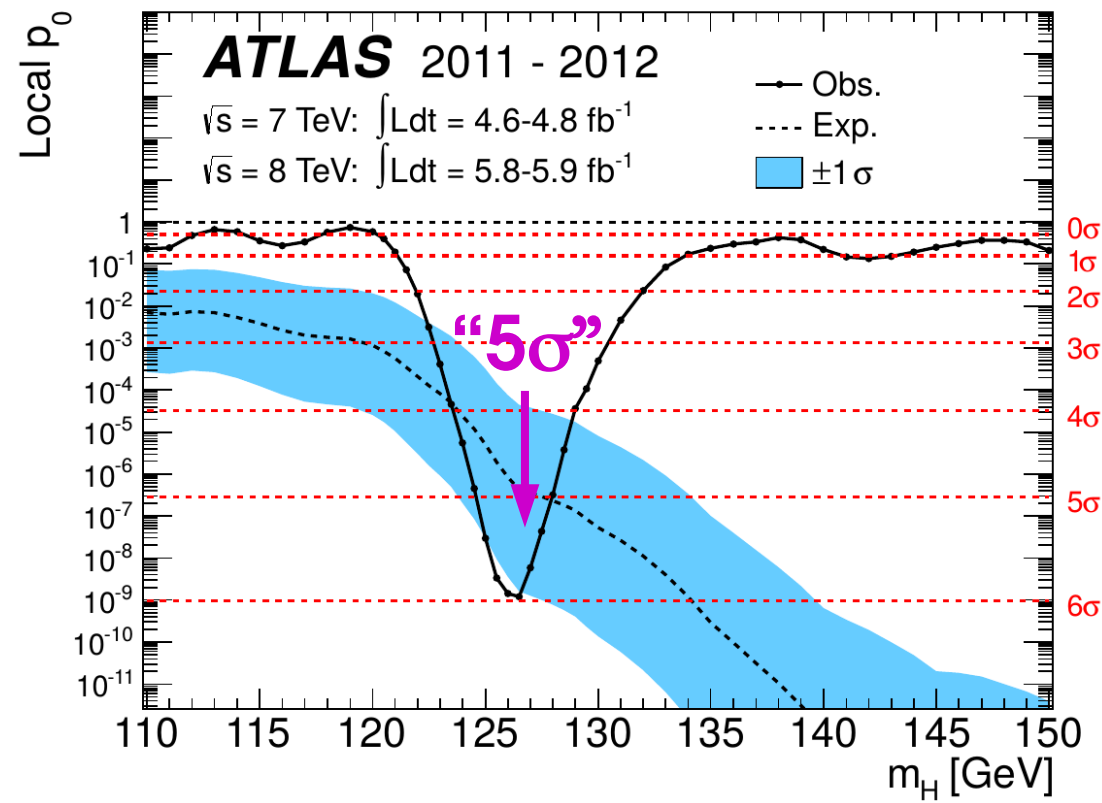
The background is a complex, abstract composition of overlapping geometric shapes in shades of yellow, green, and grey. A bright yellow starburst or explosion effect is centered in the middle. Scattered around this central point are several blue dice with white pips. The text is overlaid on this background in a large, bold, red font.

# **Statistical modeling and Systematic uncertainties in High Energy Physics**

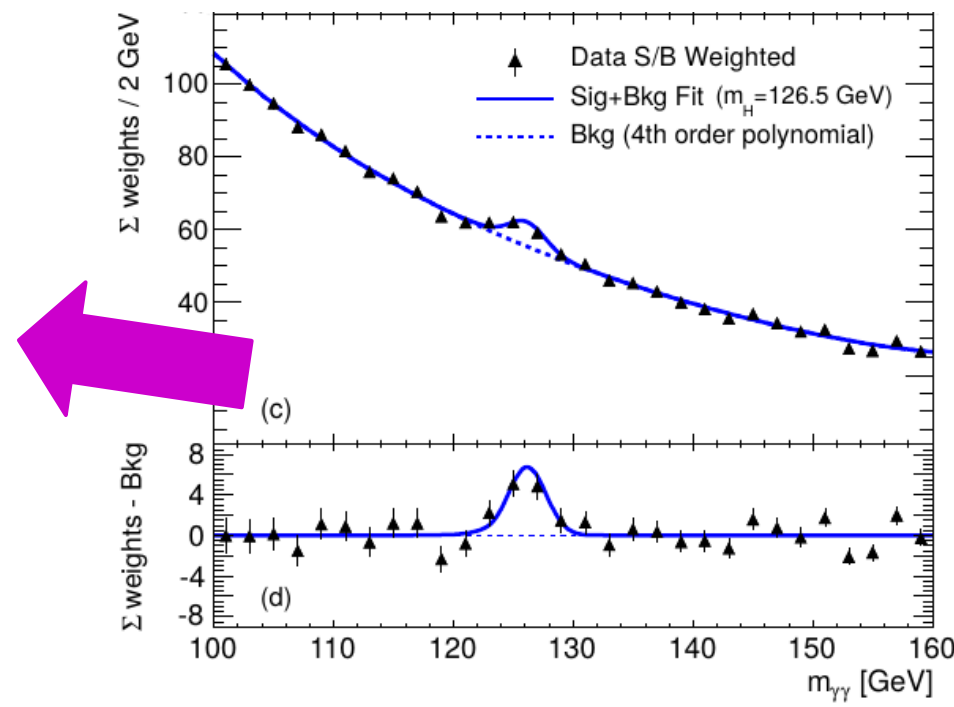
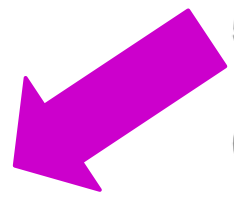
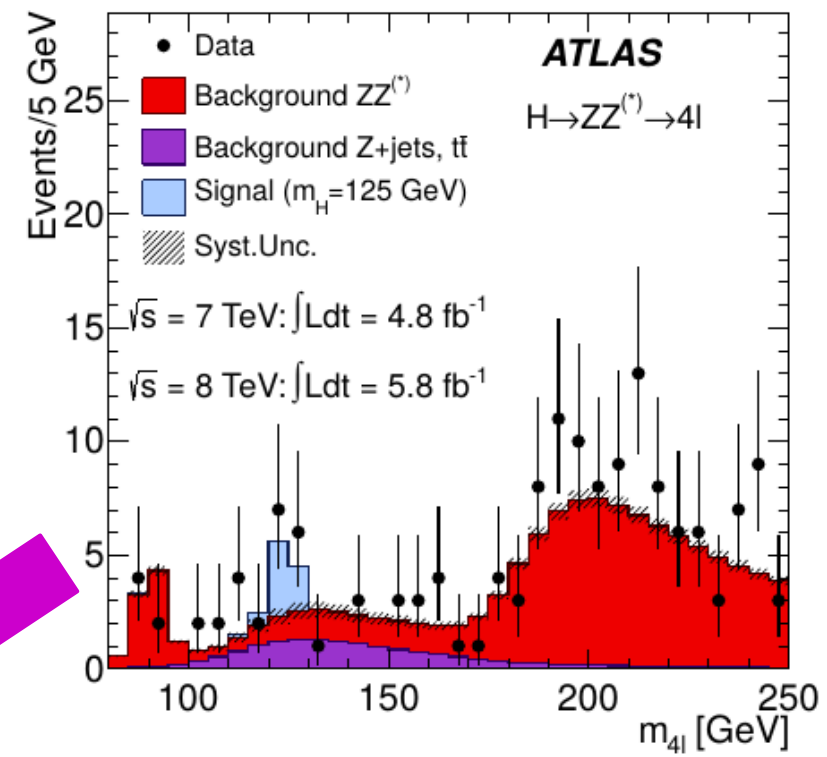
# Introduction

Statistical methods play a critical role in many areas of physics

Higgs discovery : **“We have 5 $\sigma$ ” !**

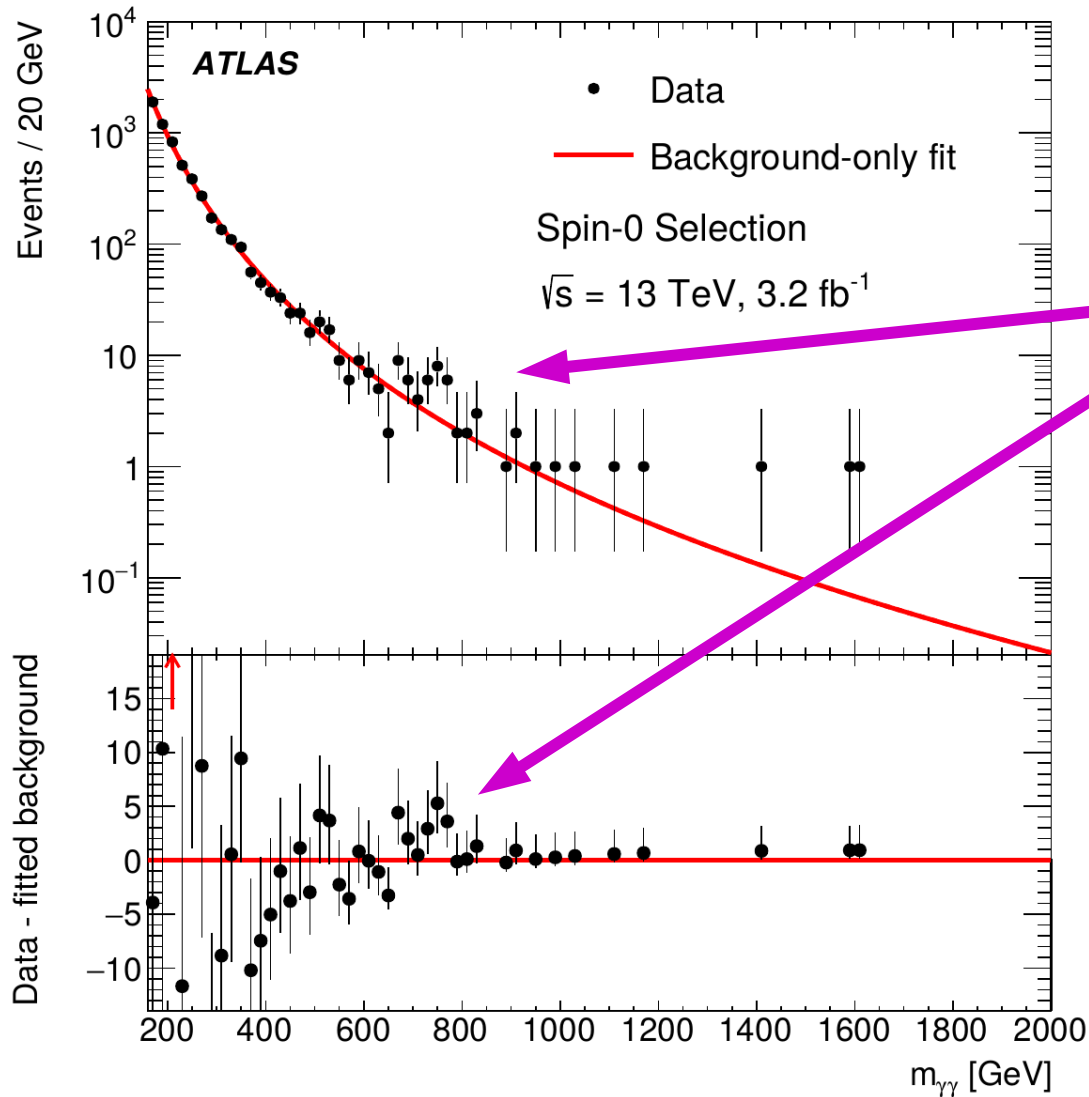


Phys. Lett. B 716 (2012) 1-29

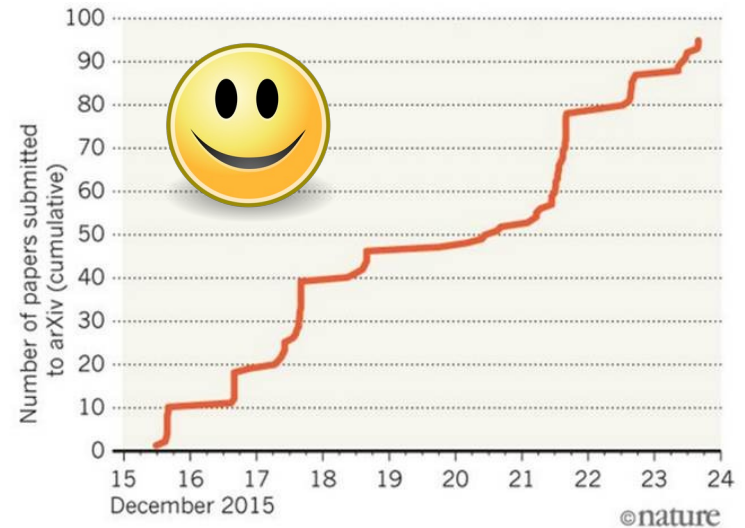


# Introduction

Sometimes difficult to distinguish a bona fide discovery  
from a **background fluctuation**...



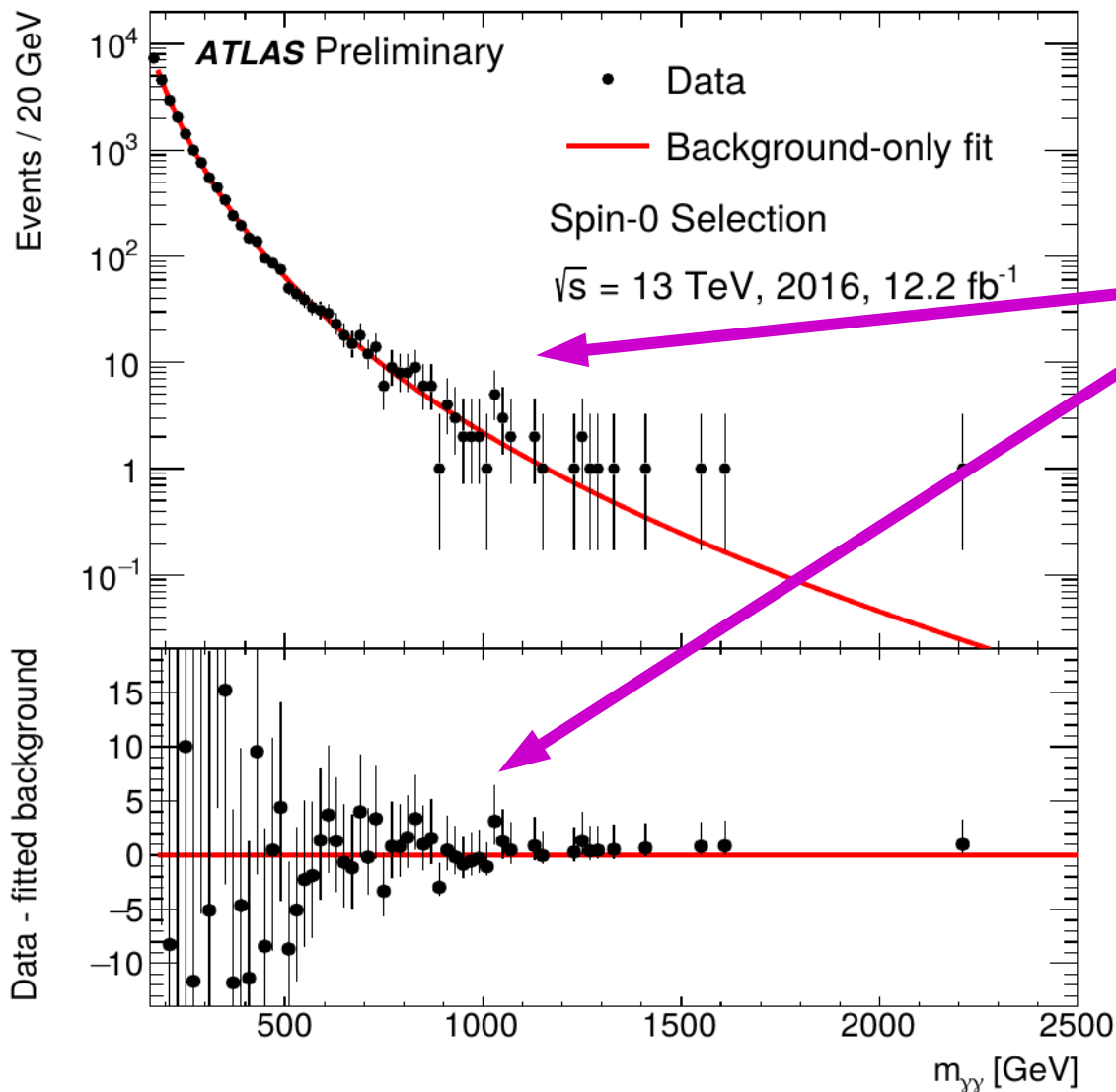
**New Physics ?**



JHEP 09 (2016) 1

# Introduction

Sometimes difficult to distinguish a bona fide discovery  
from a **background fluctuation**...



*A few months later...*

~~New Physics ?~~

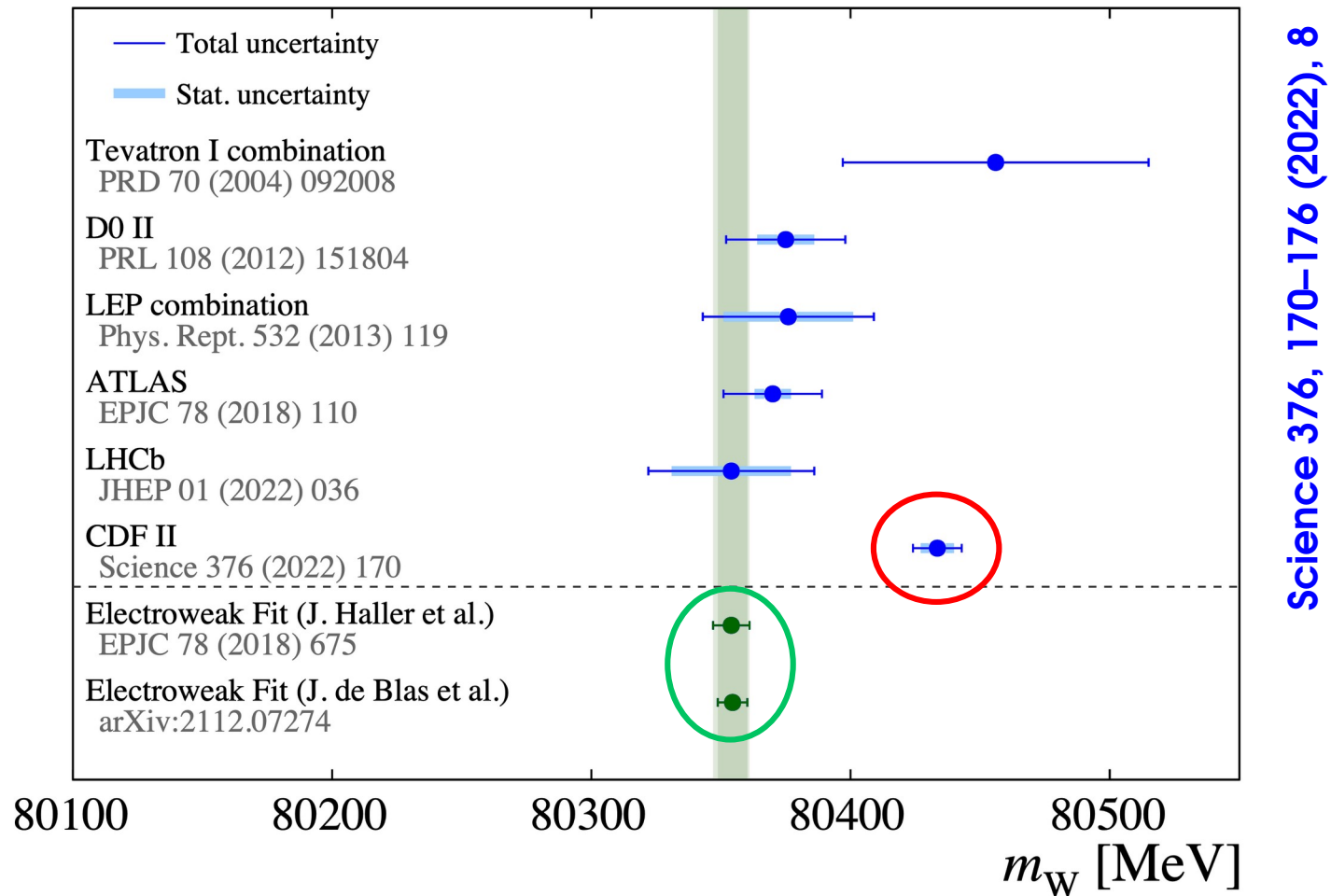


JHEP 09 (2016) 1



# Uncertainties

Many important questions answered by **precision measurements**,  
**Key point** = determination of **uncertainties**

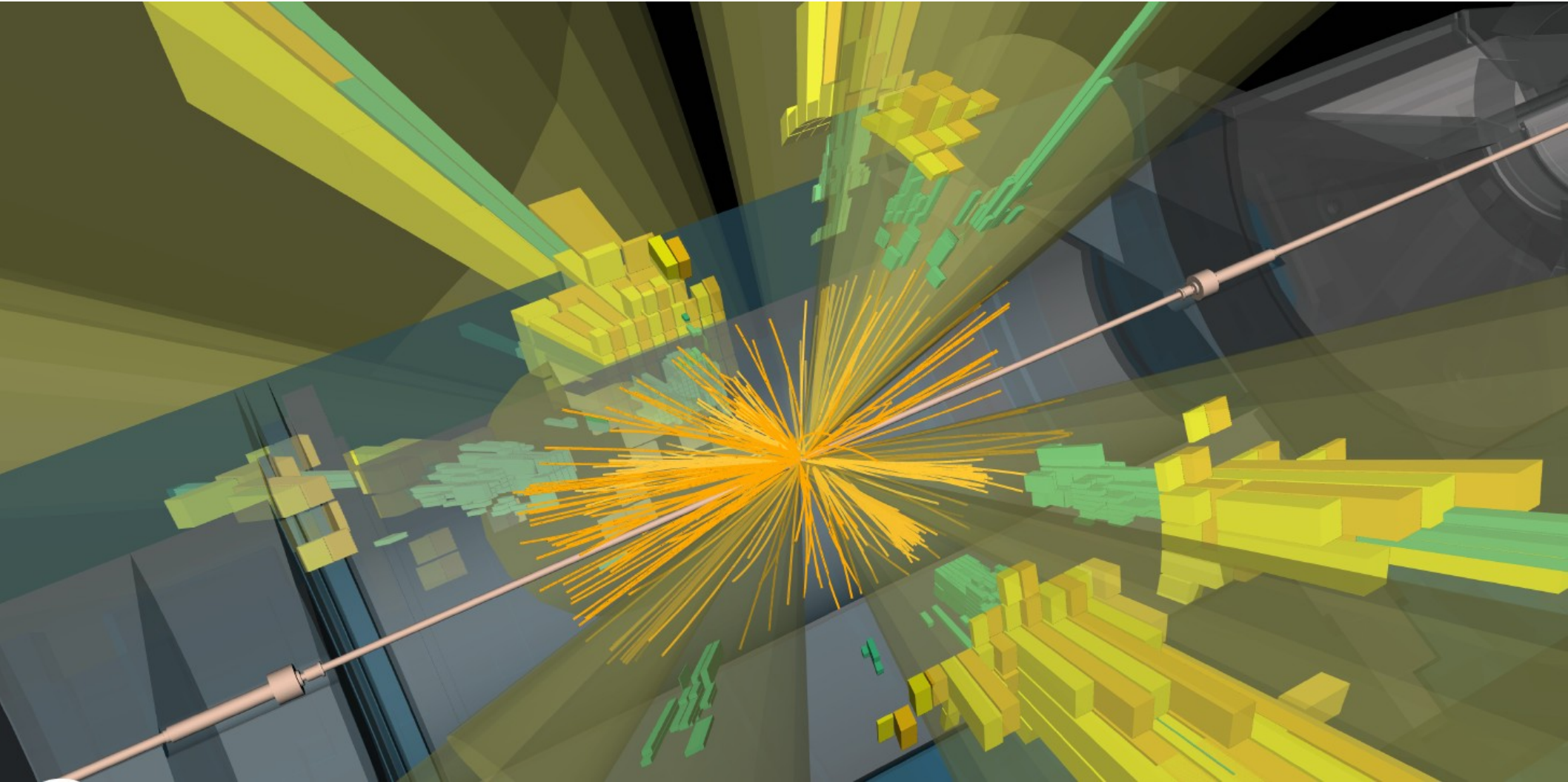


$$M_W = 80,433.5 \pm 6.4_{\text{stat}} \pm 6.9_{\text{syst}} = 80,433.5 \pm 9.4 \text{ MeV}/c^2$$

# Randomness in High-Energy Physics

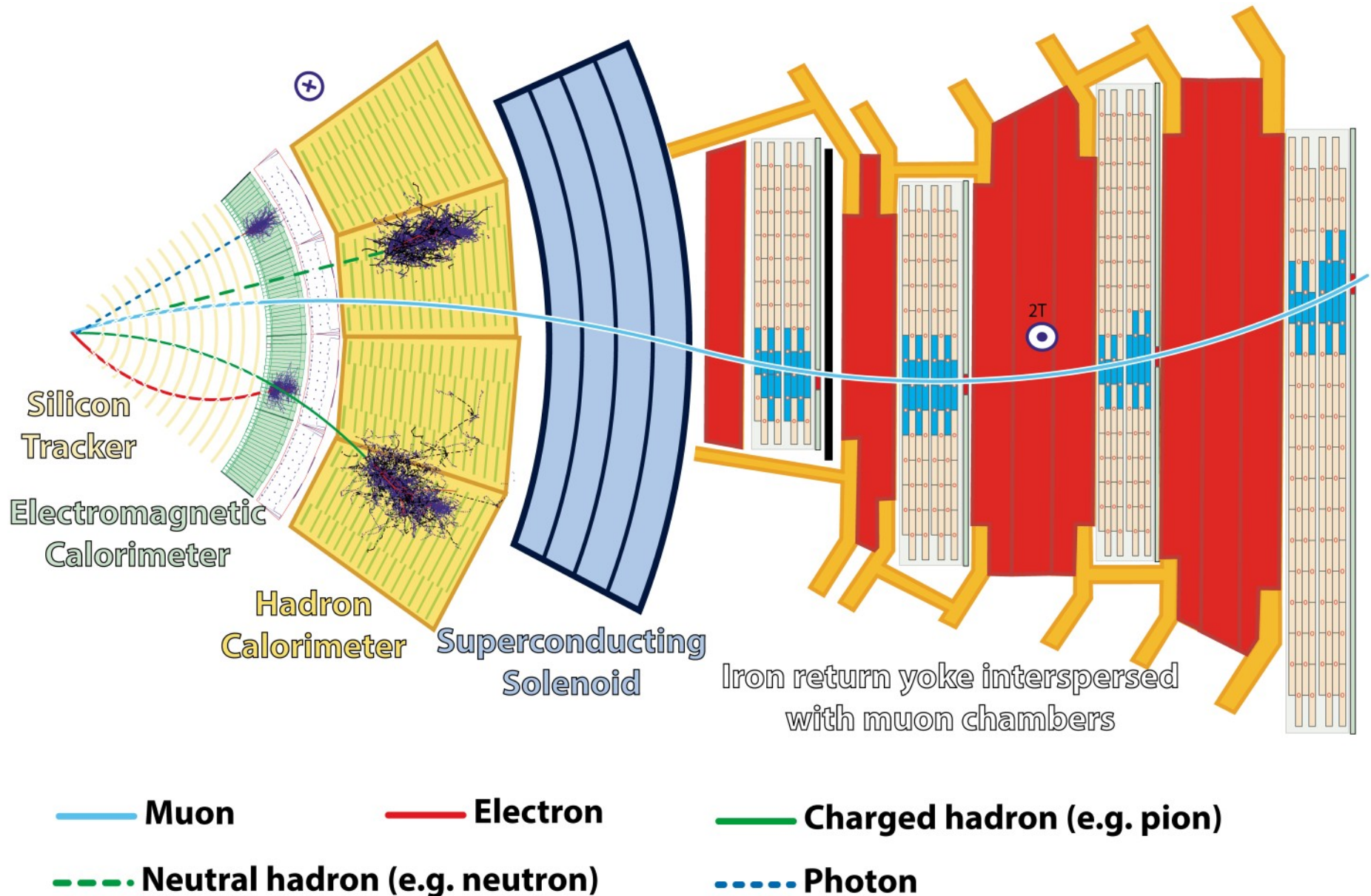
---

Experimental data is produced by incredibly complex processes



# Randomness in High-Energy Physics

Experimental data is produced by incredibly complex processes





# Randomness in High-Energy Physics

Experimental data is produced by incredibly complex processes

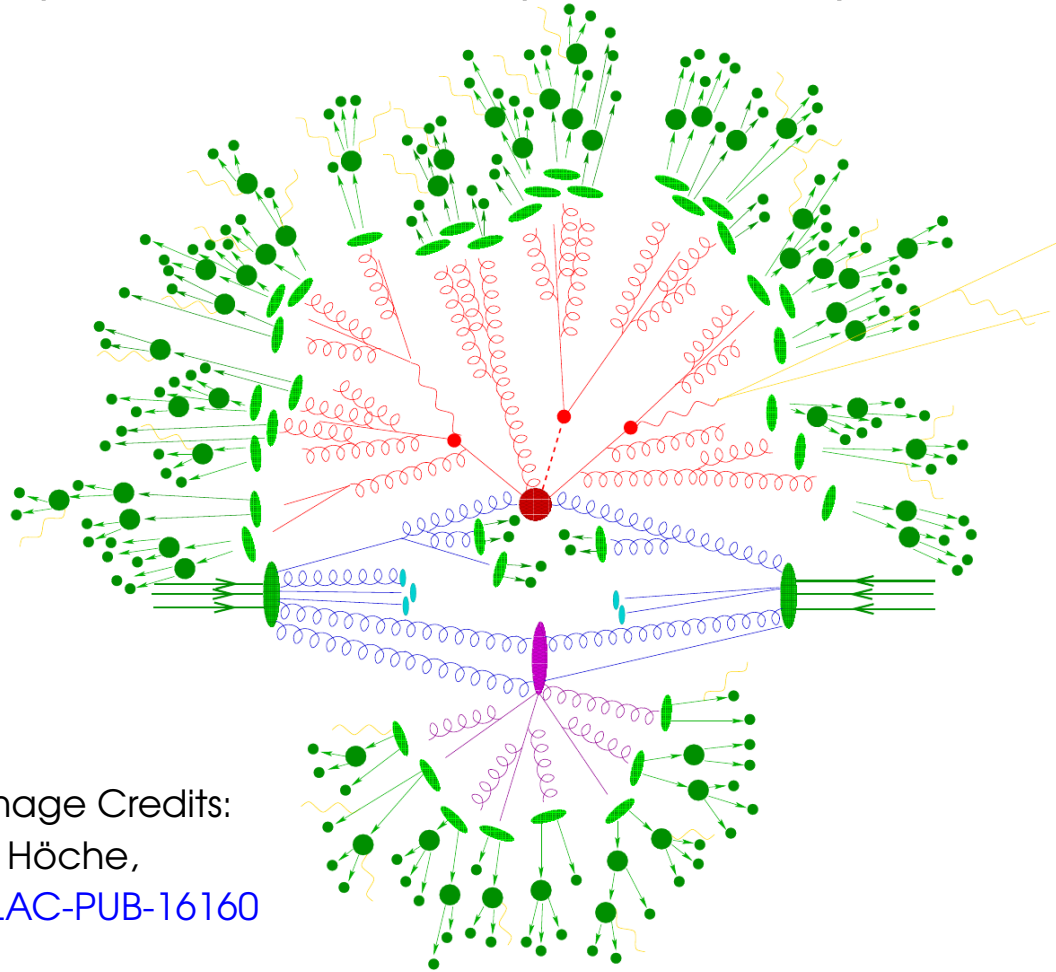


Image Credits:  
S. Höche,  
[SLAC-PUB-16160](#)

**Randomness** involved in all stages

→ **Classical** randomness: detector response

→ **Quantum** effects in particle production, decay

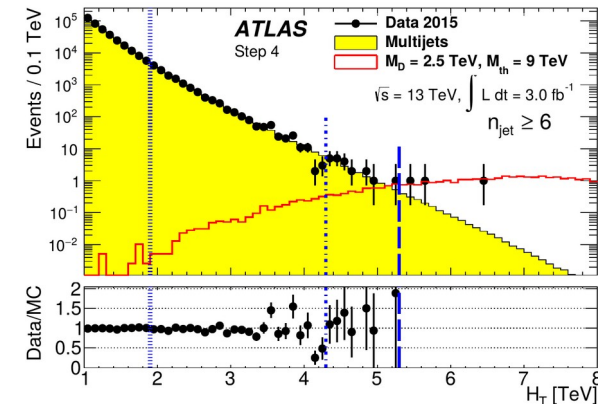
Hard scattering

PDFs, Parton shower, Pileup

Decays

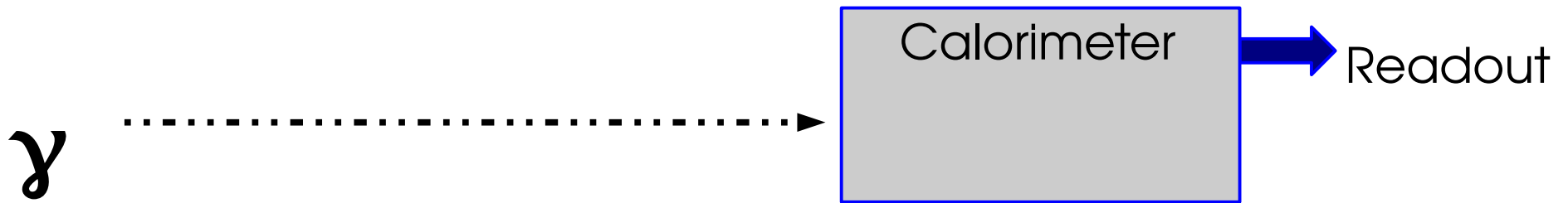
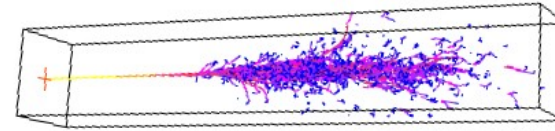
Detector response

Reconstruction



# Measurement Errors: Energy measurement

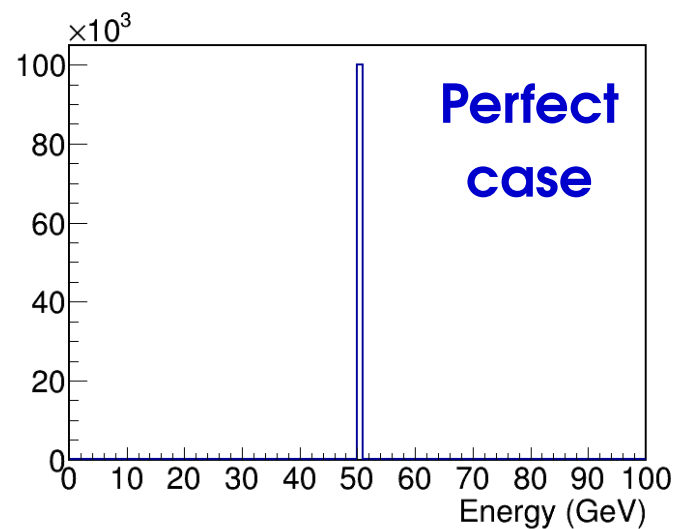
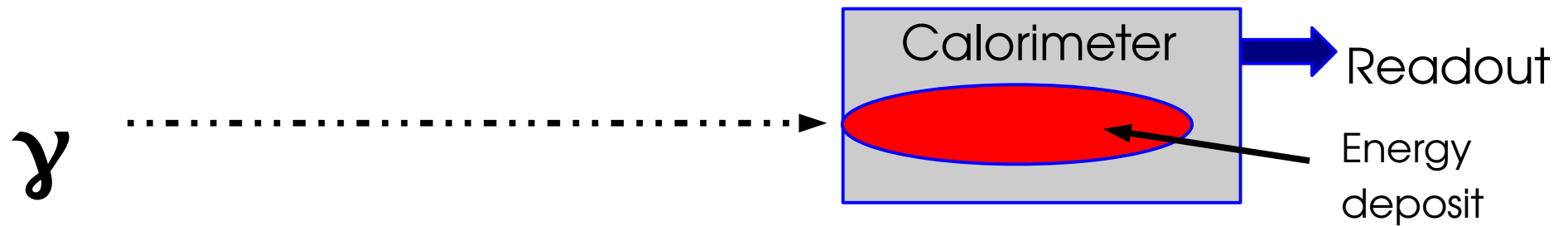
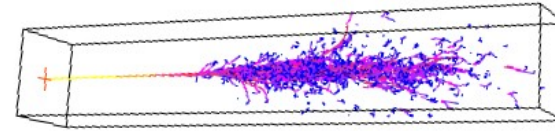
**Example:** measuring the energy of a photon in a calorimeter





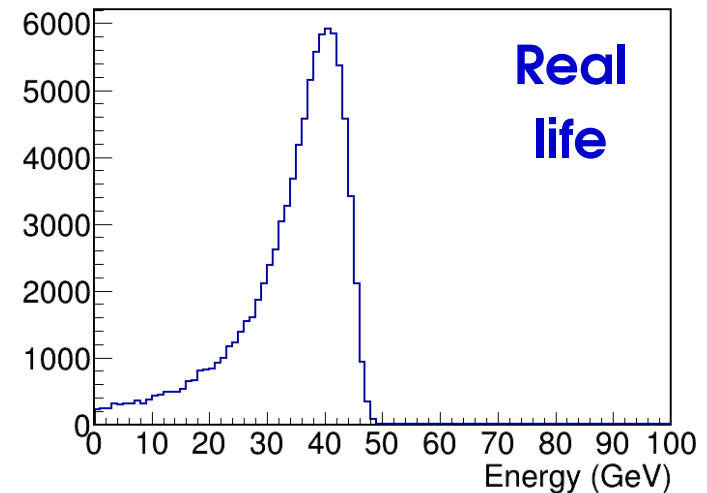
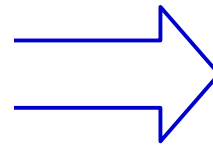
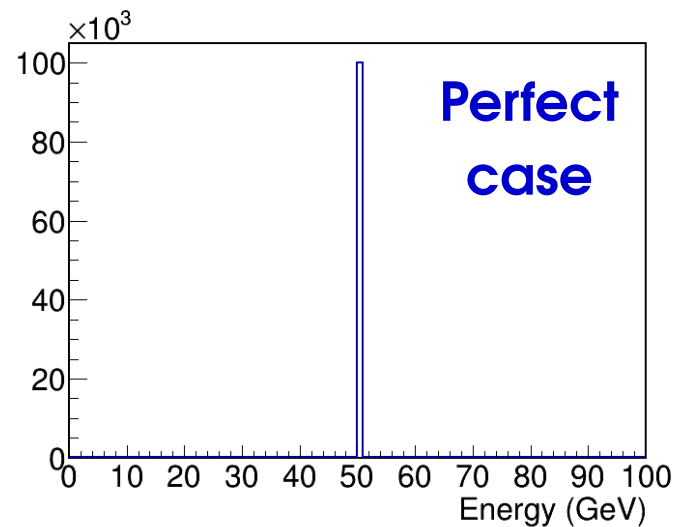
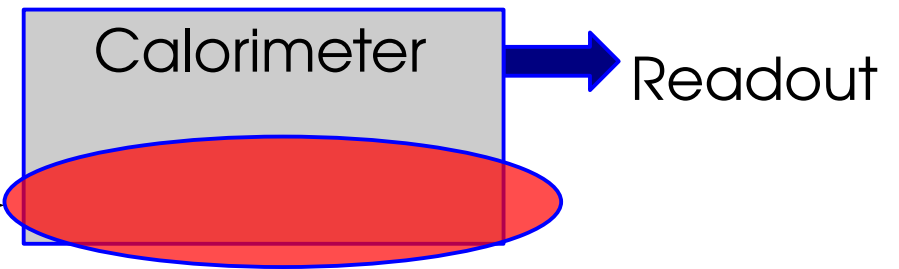
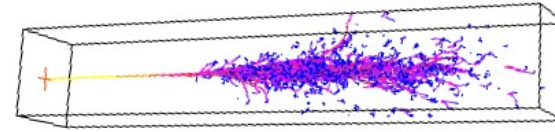
# Measurement Errors: Energy measurement

**Example:** measuring the energy  
of a photon in a calorimeter



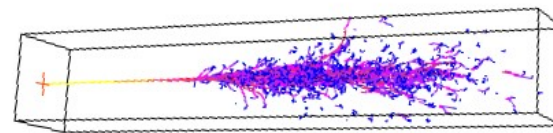
# Measurement Errors: Energy measurement

**Example:** measuring the energy  
of a photon in a calorimeter



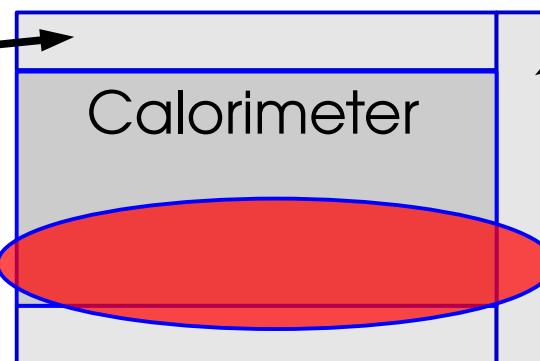
# Measurement Errors: Energy measurement

**Example:** measuring the energy  
of a photon in a calorimeter



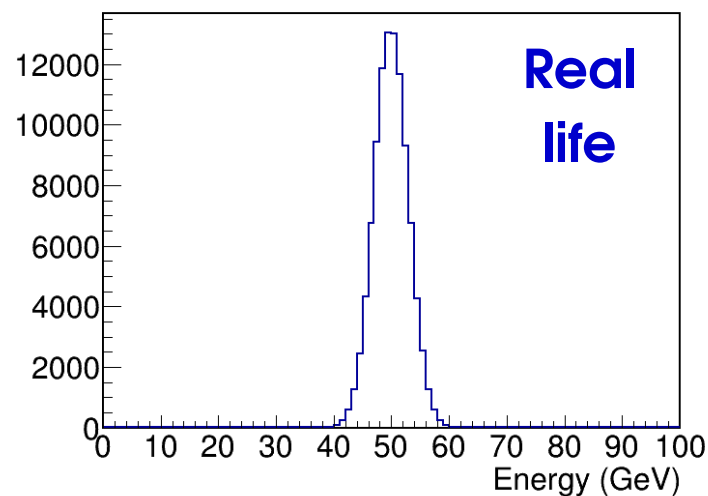
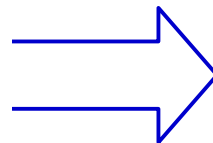
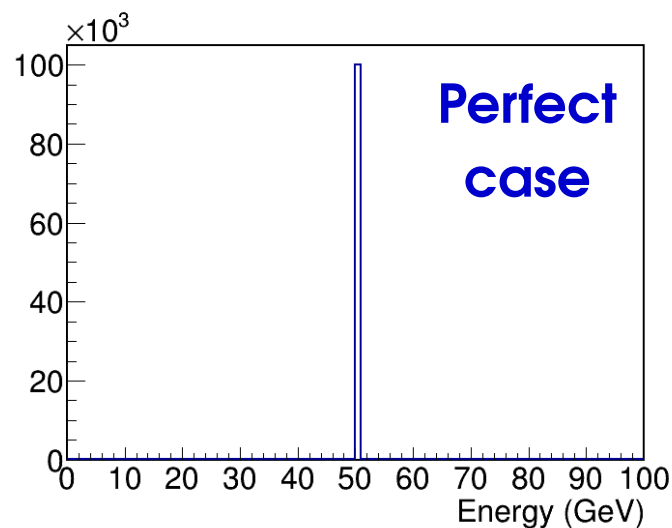
Measure leakage behind calorimeter

Measure leakage  
into neighboring cells



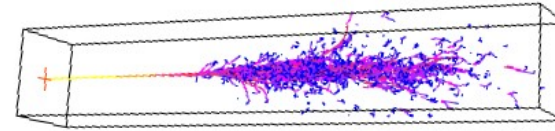
Readout

$\gamma$



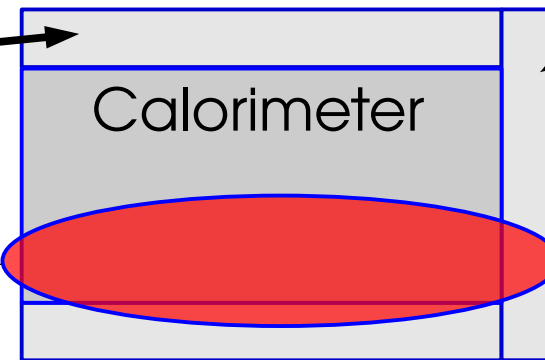
# Measurement Errors: Energy measurement

**Example:** measuring the energy  
of a photon in a calorimeter



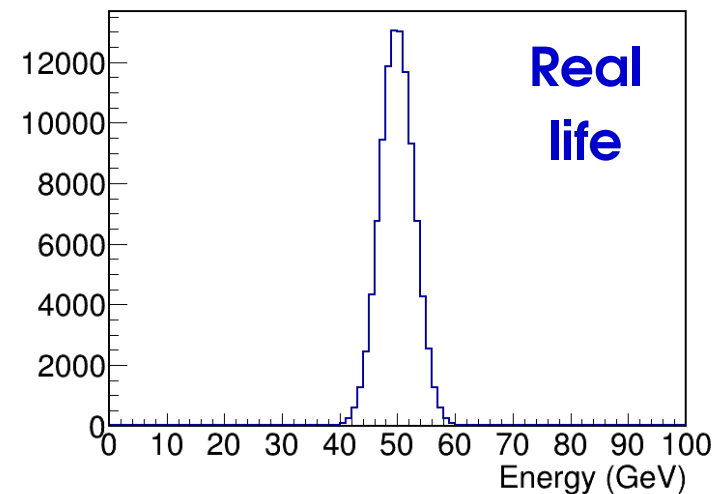
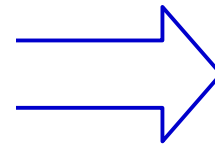
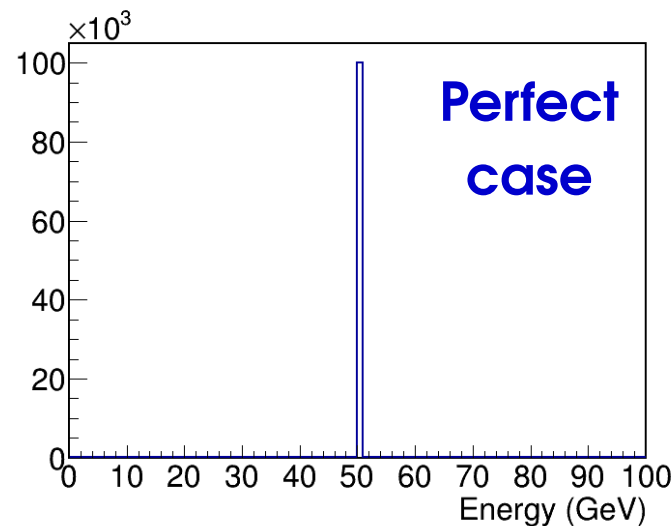
Measure leakage behind calorimeter

Measure leakage  
into neighboring cells



Readout

$\gamma$

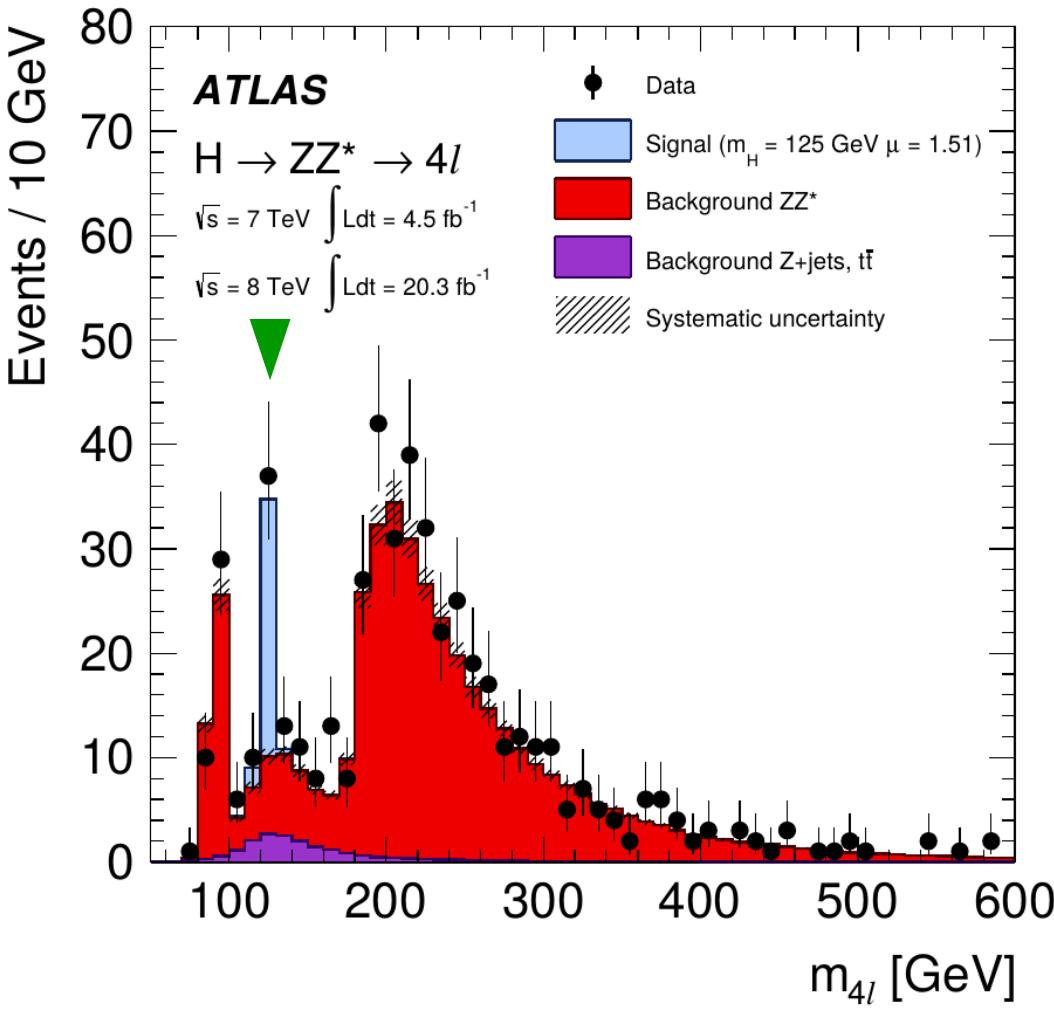


Cannot predict the measured value for a given event

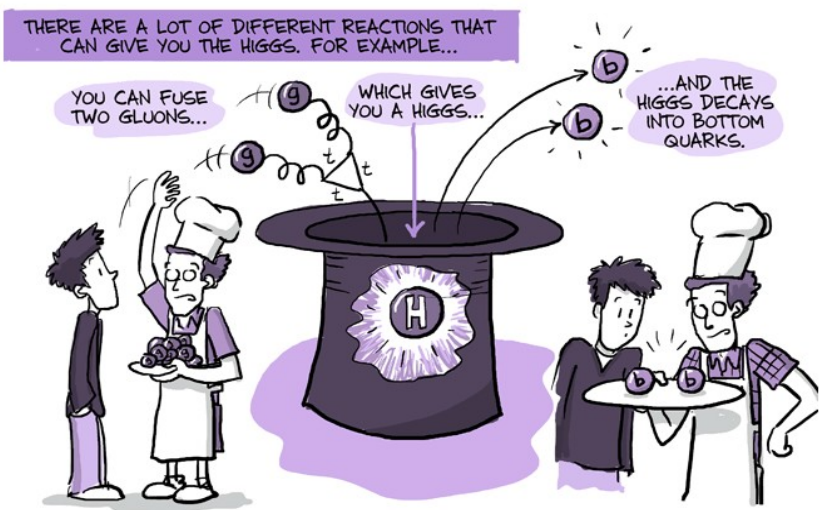
⇒ **Random process** ⇒ **Need a probabilistic description**

# Quantum Randomness: $H \rightarrow ZZ^* \rightarrow 4l$

Phys. Rev. D **91**, 012006



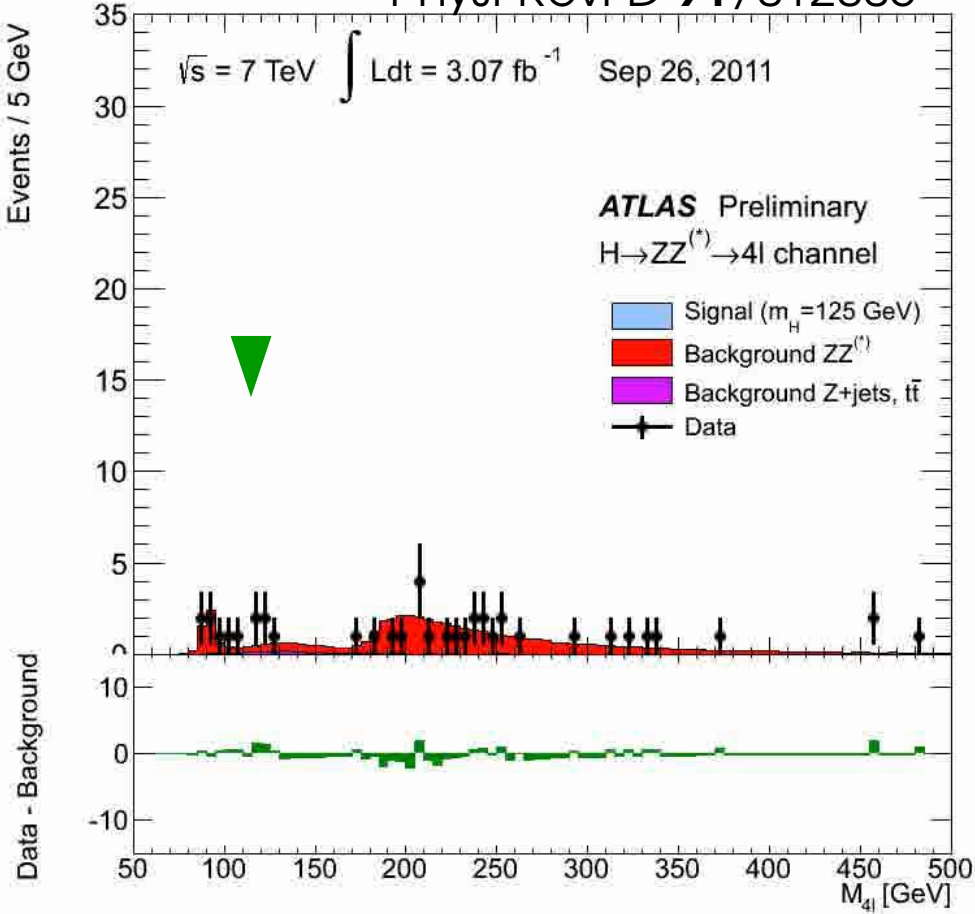
**Rare process:** Expect 1 signal event every **~6 days**



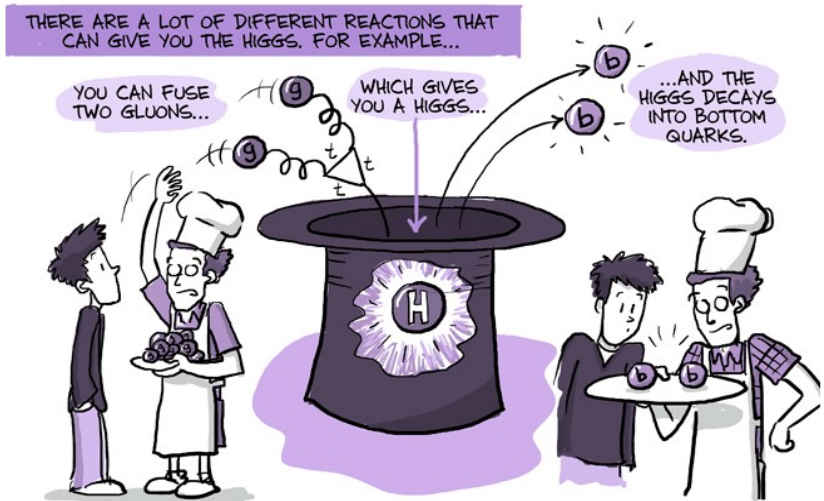


# Quantum Randomness: $H \rightarrow ZZ^* \rightarrow 4l$

Phys. Rev. D **91**, 012006



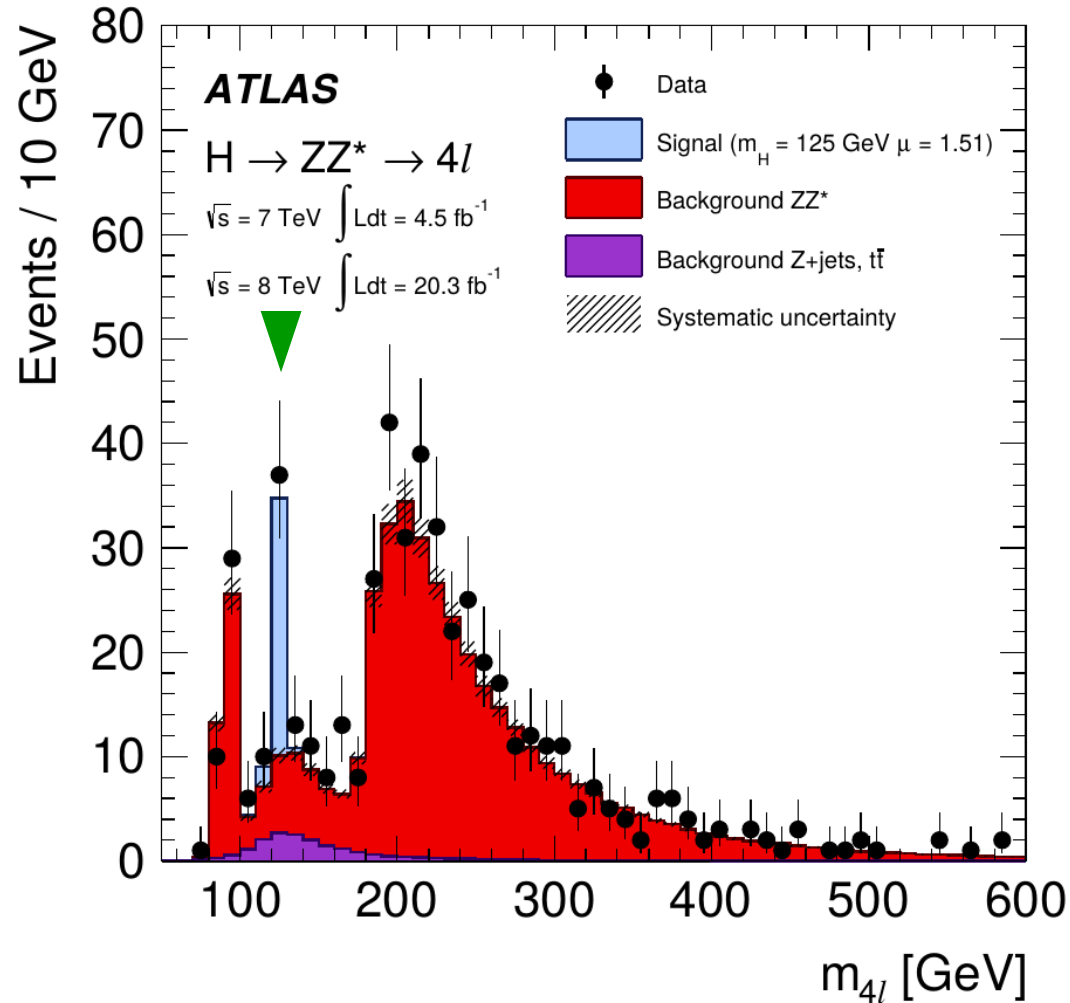
**Rare process:** Expect 1 signal event every **~6 days**



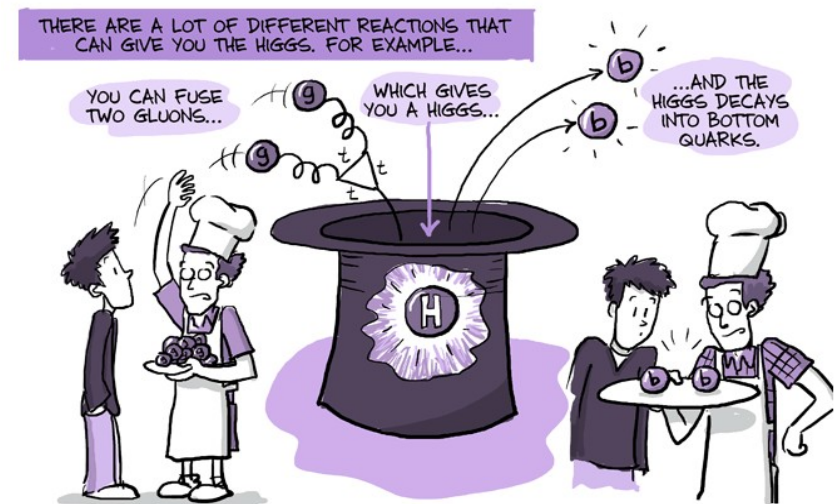
[View online](#)

# Quantum Randomness: $H \rightarrow ZZ^* \rightarrow 4l$

Phys. Rev. D **91**, 012006



**Rare process:** Expect 1 signal event every **~6 days**



“Will I get an event today ?” → only **probabilistic** answer

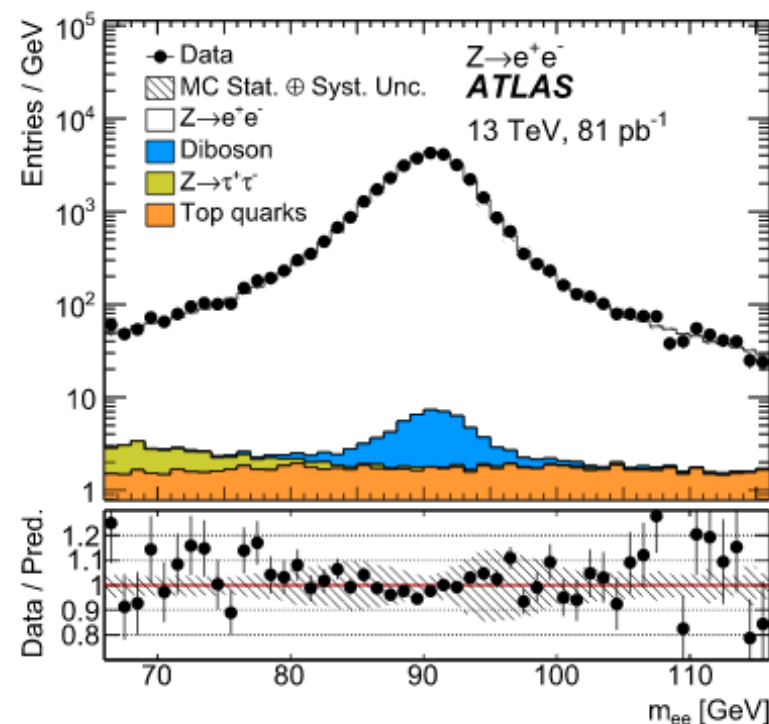
# Performing a measurement

Phys. Lett. B 759 (2016) 601

Measure the cross-section (event rate) of the  $Z \rightarrow ee$  process

$$\sigma_{fid} = \frac{n_{data} - N_{bkg}}{C_{fid} L}$$

$35000 \pm 187$  (points to  $n_{data}$ )  
 $175 \pm 8$  (points to  $N_{bkg}$ )  
 $(81 \pm 2) \text{ pb}^{-1}$  (points to  $L$ )  
 $0.552 \pm 0.006$  (points to  $C_{fid}$ )



$$\sigma^{fid} = 0.781 \pm 0.004 \text{ (stat)} \pm 0.018 \text{ (syst) nb}$$

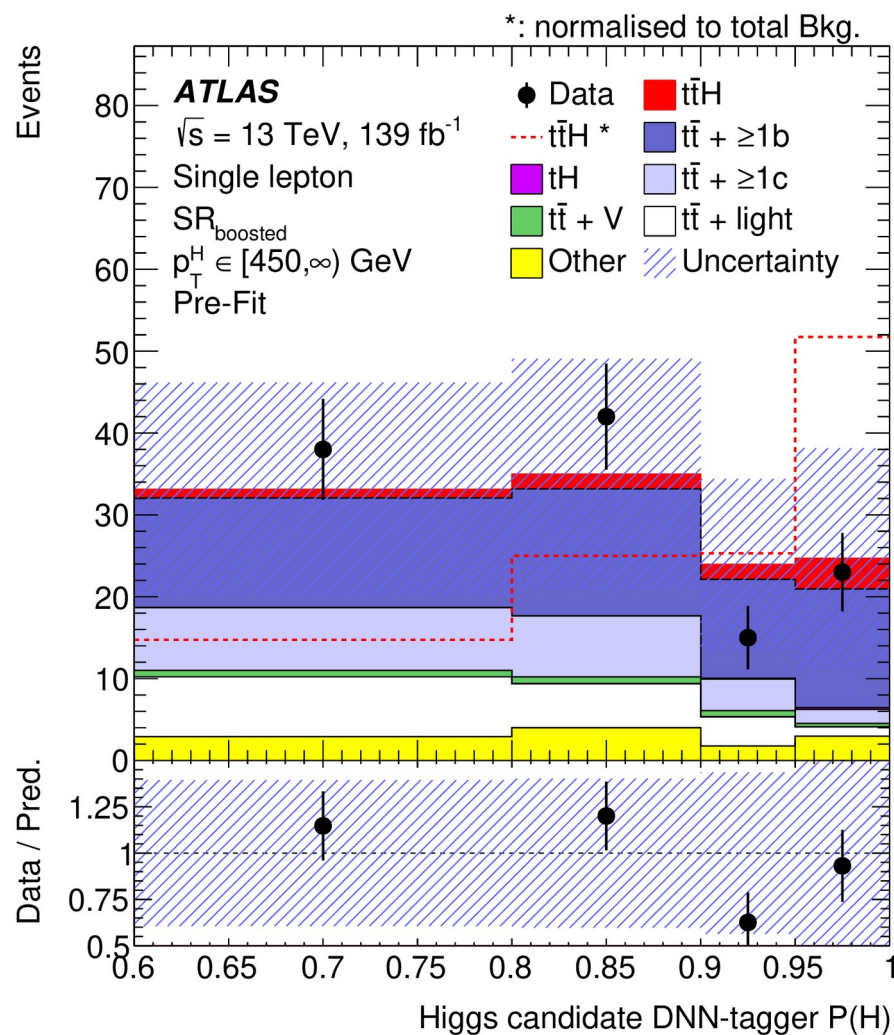
Fluctuations in the data counts

Other uncertainties (assumptions, parameter values)

“Single bin counting” : only data input is  $N_{data}$

# Example 2: $t\bar{t}H \rightarrow b\bar{b}$

arXiv:2111.06712



Event counting in different regions:

**Multiple-bin counting**

**Lots of information available**

→ Potentially higher sensitivity

→ How to make optimal use of it ?

---

# HEP Statistical Modeling



# How to count

Collider processes: produce (many) events  $N$ , select a (very) small fraction  $P$

→ In principle, binomial process

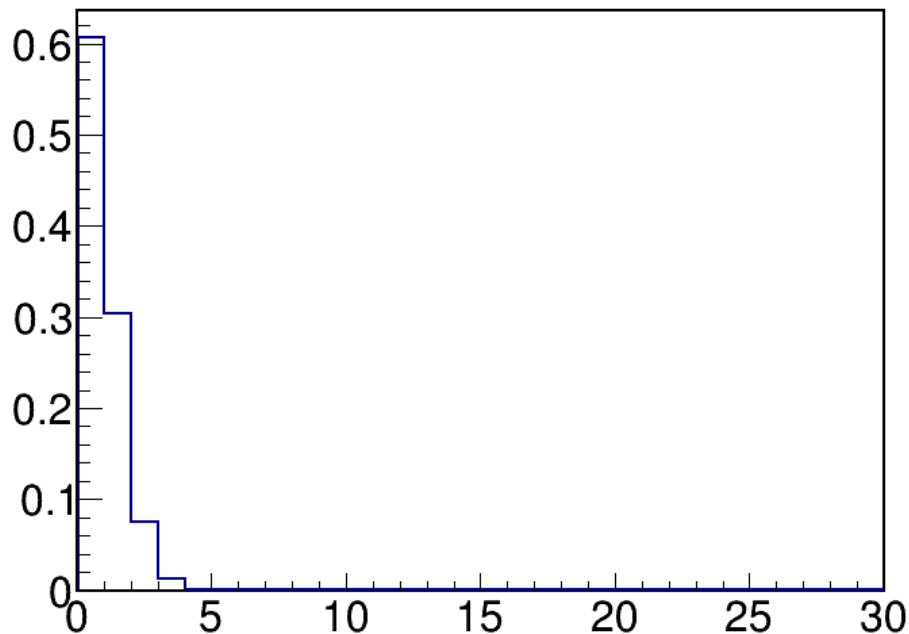
→ In practice,  $P \ll 1$ ,  $N \gg 1$ ,  $\Rightarrow$  Poisson approximation.

→ i.e. **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution:**  $P(n; \lambda) = e^{-\lambda} \frac{\lambda^n}{n!}$

$\lambda = 0.5$

$(1-P)^{N-n} \stackrel{n \ll N}{\approx} \left(1 - \frac{\lambda}{N}\right)^N \stackrel{N \gg 1}{\approx} e^{-\lambda}$



**Mean** =  $\lambda$

**Variance** =  $\lambda$

$\sigma = \sqrt{\lambda}$

Central limit theorem :

becomes **Gaussian for large  $\lambda$**  :

$$P(\lambda) \xrightarrow{\lambda \rightarrow \infty} G(\lambda, \sqrt{\lambda})$$

# How to count

Collider processes: produce (many) events  $N$ , select a (very) small fraction  $P$

→ In principle, binomial process

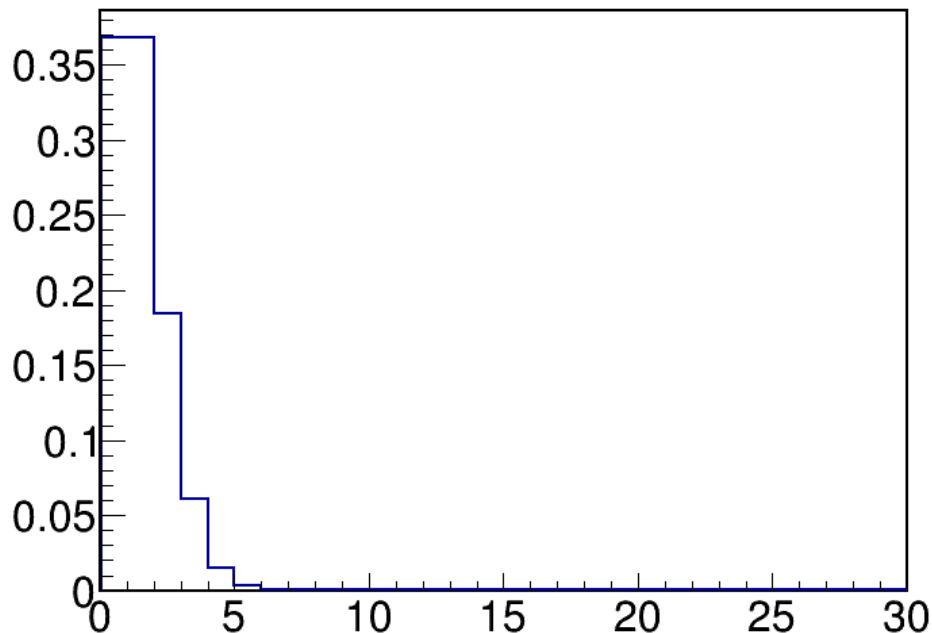
→ In practice,  $P \ll 1$ ,  $N \gg 1$ ,  $\Rightarrow$  Poisson approximation.

→ i.e. **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution:**  $P(n; \lambda) = e^{-\lambda} \frac{\lambda^n}{n!}$

$\lambda = 1$

$(1-P)^{N-n} \stackrel{n \ll N}{\approx} \left(1 - \frac{\lambda}{N}\right)^N \stackrel{N \gg 1}{\approx} e^{-\lambda}$



**Mean** =  $\lambda$

**Variance** =  $\lambda$

$\sigma = \sqrt{\lambda}$

Central limit theorem :

becomes **Gaussian for large  $\lambda$**  :

$$P(\lambda) \xrightarrow{\lambda \rightarrow \infty} G(\lambda, \sqrt{\lambda})$$

# How to count

Collider processes: produce (many) events  $N$ , select a (very) small fraction  $P$

→ In principle, binomial process

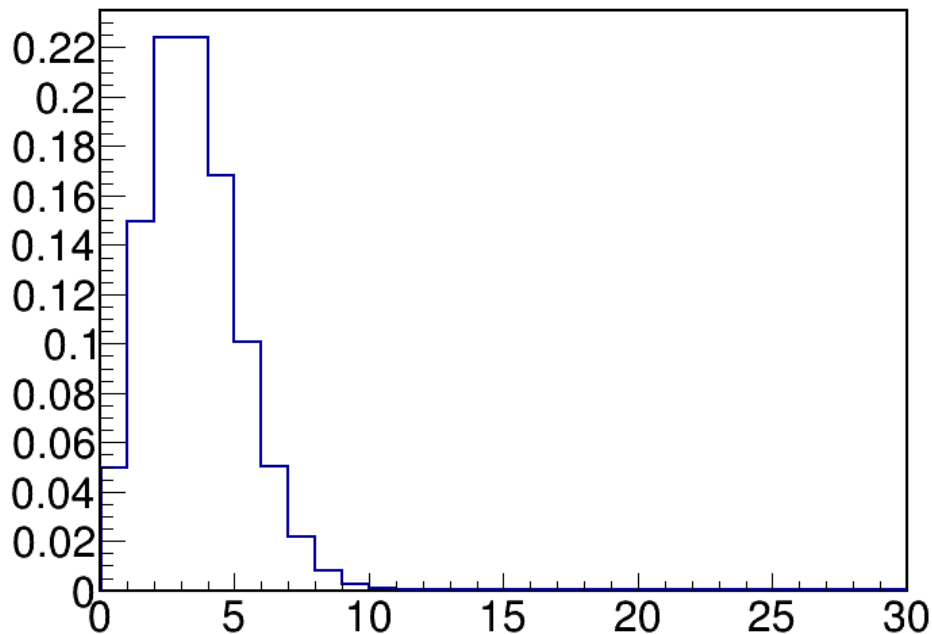
→ In practice,  $P \ll 1$ ,  $N \gg 1$ ,  $\Rightarrow$  Poisson approximation.

→ i.e. **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution:**  $P(n; \lambda) = e^{-\lambda} \frac{\lambda^n}{n!}$

$\lambda = 3$

$(1-P)^{N-n} \stackrel{n \ll N}{\approx} \left(1 - \frac{\lambda}{N}\right)^N \stackrel{N \gg 1}{\approx} e^{-\lambda}$



**Mean** =  $\lambda$

**Variance** =  $\lambda$

$\sigma = \sqrt{\lambda}$

Central limit theorem :

becomes **Gaussian for large  $\lambda$**  :

$$P(\lambda) \xrightarrow{\lambda \rightarrow \infty} G(\lambda, \sqrt{\lambda})$$

# How to count

Collider processes: produce (many) events  $N$ , select a (very) small fraction  $P$

→ In principle, binomial process

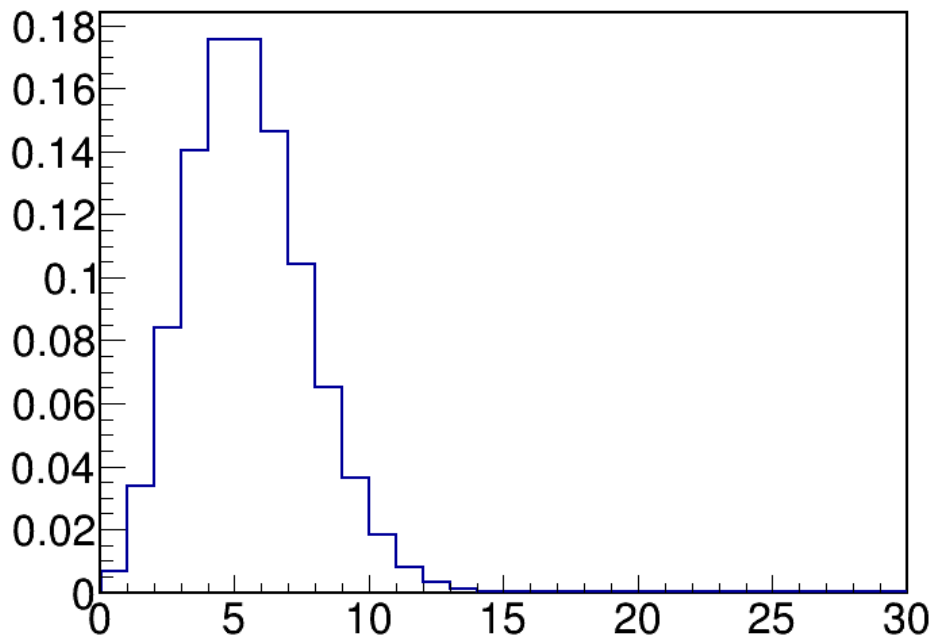
→ In practice,  $P \ll 1$ ,  $N \gg 1$ ,  $\Rightarrow$  Poisson approximation.

→ i.e. **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution:**  $P(n; \lambda) = e^{-\lambda} \frac{\lambda^n}{n!}$

$\lambda = 5$

$(1-P)^{N-n} \stackrel{n \ll N}{\approx} \left(1 - \frac{\lambda}{N}\right)^N \stackrel{N \gg 1}{\approx} e^{-\lambda}$



**Mean** =  $\lambda$

**Variance** =  $\lambda$

$\sigma = \sqrt{\lambda}$

Central limit theorem :

becomes **Gaussian for large  $\lambda$**  :

$$P(\lambda) \xrightarrow{\lambda \rightarrow \infty} G(\lambda, \sqrt{\lambda})$$

# How to count

Collider processes: produce (many) events  $N$ , select a (very) small fraction  $P$

→ In principle, binomial process

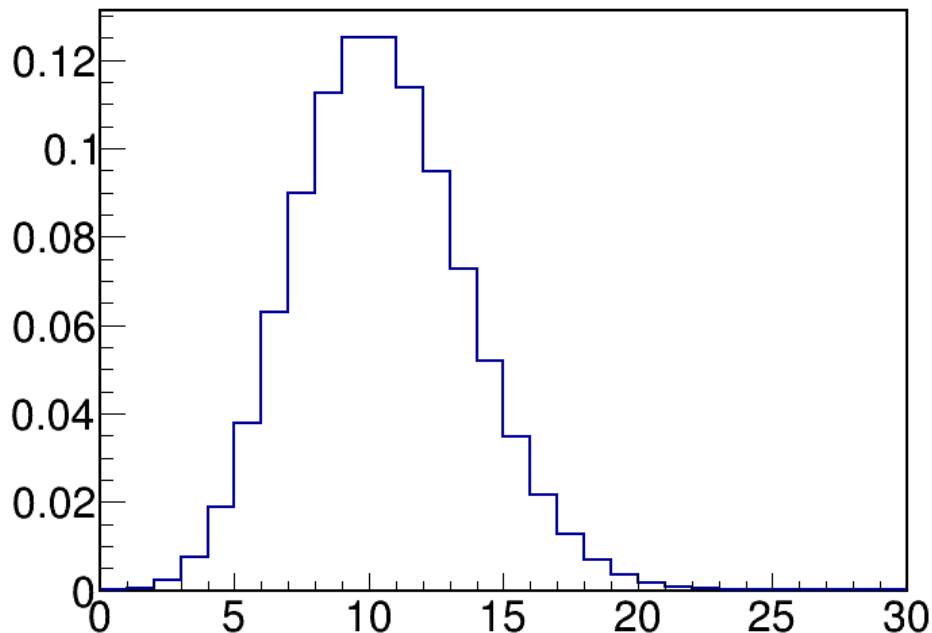
→ In practice,  $P \ll 1$ ,  $N \gg 1$ ,  $\Rightarrow$  Poisson approximation.

→ i.e. **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution:**  $P(n; \lambda) = e^{-\lambda} \frac{\lambda^n}{n!}$

$\uparrow$   $(1-P)^{N-n} \stackrel{n \ll N}{\approx} \left(1 - \frac{\lambda}{N}\right)^N \stackrel{N \gg 1}{\approx} e^{-\lambda}$

$\lambda = 10$



**Mean** =  $\lambda$

**Variance** =  $\lambda$

$\sigma = \sqrt{\lambda}$

Central limit theorem :

becomes **Gaussian for large  $\lambda$**  :

$$P(\lambda) \xrightarrow{\lambda \rightarrow \infty} G(\lambda, \sqrt{\lambda})$$



# How to count

Collider processes: produce (many) events  $N$ , select a (very) small fraction  $P$

→ In principle, binomial process

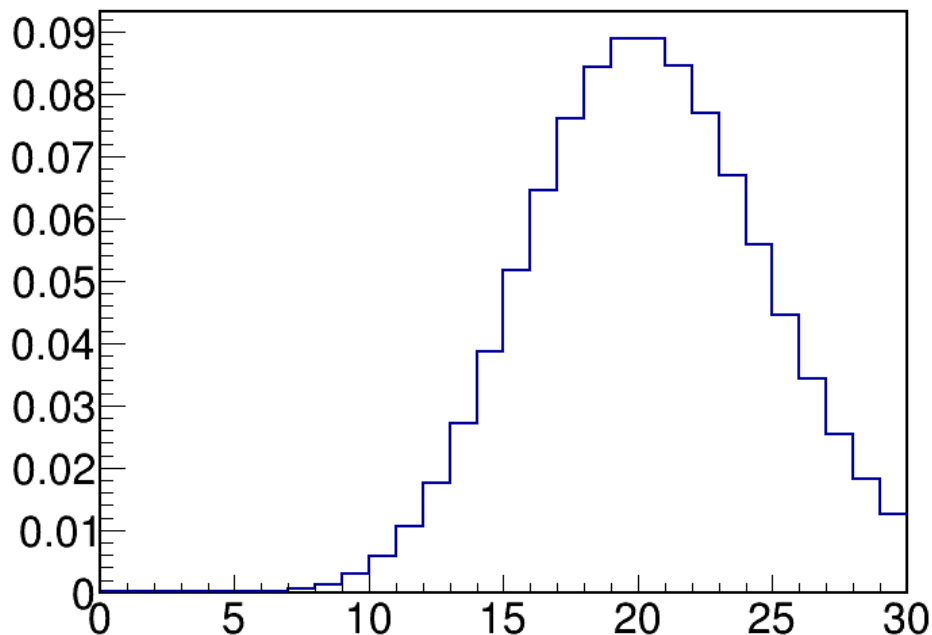
→ In practice,  $P \ll 1$ ,  $N \gg 1$ ,  $\Rightarrow$  Poisson approximation.

→ i.e. **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution:**  $P(n; \lambda) = e^{-\lambda} \frac{\lambda^n}{n!}$

$\uparrow$   $(1-P)^{N-n} \stackrel{n \ll N}{\approx} \left(1 - \frac{\lambda}{N}\right)^N \stackrel{N \gg 1}{\approx} e^{-\lambda}$

$\lambda = 20$



**Mean** =  $\lambda$

**Variance** =  $\lambda$

$\sigma = \sqrt{\lambda}$

Central limit theorem :

becomes **Gaussian for large  $\lambda$**  :

$$P(\lambda) \xrightarrow{\lambda \rightarrow \infty} G(\lambda, \sqrt{\lambda})$$

# Statistical Model for Counting

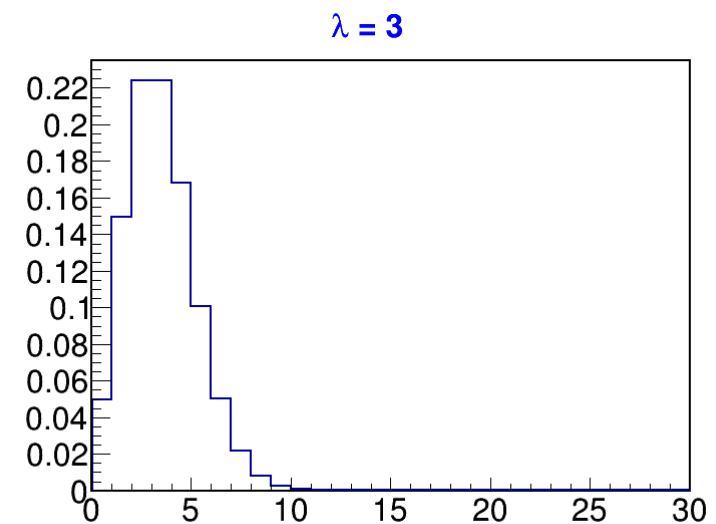
**Observable:** number of events  $n$

Typically both **S**ignal and **B**ackground present:

$$P(n; S, B) = e^{-(S+B)} \frac{(S+B)^n}{n!}$$

**S** : # of events from signal process

**B** : # of events from bkg. process(es)



Model has **parameters S** and **B**.

B can be known a priori or not (S usually not...)

→ Example: **assume B is known**, use **measured n** to find out about **S**.

# Multiple counting bins

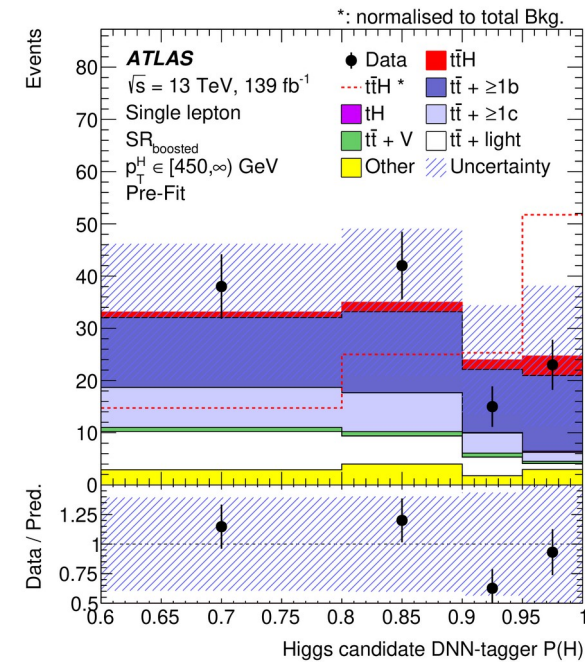
Count in bins of a variable  $\Rightarrow$  *histogram*  $n_1 \dots n_N$ .

(N : number of bins)

Per-bin fractions (=shapes)  
of Signal and Background

$$P(\{n_i\}; S, B) = \prod_{i=1}^N e^{-(Sf_{S,i} + Bf_{B,i})} \frac{(Sf_{S,i} + Bf_{B,i})^{n_i}}{n_i!}$$

Poisson distribution in each bin



**Shapes  $f$**  typically obtained from simulated events (*Monte Carlo*)

$\rightarrow$  HEP: typically excellent modeling from simulation, although some uncertainties need to be accounted for.

However not always possible to generate sufficiently large MC samples

**MC stat fluctuations** can create artefacts, especially for  $S \ll B$ .

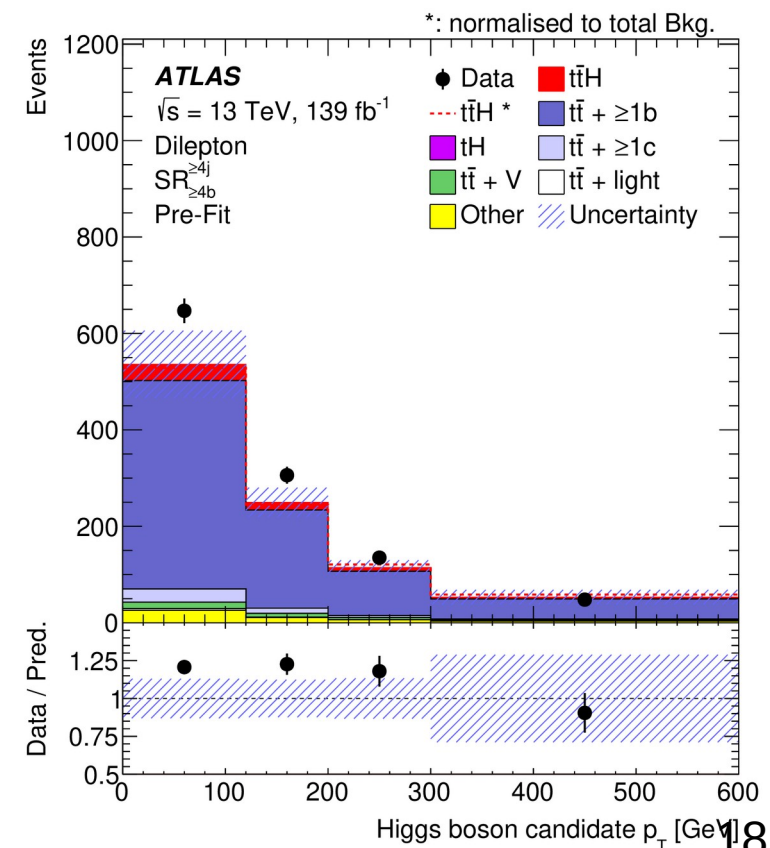
# Model Parameters

Model typically includes:

- **Parameters of interest** (POIs) : what we want to measure  
→  $S, m_W, \dots$
  - **Nuisance parameters** (NPs) : other parameters needed to define the model  
→ Background levels (**B**)  
→ For binned data,  $f_{\text{sig}}$ ,  $f_{\text{bkg}}$
- 

NPs must be either:

- **Known a priori** (within uncertainties) or
- **Constrained by the data**

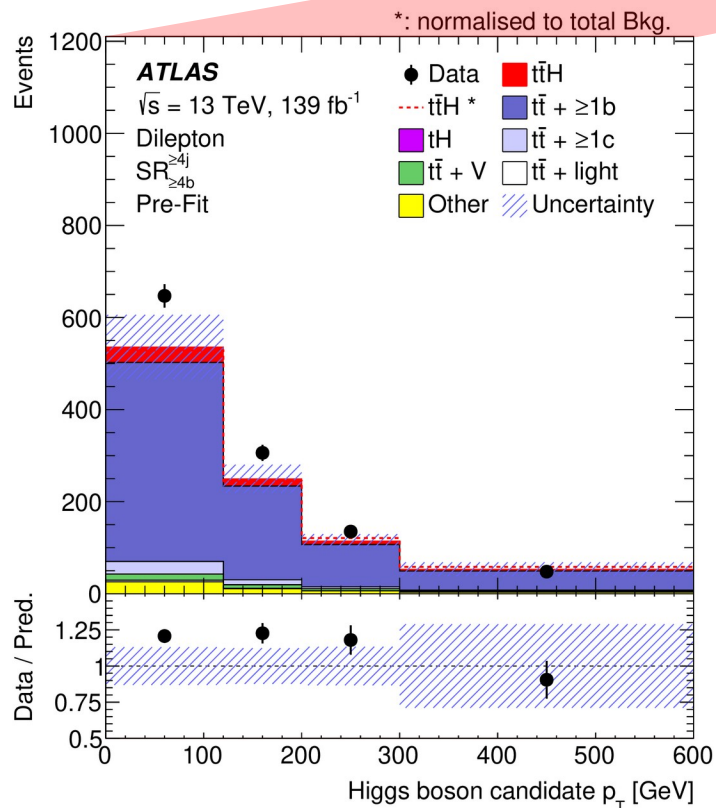
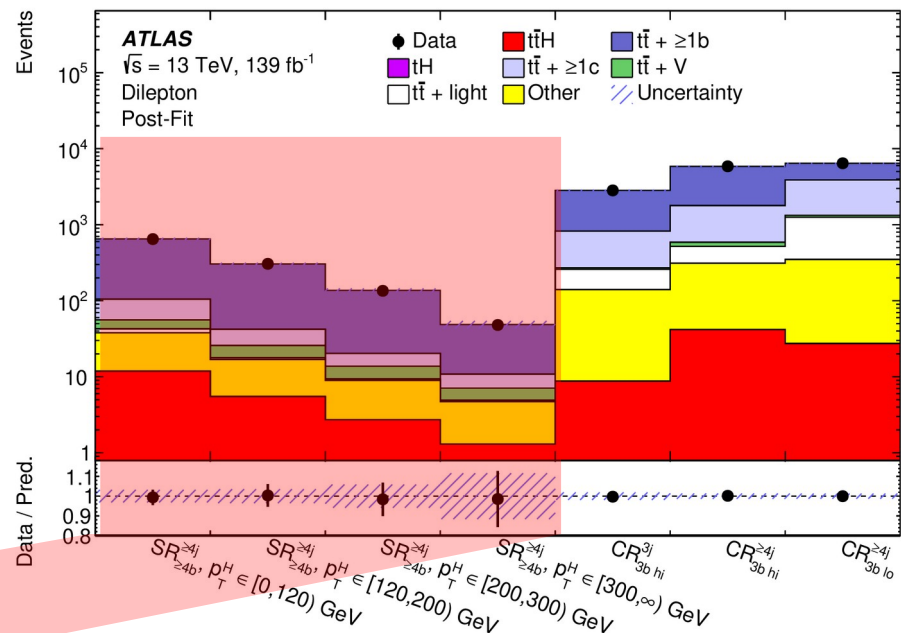


**Multiple analysis regions** often used.

→ Exploit better sensitivity in some regions

Here 7 regions:

→ 4 *Signal Regions (SR)* split in  $p_T(\text{Higgs})$



**Better sensitivity at high  $p_T$**

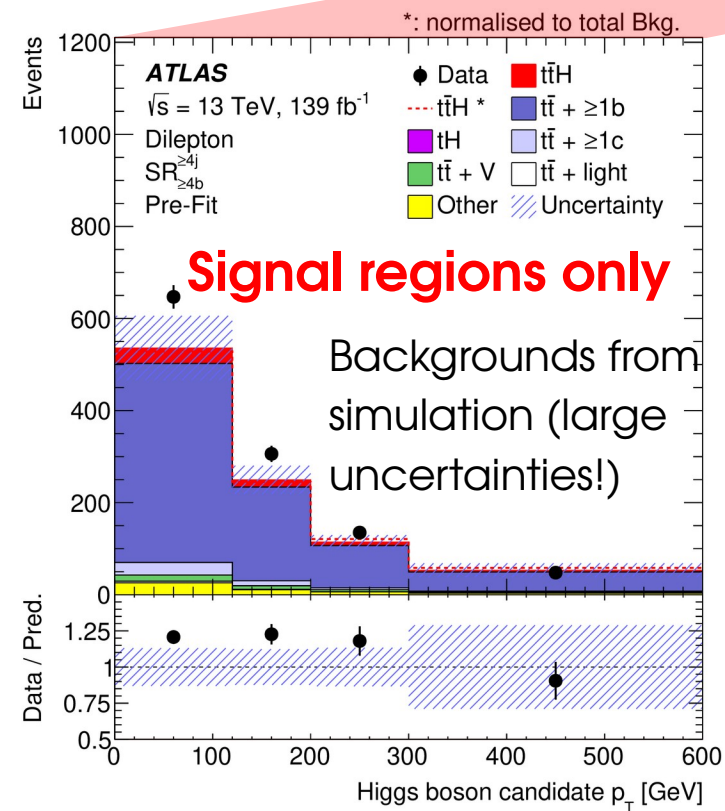
→ lower B backgrounds, higher S/B

**Backgrounds levels obtained from simulation here**

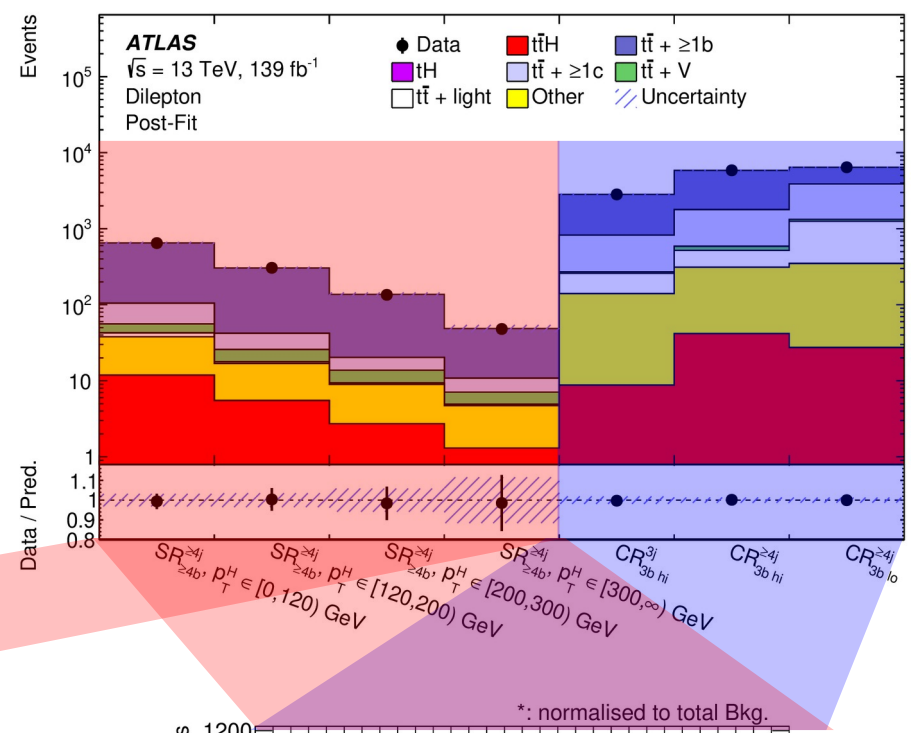
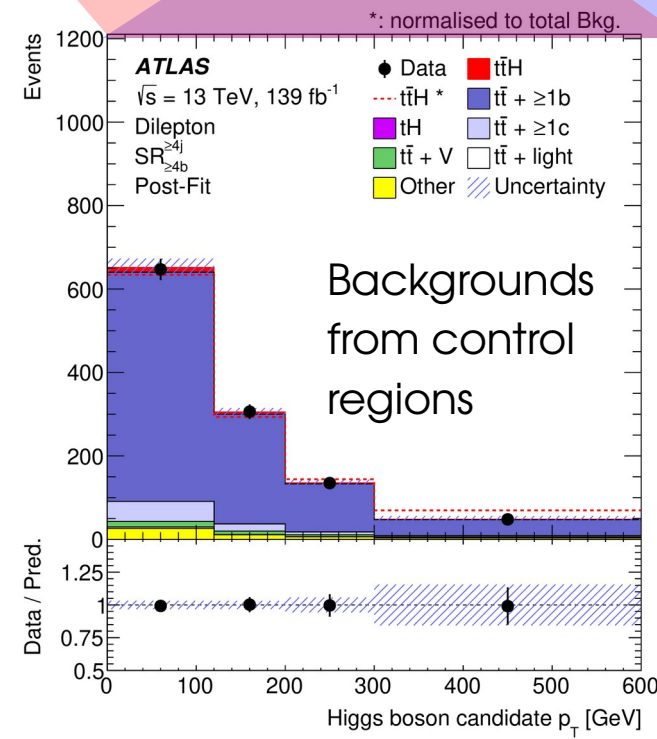
→ Large uncertainties!

- Multiple analysis regions often used.
- Exploit better sensitivity in some regions
  - Constrain NPs: **Control regions** for bkg

- Here 7 regions:
- 4 *Signal Regions* (**SR**) split in  $p_T(\text{Higgs})$
  - 3 *Background Control Regions* (**CR**)



Include  
Background CRs





**Multiple analysis regions** often used.

- Exploit better sensitivity in some regions
- Constrain NPs: **Control regions** for bkg

Here 7 regions:

- 4 *Signal Regions* (**SR**) split in  $p_T(\text{Higgs})$
- 3 *Background Control Regions* (**CR**)

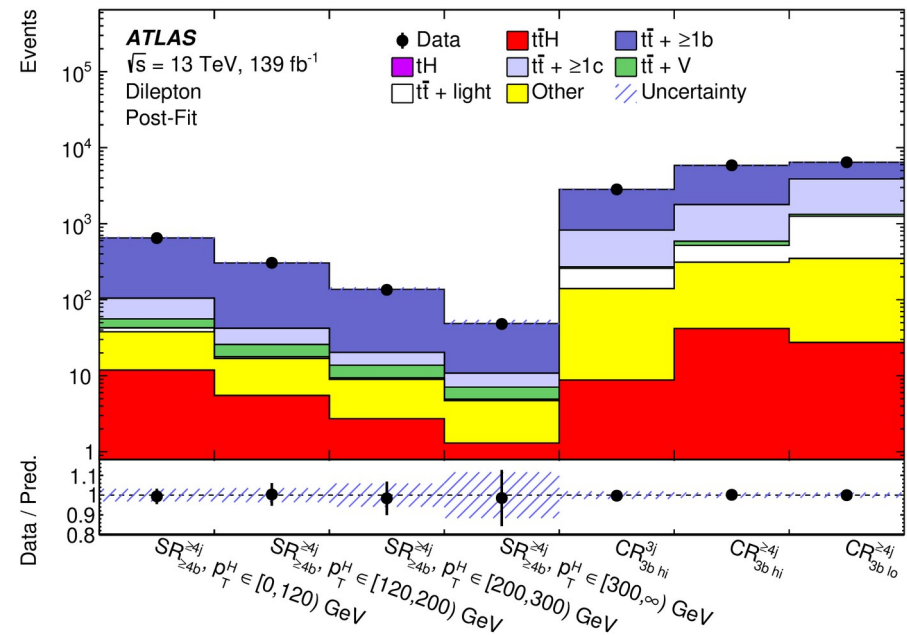
⇒ **Combined PDF** :

PDF for category k

$$P(\mathbf{S}; \{\mathbf{n}_i^{(k)}\}_{i=1 \dots n_{\text{evts}}^{(k)}}^{k=1 \dots n_{\text{cats}}}) = \prod_{k=1}^{n_{\text{cats}}} P_k(\mathbf{S}; \{\mathbf{n}_i^{(k)}\}_{i=1 \dots n_{\text{evts}}^{(k)}})$$

No overlaps between categories ⇒ No statistical correlations

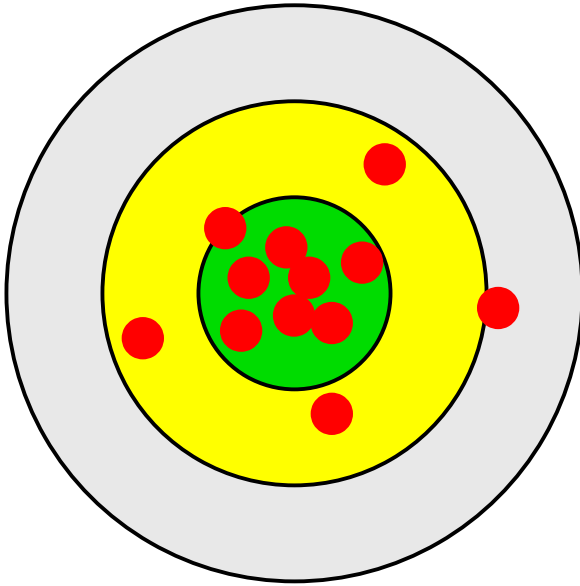
⇒ can simply take product of individual PDFs.



# Systematic Errors

---

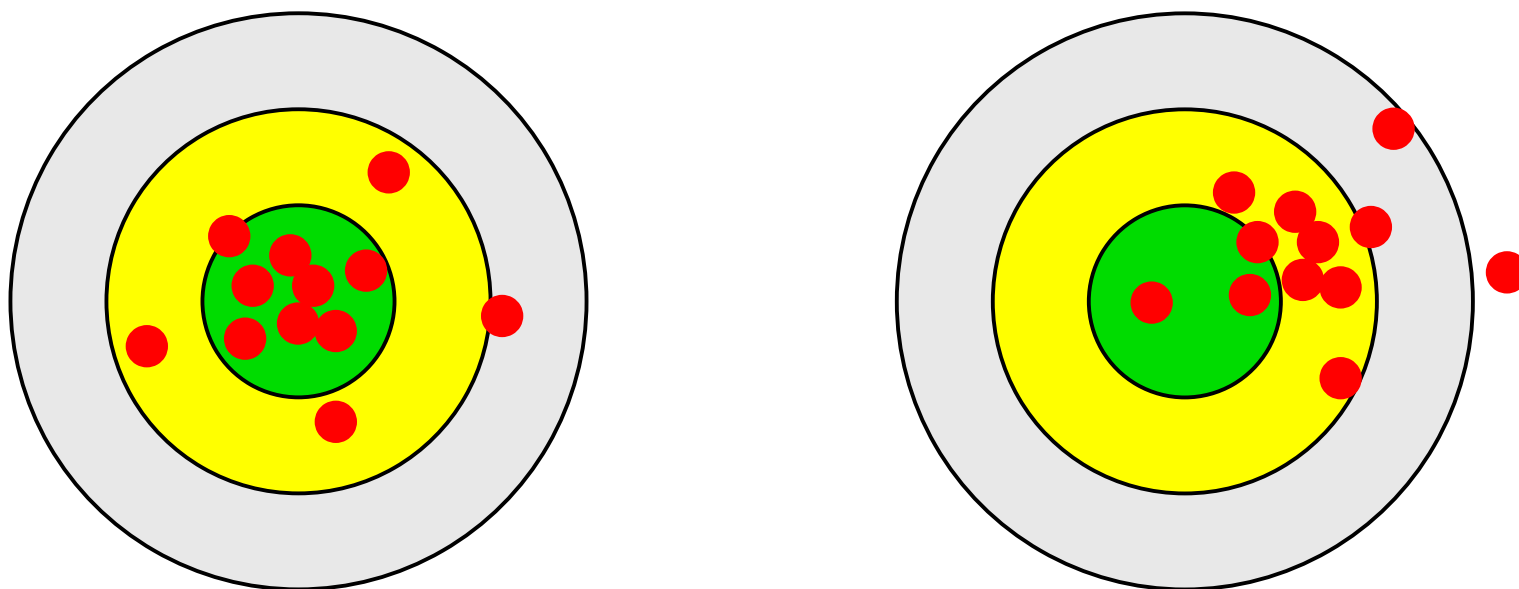
The statistical model (PDF) is a way to express **uncertainty** on the outcome of an experiment. e.g. 2D Gaussian :



These uncertainties are also called **Statistical Uncertainties** – they are the ones encoded in the model PDF.

# Systematic Errors

The statistical model (PDF) is a way to express **uncertainty** on the outcome of an experiment. e.g. 2D Gaussian :



These uncertainties are also called **Statistical Uncertainties** – they are the ones encoded in the model PDF.

However **the model itself may be wrong** : this is a *systematic error*  
→ To account for them, need a set of **Systematic uncertainties**

# Systematics

Statistical models include:

- **Parameters of interest** (POIs) :  $\mathbf{S}, \sigma \times \mathbf{B}, m_W, \dots$
- **Nuisance parameters** (NPs) : other parameters needed to define the model  
→ Ideally, **constrained by data** like the POI

**And systematics ?**

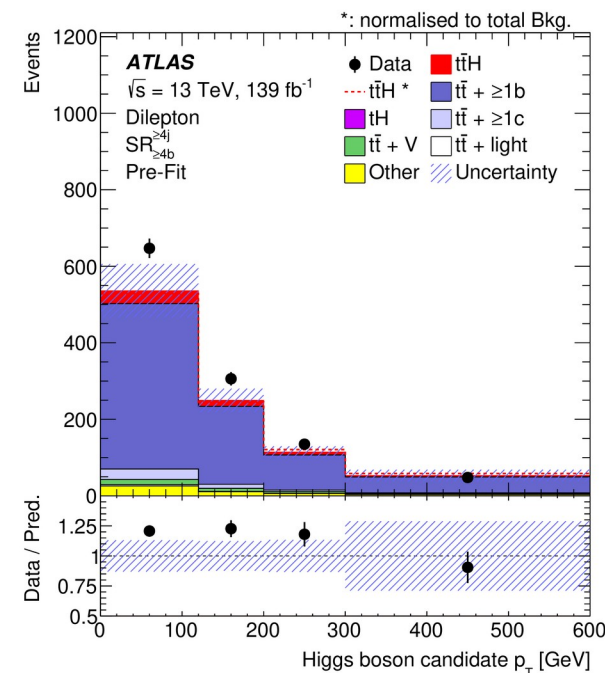
= Cover what we don't know about the random process.

⇒ **Parameterize using additional NPs**

→ Can't be constrained by the data ⇒ Add constraints in the likelihood

$$L(\underbrace{\mu}_{\text{POI}}, \underbrace{\theta}_{\text{Systematics NP}}; \text{data}) = \underbrace{L_{\text{measurement}}(\mu, \theta; \text{data})}_{\text{Measurement Likelihood}} \underbrace{C(\theta)}_{\text{NP Constraint term}}$$

$C(\theta)$  represents **external knowledge** about the NP



"Systematic uncertainty is, in any statistical inference procedure, the uncertainty due to the incomplete knowledge of the probability distribution of the observables.

G. Punzi, *What is systematics ?*

# Frequentist Systematics

**Prototype:** Systematics NP  $\rightarrow$  measured in a separate *auxiliary experiment*  
e.g. background levels

$\rightarrow$  Build the combined PDF of the main+auxiliary measurements

$$P(\mu, \theta; \text{data}) = P_{\text{main}}(\mu, \theta; \text{main data}) P_{\text{aux}}(\theta; \text{aux. data})$$

Independent measurements:  
 $\Rightarrow$  just a product

**Gaussian** form often used by default:  $P_{\text{aux}}(\theta; \text{aux. data}) = G(\theta^{\text{obs}}; \theta, \sigma_{\text{syst}})$

In the combined likelihood, **systematic NPs are constrained**

$\rightarrow$  now same as NPs constrained in data.

$\rightarrow$  Often no clear setup for auxiliary measurements  
(e.g. theory simulation uncertainties)

$\rightarrow$  **Define constraints “by hand”** (“pseudo-measurement”)

# Statistical model, the full version

$$P(\mu, \{\theta_j\}_{j=1 \dots n_{NP}}; \{n_i^{(k)}\}_{i=1 \dots n_{data}^{(k)}}^{k=1 \dots n_{cat}}, \{\theta_j^{obs}\}_{j=1 \dots n_{NP}}) =$$

Expected bin yield

$$\prod_{k=1}^{n_{cats}} P[n_i; \mu \epsilon_{i,k}(\vec{\theta}) N_{S,i,k}(\vec{\theta}) + B_{i,k}(\vec{\theta})] \prod_{j=1}^{n_{syst}} G(\theta_j^{obs}; \theta_j; 1)$$

Bin Yields or Observable values

POI

Sig/Bkg Shapes, efficiencies

NPs

Systematics

Auxiliary Data

Pseudo-experiments

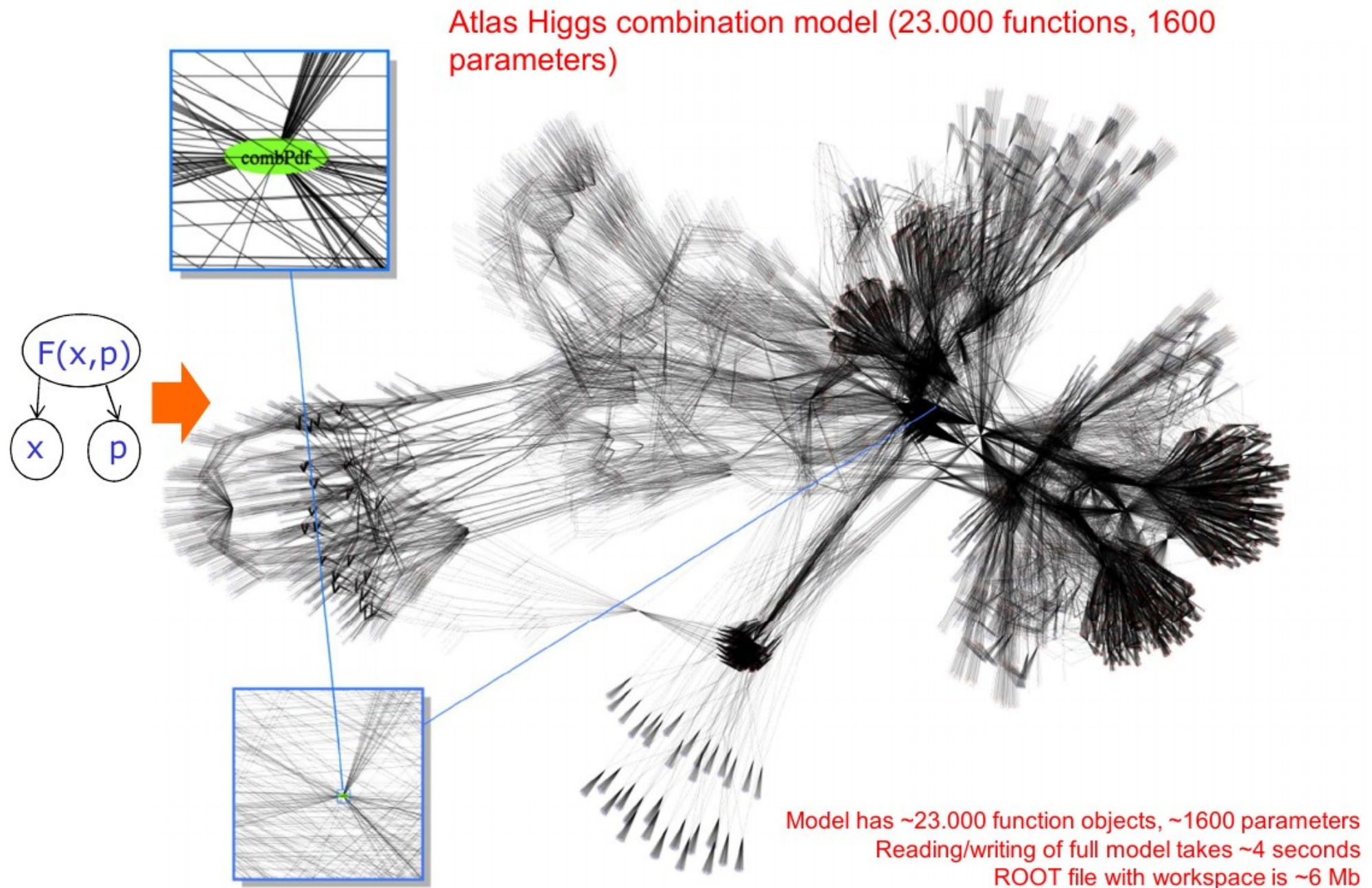
Data

MC

x number of categories!



# ATLAS Higgs Run 1 Combination Model



---

# HEP Statistical Inference : Confidence Intervals

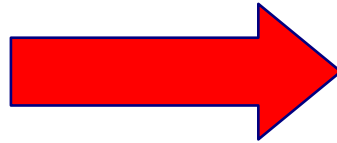


# Using the PDF

Model describes the distribution of the observable:  $P(\text{data}; \text{parameters})$

$$P(S=5)$$

provide parameter  
values



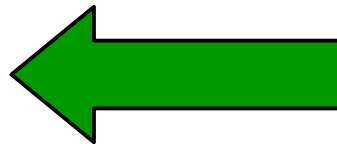
**Generate**

2, 5, 3, 7, 4, 9, ....

Each entry = a separate  
pseudo-data "experiment"

We want the **other** direction: **use data to get information on parameters**

$$P(S=?)$$



**Estimate**

2

THE observed data

# Maximum Likelihood

Define likelihood  $L(\mu) = P(\text{data}; \mu)$   
 $\Rightarrow$  Implicitly a function of the data

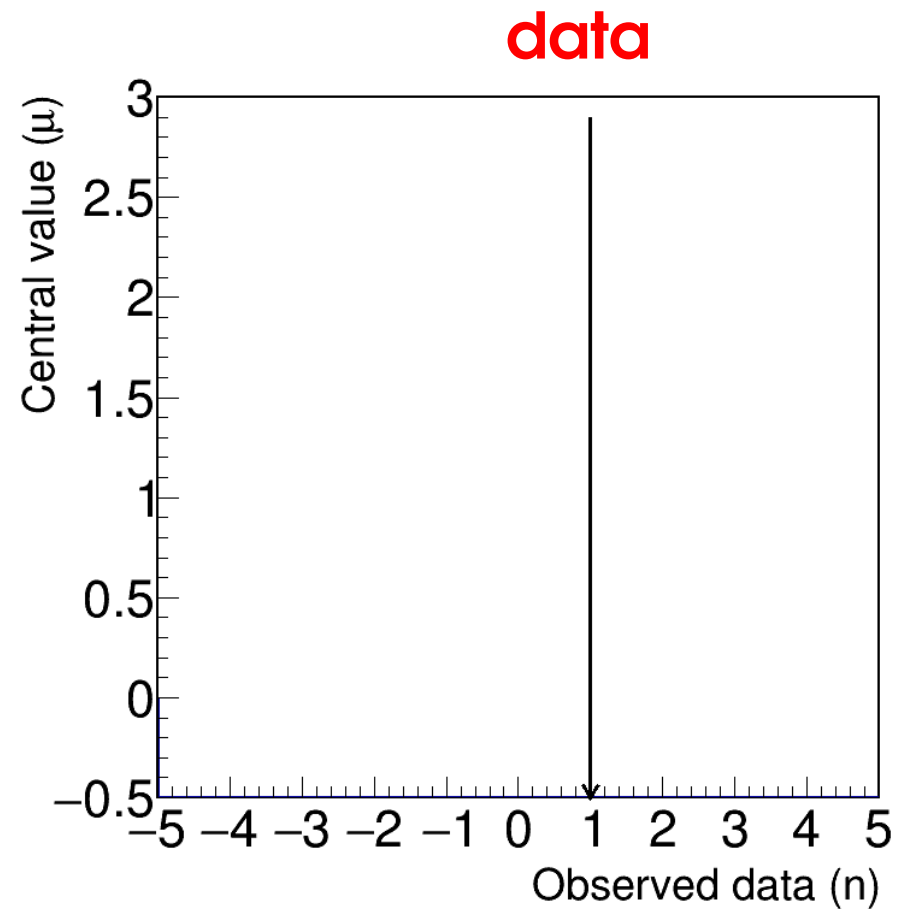
Estimate  $\mu$  as

$$\hat{\mu} = \operatorname{argmax}_{\mu} L(\mu)$$

“Best fit” of model to data

Several good properties:

- Asymptotically **Gaussian**
- Asymptotically **Unbiased**
- Asymptotically **Efficient**:  $\sigma_{\hat{\mu}}$  is the lowest possible
- Always **consistent**  $\hat{\mu} \xrightarrow{n \rightarrow \infty} \mu^*$



$$P(\hat{\mu}) \propto \exp\left(-\frac{(\hat{\mu} - \mu^*)^2}{2\sigma_{\hat{\mu}}^2}\right) \quad \text{for } n \rightarrow \infty$$

# Maximum Likelihood

Define likelihood  $L(\mu) = P(\text{data}; \mu)$   
 $\Rightarrow$  Implicitly a function of the data

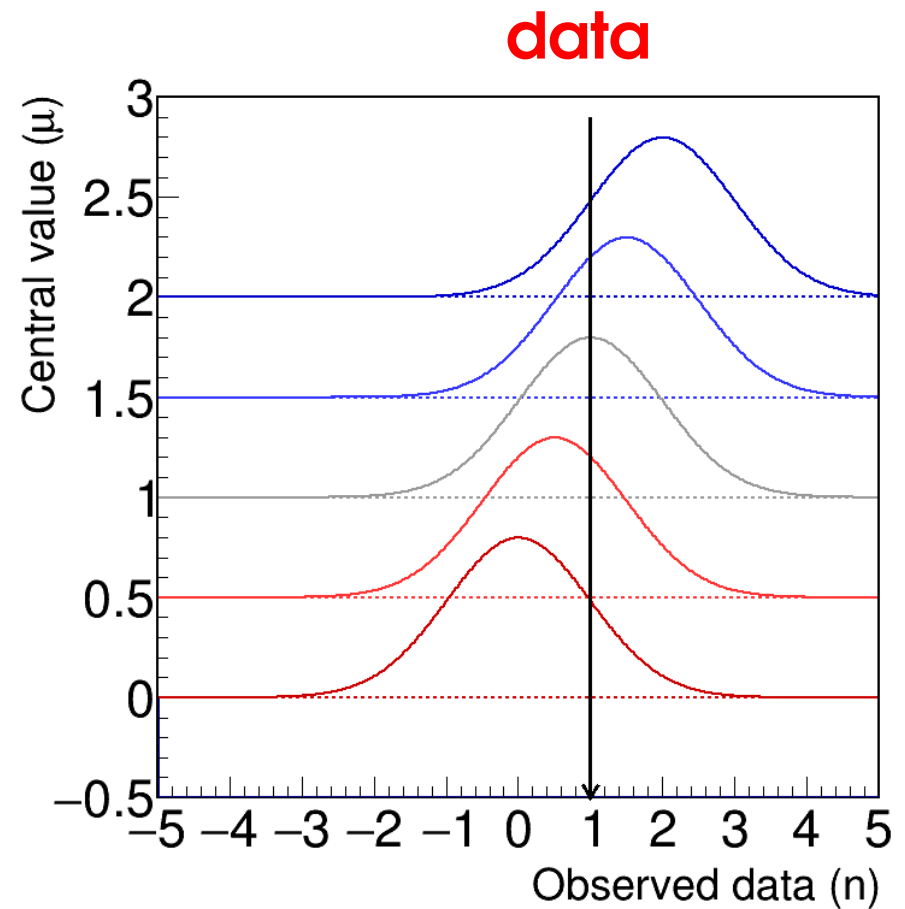
Estimate  $\mu$  as

$$\hat{\mu} = \operatorname{argmax}_{\mu} L(\mu)$$

“Best fit” of model to data

Several good properties:

- Asymptotically **Gaussian**
- Asymptotically **Unbiased**
- Asymptotically **Efficient**:  $\sigma_{\hat{\mu}}$  is the lowest possible
- Always **consistent**  $\hat{\mu} \xrightarrow{n \rightarrow \infty} \mu^*$



$$P(\hat{\mu}) \propto \exp\left(-\frac{(\hat{\mu} - \mu^*)^2}{2\sigma_{\hat{\mu}}^2}\right) \quad \text{for } n \rightarrow \infty$$

# Maximum Likelihood

Define likelihood  $L(\mu) = P(\text{data}; \mu)$   
 $\Rightarrow$  Implicitly a function of the data

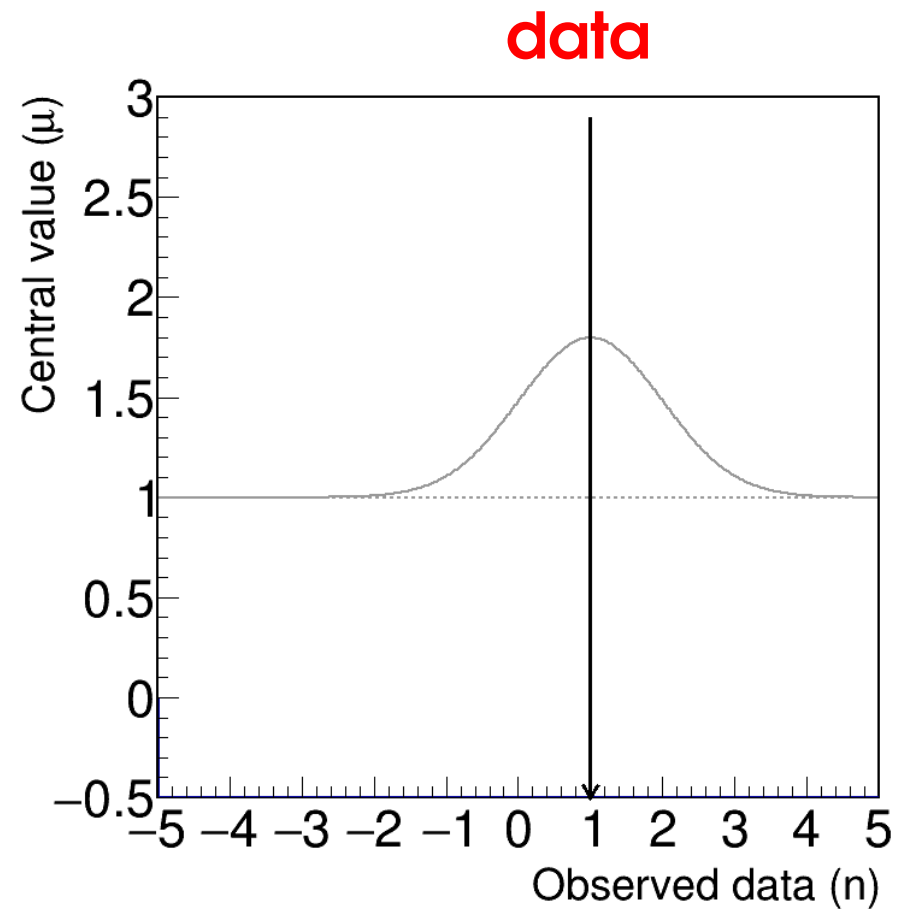
Estimate  $\mu$  as

$$\hat{\mu} = \operatorname{argmax}_{\mu} L(\mu)$$

“Best fit” of model to data

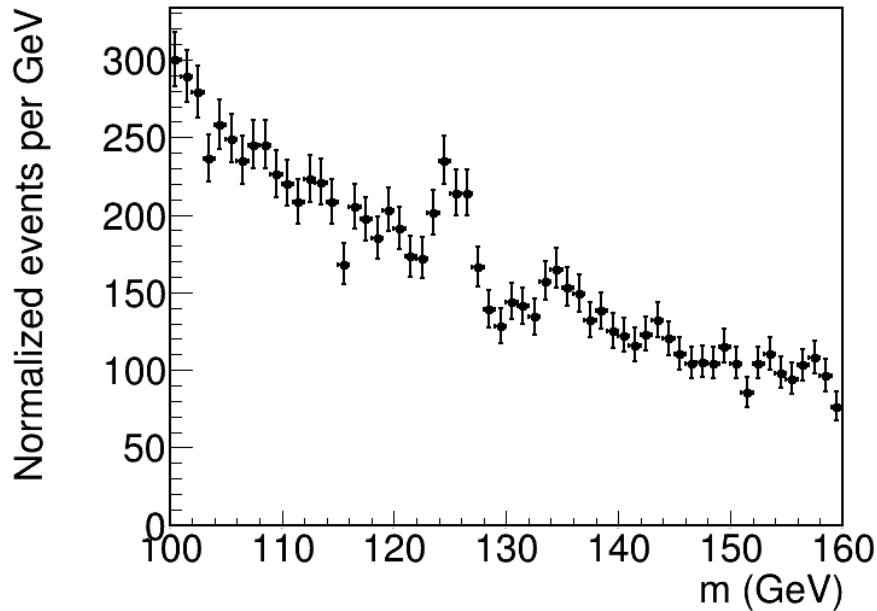
Several good properties:

- Asymptotically **Gaussian**
- Asymptotically **Unbiased**
- Asymptotically **Efficient**:  $\sigma_{\hat{\mu}}$  is the lowest possible
- Always **consistent**  $\hat{\mu} \xrightarrow{n \rightarrow \infty} \mu^*$



$$P(\hat{\mu}) \propto \exp\left(-\frac{(\hat{\mu} - \mu^*)^2}{2\sigma_{\hat{\mu}}^2}\right) \quad \text{for } n \rightarrow \infty$$

# Maximum Likelihood



Multiple Gaussian bins:

$$\lambda(\mu) = -2 \log L(\mu) = \sum_{i=1}^{N_{\text{bins}}} \left( \frac{n_i - \mu_i}{\sigma_i} \right)^2$$

**Maximum likelihood**

⇔ Minimum  $\chi^2$

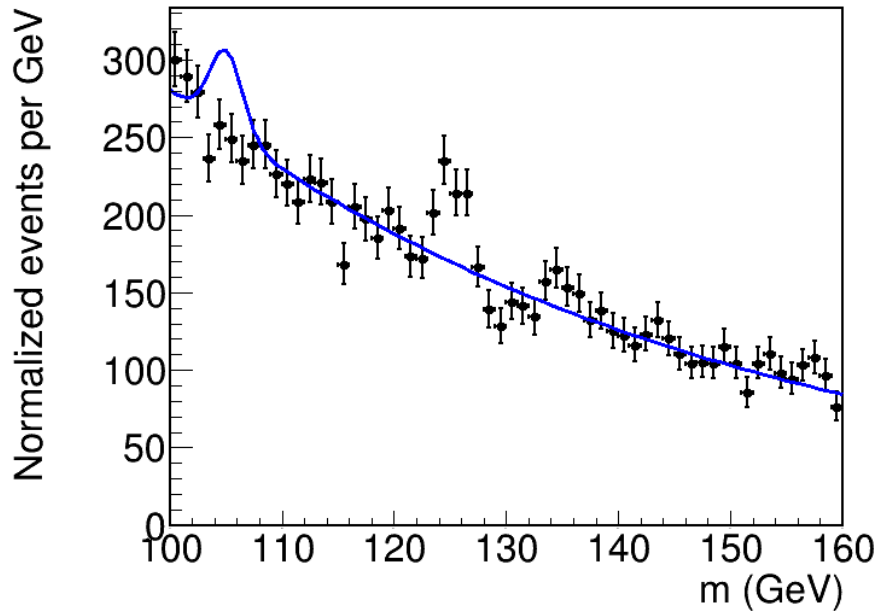
⇔ Least-squares minimization

However typically need to perform non-linear minimization.

HEP practice:

- **MINUIT** (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** – using NumPy/TensorFlow/PyTorch/... backends  
→ Usual methods – gradient-based, etc.

# Maximum Likelihood



Multiple Gaussian bins:

$$\lambda(\mu) = -2 \log L(\mu) = \sum_{i=1}^{N_{\text{bins}}} \left( \frac{n_i - \mu_i}{\sigma_i} \right)^2$$

**Maximum likelihood**

⇔ Minimum  $\chi^2$

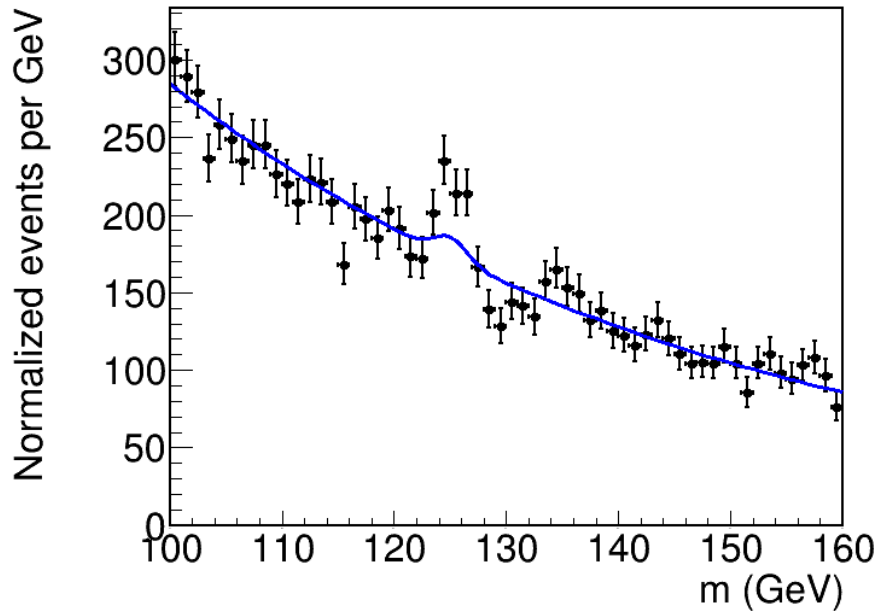
⇔ Least-squares minimization

However typically need to perform non-linear minimization.

HEP practice:

- **MINUIT** (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** – using NumPy/TensorFlow/PyTorch/... backends  
→ Usual methods – gradient-based, etc.

# Maximum Likelihood



Multiple Gaussian bins:

$$\lambda(\mu) = -2 \log L(\mu) = \sum_{i=1}^{N_{\text{bins}}} \left( \frac{n_i - \mu_i}{\sigma_i} \right)^2$$

**Maximum likelihood**

⇔ Minimum  $\chi^2$

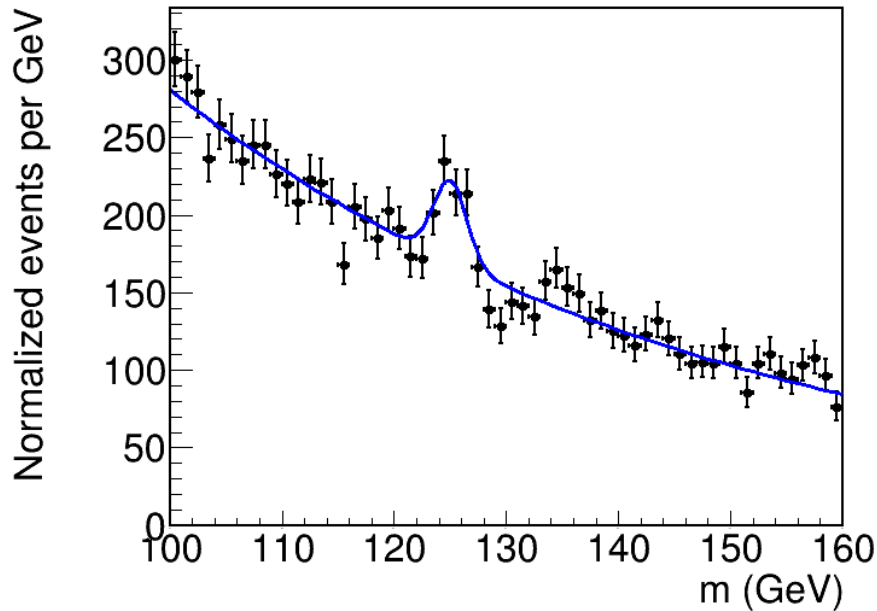
⇔ Least-squares minimization

However typically need to perform non-linear minimization.

HEP practice:

- **MINUIT** (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** – using NumPy/TensorFlow/PyTorch/... backends  
→ Usual methods – gradient-based, etc.

# Maximum Likelihood



Multiple Gaussian bins:

$$\lambda(\mu) = -2 \log L(\mu) = \sum_{i=1}^{N_{\text{bins}}} \left( \frac{n_i - \mu_i}{\sigma_i} \right)^2$$

**Maximum likelihood**

⇔ Minimum  $\chi^2$

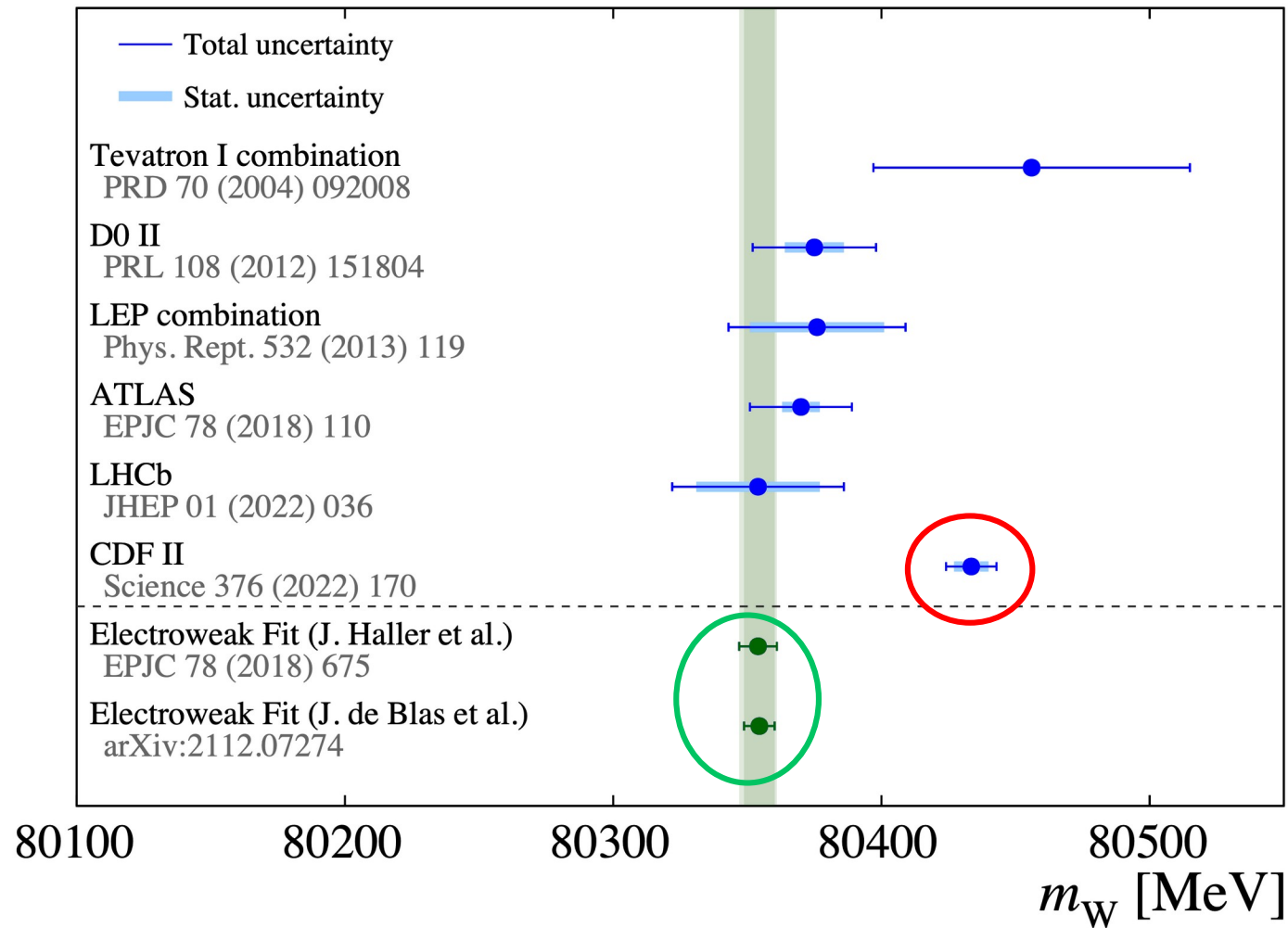
⇔ Least-squares minimization

However typically need to perform non-linear minimization.

HEP practice:

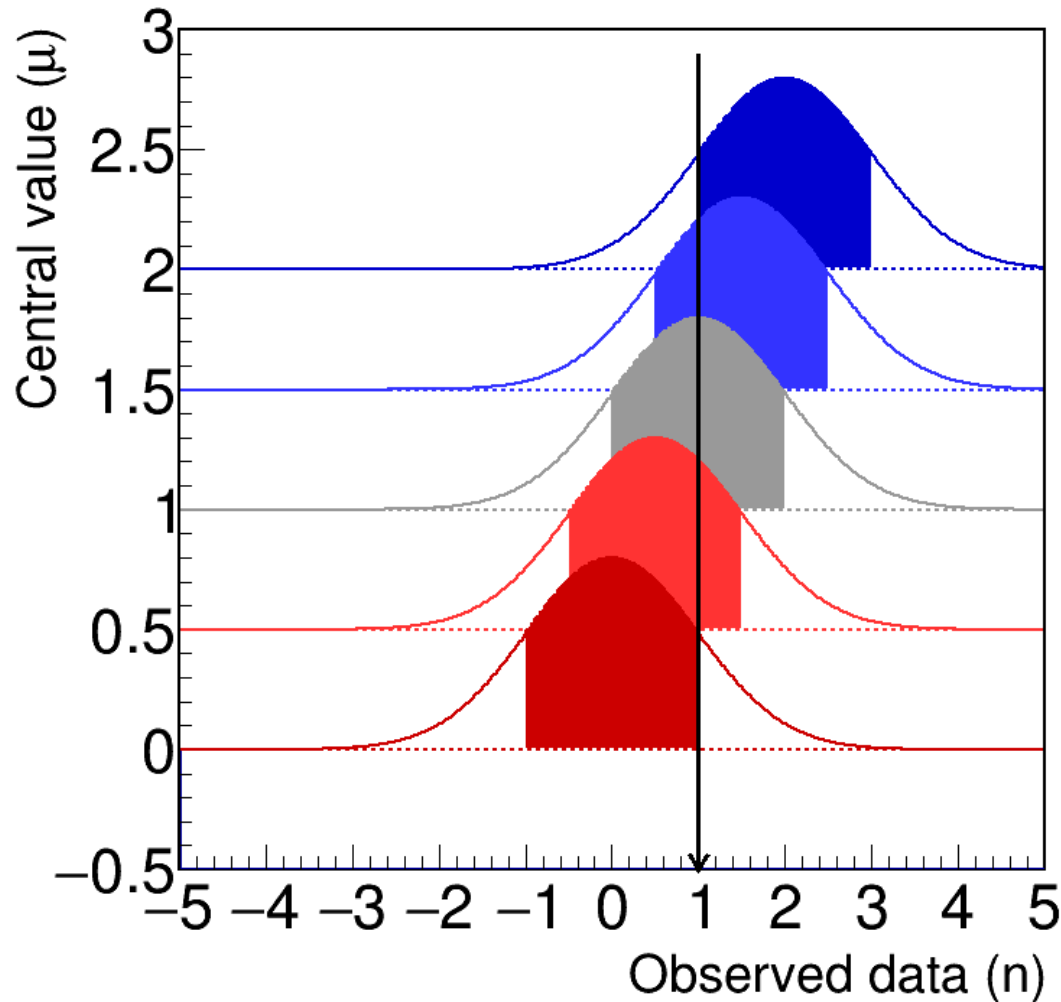
- **MINUIT** (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** – using NumPy/TensorFlow/PyTorch/... backends  
→ Usual methods – gradient-based, etc.





$$M_W = 80,433.5 \pm 6.4_{\text{stat}} \pm 6.9_{\text{syst}} = 80,433.5 \pm 9.4 \text{ MeV}/c^2$$

# Gaussian confidence intervals



Consider a Gaussian likelihood:

$$L(\mu) = \exp\left[-\frac{1}{2}\left(\frac{n-\mu}{\sigma}\right)^2\right]$$

$$P(\mu - \sigma < n < \mu + \sigma) \geq 68.3\%$$



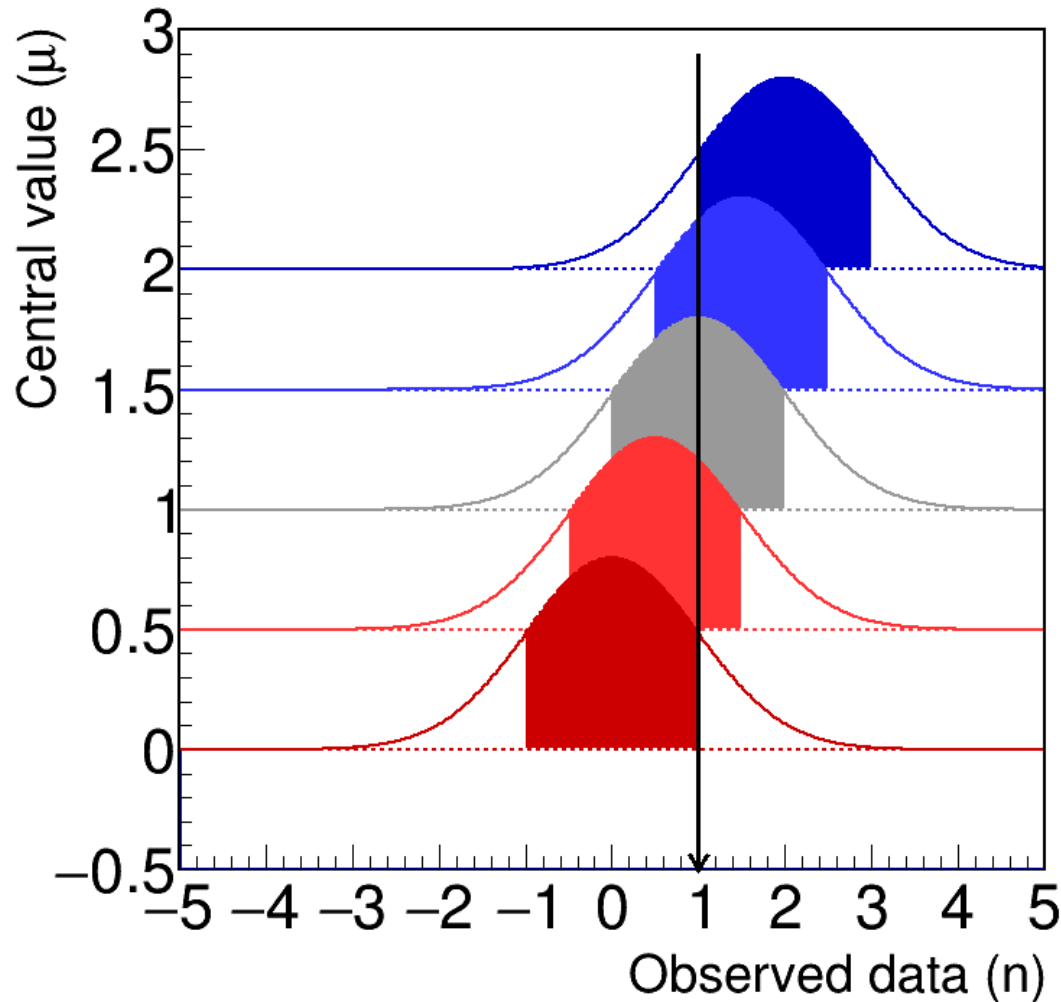
$$P(n - \sigma < \mu < n + \sigma) \geq 68.3\%$$

Still a statement on  $n$ !

$$\mu = n \pm \sigma \text{ at } 68.3\% \text{ CL}(1\sigma)$$

This interval will contain the true  $\mu$  value 68.3% of the time (“1 $\sigma$ ”)

# Gaussian confidence intervals



Consider a Gaussian likelihood:

$$L(\mu) = \exp\left[-\frac{1}{2}\left(\frac{n-\mu}{\sigma}\right)^2\right]$$

$$P(\mu - \sigma < n < \mu + \sigma) \geq 68.3\%$$



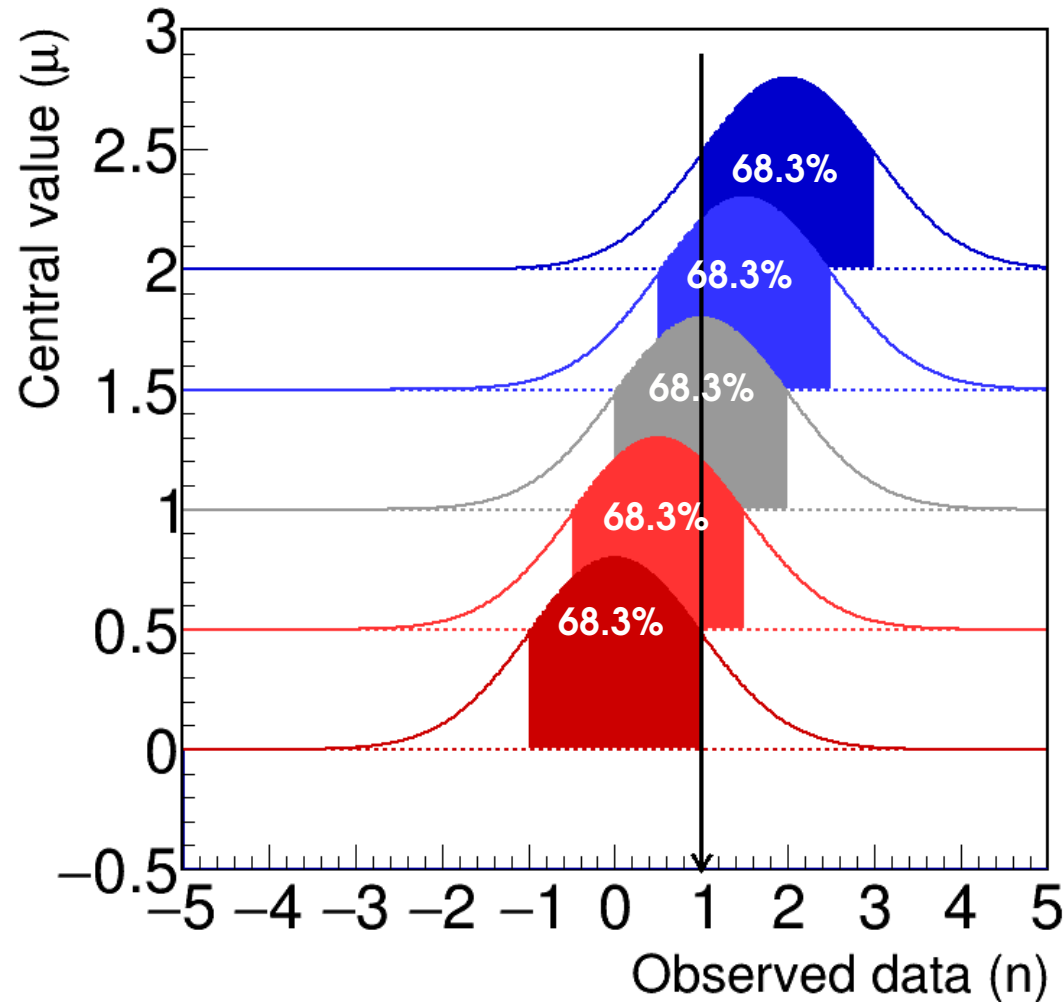
$$P(n - \sigma < \mu < n + \sigma) \geq 68.3\%$$

Still a statement on  $n$ !

$$\mu = n \pm \sigma \text{ at } 68.3\% \text{ CL}(1\sigma)$$

This interval will contain the true  $\mu$  value 68.3% of the time (“1 $\sigma$ ”)

# Gaussian confidence intervals



Consider a Gaussian likelihood:

$$L(\mu) = \exp\left[-\frac{1}{2}\left(\frac{n-\mu}{\sigma}\right)^2\right]$$

$$P(\mu - \sigma < n < \mu + \sigma) \geq 68.3\%$$



$$P(n - \sigma < \mu < n + \sigma) \geq 68.3\%$$

Still a statement on  $n$ !

$$\mu = n \pm \sigma \text{ at } 68.3\% \text{ CL}(1\sigma)$$

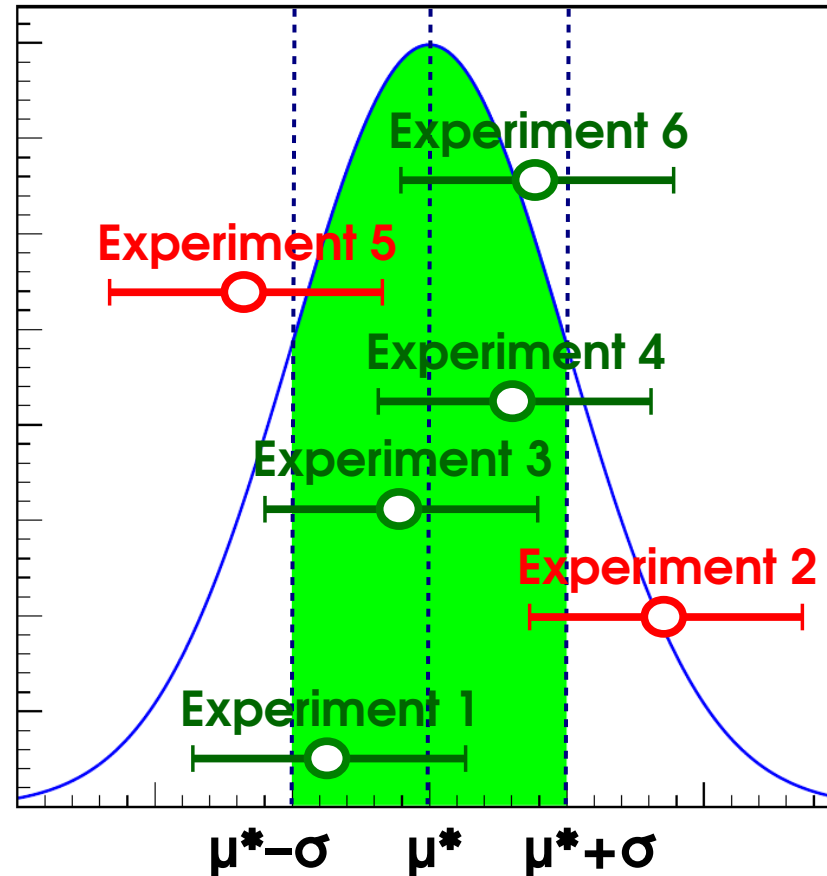
This interval will contain the true  $\mu$  value 68.3% of the time (“1 $\sigma$ ”)

# Gaussian confidence intervals

## Frequentist interpretation

If we would repeat the same experiment multiple times, with true value  $\mu^*$ , then 68.3% of the  $1\sigma$  intervals would contain  $\mu^*$ .

→ Crucially, this works even if we do not know  $\mu^*$  !



For each experiment, get the interval

$$\mu = n \pm \sigma \text{ at } 68.3\% \text{ CL}(1\sigma)$$

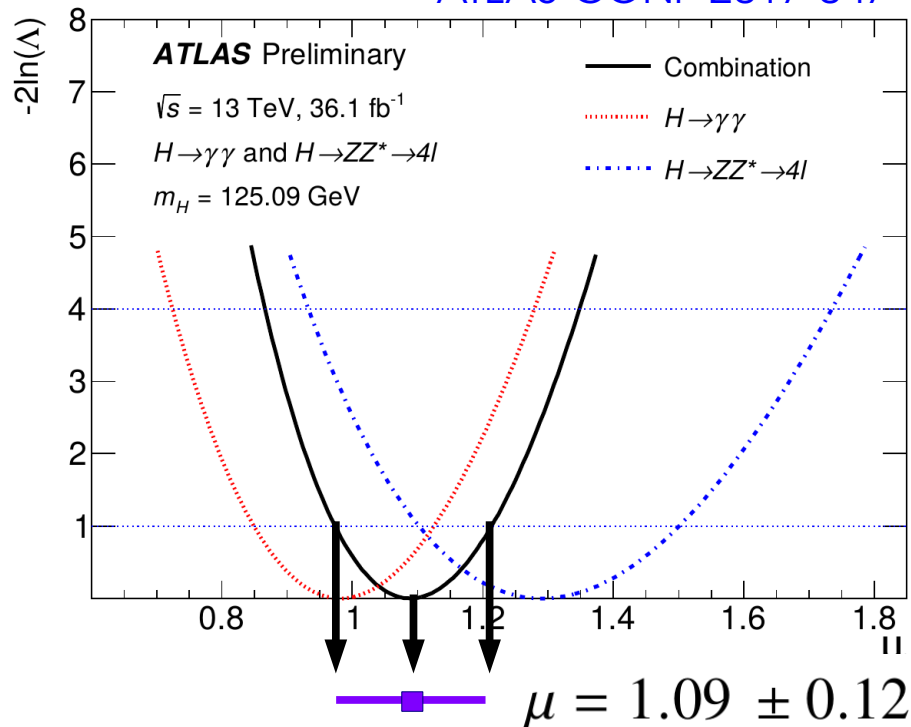
This interval will contain the true  $\mu$  value 68.3% of the time (“ $1\sigma$ ”)

# General case: Likelihood Intervals

## Confidence intervals from L:

- Test various values  $\mu$  using the **Profile Likelihood Ratio  $t(\mu)$**
- Minimum (=0) for  $\mu = \hat{\mu}$
- **Likelihood ratio** universally most powerful test for simple hypotheses (no NPs, single POI values), also used in other cases

ATLAS-CONF-2017-047



Probability to observe the data **for a given  $\mu$** . Use **conditional best-fit  $\hat{\theta}(\mu)$**  of the NPs for this  $\mu$ .

$$t(\mu) = -2 \log \frac{L(\mu, \hat{\theta}(\mu))}{L(\hat{\mu}, \hat{\theta})}$$

Probability to observe the data **for  $\hat{\mu}$** . Use **best-fit  $\hat{\theta}$**  for the NPs.

## Gaussian L( $\mu$ ):

- Parabolic in  $\mu$
- Minimum occurs at  $\mu = \hat{\mu}$
- $t(\mu)$  distributed as a  $\chi^2$
- $1\sigma$  interval  $[\mu_-, \mu_+]$  given by  $t(\mu_{\pm}) = 1$



# General case: Likelihood Intervals

## Confidence intervals from L:

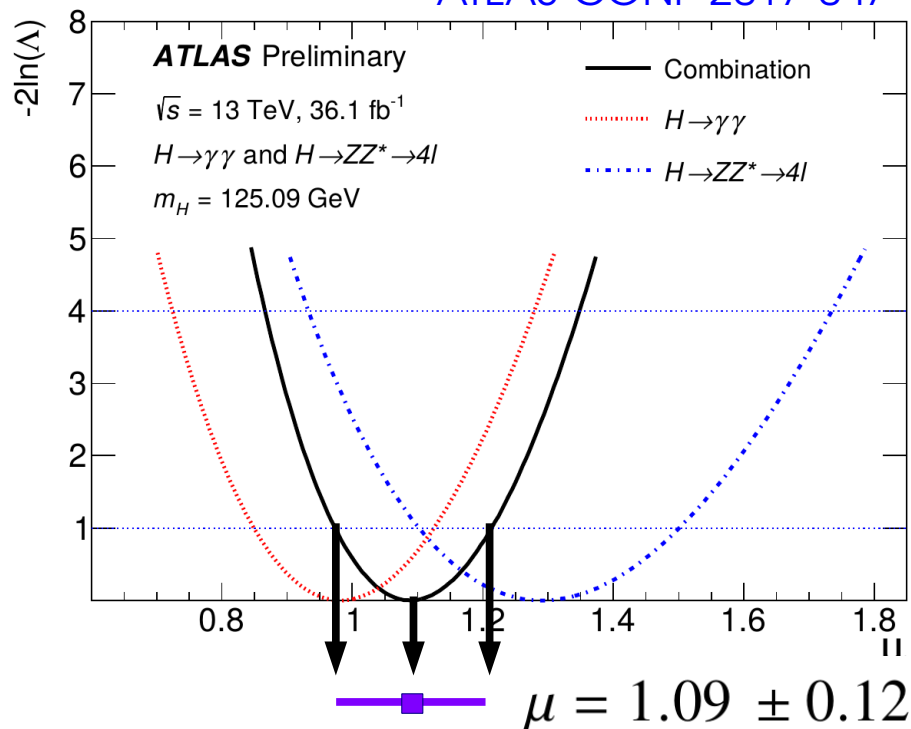
- Test various values  $\mu$  using the **Profile Likelihood Ratio**  $t(\mu)$
- Minimum (=0) for  $\mu = \hat{\mu}$
- **Likelihood ratio** universally most powerful test for simple hypotheses (no NPs, single POI values), also used in other cases

$$t(\mu) = -2 \log \frac{L(\mu, \hat{\theta}(\mu))}{L(\hat{\mu}, \hat{\theta})}$$

## General case:

- Generally not a perfect parabola
- Minimum still at  $\mu = \hat{\mu}$
- Distribution of  $t(\mu)$  ?

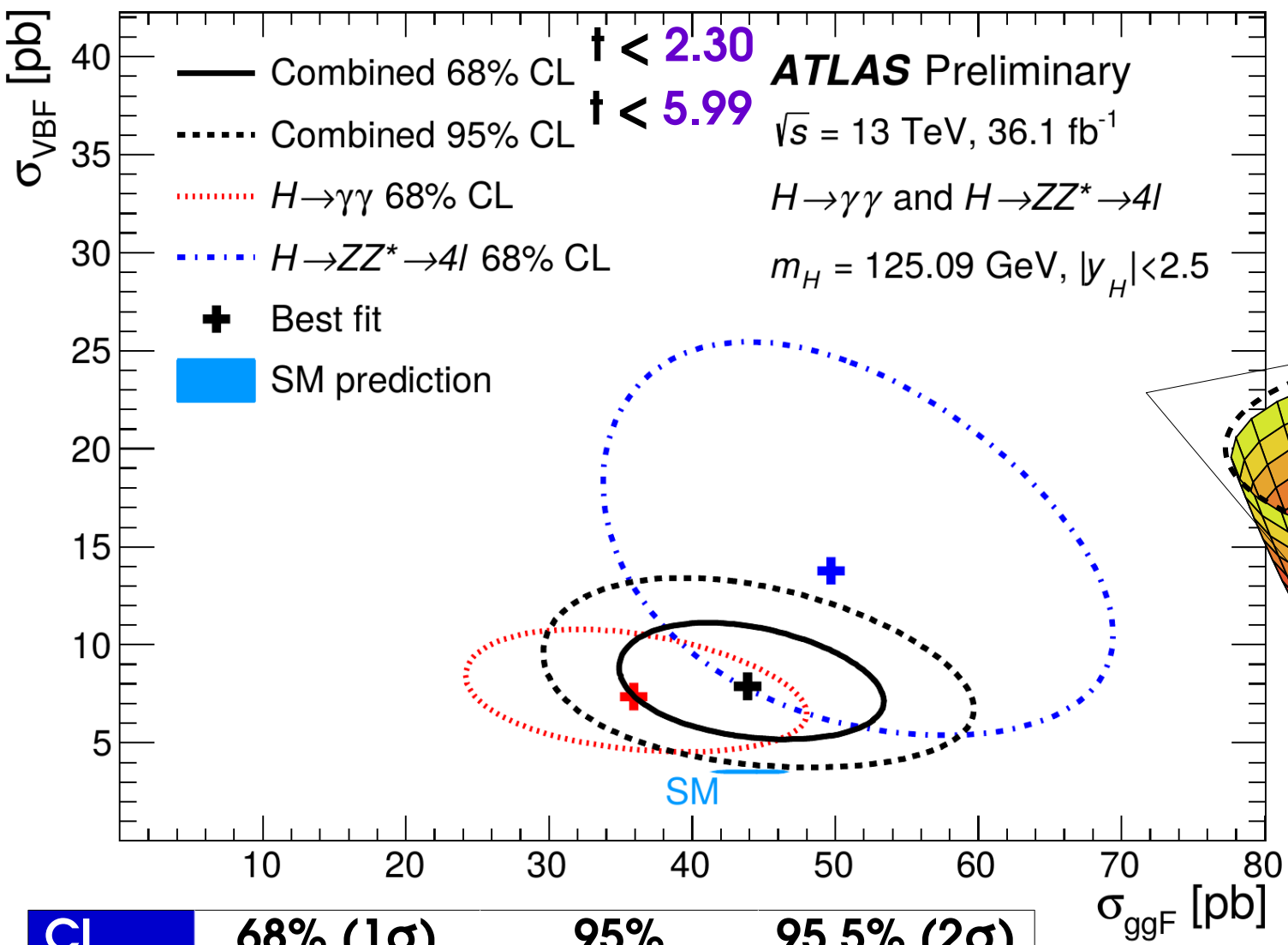
ATLAS-CONF-2017-047



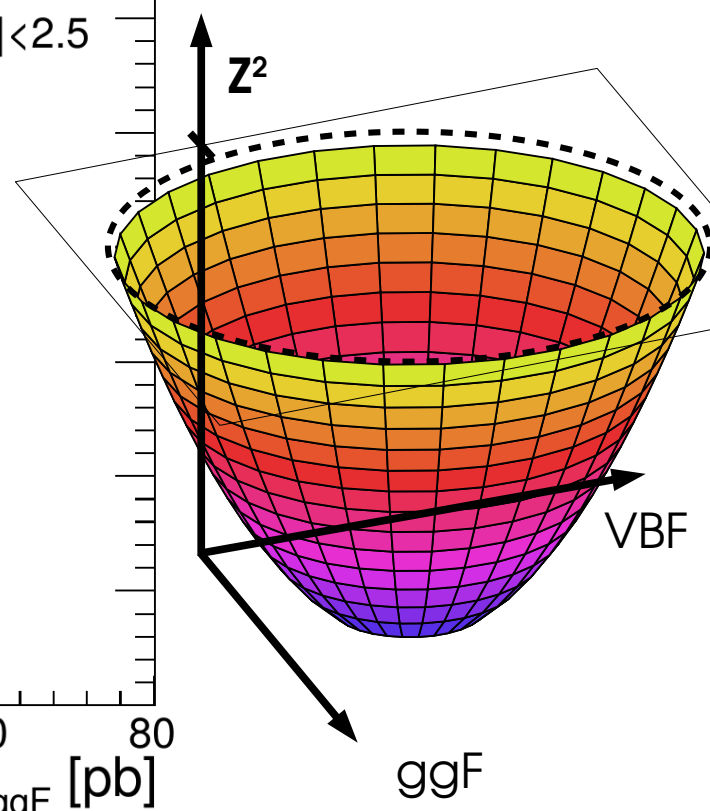
## Asymptotic approximation

- Compute  $t(\mu)$  using the exact  $L(\mu)$
- Assume  $t(\mu) \sim \chi^2$  as for Gaussian ("**Wilks' Theorem**")
- $1\sigma$  interval  $[\mu_-, \mu_+]$  given by  $t(\mu_{\pm}) = 1$
- Can also obtain exact intervals using pseudo-dataset sampling ("toys"), but generally not needed and rarely done.

# 2D Example: Higgs $\sigma_{\text{VBF}}$ vs. $\sigma_{\text{ggF}}$



$$t = -2 \log \frac{L(X_0, Y_0)}{L(\hat{X}, \hat{Y})}$$
$$\sim \chi^2(N_{\text{dof}}=2)$$



CL	68% ( $1\sigma$ )	95%	95.5% ( $2\sigma$ )
1D $Z^2$	1	3.84	4
2D $Z^2$	2.30	5.99	6.18

**Gaussian case:** elliptic paraboloid surface

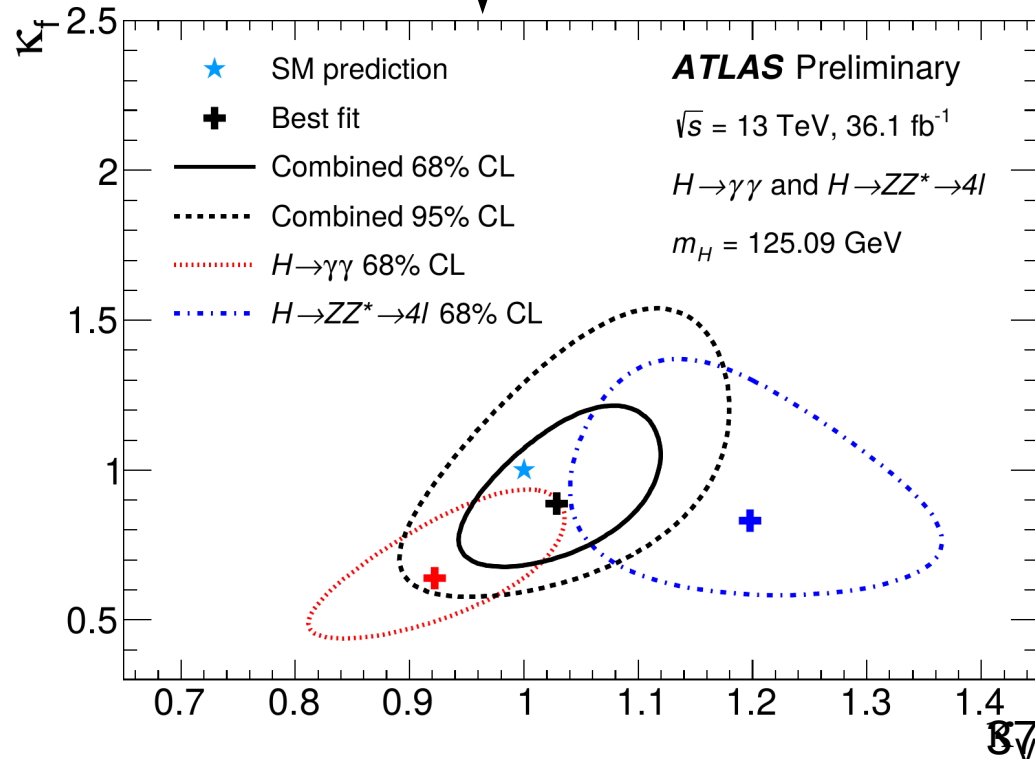
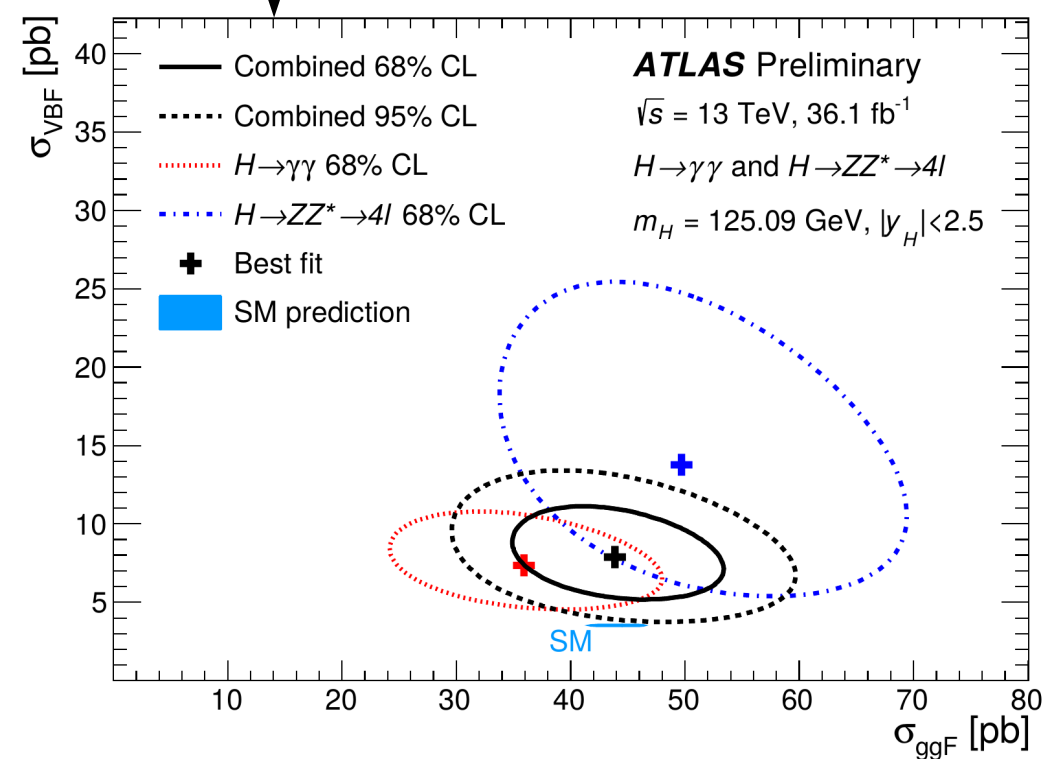
# Reparameterization

Start with basic measurement in terms of e.g.  $\sigma \times \mathbf{B}$

→ How to measure derived quantities (couplings, parameters in some theory model, etc.) ? → **just reparameterize the likelihood:**

e.g. Higgs couplings:  $\sigma_{ggF}$ ,  $\sigma_{VBF}$  sensitive to Higgs coupling modifiers  $\kappa_V$ ,  $\kappa_F$ .

$$L(\sigma_{ggF}, \sigma_{VBF}) \xrightarrow[\sigma_{VBF} \rightarrow \sigma_{VBF}(\kappa_V, \kappa_F)]{\sigma_{ggF} \rightarrow \sigma_{ggF}(\kappa_V, \kappa_F)} L(\sigma_{ggF}(\kappa_V, \kappa_F), \sigma_{VBF}(\kappa_V, \kappa_F)) \equiv L'(\kappa_V, \kappa_F)$$



# Example: Gaussian Profiling

Counting experiment with background uncertainty:  $\mathbf{n} = \mathbf{S} + \mathbf{B}$  :

$$\left. \begin{array}{l} \rightarrow \text{Signal region (SR)}: \mathbf{n}_{\text{obs}} \sim \mathbf{G}(\mathbf{S} + \mathbf{B}, \sigma_{\text{stat}}) \\ \rightarrow \text{Control region (CR)}: \mathbf{B}_{\text{obs}} \sim \mathbf{G}(\mathbf{B}, \sigma_{\text{bkg}}) \end{array} \right\} L(\mathbf{S}, \mathbf{B}) = G(\mathbf{n}_{\text{obs}}; \mathbf{S} + \mathbf{B}, \sigma_{\text{stat}}) G(\mathbf{B}_{\text{obs}}; \mathbf{B}, \sigma_{\text{bkg}})$$

**Recall:** Signal region only (fixed B):  $t_s = \left( \frac{S - n_{\text{obs}}}{\sigma_{\text{stat}}} \right)^2$   $S = (n_{\text{obs}} - B) \pm \sigma_{\text{stat}}$

→ Compute the best-fit (MLEs) for S and B

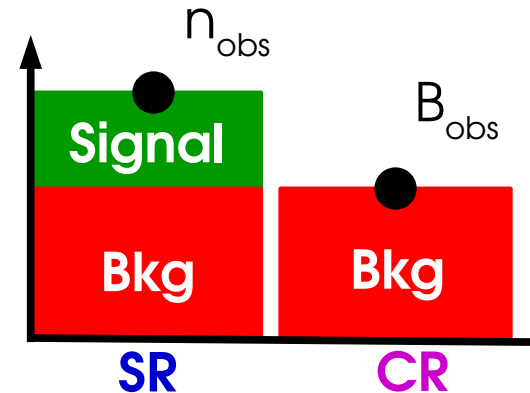
→ Show that the conditional MLE for B is

$$\hat{B}(S) = B_{\text{obs}} + \frac{\sigma_{\text{bkg}}^2}{\sigma_{\text{stat}}^2 + \sigma_{\text{bkg}}^2} (\hat{S} - S)$$

→ Compute the profile likelihood  $t_s$

→ Compute the  $1\sigma$  confidence interval on S

$$S = (n_{\text{obs}} - B_{\text{obs}}) \pm \sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{bkg}}^2} \quad \sigma_S = \sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{bkg}}^2}$$



**Stat uncertainty (on n) and systematic (on B) add in quadrature**

# Uncertainty decomposition

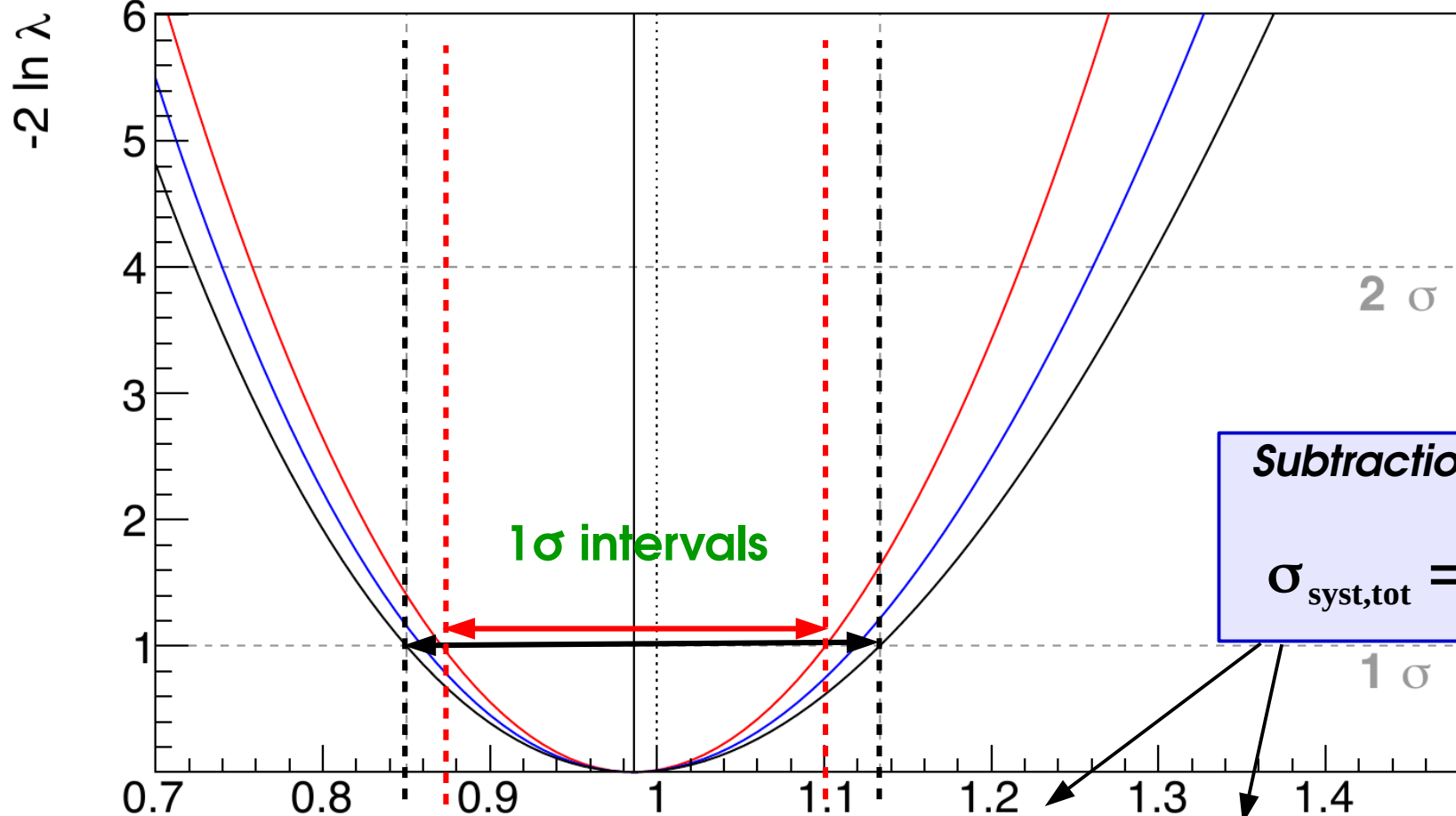
No systematics NPs included : statistical uncertainty only

All systematics NPs included: stat+syst uncertainties

**ATLAS**

$H \rightarrow \gamma\gamma, m_H = 125.09 \text{ GeV}$

— Total — Theory — Stat



$$\mu = 0.99 \pm 0.12 \text{ (stat)} \pm 0.06 \text{ (syst)} \pm 0.06 \text{ (theo)}^{\mu}$$

Systematics are described by NPs included in the fit. Define **pull** as

$$(\hat{\theta} - \theta_0) / \sigma_{\theta}$$

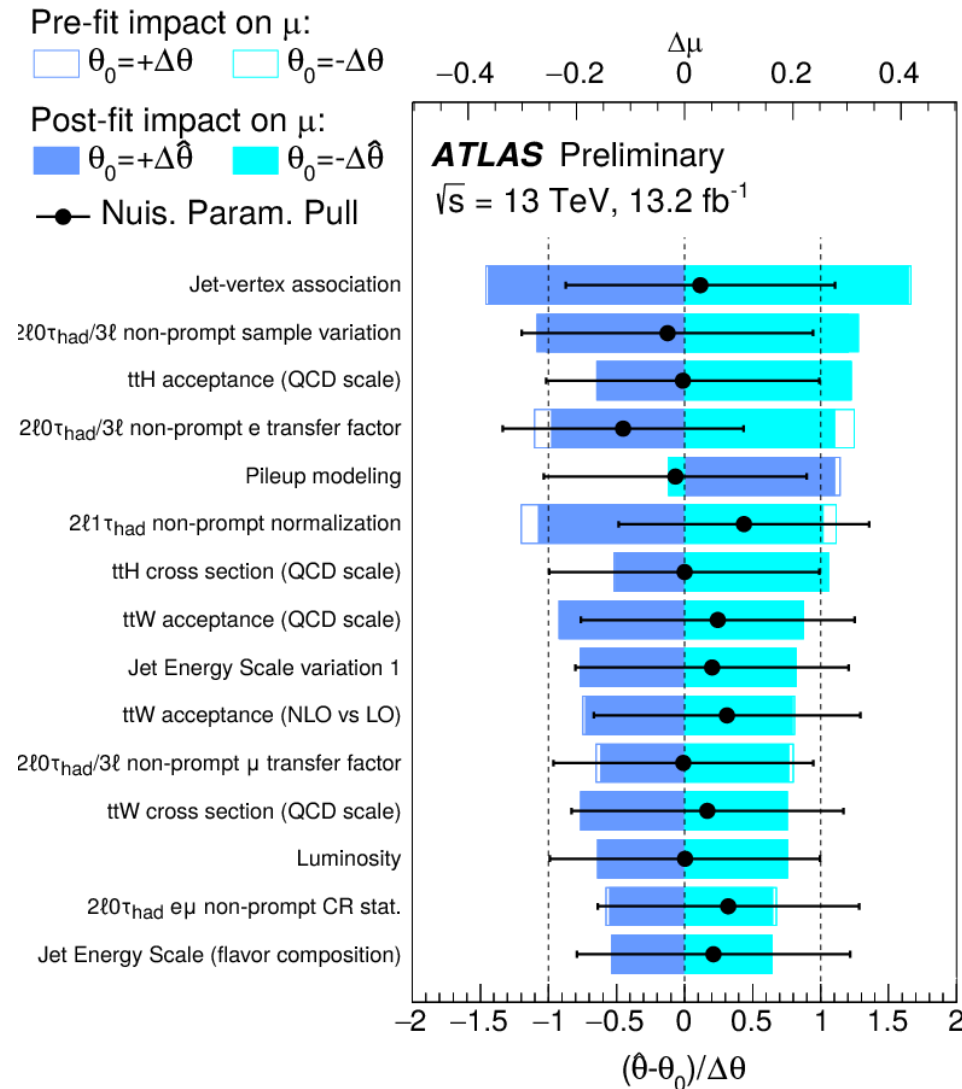
Nominally:

- **Pull = 0** : i.e. the pre-fit expectation
- **Pull uncertainty = 1** : from the Gaussian

However fit results may be different:

- **Central value  $\neq 0$** : some data feature differs from MC expectation  
⇒ Need investigation if large
- **Uncertainty  $< 1$**  : effect is *constrained* by the data ⇒ Needs checking if this legitimate or a modeling issue

→ **Impact on result** of  $\pm 1\sigma$  shift of NP allows to gauge which NPs matter most .



# Profiling issues

Systematics are described by NPs included in the fit. Define **pull** as

$$(\hat{\theta} - \theta_0) / \sigma_{\theta}$$

Nominally:

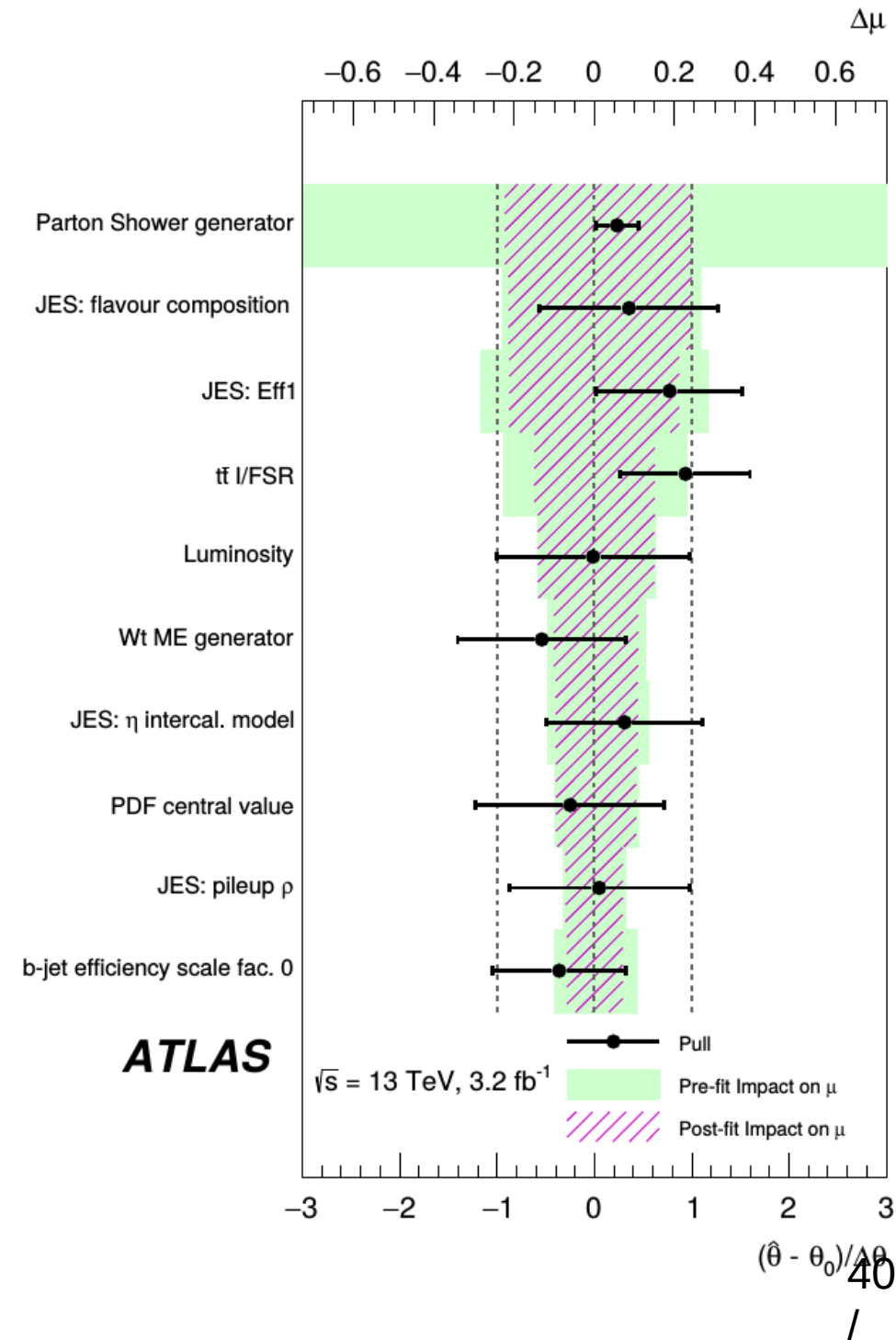
- **Pull = 0** : i.e. the pre-fit expectation
- **Pull uncertainty = 1** : from the Gaussian

However fit results may be different:

- **Central value  $\neq 0$** : some data feature differs from MC expectation  
⇒ Need investigation if large
- **Uncertainty  $< 1$**  : effect is *constrained* by the data ⇒ Needs checking if this legitimate or a modeling issue

→ **Impact on result** of  $\pm 1\sigma$  shift of NP allows to gauge which NPs matter most .

13 TeV single- $t$  XS ([arXiv:1612.07231](https://arxiv.org/abs/1612.07231))



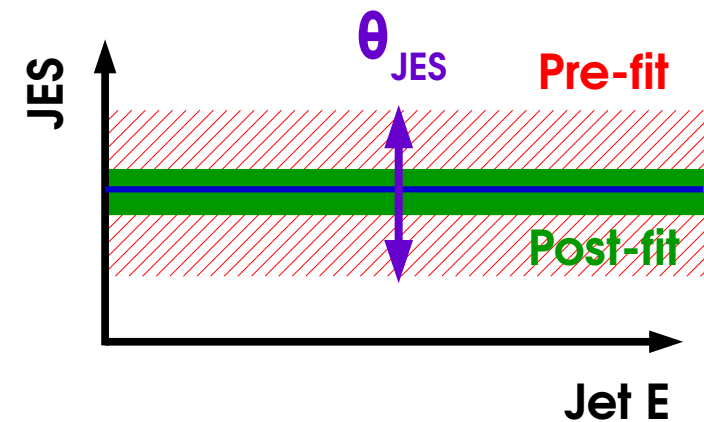


# Profiling issues

**Too simple modeling** can have unintended effects

→ e.g. single Jet E scale parameter:

⇒ Low-E jets calibrate high-E jets – intended ?

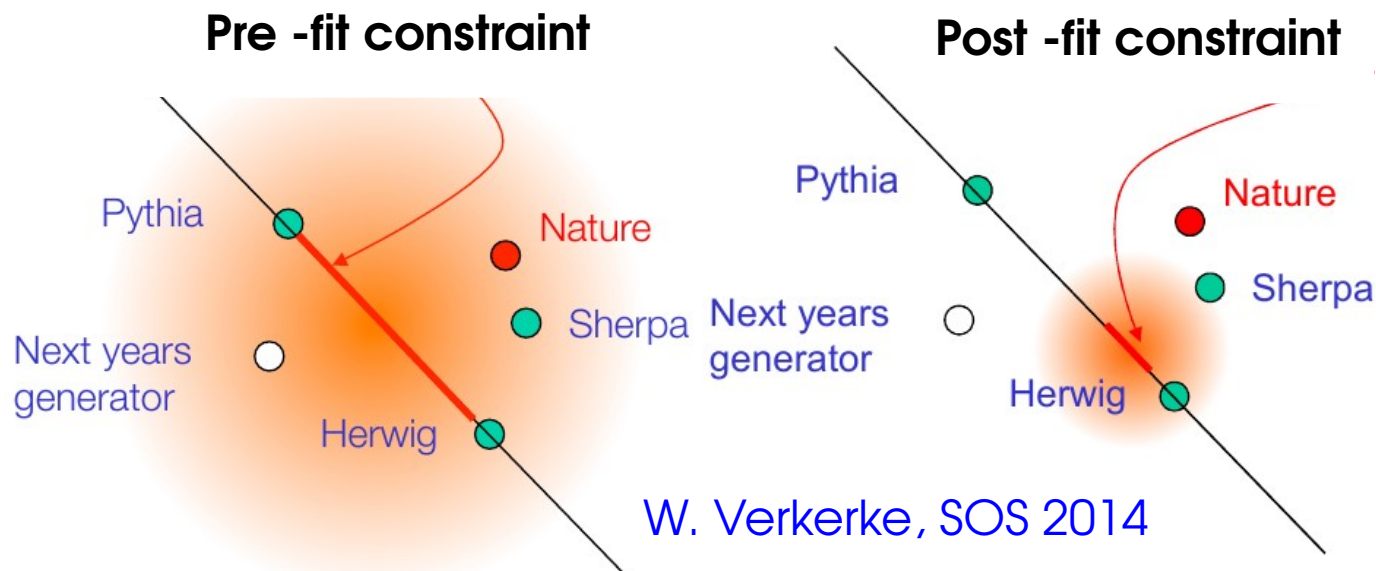


**Two-point uncertainties:**

→ Interpolation may not cover full configuration space

⇒ Can lead to too-strong constraints

**Typical examples:** simulation uncertainties (“PYTHIA vs. HERWIG”)

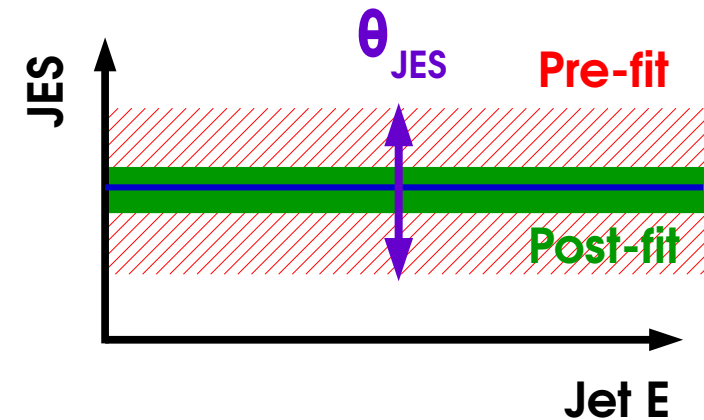


# Profiling issues

**Too simple modeling** can have unintended effects

→ e.g. single Jet E scale parameter:

⇒ Low-E jets calibrate high-E jets – intended ?

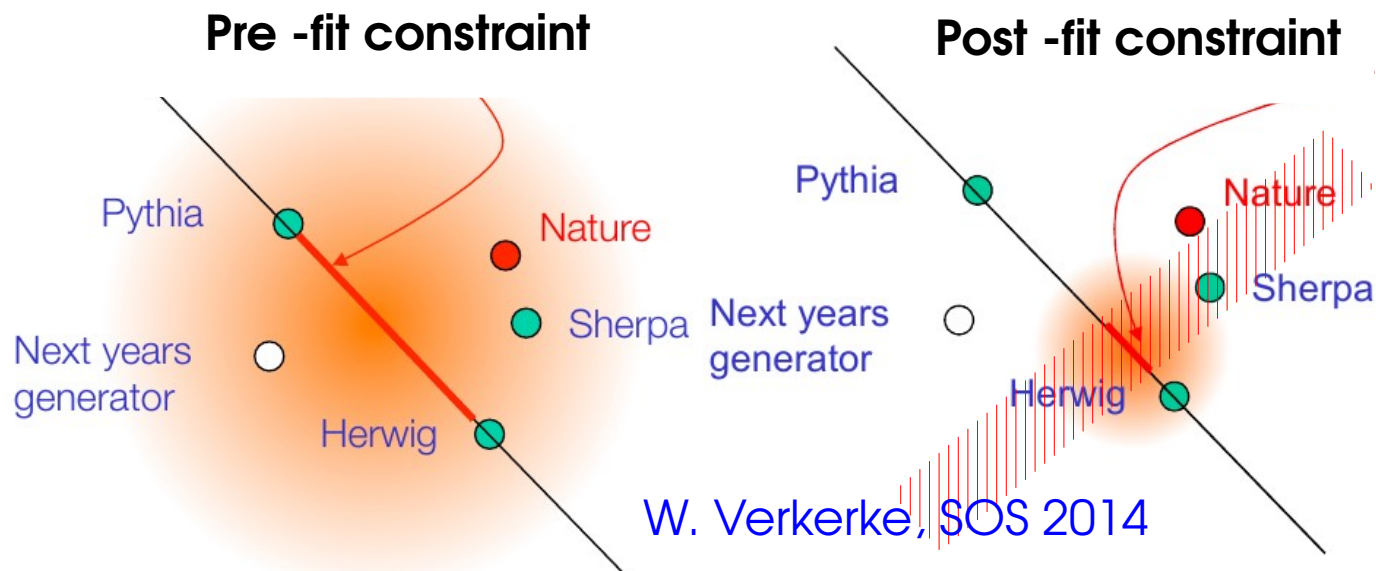


**Two-point uncertainties:**

→ Interpolation may not cover full configuration space

⇒ Can lead to too-strong constraints

**Typical examples:** simulation uncertainties (“PYTHIA vs. HERWIG”)



---

# Systematics

# Impact of Systematics

$$L(\mu, \{\theta_j\}_{j=1 \dots n_{NP}}; \{n_i^{(k)}\}_{i=1 \dots n_{data}^{(k)}}^{k=1 \dots n_{cat}}, \{\theta_j^{obs}\}_{j=1 \dots n_{NP}}) =$$

Expected  
bin yield

$$\prod_{k=1}^{n_{cat}} P[n_i; \mu \epsilon_{i,k}(\vec{\theta}) N_{S,i,k}(\vec{\theta}) + B_{i,k}(\vec{\theta})] \prod_{j=1}^{n_{syst}} G(\theta_j^{obs}; \theta_j; 1)$$

Bin Yields or  
Observable  
values

POI

NPs

Systematics

Sig/Bkg Shapes,  
efficiencies

Pseudo-  
experiments

Data

MC

Auxiliary  
Data

× number of categories!

# Impact of Systematics

$$L(\mu, \{\theta_j\}_{j=1 \dots n_{NP}}; \{c n_i^{(k)}\}_{i=1 \dots n_{data}^{(k)}}, \{e \theta_j^{obs}\}_{j=1 \dots n_{NP}}) =$$

Expected  
bin yield

$$\prod_{k=1}^{n_{cat}} P[n_i; \mu, \epsilon_{i,k}(\vec{\theta}) N_{S,i,k}(\vec{\theta}) + B_{i,k}(\vec{\theta})] \prod_{j=1}^{n_{syst}} G(\theta_j^{obs}; \theta_j; 1)$$

Bin Yields or  
Observable  
values

POI

Key ingredient: impact of systematics  
on signal and background yields

$$N_{S,i,k}(\theta_j) = N_{S,i,k}^0 \prod_j (1 + \delta_{i,j,k} \theta_j)$$

Systematics

Auxiliary  
Data

Pseudo-  
experiments

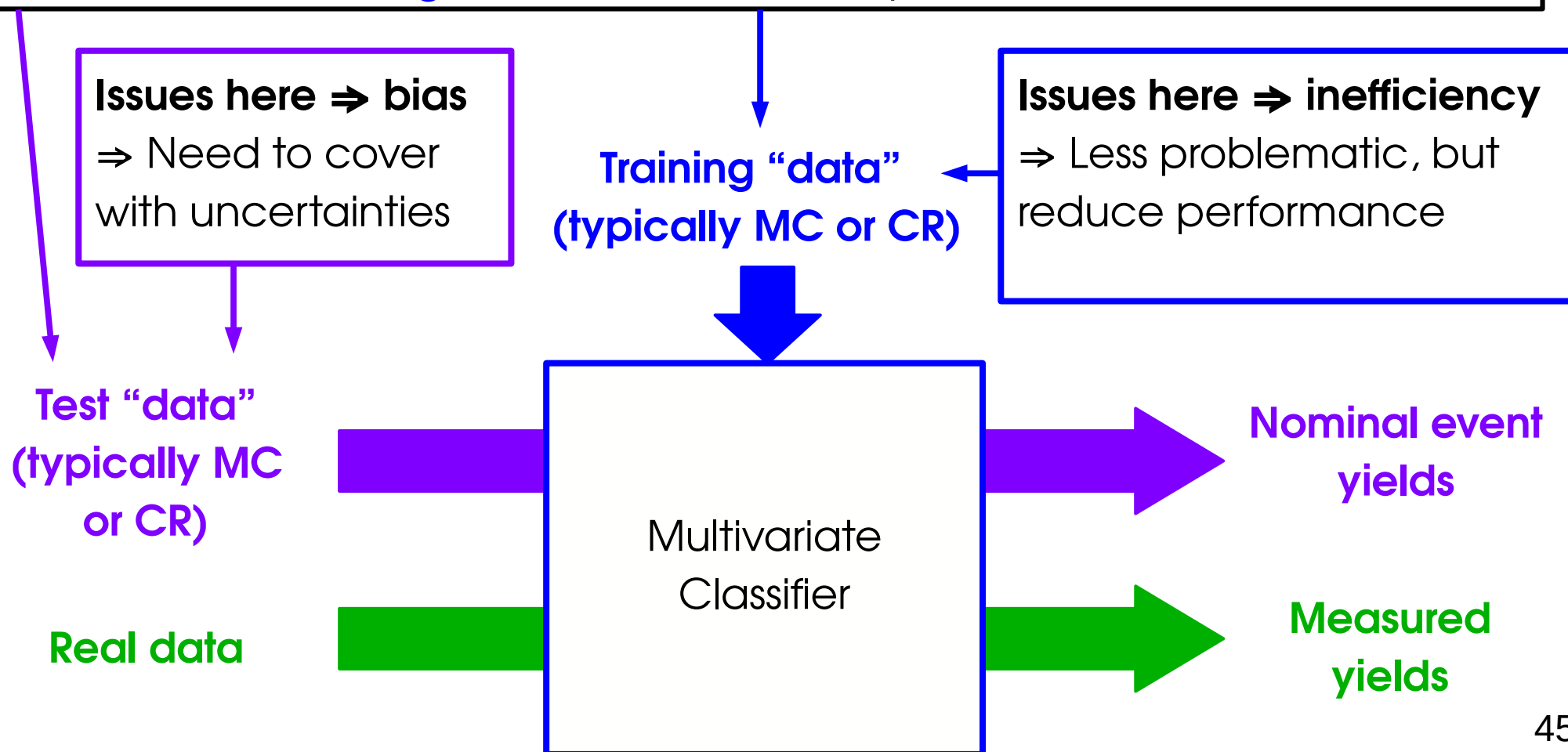
Data

MC

× number of categories!

## Effects to account for:

- **Training/test data not representative of real data (mismodeling)**  
→ Typically covered by MC systematics: variations in the MC sample describing a range of possible models.
- **Limited size of training dataset** → Covered by “MC stat” uncertainties



# Randomness in High-Energy Physics

Experimental data is produced by incredibly complex processes

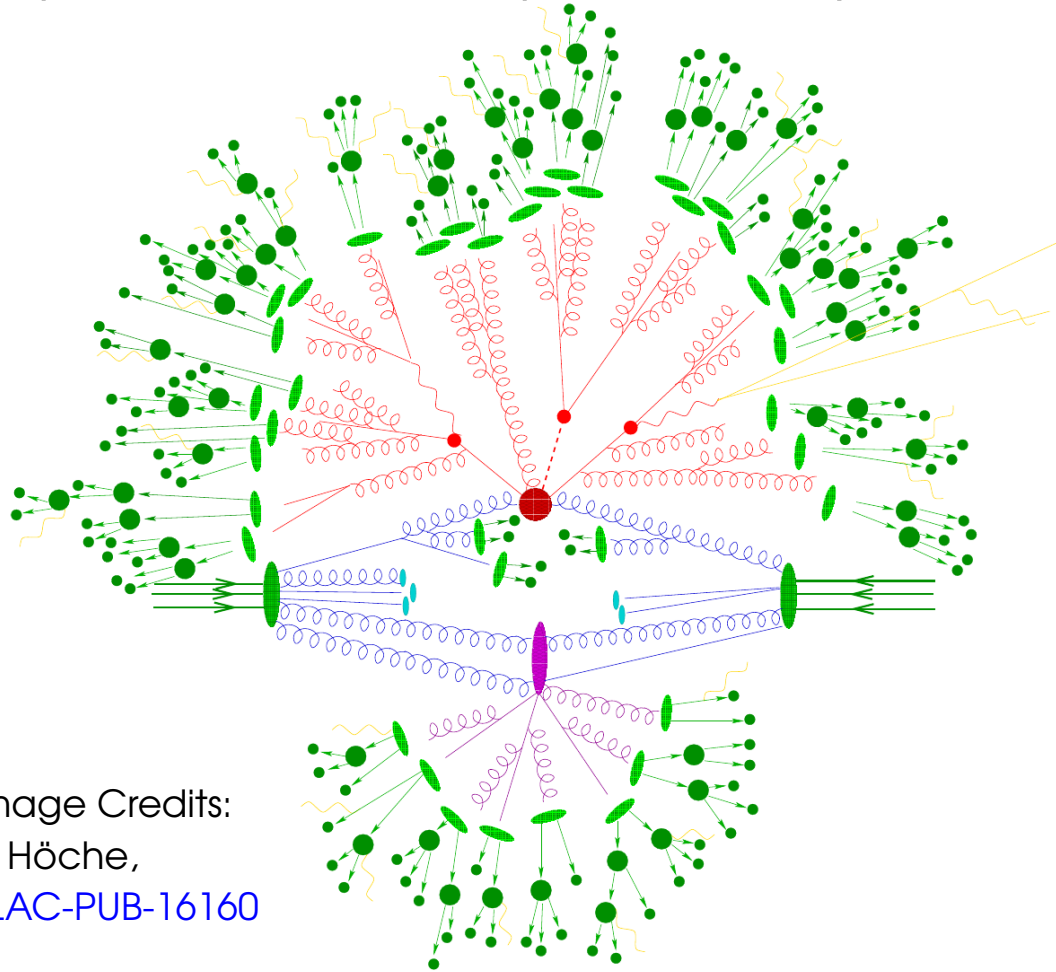


Image Credits:  
S. Höche,  
[SLAC-PUB-16160](#)

**Randomness** involved in all stages

→ **Classical** randomness: detector response

→ **Quantum** effects in particle production, decay

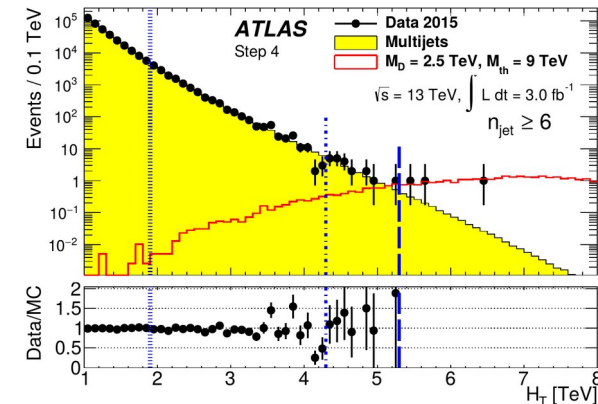
Hard scattering

PDFs, Parton shower, Pileup

Decays

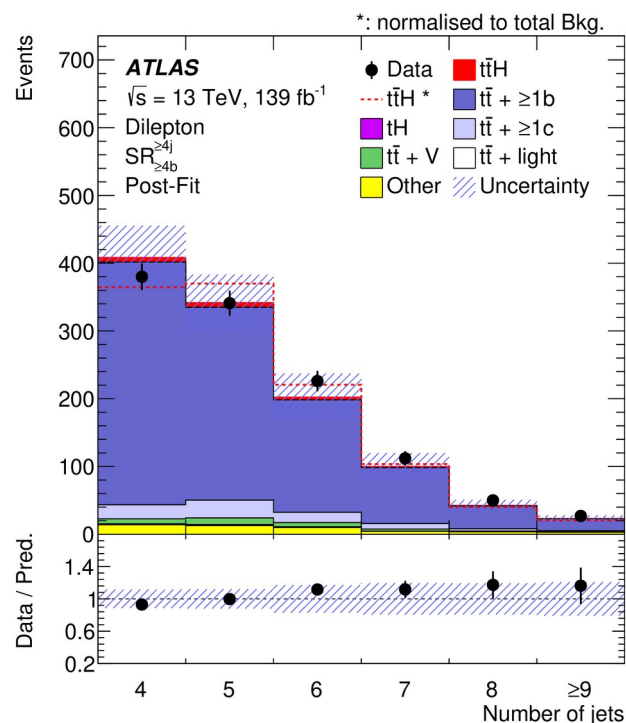
Detector response

Reconstruction





# Modeling Systematics



**Some distributions not predicted with sufficient accuracy:**

- MC modeling
- Detector response
- CR statistics, CR  $\rightarrow$  SR extrapolation

**Error band:** combination of above  
 Typically described by many NPs

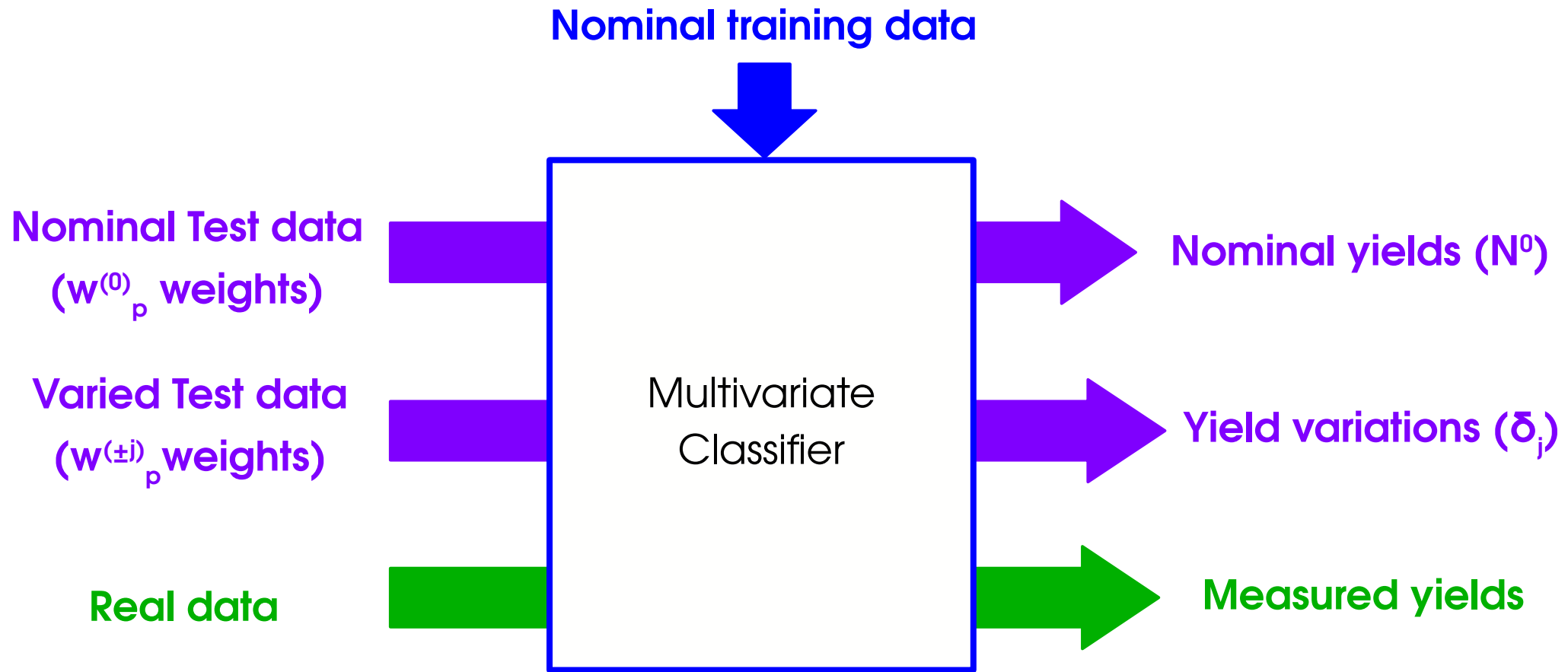
**Modeling variations typically implemented through event weights:**

- **Nominal modeling**  $\rightarrow$  nominal event weight  $w_p^{(0)}$ .
- **Each variation  $\theta_j = \pm 1$**   $\rightarrow$  associated event weight  $w_p^{(\pm j)}$ .

Distributions for each case obtained by applying the appropriate weights.

Ultimately, need impact on yields: 
$$N_{s,i}(\theta_j) = N_{s,i}^0 \prod_j (1 + \delta_{i,j} \theta_j)$$

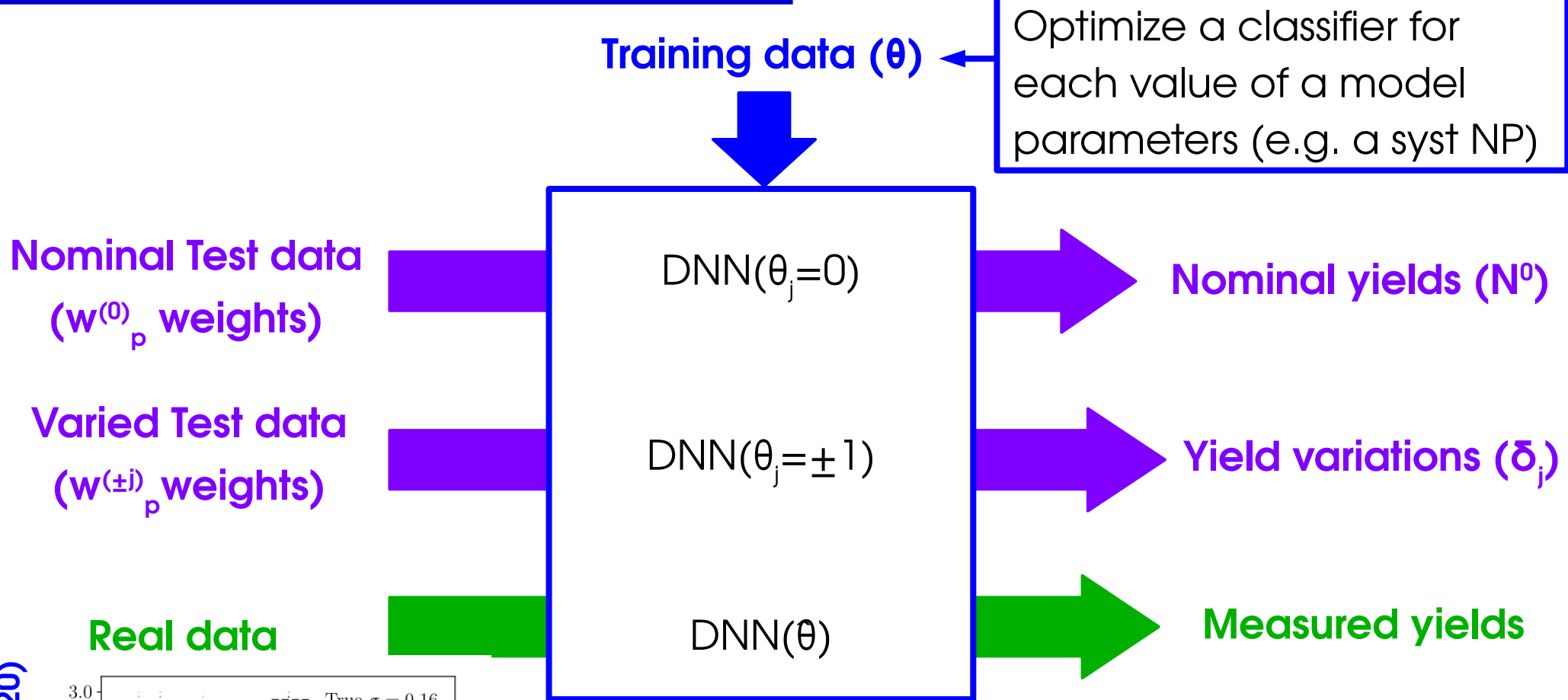
# Treating ML Systematics



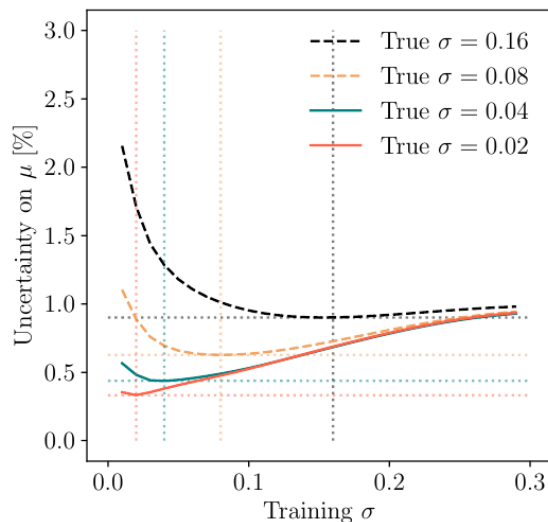
- **“Propagate uncertainties through the DNN”**
- MC stat uncertainties can be treated similarly using resampling  
→ Allows to properly cover for uncertainties, but optimal performance only in nominal case (since used in training).

# Parameterized classifiers

Eur. Phys. J. C (2016) 76:235



- ⊕ Use the optimal classifier for each NP value  
⇒ Retain optimal performance in each case
- ⊖ Scaling with  $|\theta|$  ?  
(in practice  $|\theta|$  can be  $10^{2-3}$  ...)



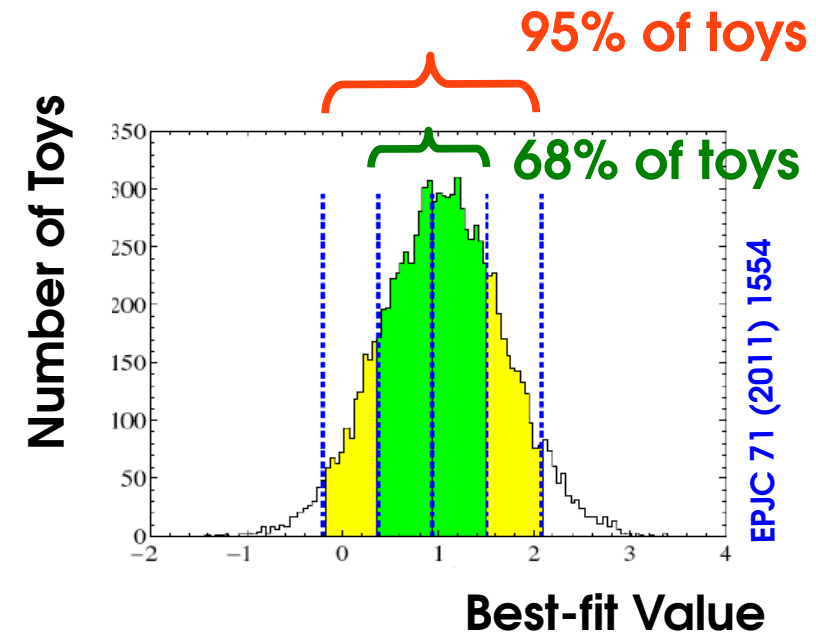
# Brute force approach

Generate pseudo-experiments (“toys”) and repeat best-fit for each case

→ **Statistics**: resample observed dataset

→ **Systematics**: randomize auxiliary obs.  $\theta_j^{obs}$

Obtain intervals from quantiles of the distribution of results



$$\prod_{k=1}^{n_{cat}} P\left[n_i; \mu \epsilon_{i,k}(\vec{\theta}) N_{S,i,k}(\vec{\theta}) + B_{i,k}(\vec{\theta})\right] \prod_{j=1}^{n_{syst}} G(\theta_j^{obs}; \theta_j; 1)$$

⊕ No reliance on asymptotic formulas

- ⊖ High CPU requirements (need a fit for each of  $O(1000)$  toys)
- ⊖ As before, changing syst NPs  $\Rightarrow$  non-optimal classifier performance
- ⊖ **Optimal case**: need to retrain classifier for each toy ?

# Other approaches

---

- **Inference-aware NN** (De Castro, Dorigo, [Com. Phys. Comm. 244 \(2019\), 170-179](#))  
→ Design a NN to directly minimize the width of the confidence interval on the target POI
- **Likelihood-free inference** (Cranmer, Pavez, Louppe, [arXiv:506.02169](#)).
  - Typically, trained classifiers asymptotically learn the likelihood ratio  $p(x|S)/p(x|B)$ , e.g. when using cross-entropy loss.
  - Parameterized classifiers can estimate POIs without computing  $L$ .⇒ Bypass the profile likelihood construction, get intervals from toys ?
- ... ?

---

**(Further) Discussion,  
Questions, Comments ?**

---

# Backup



# Collider processes

**HEP** : Poisson approximation almost always valid:

**ATLAS** :

- **Event rate  $\sim 1$  GHz**  
( $L \sim 10^{34} \text{ cm}^{-2}\text{s}^{-1} \sim 10 \text{ nb}^{-1}/\text{s}$ ,  $\sigma_{\text{tot}} \sim 10^8 \text{ nb}$ , )

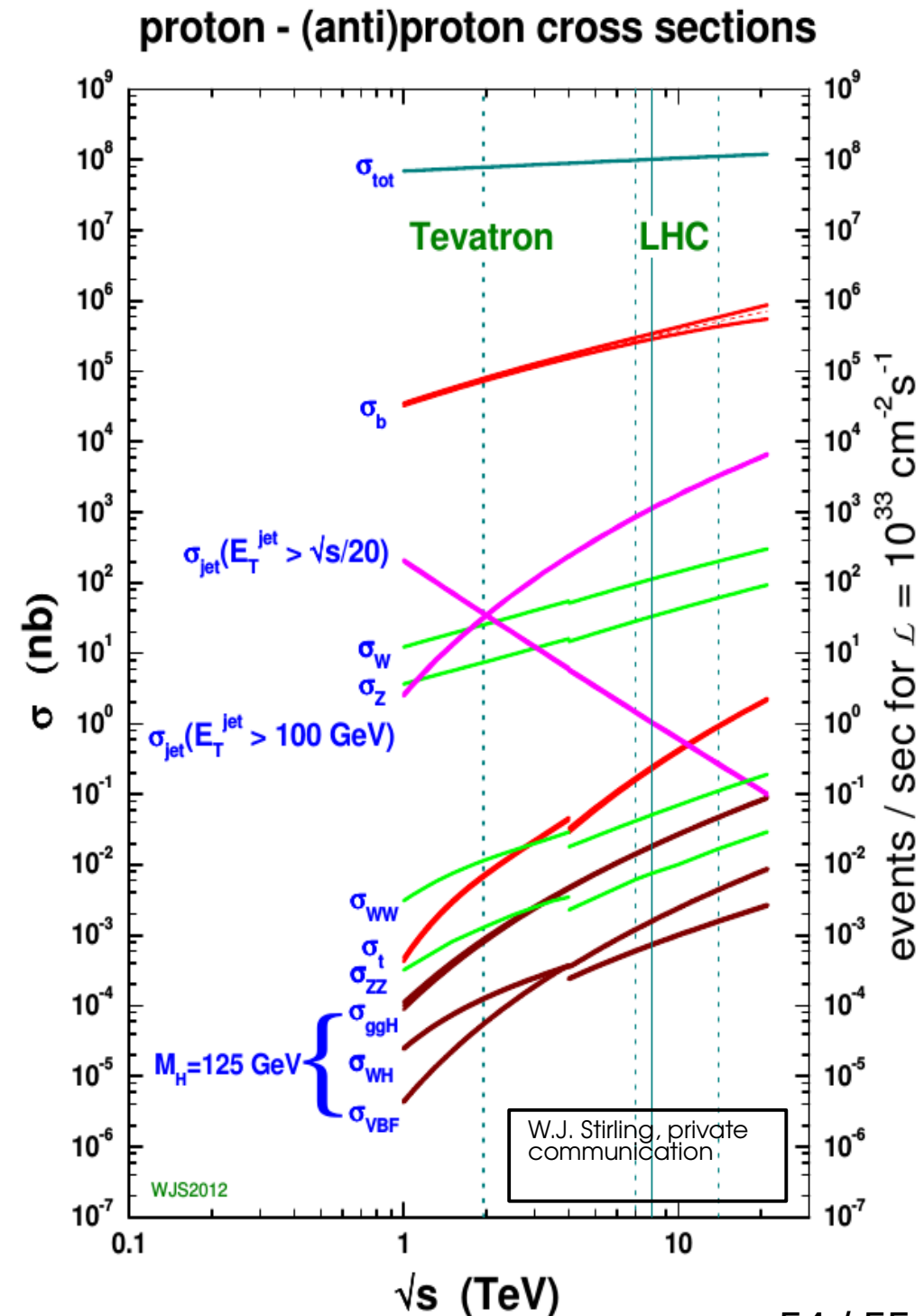
- **Trigger rate  $\sim 1$  kHz**  
(Higgs rate  $\sim 0.1$  Hz)

$\Rightarrow p \sim 10^{-6} \ll 1$  ( $p_{H \rightarrow \gamma\gamma} \sim 10^{-13}$ )

A day of data:  **$N \sim 10^{14} \gg 1$**

$\Rightarrow$  Poisson regime! Similarly true in many other physics situations.

Large  $N$  = design requirement, to get not-too-small  $\lambda = Np$ ...



# Bayesian methods

**Probability distribution** (= likelihood) :

→ Same as frequentist case, but treat systematics by **marginalization**, i.e. **integrating over priors**, instead of profiling:

→ Integrate out  $\theta$  to get  $P(\mu)$  : 
$$P(\mu) = \int P(\mu, \theta) C(\theta) d\theta$$

→ Use probability distribution  $P(\mu)$  directly for limits & intervals

e.g. 68% CL (“Credibility Level”) interval  $[A, B]$  is: 
$$\int_A^B P(\mu) \pi(\mu) d\mu = 68\%$$

where  $\pi(\mu)$  is the prior on  $\mu$ . Uses **Bayes’ Theorem**: 
$$P(\mu | n) = P(n | \mu) \frac{P(\mu)}{P(n)}$$

- ➊ No simple way to test for discovery
- ➋ Integration over NPs can be CPU-intensive (but can use MCMC methods)

**Priors** : most analyses use flat priors in the analysis variable(s)

⇒ **Parameterization-dependent**: if flat in  $\sigma \times B$ , then not flat in couplings....

→ Can use the Jeffreys’ or reference priors, but difficult in practice