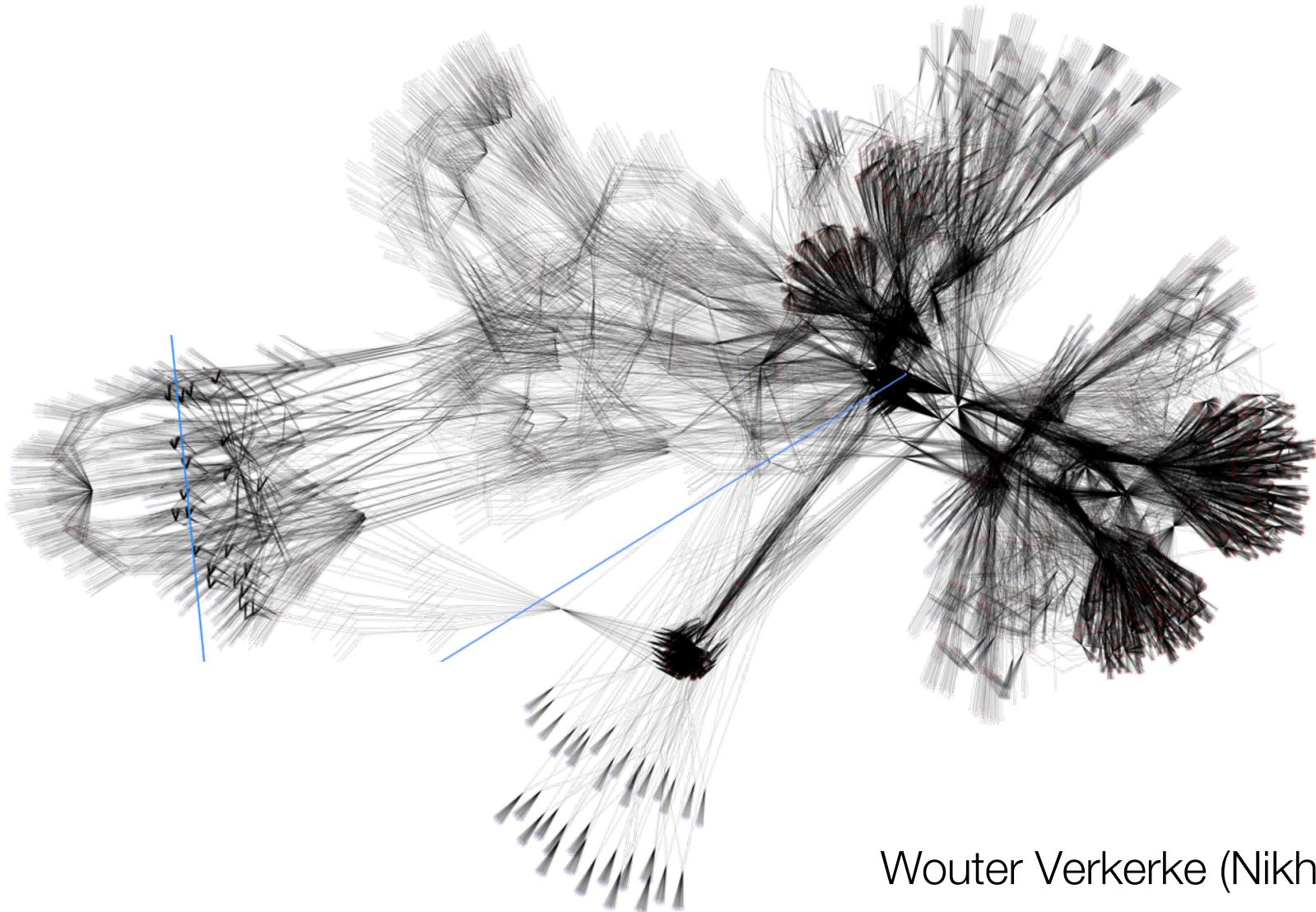


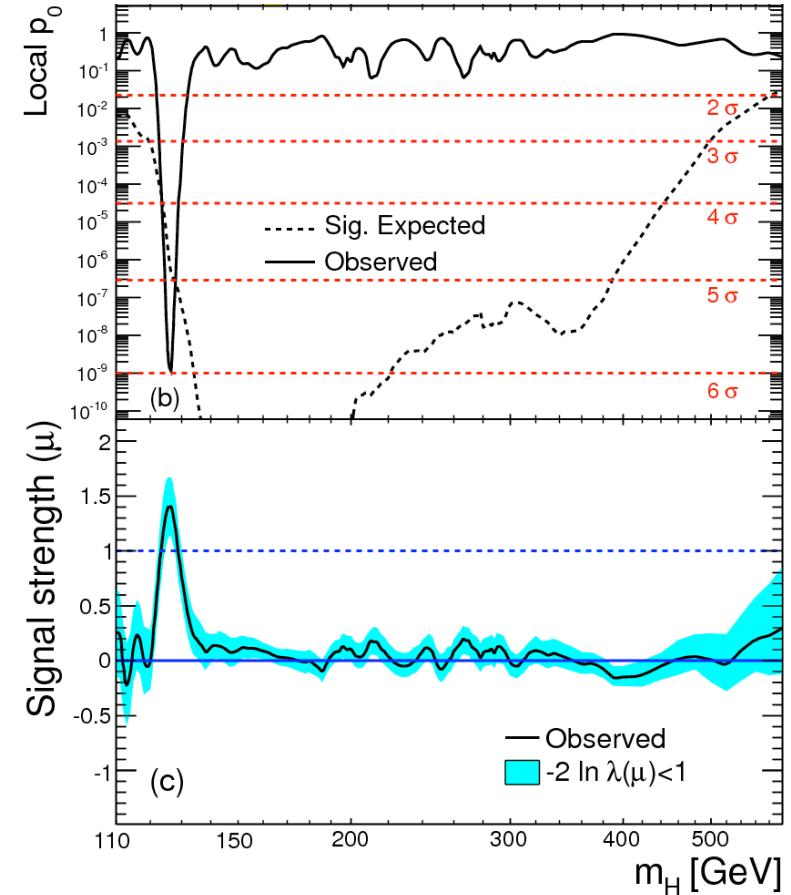
# Statistical analysis with RooFit/RooStats/HistFactory, a brief overview



# What do you want to know?

- Physics questions we have...
  - Does the (SM) Higgs boson exist?
  - What is its production cross-section?
  - What is its boson mass?
- Statistical tests construct probabilistic statements:  $p(\text{theo}|\text{data})$ , or  $p(\text{data}|\text{theo})$ 
  - Hypothesis testing (discovery)
  - (Confidence) intervals
  - Measurements & uncertainties
- Result: *Decision based on tests*

“As a layman I would now say: I think we have it”

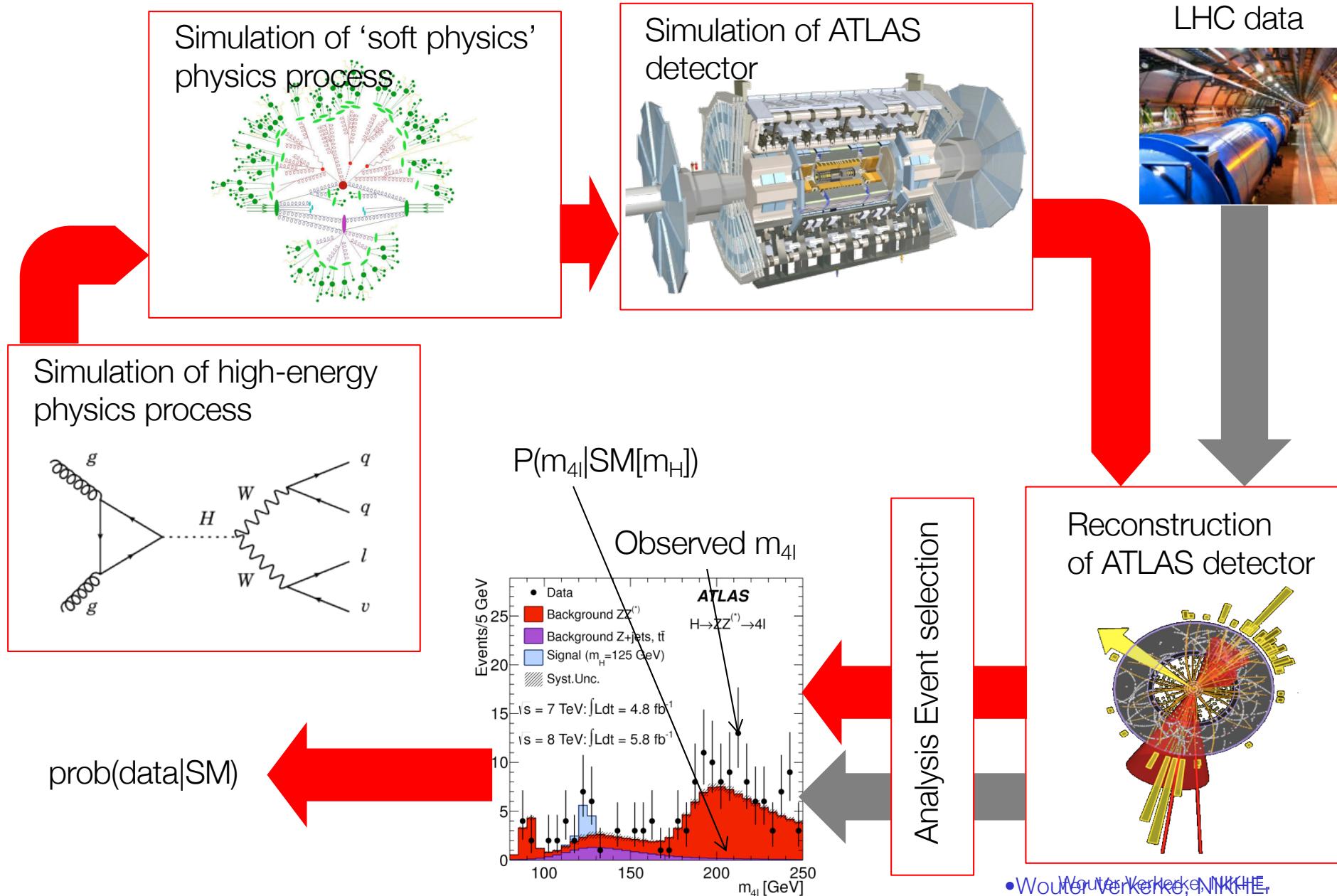


Wouter Verkerke, NIKHEF

## All experimental results start with the formulation of a model

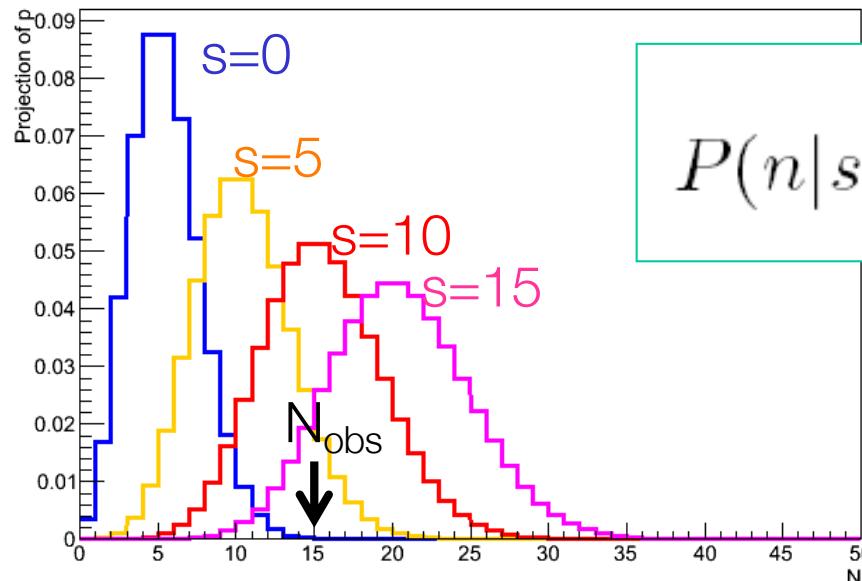
- Examples of HEP **physics** models being tested
  - SM with  $m(\text{top})=172,173,174 \text{ GeV}$  → Measurement top quark mass
  - SM with/without Higgs boson → Discovery of Higgs boson
  - SM with composite fermions/Higgs → Measurement of Higgs coupling properties
- Via chain of physics simulation, showering MC, detector simulation and analysis software, a physics model is reduced to a **statistical** model

# The HEP analysis workflow illustrated



# All experimental results start with the formulation of a model

- Examples of HEP **physics** models being tested
  - SM with  $m(\text{top})=172, 173, 174 \text{ GeV} \rightarrow$  Measurement top quark mass
  - SM with/without Higgs boson  $\rightarrow$  Discovery of Higgs boson
  - SM with composite fermions/Higgs  $\rightarrow$  Measurement of Higgs coupling properties
- Via chain of physics simulation, showering MC, detector simulation and analysis software, a physics model is reduced to a statistical model
- A **statistical** model defines  $p(\text{data}|\text{theory})$  for all observable outcomes
  - Example of a statistical model for a counting measurement with a known background



$$P(n|s+b) = \frac{(s+b)^n}{n!} e^{-(s+b)}$$

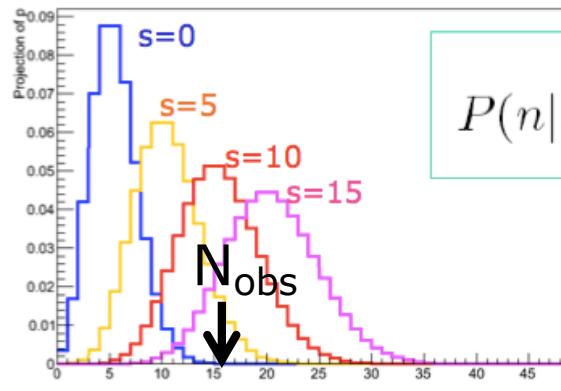
*NB:  $b$  is a constant in this example*

**Definition: the Likelihood  
is  $P(\text{observed data}|\text{theory})$**

Wouter Verkerke, NIKHEF

# Everything starts with the likelihood

- **All** fundamental statistical procedures are based on the likelihood function as ‘description of the measurement’



$$P(n|s+b) = \frac{(s+b)^n}{n!} e^{-(s+b)}$$

NB: b is a constant in this example

Definition: the Likelihood  
is  $P(\text{observed data}|\text{theory})$

e.g.  $L(15|s=0)$   
e.g.  $L(15|s=10)$

Frequentist statistics

Bayesian statistics

Maximum Likelihood

Confidence interval on s

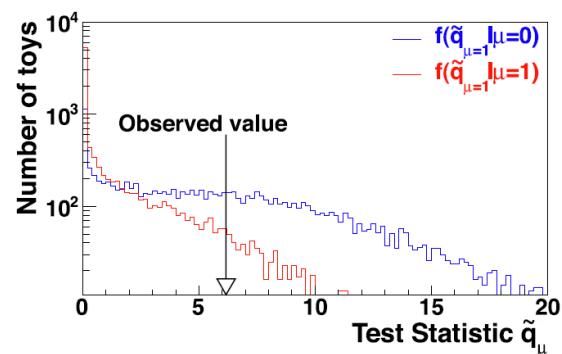
Posterior on s

$s = x \pm y$

# Everything starts with the likelihood

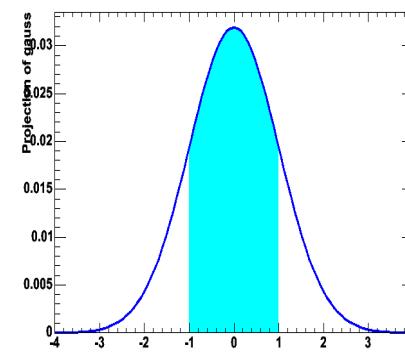
Frequentist statistics

$$\lambda_\mu(\vec{N}_{obs}) = \frac{L(\vec{N} | \mu)}{L(\vec{N} | \hat{\mu})}$$



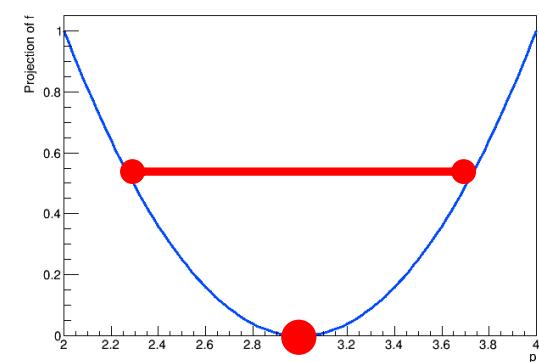
Bayesian statistics

$$P(\mu) \propto L(x | \mu) \cdot \pi(\mu)$$



Maximum Likelihood

$$\left. \frac{d \ln L(\vec{p})}{d \vec{p}} \right|_{p_i = \hat{p}_i} = 0$$



Confidence interval  
or p-value

Posterior on s  
or Bayes factor

$s = x \pm y$

## Complications start when your model uncertain

---

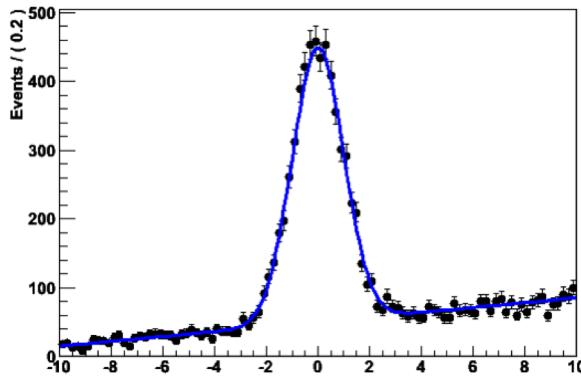
- Canonical problem: what happens to your inference on  $s$  when  $b$  is not exactly known, but has some uncertainty

$$P(n|s+b) = \frac{(s+b)^n}{n!} e^{-(s+b)}$$

- Uncertain models are the norm in HEP (predictions from the simulation chain pick up a multitude of uncertain parameters from each step)
- Not a problem for modeling per se (e.g. for Poisson formula is the same, with  $b$  a floating parameter now) but extra challenge for inference
  - Must propagate uncertainty on  $b$  in inference in  $s$
  - Procedure depends on chosen inference technique, even then multiple choices possible. (Common choice in ATLAS is the profile likelihood ratio)
  - Terminology:  $s$  is **parameter of interest** (POI),  $b$  is **nuisance parameter**
  - Will not dig deeply into this now, but procedure can be quite involved depending on circumstances (special tools exist → RooStats)

# How is Higgs discovery different from a simple fit?

*Gaussian + polynomial*

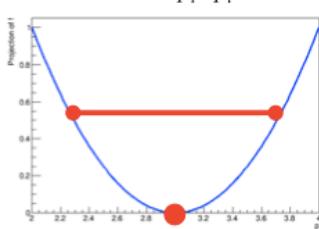


ROOT TH1

ROOT TF1

$$L(\vec{N} \mid \mu, \vec{\theta}) = \prod_i \text{Poisson}(N_i \mid f(x_i, \mu, \vec{\theta}))$$

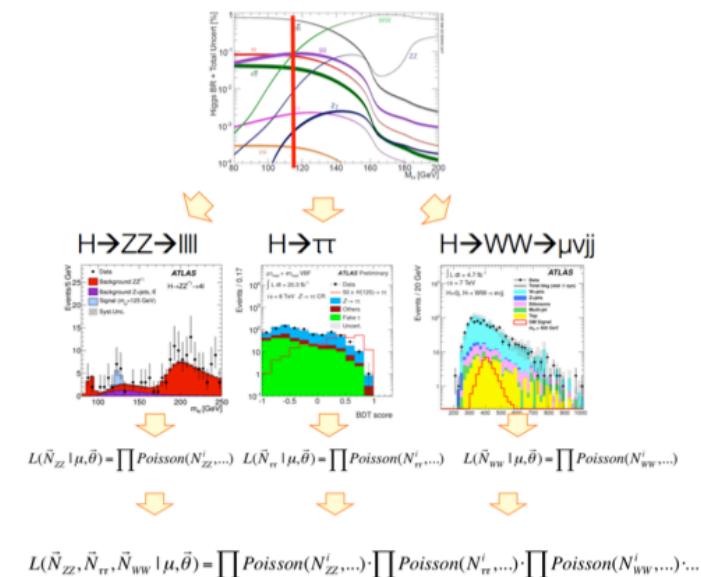
*"inside ROOT"*



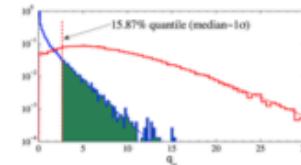
ML estimation of parameters  $\mu, \theta$  using MINUIT  
(MIGRAD, HESSE, MINOS)

$$\mu = 5.3 \pm 1.7$$

*Higgs combination model*



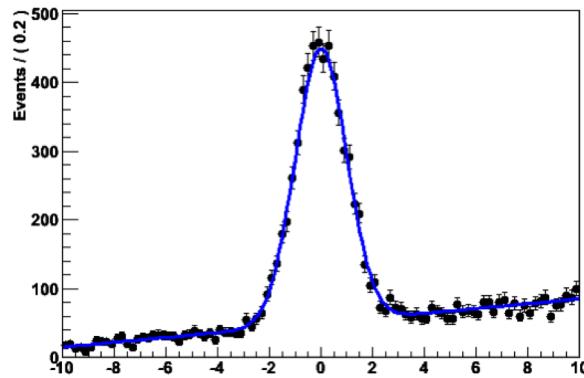
$$\lambda_\mu(\vec{N}_{ZZ}, \vec{N}_{WW}, \vec{N}_{\pi\pi} \mid \mu, \hat{\vec{\theta}}) = \frac{L(\vec{N}_{ZZ}, \vec{N}_{WW}, \vec{N}_{\pi\pi} \mid \mu, \hat{\vec{\theta}})}{L(\vec{N}_{ZZ}, \vec{N}_{WW}, \vec{N}_{\pi\pi} \mid \hat{\mu}, \hat{\vec{\theta}})}$$



$$p(H_\mu) = \int_{\lambda_{\text{obs}}}^{\infty} f(\lambda \mid H_\mu) d\lambda = \dots$$

# How is Higgs discovery different from a simple fit?

Gaussian + polynomial

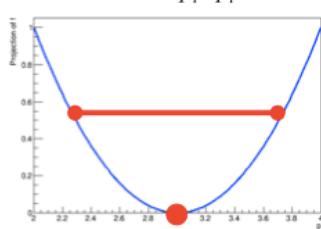


ROOT TH1

ROOT TF1

$$L(\vec{N} \mid \mu, \vec{\theta}) = \prod_i \text{Poisson}(N_i \mid f(x_i, \mu, \vec{\theta}))$$

"inside ROOT"



ML estimation of parameters  $\mu, \theta$  using MINUIT  
(MIGRAD, HESSE, MINOS)

$$\mu = 5.3 \pm 1.7$$

Likelihood Model

orders of magnitude more complicated. Describes

- O(100) signal distributions
- O(100) control sample distr.
- O(1000) parameters representing syst. uncertainties

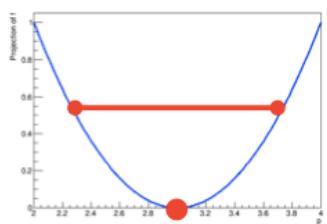
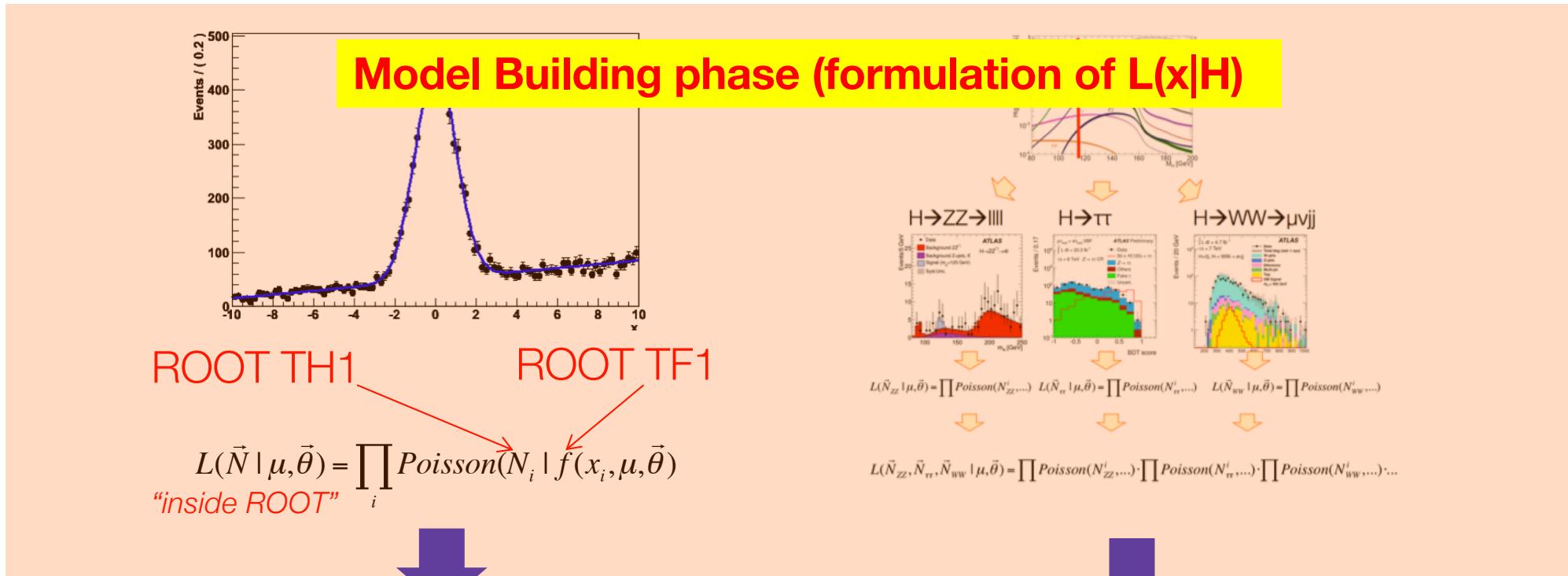
$$L(\vec{N}_{ZZ}, \vec{N}_{\tau\tau}, \vec{N}_{WW} \mid \mu, \vec{\theta}) = \prod \text{Poisson}(N_{ZZ}^i, \dots) \cdot \prod \text{Poisson}(N_{\tau\tau}^i, \dots) \cdot \prod \text{Poisson}(N_{WW}^i, \dots) \cdot \dots$$

Frequentist confidence interval construction and/or p-value calculation not available as 'ready-to-run' algorithm in ROOT

# How is Higgs discovery different from a simple fit?

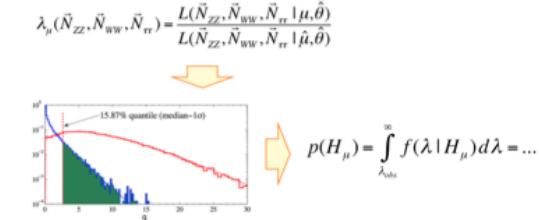
Gaussian + polynomial

Higgs combination model



ML estimation of parameters  $\mu, \theta$  using MINUIT  
(MIGRAD, HESSE, MINOS)

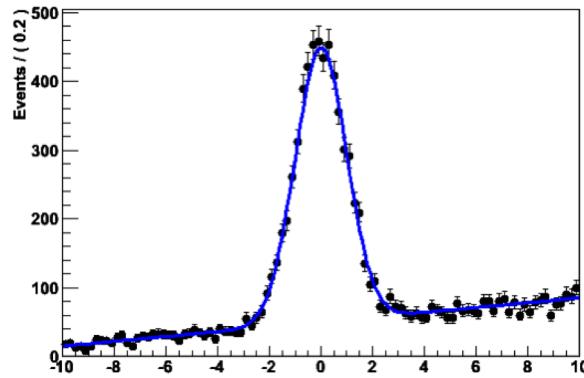
$$\mu = 5.3 \pm 1.7$$



Wouter Verkerke, NIKHEF

# How is Higgs discovery different from a simple fit?

*Gaussian + polynomial*

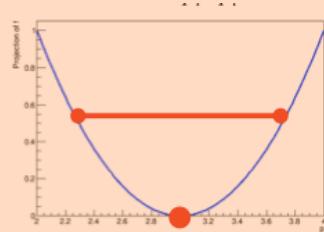


ROOT TH1

ROOT TF1

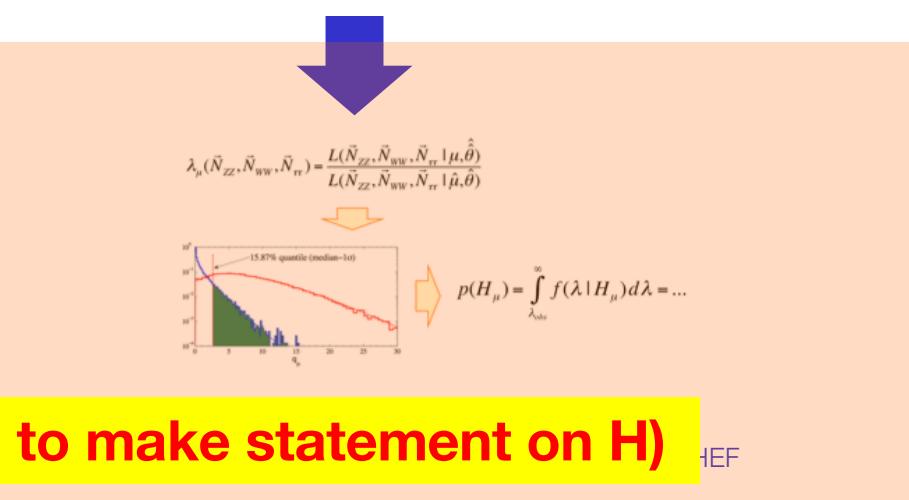
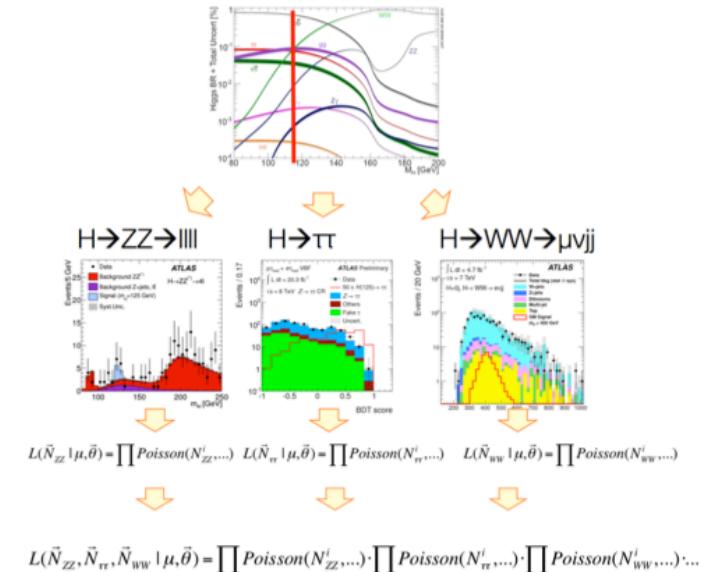
$$L(\vec{N} \mid \mu, \vec{\theta}) = \prod_i \text{Poisson}(N_i \mid f(x_i, \mu, \vec{\theta}))$$

*"inside ROOT"*



ML estimation of parameters  $\mu, \theta$  using MINUIT  
(MIGRAD, HESSE, MINOS)

*Higgs combination model*



**Model Usage phase (use  $L(x|H)$  to make statement on  $H$ )**

HEF

# How is Higgs discovery different from a simple fit?

Gaussian + polynomial

Higgs combination model

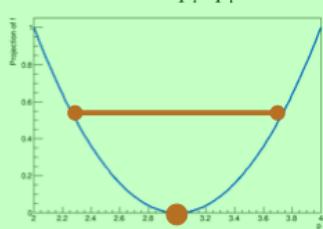
Design goal:

Separate **building of Likelihood model** as much as possible  
from statistical analysis **using the Likelihood model**

- More modular software design
- ‘Plug-and-play with statistical techniques
- Factorizes work in collaborative effort

Re

"in

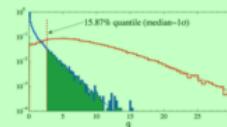


ML estimation of  
parameters  $\mu, \theta$  using MINUIT  
(MIGRAD, HESSE, MINOS)



$$\mu = 5.3 \pm 1.7$$

$$\lambda_\mu(\vec{N}_{ZZ}, \vec{N}_{WW}, \vec{N}_{\tau\tau}) = \frac{L(\vec{N}_{ZZ}, \vec{N}_{WW}, \vec{N}_{\tau\tau} | \mu, \hat{\theta})}{L(\vec{N}_{ZZ}, \vec{N}_{WW}, \vec{N}_{\tau\tau} | \hat{\mu}, \hat{\theta})}$$



$$p(H_\mu) = \int_{\lambda_{obs}}^{\infty} f(\lambda | H_\mu) d\lambda = \dots$$

## The idea behind the design of RooFit/RooStats/HistFactory

---

- Modularity, Generality and flexibility
- Step 1 – Construct the likelihood function  $L(x|p)$

## RooFit, or RooFit+HistFactory

- Step 2 – Statistical tests on parameter of interest  $p$

Procedure can be Bayesian, Frequentist, or Hybrid),  
but always based on  $L(x|p)$

## RooStats

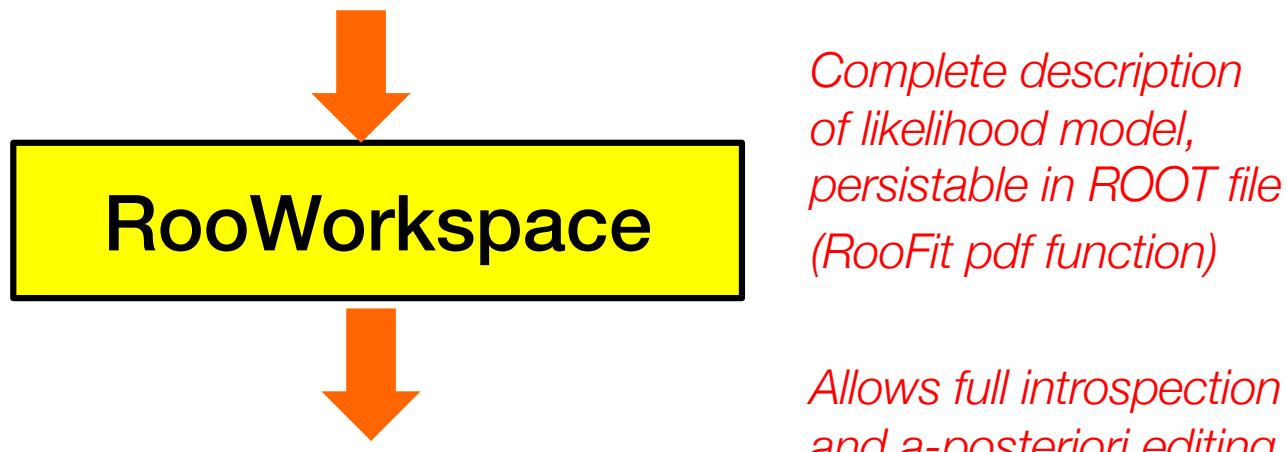
- Steps 1 and 2 are conceptually separated,  
and in Roo\* suit also implemented separately.

## The idea behind the design of RooFit/RooStats/HistFactory

---

- Steps 1 and 2 can be ‘physically’ separated (in time, or user)
- **Step 1** – Construct the likelihood function  $L(x|p)$

### RooFit, or RooFit+HistFactory



- **Step 2** – Statistical tests on parameter of interest  $p$

### RooStats

## The benefits of modularity

---

- Perform different statistical test on exactly the same model

**RooFit, or RooFit+HistFactory**



**RooWorkspace**



**“Simple fit”**

(ML Fit with  
HESSE or  
MINOS)



**RooStats**

**(Frequentist  
with toys)**



**RooStats**

**(Frequentist  
asymptotic)**



**RooStats**

**Bayesian  
MCMC**

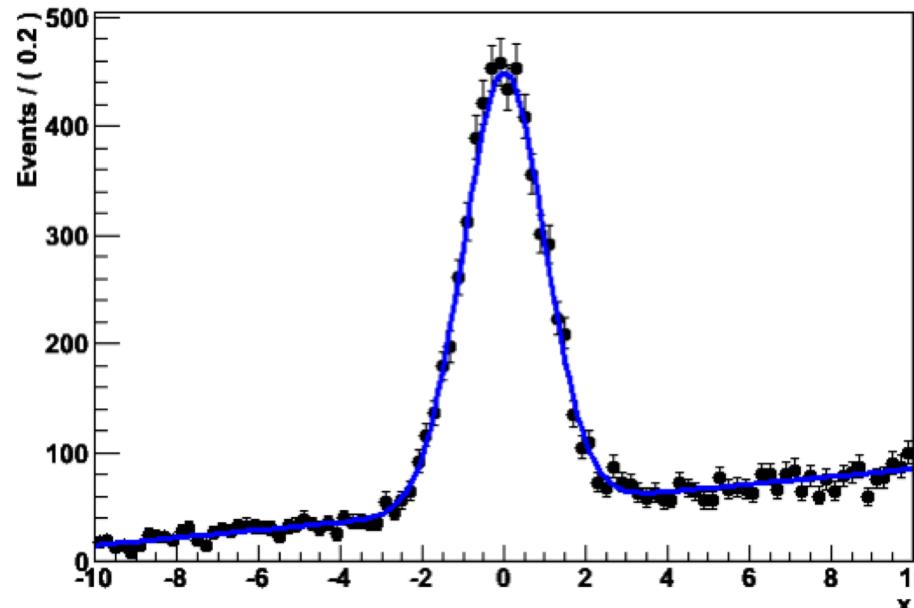
# RooFit

WV + D. Kirkby - 1999

## RooFit – Focus: coding a probability density function

---

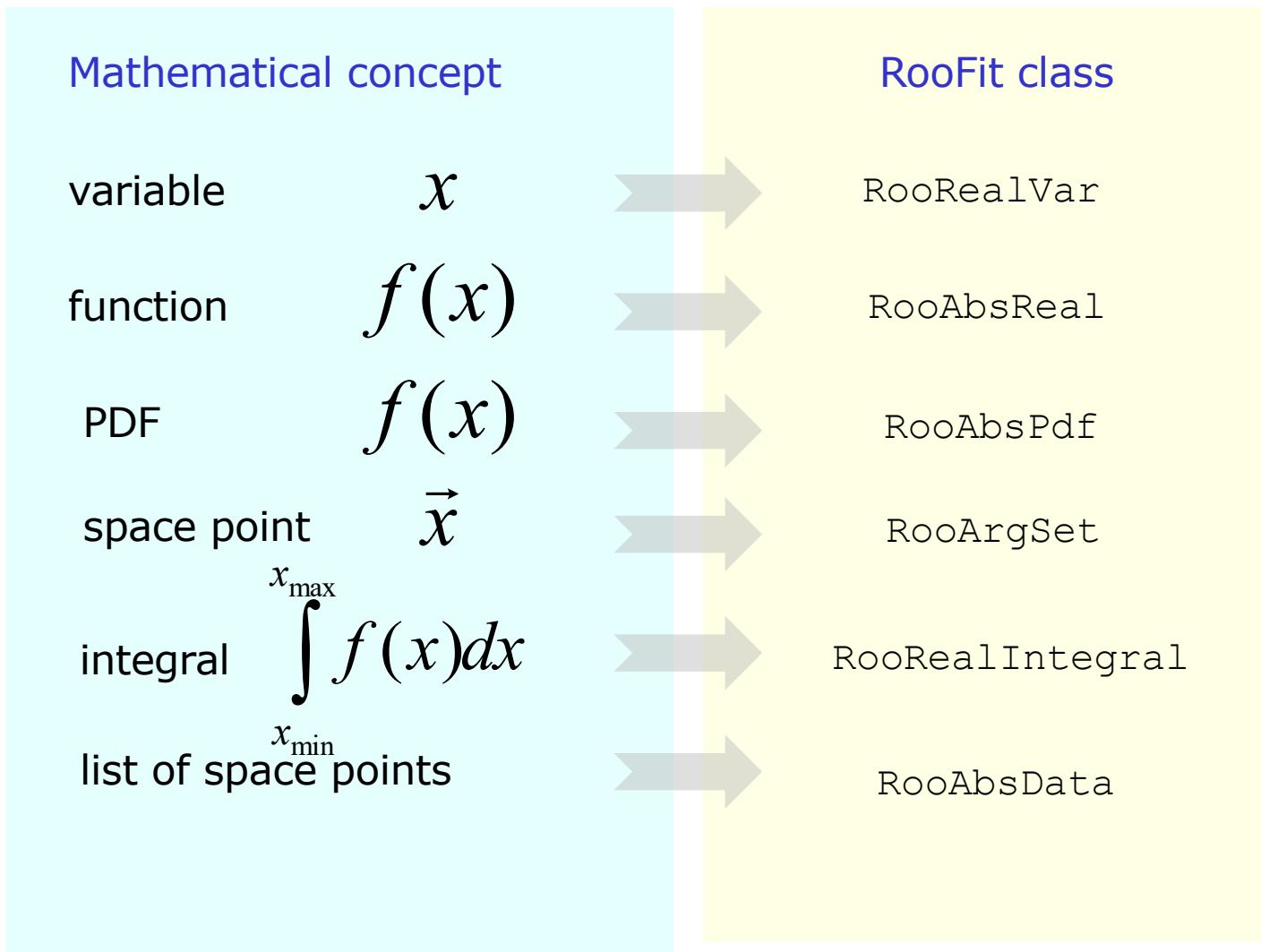
- Focus on one practical aspect of many data analysis in HEP:  
**How do you formulate your p.d.f. in ROOT**
  - For ‘simple’ problems (gauss, polynomial) this is easy
  - But if you want to do unbinned ML fits, use non-trivial functions, or work with multidimensional functions you quickly find that you need some tools to help you



- The RooFit project started in 1999 for data modeling needs for BaBar collaboration initially, publicly available in ROOT since 2003

# RooFit core design philosophy

- Mathematical objects are represented as C++ objects



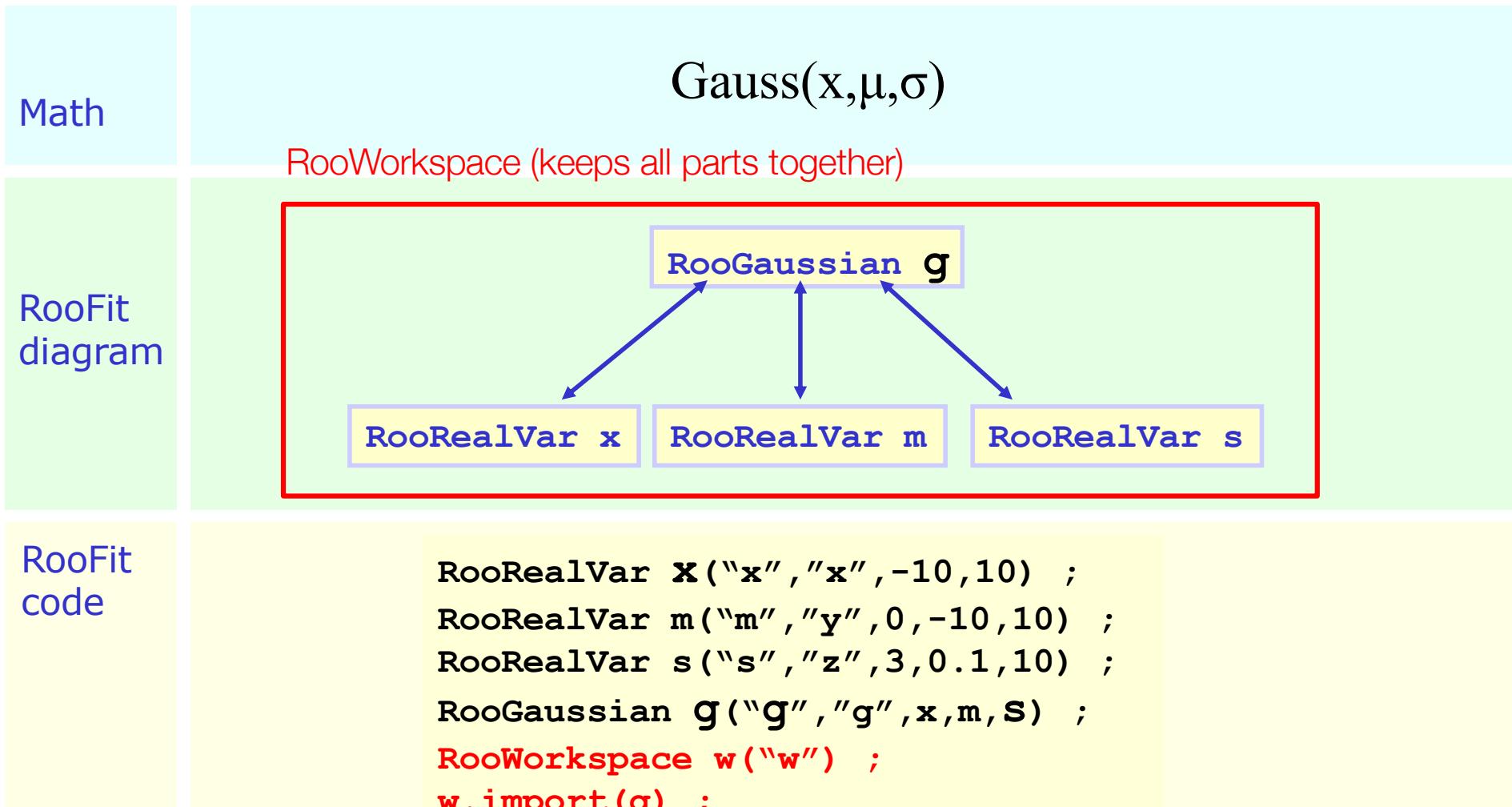
## Data modeling – Constructing composite objects

- Straightforward correlation between mathematical representation of formula and RooFit code

Math	$gauss(x, m, \sqrt{s})$
RooFit diagram	<pre>graph TD; x[RooRealVar x] -- 1 --&gt; g[RooGaussian g]; m[RooRealVar m] -- 2 --&gt; g; s[RooRealVar s] -- 3 --&gt; sqrt[sqrt[sqrts]]; sqrt -- 4 --&gt; g; sqrts[RooFormulaVar sqrts] -- 5 --&gt; g;</pre>
RooFit code	<pre>① RooRealVar x("x","x",-10,10) ; ② RooRealVar m("m","mean",0) ; ③ RooRealVar s("s","sigma",2,0,10) ; ④ RooFormulaVar sqrts("sqrt","sqrt(s)",s) ; ⑤ RooGaussian g("g","gauss",x,m,sqrts) ;</pre>

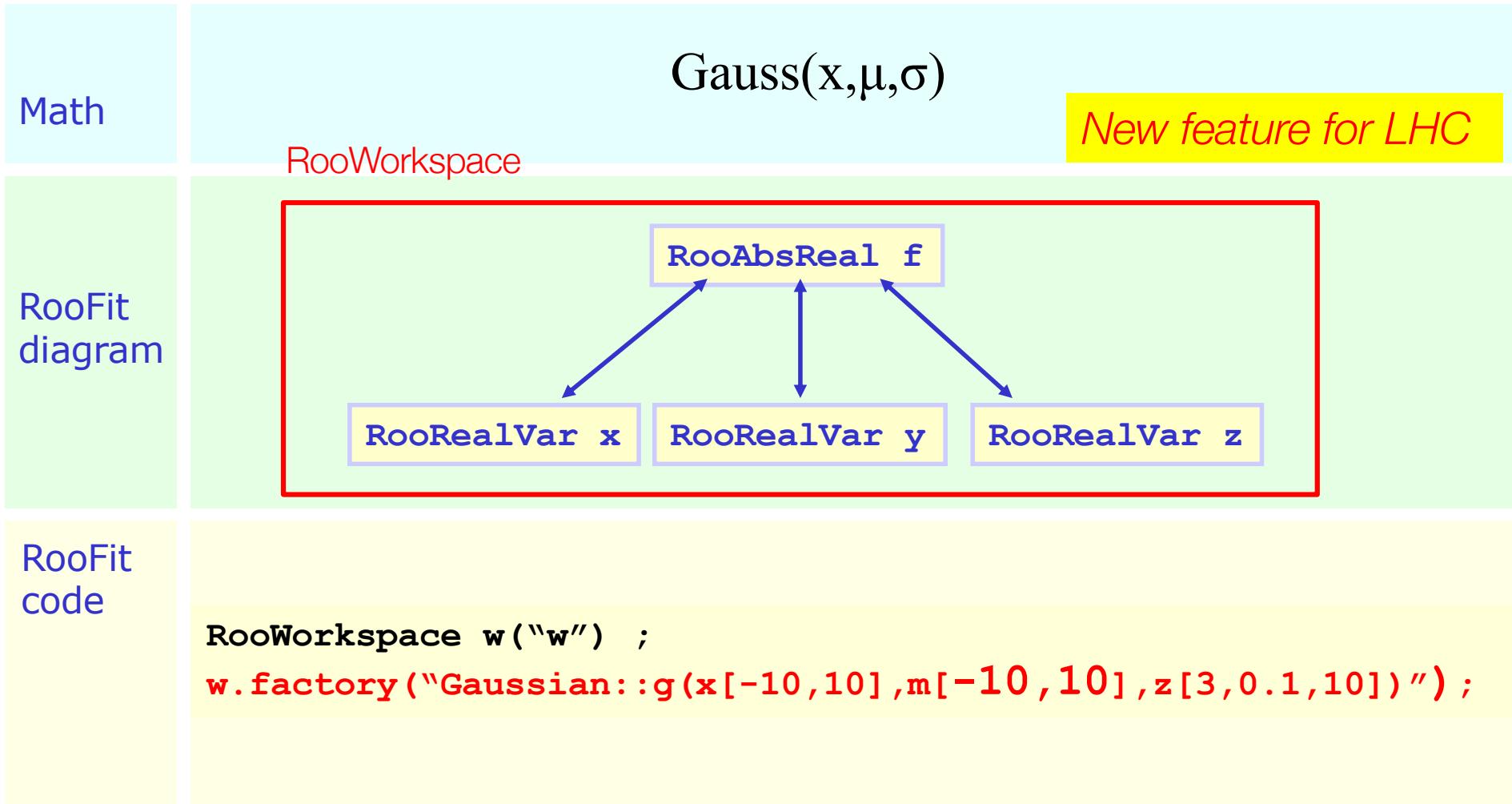
# RooFit core design philosophy

- A special container class owns all objects that together build a likelihood function



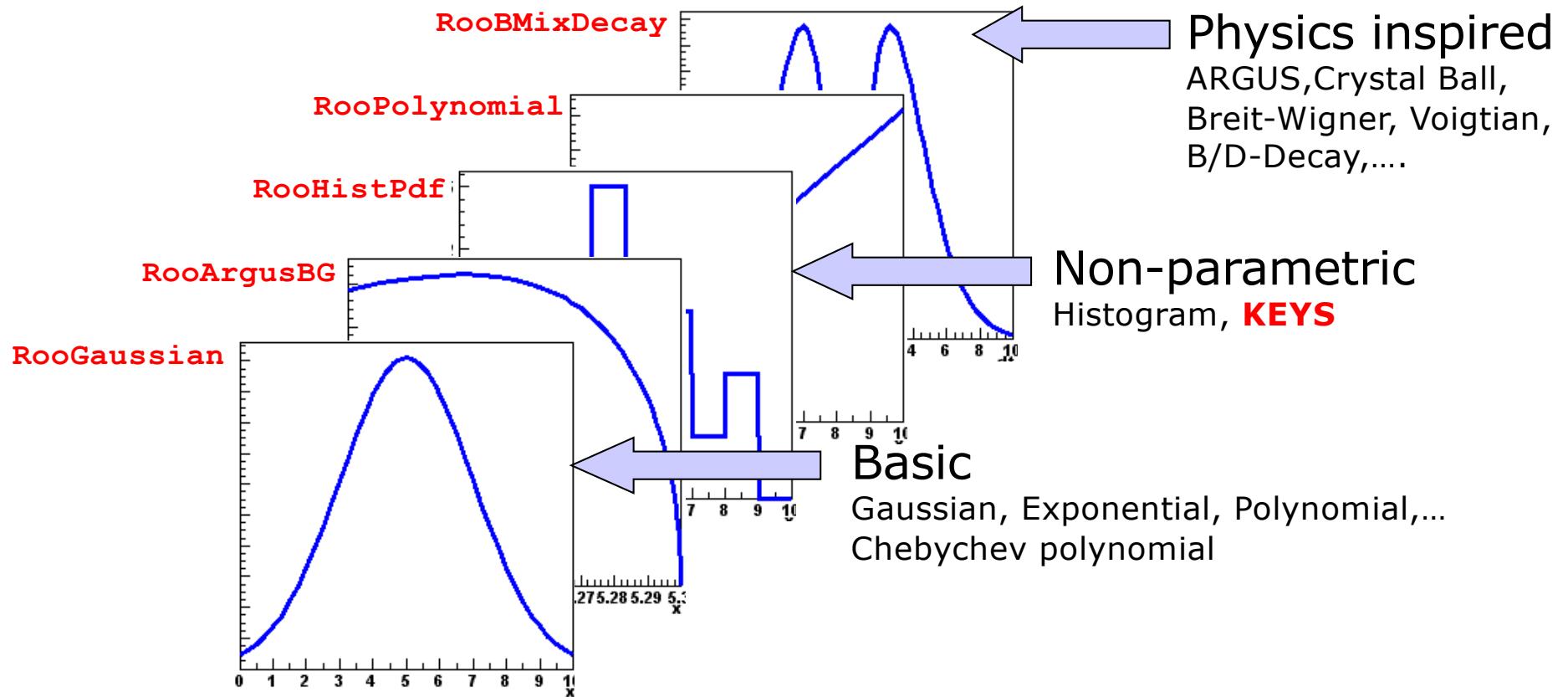
## Populating a workspace the easy way – “the factory”

- The **factory** allows to fill a workspace with pdfs and variables using a simplified scripting language



## Model building – (Re)using standard components

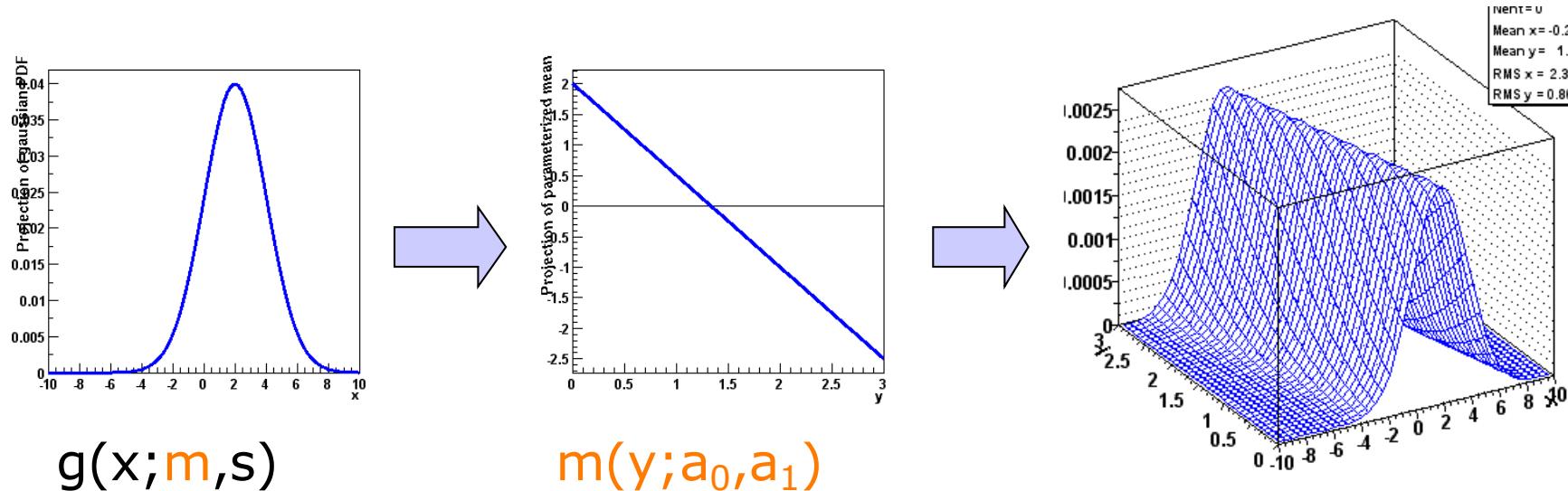
- RooFit provides a collection of compiled standard PDF classes



***Easy to extend the library: each p.d.f. is a separate C++ class***

## Model building – (Re)using standard components

- Library p.d.f.s can be adjusted on the fly.
  - Just plug in *any function expression* you like as input variable
  - Works universally, even for classes you write yourself



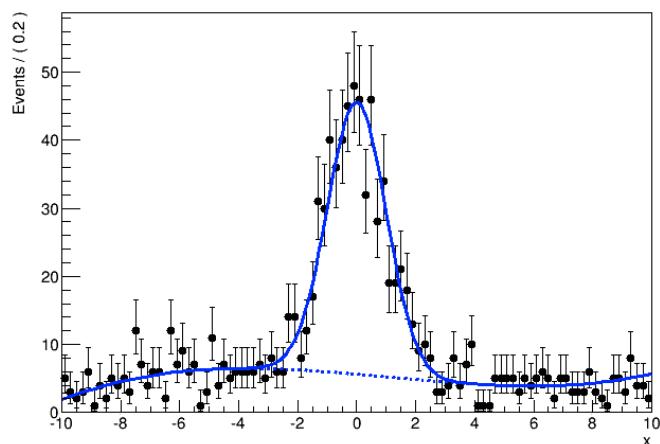
```
RooPolyVar m("m",y,RooArgList(a0,a1)) ;  
RooGaussian g("g","gauss",x,m,s) ;
```

- Maximum flexibility of library shapes keeps library small

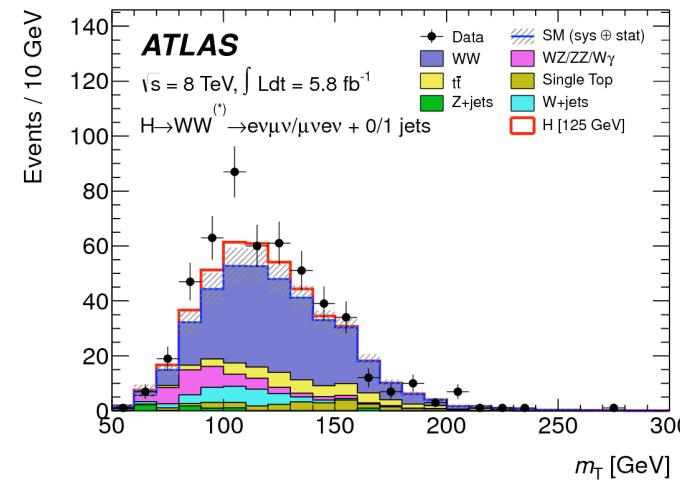
## From empirical probability models to simulation-based models

- Large difference between B-physics and LHC hadron physics is that for the latter distributions usually don't follow simple analytical shapes

*Unbinned analytical probability model*

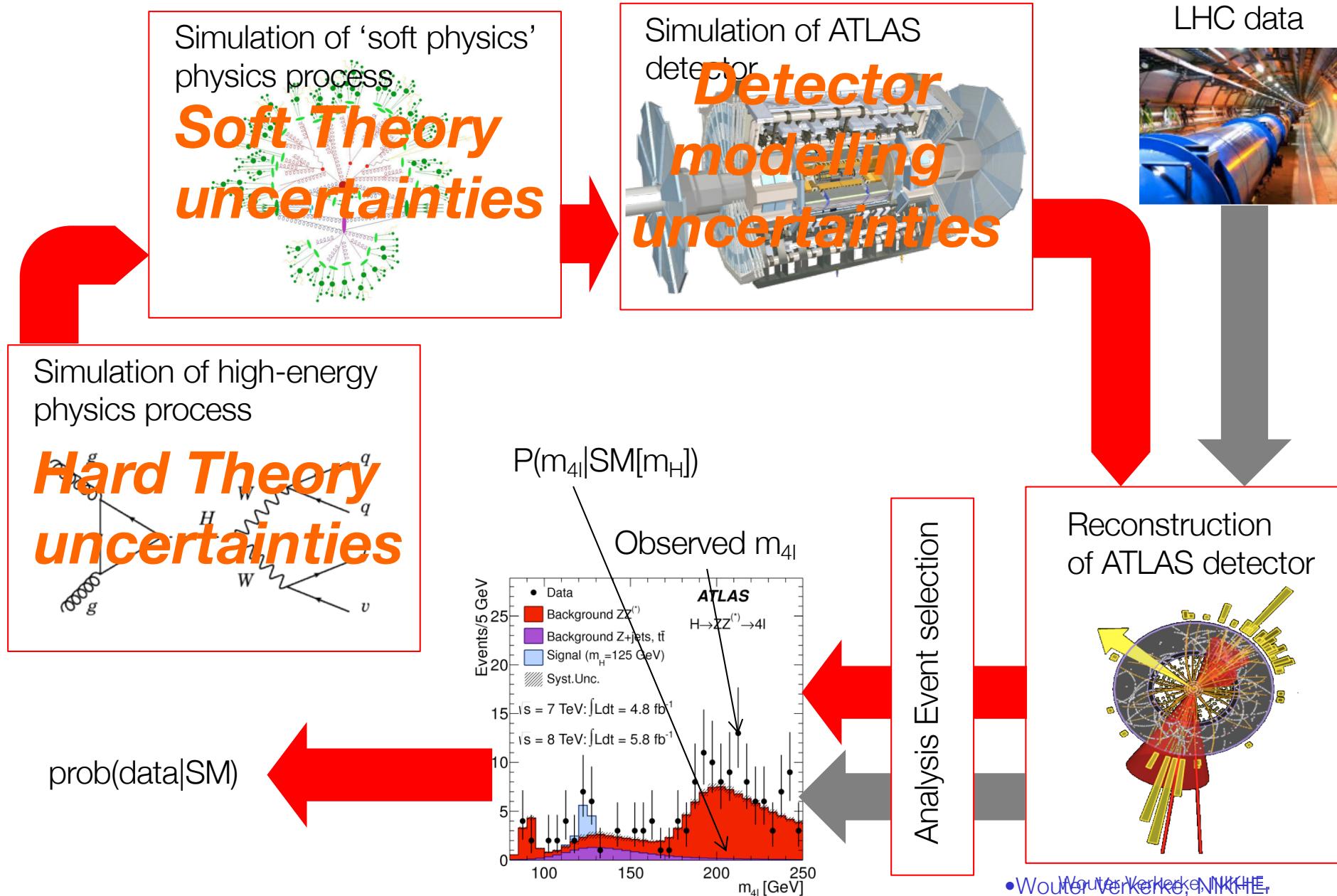


*(Geant) Simulation-driven binned template model*



- But concept of simulation-driven template models can also extend to systematic uncertainties. Instead of empirically chosen ‘nuisance parameters’ (e.g. polynomial coeffs) construct degrees of freedom that correspond to known systematic uncertainties

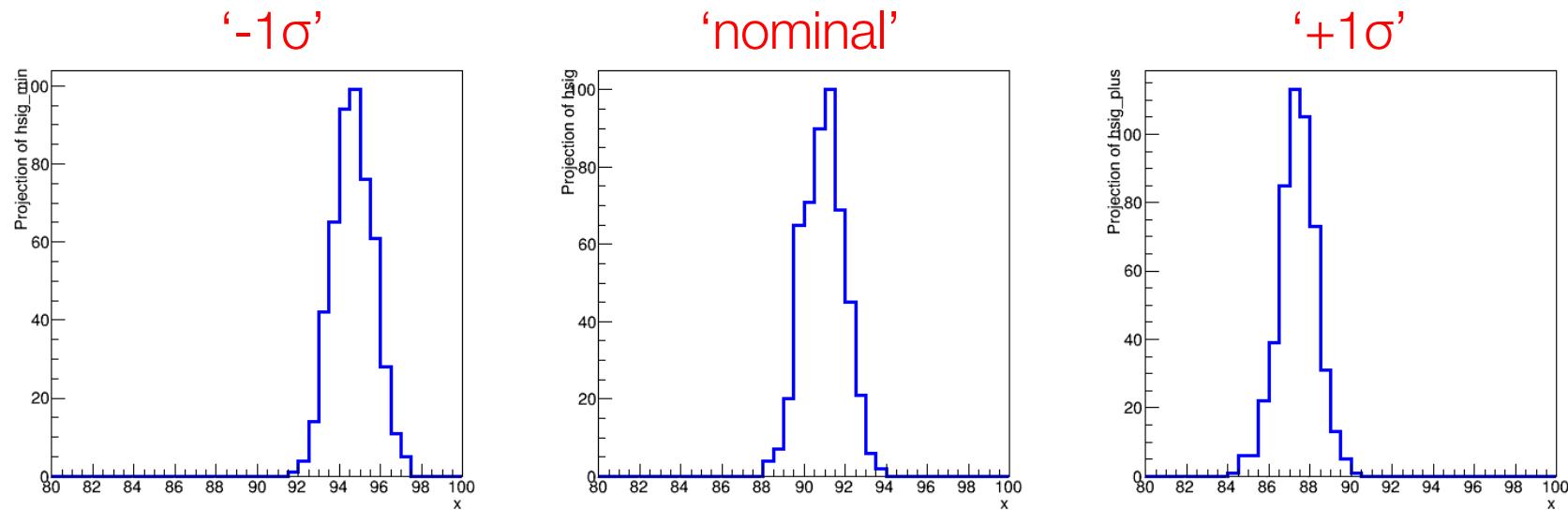
# The HEP analysis workflow illustrated



# Modeling of shape systematics in the likelihood

- Effect of *any* systematic uncertainty that affects the shape of a distribution can in principle be obtained from MC simulation chain
  - Obtain histogram templates for distributions at ‘ $+1\sigma$ ’ and ‘ $-1\sigma$ ’ settings of systematic effect

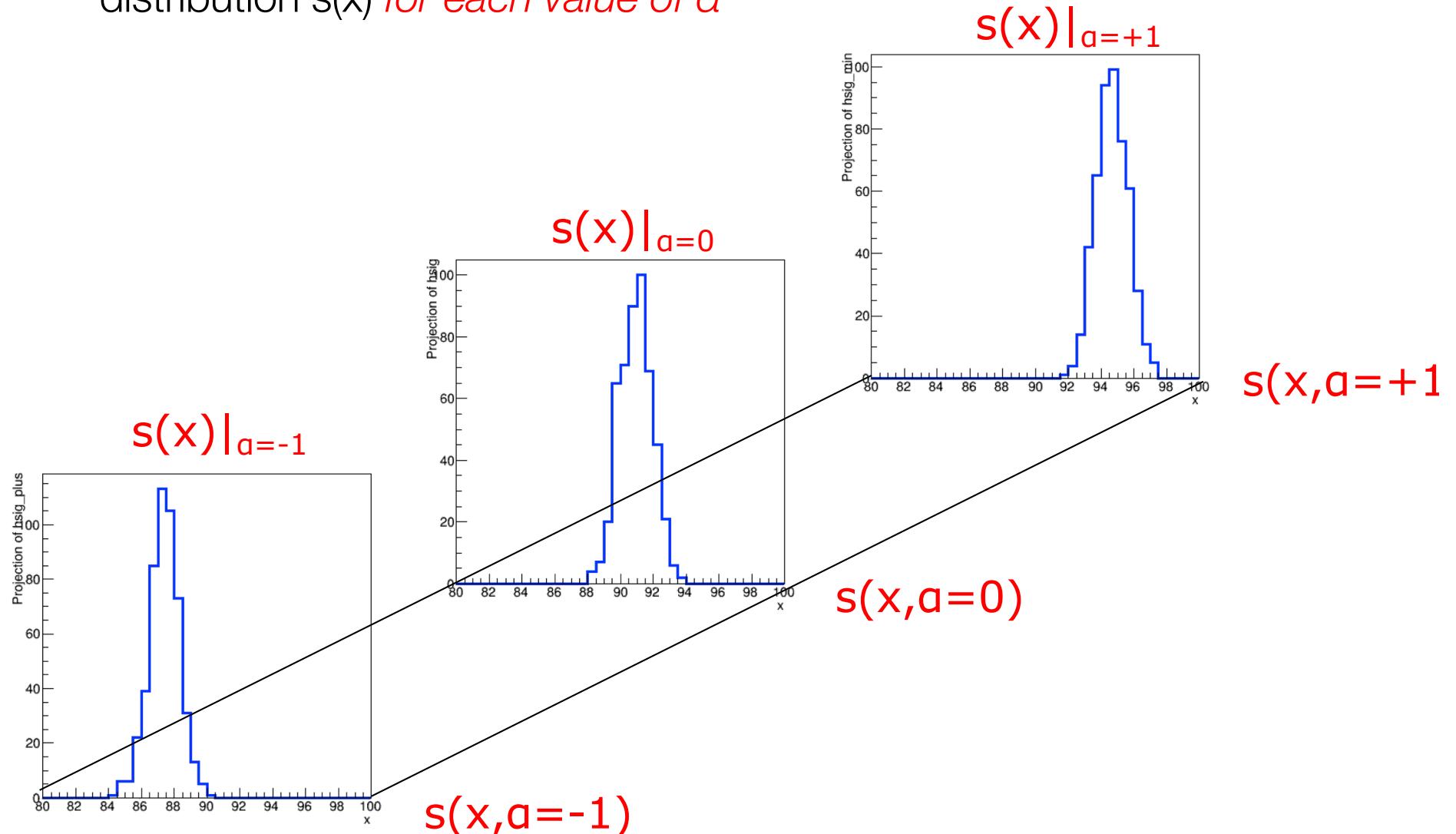
“Jet Energy Scale”



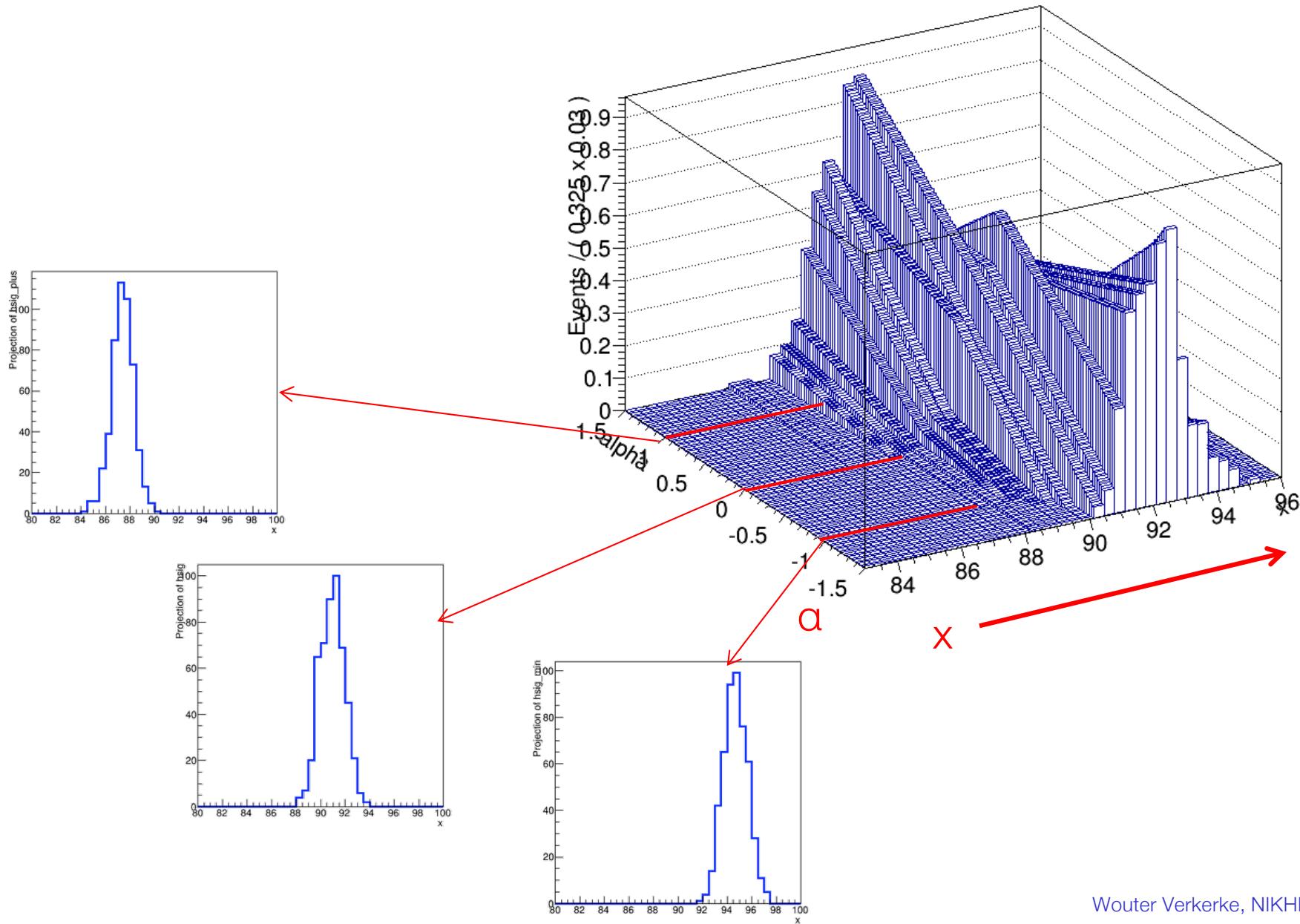
- Challenge: **construct an empirical response function based on the interpolation of the shapes of these three templates.**

## Need to interpolate between template models

- Need to define ‘morphing’ algorithm to define distribution  $s(x)$  *for each value of  $a$*



# Visualization of bin-by-bin linear interpolation of distribution

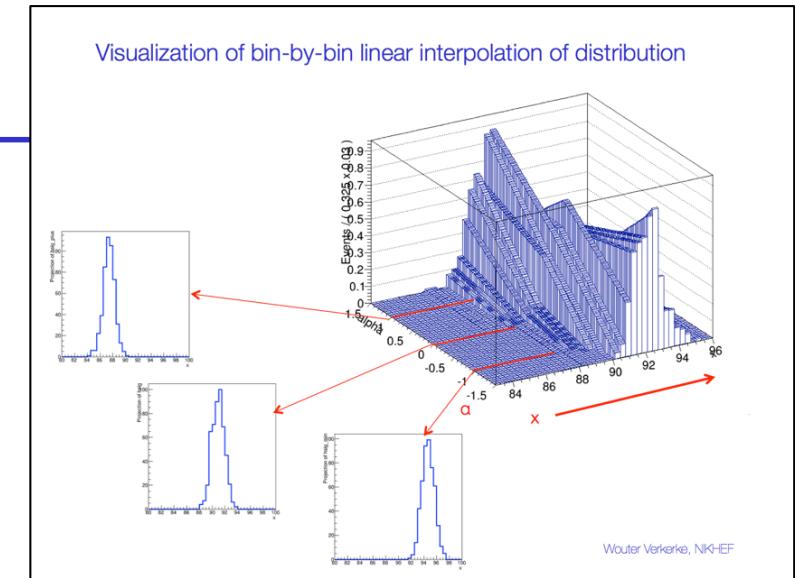


## Example 2 : binned L with syst

- Example of template morphing systematic in a binned likelihood

$$s_i(\alpha, \dots) = \begin{cases} s_i^0 + \alpha \cdot (s_i^+ - s_i^0) & \forall \alpha > 0 \\ s_i^0 + \alpha \cdot (s_i^0 - s_i^-) & \forall \alpha < 0 \end{cases}$$

$$L(\vec{N} | \alpha, \vec{s}^-, \vec{s}^0, \vec{s}^+) = \prod_{bins} P(N_i | s_i(\alpha, s_i^-, s_i^0, s_i^+)) \cdot G(0 | \alpha, 1)$$



```
// Import template histograms in workspace
w.import(hs_0, hs_p, hs_m) ;

// Construct template models from histograms
w.factory("HistFunc::s_0(x[80,100],hs_0)" );
w.factory("HistFunc::s_p(x,hs_p)" );
w.factory("HistFunc::s_m(x,hs_m)" );

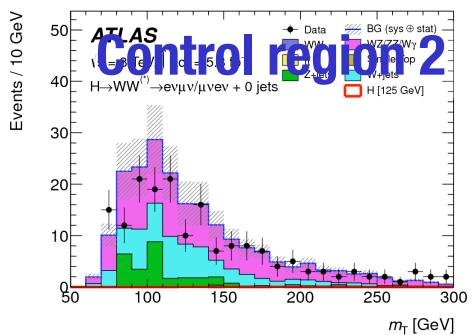
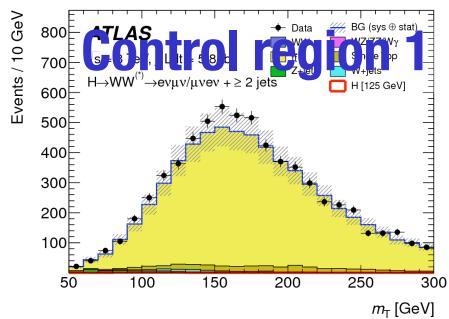
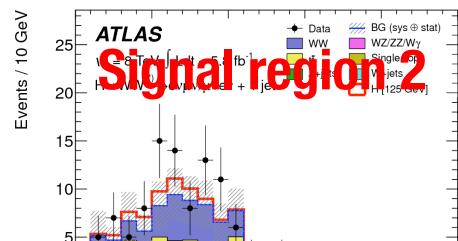
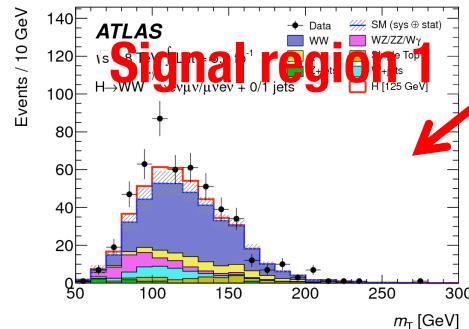
// Construct morphing model
w.factory("PiecewiseInterpolation::sig(s_0,s_,m,s_p,alpha[-5,5])" );

// Construct full model
w.factory("PROD::model(ASUM(sig,bkg,f[0,1]),Gaussian(0,alpha,1))" );
```

# The structure of an (Higgs) profile likelihood function

- Likelihood describing Higgs samples have following structure

$$L_{H \rightarrow X}(x | \mu, \vec{\theta}) = \prod_{i=0 \dots n} L_{phys}(x | \mu, \vec{\theta}) \cdot \prod_{i=0 \dots n} L_{control}(x | \mu, \vec{\theta}) \cdot L_{sub}(\theta_1) \cdot L_{sub}(\theta_1) \cdots L_{sub}(\theta_n)$$



*Strength of systematic uncertainties*

'Constraint  $\theta_1$ '

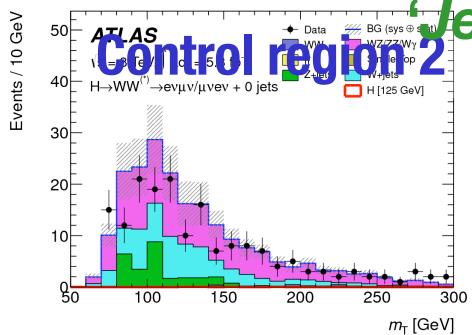
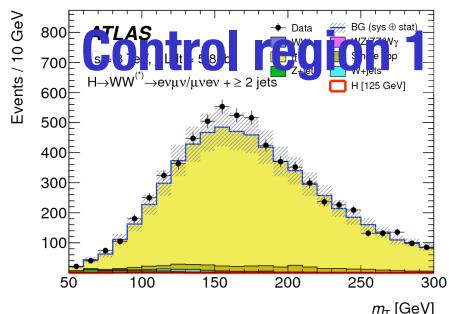
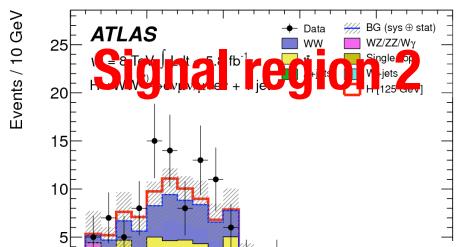
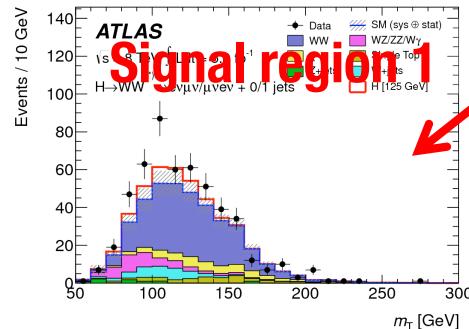
'Constraint  $\theta_n$ '

'Constraint  $\theta_n$ '

# The structure of an (Higgs) profile likelihood function

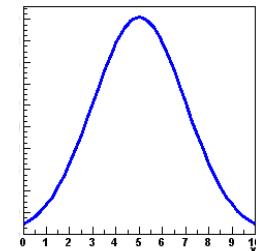
- A simultaneous fit of physics samples and (simplified) performance measurements

$$L_{H \rightarrow X}(x | \mu, \vec{\theta}) = \prod_{i=0 \dots n} L_{phys}(x | \mu, \vec{\theta}) \cdot \prod_{i=0 \dots n} L_{control}(x | \mu, \vec{\theta}) \cdot L_{sub}(\theta_1) \cdot L_{sub}(\theta_2) \cdots L_{sub}(\theta_n)$$



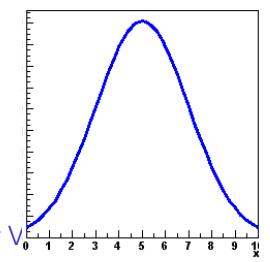
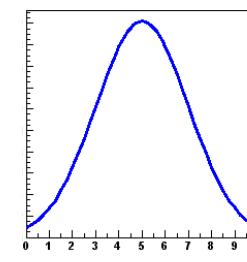
'Simplified Likelihood of Factorization scale  
a measurement related  
to systematic uncertainties'

'Subsidiary  
measurement n'



'Subsidiary  
measurement 1'      'Subsidiary  
measurement 2'

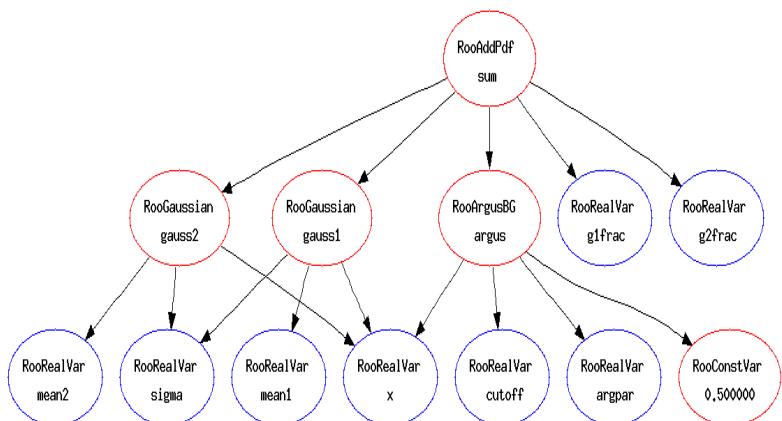
'jet Energy scale'  
'B-tagging eff'



# The Workspace

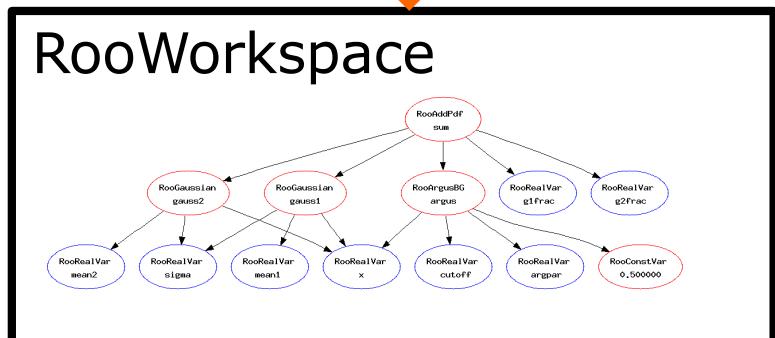
# The workspace

- The workspace concept has revolutionized the way people share and combine analysis
  - Completely factorizes process of building and using likelihood functions
  - You can give somebody an analytical likelihood of a (potentially very complex) physics analysis in a way to the easy-to-use, provides introspection, and is easy to modify.

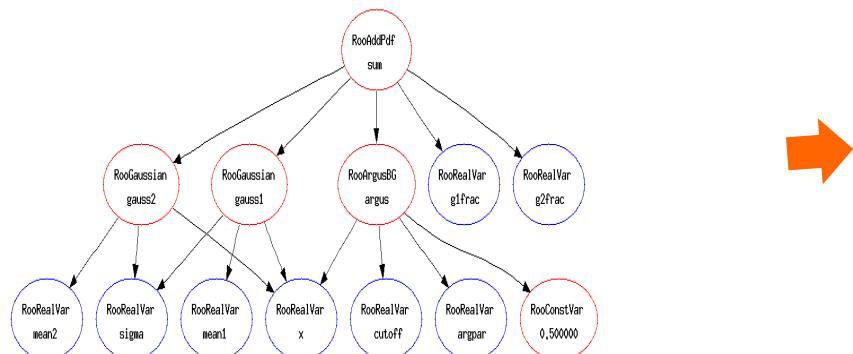
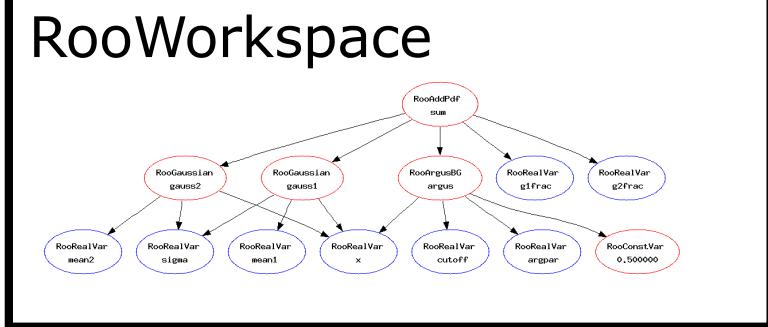


```
RooWorkspace w("w");  
w.import(sum);  
w.writeToFile("model.root");
```

model.root

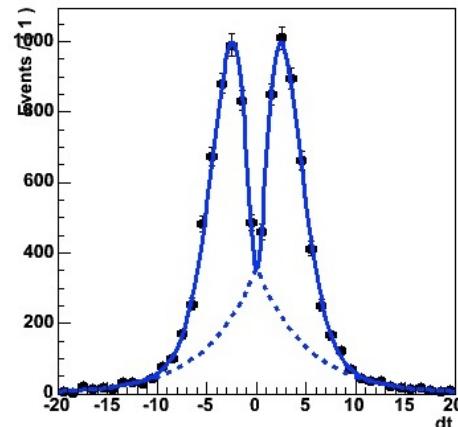


# Using a workspace



```
// Resurrect model and data
TFile f("model.root") ;
RooWorkspace* w = f.Get("w") ;
RooAbsPdf* model = w->pdf("sum") ;
RooAbsData* data = w->data("xxx") ;
```

```
// Use model and data
model->fitTo(*data) ;
RooPlot* frame =
    w->var("dt")->frame() ;
data->plotOn(frame) ;
model->plotOn(frame) ;
```



## The idea behind the design of RooFit/RooStats/HistFactory

---

- **Step 1** – Construct the likelihood function  $L(x|p)$

```
RooWorkspace w("w") ;  
w.factory("Gaussian::sig(x[-10,10],m[0],s[1])") ;  
w.factory("Chebychev::bkg(x,a1[-1,1])") ;  
w.factory("SUM::model(fsig[0,1]*sig,bkg)") ;  
w.writeToFile("L.root") ;
```



**RooWorkspace**



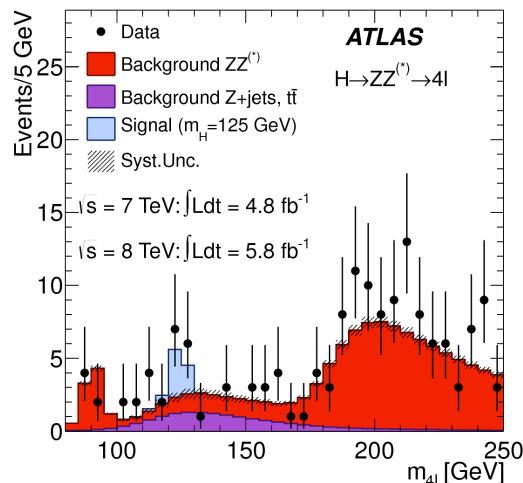
*Complete description  
of likelihood model,  
persistable in ROOT file  
(RooFit pdf function)  
Allows full introspection  
and a-posteriori editing*

- **Step 2** – Statistical tests on parameter of interest  $p$

```
RooWorkspace* w=TFile::Open("L.root")->Get("w") ;  
RooAbsPdf* model = w->pdf("model") ;  
pdf->fitTo(data) ;
```

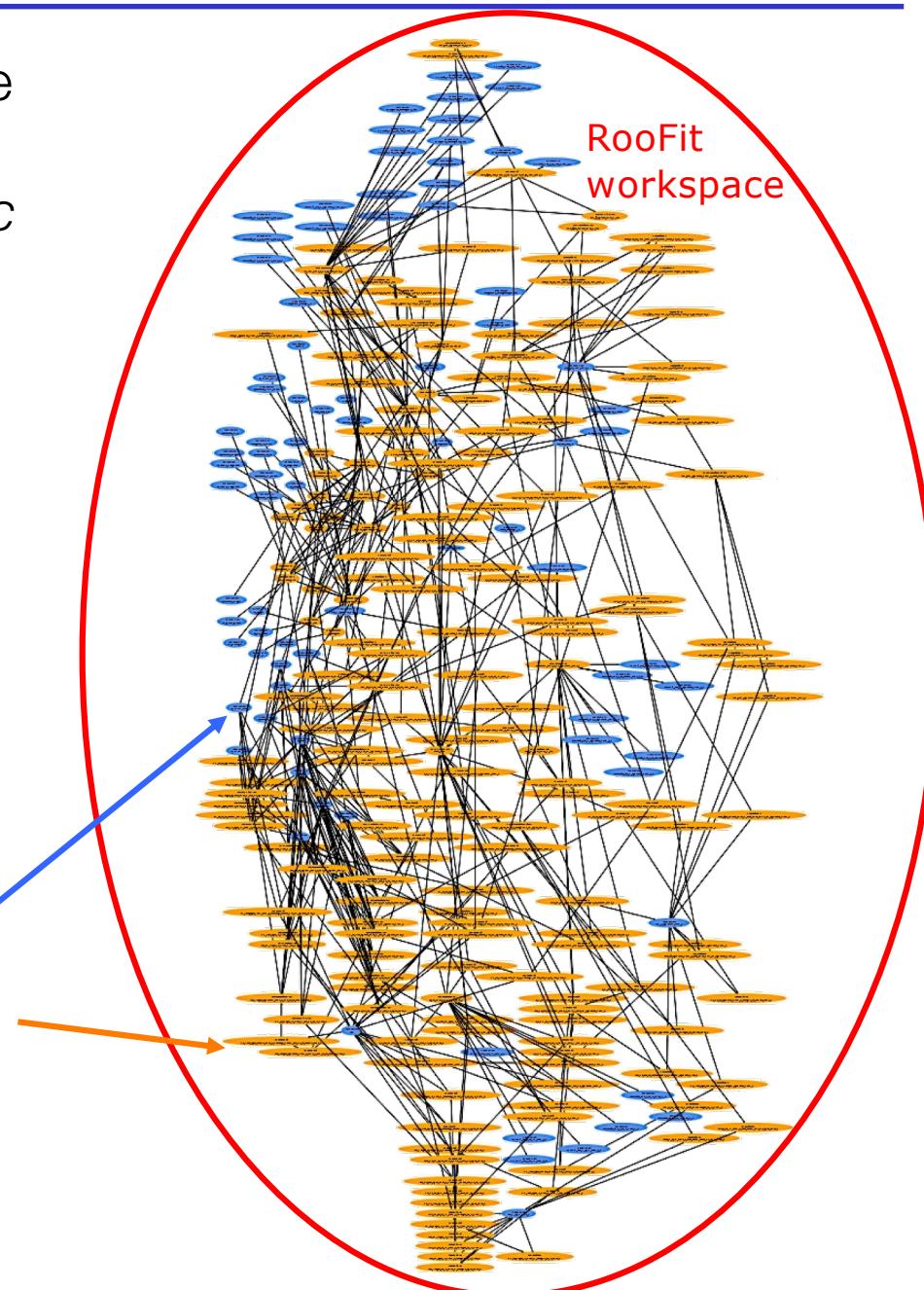
## Example RooFit component model for realistic Higgs analysis

Likelihood model describing the ZZ invariant mass distribution *including all possible systematic uncertainties*



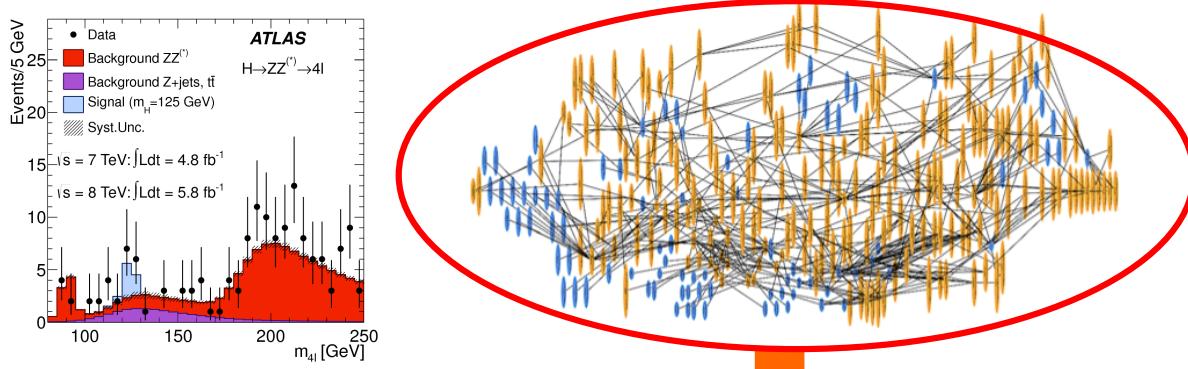
variables  
function objects

Graphical illustration of function components that call each other



# Analysis chain identical for highly complex (Higgs) models

- **Step 1** – Construct the likelihood function  $L(x|p)$



*Complete description  
of likelihood model,  
persistable in ROOT file  
(RooFit pdf function)  
Allows full introspection  
and a-posteriori editing*

- **Step 2** – Statistical tests on parameter of interest  $p$

```
RooWorkspace* w=TFile::Open("L.root")->Get("w") ;  
RooAbsPdf* model = w->pdf("model") ;  
pdf->fitTo(data,  
            GlobalObservables(w->set("MC_G10bs")),  
            Constrain(*w->st("MC_NuisParams")) ;
```

# Workspaces power collaborative statistical modelling

---

- Ability to persist complete<sup>(\*)</sup> Likelihood models has profound implications for HEP analysis workflow
  - (\*) Describing signal regions, control regions, and including nuisance parameters for all systematic uncertainties)
- **Anyone with ROOT (and one ROOT file with a workspace) can re-run any entire statistical analysis out-of-the-box**
  - About 5 lines of code are needed
  - Including estimate of systematic uncertainties
- Unprecedented new possibilities for cross-checking results, in-depth checks of structure of analysis
  - Trivial to run variants of analysis (what if ‘Jet Energy Scale uncertainty’ is 7% instead of 4%). Just change number and rerun.
  - But can also make structural changes a posteri. For example, rerun with assumption that JES uncertainty in forward and barrel region of detector are 100% correlated instead of being uncorrelated.

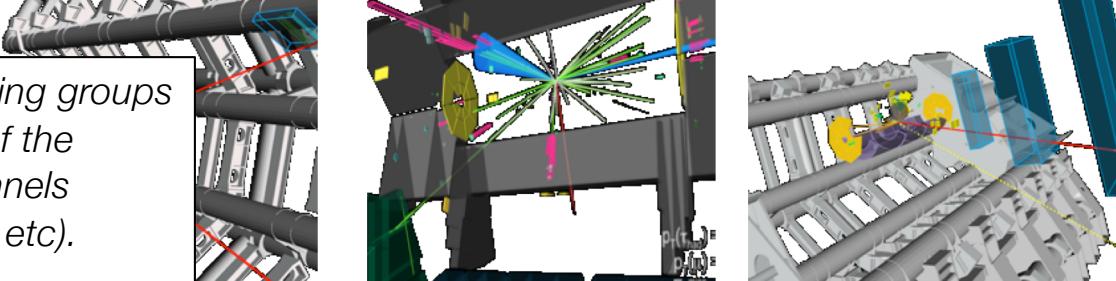
# Collaborative statistical modelling

---

- As an experiment, you can effectively **build a library of measurements**, of which the full likelihood model is preserved for later use
  - Already done now, experiments have such libraries of workspace files,
  - Archived in AFS directories, or even in SVN....
  - Version control of SVN, or numbering scheme in directories allows for easy validation and debugging as new features are added
- Building of combined likelihood models greatly simplified.
  - Start from persisted components. No need to (re)build input components.
  - No need to know how individual components were built, or are internally structured. Just need to know meaning of parameters.
  - Combinations can be produced (much) later than original analyses.
  - Even analyses that were never originally intended to be combined with anything else can be included in joint likelihoods at a later time

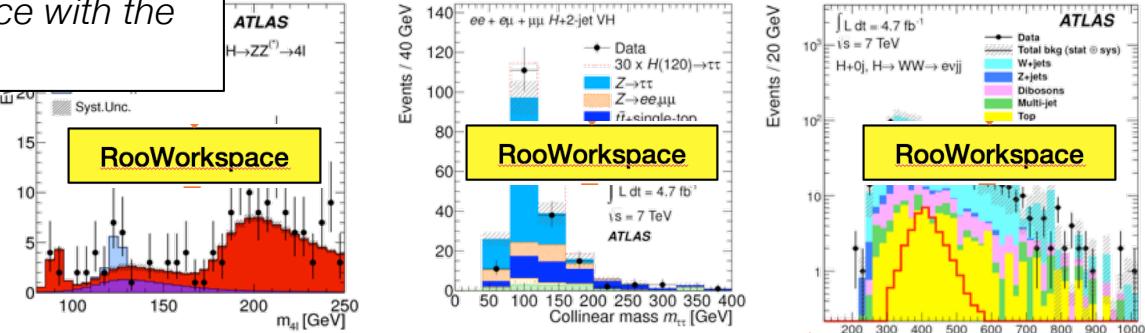
# Higgs discovery strategy – add everything together

$H \rightarrow ZZ \rightarrow llll$        $H \rightarrow \tau\tau$        $H \rightarrow WW \rightarrow \mu\nu jj$

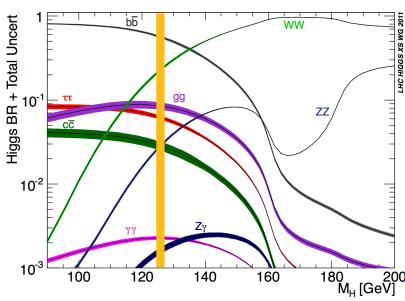


*Dedicated physics working groups define search for each of the major Higgs decay channels ( $H \rightarrow WW$ ,  $H \rightarrow ZZ$ ,  $H \rightarrow \tau\tau$  etc).*

*Output is physics paper or note, and a RooFit workspace with the full likelihood function*



**Assume SM rates**



$L(\mu, \vec{\theta}) = L_{H \rightarrow WW}(\mu_{WW}, \vec{\theta}) \cdot L_{H \rightarrow \gamma\gamma}(\mu_{\gamma\gamma}, \vec{\theta}) \cdot L_{H \rightarrow ZZ}(\mu_{ZZ}, \vec{\theta}) \cdot \dots$

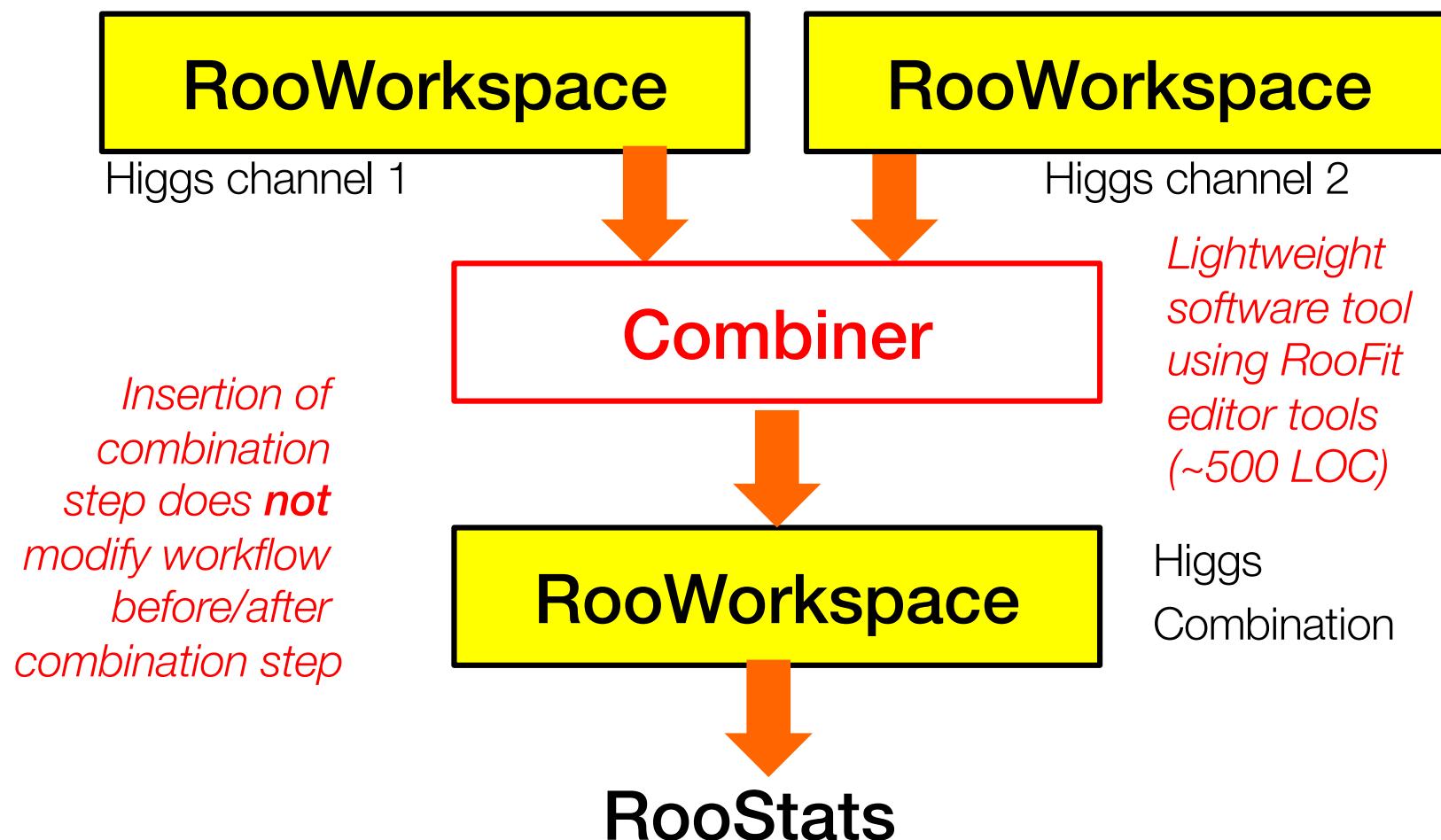
*A small dedicated team of specialists builds a combined likelihood from the inputs.  
Major discussion point: naming of parameters, choice of parameters for systematic uncertainties (a physics issue, largely)*

## The benefits of modularity

---

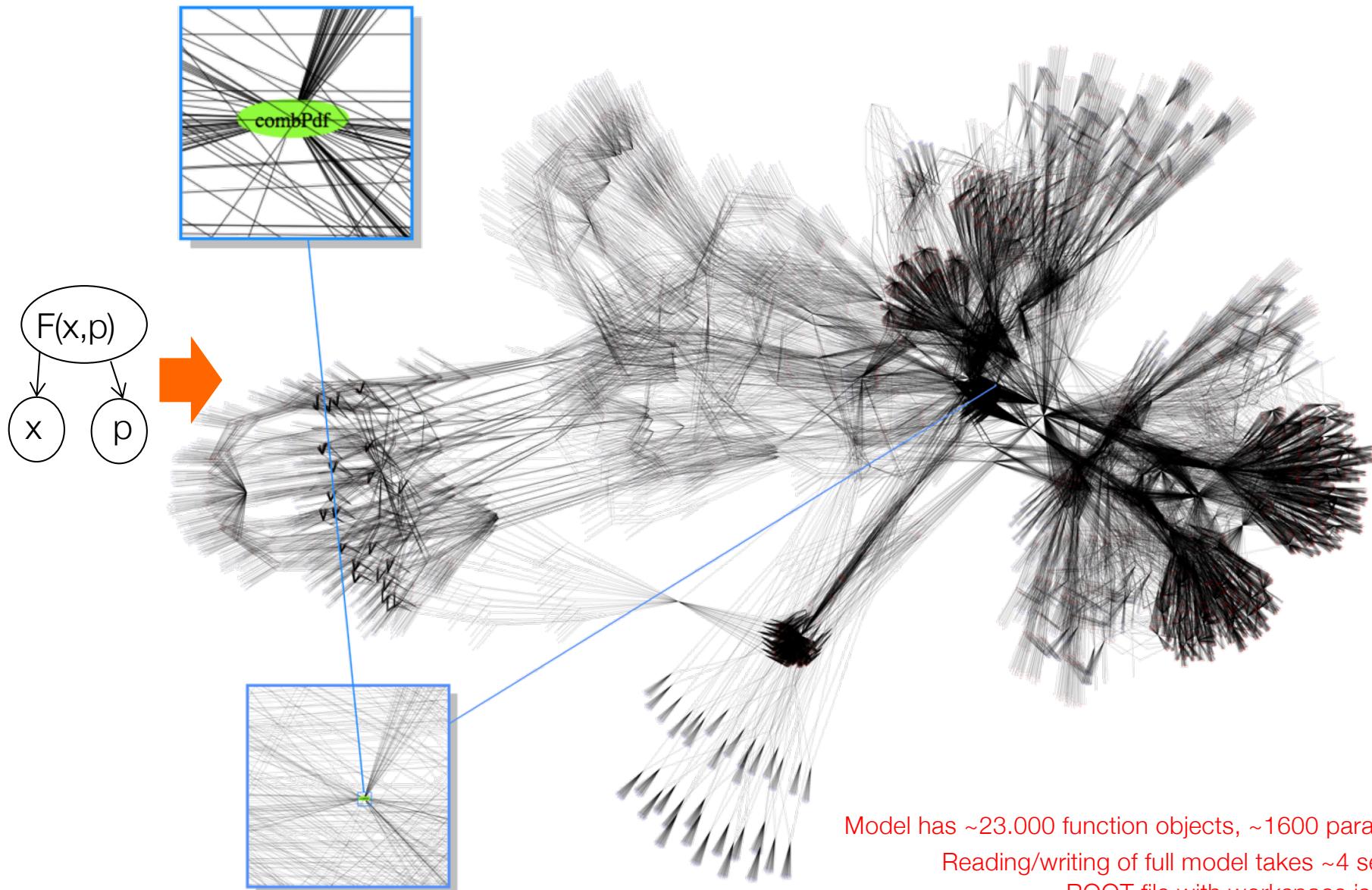
- Technically very straightforward to combine measurements

### RooFit, or RooFit+HistFactory



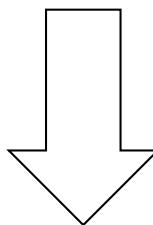
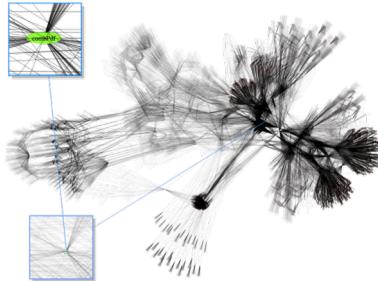
## Workspace persistence of *really* complex models works too!

Atlas Higgs combination model (23.000 functions, 1600 parameters)

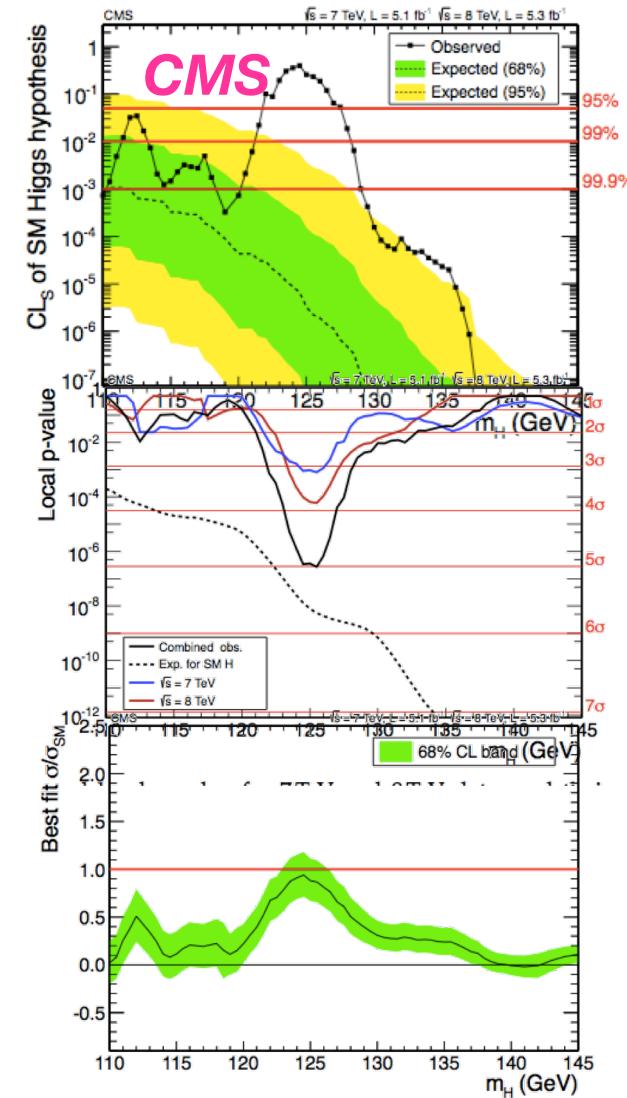
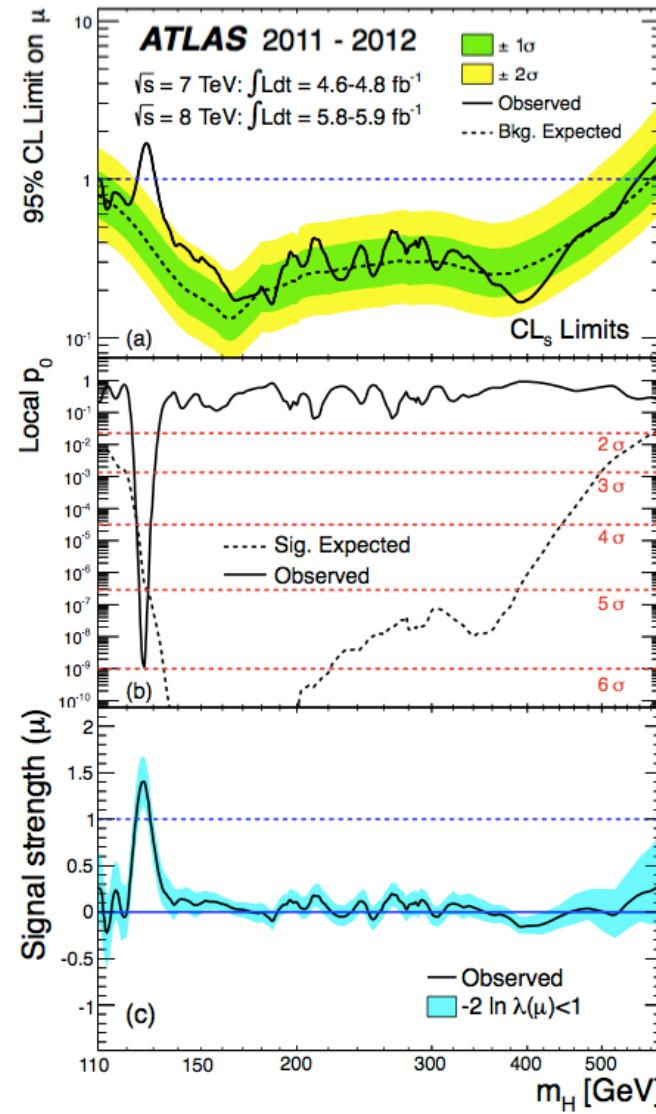
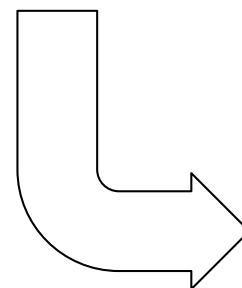


With these combined models the Higgs discovery plots were produced...

$$L_{\text{ATLAS}}(\mu, \theta) =$$



Neyman construction  
with profile likelihood  
ratio test

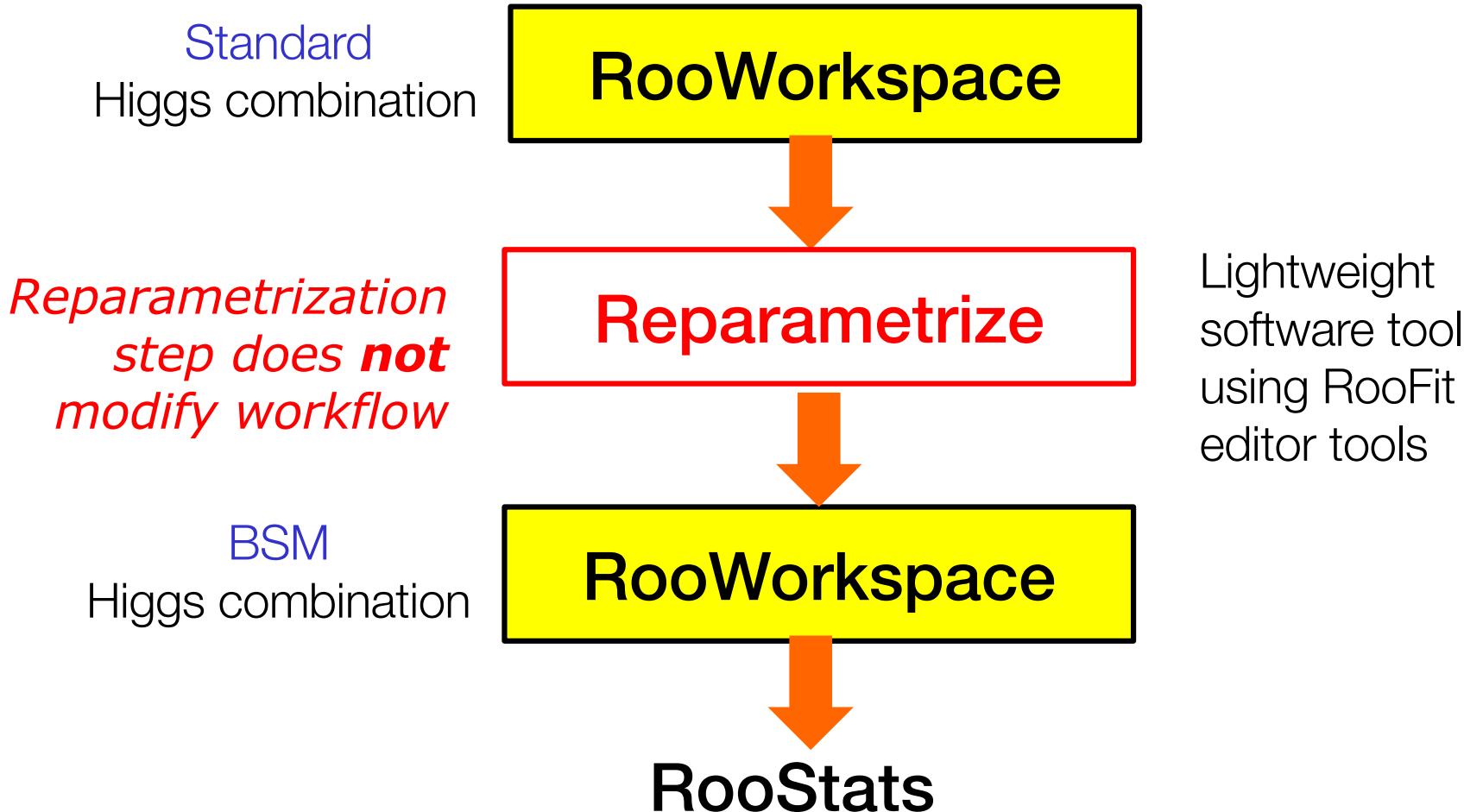


## More benefits of modularity

---

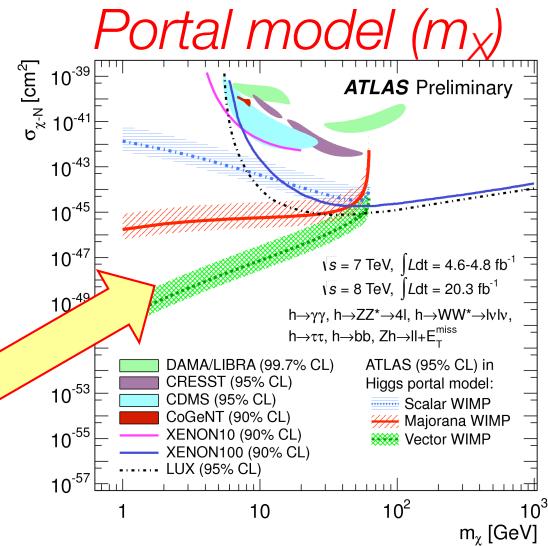
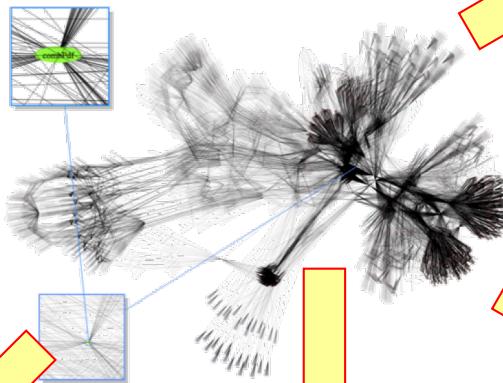
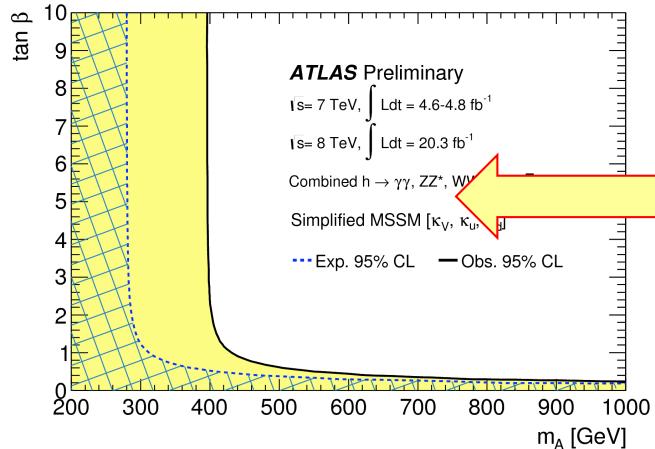
- Technically very straightforward to reparametrize measurements

## RooFit, or RooFit+HistFactory

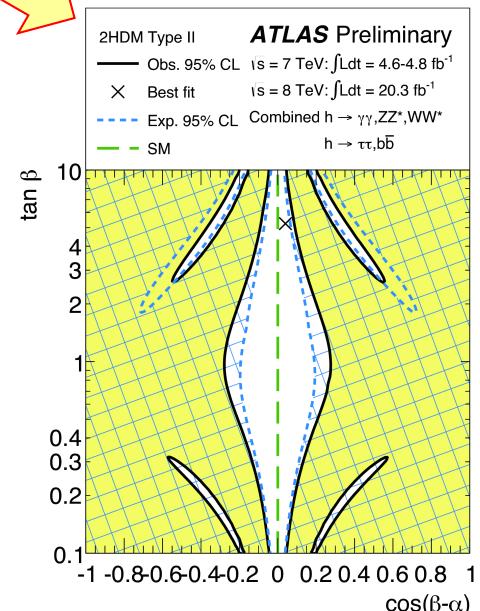
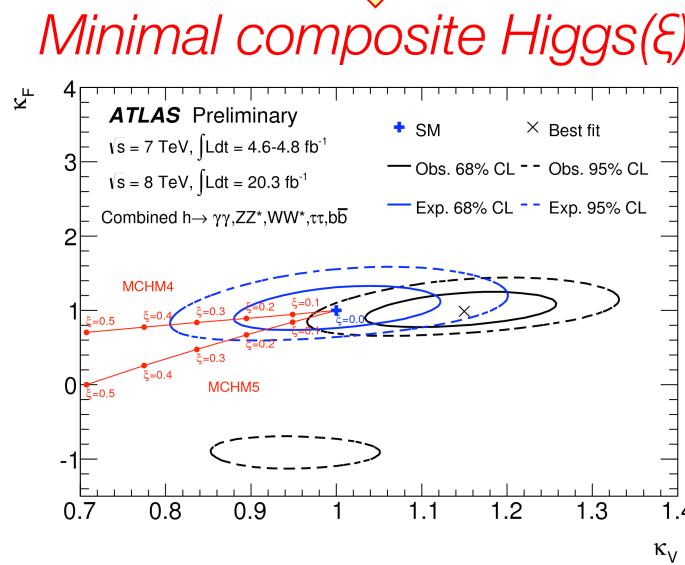
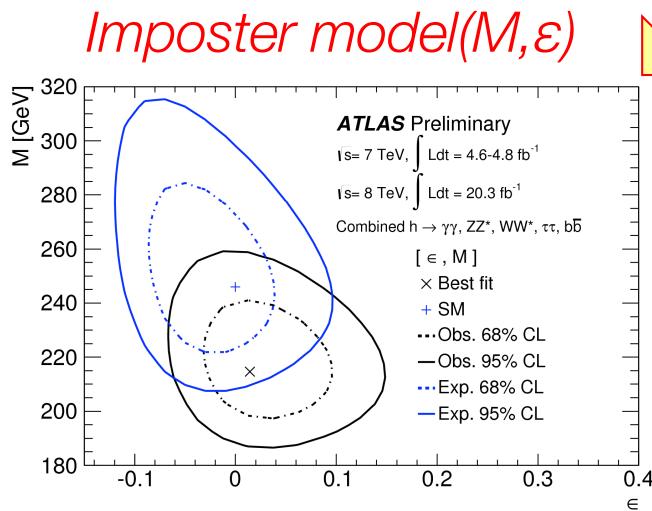


# BSM Higgs constraints from reparametrization of SM Higgs Likelihood model

## Simplified MSSM ( $\tan\beta, m_A$ )



## Two Higgs Double Model ( $\tan\beta, \cos(\alpha-\beta)$ )



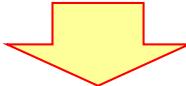
(ATLAS-CONF-2014-010)

Wouter Verkerke, NIKHEF

## An excursion – Collaborative analyses with workspaces

---

- *How can you reparametrize existing Higgs likelihoods in practice?*
- Write functions expressions corresponding to new parameterization

$$\sigma(gg \rightarrow H) * \text{BR}(H \rightarrow \gamma\gamma) \sim \frac{\kappa_F^2 \cdot \kappa_\gamma^2(\kappa_F, \kappa_V)}{0.75 \cdot \kappa_F^2 + 0.25 \cdot \kappa_V^2}$$


```
w.factory("expr::mu_gg_func(' (KF2*Kg2)/  
                           (0.75*KF2+0.25*KV2)',  
                           KF2,Kg2,KV2) ;
```

- Import transformation in workspace, edit *existing* model

```
w.import(mu_gg_func) ;  
w.factory("EDIT::newmodel(model,mu_gg=mu_gg_func)");
```

# HistFactory

K. Cranmer, A. Shibata, G. Lewis, L. Moneta, W. Verkerke (2010)

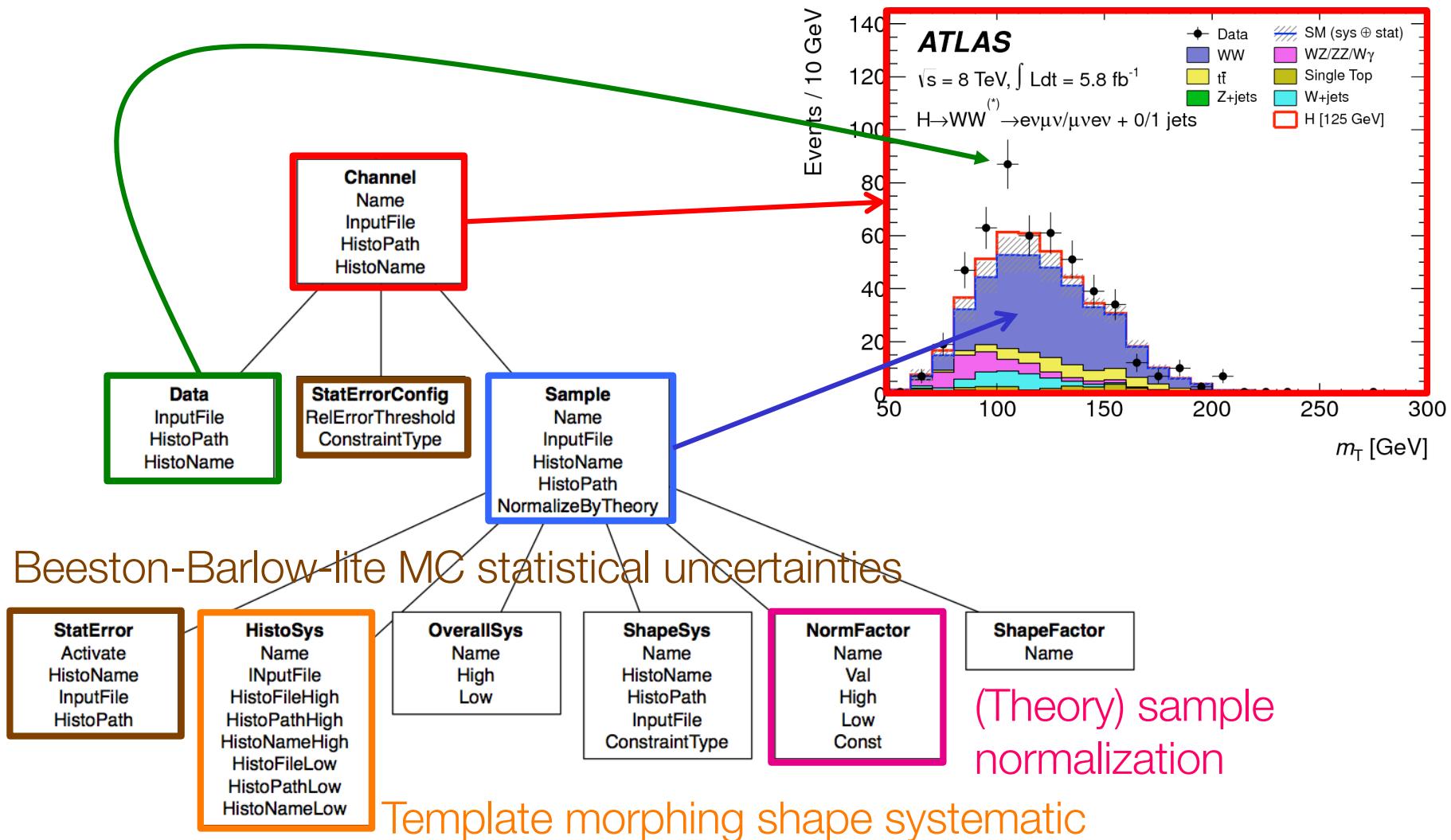
## HistFactory – structured building of binned template models

---

- RooFit modeling building blocks allow to easily construct likelihood models that model shape and rate systematics with one or more nuisance parameter
  - Only few lines of code per construction
- Typical LHC analysis required modeling of 10-50 systematic uncertainties in  $O(10)$  samples in anywhere between 2 and 100 channels → **Need structured formalism to piece together model from specifications.** This is the purpose of HistFactory
- HistFactory conceptually similar to workspace factory, but has much higher level semantics
  - Elements represent physics concepts (channels, samples, uncertainties and their relation) rather than mathematical concepts
  - Descriptive elements are represented by C++ objects (like roofit), and can be configured in C++, or alternatively from an XML file
- **HistFactory builds a RooFit (mathematical) model from a physics model.**

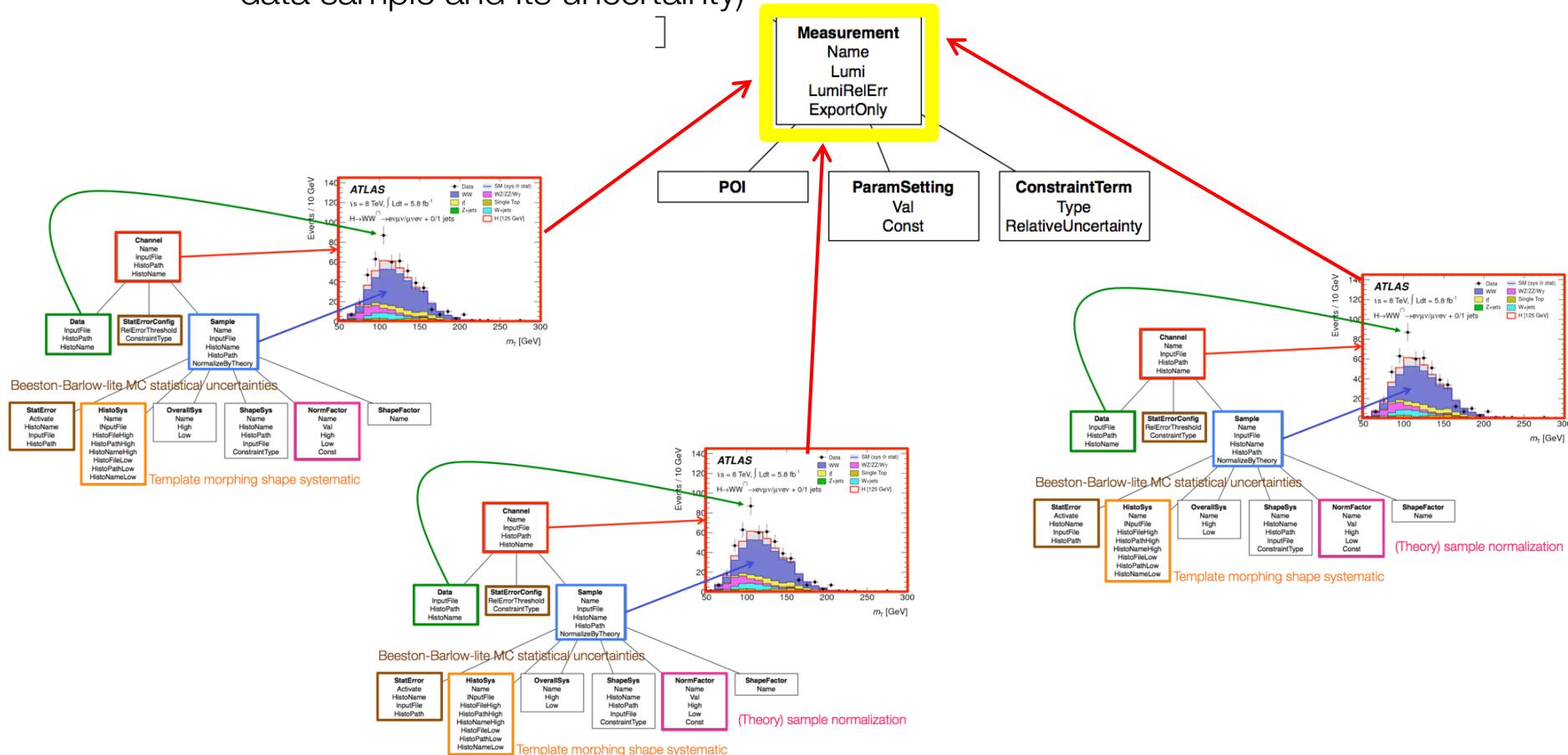
# HistFactory elements of a channel

- Hierarchy of concepts for description of one measurement channel



# HistFactory elements of measurement

- One or more **channels** are combined to form a **measurement**
  - Along with some extra information (declaration of the POI, the luminosity of the data sample and its uncertainty)

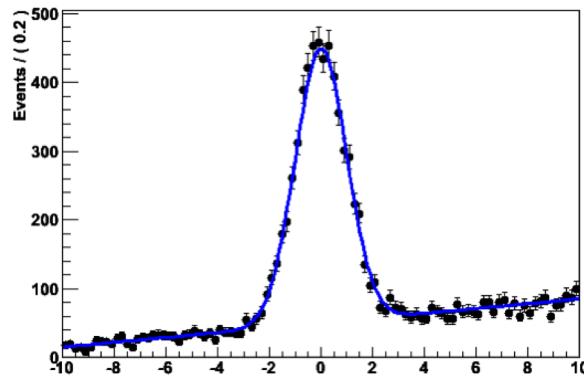


Once physics model is defined, one line of code will turn it into a RooFit likelihood

Wouter Verkerke, NIKHEF

# How is Higgs discovery different from a simple fit?

Gaussian + polynomial



ROOT TH1

ROOT TF1

$$L(\vec{N} \mid \mu, \vec{\theta}) = \prod_i \text{Poisson}(N_i \mid f(x_i, \mu, \vec{\theta}))$$

"inside ROOT"

Maximum Likelihood estimation of parameters  $\mu, \theta$  using MINUIT (MIGRAD, HESSE, MINOS)

$$\mu = 5.3 \pm 1.7$$

## Likelihood Model

orders of magnitude more complicated. Describes

- O(100) signal distributions
- O(100) control sample distr.
- O(1000) parameters representing syst. uncertainties

$$L(\vec{N}_{ZZ}, \vec{N}_{\tau\tau}, \vec{N}_{WW} \mid \mu, \vec{\theta}) = \prod \text{Poisson}(N_{ZZ}^i, \dots) \cdot \prod \text{Poisson}(N_{\tau\tau}^i, \dots) \cdot \prod \text{Poisson}(N_{WW}^i, \dots) \cdot \dots$$

Frequentist confidence interval construction and/or p-value calculation not available as 'ready-to-run' algorithm in ROOT

# RooStats

K. Cranmer, L. Moneta, S. Kreiss, G. Kukartsev, G. Schott, G. Petrucciani, WV - 2008

## The benefits of modularity

---

- Perform different statistical test on exactly the same model

**RooFit, or RooFit+HistFactory**



**RooWorkspace**



**“Simple fit”**

(ML Fit with  
HESSE or  
MINOS)

**RooStats**  
**(Frequentist  
with toys)**

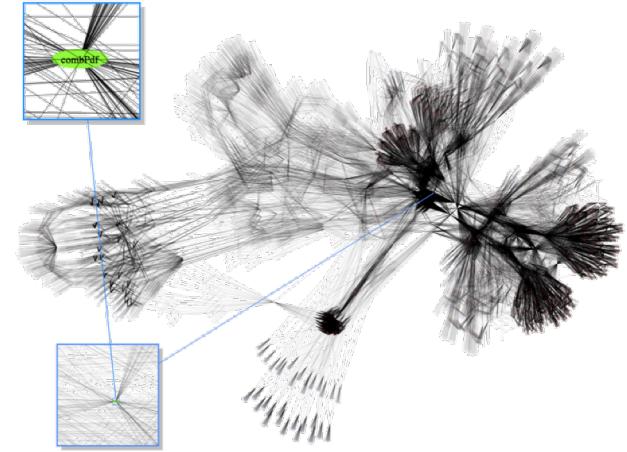
**RooStats**  
**(Frequentist  
asymptotic)**

**RooStats**  
**Bayesian  
MCMC**

# Maximum Likelihood estimation as simple statistical analysis

- **Step 1** – Construct the likelihood function  $L(x|p)$

```
RooWorkspace w("w") ;  
w.factory("Gaussian::sig(x[-10,10],m[0],s[1])";  
w.factory("Chebychev::bkg(x,a1[-1,1])");  
w.factory("SUM::model(fsig[0,1]*sig,bkg)");  
w.writeToFile("L.root") ;
```



**RooWorkspace**

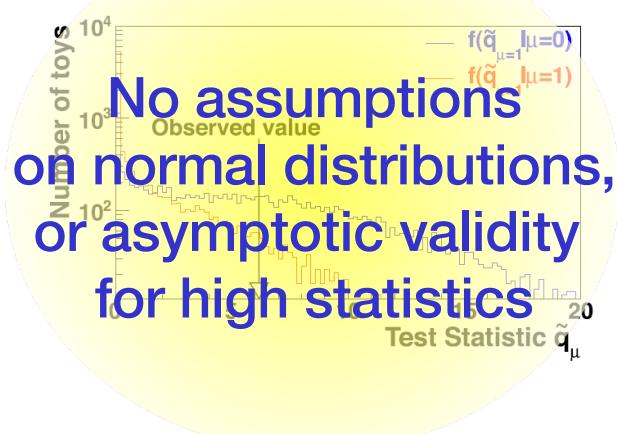
- **Step 2** – Statistical tests on parameter of interest  $p$

```
RooWorkspace* w=TFile::Open("L.root")->Get("w") ;  
RooAbsPdf* model = w->pdf("model") ;  
pdf->fitTo(data) ;
```

# The need for fundamental statistical techniques

Frequentist statistics

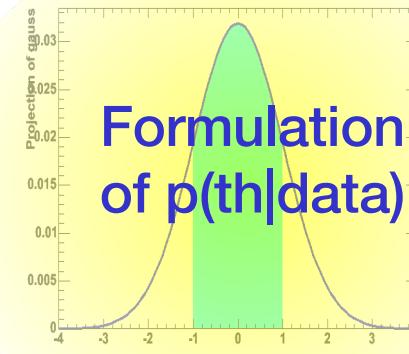
$$\lambda_\mu(\vec{N}_{obs}) = \frac{L(\vec{N} | \mu)}{L(\vec{N} | \hat{\mu})}$$



Confidence interval or p-value

Bayesian statistics

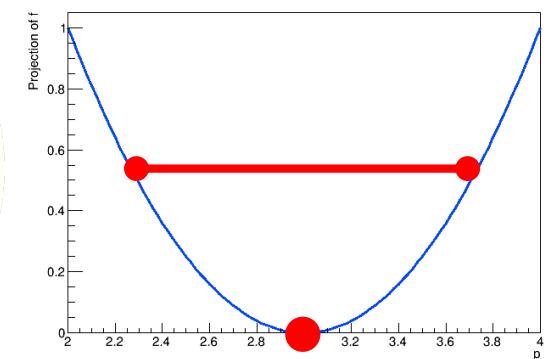
$$P(\mu) \propto L(x | \mu) \cdot \pi(\mu)$$



Posterior or Bayes factor

Maximum Likelihood

$$\left. \frac{d \ln L(\vec{p})}{d\vec{p}} \right|_{p_i=\hat{p}_i} = 0$$

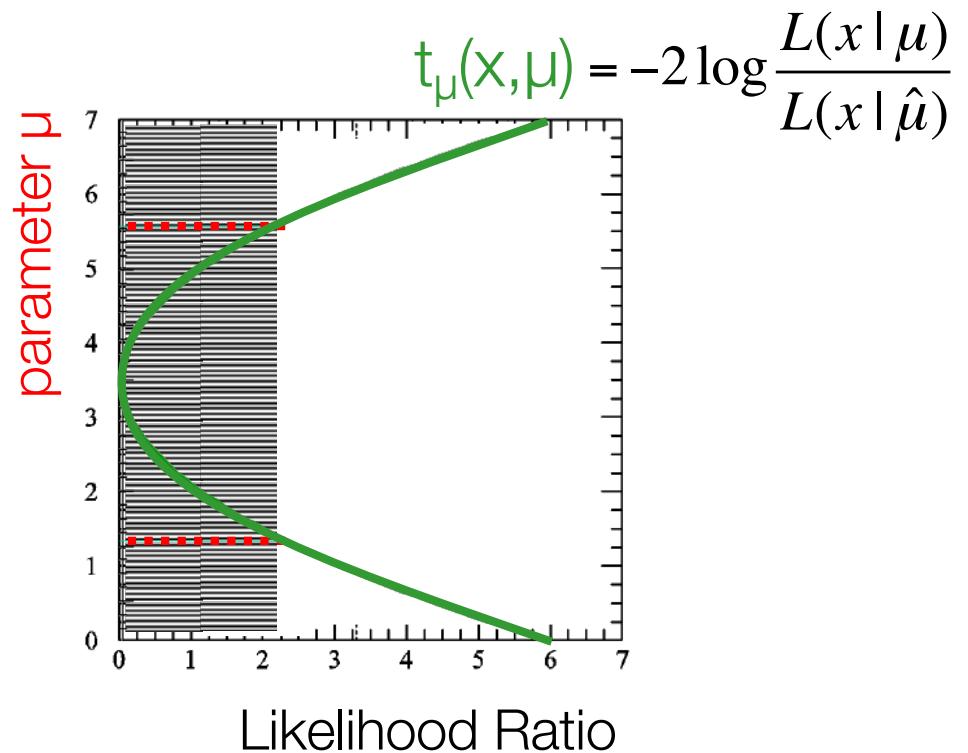
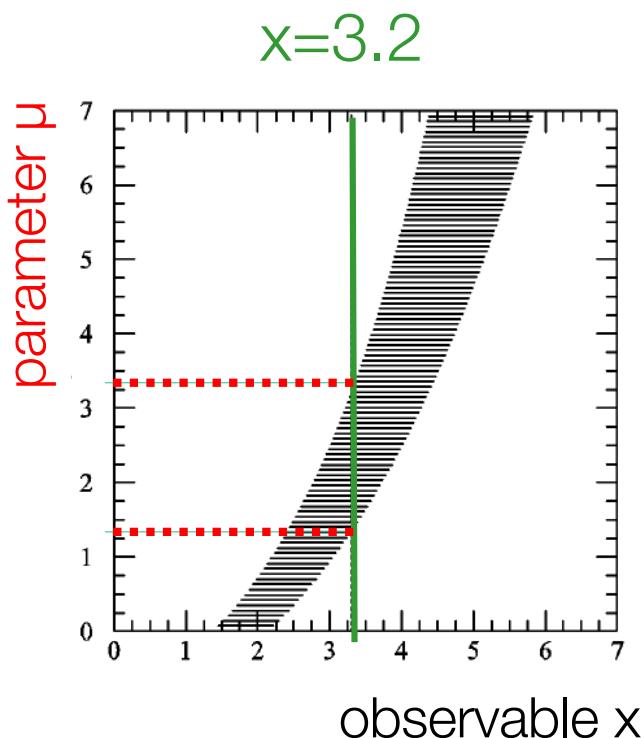


$s = x \pm y$

But fundamental techniques can be complicated to execute...

---

- Example of confidence interval calculation with Neyman construction
  - Need to construct ‘confidence belt’ using toy MC. Intersection observed data with belt defined interval in POI with guaranteed coverage



- Expensive, complicated procedure, but completely procedural once Likelihood and parameter of interest are fixed  
→ Can be wrapped in a tool that runs effectively ‘out-of-the-box’

# Running RooStats interval calculations ‘out-of-the-box’

- Confidence intervals calculated with model

- ‘Simple Fit’

```
RooAbsReal* nll = myModel->createNLL(data) ;  
RooMinuit m(*nll) ;  
m.migrad() ;  
m.hesse() ;
```

- Feldman Cousins (Frequentist Confidence Interval)

```
FeldmanCousins fc;  
fc.SetPdf(myModel);  
fc.SetData(data); fc.SetParameters(myPOU);  
fc.UseAdaptiveSampling(true);  
fc.FluctuateNumDataEntries(false);  
fc.SetNBins(100); // number of points to test per parameter  
fc.SetTestSize(.1);  
ConfInterval* fcint = fc.GetInterval();
```

- Bayesian (MCMC)

```
UniformProposal up;  
MCMCCalculator mc;  
mc.SetPdf(w::PC);  
mc.SetData(data); mc.SetParameters(s);  
mc.SetProposalFunction(up);  
mc.SetNumIters(100000); // steps in the chain  
mc.SetTestSize(.1); // 90% CL  
mc.SetNumBins(50); // used in posterior histogram  
mc.SetNumBurnInSteps(40);  
ConfInterval* mcmcint = mc.GetInterval();
```

## But you can also look ‘in the box’ and build your own

Tool to calculate p-values for a given hypothesis

$$\int_{q_{\mu, \text{obs}}}^{\infty} f(q_{\mu} | \mu') dq_{\mu}$$

```
// create first HypoTest calculator (N.B null is s+b model)
FrequentistCalculator fc(*data, *bModel, *sbModel);

// configure ToyMCSampler and set the test statistics
ToyMCSampler *toymcs = (ToyMCSampler*)fc.GetTestStatSampler();

ProfileLikelihoodTestStat profll(*sbModel->GetPdf());
// for CLs (bounded intervals) use one-sided profile likelihood
profll.SetOneSided(true);
toymcs->SetTestStatistic(&profll);

HypoTestInverter calc(*fc);
calc.UseCLs(true);

// configure and run the scan
calc.SetFixedScan(npoints, poimin, poimax);
HypoTestInverterResult * r = calc.GetInterval();

// get result and plot it
double upperLimit = r->UpperLimit();
double expectedLimit = r->GetExpectedUpperLimit(0);

HypoTestInverterPlot *plot = new HypoTestInverterPlot("hi","","",r);
plot->Draw();
```

$f(q_{\mu} | \mu')$

Tool to construct  
test statistic  
distribution

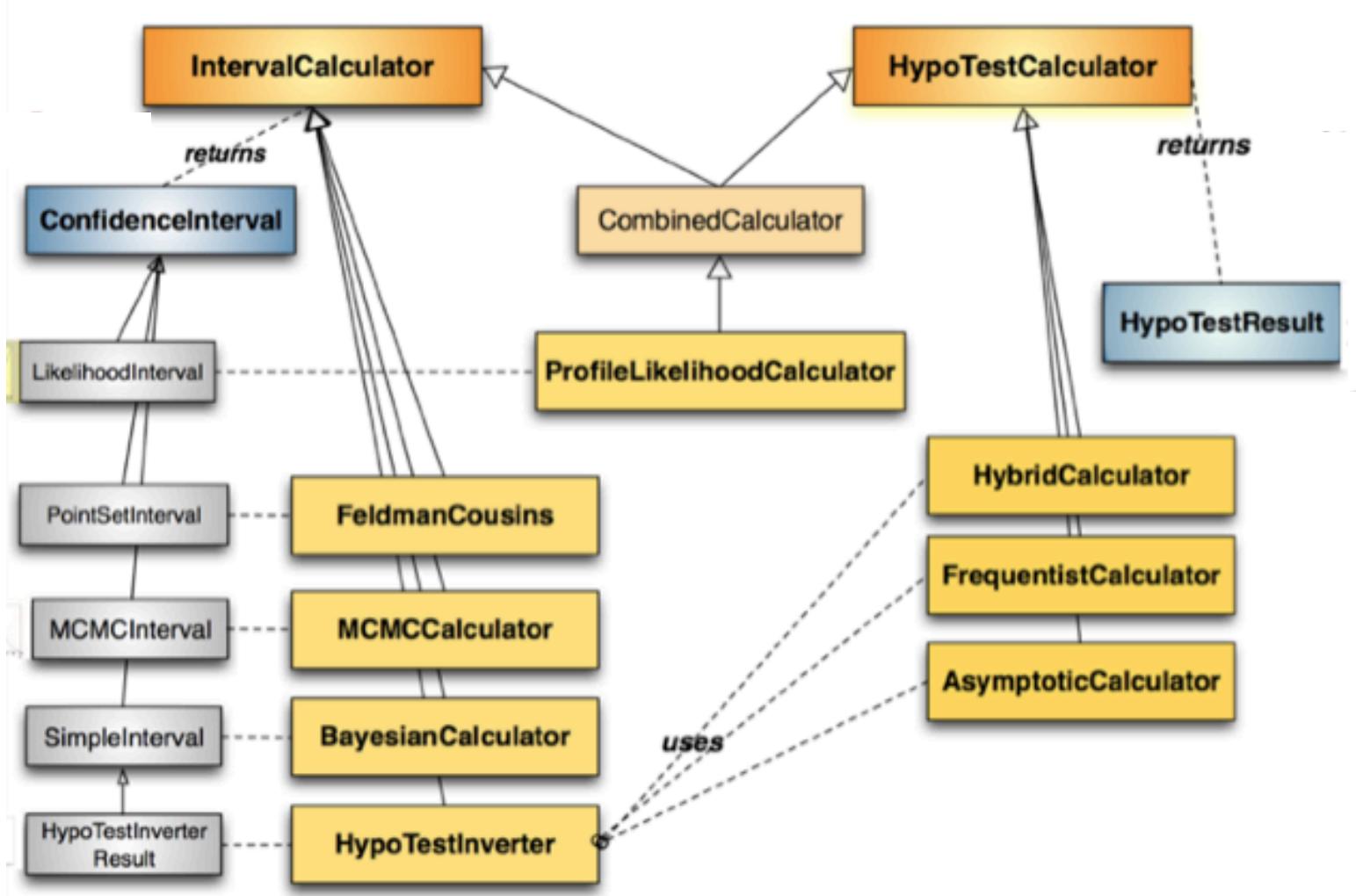
$q_{\mu}(\mu')$

The test statistic  
to be used for  
the calculation  
of p-values

Tool to construct  
interval from  
hypo test results

*Offset advanced control over details of statistical  
procedure (use of CLS, choice of test statistic, boundaries...)*

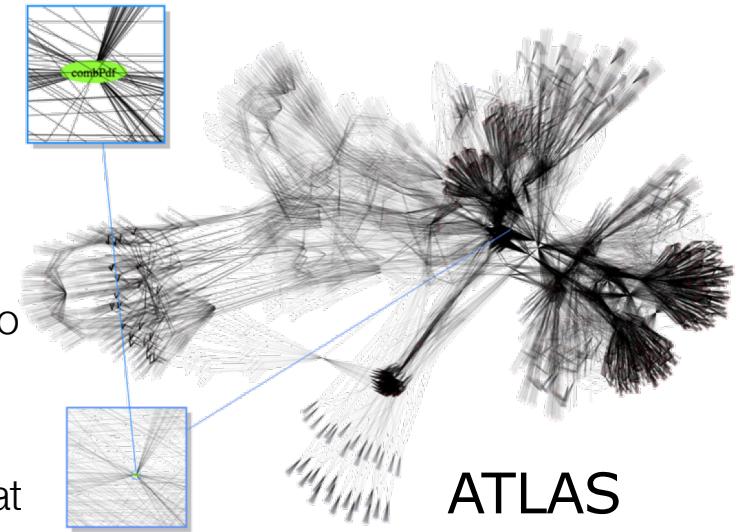
# RooStats class structure



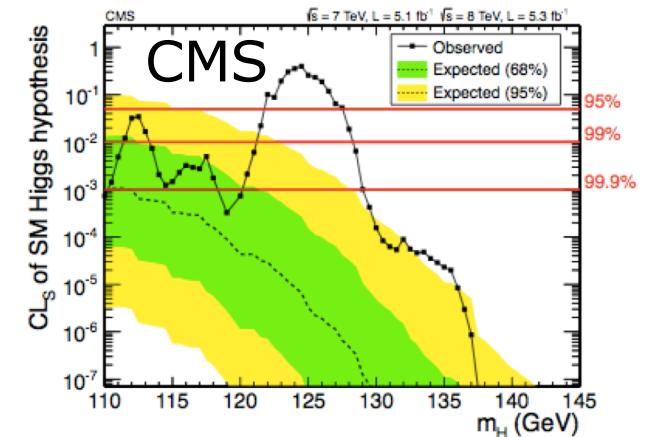
# Summary

---

- **RooFit** and **RooStats** allow you to perform advanced statistical data analysis
  - LHC Higgs results a prominent example
- **RooFit** provides (almost) limitless model building facilities
  - Concept of persistable model workspace allows to separate model building and model interpretation
  - **HistFactory** package introduces structured model building for binned likelihood template models that are common in LHC analyses
- Concept of RooFit **Workspace** has completely restructured HEP analysis workflow with ‘collaborative modeling’
- **RooStats** provide a wide set of statistical tests that can be performed on RooFit models
  - Bayesian, Frequentist and Likelihood-based test concepts



ATLAS



Wouter Verkerke, NIKHEF

## HistFitter

---

- RooFit/RooStats/Histfitter provide large amount of flexibility, modeling possibilities & statistical tools
- HistFitter is cased on RooFit/RooStats and provides steering for the full analysis chain for analyses that follow a certain pattern

- **Step 0** – Definition of signal/control/validation regions

- **Step 1** – Construct the likelihood function  $L(x|p)$

RooFit, or RooFit+HistFactory



RooWorkspace



- **Step 2** – Statistical tests on parameter of interest  $p$

RooStats

