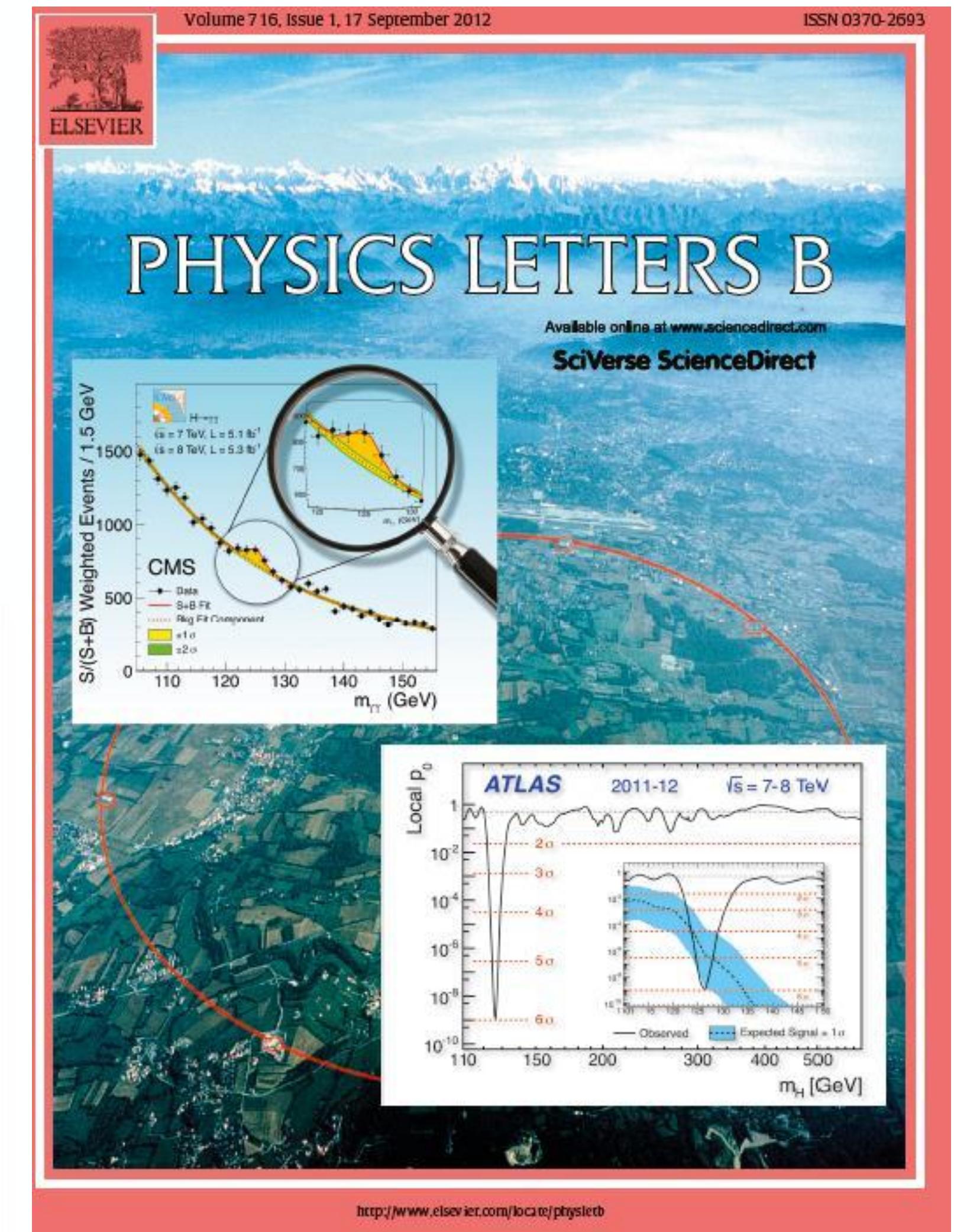
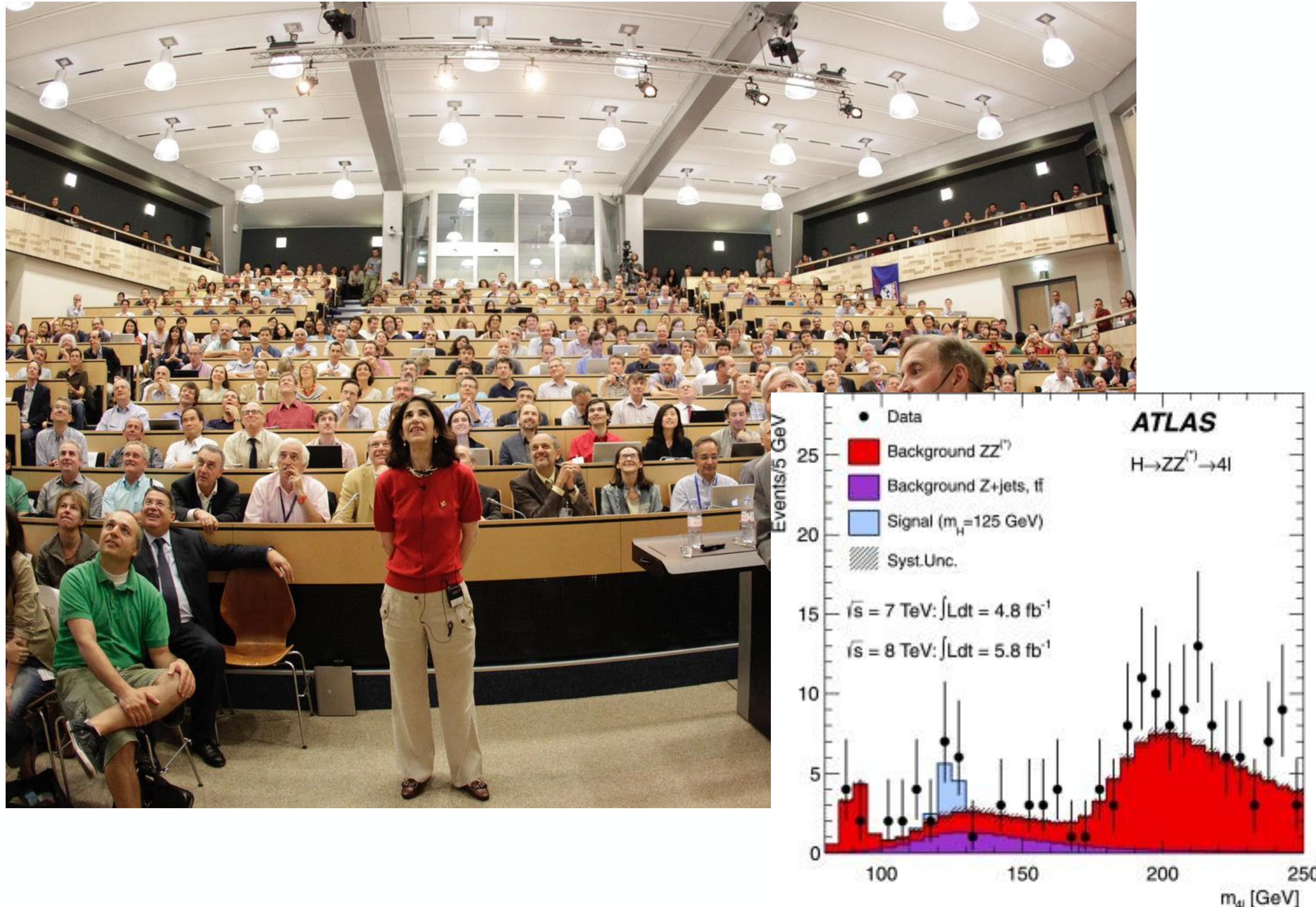


Systematics in HEP

Event Selection, Limits, Discovery

Lukas Heinrich, TUM - STAMPS

10 years ago

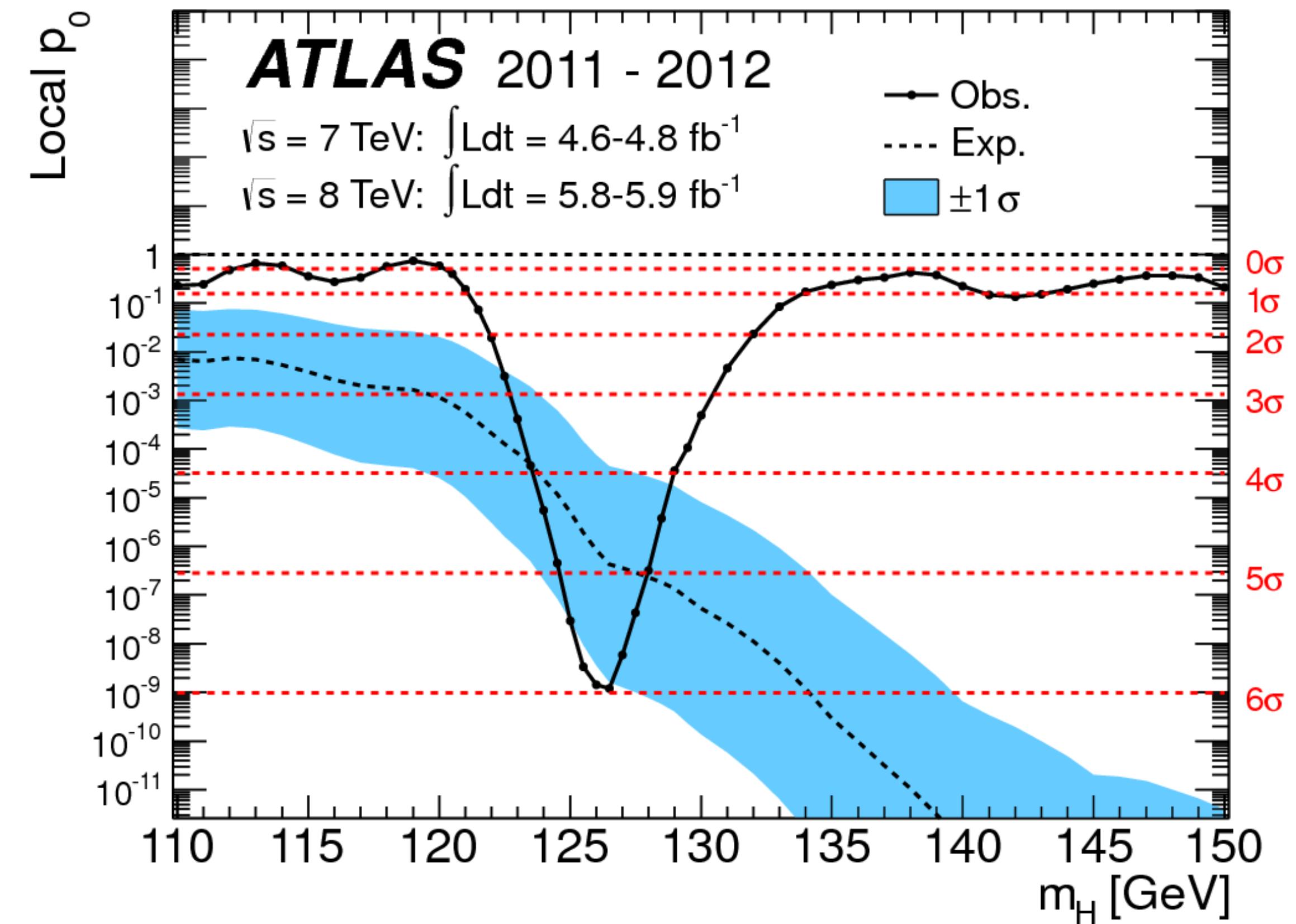


10 years ago

**Not only a significant moment
for physics, but also:**

LHC reached a **new level of statistical
sophistication** in producing these
results

This Talk: will try to give a feel
for how these results come about



Goals for the talk

I will try to summarize

- how we **build the statistical model** at the LHC
esp. **how we incorporate nuisance parameters**
- describe the current practice for **statistical tests**
- discuss **some frequent assumptions** made during inference
- focus on **frequentist viewpoint** as its dominant at LHC

Assuming familiarity w/ stats but less with HEP

The Large Hadron Collider

CMS

CERN

LHC

[Large Hadron Collider]

ALICE

LHCb

ATLAS



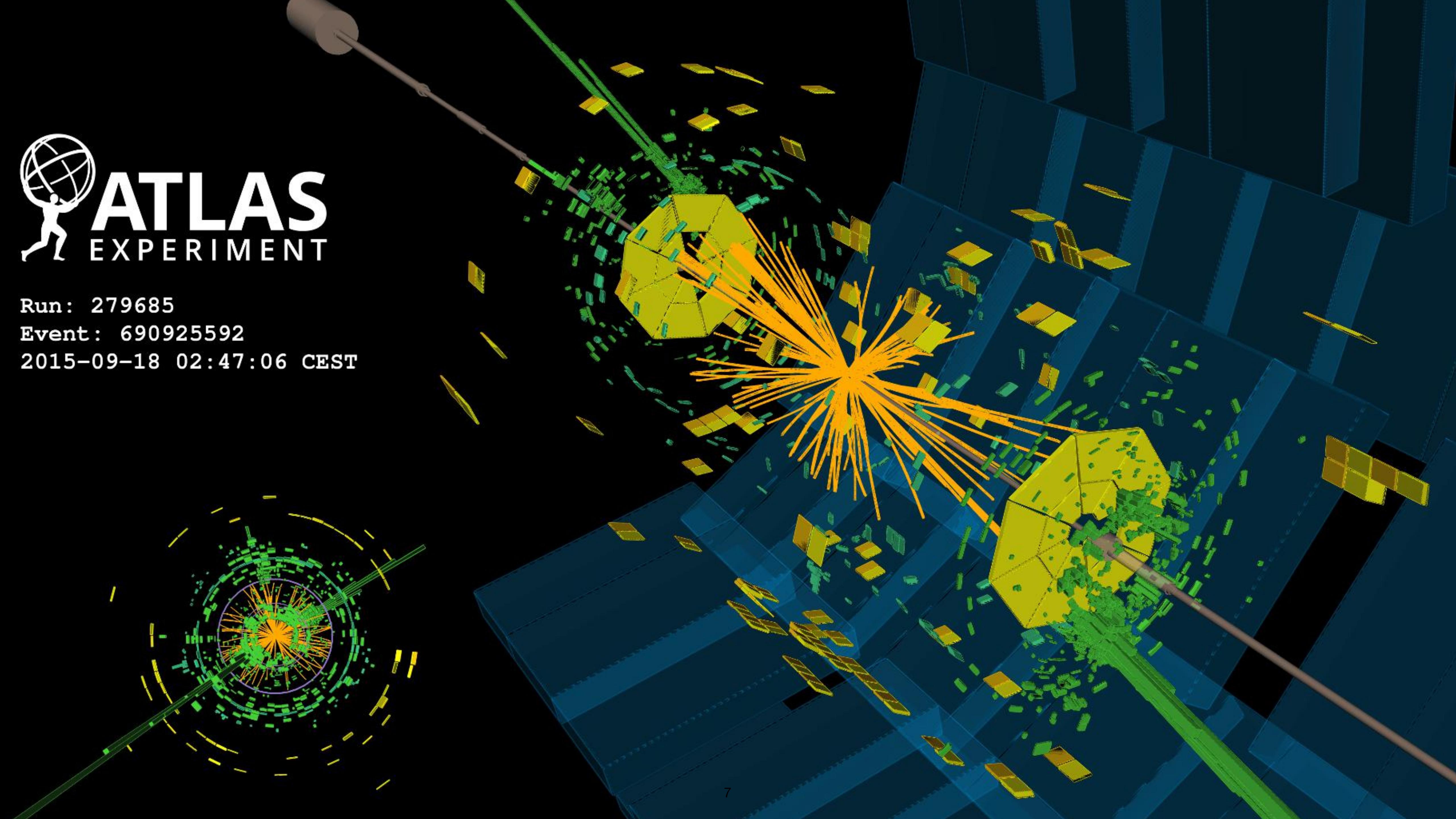


ATLAS EXPERIMENT

Run: 279685

Event: 690925592

2015-09-18 02:47:06 CEST

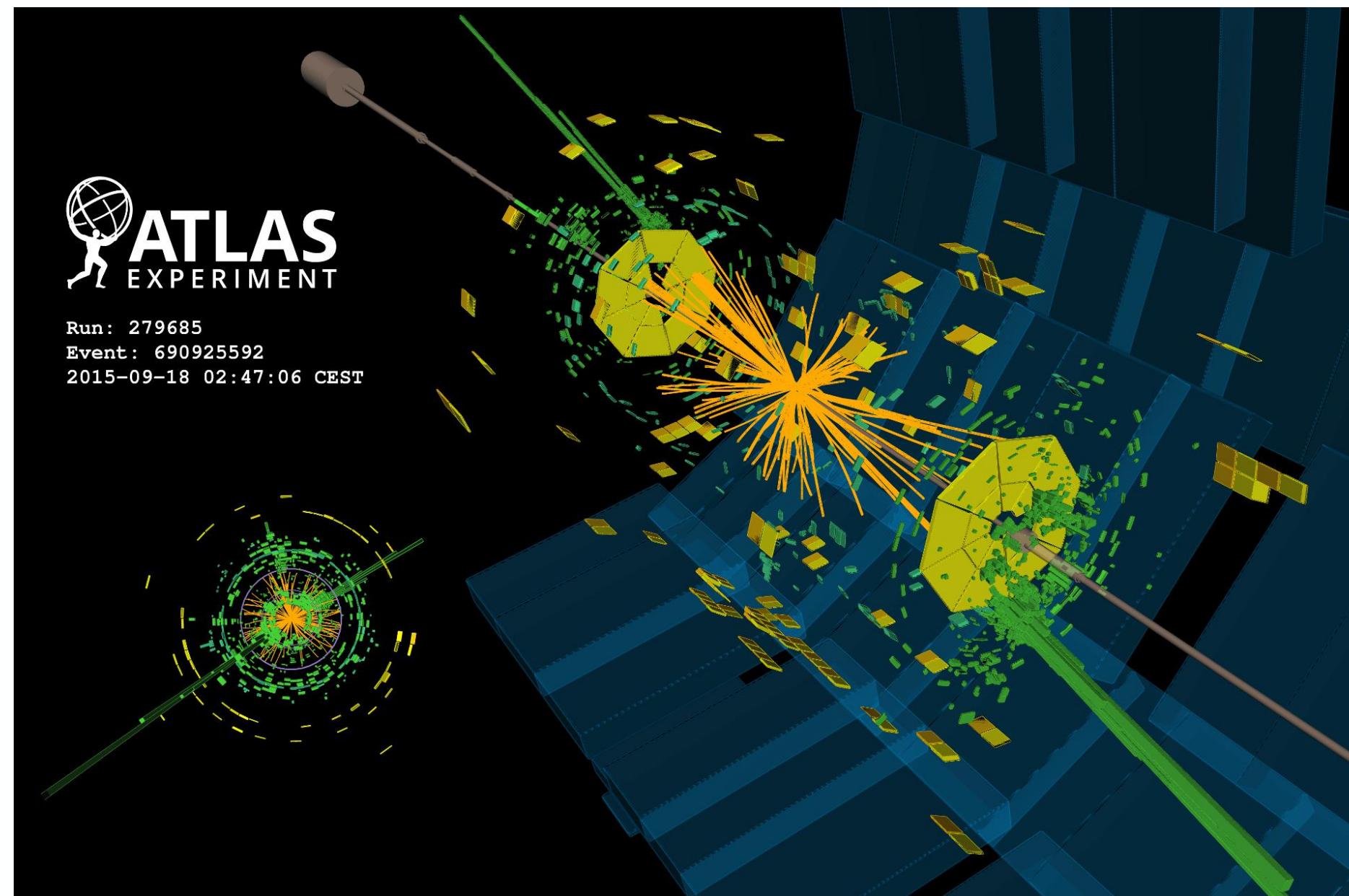


What we would like to do

$$p(\text{theory} \mid \text{data}) = \frac{p(\text{data} \mid \text{theory})p(\text{theory})}{p(\text{data})}$$

How to describe the data?

The forward model of the data is pretty complicated



data

$p(\text{data} \mid \text{theory})?$



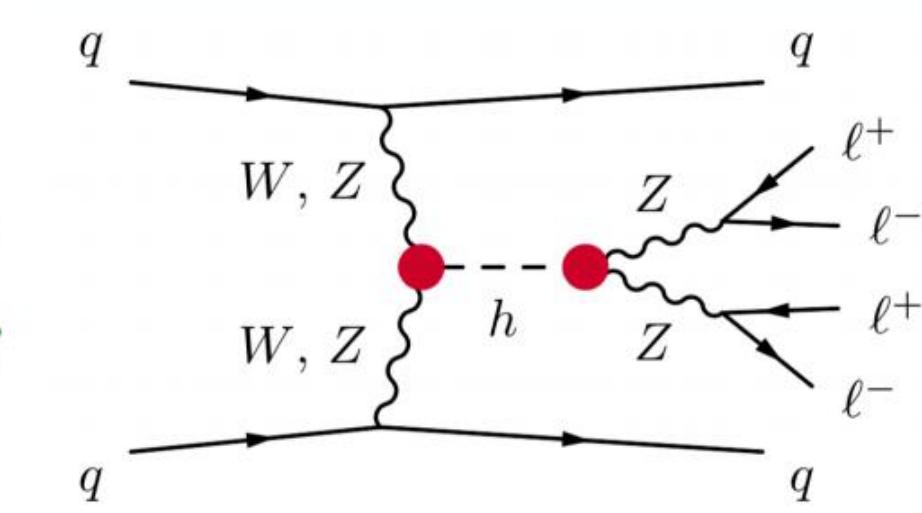
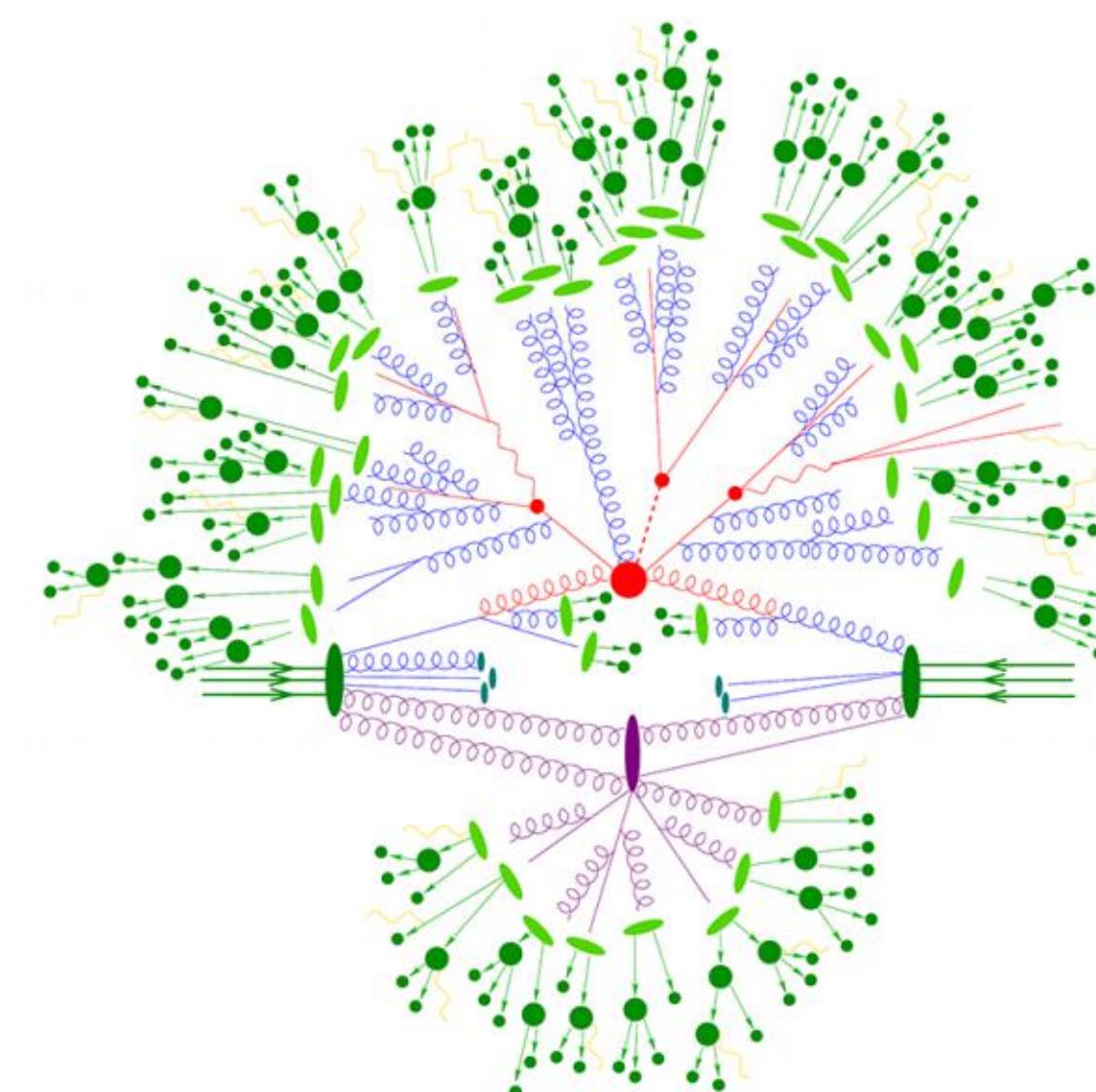
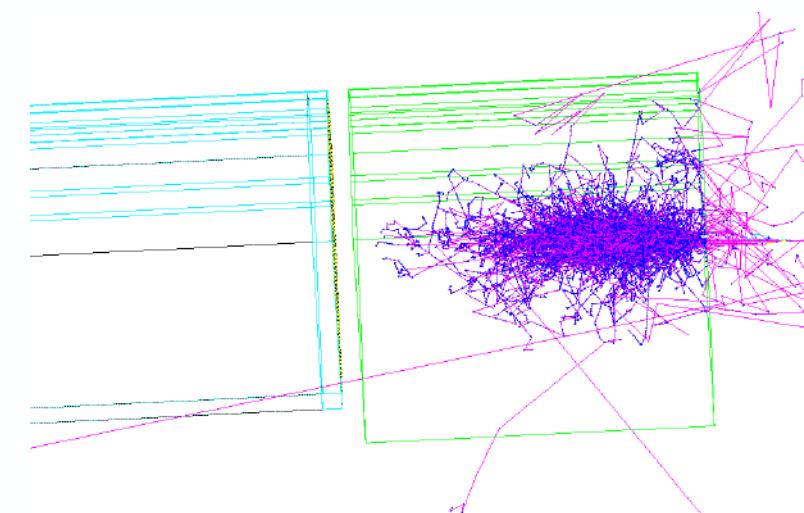
$$\begin{aligned}\mathcal{L} = & -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} \\ & + i \bar{\psi} \not{D} \psi + h.c. \\ & + \bar{\Psi}_i \gamma_{ij} \Psi_j \phi + h.c. \\ & + D_\mu \phi l^2 - V(\phi)\end{aligned}$$

theory

How to describe the data?

Fortunately: particle physics has very strong foundational modeling:

- can describe generative process as a (very deep) Markov Model



$$\begin{aligned}\mathcal{L} = & -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} \\ & + i \bar{\psi} \not{D} \psi + h.c. \\ & + \bar{\psi}_i y_{ij} \psi_j \phi + h.c. \\ & + \not{D}_\mu \phi /^2 - V(\phi)\end{aligned}$$

$$p(x|z_d)$$

*observed voltages, etc
in read-out electronics*

$$p(z_d|z_h)$$

*particle interaction
with dense material*

$$p(z_h|z_p)$$

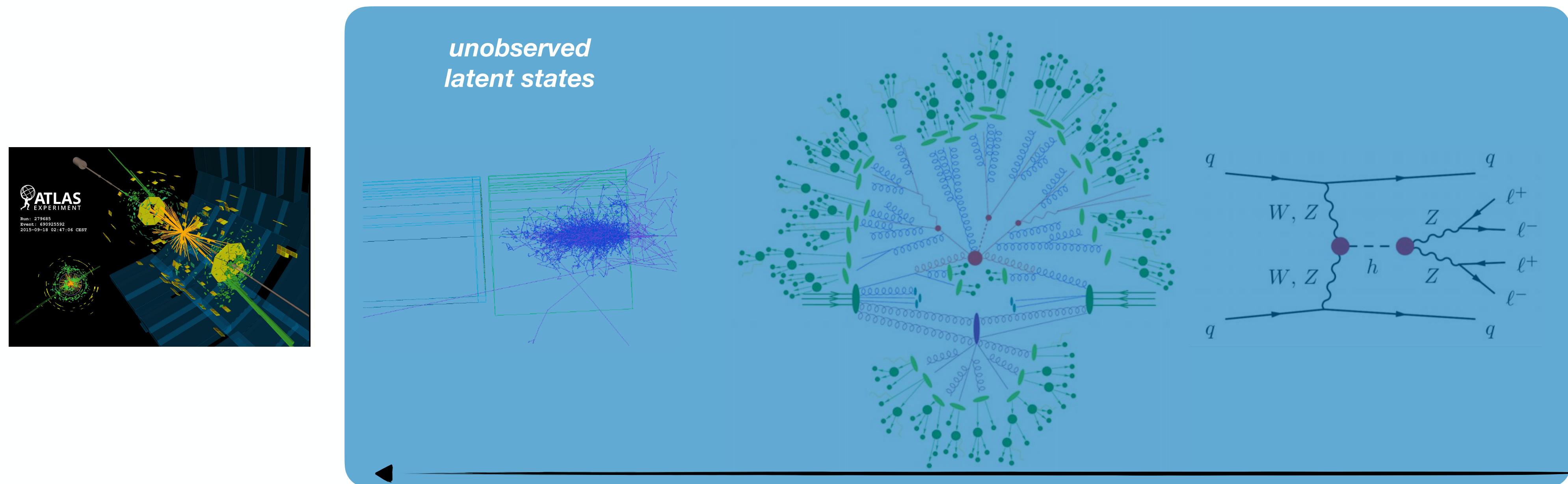
*particle evolution
(decays, radiation)*

$$p(z_p|\theta)$$

*core physics
process*

How to describe the data?

Unfortunately, all this latent evolution is not observable
all we get is the final data x : *integrate out all latent states*



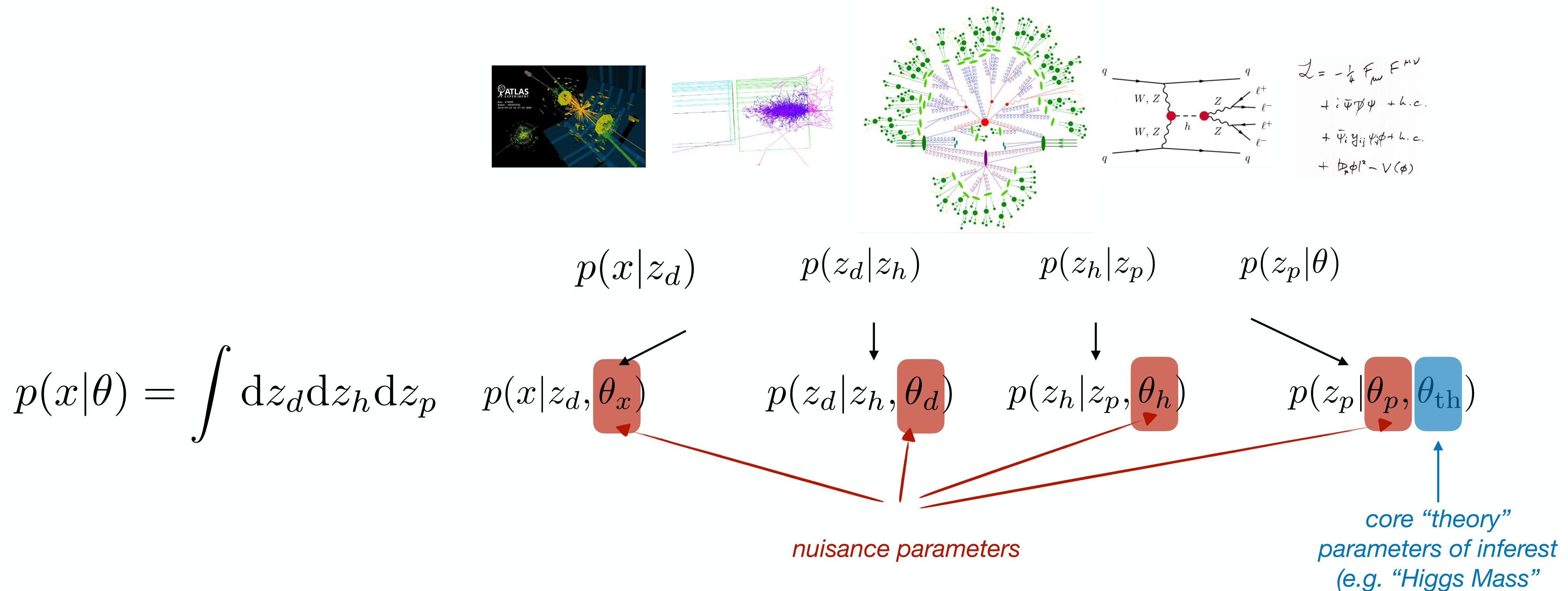
$$\begin{aligned} \mathcal{L} = & -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} \\ & + i \bar{\psi} \not{D} \psi + h.c. \\ & + \bar{\psi}_i y_{ij} \psi_j \phi + h.c. \\ & + \not{D}_\mu \phi /^2 - V(\phi) \end{aligned}$$

$$p(x|\theta) = \int dz \, p(x, z) = \int dz_d dz_h dz_p \, p(x|z_d) p(z_d|z_h) p(z_h|z_p) p(z_p|\theta)$$

The parameters θ

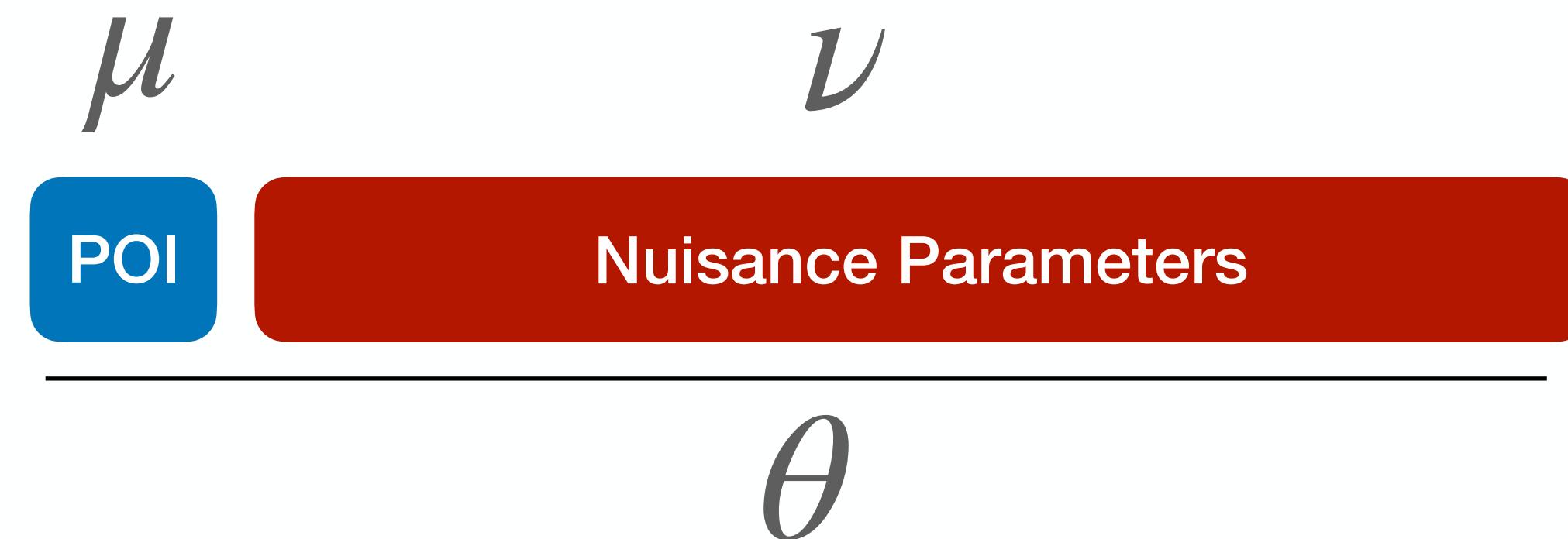
One more issue: the “theory” space is not the only thing effecting the data

- **every step of the forward process comes with its own parameters**
(we understand the process generally but need additional knobs to model the data)



Inference Goals

The nuisance parameters are necessary but seriously complicate the model



But not what we are interested in producing results on them.
Need to reflect this in inference methods

The Saving Grace: Simulators

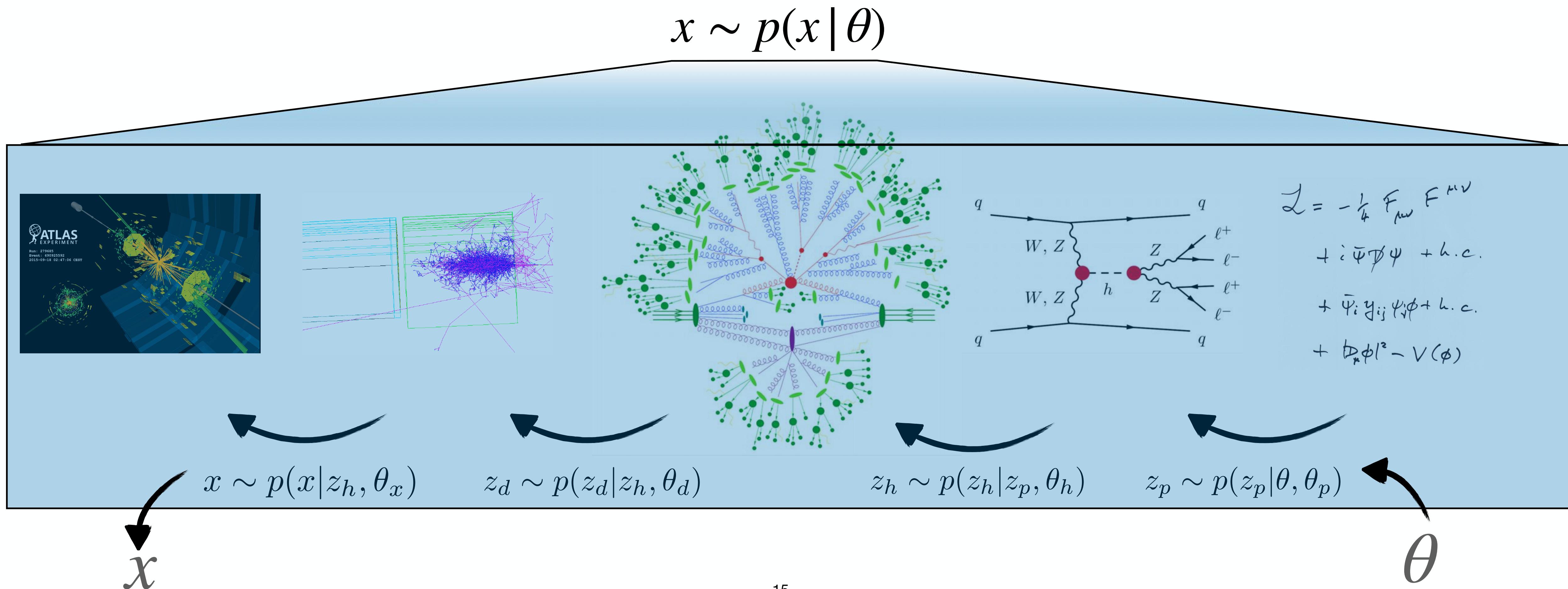
So far a pretty bleak picture: but we still do science, what gives?

Simulators!



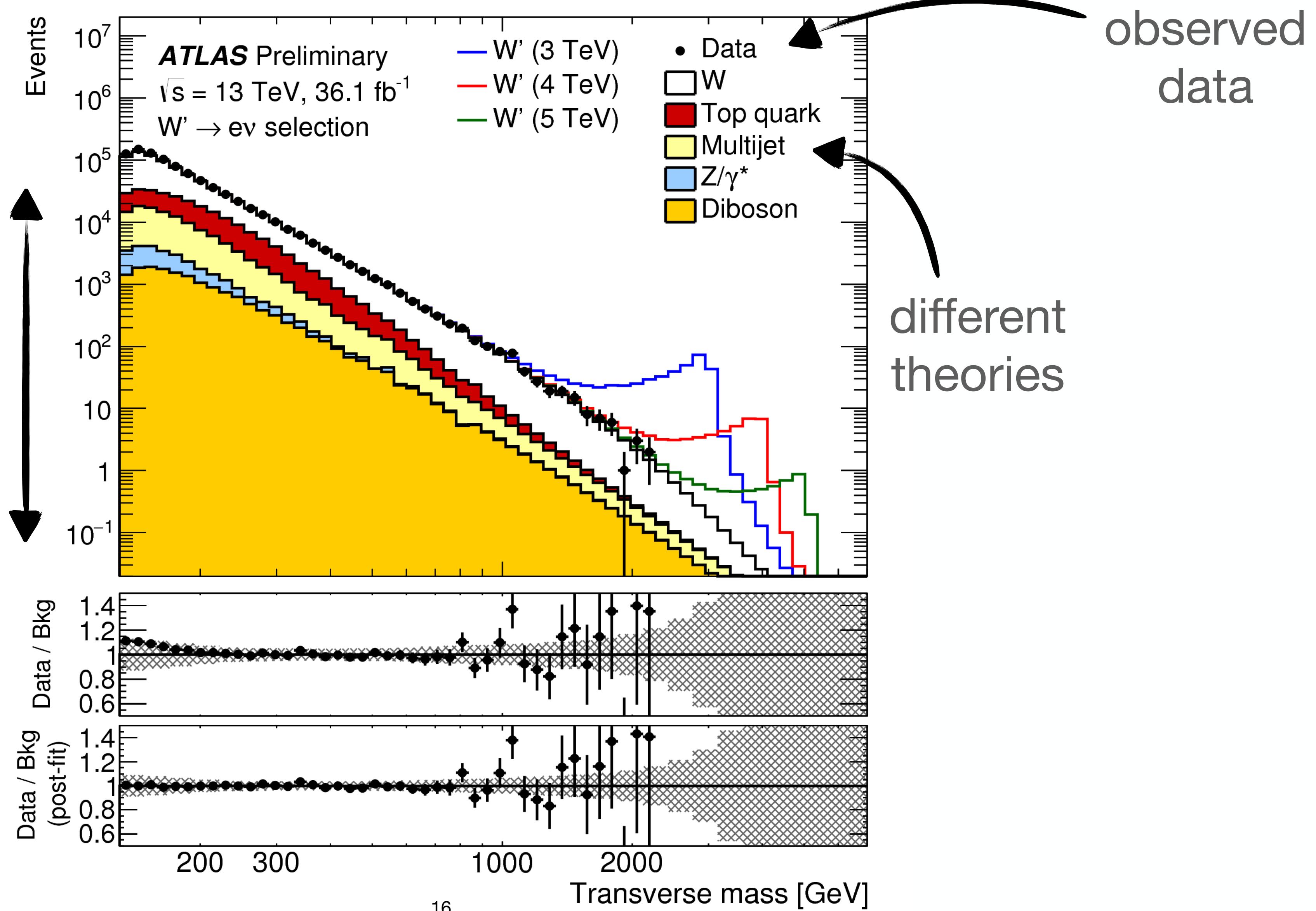
The Saving Grace: Simulators

While we cannot *evaluate the likelihood*, we can *sample it*.



Extremely High-Quality Simulation

7 orders
of magnitude
in density



Frequentist Statistics and Likelihood-Free Inference

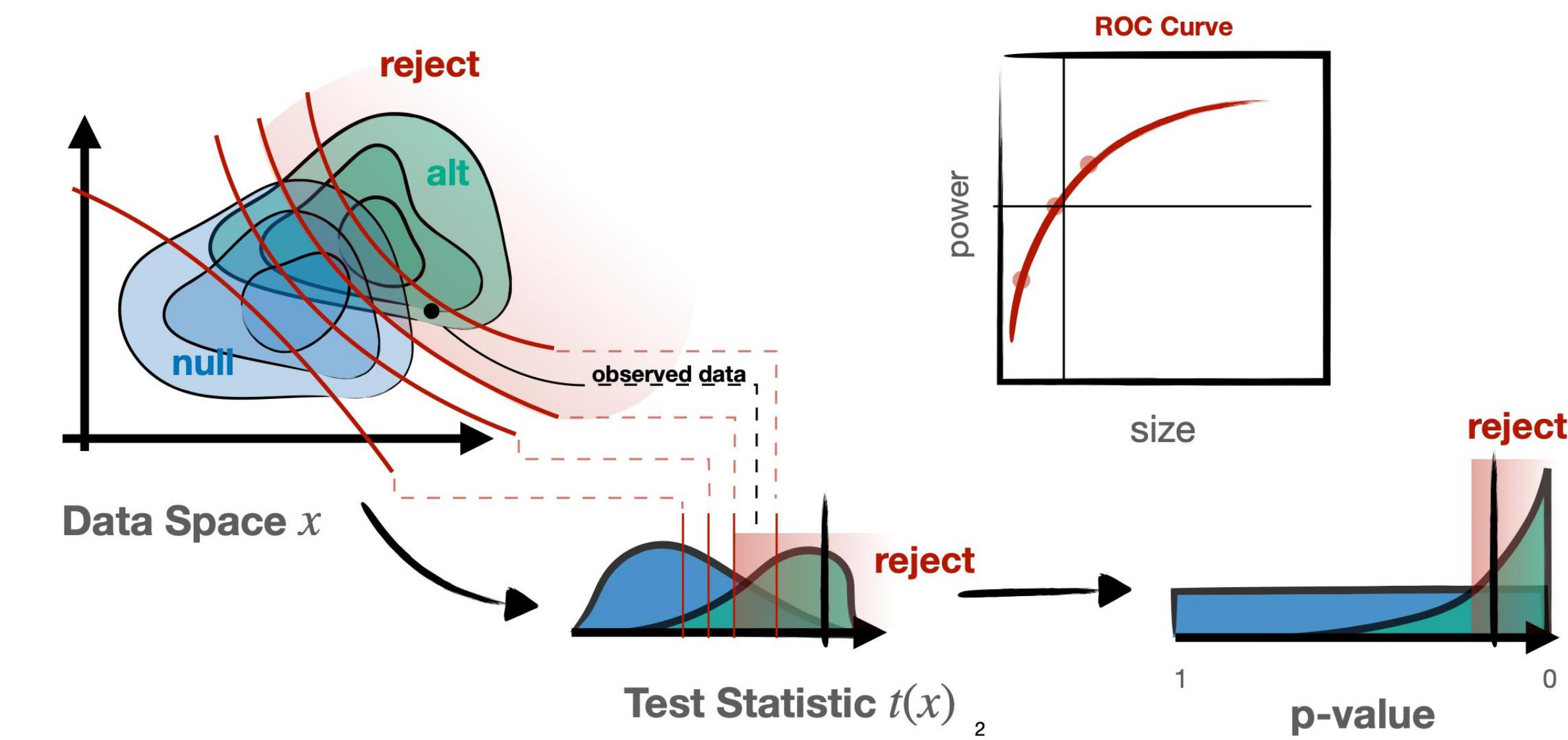
HEP one of major domains where frequentist statistics is used most results

One (possible) reason:

Freq. Stats is fundamentally more amenable to likelihood-free inference

Example: Hypothesis Tests

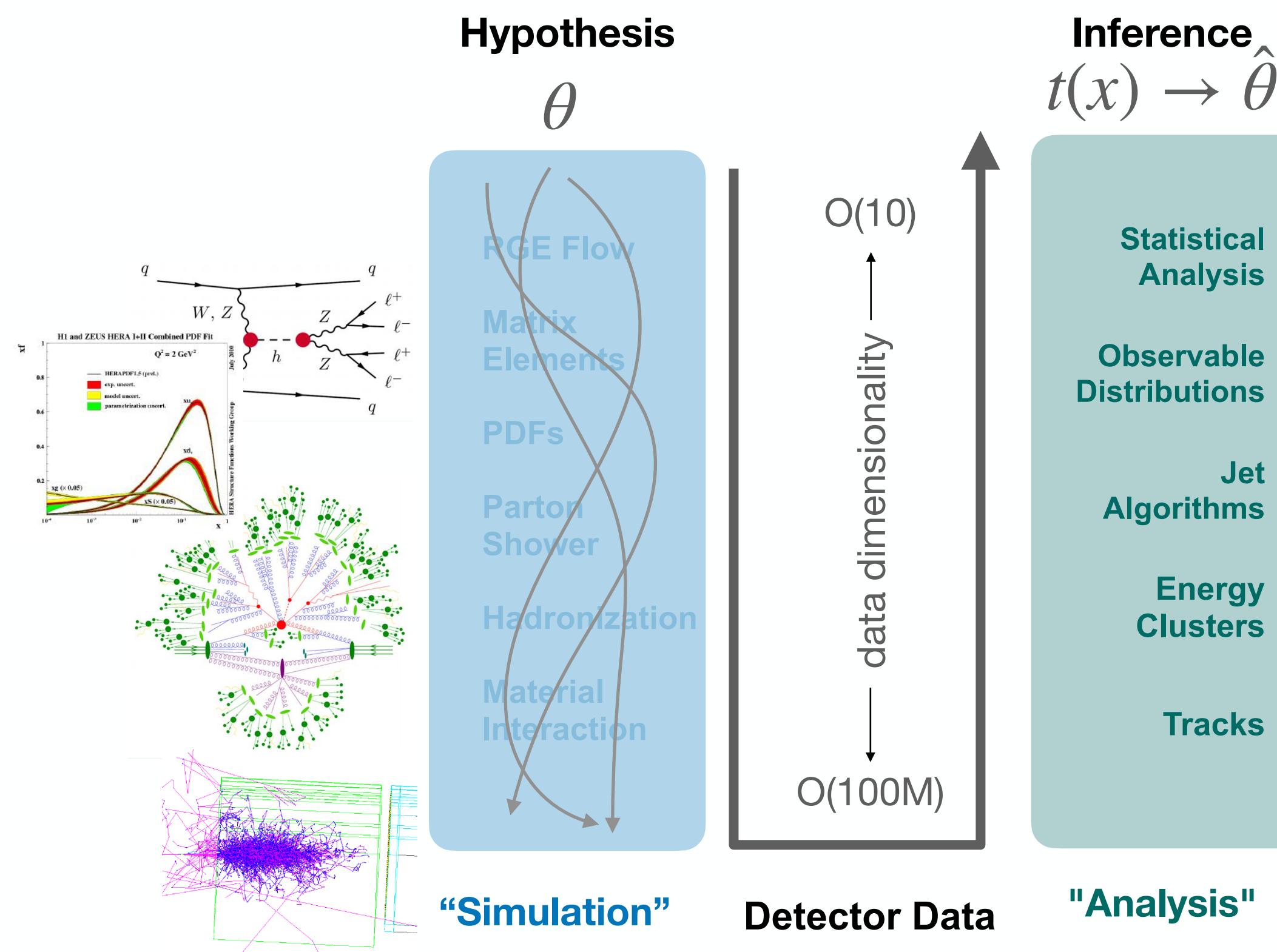
- just need be able to sample
 $x \sim p(x | \theta)$
- and be able to evaluate test statistic
 $x \rightarrow t(x)$



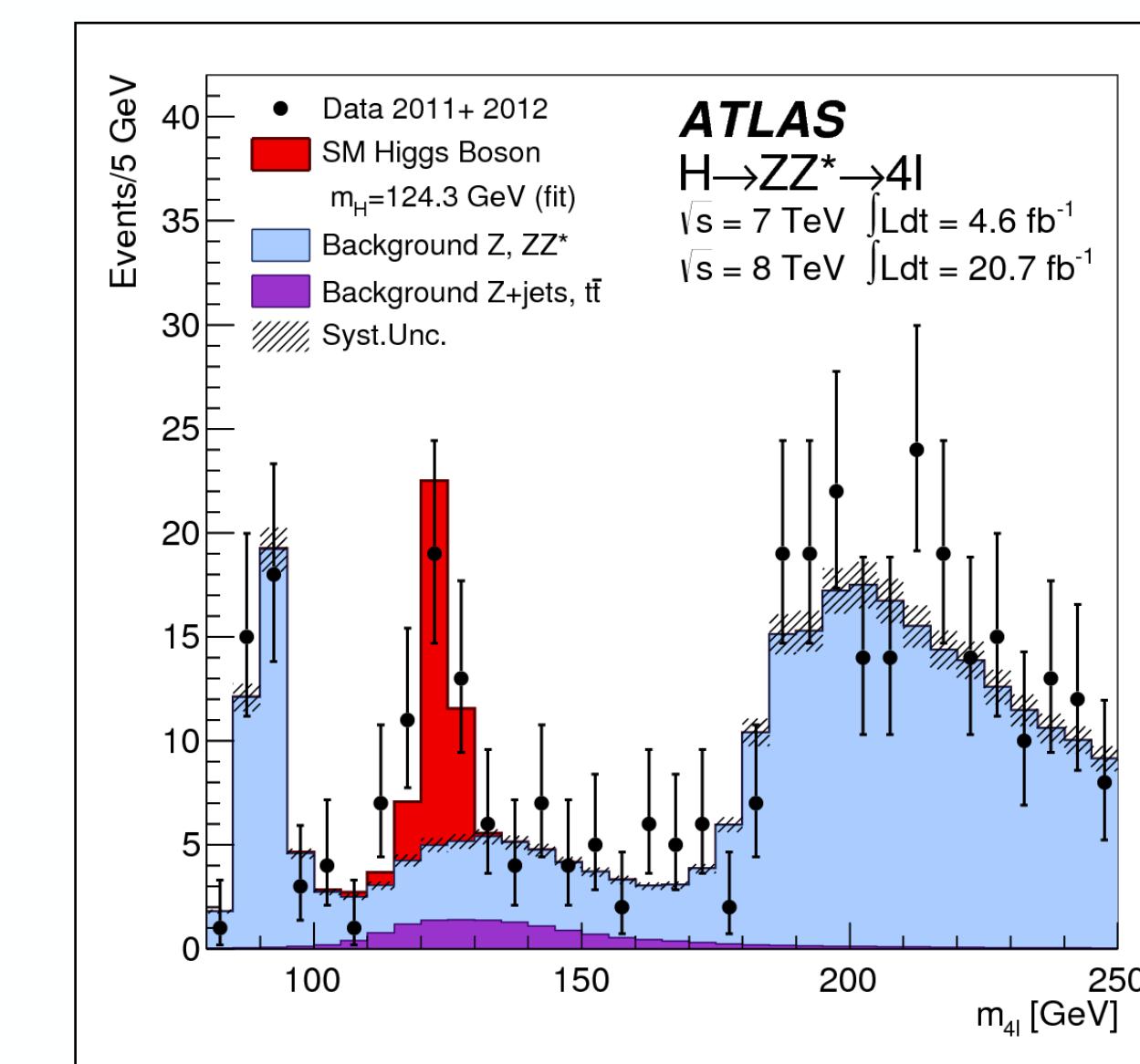
Compressing the data

Given the freedom of frequentist analysis, most of LHC analysis is about finding good data representations

Approximately follows a step-by-step inference / inversion of the data



In some sense just a very sophisticated way to build summary statistics of the data

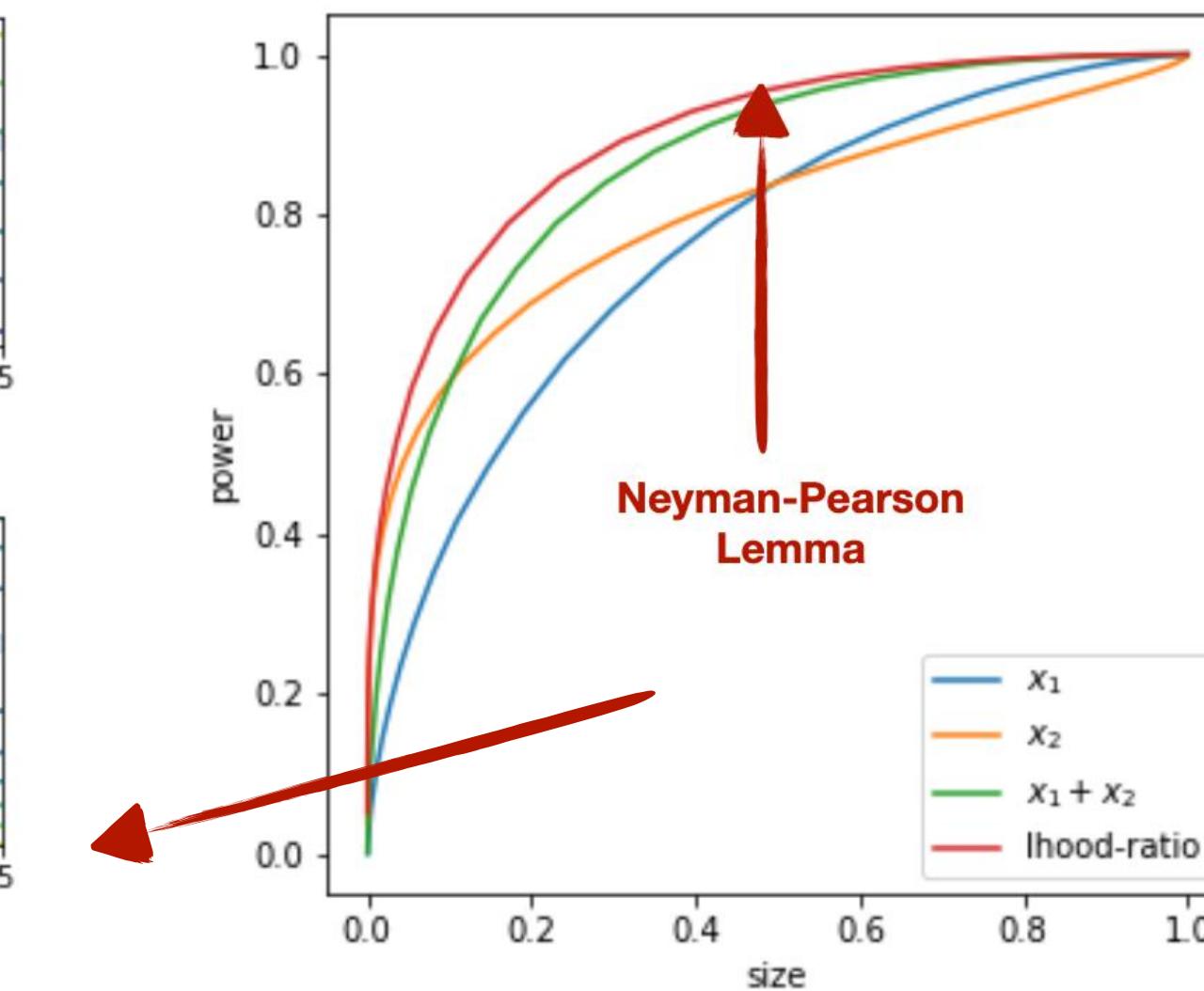
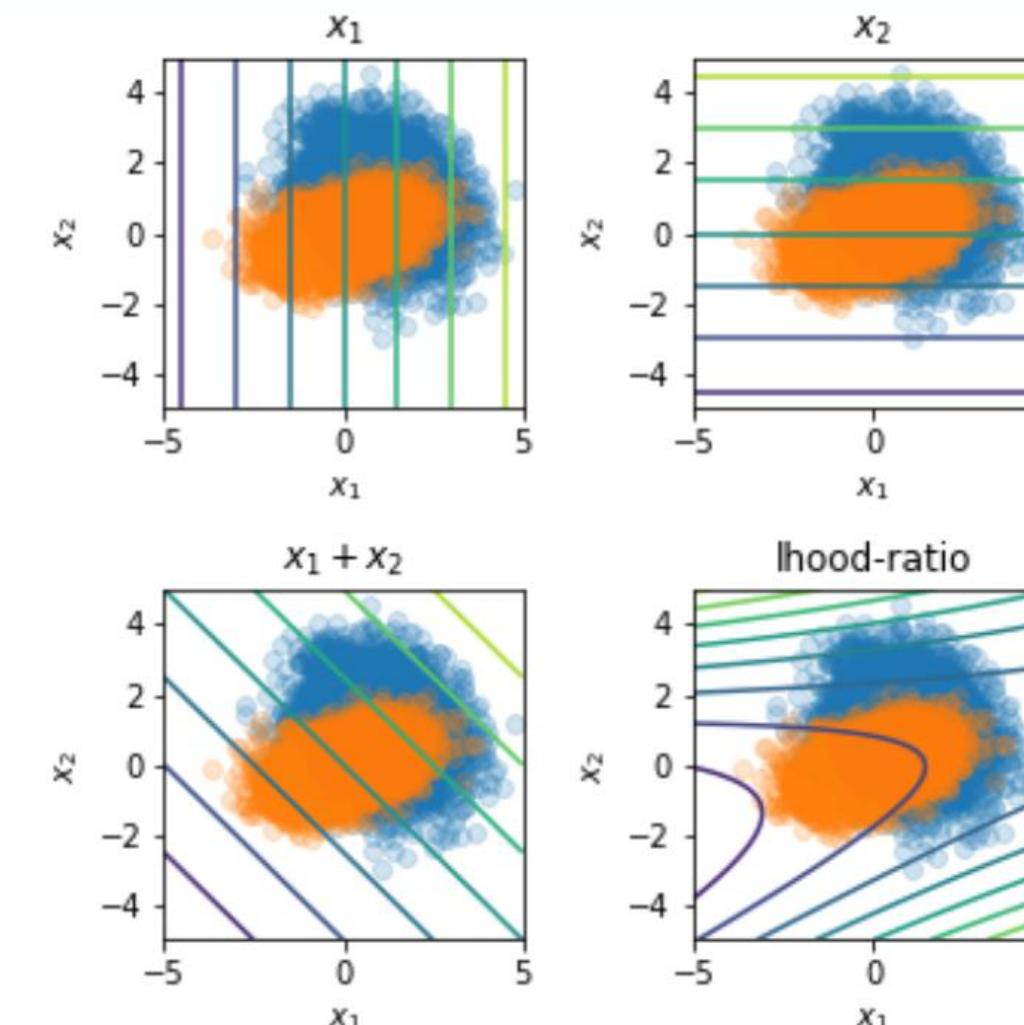
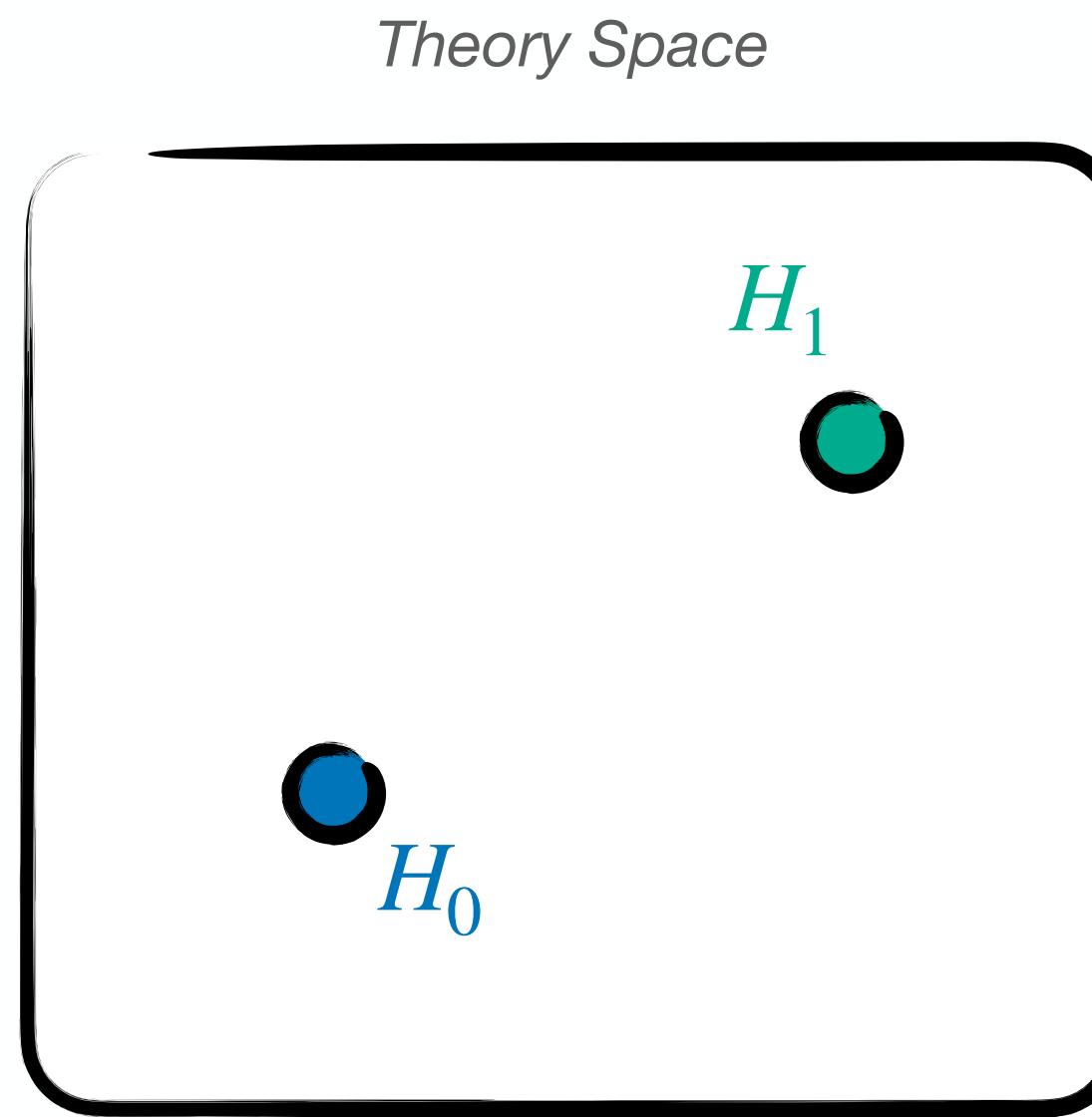


Likelihood are still useful

Even though frequentist analysis does not need likelihoods,
they are still useful in order to find good test statistics

Classic Neyman-Pearson Lemma: optimal binary test
requires evaluating the likelihood (ratio)

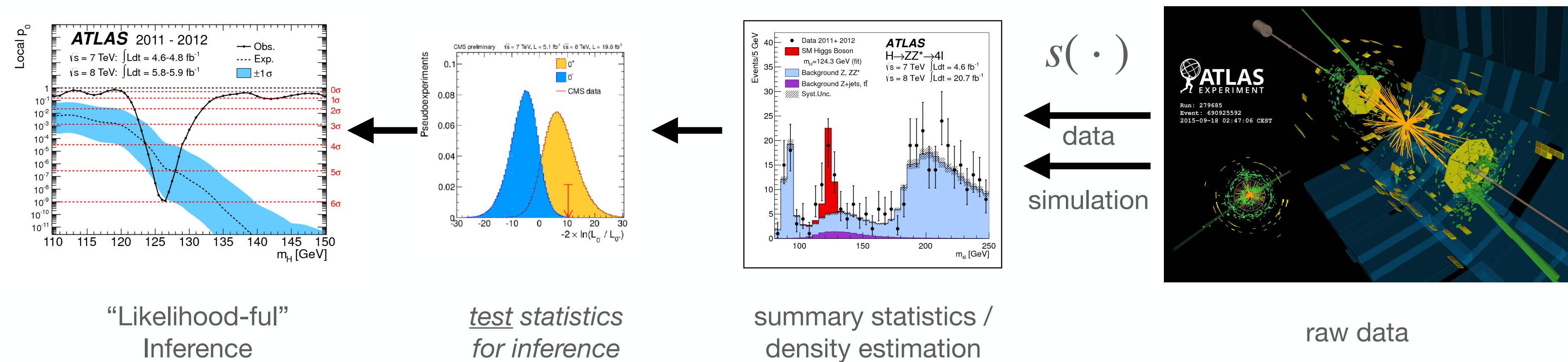
$$t(x) = -2 \log \frac{p(x|H_0)}{p(x|H_1)}$$



The usual workflow

In light of this, ***the usual strategy in HEP*** uses “approximate likelihoods” designed to model densities of summary statistics by using samples from simulation (e.g. histograms, kernel densities, ...)

$$p(s | \theta) = p(s(x) | \theta)$$



Most work goes into carefully building this “approximate likelihood”
(recently lots of interest in “more likelihood-free” methods, cf. later)

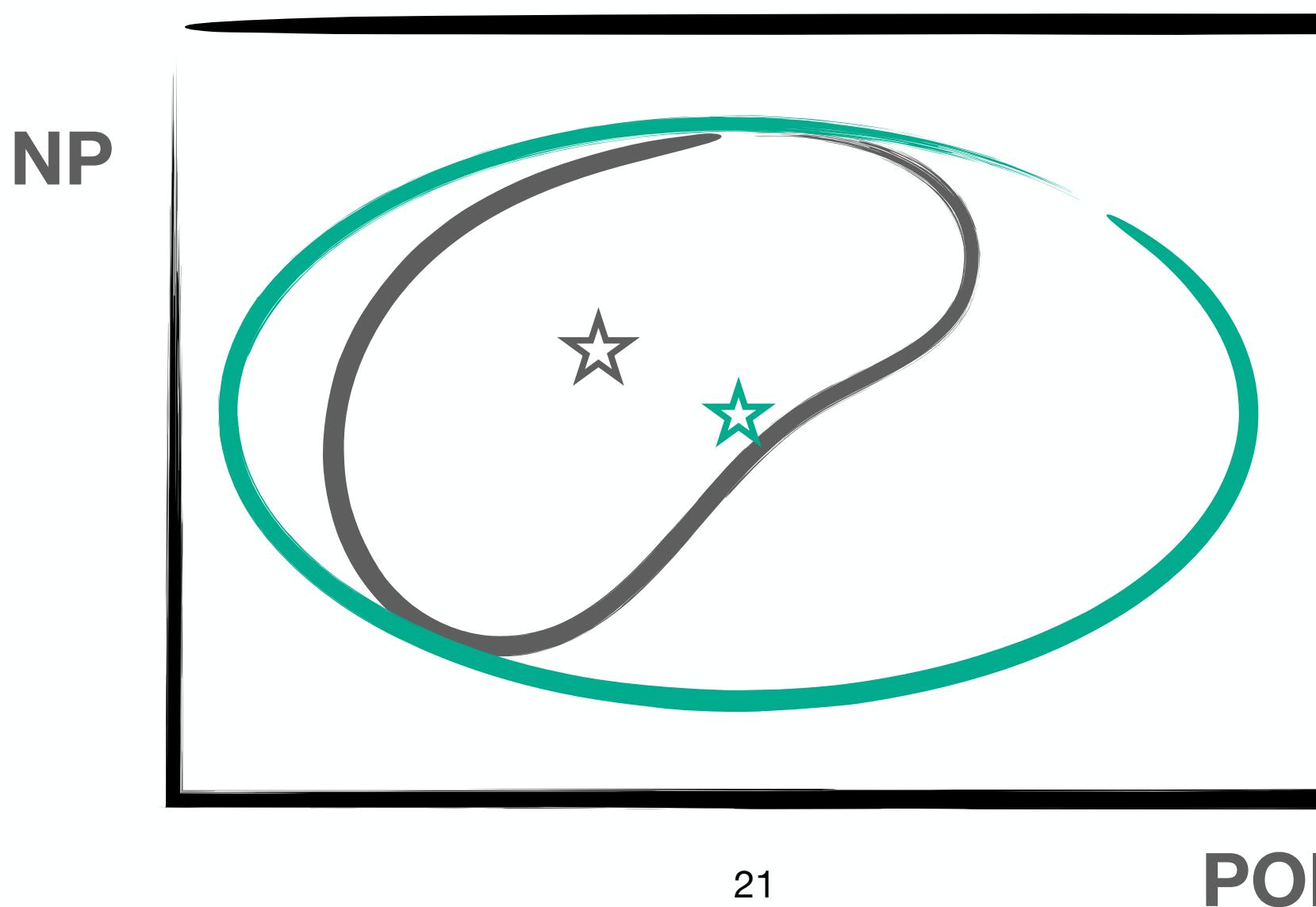
Advantages & Disadvantages of this workflow

Analysis via summary statistics: $p(x = t(\text{data}), \theta)$

Curse: by definition lose sensitivity (data processing inequality)

Blessing: by the same token gives freedom to tune for desiderata (“planned degradation”)

- e.g. design $t(x)$ to be robustness against NPs, remove correlations

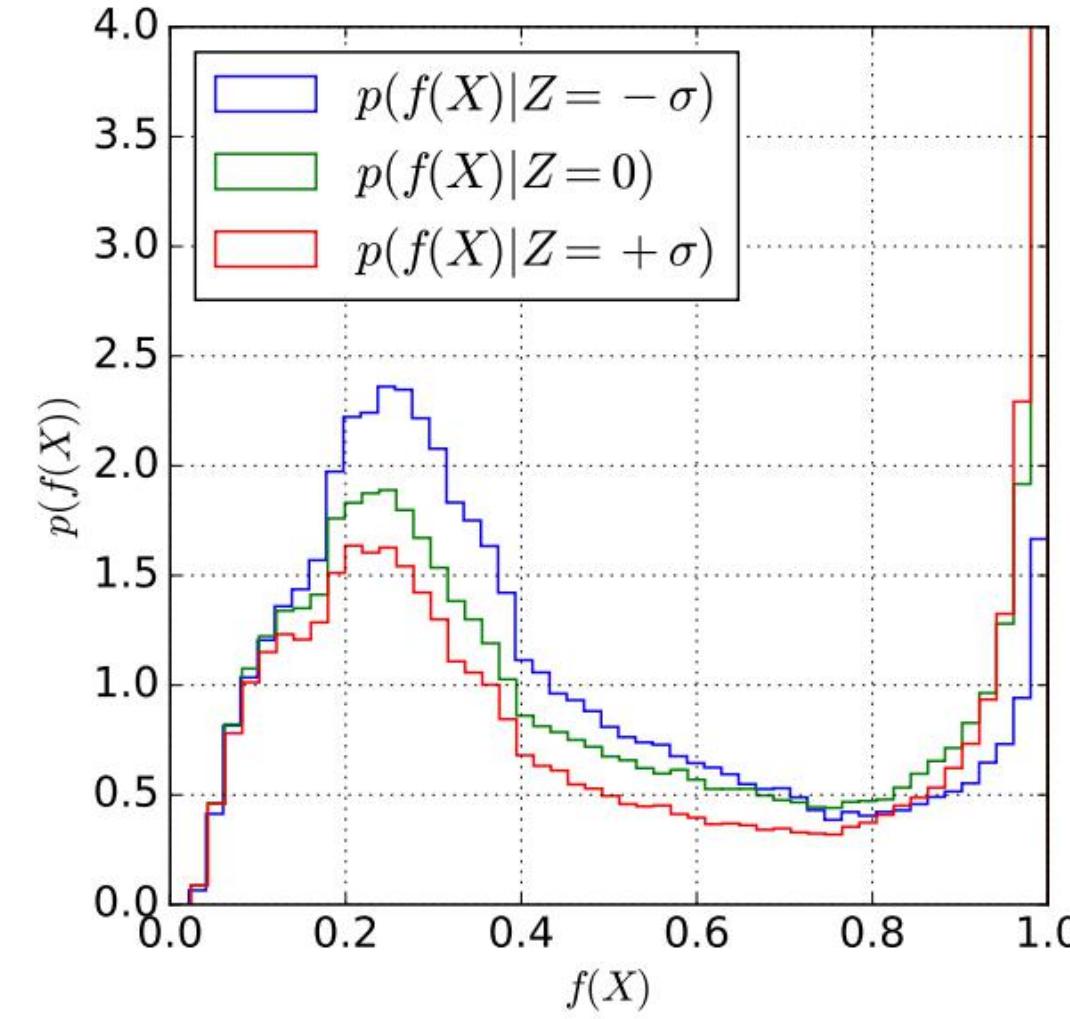


ML for controlling systematics

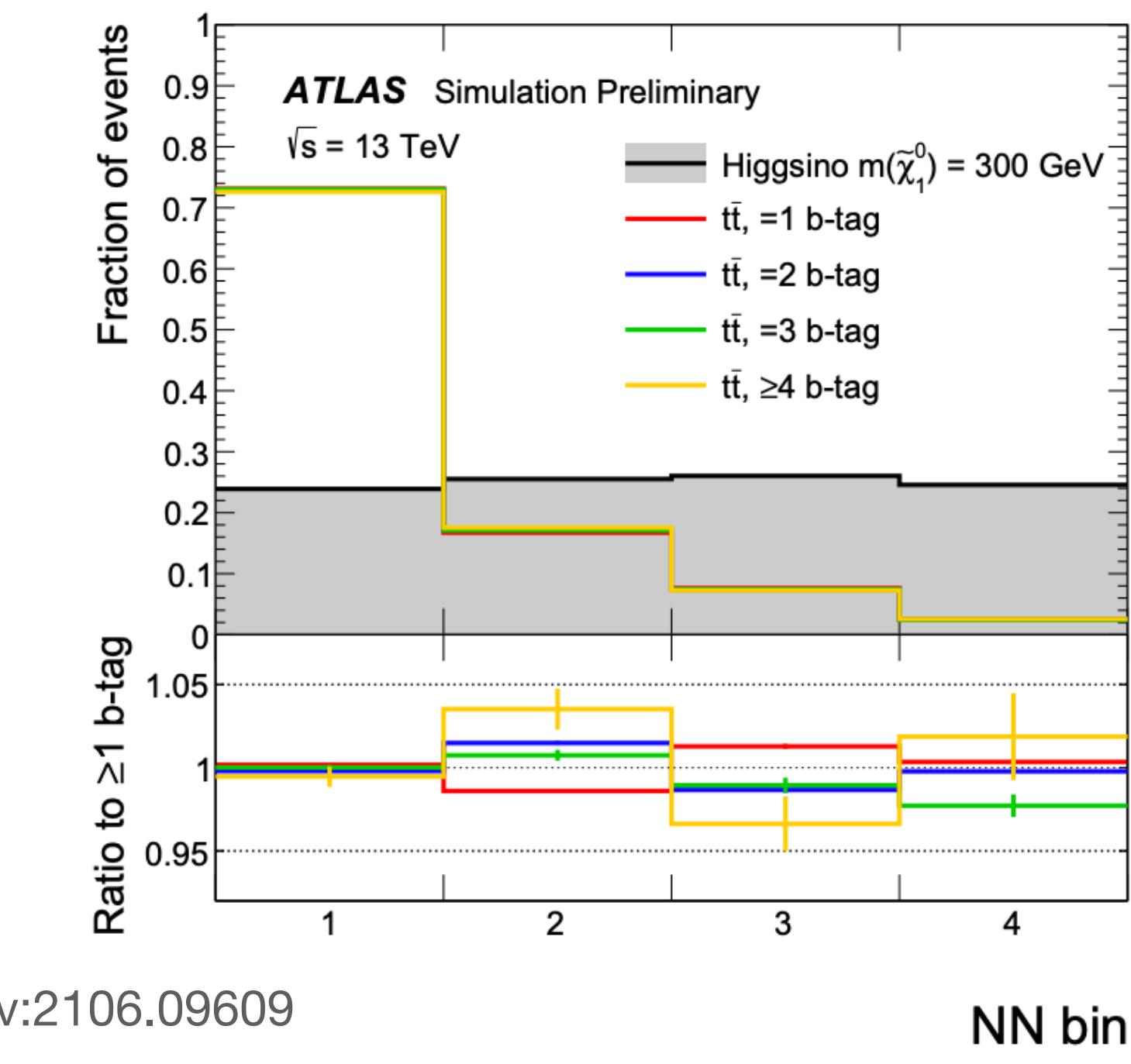
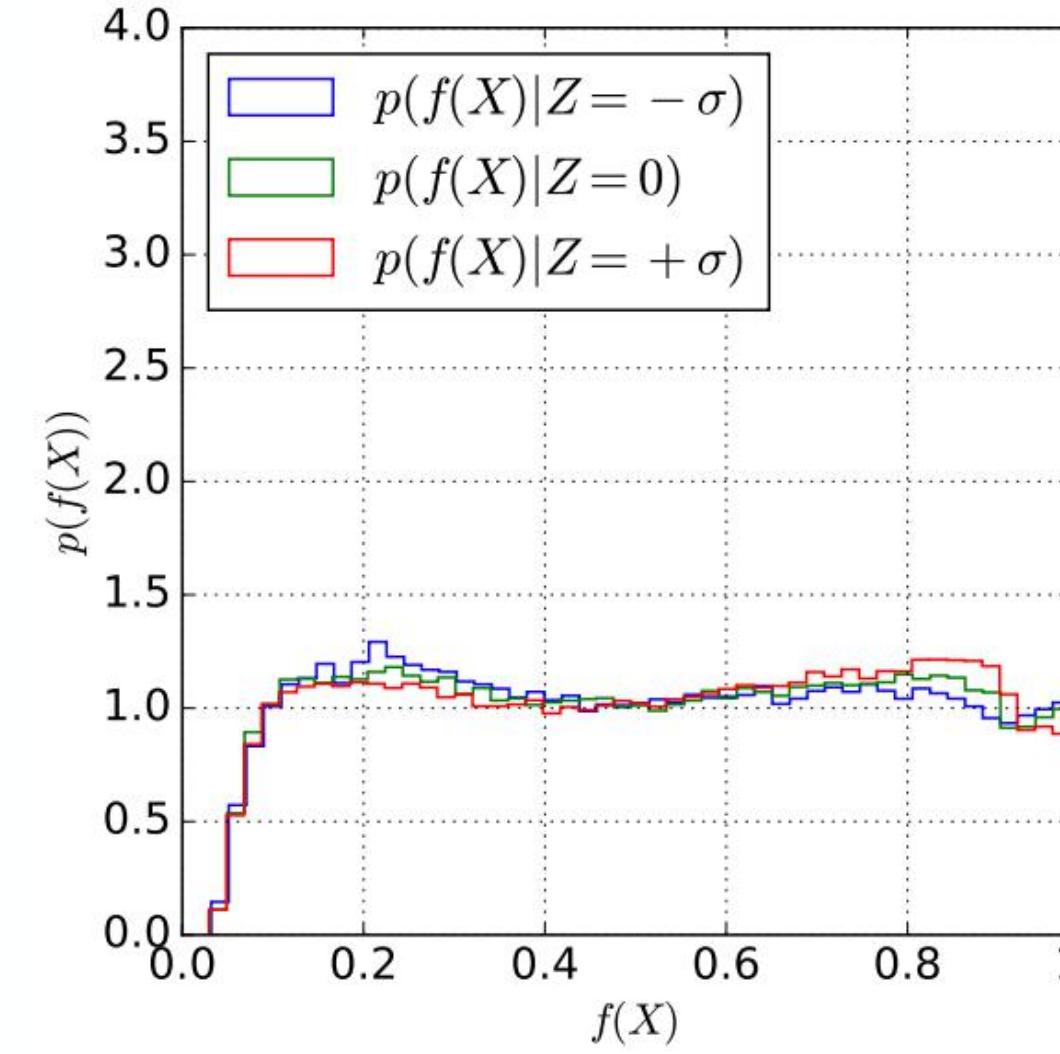
Using ML and “planned degradation” to achieve pivotal projections

Goal: representations of the data that

- **good at separating sig v. bkg**
- **bad at measuring NPs (i.e. invariant)**



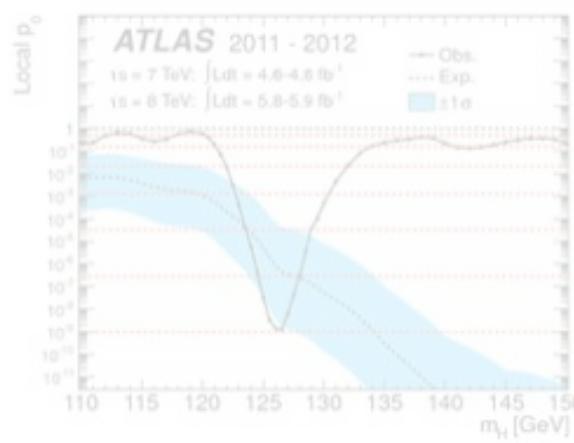
adversarial training



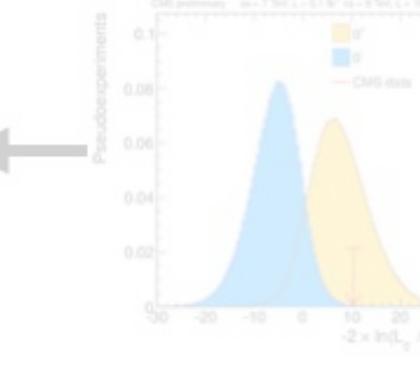
arxiv:2106.09609

arxiv:2001.05310 distance correlation penalty

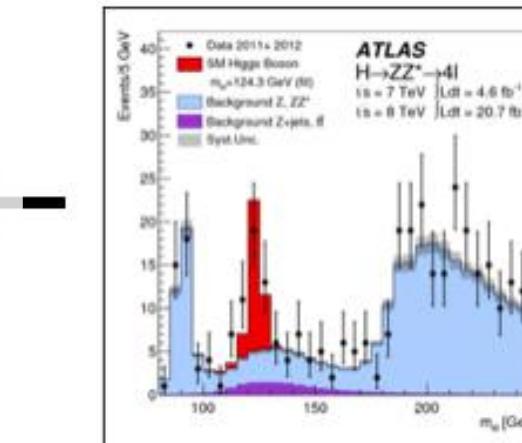
Building the model



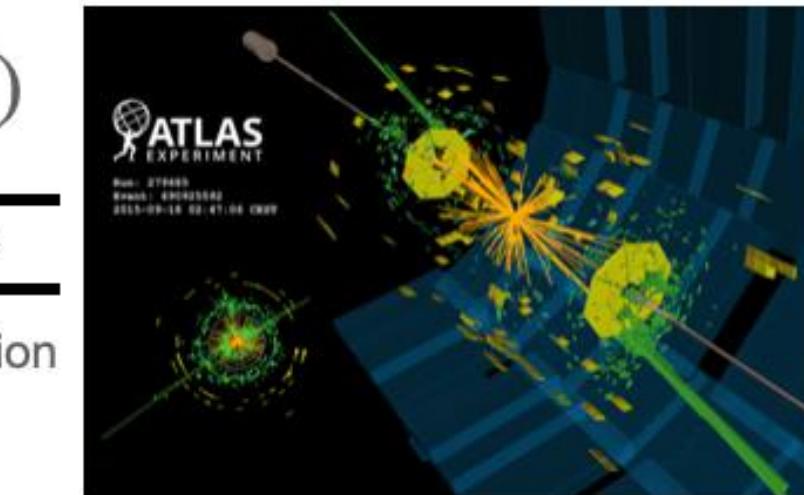
"Likelihood-fu!"
Inference



test statistics
for inference



summary statistics /
density estimation

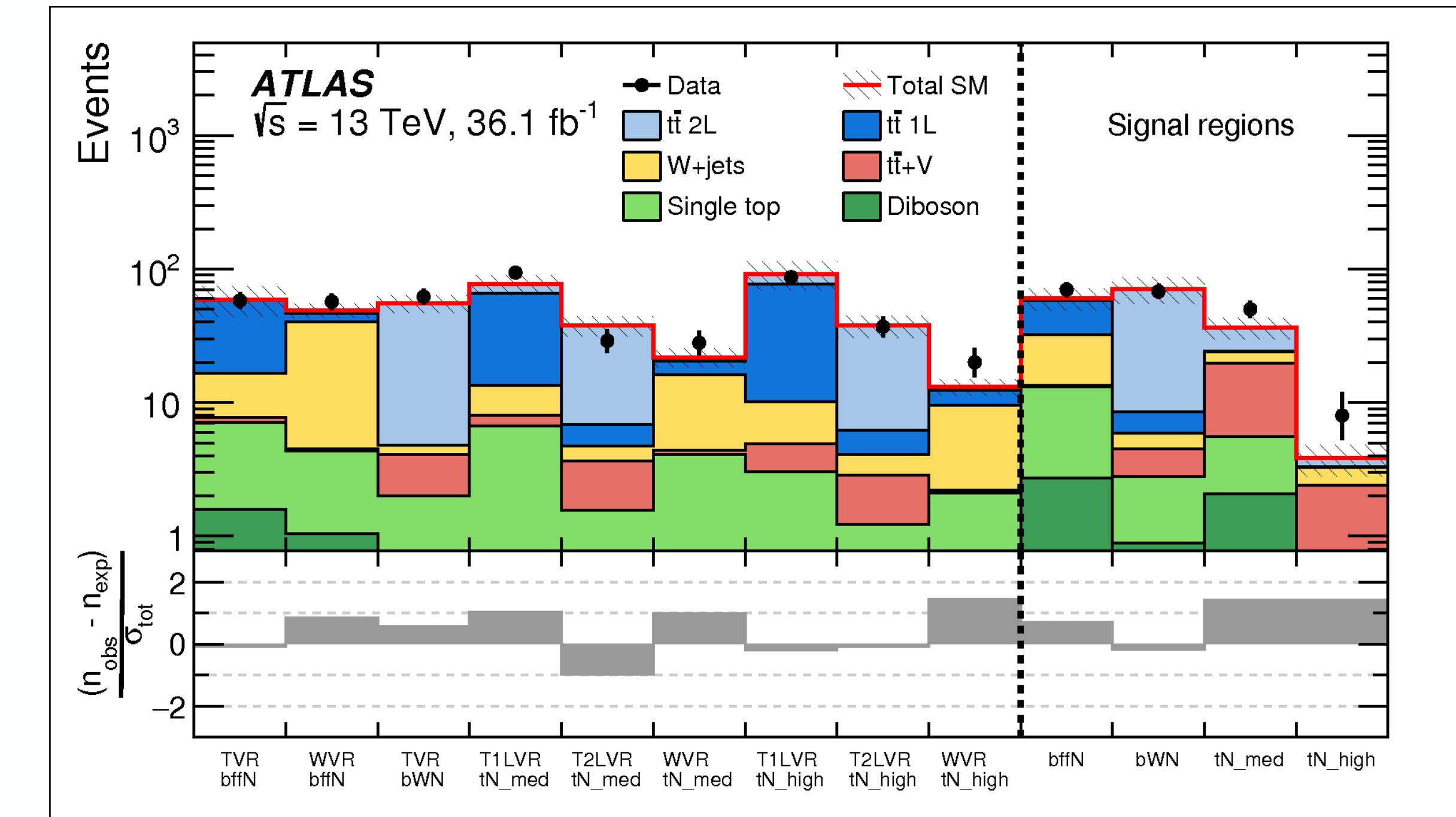
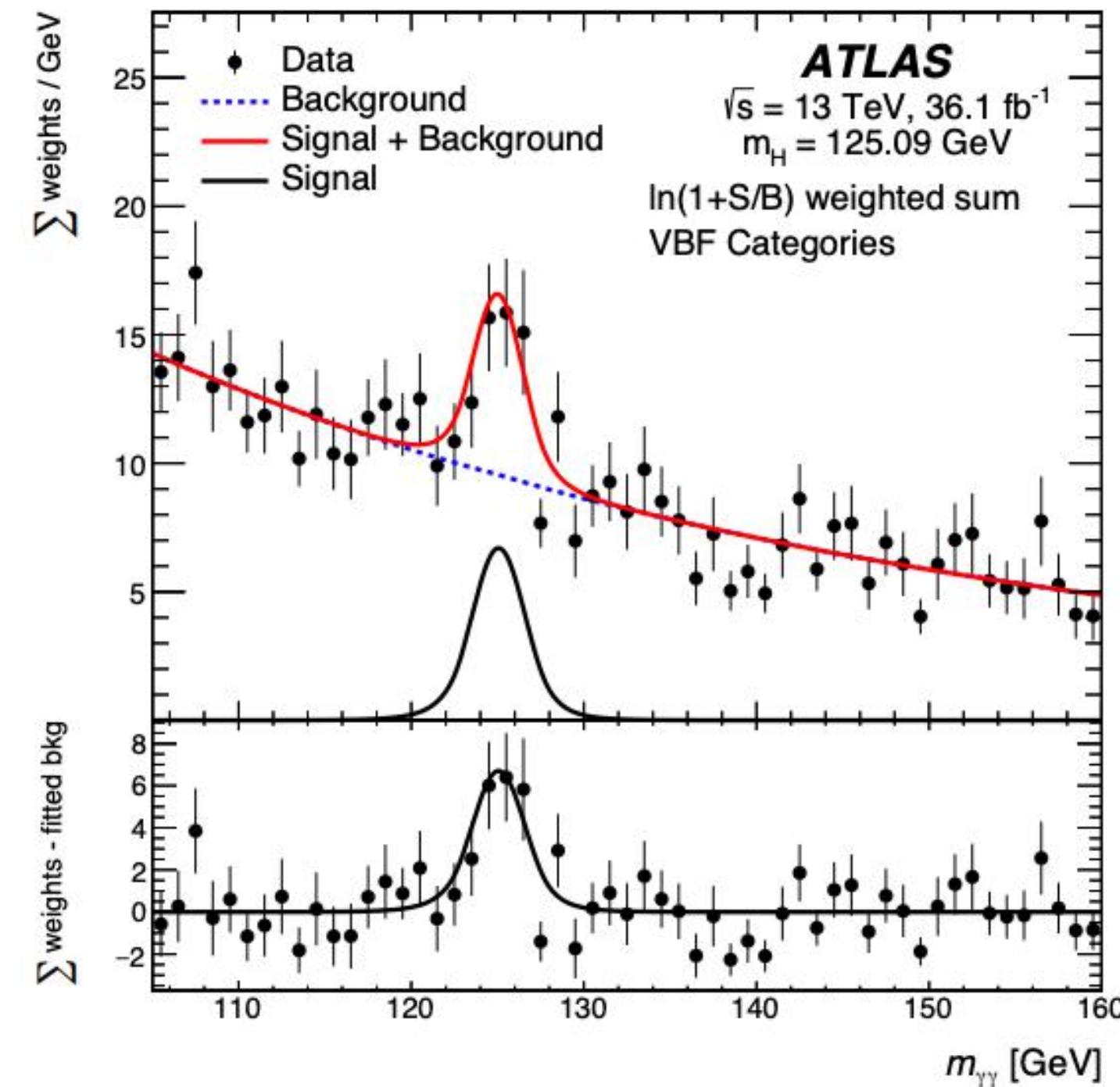


raw data

$s(\cdot)$
data
simulation

Focus on building the low-dim. model

Two broad styles:



$$p(\{x_i\} | \theta) = \text{Pois}(N | \lambda(\theta)) \prod_{i=0}^N p(x_i | \theta)$$

“Unbinned”

$$p(\{n_b\} | \theta) = \prod_b \text{Pois}(n_b | \lambda_b(\theta))$$

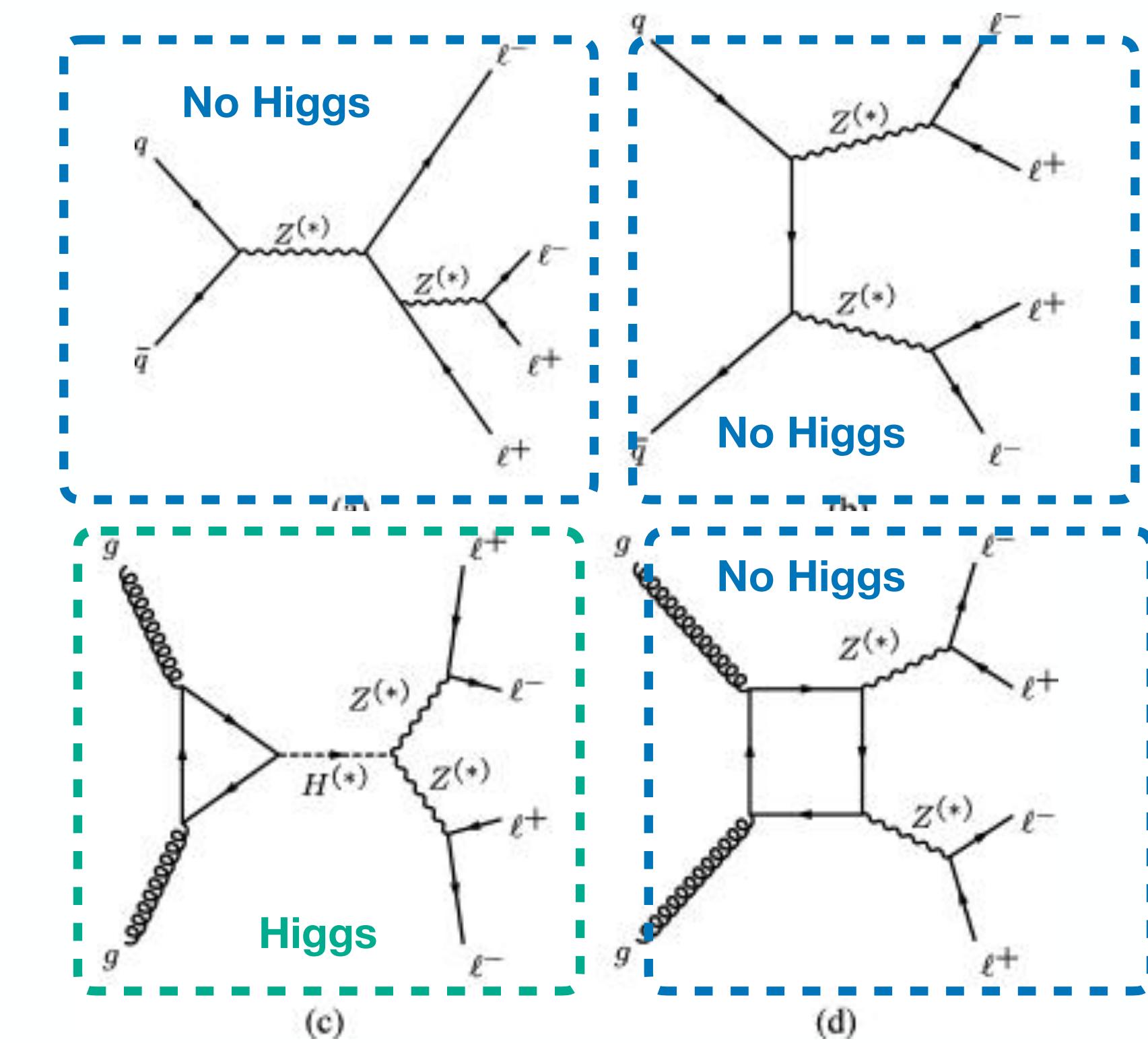
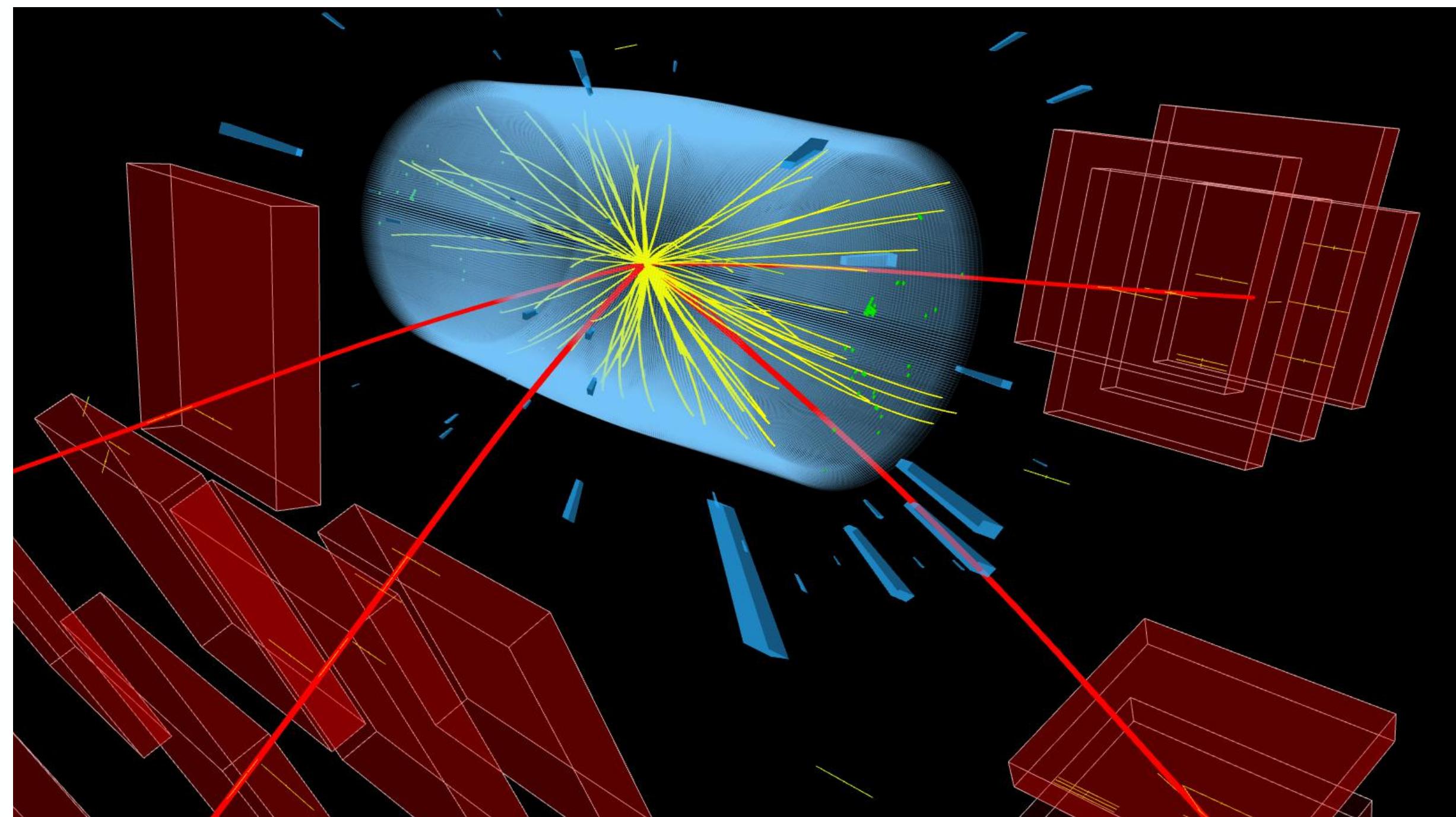
“Binned”

(equivalent to unbinned w/ step-wise constant)

Building the model: Mixture Models

Fundamentally, Quantum Mechanics doesn't allow us to predict, exactly what “process” happens: Particle physics will always be a mixture model.

$$p(x | \theta) = \sum_{\text{proc}} c_i(\theta) p(x | \text{proc}, \theta)$$



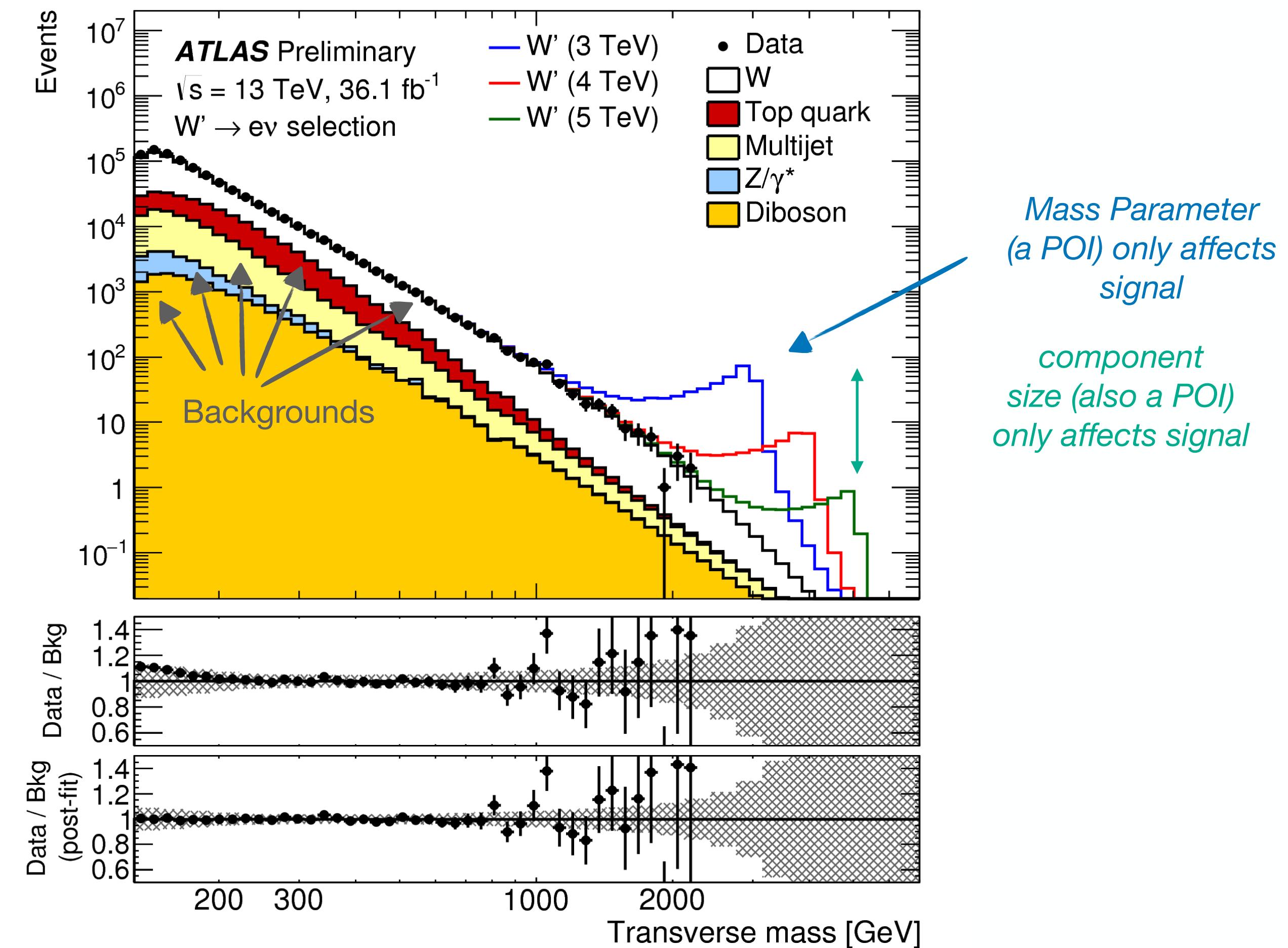
Building the model: Mixture Models

Useful categorization: “signal” and “background” components

$$\text{All} = c_B(\nu)\text{Bkg}(\nu) + c_S(\theta)\text{Sig}(\theta, \nu)$$

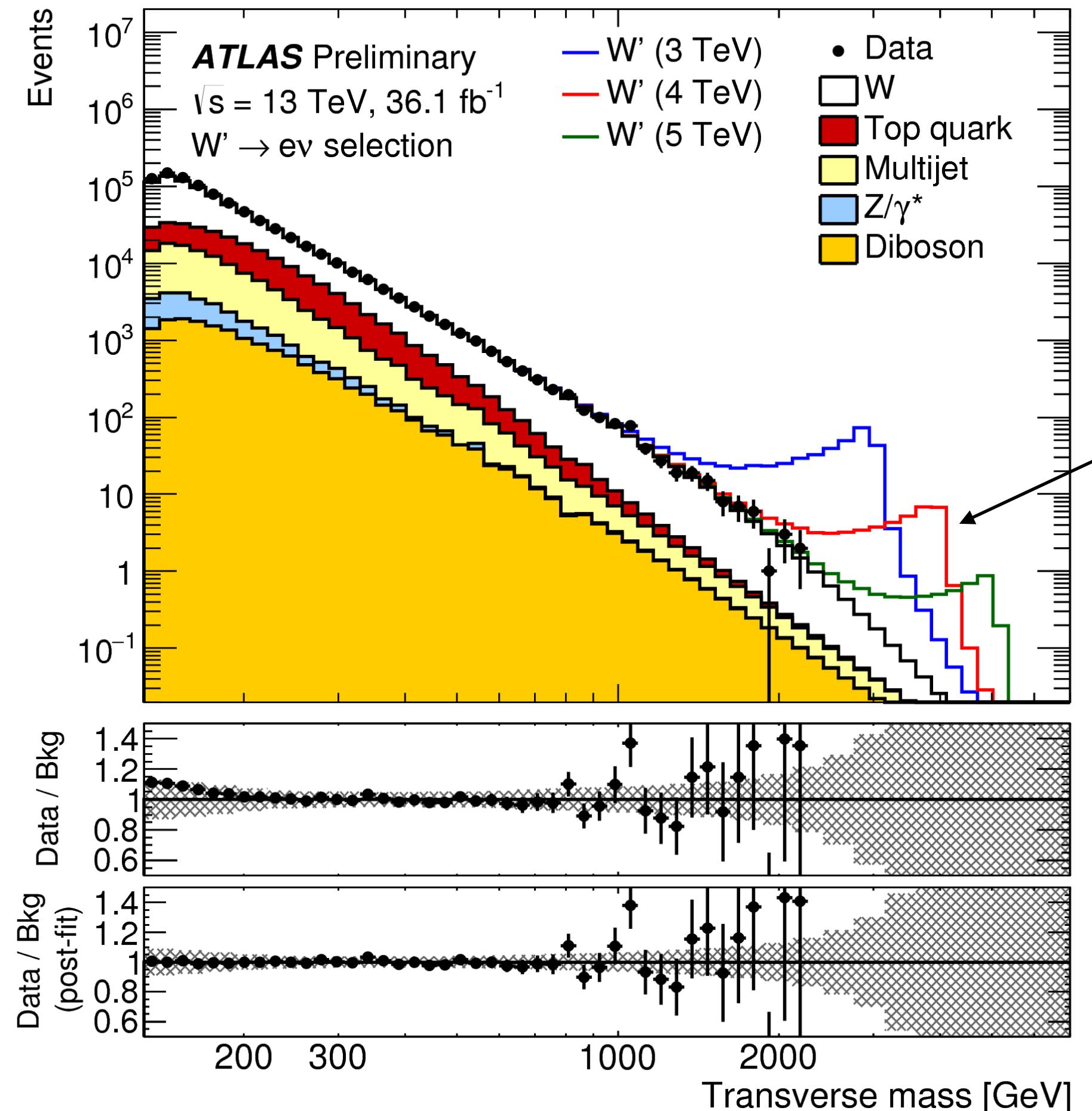
Background:
only affected by nuisance parameter

Signal:
affected by nuisance and POIs



Interpolation

Simulators are very good, but very expensive as well



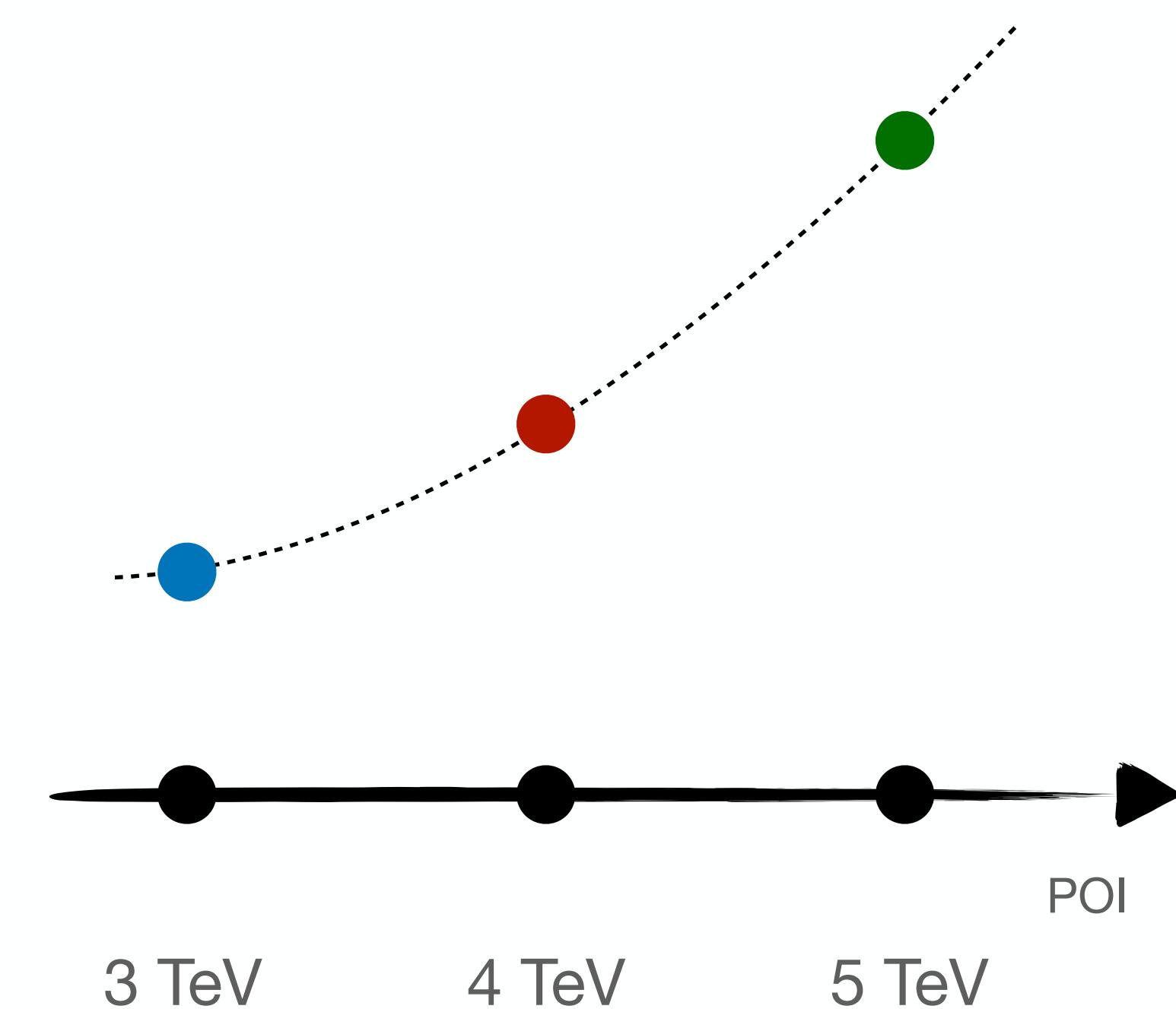
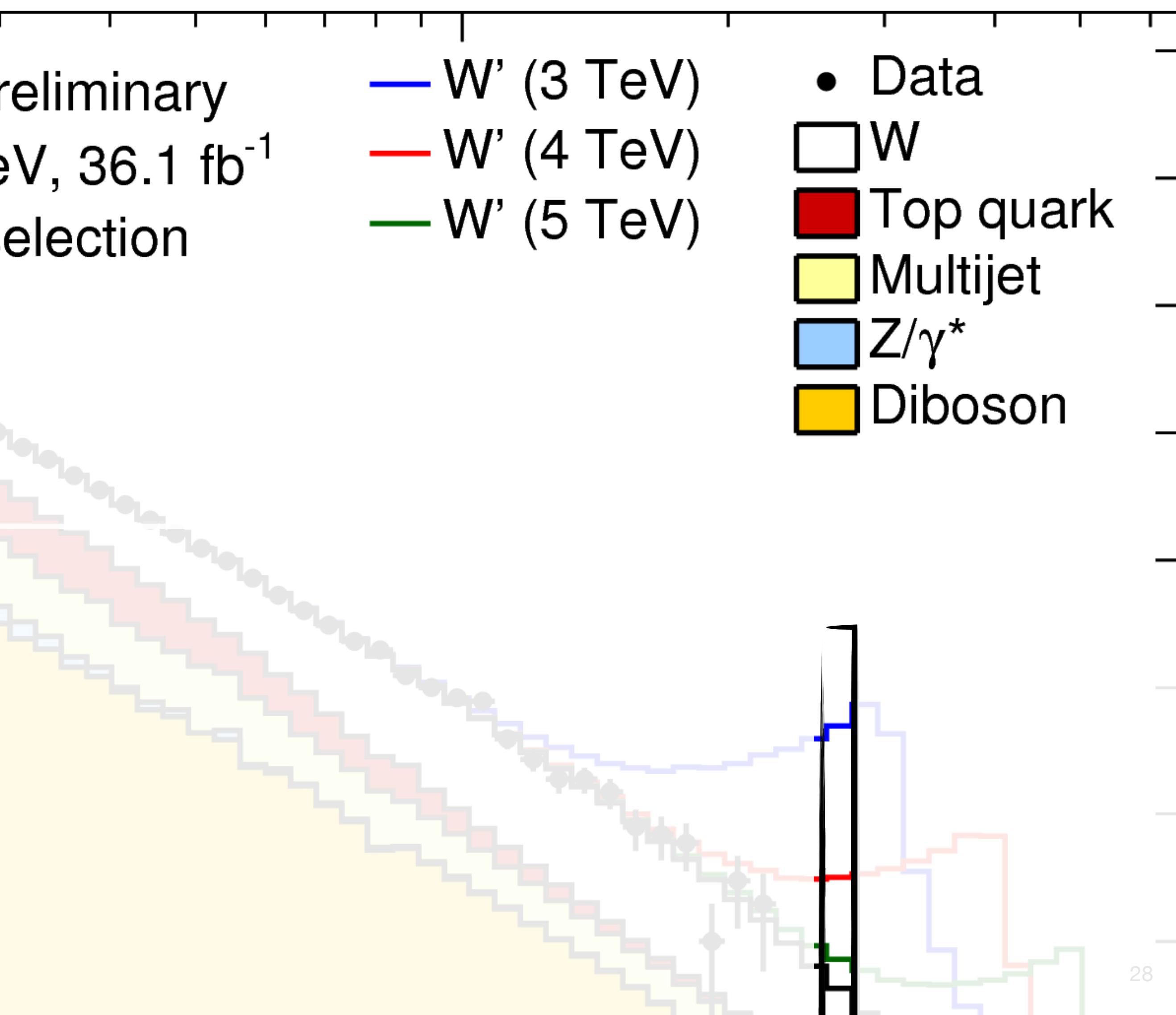
Getting expected rates $\lambda(\theta)$ O(10k) CPUh!

$$\text{Pois}(n_i | \lambda_i(\theta))$$

Cannot be used to during inference
(e.g. MLE fits). Need a faster option.

Interpolation

Simulators are very good, but very expensive as well



Incorporating Prior Information

As HEP analysis is so dominated by NPs, we must incorporate prior info.

- But methodology designed to allow Bayesian OR Frequentist treatment

where does this
prior come from?



$$p(x | \mu, \nu)p(\nu)$$

Incorporating Prior Information

As HEP analysis is so dominated by NPs, we must incorporate prior info.

- But methodology designed to allow Bayesian OR Frequentist treatment

where does this
prior come from?

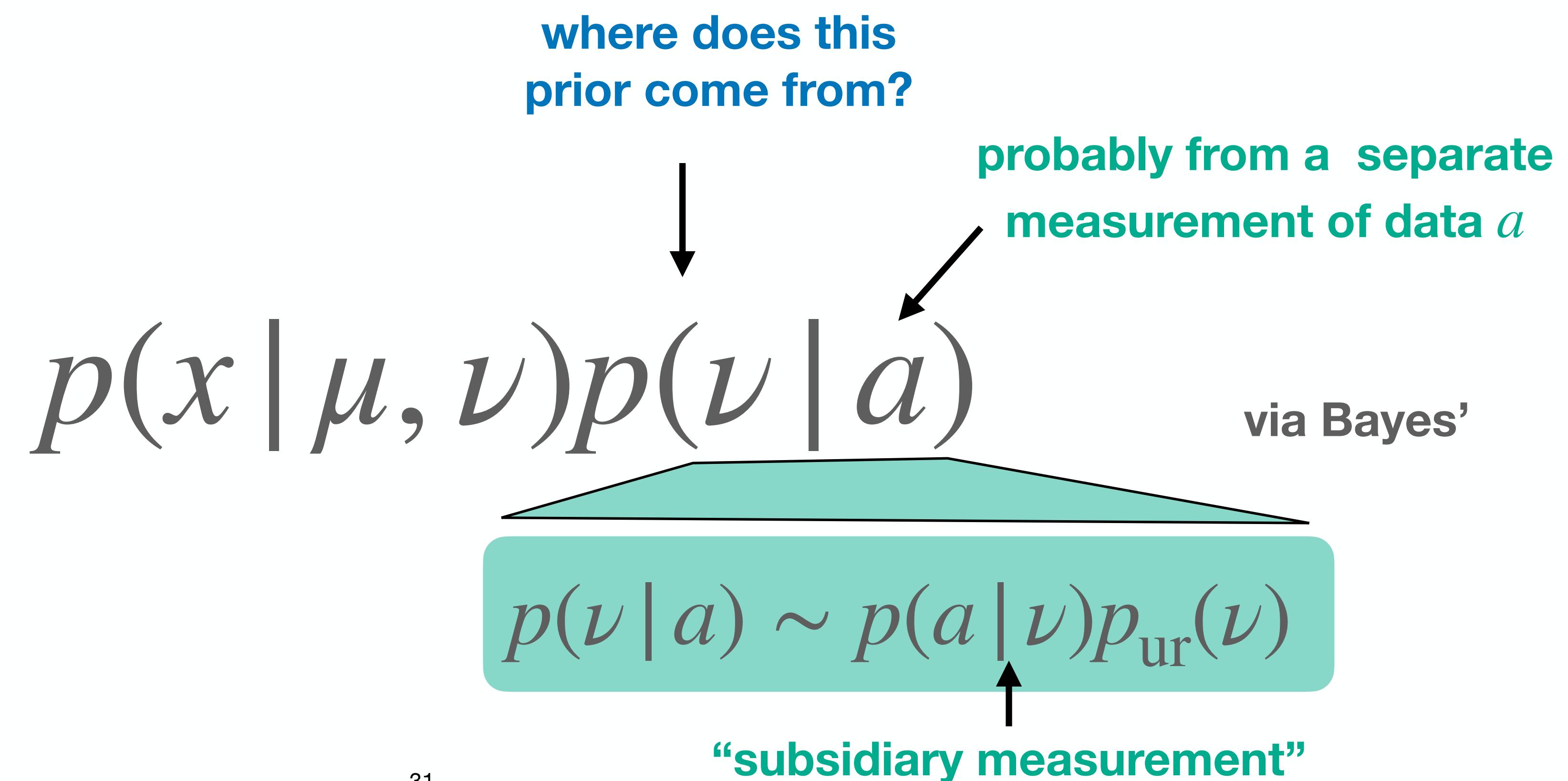
probably from a separate
measurement a

$$p(x | \mu, \nu)p(\nu | a)$$

Incorporating Prior Information

As HEP analysis is so dominated by NPs, we must incorporate prior info.

- But methodology designed to allow Bayesian OR Frequentist treatment



Incorporating Prior Information

Can incorporate same prior information by collecting
“subsidiary measurements” $p(a | \nu)$

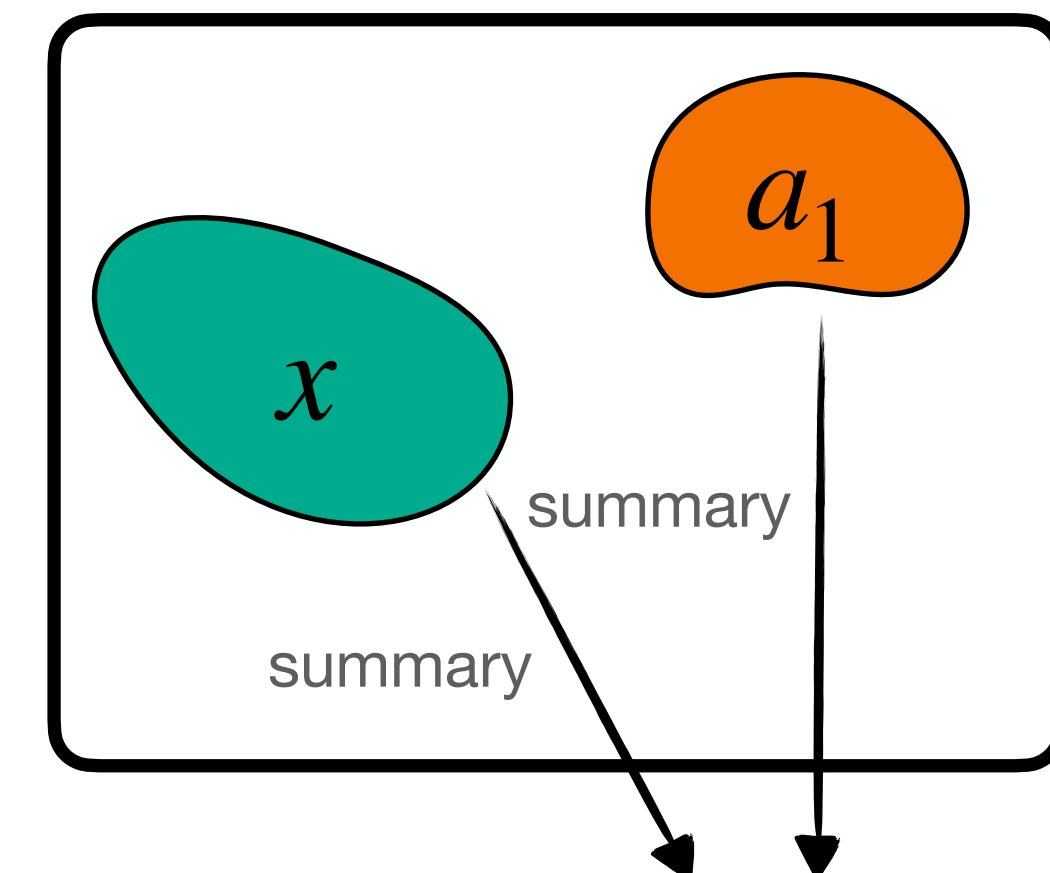
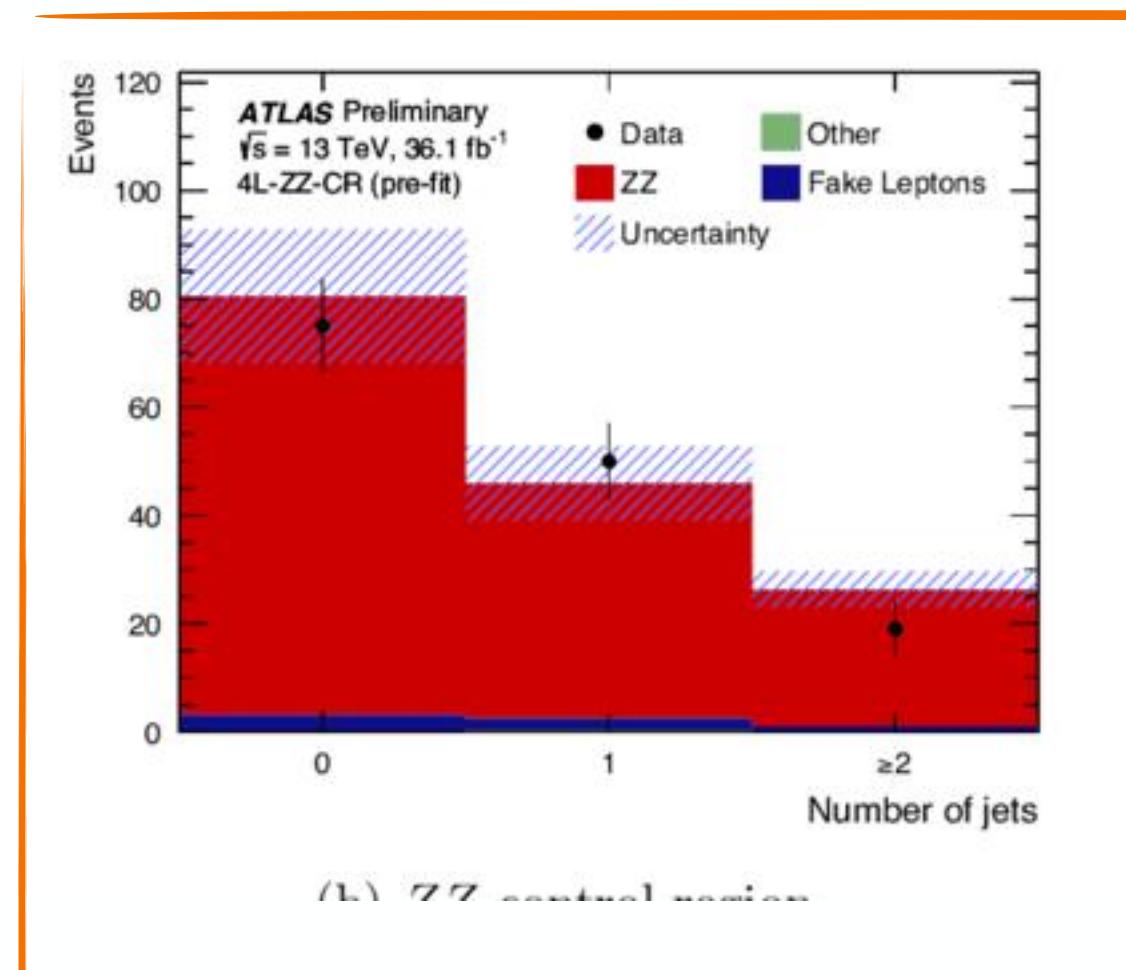
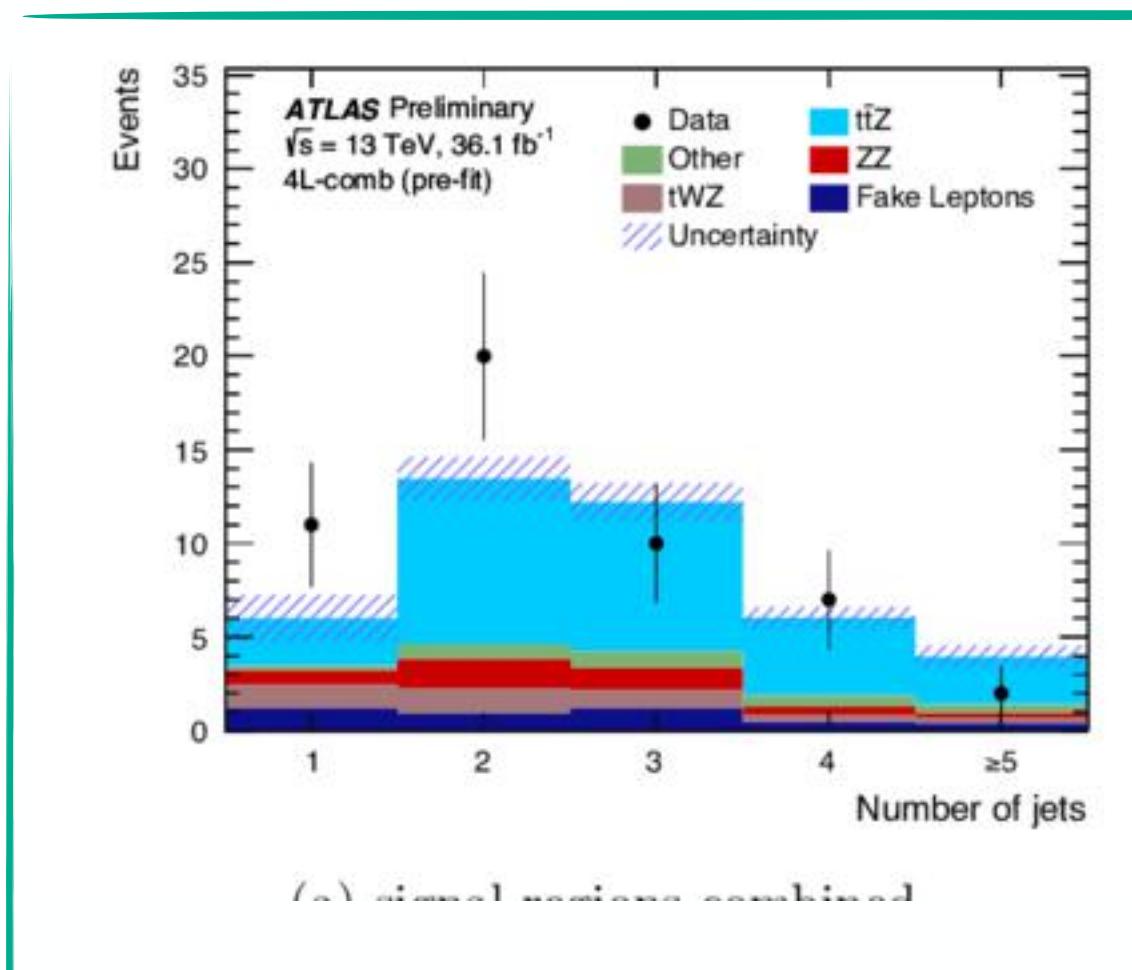
$$p(x | \theta) \rightarrow p(x, a | \theta) = p(x | \mu, \nu)p(a | \nu)$$

Next Question: what are these subsidiary likelihoods?

Option 1: “The Gold Standard”

Perform a full separate measurement (i.e. derive a full fledged model)

Example: constaining background
by measuring a “control region”



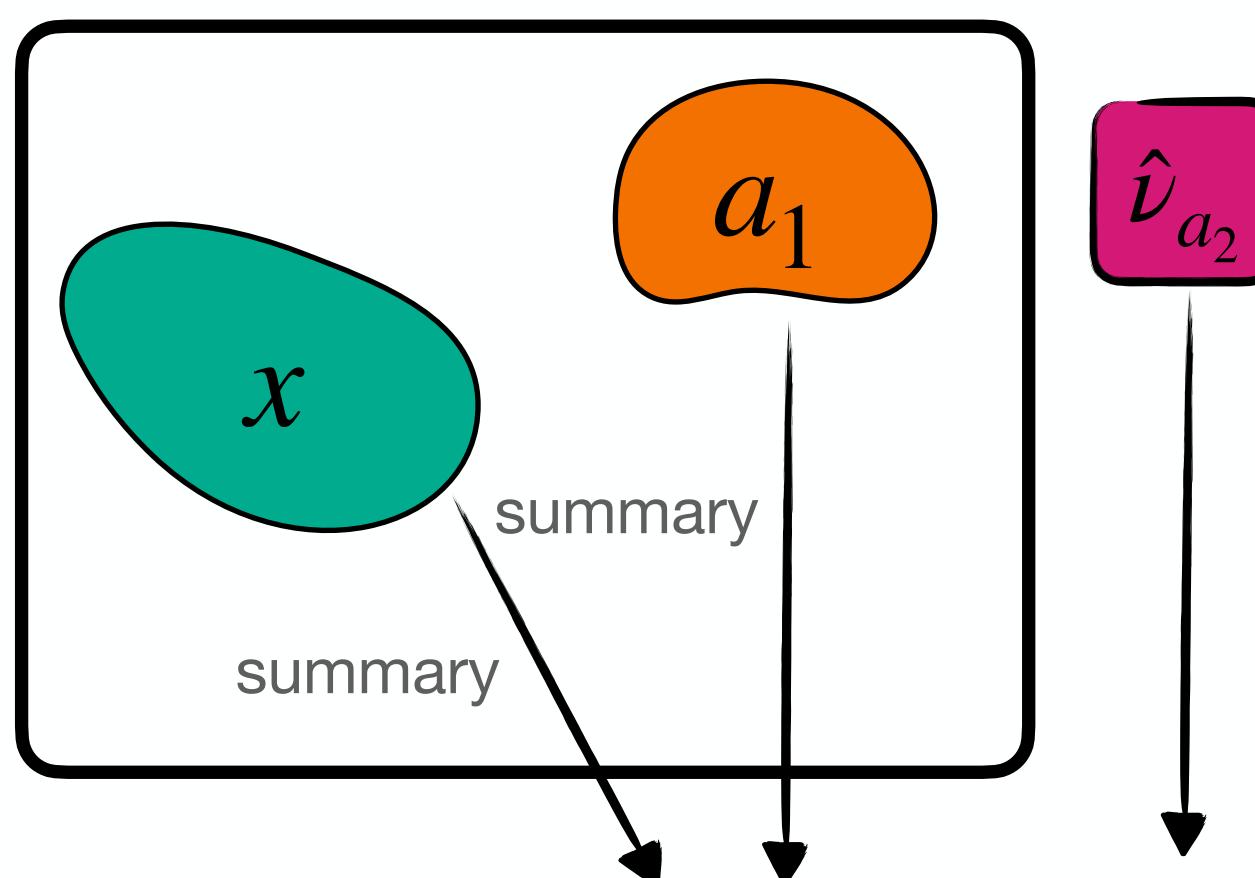
$$p(x, a | \theta) = p(x | \mu, \nu)p(a_1 | \nu)$$

A lot of work to get a full $p(a | \nu)$

Option 2: “Simplified Subsidiaries”

If no $p(a | \nu)$ is not available, or too complex, approximate contribution via

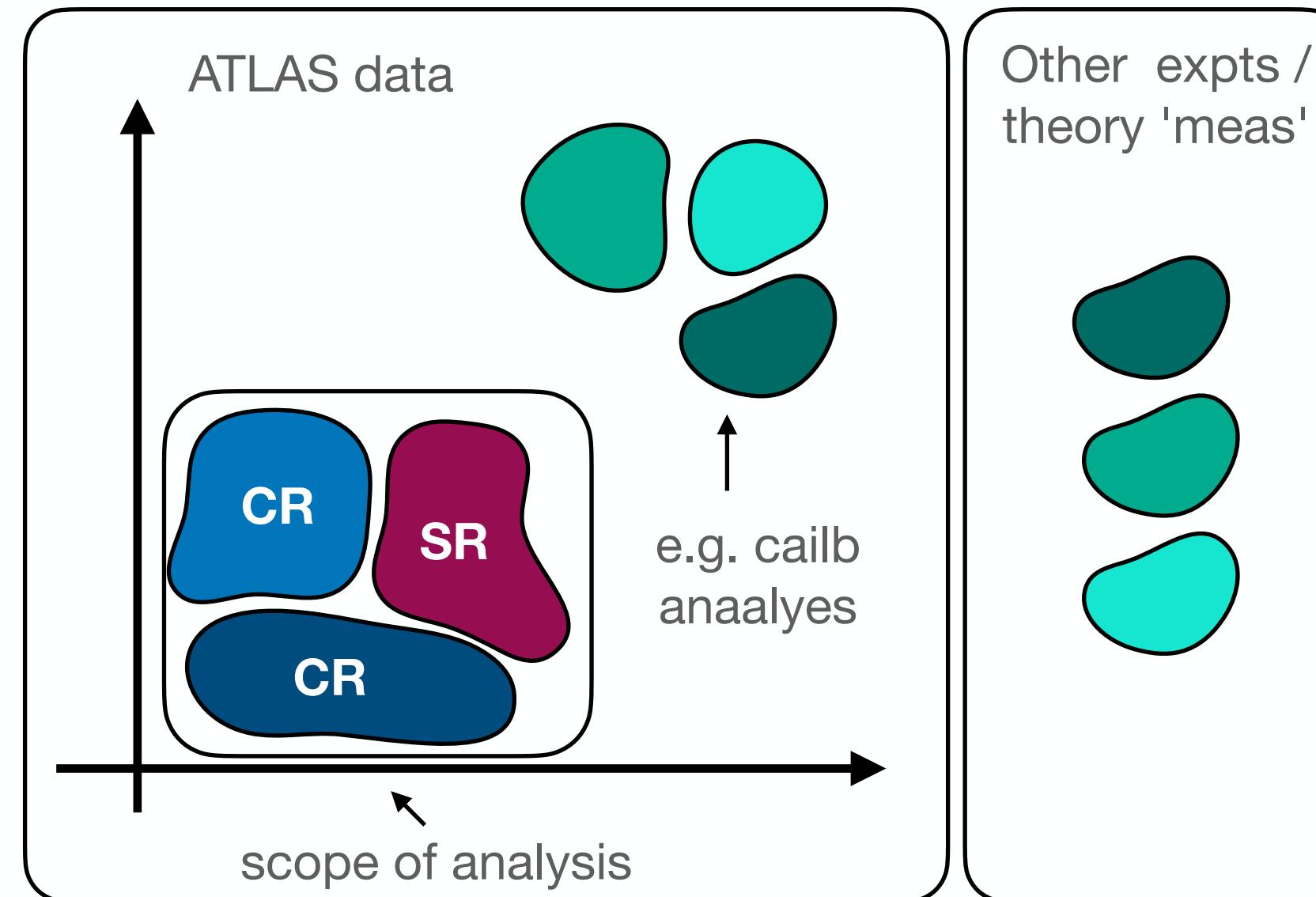
$$p(a | \nu) \sim p(\hat{\nu}(a) | \nu) \sim \text{Gaus}(\hat{\nu}_a | \nu, \hat{\sigma})$$



$$p(x, a | \theta) = p(x | \mu, \nu)p(a_1 | \nu)G(\hat{\nu}_{a_2} | \nu)$$

Putting it all together:

Realistic Analyses use a mixture. A lot of analyses covered by this approach.



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

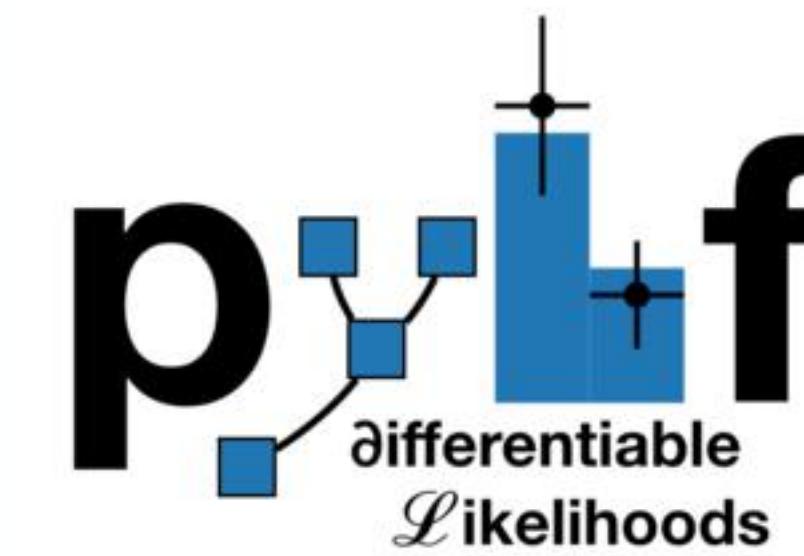
$$\nu_{cb}(\phi) = \sum_{s \in \text{samples}} \nu_{scb}(\boldsymbol{\eta}, \chi) = \sum_{s \in \text{samples}} \left(\underbrace{\prod_{\kappa \in \kappa} \kappa_{scb}(\boldsymbol{\eta}, \chi)}_{\text{multiplicative modifiers}} \right) \left(\underbrace{\nu_{scb}^0(\boldsymbol{\eta}, \chi) + \sum_{\Delta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

sum over components in mixture mode

"parametrized interpolation"

nominal rates

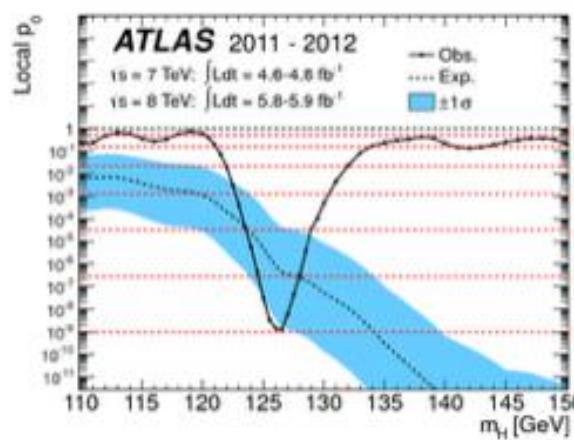
Public Tools available:



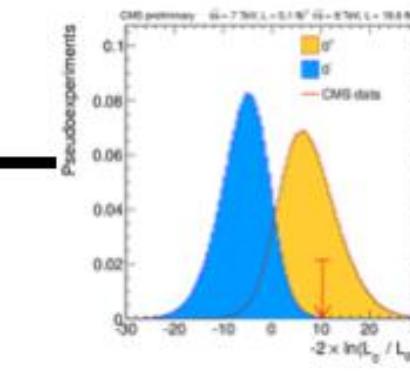
[LH, M. Feickert, G. Stark]

github.com/scikit-hep/pyhf

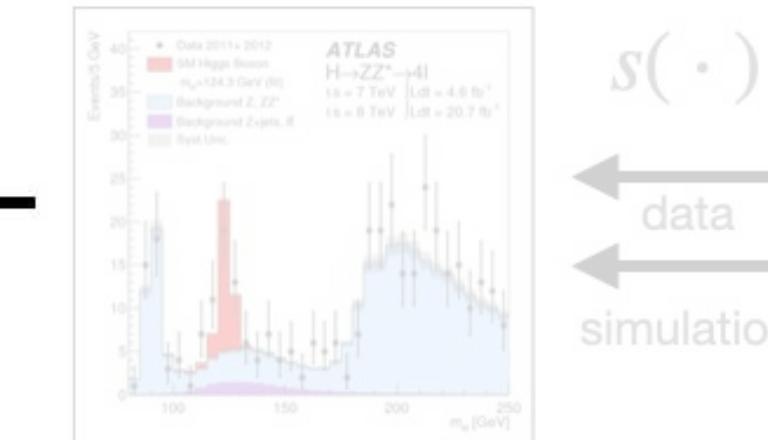
Statistical Inference



“Likelihood-ful”
Inference



test statistics
for inference



summary statistics /
density estimation



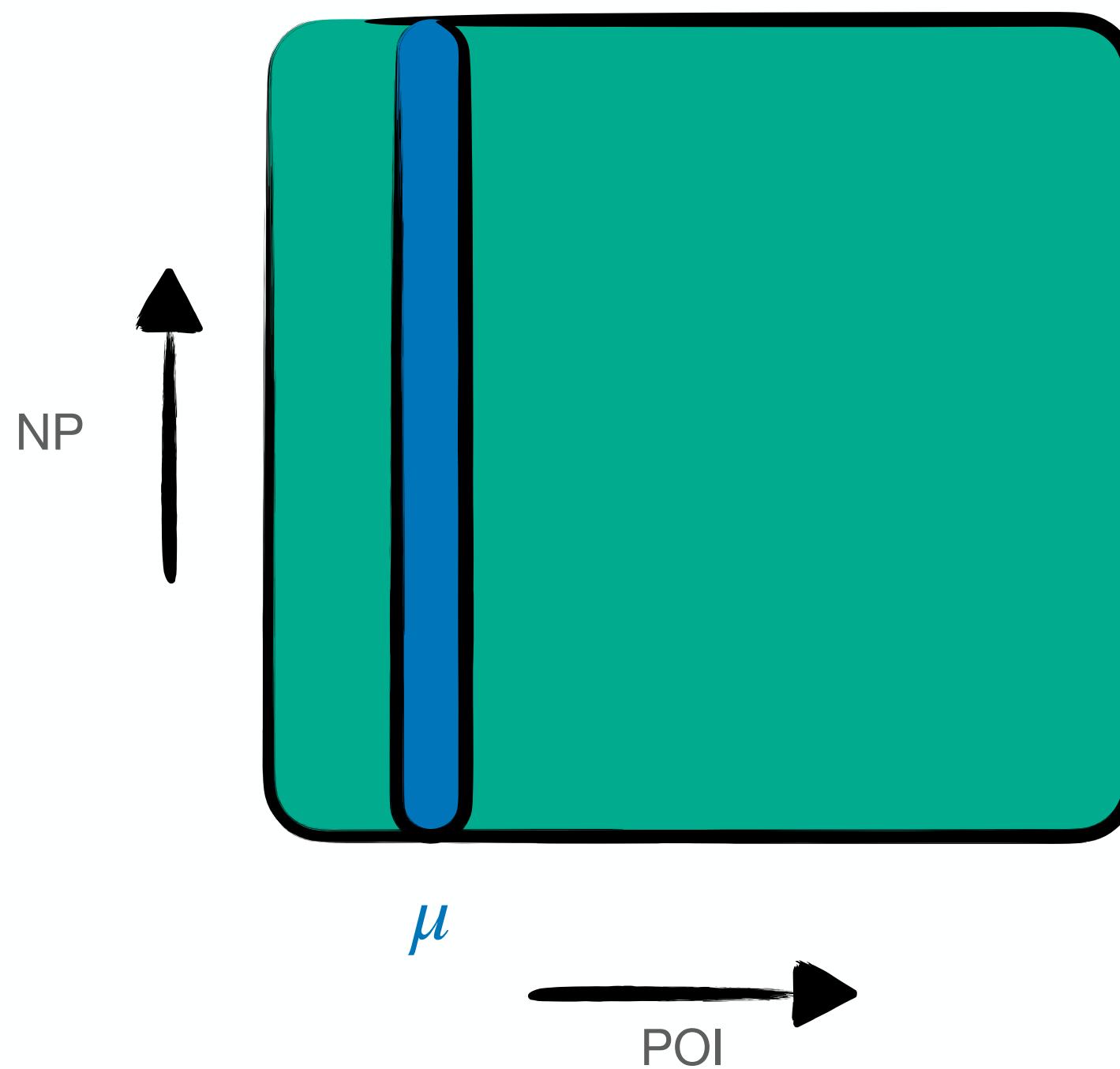
raw data

$s(\cdot)$
data
simulation

Nested Models

Split out parameter space into hypothesis spaces.

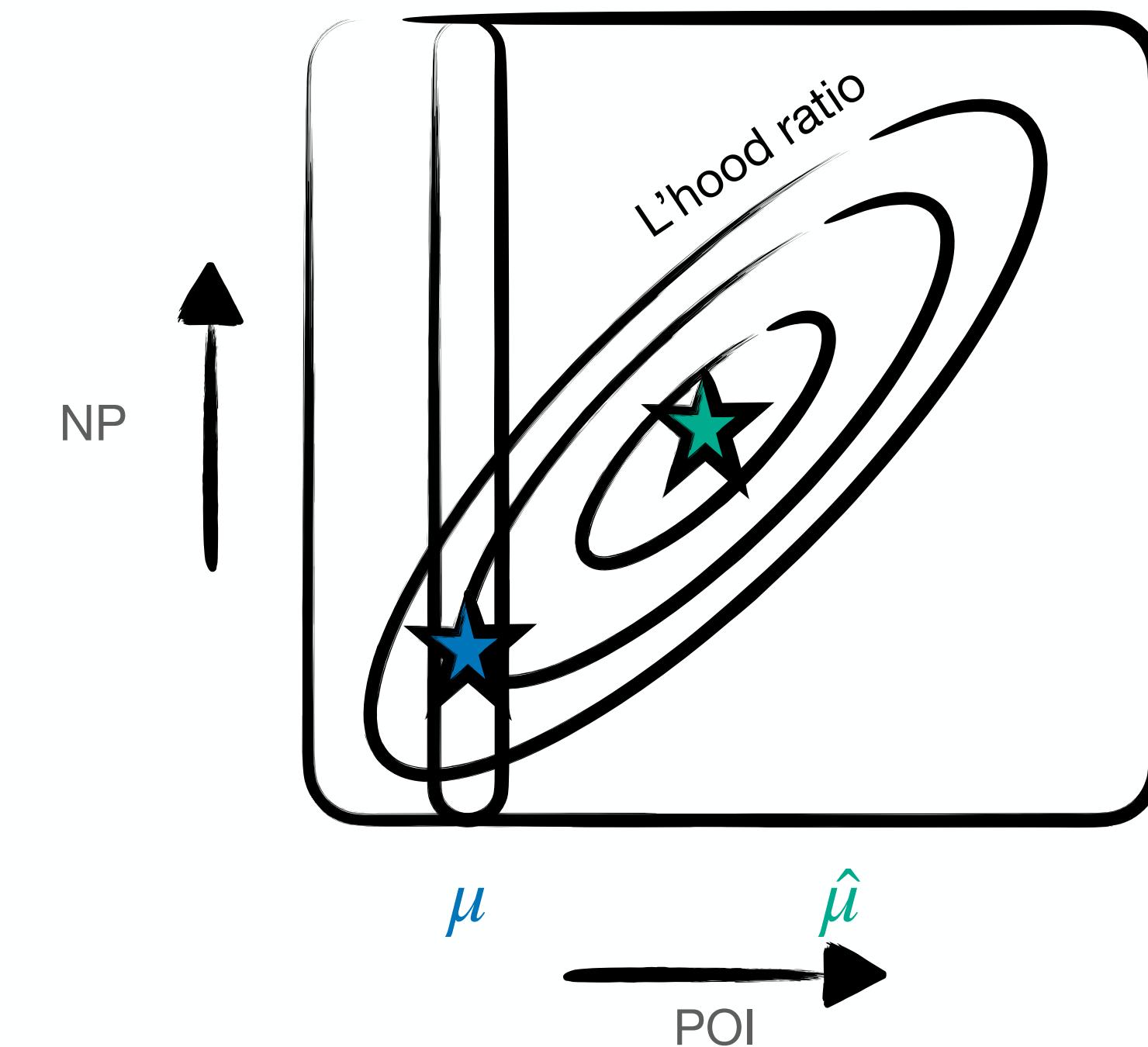
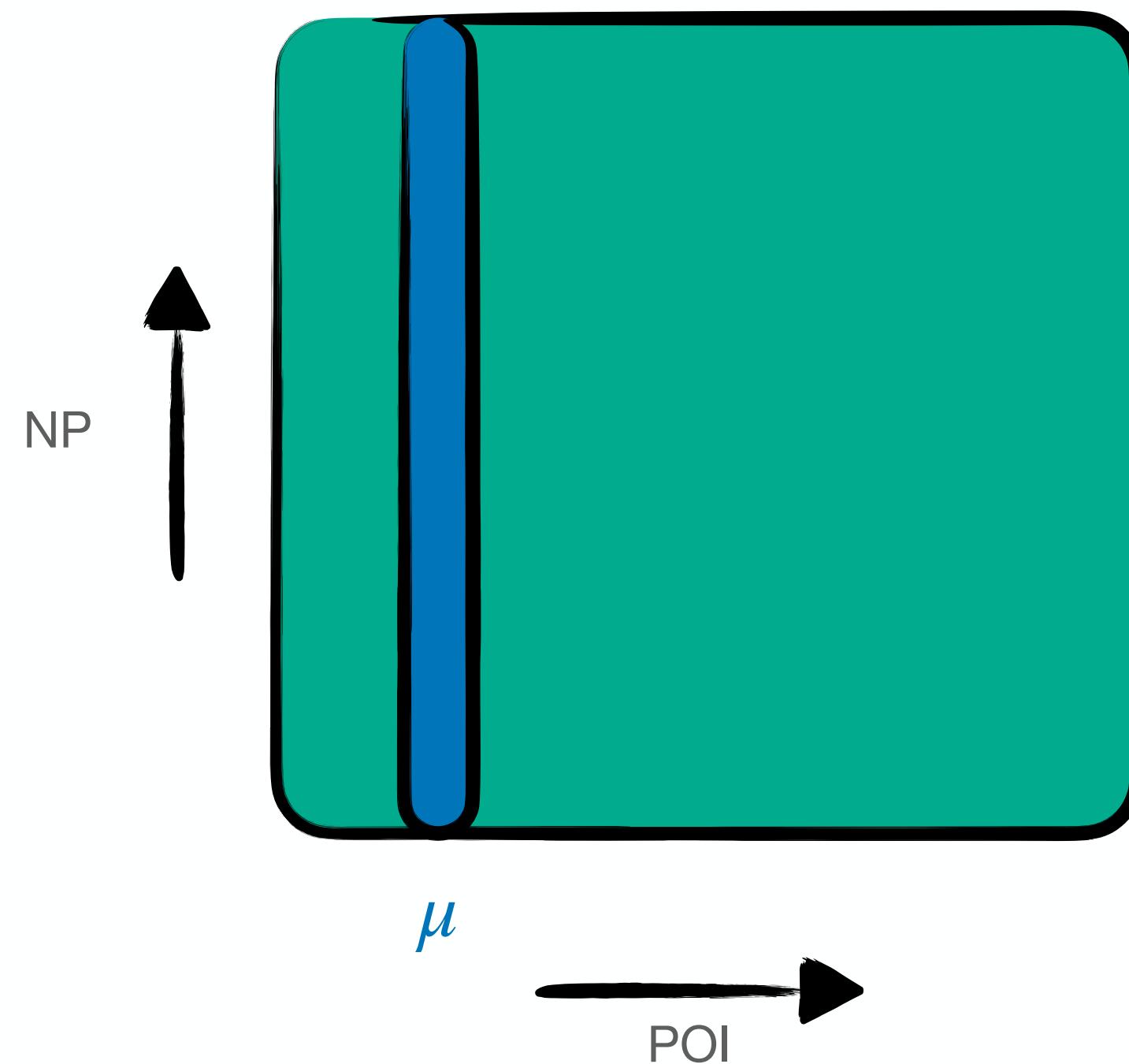
- Reminder: not interested in inference on NPs
- Null hypotheses only indexed by POIs



Nested Models

After decades of iteration, HEP landed on classic Likelihood Ratio Test

- “profile likelihood ratio”



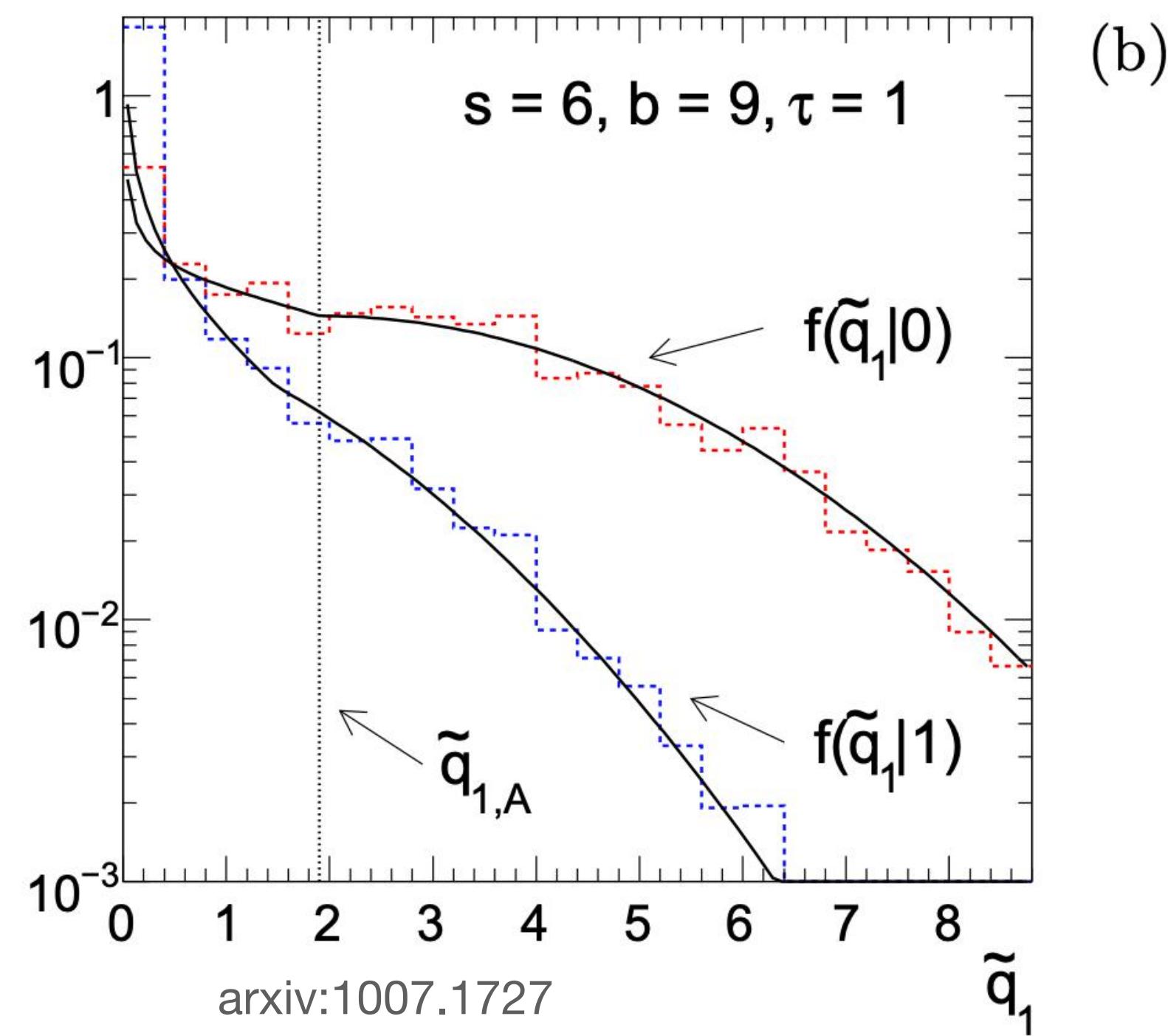
**More Stats <>> Physics connections
would have helped as a lot!**

$$t(x|\mu) = -2 \log \frac{p(x|\mu, \hat{\nu})}{p(x|\hat{\mu}, \hat{\nu})}$$

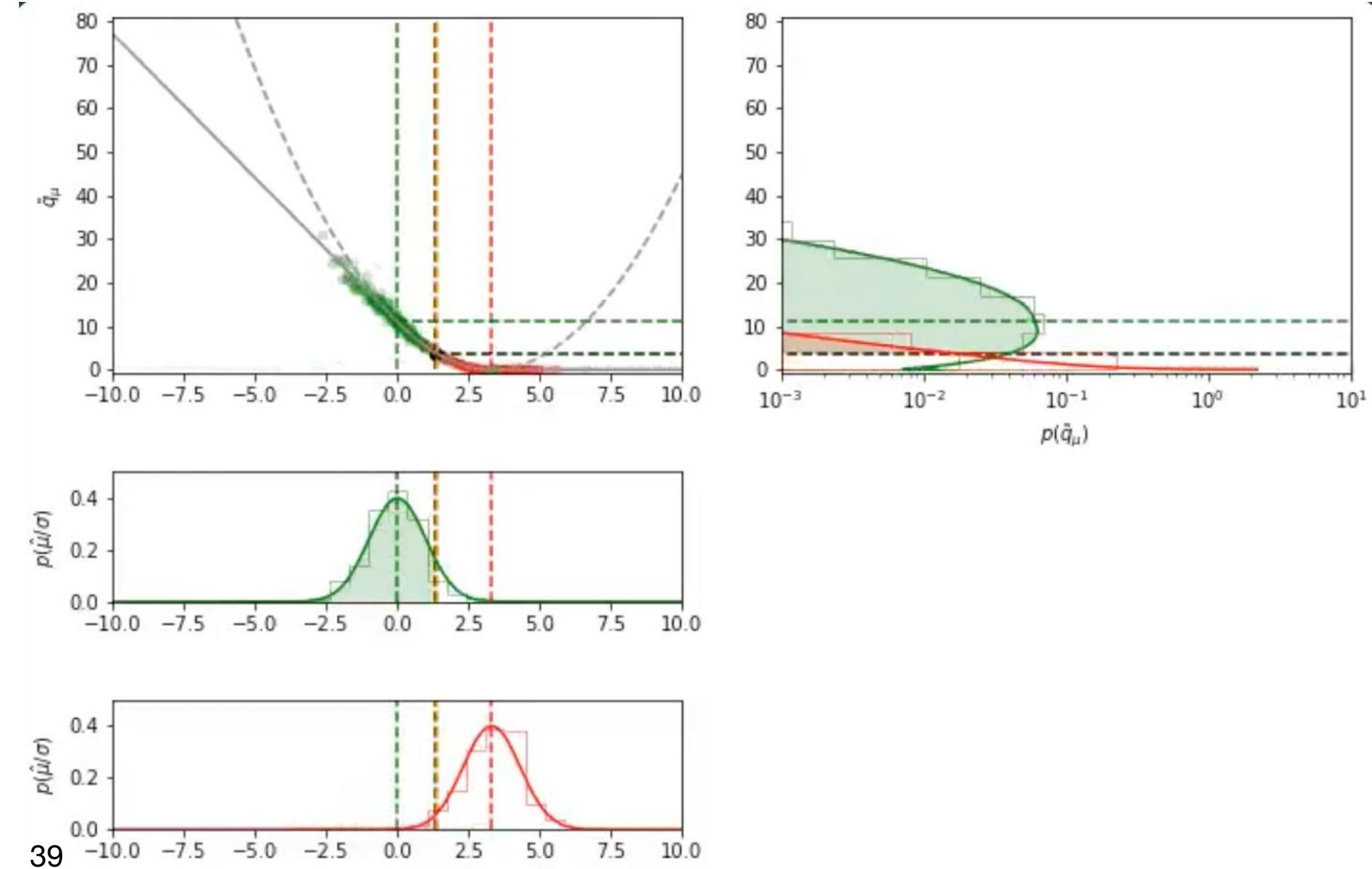
Use “best-in-class” parameter
for each hypothesis set for ratio

Asymptotic Calculations

- Many analyses rely heavily on asymptotic shapes of test stat distributions
- Assumption: Wilks & Wald hold (cf. see next slide)



(b)



Interval construction

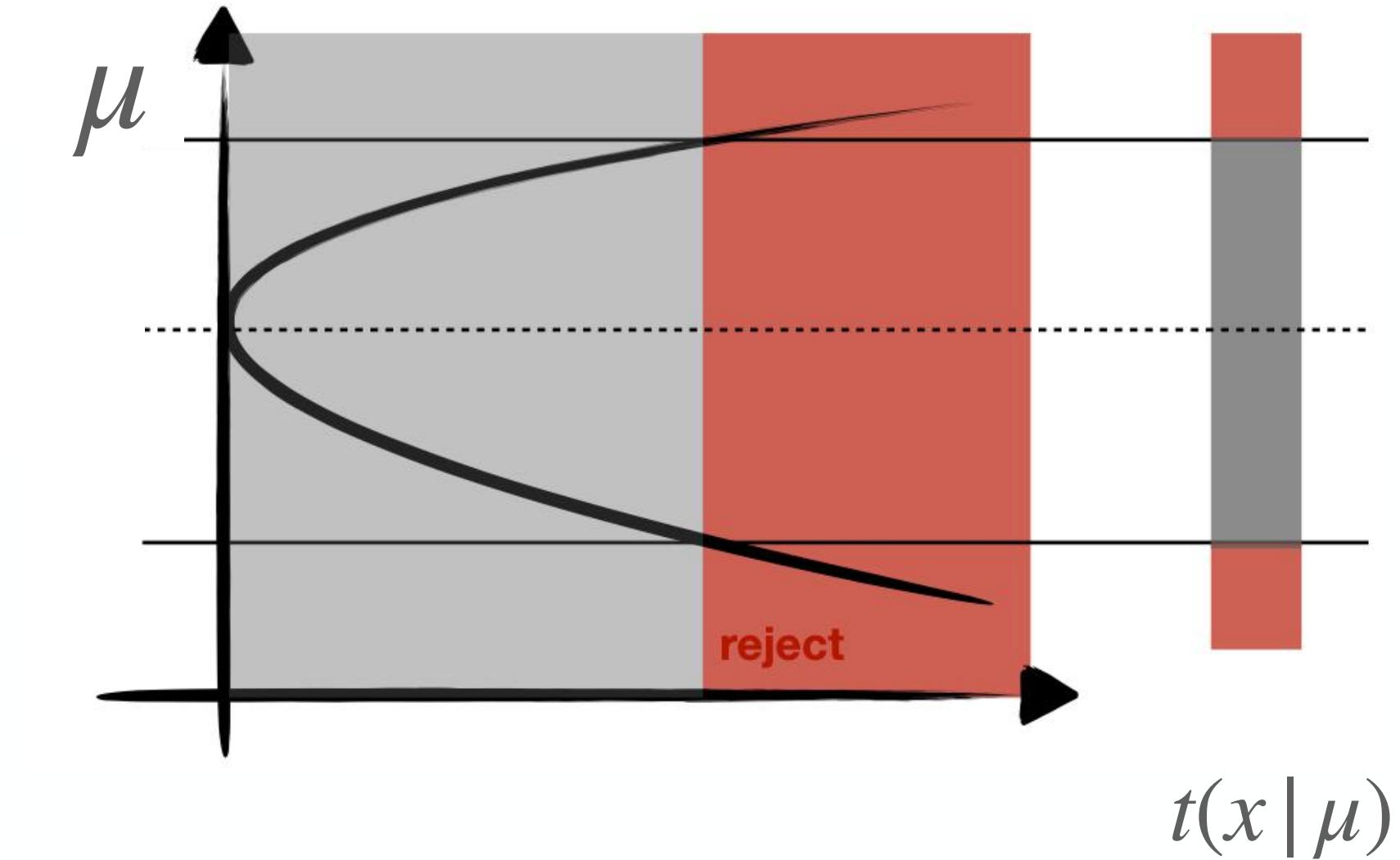
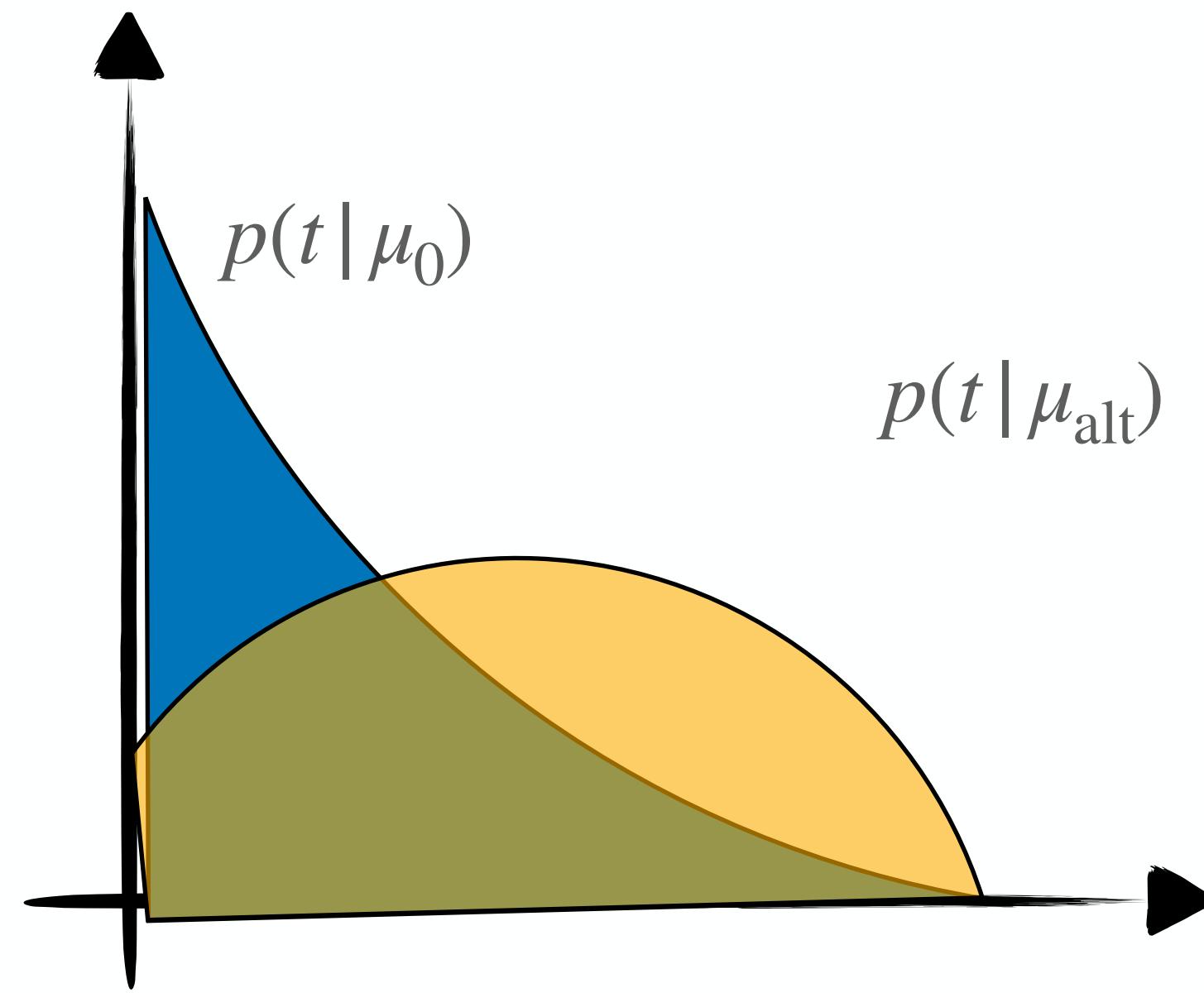
HEP fairly principled: Intervals strictly via exact Neyman-Construction.

Widely assumed, rarely checked

$$p(t | \mu, \nu) = p(t | \mu)$$

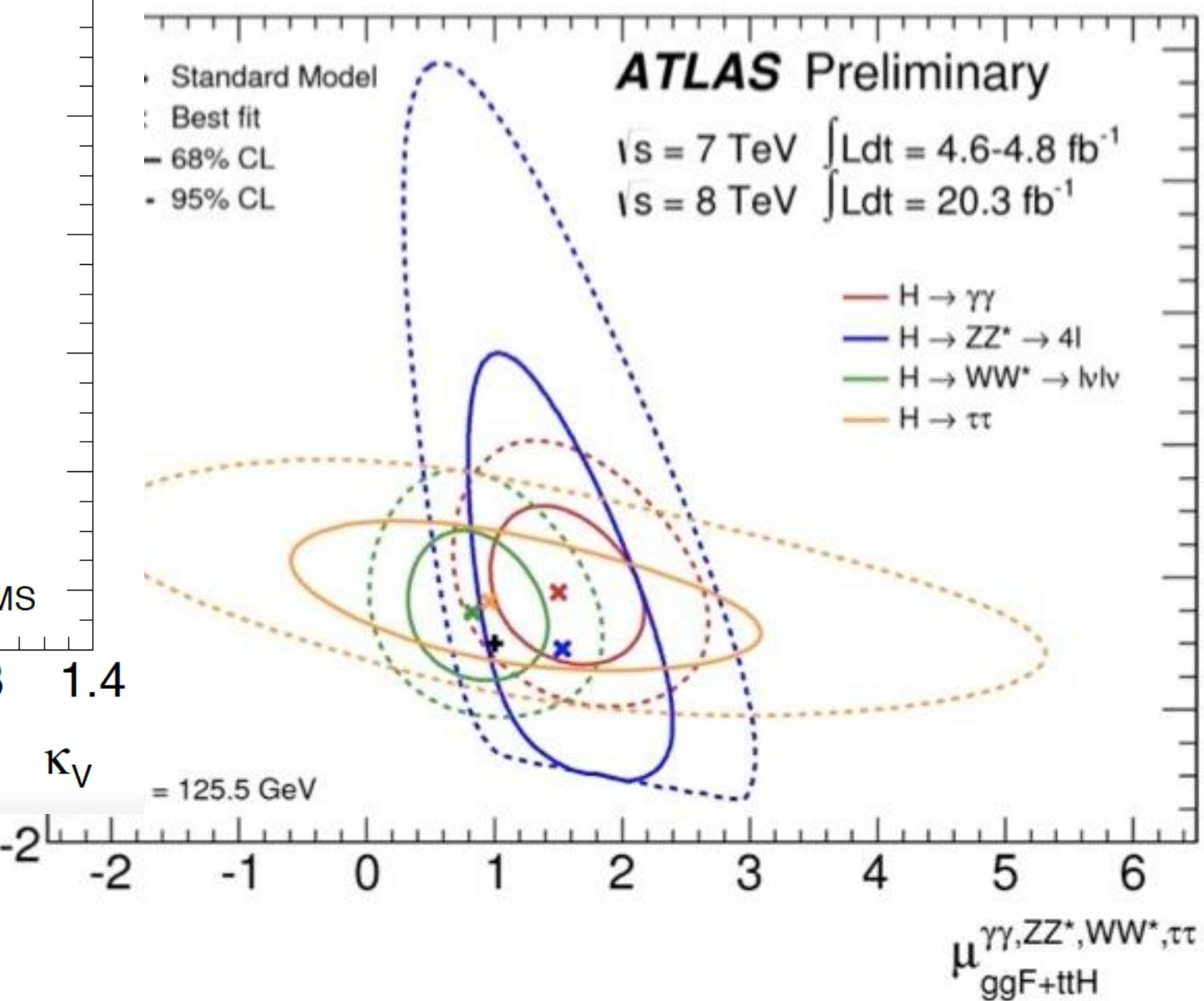
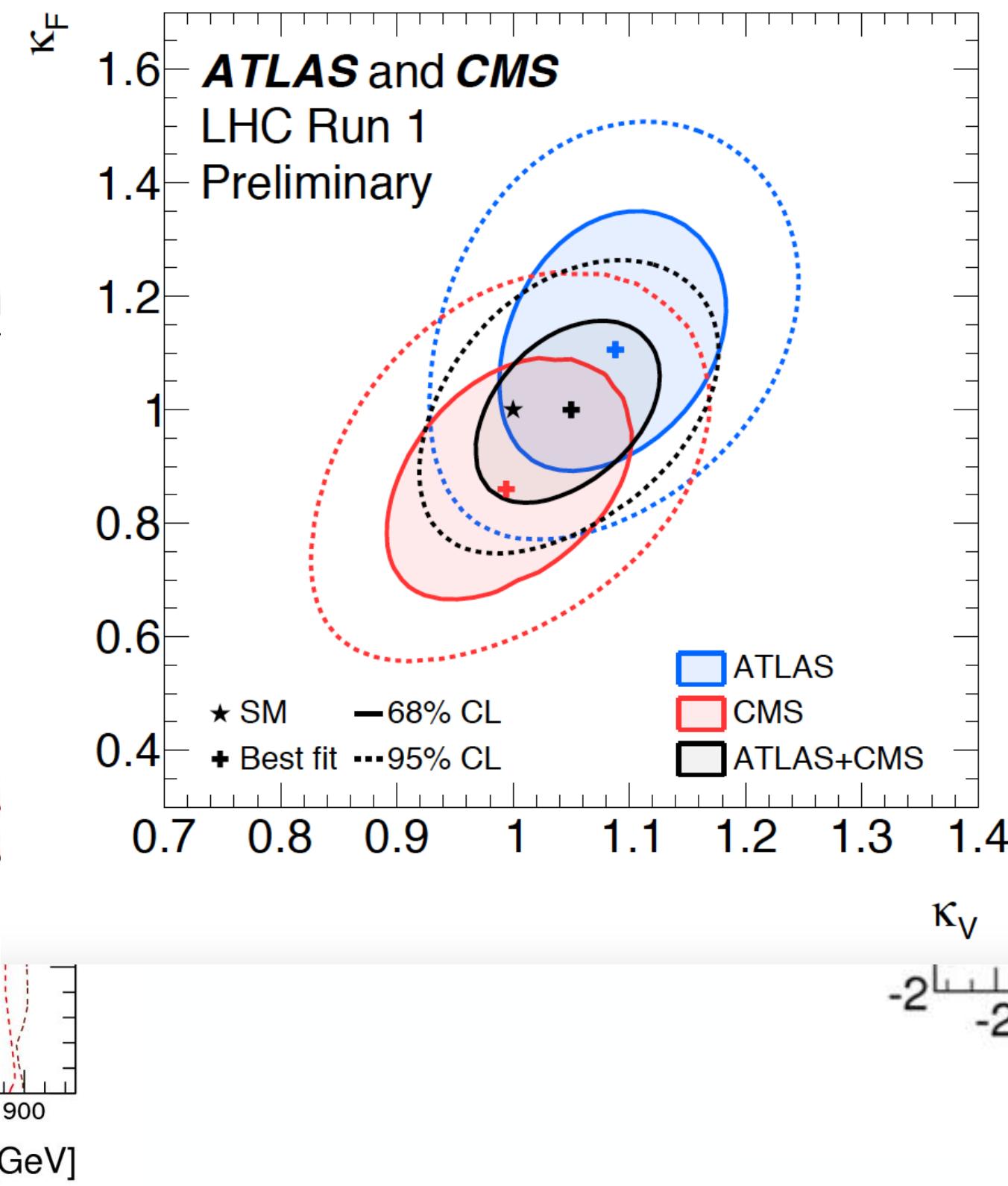
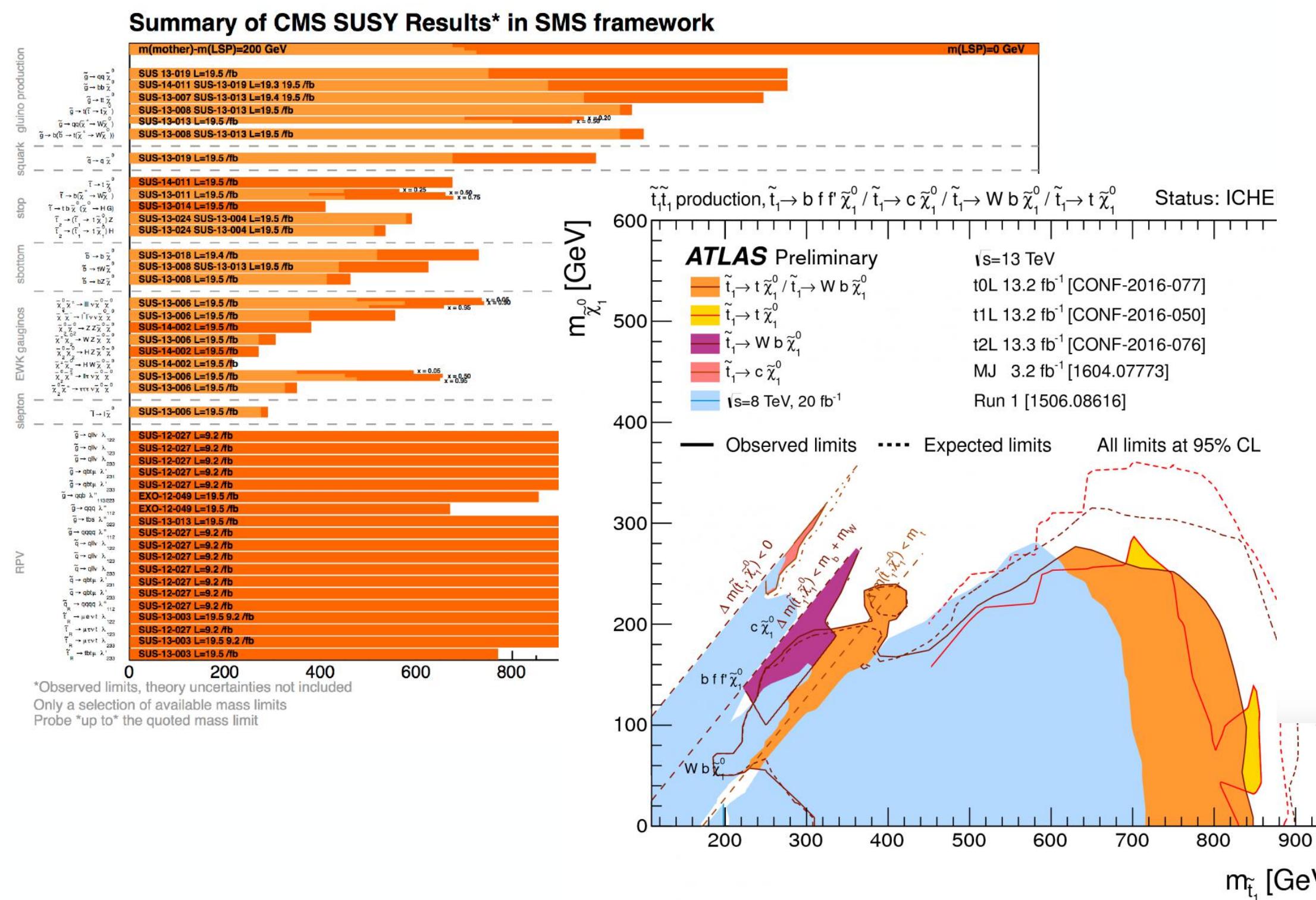
If holds, can do Intervals only POI-space without loss of coverage

- removed NP from quoted results
- saved a lot of compute $O(1) \ll O(1000)$



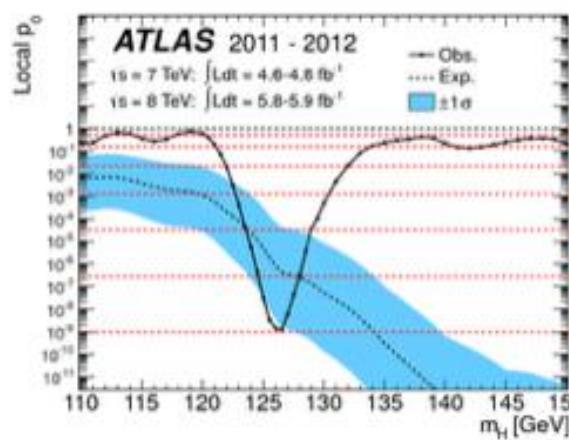
Looking Back and Forward

The methodology of simulation-driven frequentist inference is the result of many decades of development. Supported us for Higgs discovery and beyond

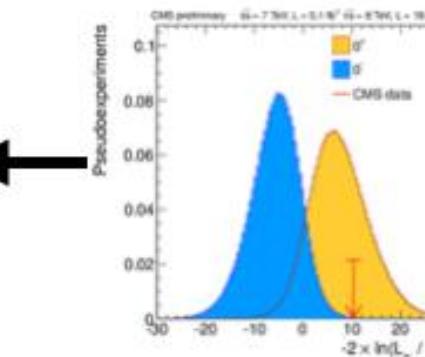


Are we done with Stats methodology?

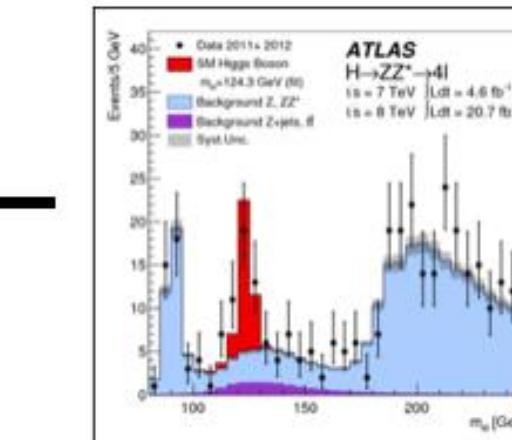
Going beyond



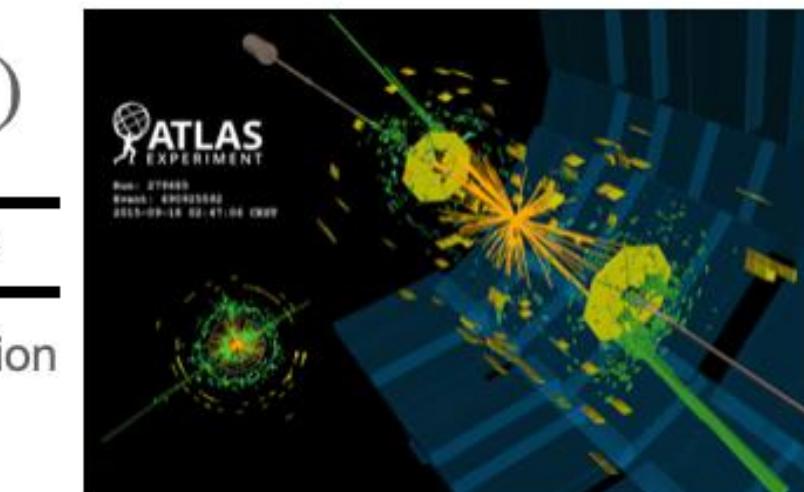
“Likelihood-ful”
Inference



test statistics
for inference



summary statistics /
density estimation



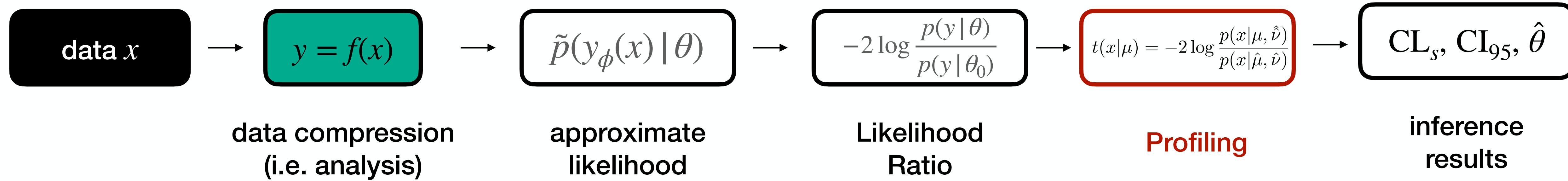
raw data

$s(\cdot)$
data
simulation

Back to basics:

In frequentist analysis we do not need the likelihood

- we choose to run likelihood-free inference with an approximate l'hood



But ultimately all these intermediate step are optional

Recently: A lot of development in exploring ML-driven “direct” likelihood-free inference

What's the connection?

Key Stat. Concept: Likelihood-Ratios

Why? Among all functions of the data ***they are optimal*** when used as a test statistic

(Neyman-Pearson)

Machine Learning:

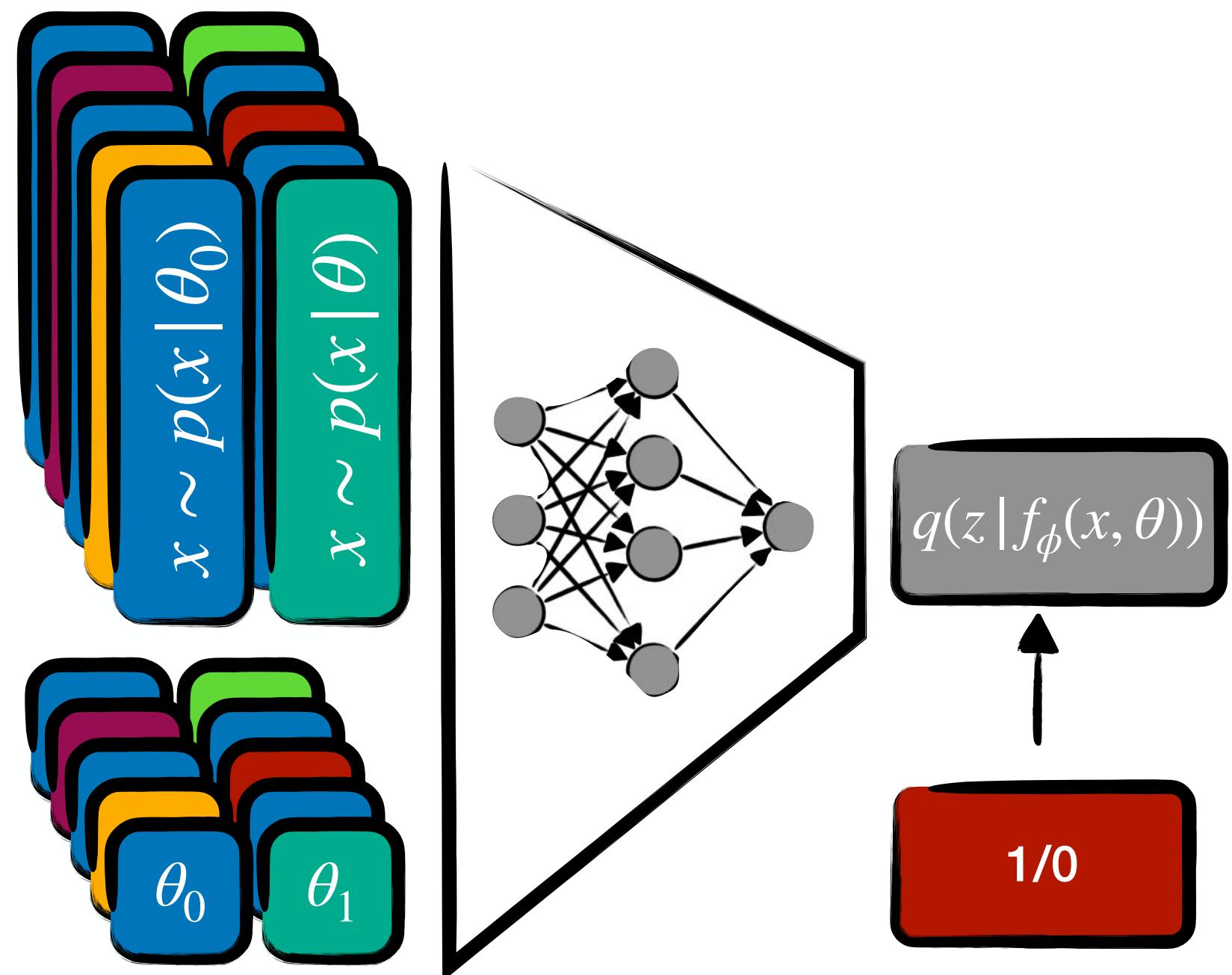
extremely good at searching through functions that are optimal in one sense of another

(based on samples)

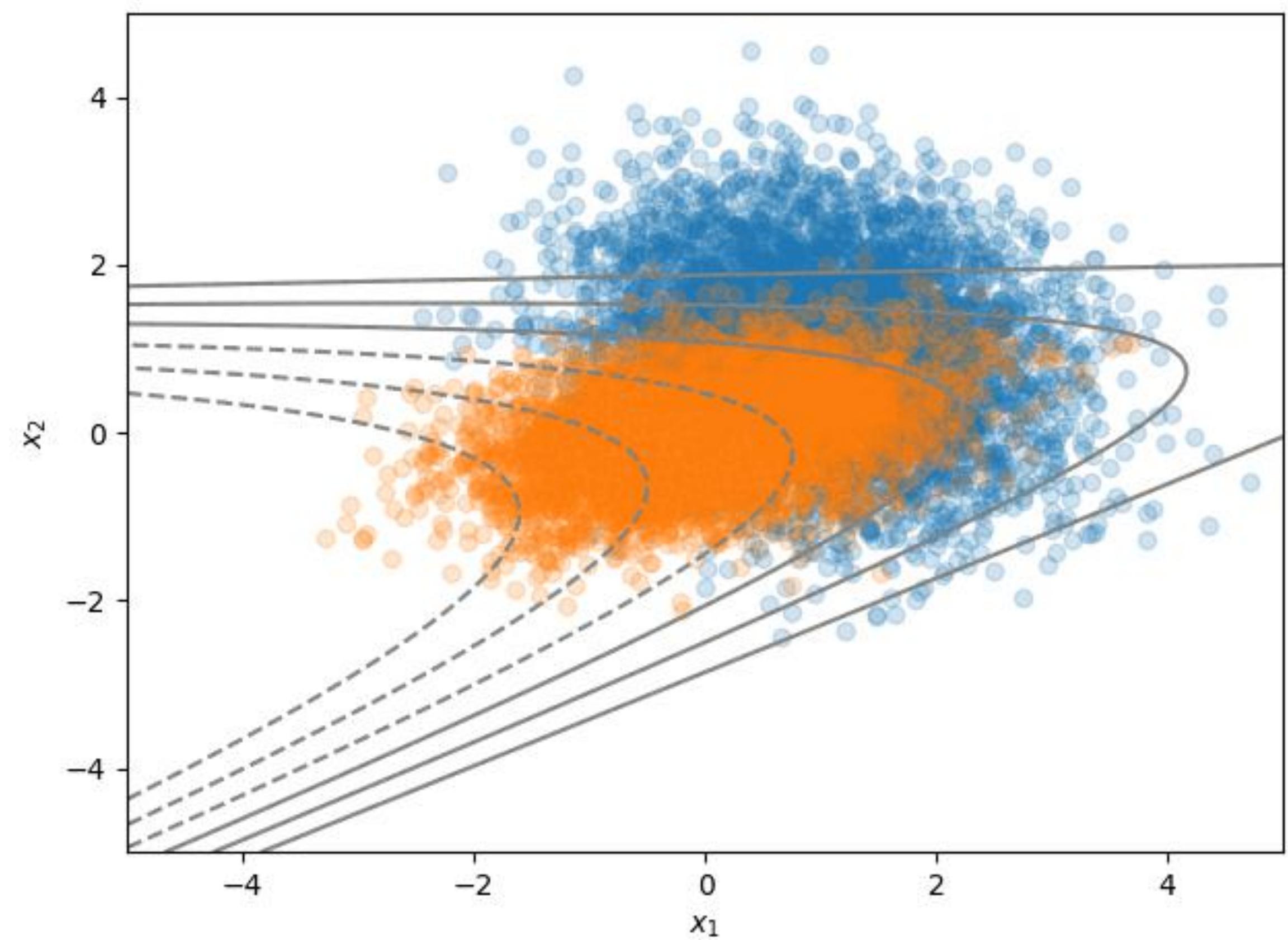
seems like one should be able to use ML to do frequentist statistics

Likelihood Ratio Trick

Training neural networks for classification produces likelihood ratios, purely from samples without ever having to specify $p(x | \theta)$

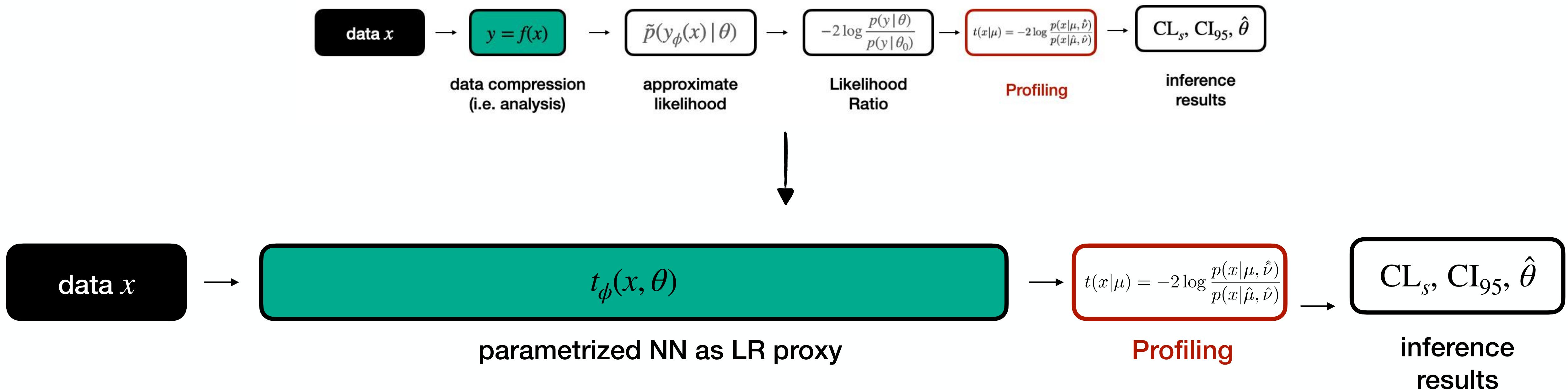


$$t_\phi(x, \theta) \rightarrow \frac{p(x | \theta)}{p(x | \theta_0)}$$



Back to basics:

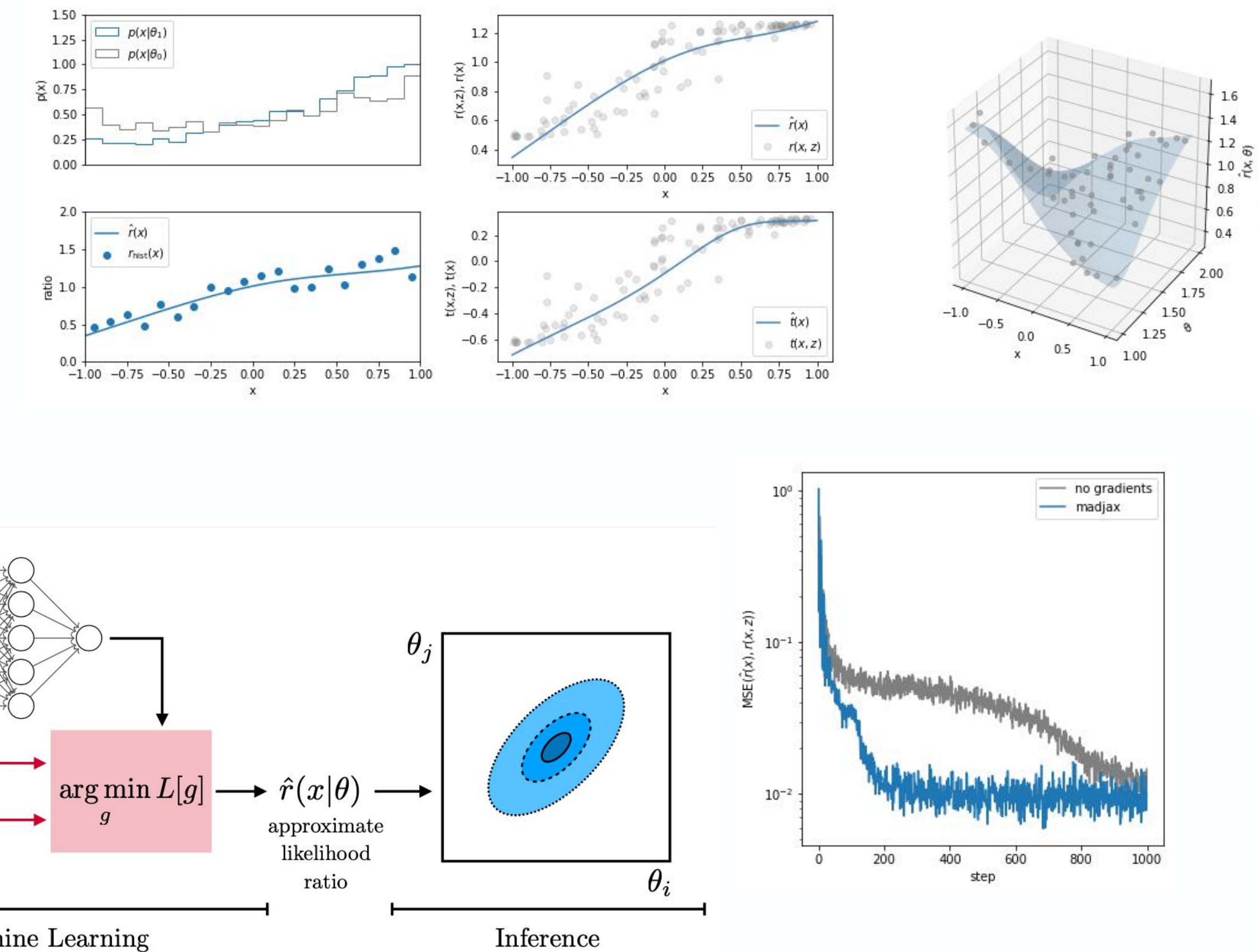
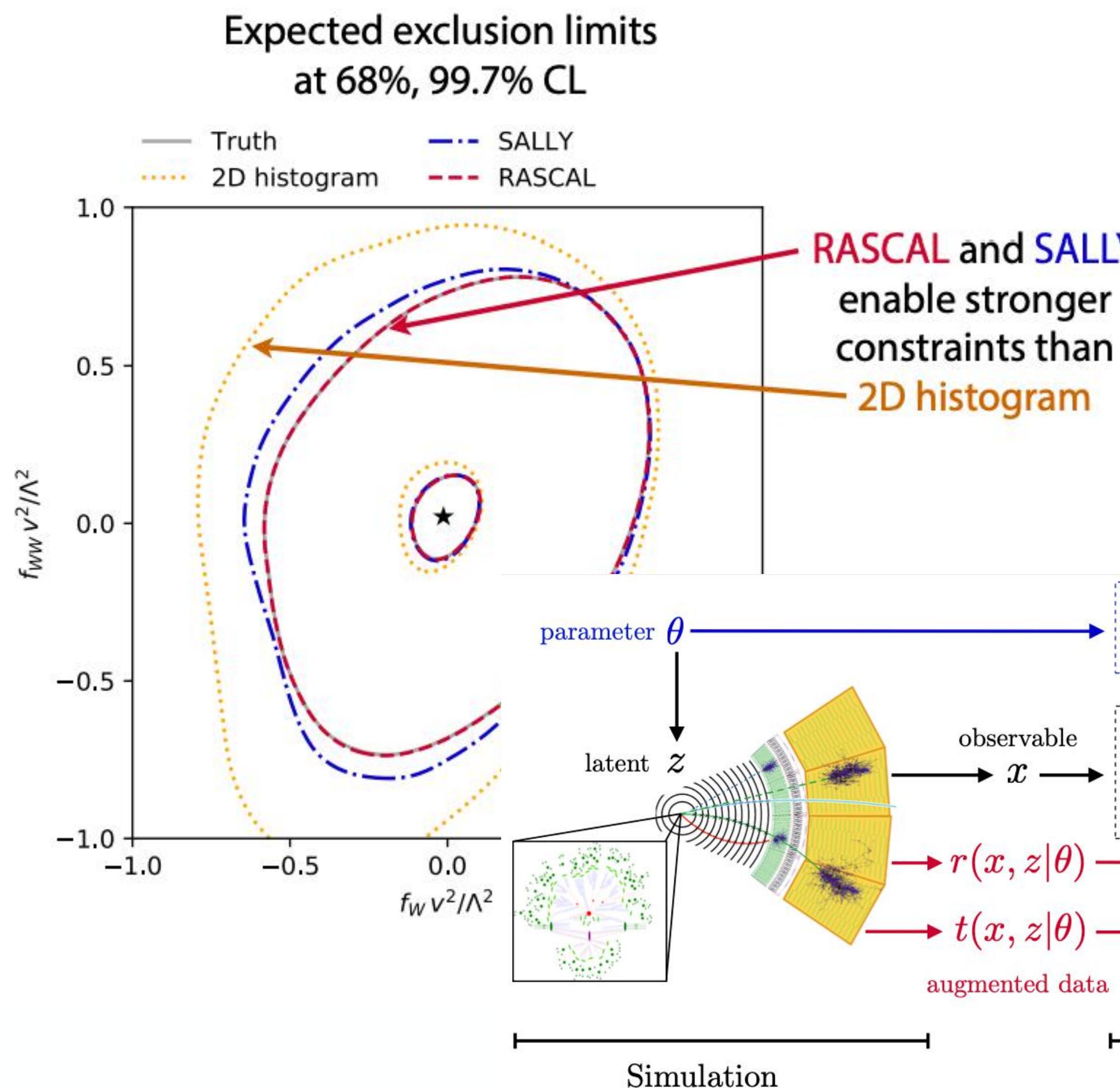
We can start to replace a lot of our pipeline with ML components that directly output likelihood ratios



**Key: this also removes the (lossy) compression we did to reduce dimensionality
→ increased sensitivity**

Back to basics:

Completely new workflow, but with promise: New generation of tooling to train these “likelihood ratio” networks



Going even Further

Now that we've stared? Just how much can we replace with ML?



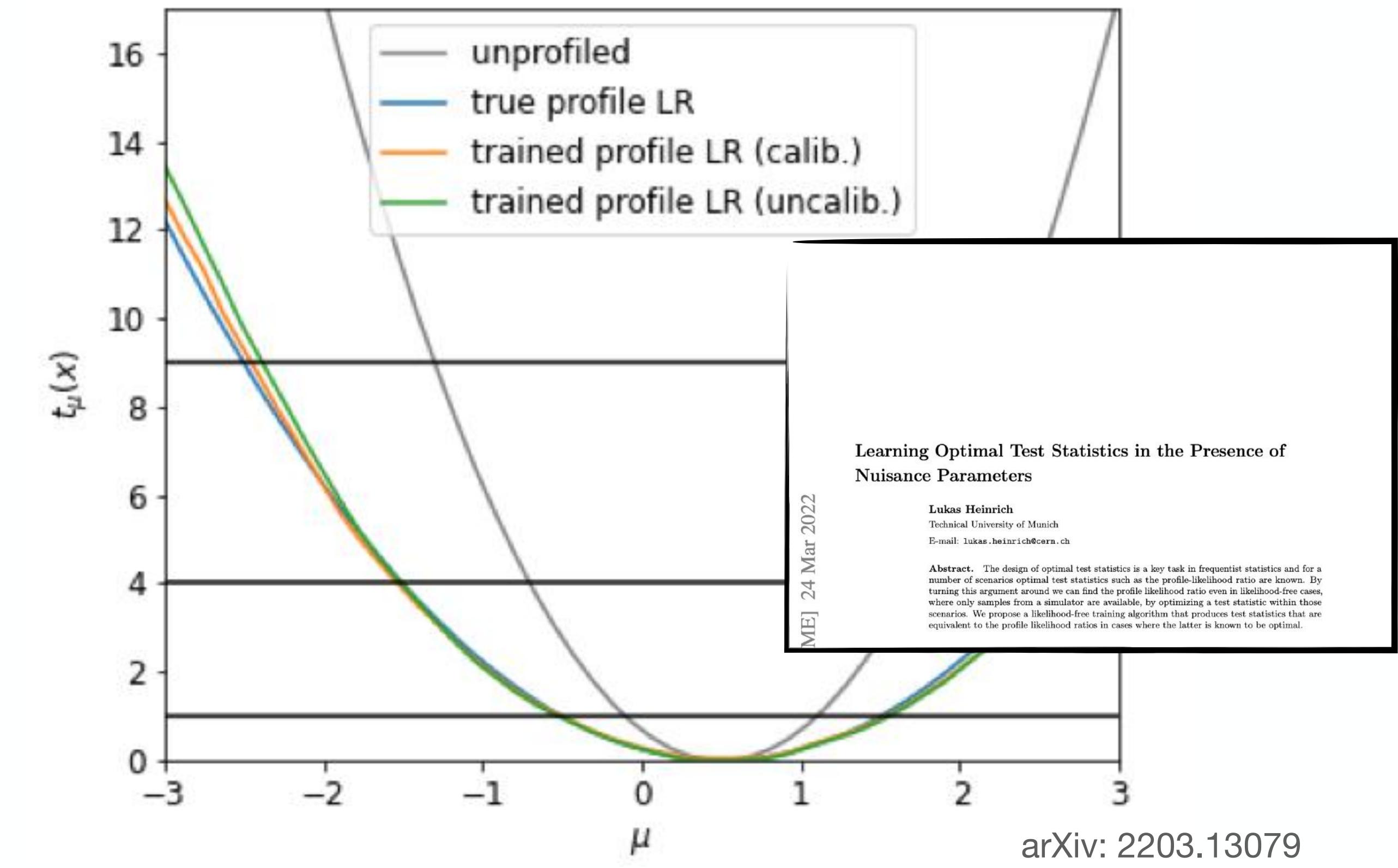
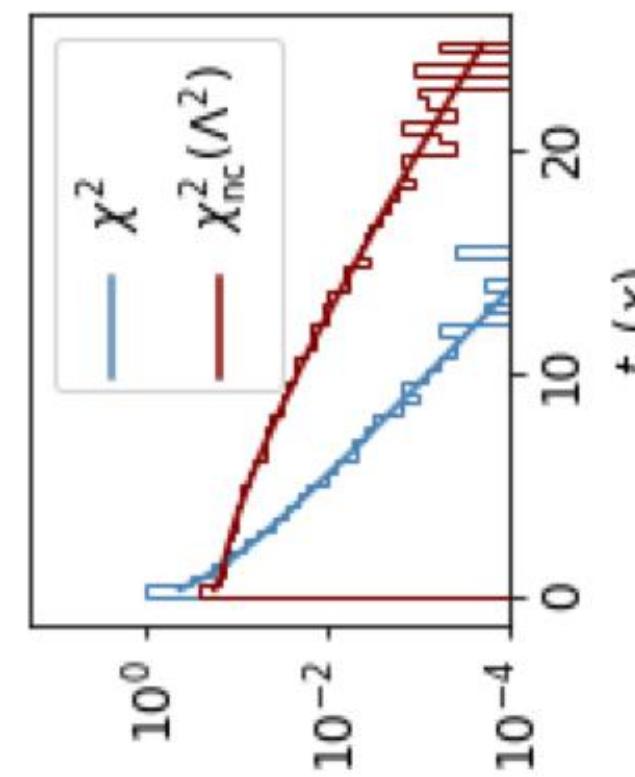
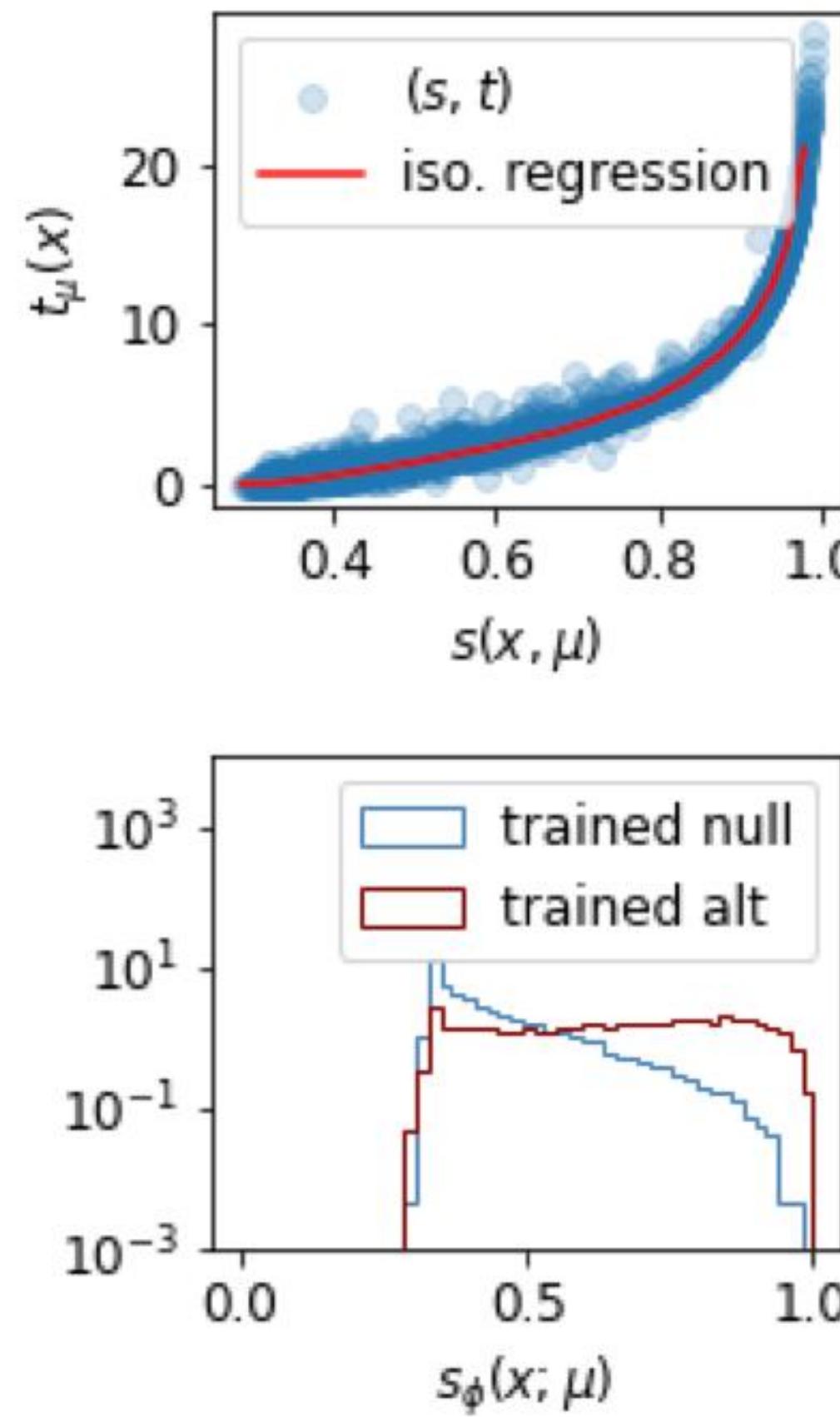
Profile Likelihood is also just used because it is optimal in a certain sense

Can we learn a network that directly outputs profile likelihood?

$$t(x|\mu) = -2 \log \frac{p(x|\mu, \hat{\nu})}{p(x|\hat{\mu}, \hat{\nu})}$$

Going even Further

Yes, if the profile likelihood is optimal, we can find it through optimization



Learning Optimal Test Statistics in the Presence of Nuisance Parameters

Lukas Heinrich
Technical University of Munich
E-mail: lukas.heinrich@cern.ch

Abstract: The design of optimal test statistics is a key task in frequentist statistics and for a number of scenarios optimal test statistics with the profile-likelihood ratio are known. By turning this argument around we can find the profile likelihood ratio even in likelihood-free cases, where only samples from a simulator are available, by optimizing a test statistic within those scenarios. We propose a likelihood-free training algorithm that produces test statistics that are equivalent to the profile likelihood ratios in cases where the latter is known to be optimal.

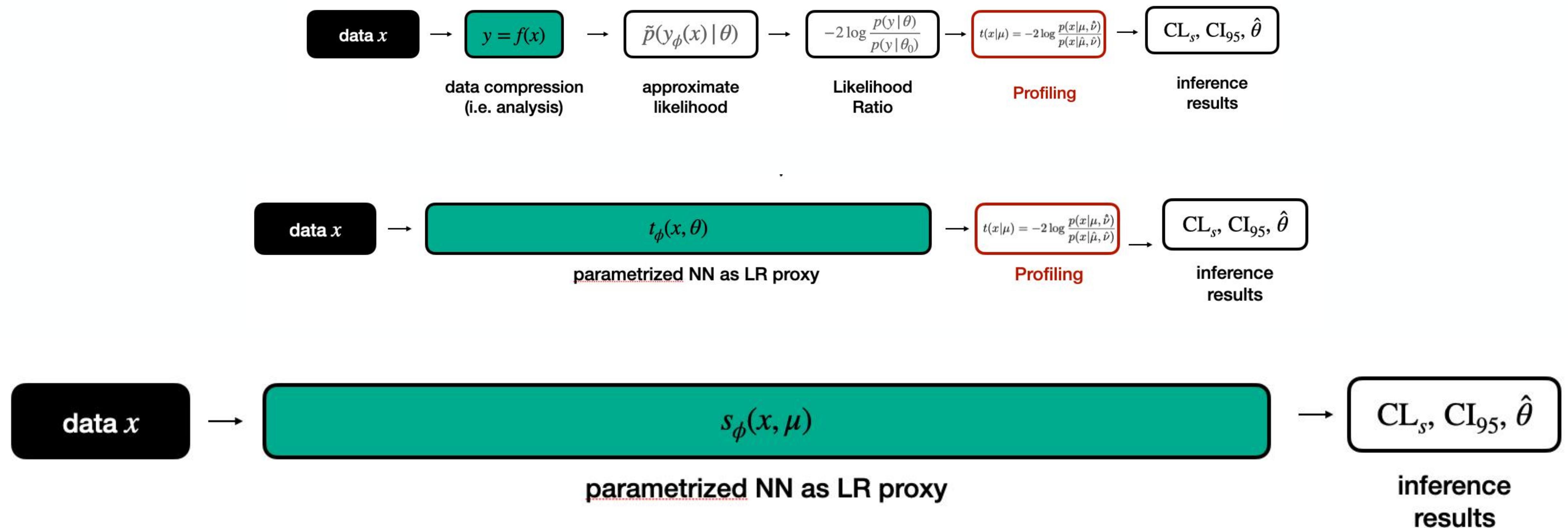
MEJ 24 Mar 2022

arXiv: 2203.13079

Active area of research: happy to chat more offline

Going even Further

It seems within reach, that one could target even more of the classical frequentist statistics pipeline with an end-to-end ML workflow



Summary

LHC Statistics is a poster child of frequentist analysis in a complex setting:

- likelihood-free (and even then expensive sampling)
- high-dimensional data and parameter spaces

Many decades of development culminated in a well-oiled machine with fairly principled approach to inference.

Classic Approach: proxy likelihood and focus on good modelling with incorporation of prior information for hundreds of nuisance parameters

Emerging: Machine Learning driving a resurgence in likelihood-free techniques that aim to eliminate some of the compromises we usually make

Backup

Summary: Many things work but also lots of interesting questions

Exciting development: public real-world probability models

**Not only opportunity to push new science but
also new stats methodology / answer "what if's**

A new public workhorse / platform akin to e.g.
on-off problem, $\text{Pois}(n | \mu s + b)$, ...

Examples:

- Bayesian Workflow on LHC models
- Coverage & Asymptotics Studies
- Signal interpolation strategies
- Reinterpretation Tools
- Model distillation ("simplified likelihoods")



Updates > News > New open release streamlines interactions with theoretical physicists

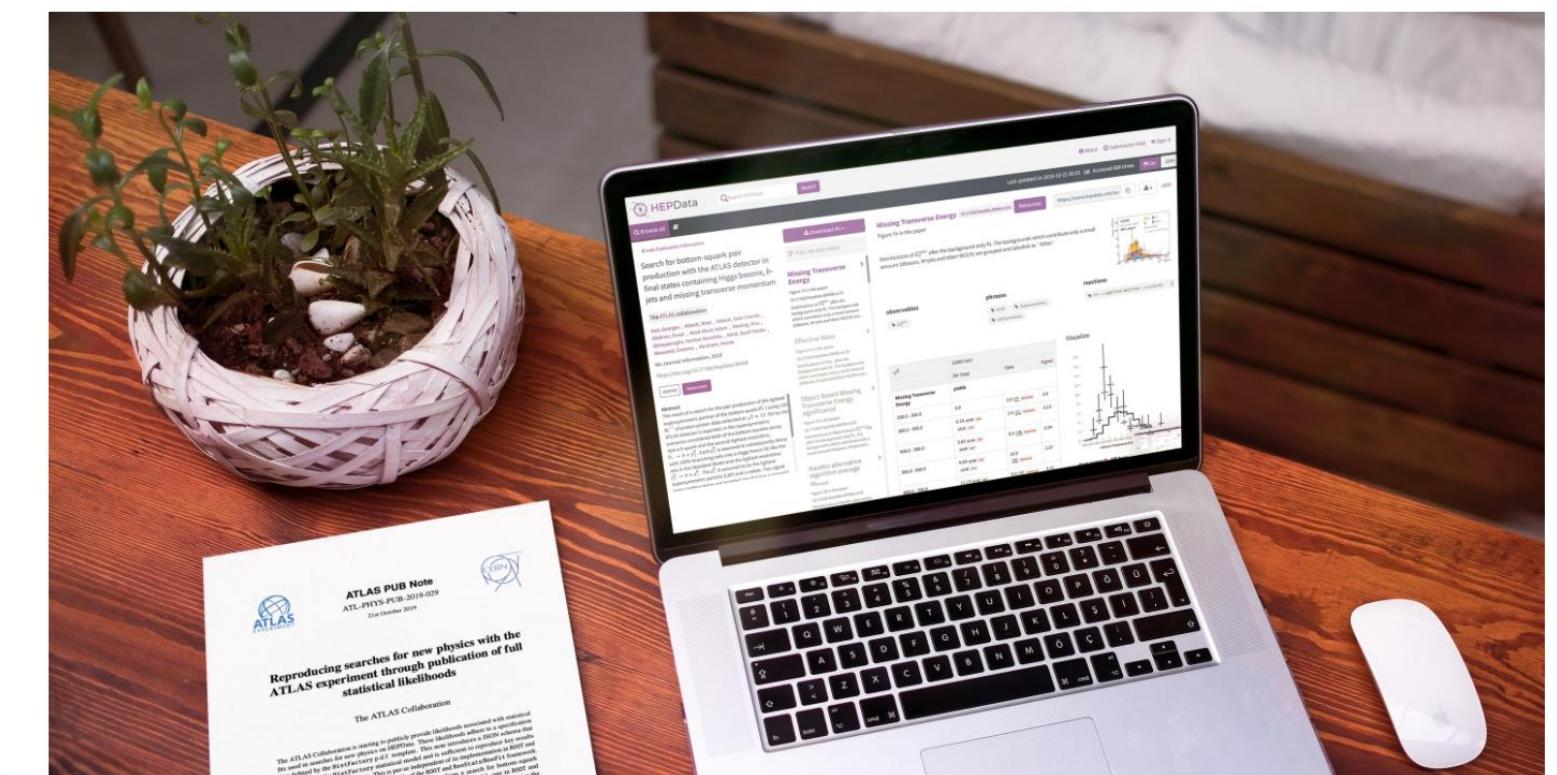
News

Tags:
[open data](#)

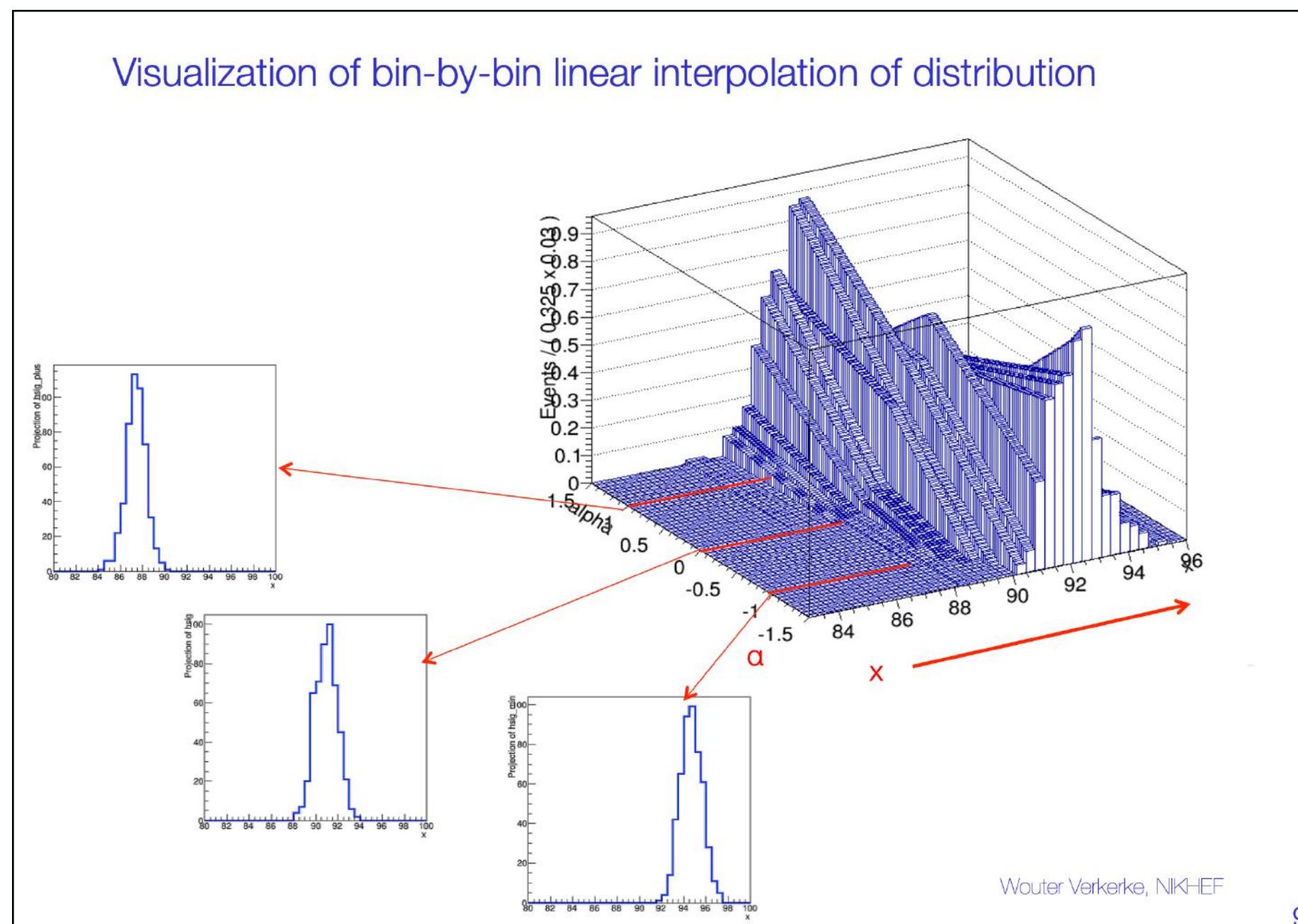
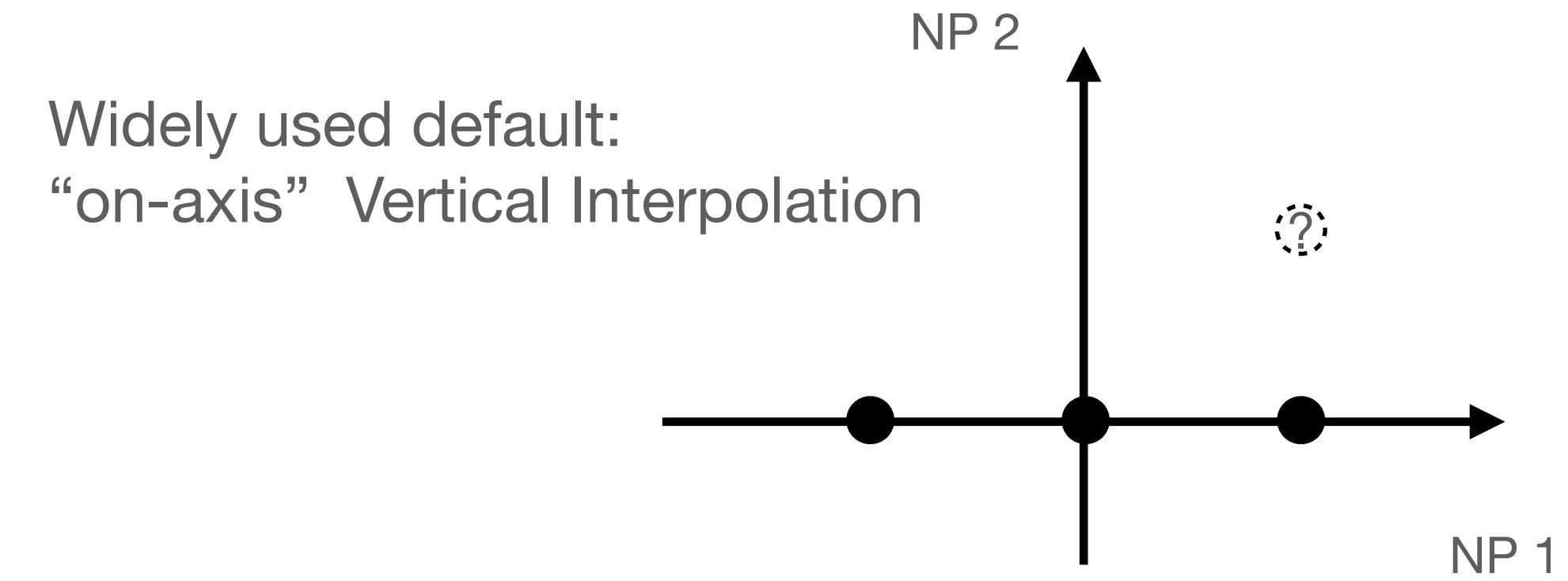
New open release streamlines interactions with
theoretical physicists — **statisticians**

The ATLAS Collaboration has released the first open likelihoods from
an LHC experiment.

12th December 2019 | By [Katarina Anthony](#)



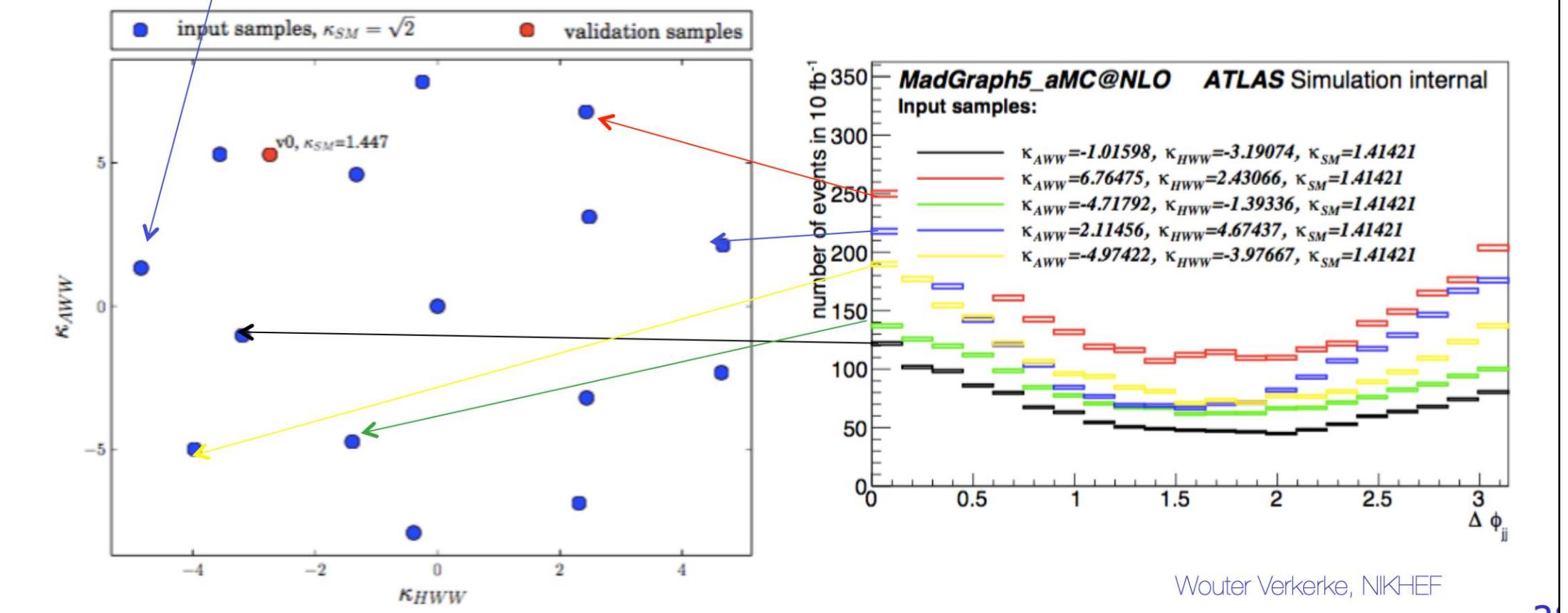
Example Interpolations



More Complex Interpolation Schemes
with “off-axis” input distributions

Truth-level validation study on simulation samples

- Procedure
 - VBF H \rightarrow WW process with SM (g_{SM}) and 2 BSM operators (g_{HWW} , g_{AWW})
50k events generated. Kinematic observable used: $\Delta\phi_{jj}$, **Only signal considered**
 - 15 samples with different parameter settings used to construct EFT morphing model



Two-dimensional interpolation

[W. Verkerke]

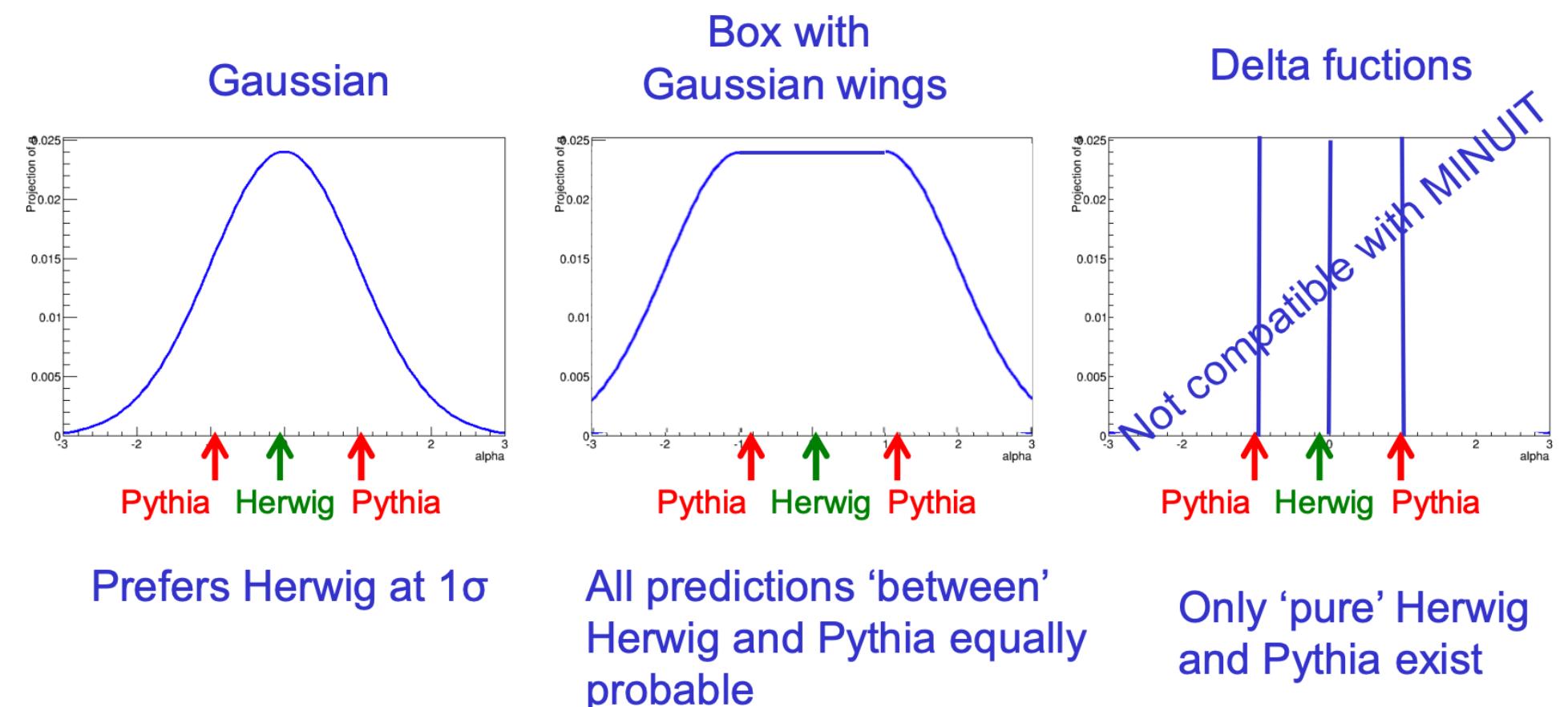
A problem for non-asymptotic cases

Sometimes, NPs relate to (categorial) choice of simulator itself

- lots of discussions on appropriate “simplified subsidiary”

Specific issues with theory uncertainties

- Subsidiary measurement of a theoretical 2-point uncertainty effectively quantifies the ‘knowledge’ on these models
 - Extra difficult to make meaningful statement about this*, since meaning of parameter is not well embedded in underlying theory model
 - But again, all procedures need to assume some distribution... Profiling requires you to spell it out
- Some options and their effects



Wouter Verkerke, NIKHEF

Decision Theory

Key point in inference (intervals, tests, ...): reject / accept.

Here HEP differs from orthodox statistics.

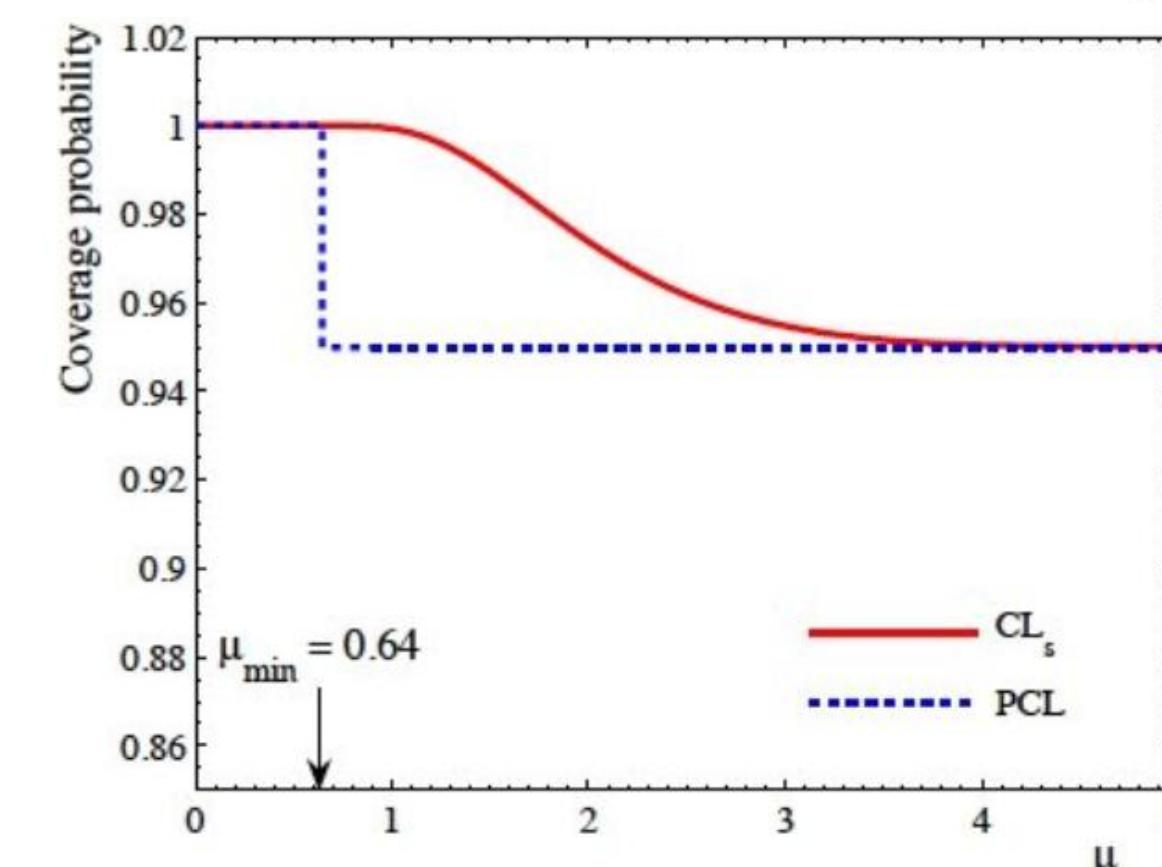
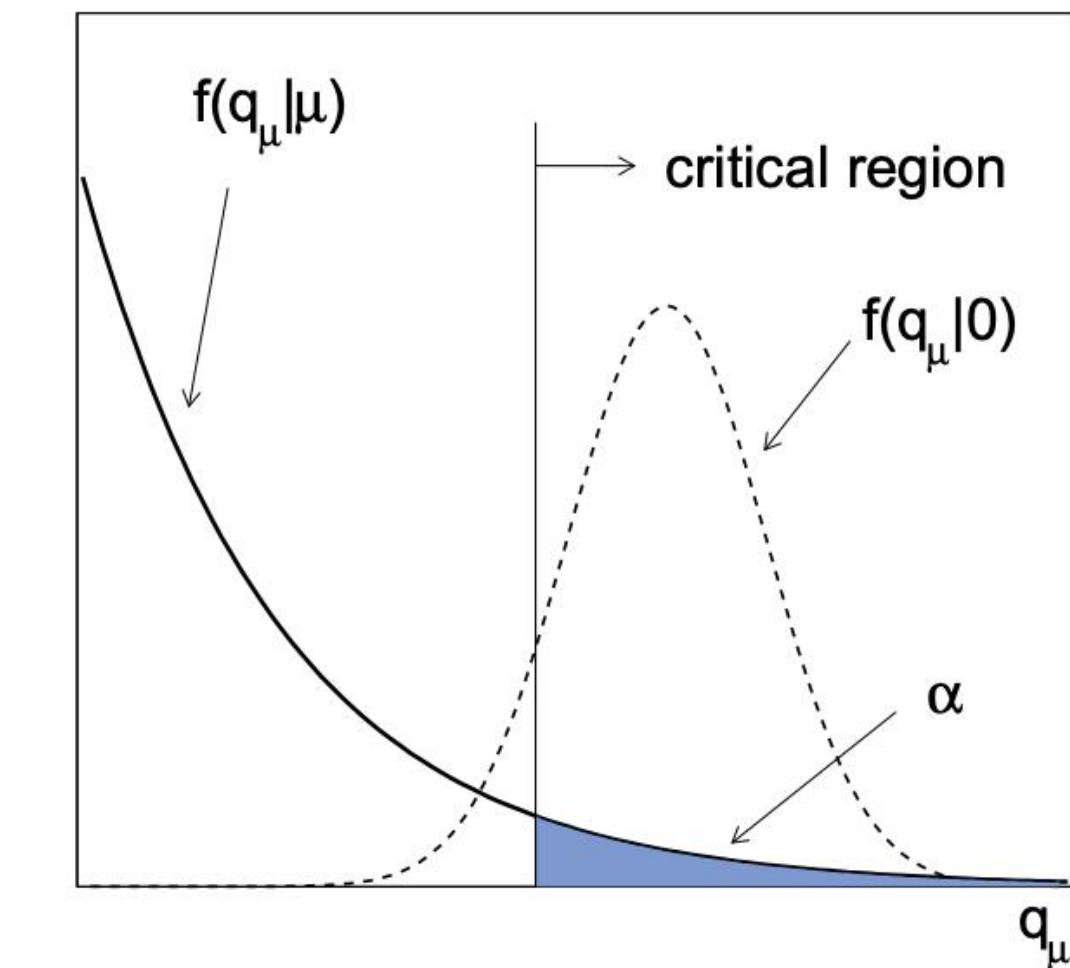
Incorporate power of alternative into decision function of rejecting null → no power: never reject

Reason: huge consequence of rejecting null
(only reject if good alternative available)

HEP Jargon: “CLs” method

$$CL_s = p/\beta$$

↑
power at p-value
↓
p-value



Boundary Effects

Common occurrence in HEP:

- null is on boundary of parameter space: → adjust Wilks'
- extended region in alternative that's indistinguishable from null
→ input from statisticians interesting

