

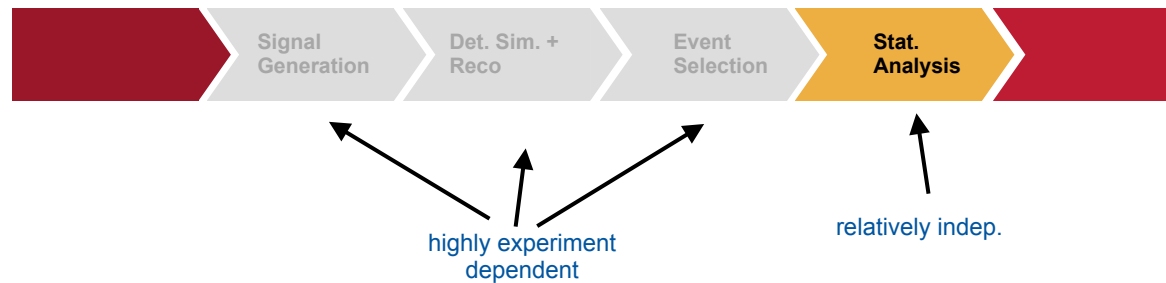
pyhf

Matthew Feickert, [Lukas Heinrich](#), Giordon Stark

IRIS-HEP

pyhf in the HEP analysis workflow:

Overall Goal of an analysis: inference on some parameters of nature (masses, couplings, ...)...



Statistical Analysis:

- final step of analysis
- you've collected all events, what do they tell you

Two broad categories

Decision: **aggregate** in the events before statistical analysis or not?

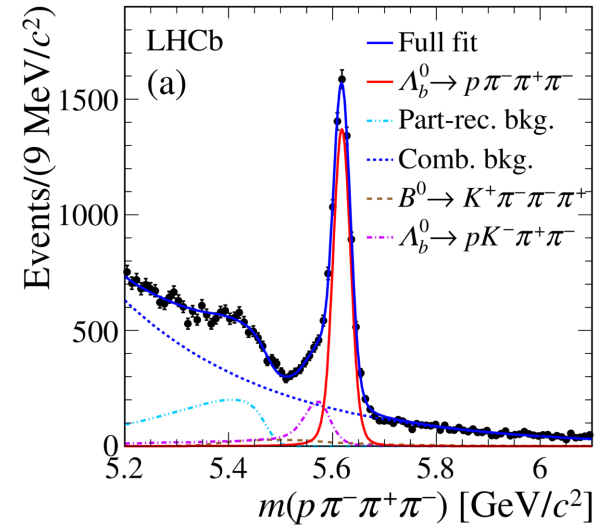
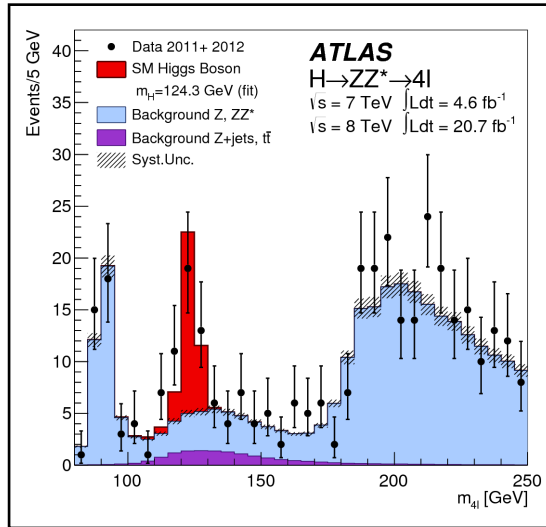
Statistical Analysis

yes

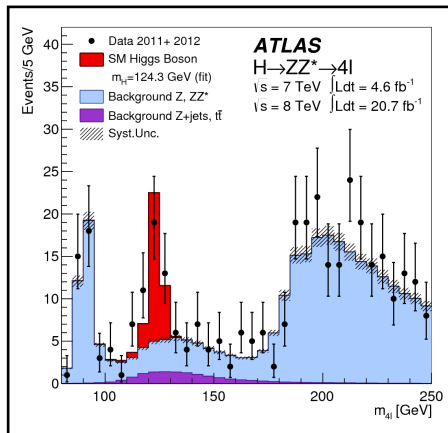
no

Binned

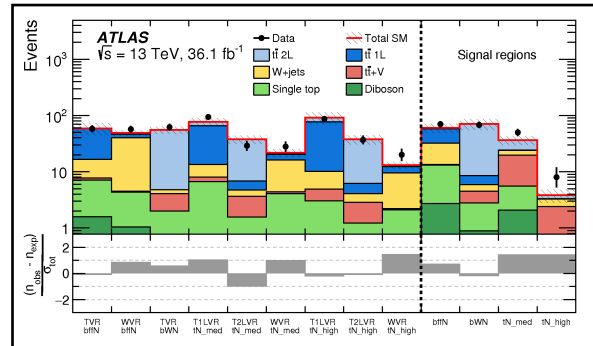
Unbinned



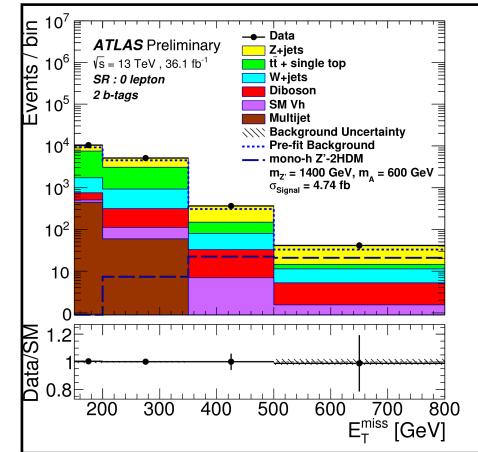
HistFactory is one of the most widely used families of stat. models used in HEP, esp. ATLAS: Majority binned template analyses formulated within this **one model**.



SM



SUSY



Exotics

pyhf

- Simultaneous Fit to multiple channels, each with multiple samples.
- Sample yields estimated function of nominal rate, scale factors and systematics.
- Systematics imply constraint term on the pdf.

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

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

Straight forward math:

but sofar only implementation in ROOT/RooStats/RooFit

HistFactory: A tool for creating statistical models for use with
RooFit and RooStats

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

June 20, 2012

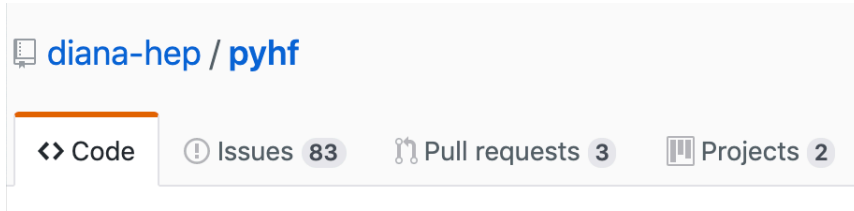
Contents

1 Introduction

2

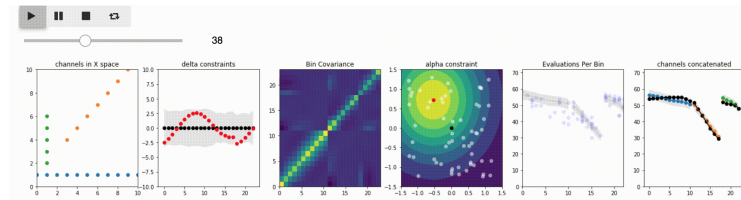
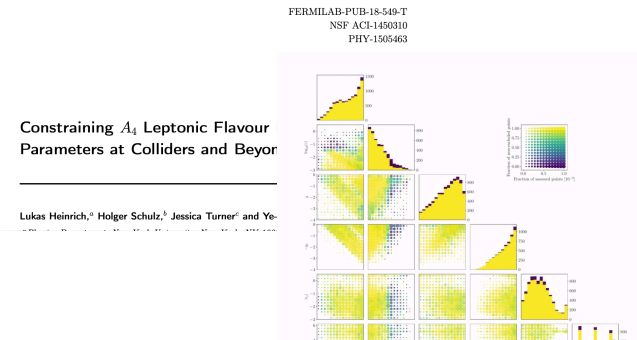
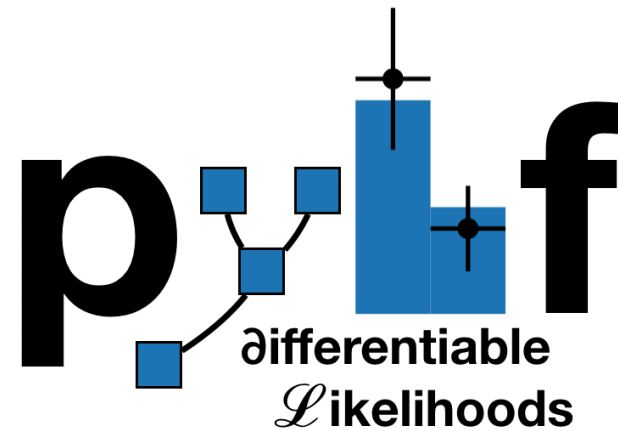
... but in principle implementable in other languages (of course)

pyhf



High Level Goals of pyhf:

- Python-based implementation
 - unlock python / data science eco-system
 - (new systematic types Gaussian Processes)
 - differentiable formulation
 - performance
 - lower barrier of entry to use HistFactory (e.g. phenomenologists..)
- Likelihood Preservation
 - side benefit: find language-independent spec
 - **likelihoods more important data product of an analysis**



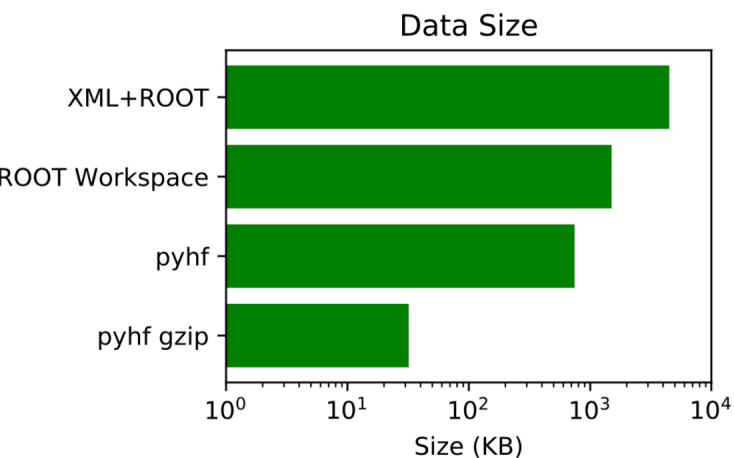
```
$> pip install pyhf
$> pyhf cls workspace.json
{
  "CLs_exp": [
    0.008897411763217407,
    0.03524468002619176,
    0.1243148689002353,
    0.3514186235832989,
    0.6941411699405086
  ],
  "CLs_obs": 0.03607409335946063
}
```

```
$> curl http://url-to-json/workspace.json|pyhf cls
{
  "CLs_exp": [
    0.002606408505279359,
    0.013820656047622592,
    0.0644552079856191,
    0.23526102499555396,
    0.573041803728844
  ],
  "CLs_obs": 0.05290116065118097
}
```


JSON Format:

Idea: remove "split brain" from XML + ROOT and inline all data into a single JSON document. For binned data, this should be fine.

(Should be find for very large binned likelihoods, but can use pointers into external storage if needed)



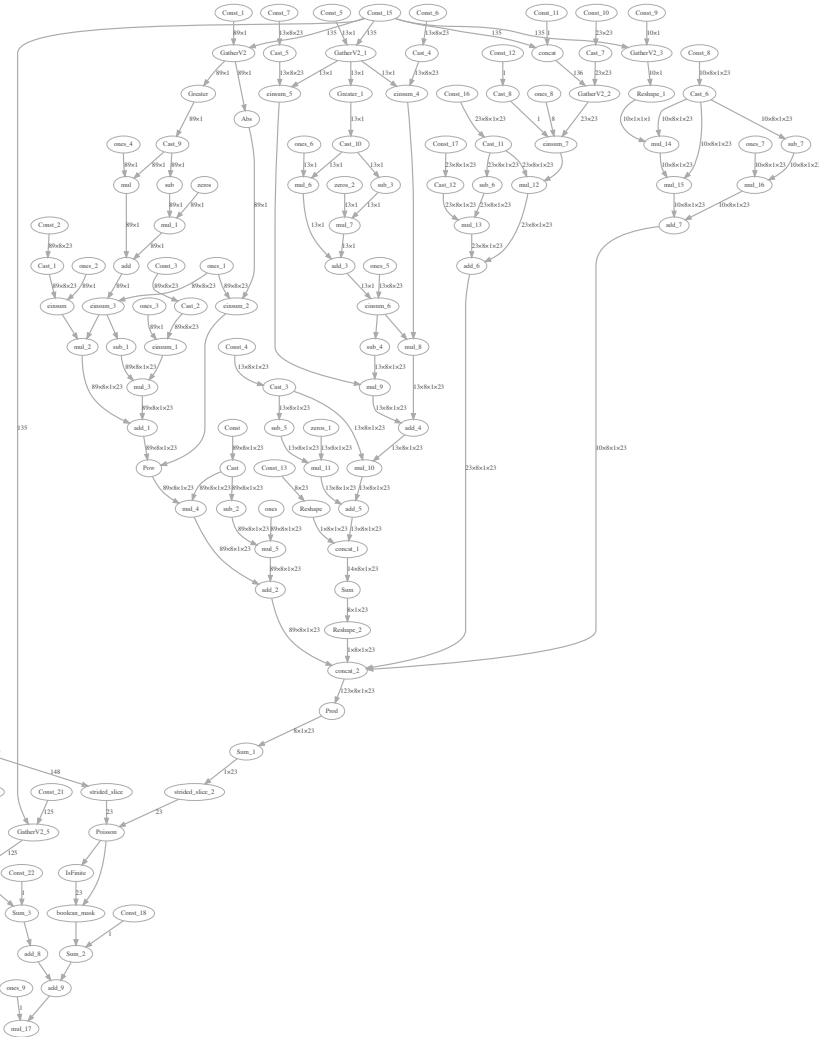
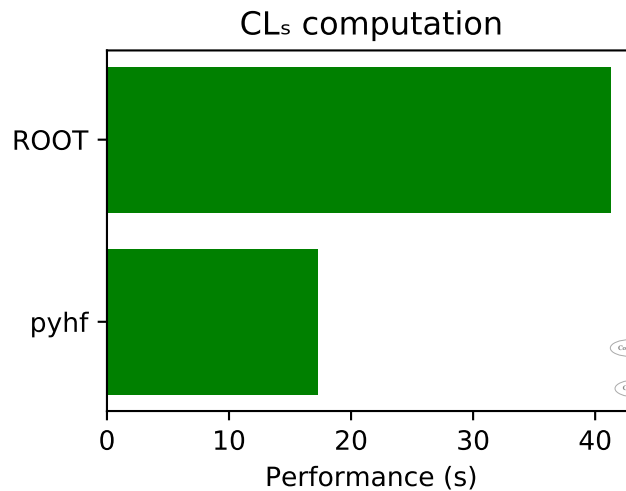
```
{
  "channels": [
    { "name": "singlechannel",
      "samples": [
        { "name": "signal",
          "data": [7.0, 2.0],
          "modifiers": [ { "name": "mu", "type": "normfactor", "data": null } ]
        },
        { "name": "background",
          "data": [50.0, 60.0],
          "modifiers": [ { "name": "uncorr_bkguncrt", "type": "shapesys", "data": [5.0,12.0] } ]
        }
      ]
    }
  ],
  "data": {
    "singlechannel": [50, 60]
  },
  "measurements": [
    { "name": "Measurement", "config": { "poi": "mu", "parameters": [] } }
  ]
}
```

Fully vectorized computation



Use shim to make backend agnostic

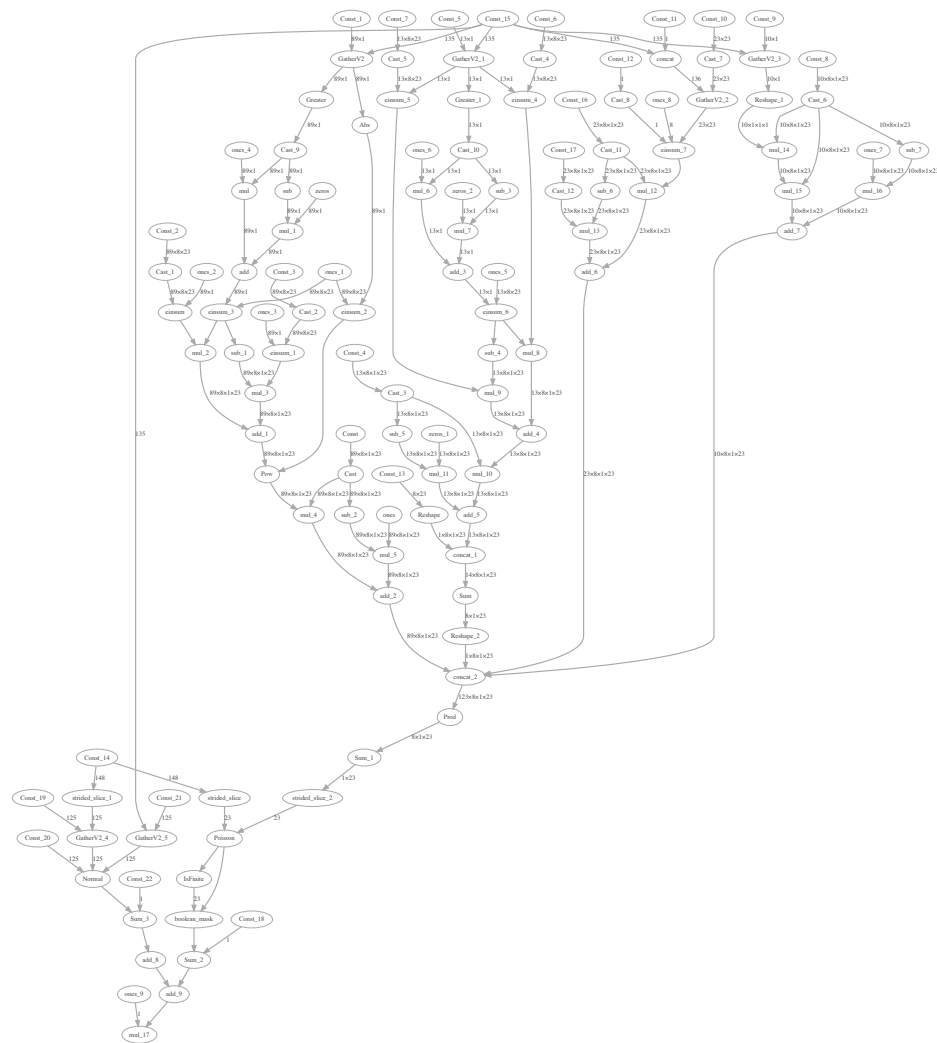
- NumPy (default)
- Tensorflow
- PyTorch
- (MXnet)
- (jax)
- (Dask)



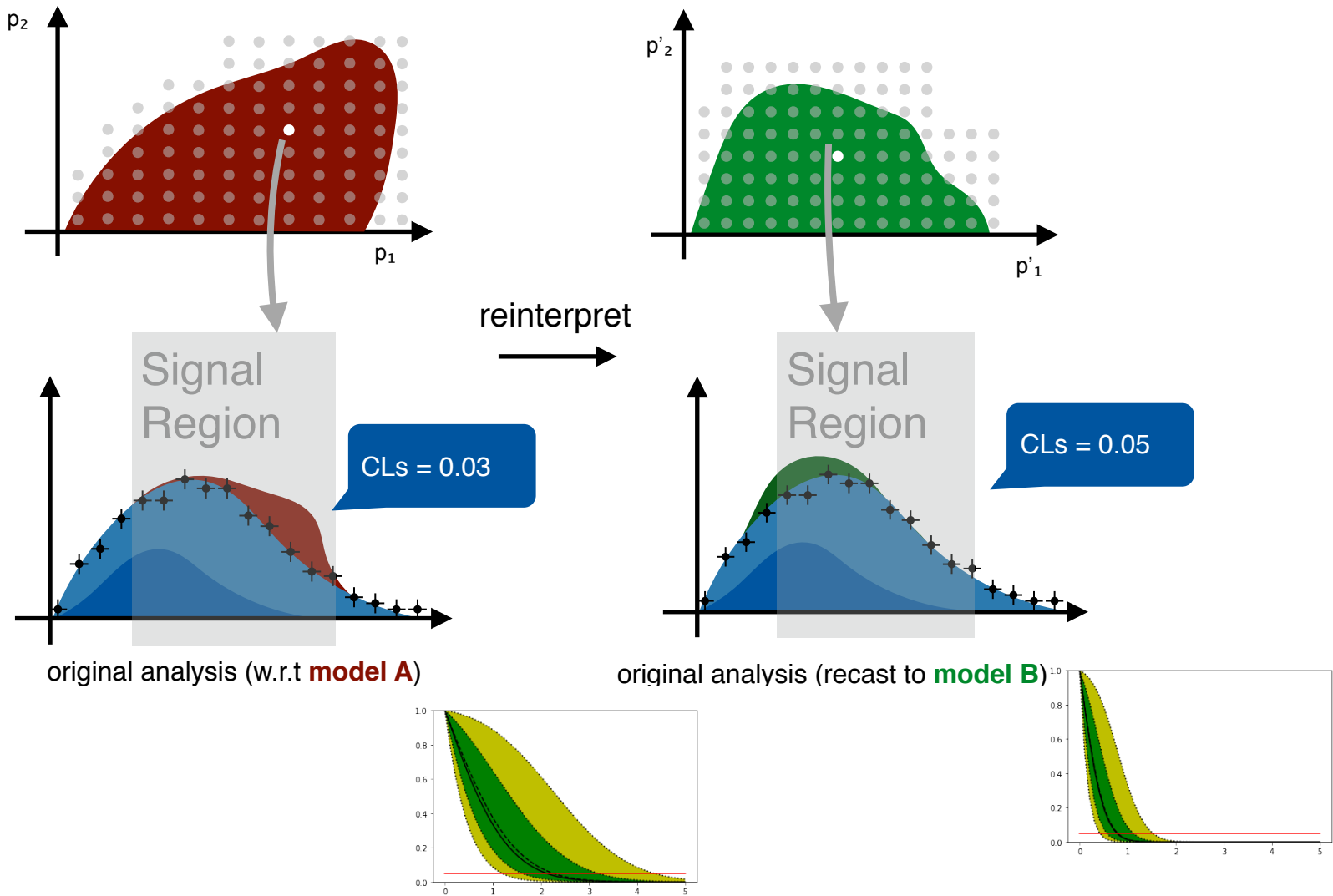
Advantage of non-Numpy backends



- Distribution across multiple machines (Dask)
- Hardware Acceleration (ML backends)
- Improved Fitting through Automatic Differentiation (ML Backends)



Easy RECAST



Easy RECAST

Advantages: very easy to modify likelihoods using **JSONPatch**
Important for e.g. Reinterpretations (think: pMSSM scan)

```
{
  "channels": [
    {
      "name": "singlechannel",
      "samples": [
        {
          "name": "signal",
          "data": [7.0, 2.0],
          "modifiers": [ { "name": "mu", "type": "normfactor", "data": null } ]
        },
        {
          "name": "background",
          "data": [50.0, 60.0],
          "modifiers": [ { "name": "uncorr_bkguncrt", "type": "shapesys", "data": [5.0,12.0] } ]
        }
      ]
    }
  ],
  "data": {
    "singlechannel": [50, 60]
  },
  "measurements": [
    { "name": "Measurement", "config": { "poi": "mu", "parameters": [] } }
  ]
}
```

+

```
[{
  "op": "replace",
  "path": "/channels/0/samples/0/data",
  "value": [7.0, 2.0]
}]
```

=

```
jsonpatch likelihood.json patch.json > new.json
pyhf cls new.json
```

```
{
  "channels": [
    {
      "name": "singlechannel",
      "samples": [
        {
          "name": "signal",
          "data": [7.0, 2.0],
          "modifiers": [ { "name": "mu", "type": "normfactor", "data": null } ]
        },
        {
          "name": "background",
          "data": [50.0, 60.0],
          "modifiers": [ { "name": "uncorr_bkguncrt", "type": "shapesys", "data": [5.0,12.0] } ]
        }
      ]
    }
  ],
  "data": {
    "singlechannel": [50, 60]
  },
  "measurements": [
    { "name": "Measurement", "config": { "poi": "mu", "parameters": [] } }
  ]
}
```

new yields injected

Connecting upstream

Input to HistFactory
is (structured) set of Histograms
(.. what used to be histbook?)

Histogram = result of
a functional DL operating
on event-wise data

Example:

- simple TTree::Draw-like syntax
- coffea / uproot / awkward ...

```
31 regions:
32 - name: regionA
33   binning: [-0.5,15.5,17]
34   filter: 'n_jet < 7 && n_mu > 0'
35   observable: 'n_jet'
36 - name: regionB
37   binning: [-0.5,15.5,17]
38   filter: 'n_jet > 8 && n_mu > 0'
39   observable: 'n_jet'
40
```

```
1 samples:
2 - name: 'data'
3   variations:
4   - name: nominal
5     tree: data16
6     data: true
7     glob: 'exported/data16/*.root'
8
9 - name: 'wjets_mc16a'
10  variations:
11  - name: nominal
12    tree: wjets_mc16a_Nom
13    glob: 'exported/wjets_mc16a/*.root'
14    weight: 'xs_weight * weight'
15    lumi: 36207.66
16 - name: 'zjets_mc16a'
17  variations:
18  - name: nominal
19    tree: zjets_mc16a_Nom
20    glob: 'exported/zjets_mc16a/*.root'
21    weight: 'xs_weight * weight'
22    lumi: 36207.66
23 - name: 'ttbar_mc16a'
24  variations:
25  - name: nominal
26    tree: ttbar_mc16a_Nom
27    glob: 'exported/ttbar_mc16a/*.root'
28    weight: 'xs_weight * weight'
29    lumi: 36207.66
30
31 regions:
32 - name: regionA
33   binning: [-0.5,15.5,17]
34   filter: 'n_jet < 7 && n_mu > 0'
35   observable: 'n_jet'
36 - name: regionB
37   binning: [-0.5,15.5,17]
38   filter: 'n_jet > 8 && n_mu > 0'
39   observable: 'n_jet'
```

Search for bottom-squark pair production with the ATLAS detector in final states containing Higgs bosons, b -jets and missing transverse momentum

The ATLAS Collaboration

