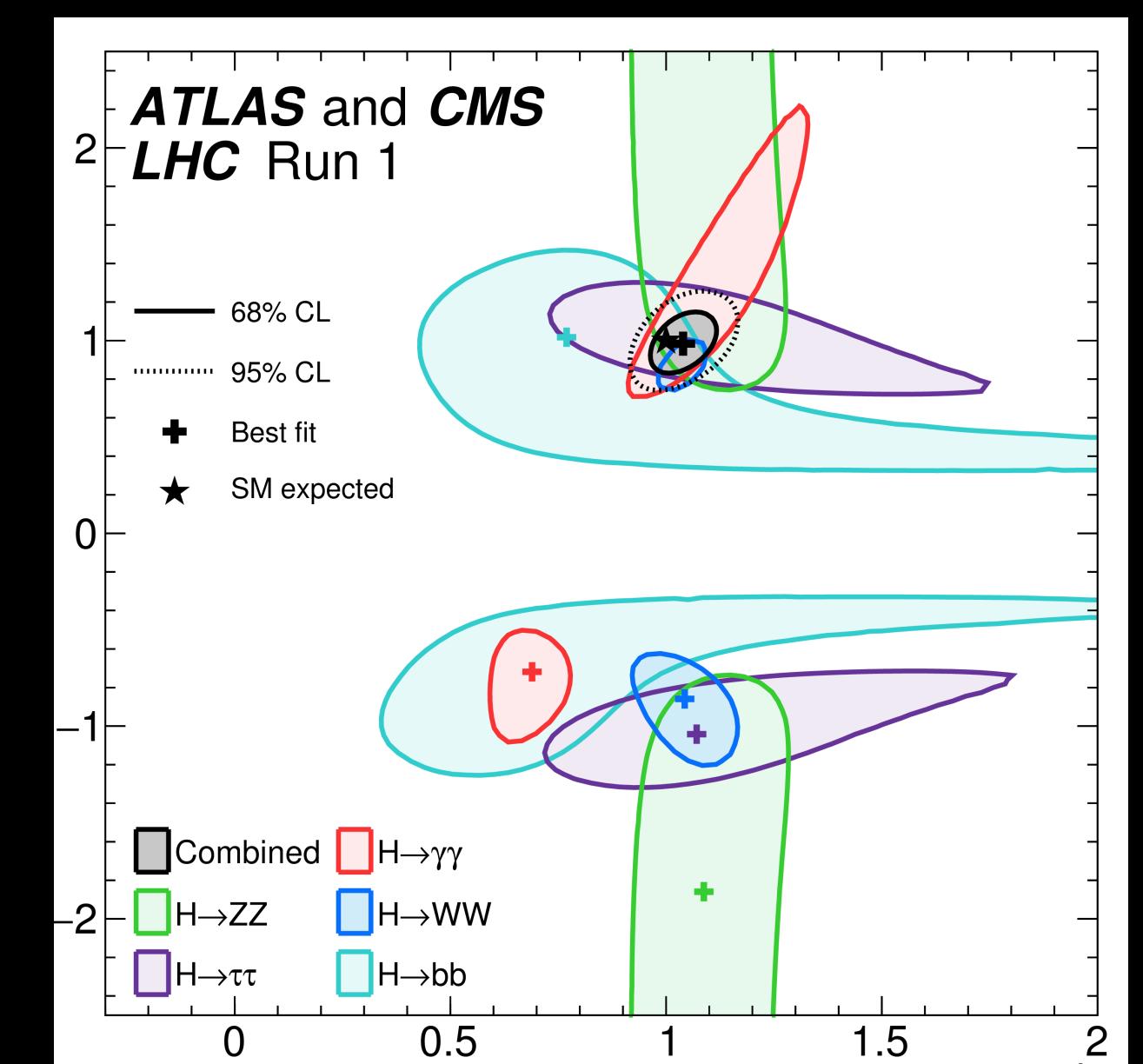
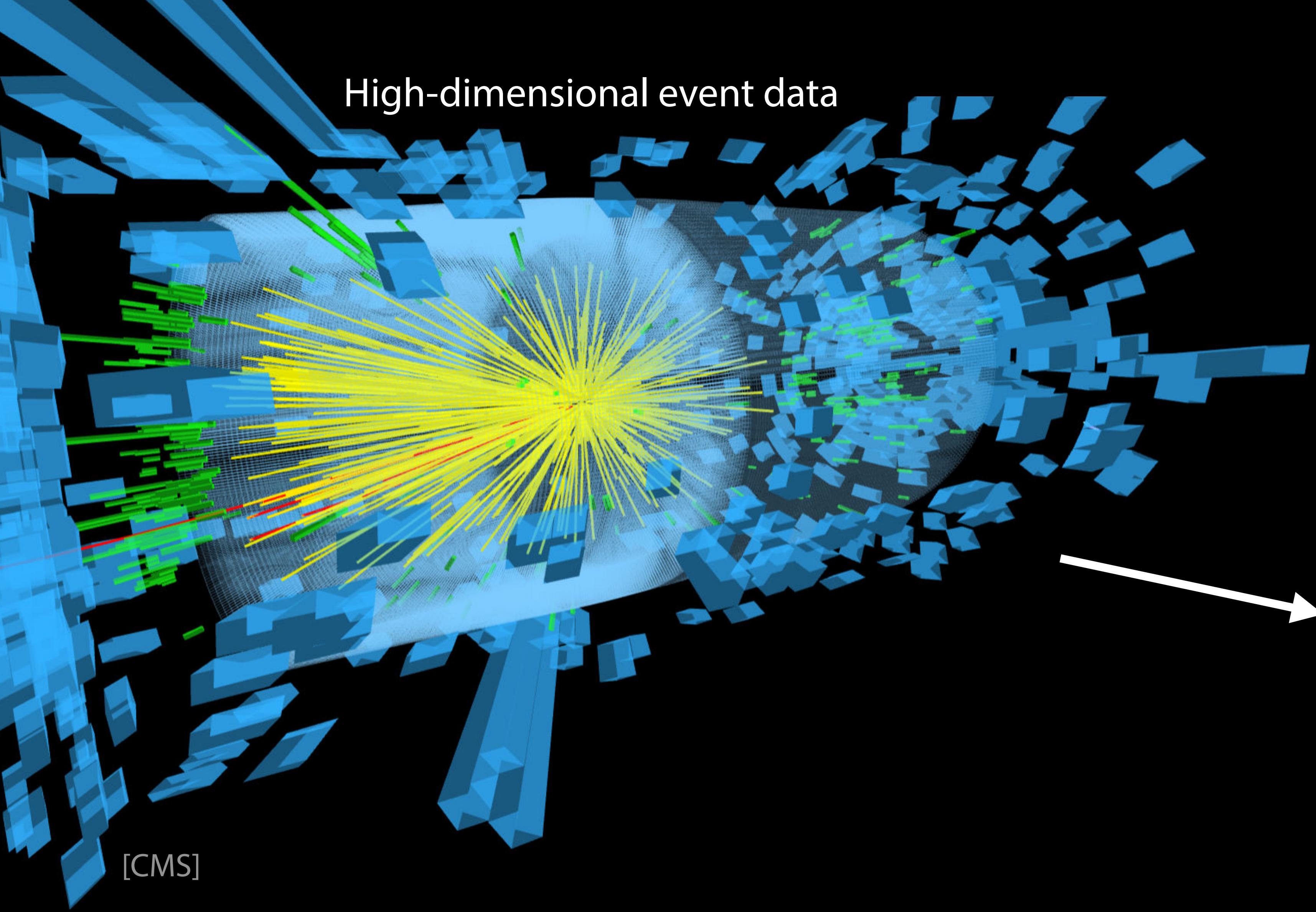
- We brought together experts from different workflow languages to understand the state of the art and needs of particle physics community



COMMON
WORKFLOW
LANGUAGE



Simulation-based Inference



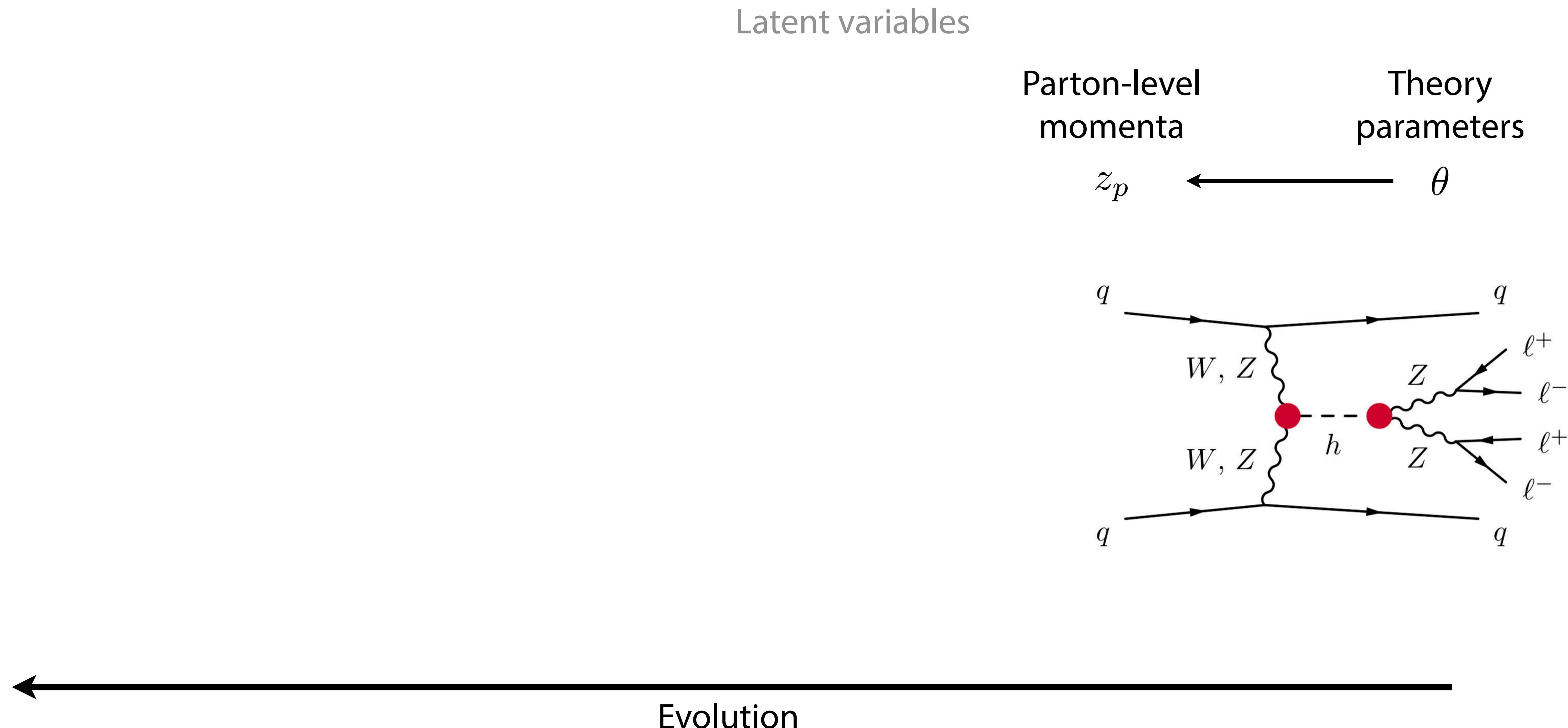
Precision constraints on
new physics

Simulating particle physics processes

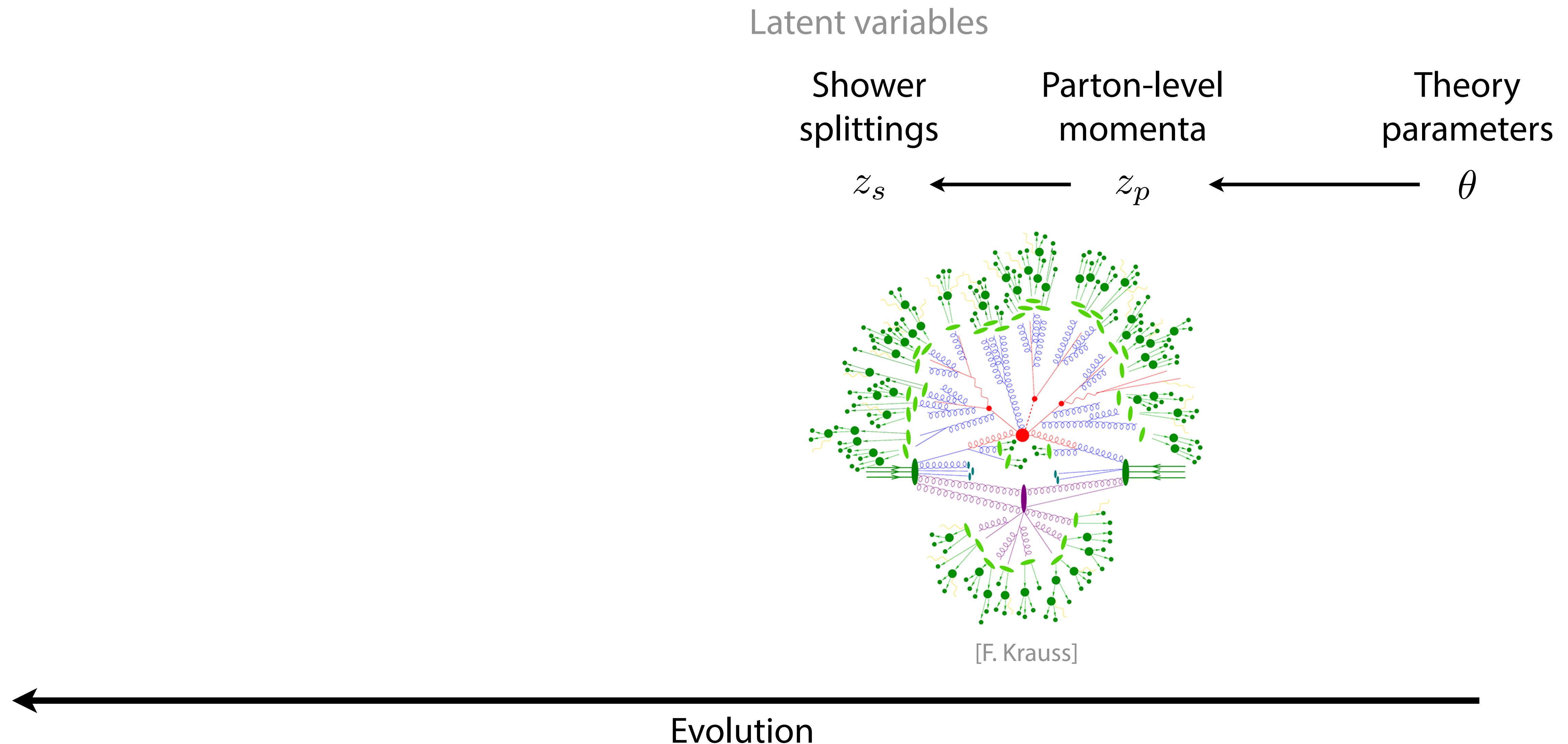
Theory
parameters
 θ



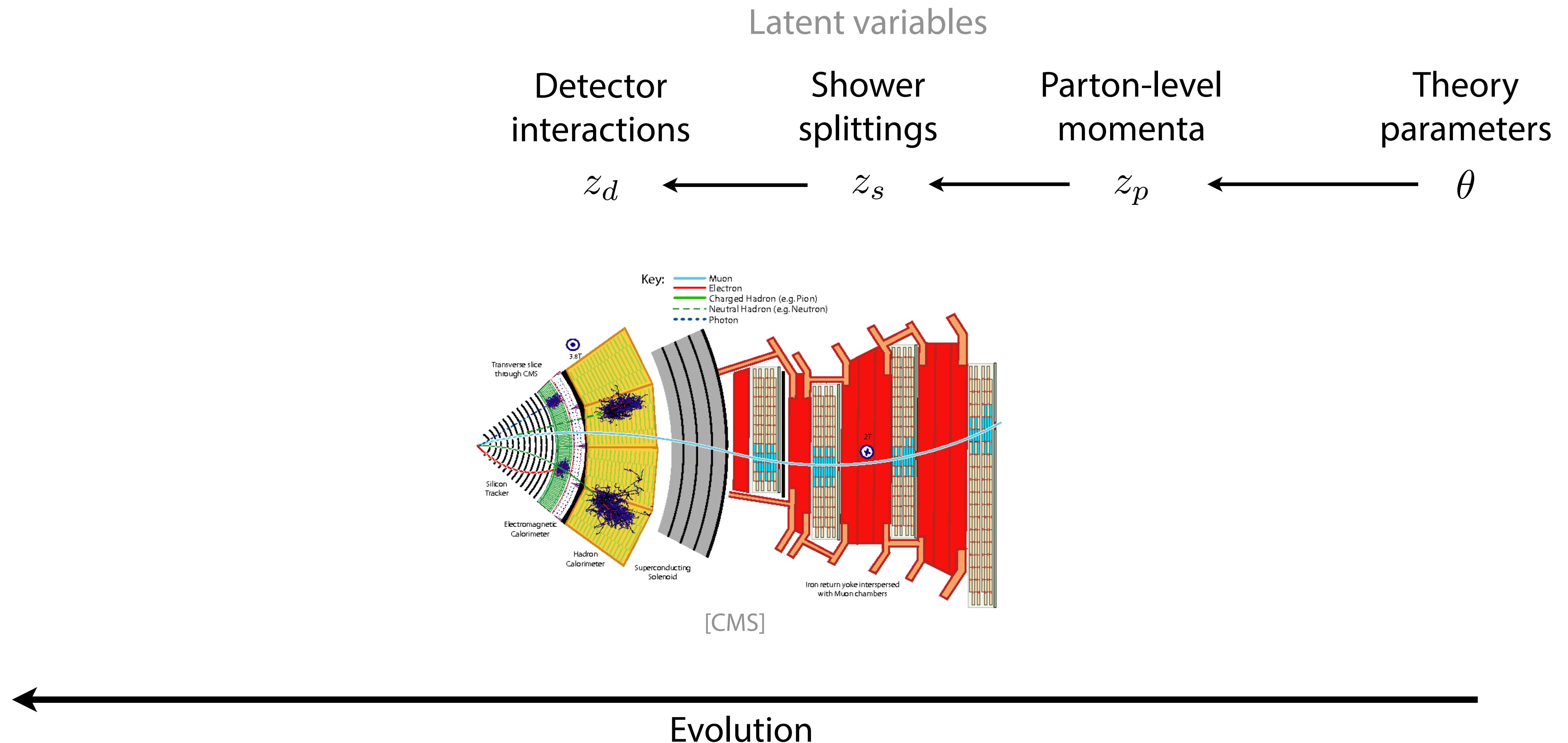
Simulating particle physics processes



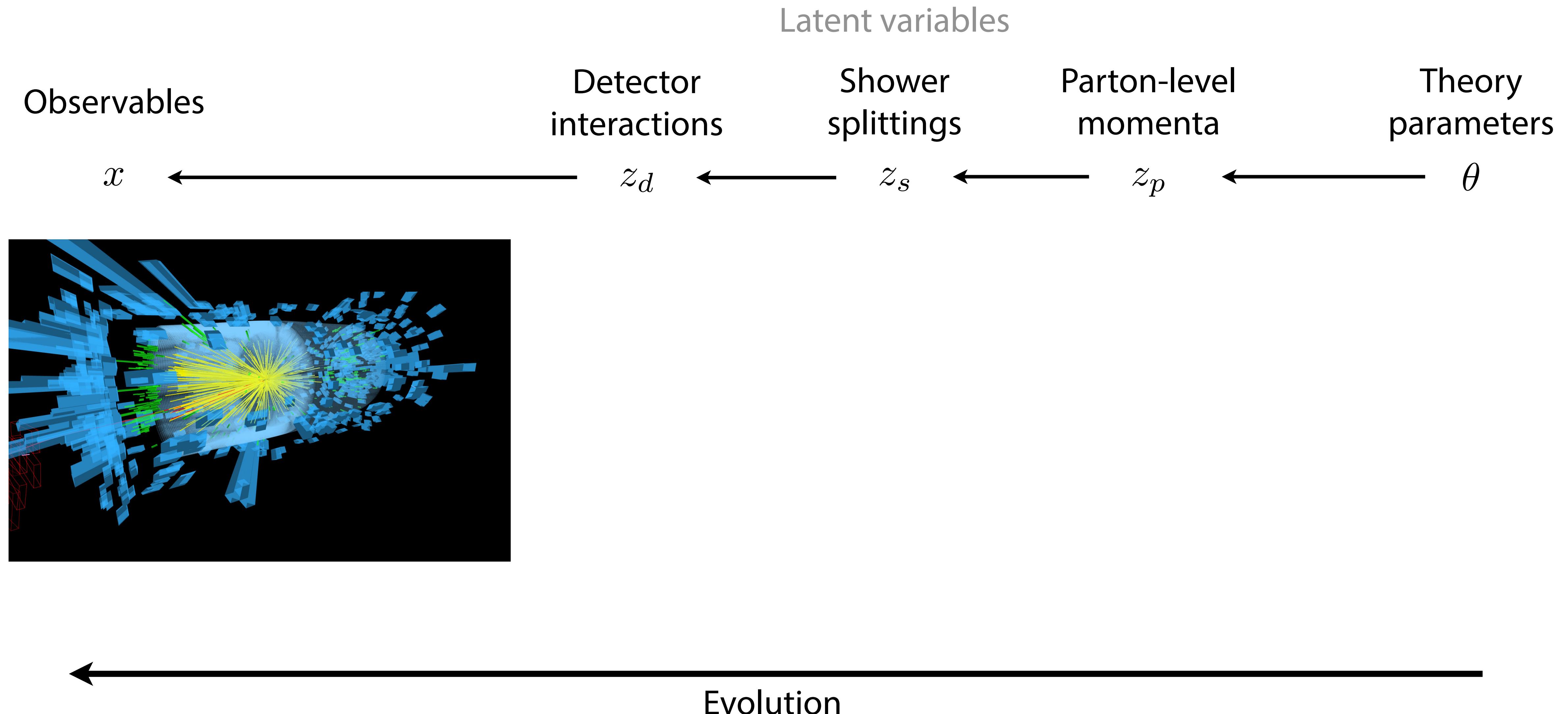
Simulating particle physics processes



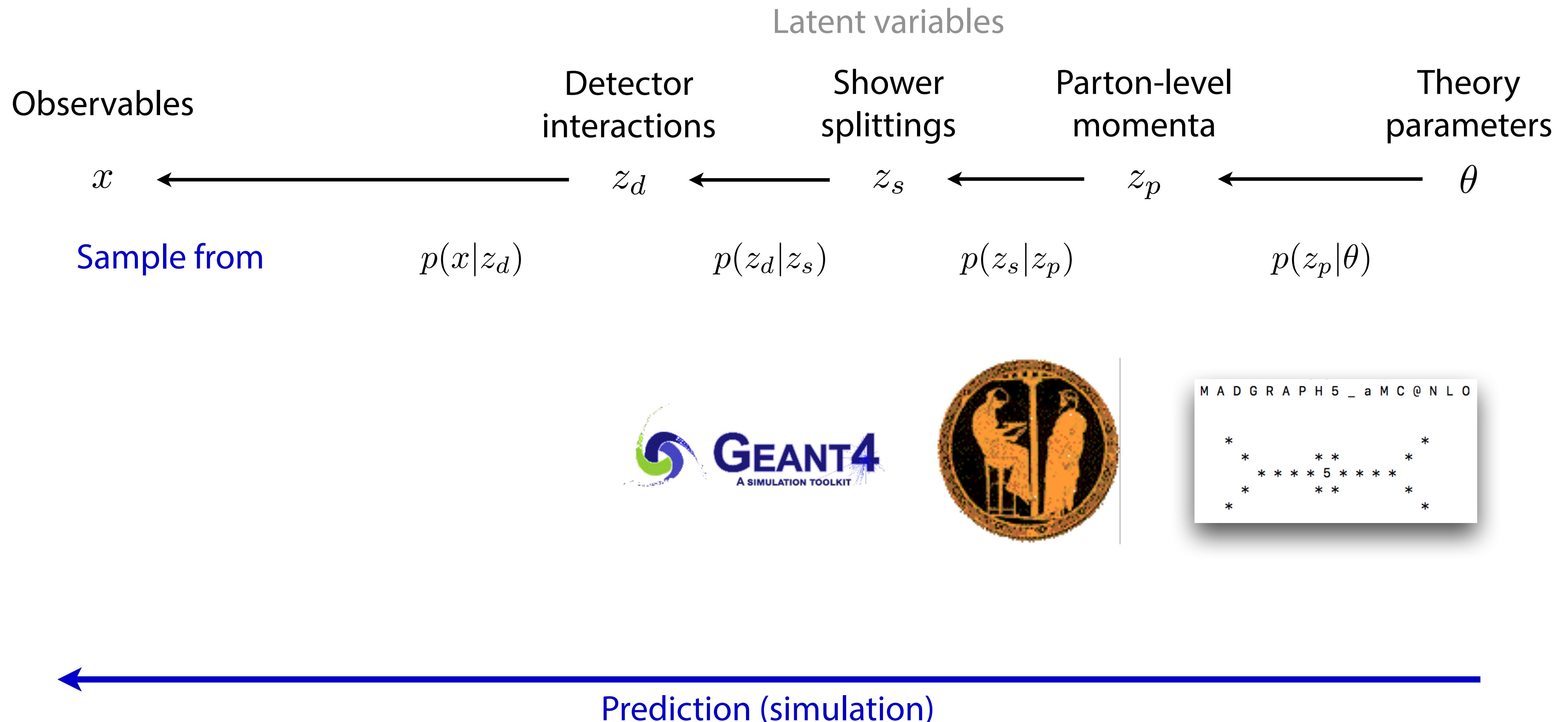
Simulating particle physics processes



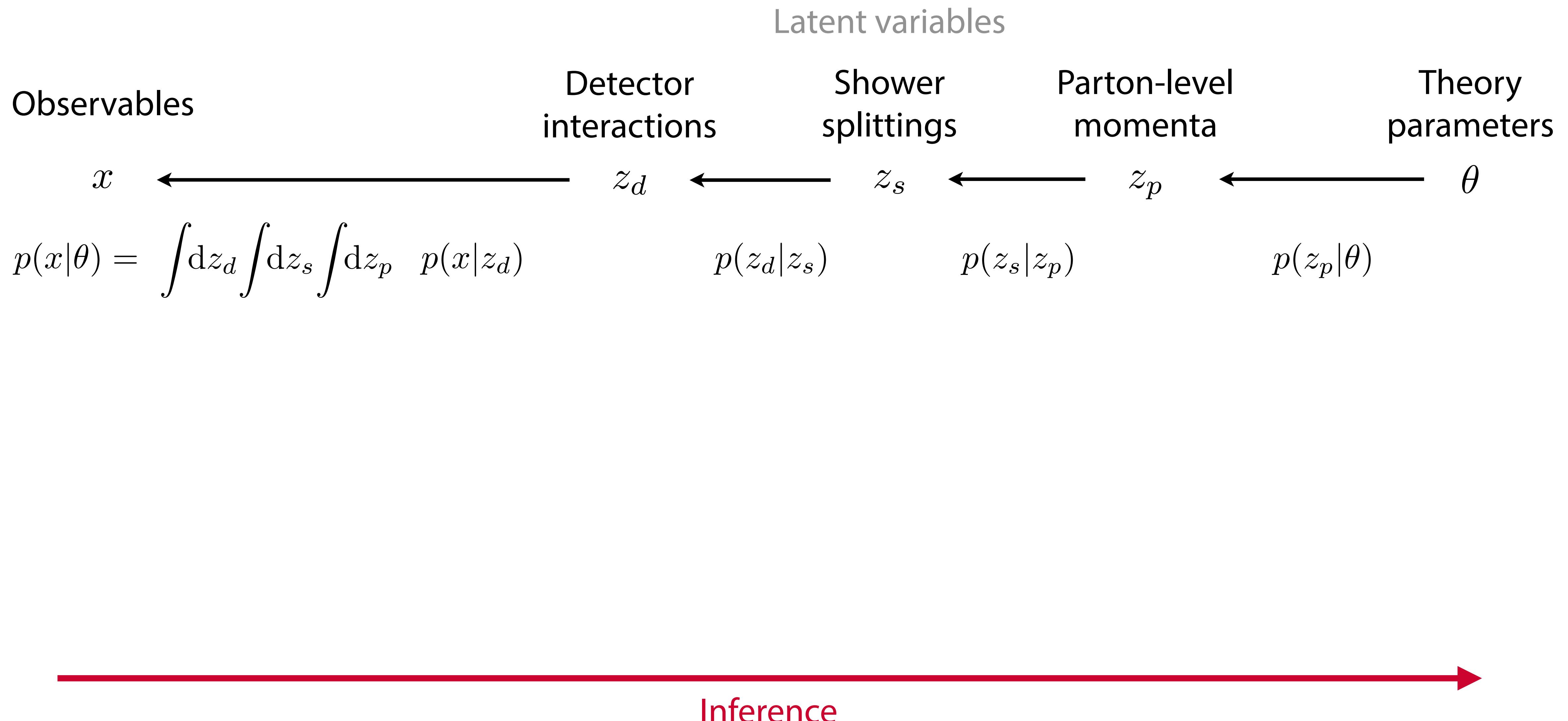
Simulating particle physics processes



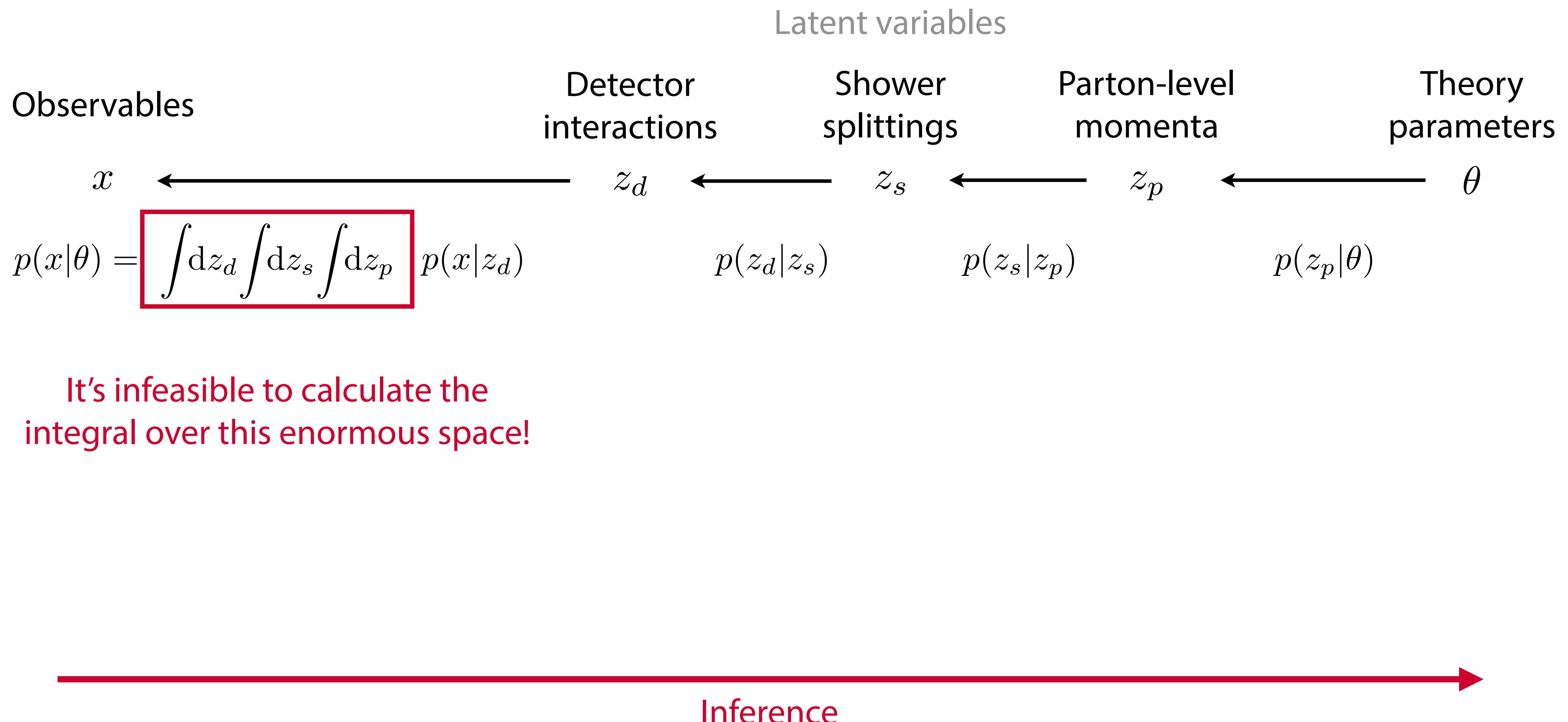
Simulating particle physics processes



Simulating particle physics processes



Simulating particle physics processes



Simulation-based Inference

This motivates a class of inference methods for a stochastic simulator where

- evaluating the **likelihood is intractable**, but
- it is **possible to sample** synthetic data $x \sim p(x | \theta)$

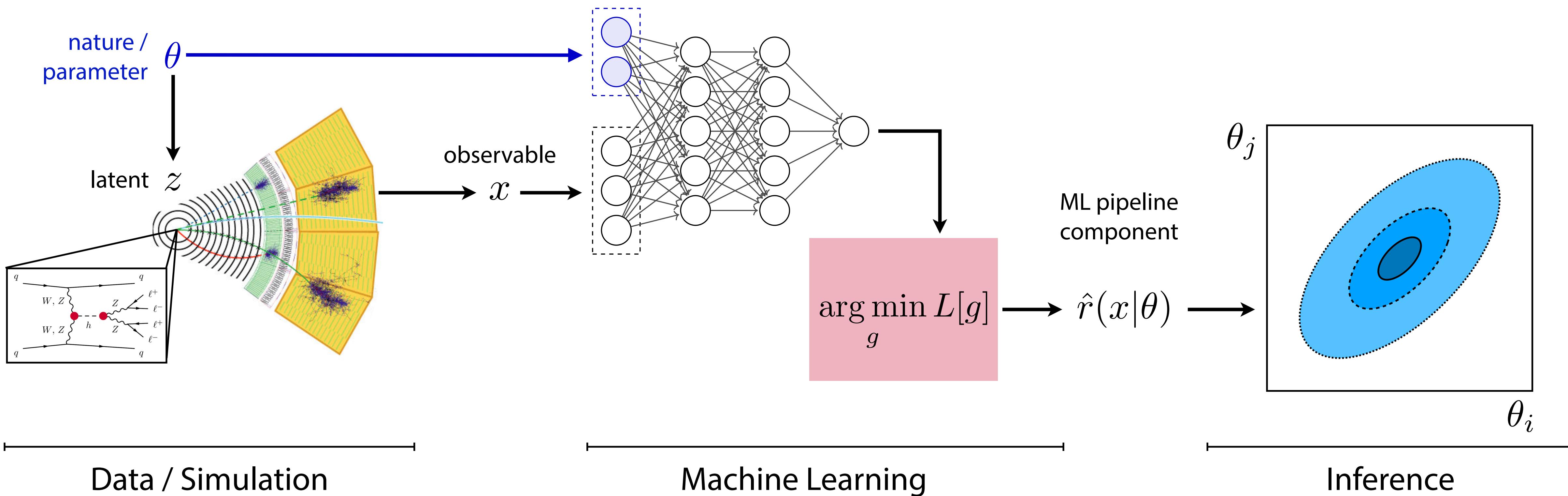
This setting is often referred to as **likelihood-free inference**, but I prefer the term **simulation-based inference** because usually one approximates the likelihood (or likelihood ratio) and then use established inference techniques

- applies to both Bayesian or Frequentist inference

Simulation-Based Inference

Deep learning and neural density estimation are effective at learning approximate surrogates for the likelihood and posterior, **revolutionizing principled statistical inference in science!**

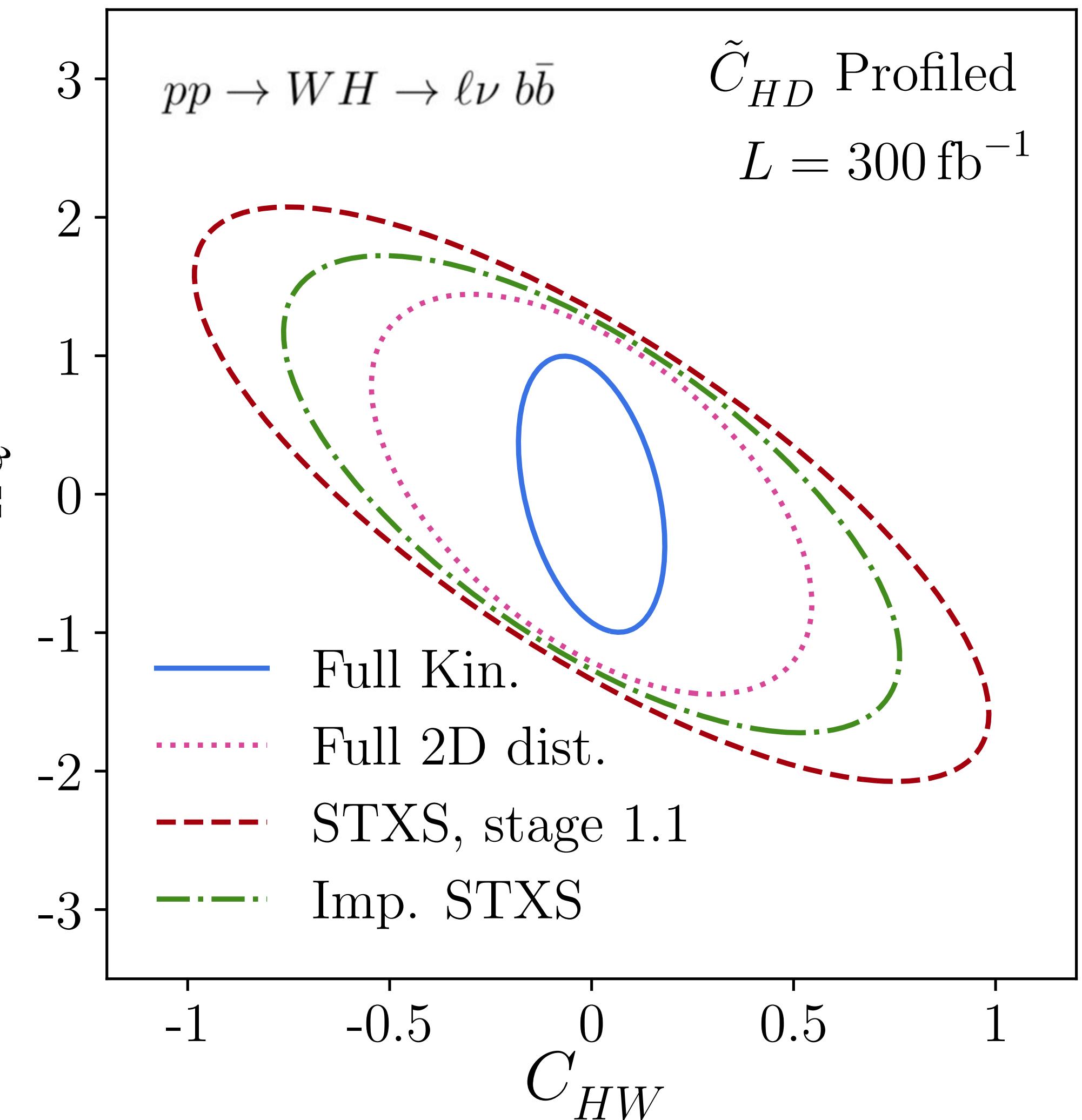
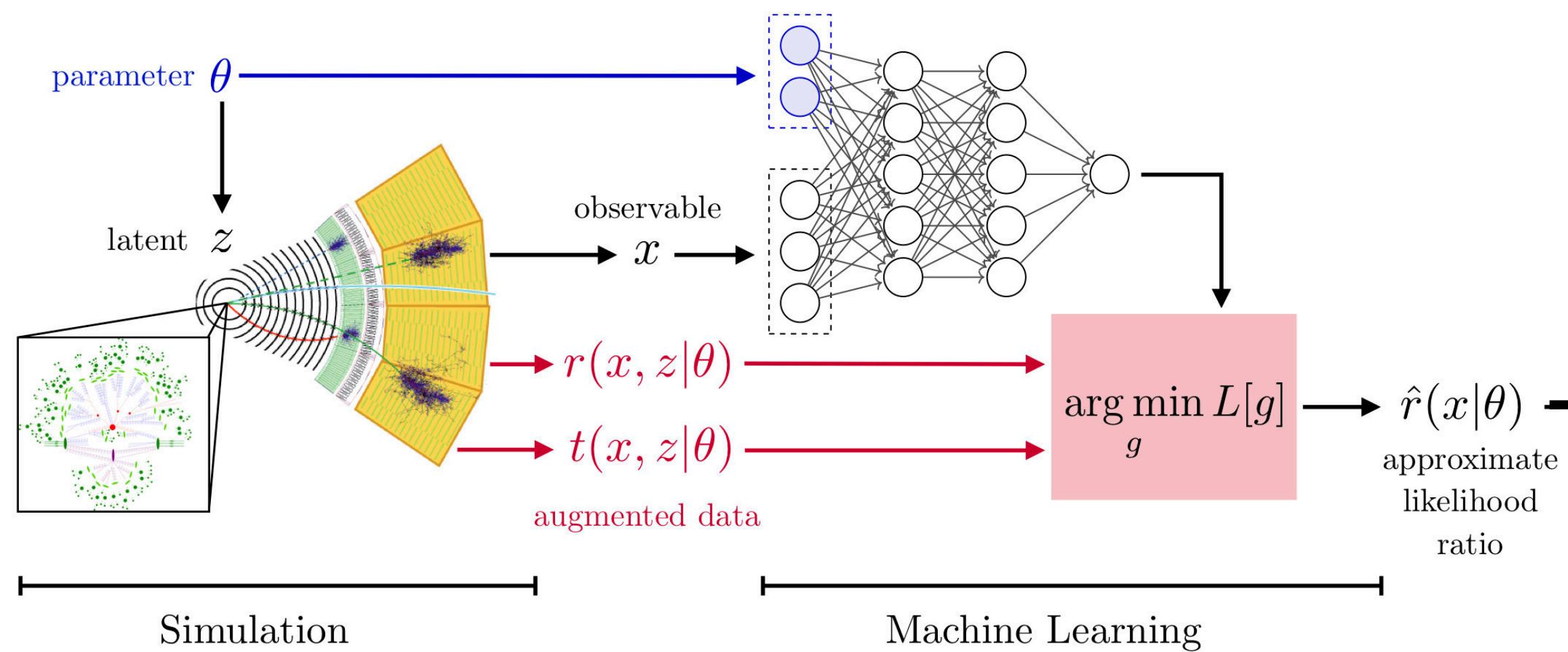
- Removes the need for hand-engineered summary statistics that sacrifice power
- Expert knowledge in simulator is transferred to surrogate via learning



Impact on Studies of The Higgs Boson

Potential for massive gains in precision of a flagship measurement at the LHC !

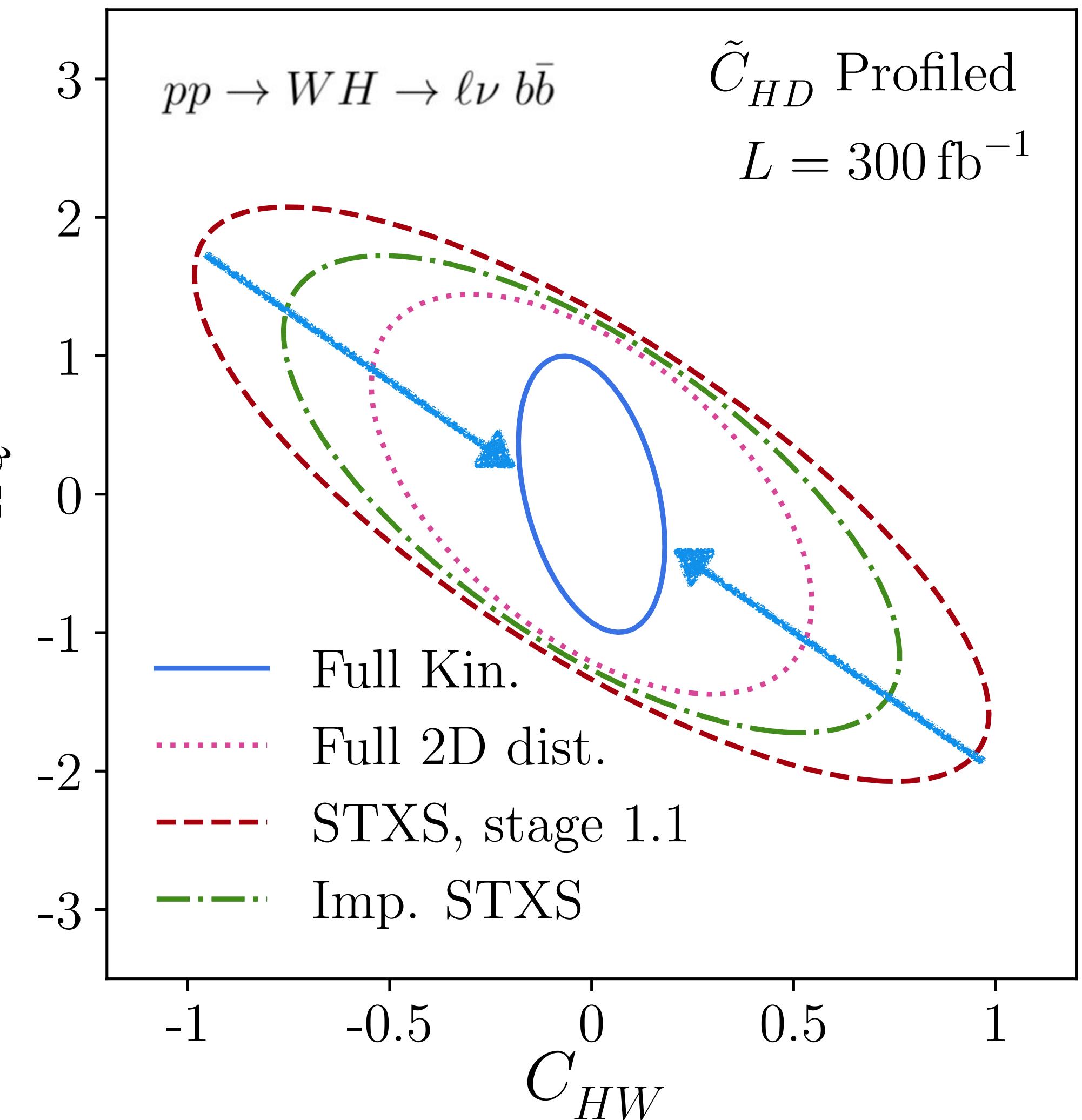
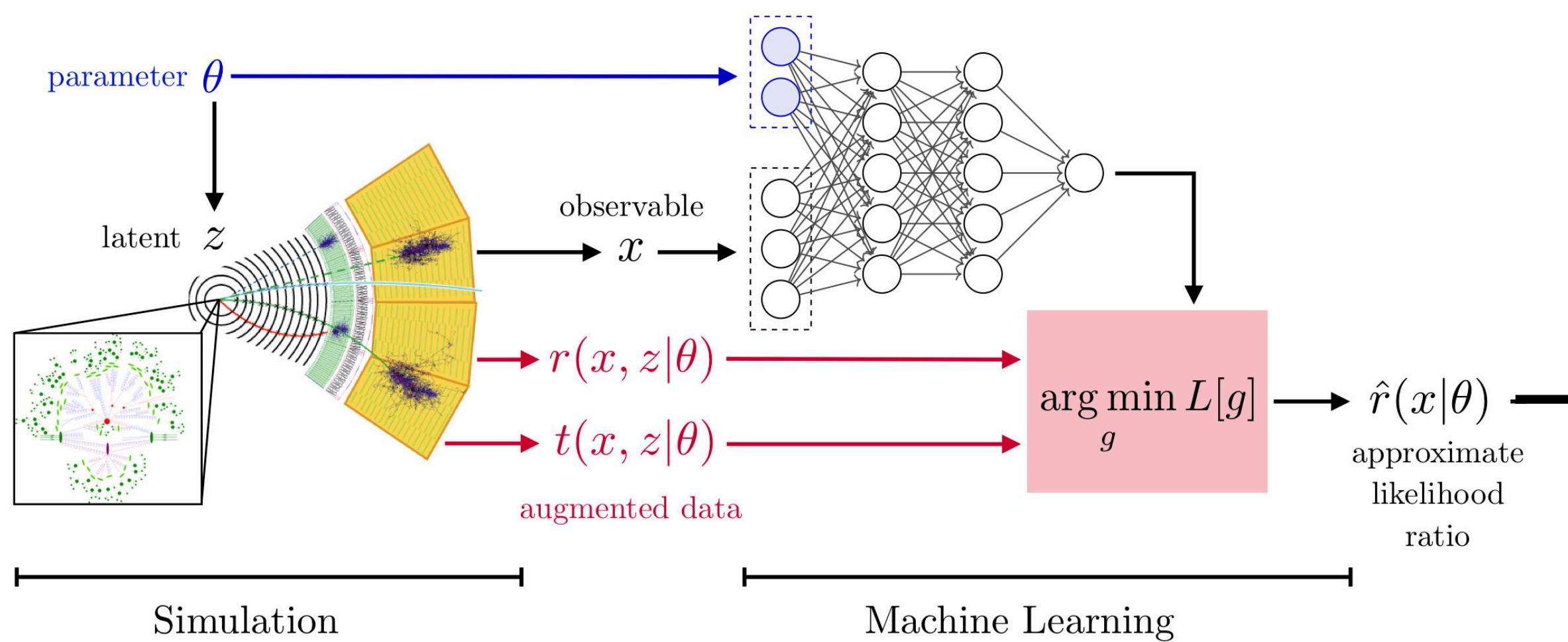
Equivalent increasing data collected by LHC by several factors



Impact on Studies of The Higgs Boson

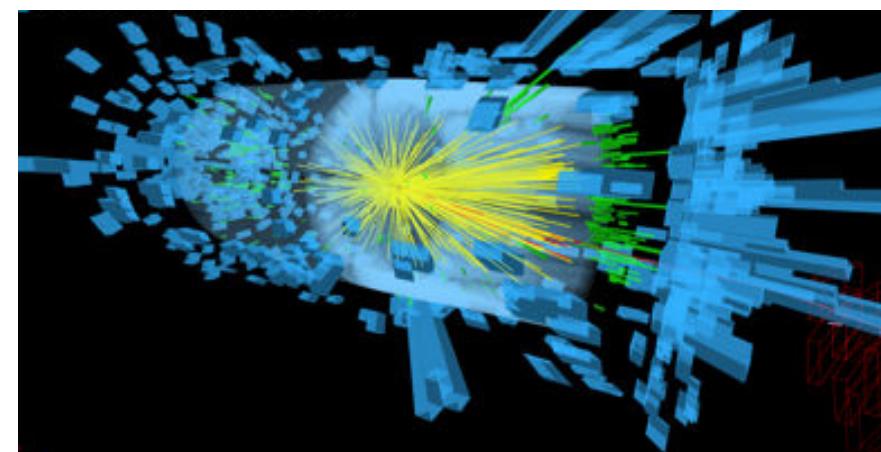
Potential for massive gains in precision of a flagship measurement at the LHC !

Equivalent increasing data collected by LHC by several factors

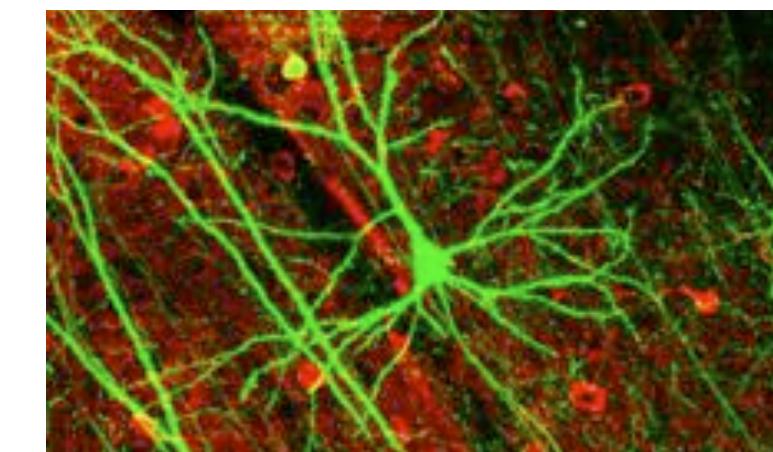


Science is replete with high-fidelity simulators

Simulation-based inference is a major evolution in statistical inference for science!



Particle
colliders



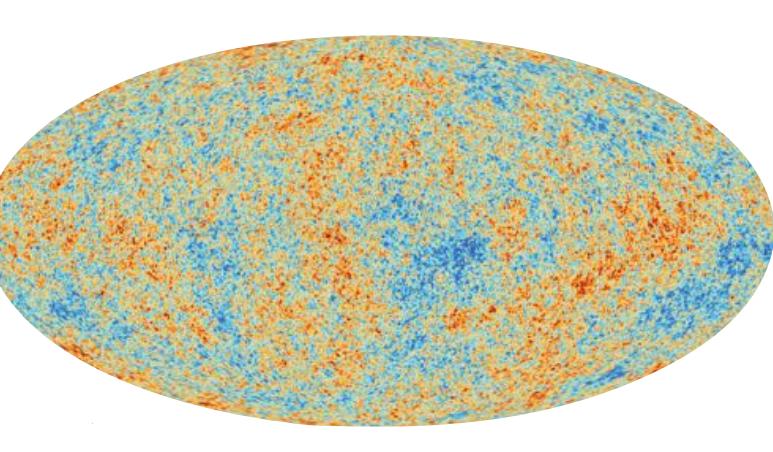
Neuron
activity



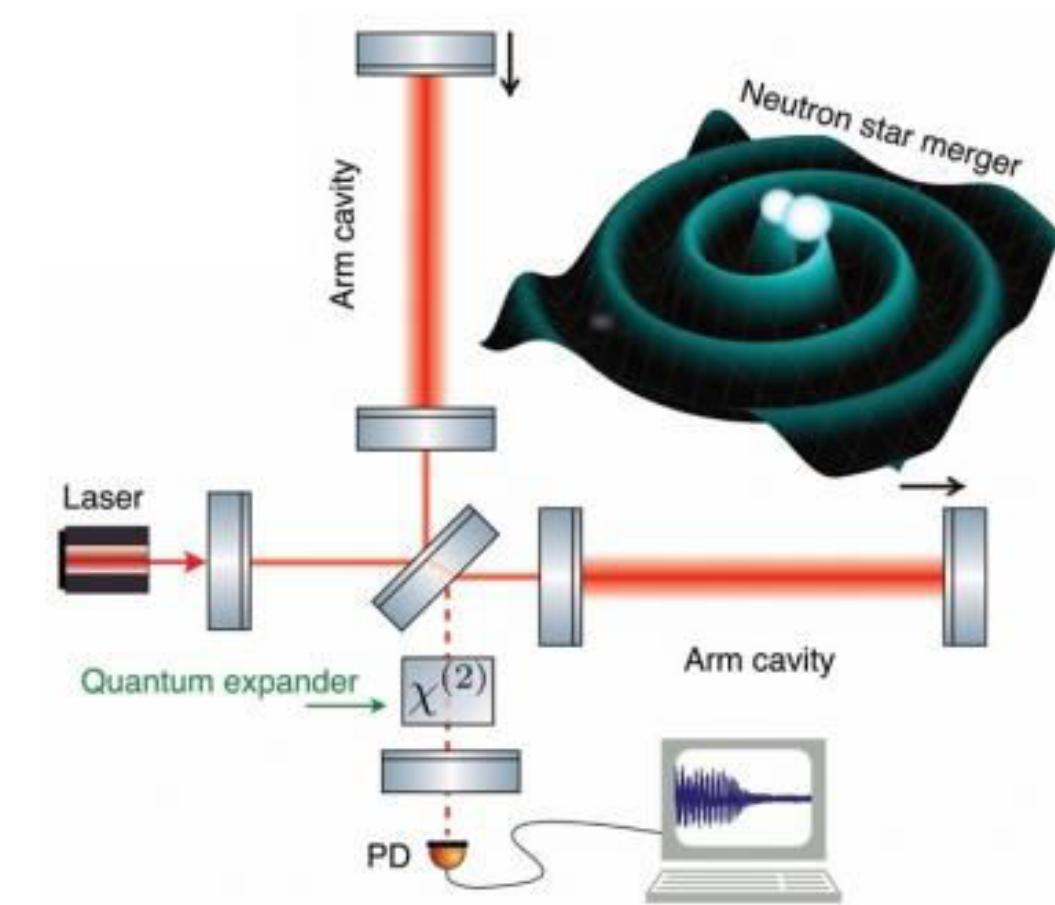
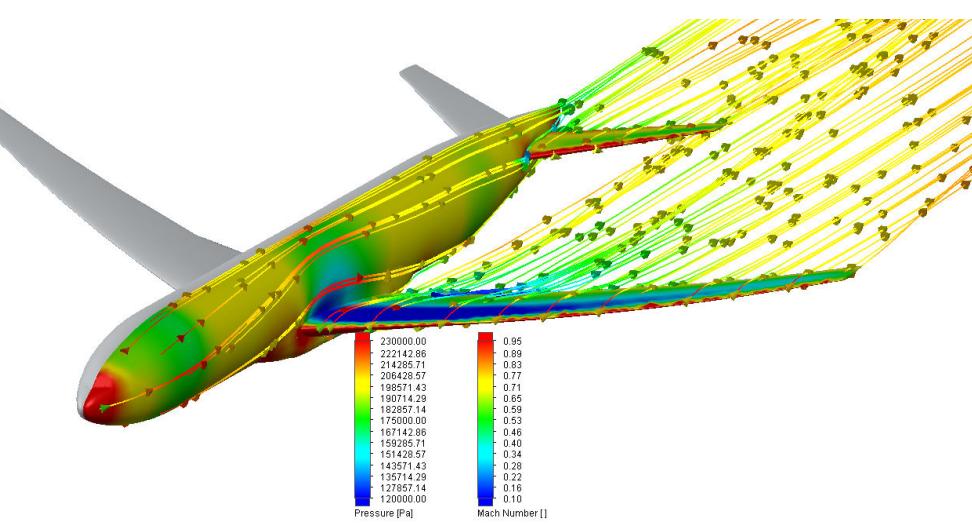
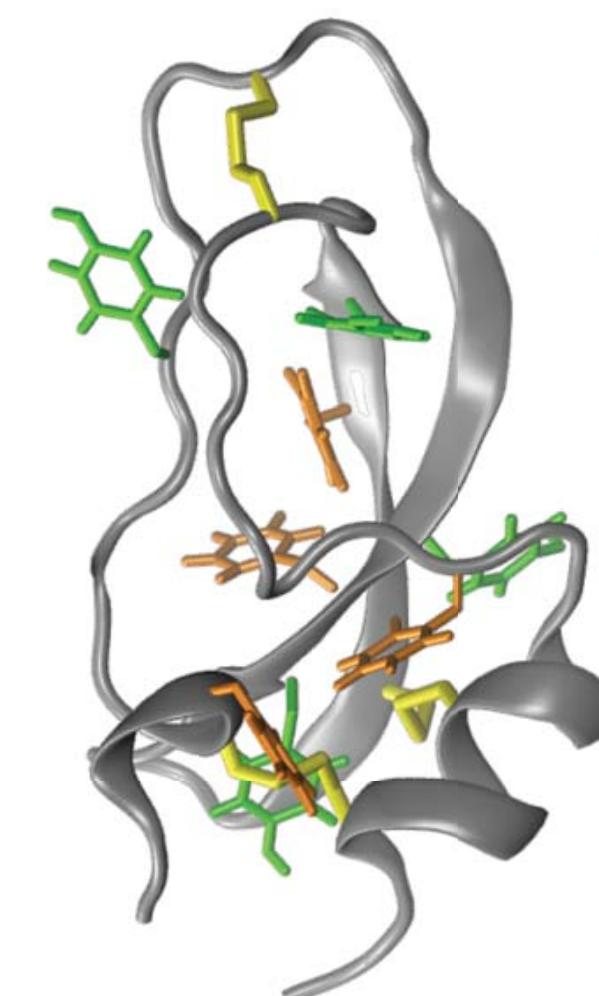
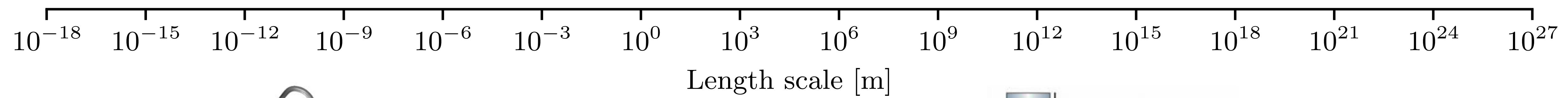
Epidemics



Gravitational
lensing



Evolution of
the Universe



Papers

The list is automatically compiled each day. Should you observe any inaccuracies or concerns, kindly [bring them to our attention](#).

Additionally, if you believe a new paper aligns with the topic, feel free to [submit it](#).

[Visualize the annual growth in the number of publications.](#)

 Sort by Category

 Sort by Year

 Sort by

Total (840)
Statistics (208)
Computer Science (116)

Astrophysics (100)
Mathematics (58)
Education (53)
Physics (48)
Economics (46)
Quantitative Biology (32)

Neuroscience (30)
Quantitative Finance (21)
Astronomy (18)
Engineering (14)

Genetics (13)
Epidemiology (11)
Medicine (11)
Geography (8)

Social Science (7)
Ecology (6)
Evolutionary biology (6)

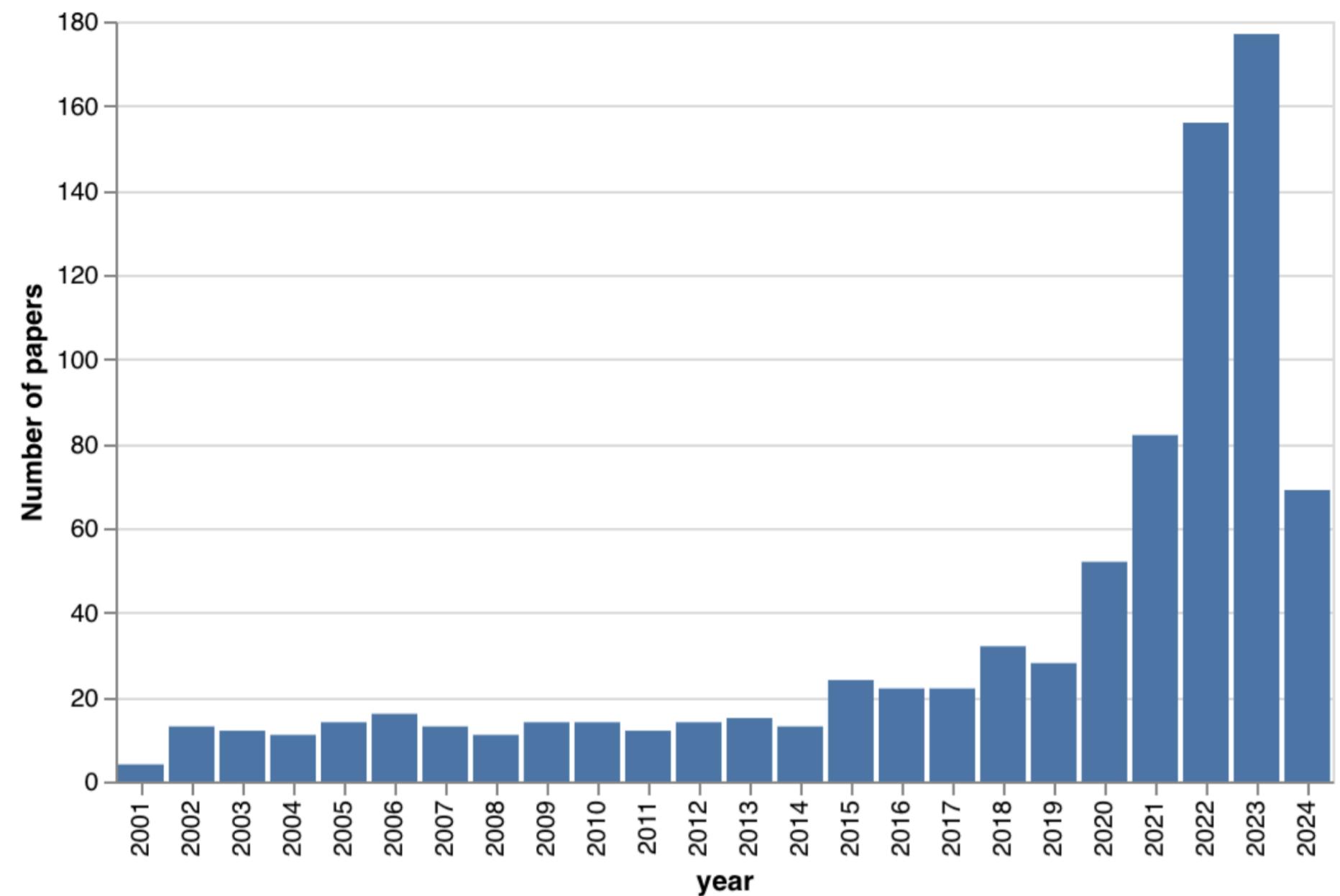
Environmental Science (4)
Cognitive Science (4)
Robotics (4)
Systems biology (4)

Electrical Engineering and
Computer Science (3)

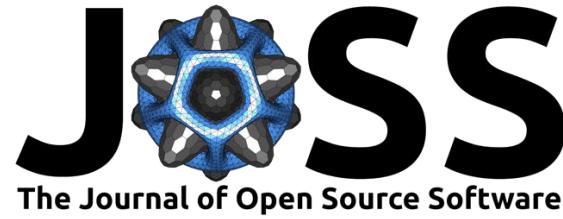
Statistics

- Addressing Misspecification in Simulation-based Inference through Data-driven Calibration, Gamella, O Sener, J Behrmann... - arXiv preprint arXiv ..., 2024 - arxiv.org
- Modelling Sampling Distributions of Test Statistics with Autograd, AA Kadhim, HB Prosper - arXiv preprint arXiv:2405.02488, 2024 - arxiv.org
- Preconditioned Neural Posterior Estimation for Likelihood-free Inference, X Wang, RP Kelly, DJ Warne, C Drovandi - arXiv preprint arXiv ..., 2024 - arxiv.org
- A variational neural Bayes framework for inference on intractable posterior distributions, E Maceda, EC Hector, A Lenzi, BJ Reich - arXiv preprint arXiv:2404.10899, 2024 - arxiv.org
- Increased perceptual reliability reduces membrane potential variability in cortical neurons, B von Hünerbein, J Jordan, M Oude Lohuis... - bioRxiv, 2024 - biorxiv.org
- How much information can be extracted from galaxy clustering at the field level?, NM Nguyen, F Schmidt, B Tucci, M Reinecke... - arXiv preprint arXiv ..., 2024 - arxiv.org
- Evolution of Analysis Techniques and Statistical Treatment, A Held - Bulletin of the American Physical Society, 2024 - APS
- Simulation-Based Inference with Quantile Regression, H Jia - arXiv preprint arXiv:2401.02413, 2024 - arxiv.org
- Direct Amortized Likelihood Ratio Estimation, AD Cobb, B Matejek, D Elenius, A Roy... - arXiv preprint arXiv ..., 2023 - arxiv.org
- On simulation-based inference for implicitly defined models, J Park - arXiv preprint arXiv:2311.09446, 2023 - arxiv.org
- Machine Learning for Mechanistic Models of Metapopulation Dynamics, J Li, EL Ionides, AA King, M Pascual, N Ning - arXiv preprint arXiv ..., 2023 - arxiv.org
- Inference on spatiotemporal dynamics for networks of biological populations, J Li, EL Ionides, AA King, M Pascual, N Ning - arXiv preprint arXiv ..., 2023 - arxiv.org
- Optimal simulation-based Bayesian decisions, J Alsing, TDP Edwards, B Wandelt - arXiv preprint arXiv:2311.05742, 2023 - arxiv.org

Number of Simulation-based Inference Papers by Year



Simulation-based Inference software packages



The Journal of Open Source Software

sbi: A toolkit for simulation-based inference

Alvaro Tejero-Cantero^{e, 1}, Jan Boelts^{e, 1}, Michael Deistler^{e, 1},
Jan-Matthis Lueckmann^{e, 1}, Conor Durkan^{e, 2}, Pedro J. Gonçalves^{1, 3},
David S. Greenberg^{1, 4}, and Jakob H. Macke^{1, 5, 6}

^e Equally contributing authors ¹ Computational Neuroengineering, Department of Electrical and Computer Engineering, Technical University of Munich ² School of Informatics, University of Edinburgh ³ Neural Systems Analysis, Center of Advanced European Studies and Research (caesar), Bonn ⁴ Model-Driven Machine Learning, Centre for Materials and Coastal Research, Helmholtz-Zentrum Geesthacht ⁵ Machine Learning in Science, University of Tübingen ⁶ Empirical Inference, Max Planck Institute for Intelligent Systems, Tübingen

DOI: [10.21105/joss.02505](https://doi.org/10.21105/joss.02505)

Software

- [Review ↗](#)
- [Repository ↗](#)
- [Archive ↗](#)

Summary

Scientists and engineers employ stochastic numerical simulators to model empirically observed phenomena. In contrast to purely statistical models, simulators express scientific principles that provide powerful inductive biases, improve generalization to new data or scenarios and allow for fewer, more interpretable and domain-relevant parameters. Despite these advantages, tuning a simulator's parameters so that its outputs match data is challenging. Simulation-based inference (SBI) seeks to identify parameter sets that a) are compatible with prior knowledge and b) match empirical observations. Importantly, SBI does not seek to recover a single 'best' data-compatible parameter set, but rather to identify all high probability regions of parameter space that explain observed data, and thereby to quantify parameter uncertainty. In Bayesian terminology, SBI aims to retrieve the posterior distribution over the parameters of interest. In

BayesFlow

Tests passing License MIT JOSS 10.21105/joss.05702 contributions welcome

Welcome to our BayesFlow library for efficient simulation-based Bayesian workflows! Our library enables users to create specialized neural networks for *amortized Bayesian inference*, which repay users with rapid statistical inference after a potentially longer simulation-based training phase.



MadMiner: Machine learning–based inference for particle physics

By Johann Brehmer, Felix Kling, Irina Espejo, and Kyle Cranmer

pypi package 0.6.3 build passing docs failing chat on gitter code style black License MIT DOI 10.5281/zenodo.1489147
 arXiv 1907.10621

Swyft



Swyft is a system for scientific simulation-based inference at scale.

pypi package 0.4.5 codecov 26% JOSS 10.21105/joss.04205 DOI 10.5281/zenodo.5752734

Swyft is the official implementation of Truncated Marginal Neural Ratio Estimation (TMNRE), a hyper-efficient, simulation-based inference technique for complex data and expensive simulators.

Wrap it up

The next 10 years



6.7

Software, Computing, and Data Science

Software and computing play a critical role in maximizing the science output of particle physics experiments. They are an integral component of experimental design, trigger and data acquisition, simulation, reconstruction, and data analysis. They also underlie simulation and design of particle accelerators. The complexity, computational needs and data volumes of particle physics experiments are expected to increase dramatically in the next decade or so. Advances in software and computing, including AI/ML, will be key for

6.8

Technology Innovations and Impact on Society

Particle physicists have a long history of recognizing, embracing, and fostering emerging technologies. The devices and tools that particle physics develops enable applications beyond fundamental physics in fields as different as medicine and aerospace. This section outlines the many ways particle physics spurs innovation, with example impacts.

6.8.1 – Computing

Particle physicists are the stewards of some of the world's largest datasets and therefore have a special connection to computing and data applications. Particle physicists were among the first to transition from analog to digitized data systems and were early adopters of the use of machine learning algorithms. The world wide web and the web browser were invented at CERN as an efficient information sharing system for particle physicists. To cope with enormous datasets, particle physicists played leading roles in the development of global-scale distributed computing technologies and high-performance networking. The ambitious portfolio of projects proposed for the next 10 years will drive innovation that addresses the scientific needs of the particle physics community and will broadly impact society.

6.8.2 – Artificial Intelligence and Machine Learning

The unique challenges encountered in particle physics have proven to be fertile ground for innovation of AI/ML. Importing new techniques into the particle physics context requires a process of translation, through which those techniques are stress-tested, generalized, and improved. Those improvements feed back into the wider AI/ML research community, completing a cycle of use-inspired research.

Thank you!

Questions?

“New directions in science are launched by new tools much more often than by new concepts. The effect of a concept-driven revolution is to explain old things in new ways. The effect of a tool-driven revolution is to discover new things that have to be explained.”

– FREEMAN DYSON

