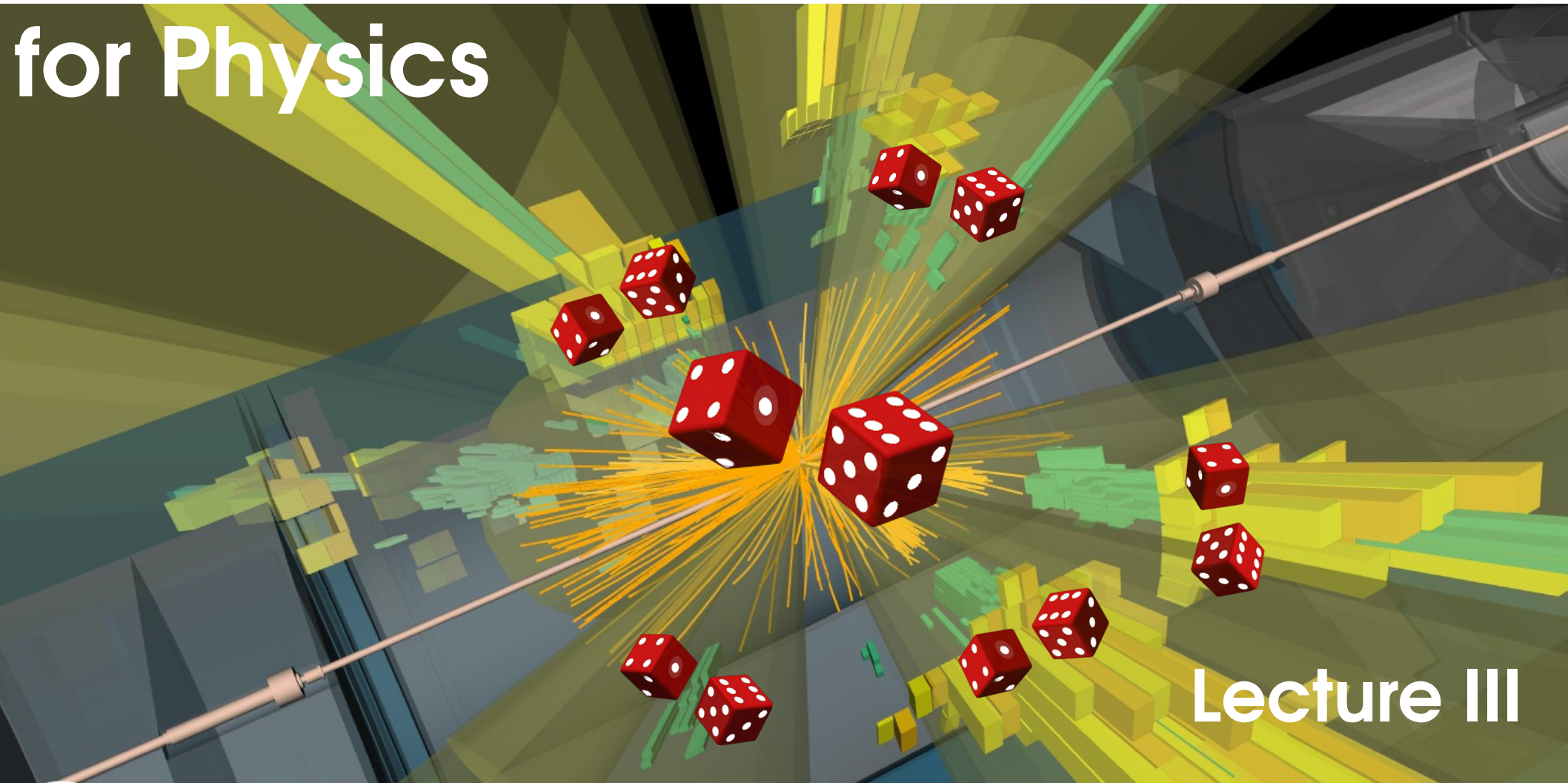
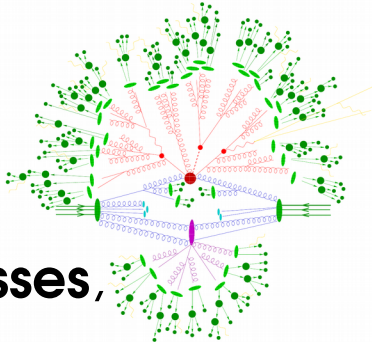

Statistical analysis methods for Physics



Lecture III

Nicolas Berger (LAPP Annecy)

Reminders From Lecture I



Physics measurement data are produced through **random processes**,
 Need to be described using a statistical model:

Description	Observable	Likelihood
Counting	n	Poisson $P(n; S, B) = e^{-(S+B)} \frac{(S+B)^n}{n!}$
Binned shape analysis	$n_i, i=1..N_{bins}$	Poisson product $P(n_i; S, B) = \prod_{i=1}^{n_{bins}} e^{-(S f_i^{sig} + B f_i^{bkg})} \frac{(S f_i^{sig} + B f_i^{bkg})^{n_i}}{n_i!}$
Unbinned shape analysis	$m_i, i=1..n_{evts}$	Extended Unbinned Likelihood $P(m_i; S, B) = \frac{e^{-(S+B)}}{n_{evts}!} \prod_{i=1}^{n_{evts}} S P_{sig}(m_i) + B P_{bkg}(m_i)$

Model can include multiple **categories**, each with a separate description
 Includes **parameters of interest** (POIs) but also **nuisance parameters** (NPs)

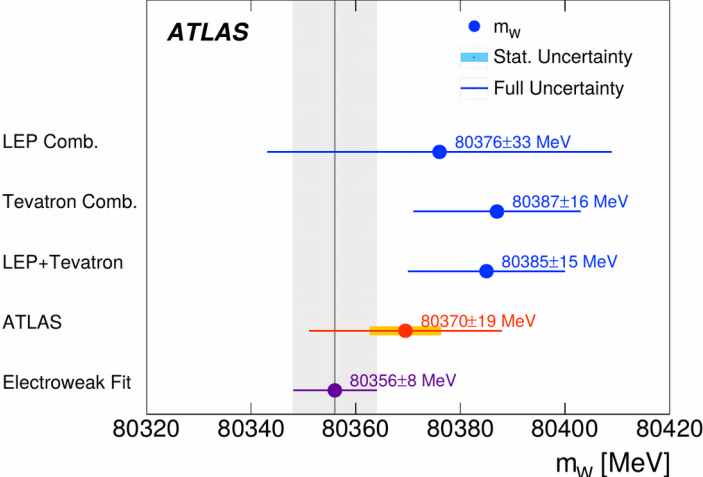
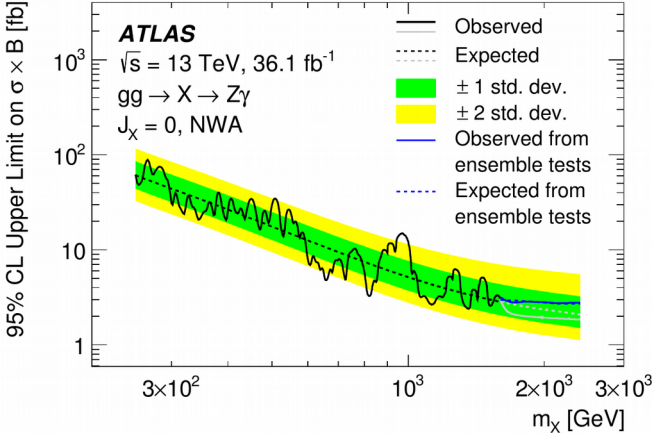
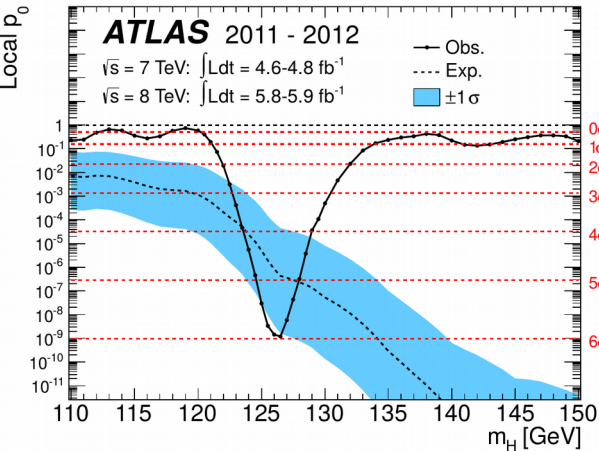
Reminders From Lecture I

To **estimate a parameter value**, use the Maximum-likelihood estimate (MLE), a.k.a. Best-fit value of the parameter,

Today, further results:




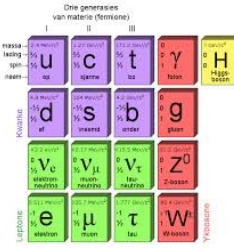
- **Discovery:** we see an excess – is it a (new) signal, or a background fluctuation ?
- **Upper limits:** we don't see an excess – if there is a signal present, how small must it be ?
- **Parameter measurement:** what is the allowed range (“confidence interval”) for a model parameter ?

→ The Statistical Model already contains all the needed information – how to use it ?

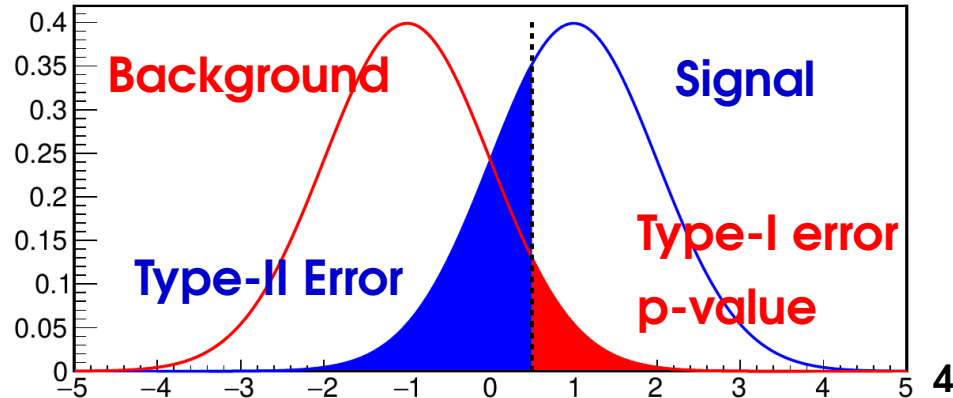


Reminders from Lecture II: Hypothesis Testing

Hypothesis: assumption on model parameters, say value of S (e.g. $H_0 : \mathbf{S}=\mathbf{0}$)
 → Goal : determine if H_0 is true or false using a test based on the data




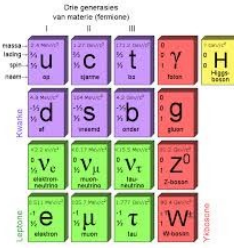
Possible outcomes:	Data disfavors H_0 (Discovery claim)	Data favors H_0 (Nothing found)
H_0 is false (New physics!)	Discovery! 	Missed discovery Type-II error (1 - Power) 
H_0 is true (Nothing new)	False discovery claim Type-I error (→ p-value, significance) 	No new physics, none found 

Stringent discovery criteria
 ⇒ lower Type-I errors, higher Type-II errors
 → Goal: test that minimizes Type-II errors for given level of Type-I error.

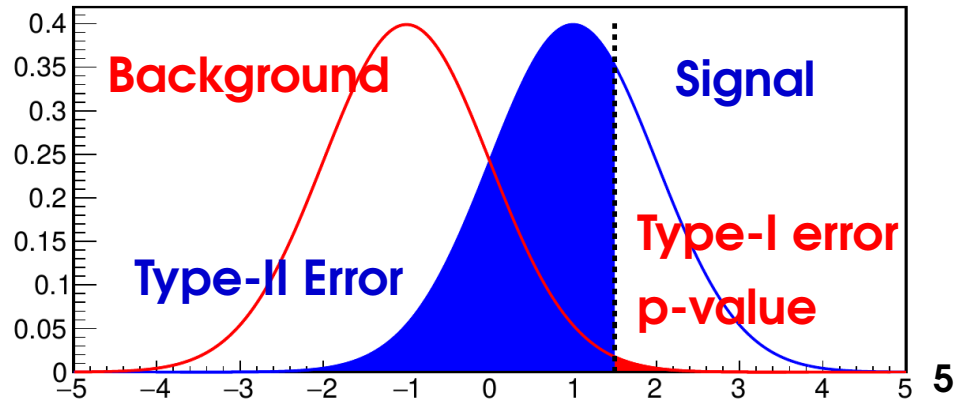


Reminders from Lecture II: Hypothesis Testing

Hypothesis: assumption on model parameters, say value of S (e.g. $H_0 : S=0$)
 → Goal : determine if H_0 is true or false using a test based on the data

Possible outcomes:	Data disfavors H_0 (Discovery claim)	Data favors H_0 (Nothing found)
H_0 is false (New physics!)	Discovery! 	Missed discovery Type-II error (1 - Power) 
H_0 is true (Nothing new)	False discovery claim Type-I error (→ p-value, significance) 	No new physics, none found 

Stringent discovery criteria
 ⇒ lower Type-I errors, higher Type-II errors
 → Goal: test that minimizes Type-II errors for given level of Type-I error.



Reminders from Lecture II: Discovery Significance

Given a statistical model $P(\text{data}; \mu)$, define likelihood $L(\mu) = P(\text{data}; \mu)$

To estimate a parameter, use value $\hat{\mu}$ that maximizes $L(\mu)$.

To decide between hypotheses H_0 and H_1 , use the likelihood ratio $\frac{L(H_0)}{L(H_1)}$

To test for **discovery**, use $q_0 = \begin{cases} -2 \log \frac{L(S=0)}{L(\hat{S})} & \hat{S} \geq 0 \\ +2 \log \frac{L(S=0)}{L(\hat{S})} & \hat{S} < 0 \end{cases}$

For large enough datasets ($n > 5$), $Z = \sqrt{q_0}$

For a **Gaussian** measurement, $Z = \frac{\hat{S}}{\sqrt{B}}$

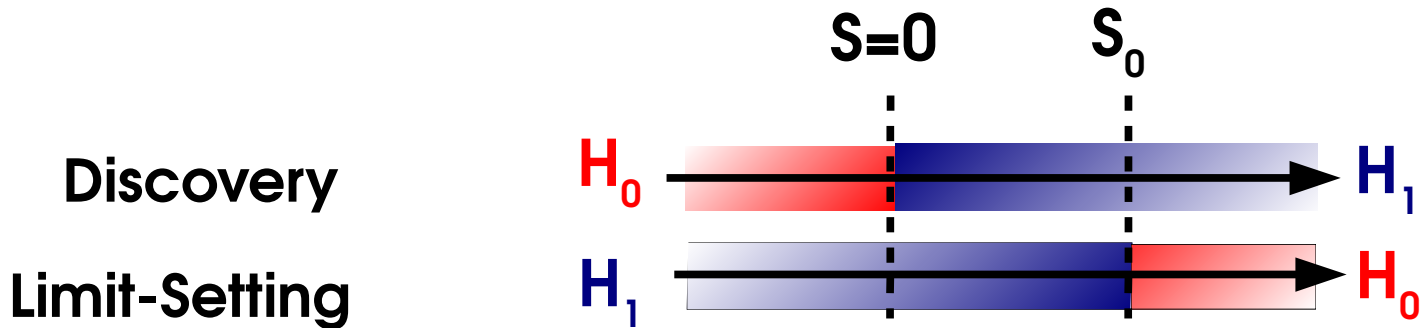
For a **Poisson** measurement, $Z = \sqrt{2 \left[(\hat{S} + B) \log \left(1 + \frac{\hat{S}}{B} \right) - \hat{S} \right]}$

Reminders from Lecture II: Test Statistic for Limits

For upper limits, alternate is $H_1 : S < \mu_0$:

→ If **large** signal observed ($\hat{S} \gg S_0$), does not favor H_1 over H_0

→ Only consider $\hat{S} < S_0$ for H_1 , and include $\hat{S} \geq S_0$ in H_0 .



⇒ Set $q_{s_0} = 0$ for $\hat{S} > S_0$ – only small signals ($\hat{S} < S_0$) help lower the limit.

→ Also treat separately the case $S < 0$ to avoid technical issues in $-2\log L$ fits.

Asymptotics:

$q_{s_0} \sim \text{“}1/2\chi^2\text{”}$ under $H_0(S=S_0)$, same as q_0 , except for special treatment of $\hat{S} < 0$.

$$\tilde{q}_{s_0} = \begin{cases} 0 & \hat{S} \geq S_0 \\ -2 \log \frac{L(S=S_0)}{L(\hat{S})} & 0 \leq \hat{S} \leq S_0 \\ -2 \log \frac{L(S=S_0)}{L(S=0)} & \hat{S} < 0 \end{cases}$$

$$p_0 = 1 - \Phi\left(\sqrt{q_{s_0}}\right)$$

Reminders from Lecture II: Limit Inversion

Procedure

→ Consider $H_0 : H(S=S_0)$ – alternative $H_1 : H(\hat{S} < S_0)$

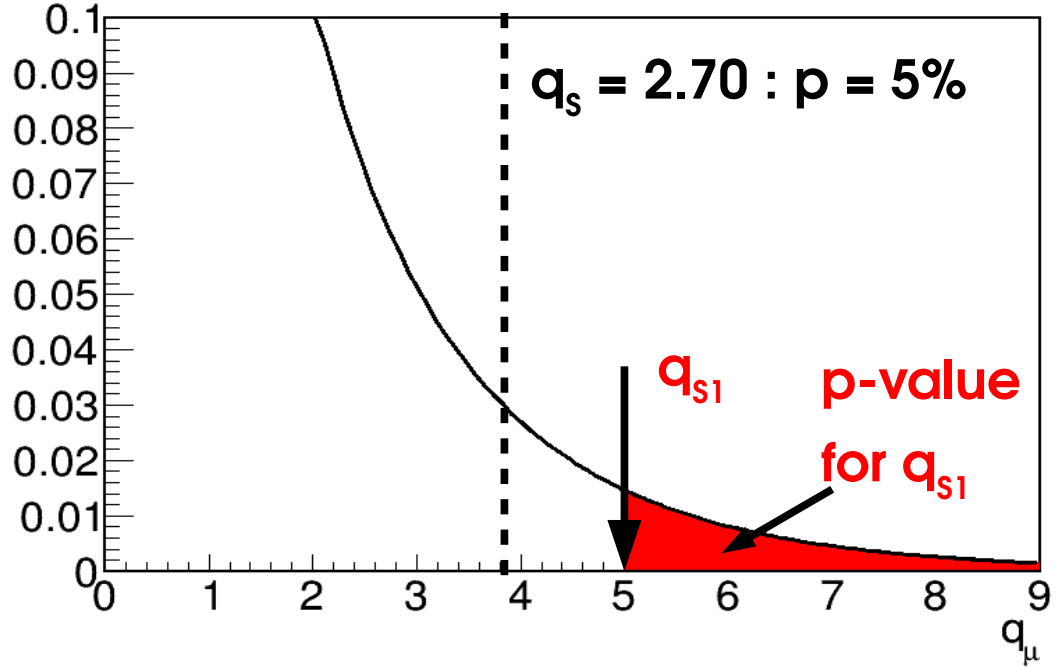
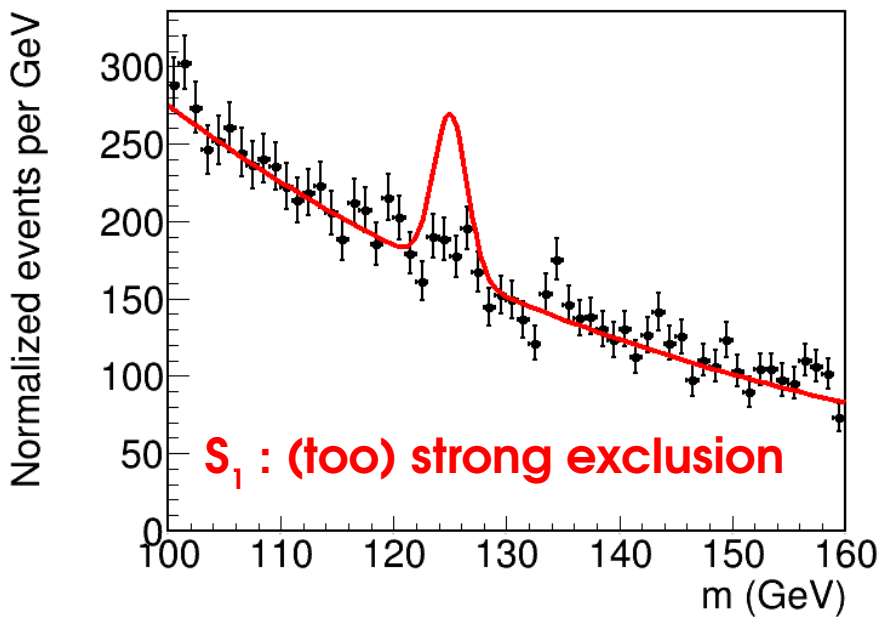
→ Compute q_{S_0} , get **exclusion p-value** p_{S_0}

→ **Adjust S_0 until 95% CL exclusion ($p_{S_0} = 5\%$) is reached**

Asymptotics: set target in terms of $q_{S_0} : \sqrt{q_{S_0}} = \Phi^{-1}(1 - p_0)$

Asymptotics

CL	Region
90%	$q_s > 1.64$
95%	$q_s > 2.70$
99%	$q_s > 5.41$



Reminders from Lecture II: Limit Inversion

Procedure

→ Consider $H_0 : H(S=S_0)$ – alternative $H_1 : H(\hat{S} < S_0)$

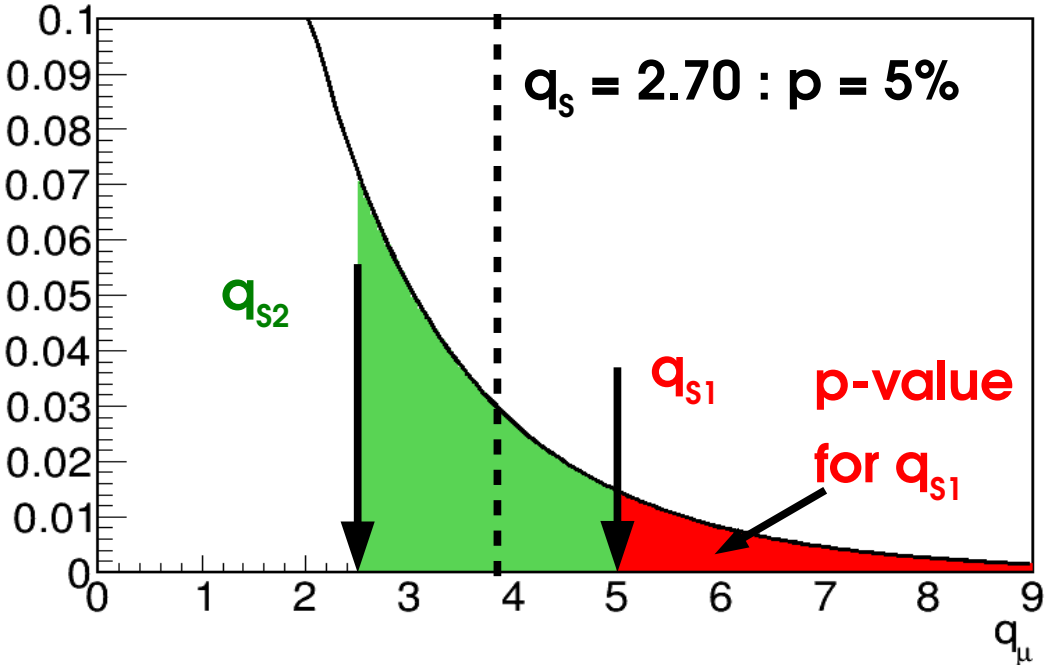
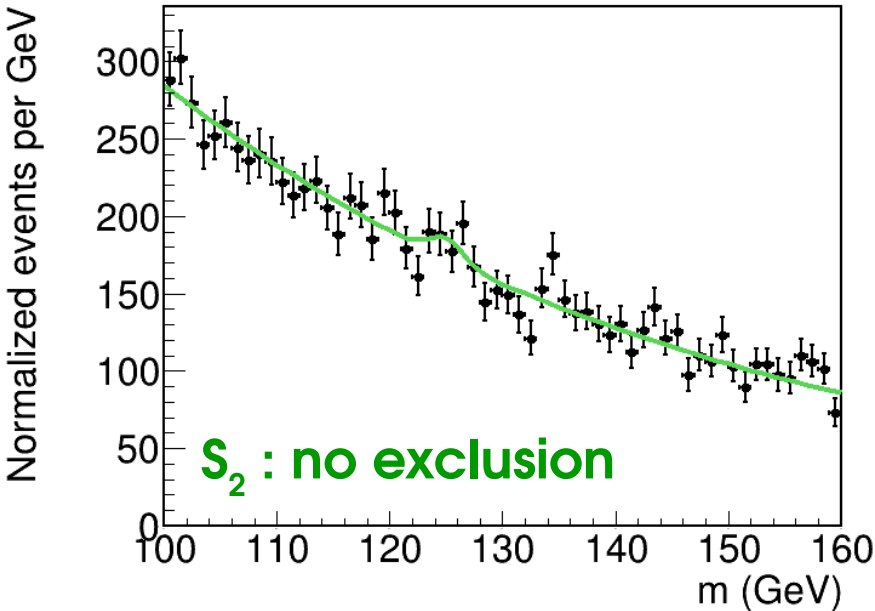
→ Compute q_{S_0} , get **exclusion p-value** p_{S_0}

→ **Adjust S_0 until 95% CL exclusion ($p_{S_0} = 5\%$) is reached**

Asymptotics: set target in terms of $q_{S_0} : \sqrt{q_{S_0}} = \Phi^{-1}(1 - p_0)$

Asymptotics

CL	Region
90%	$q_s > 1.64$
95%	$q_s > 2.70$
99%	$q_s > 5.41$



Reminders from Lecture II: Limit Inversion

Procedure

→ Consider $H_0 : H(S=S_0)$ – alternative $H_1 : H(\hat{S} < S_0)$

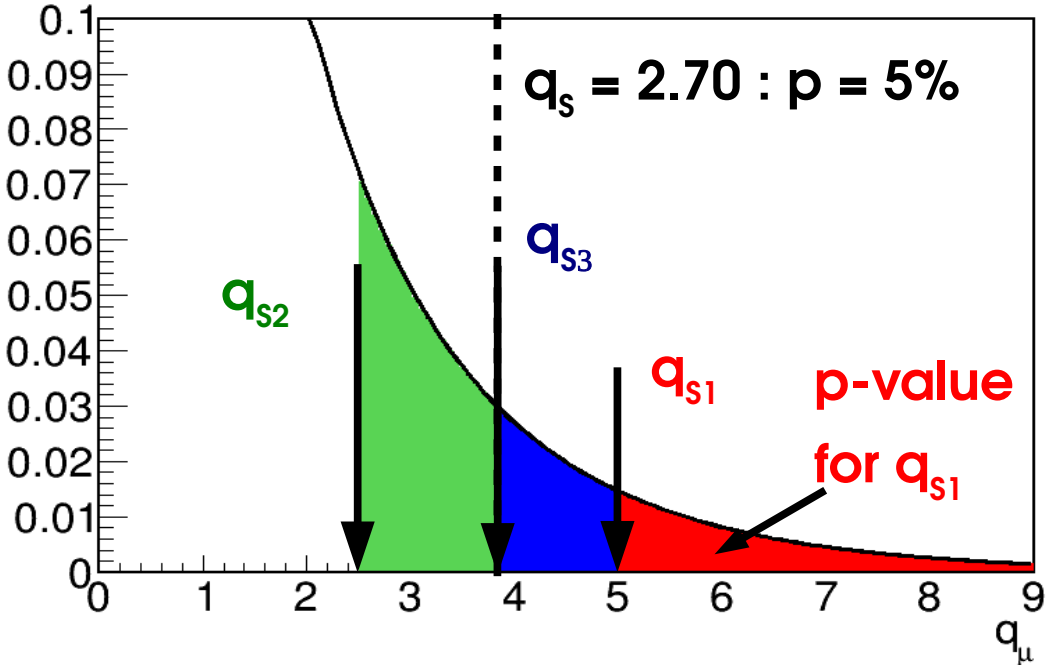
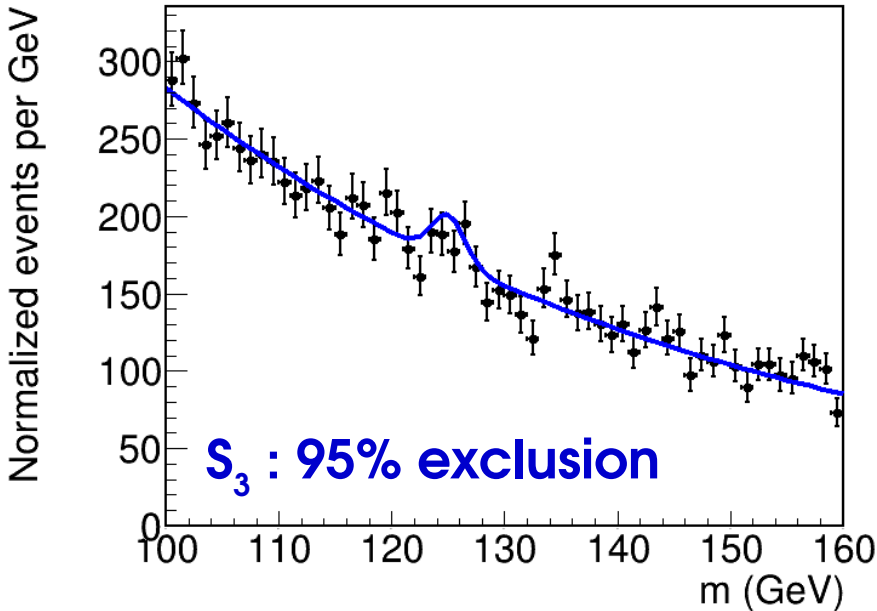
→ Compute q_{S_0} , get **exclusion p-value** p_{S_0}

→ **Adjust S_0 until 95% CL exclusion ($p_{S_0} = 5\%$) is reached**

Asymptotics: set target in terms of $q_{S_0} : \sqrt{q_{S_0}} = \Phi^{-1}(1 - p_0)$

Asymptotics

CL	Region
90%	$q_s > 1.64$
95%	$q_s > 2.70$
99%	$q_s > 5.41$



Reminders from Lecture II: CL_s

How to avoid negative limits ? in HEP, use : CL_s .

→ Compute modified p-value

• p_{S_0} is the usual p-value (5%)

• p_0 is the p-value computed under $H(S=0)$.

⇒ **Rescale** exclusion at S_0 by exclusion at $S=0$.

→ Somewhat ad-hoc, but good properties...

Good case : $p_0 \sim O(1)$

$p_{CL_s} \sim p_{S_0} \sim 5\%$, no change.

Pathological case : $p_0 \ll 1$

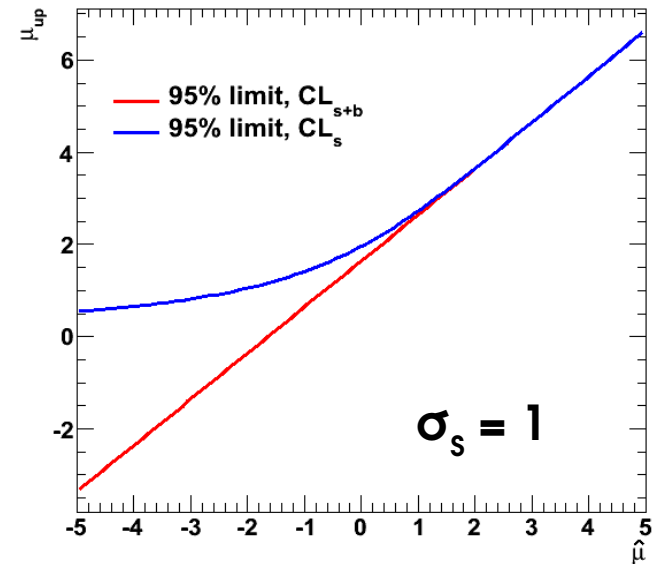
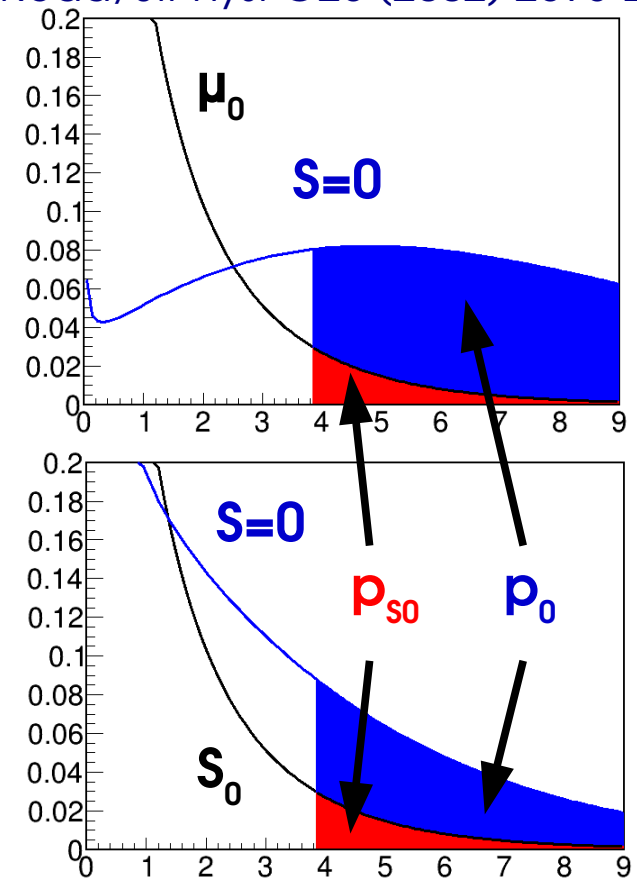
$p_{CL_s} \sim p_{S_0}/p_0 \gg 5\%$

→ no exclusion ⇒ worse limit, usually >0 as desired

Drawback: overcoverage

→ limit is actually $>95\%$ CL for small p_0 .

$$p_{CL_s} = \frac{p_{S_0}}{p_0}$$



Outline

Computing Statistical Results

Limits, continued

Confidence Intervals

Profiling

Look-Elsewhere Effect

Bayesian methods

Statistical modeling in practice

BLUE

CL_s : Gaussian Example

Usual Gaussian counting example with known B:

$$\lambda(S) = \left(\frac{n - (S + B)}{\sigma_S} \right)^2$$

Reminder

Best fit signal : $\hat{S} = n - B$

CL_{s+b} limit: $S_{\text{up}} = \hat{S} + 1.64 \sigma_S$ at 95% CL

CL_s upper limit : still have

so need to solve

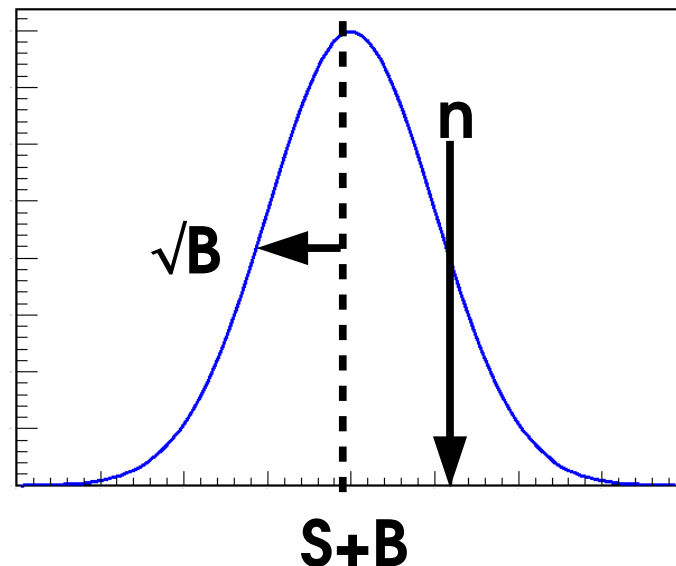
$$q_{S_0} = \left(\frac{S_0 - \hat{S}}{\sigma_S} \right)^2 \quad (\text{for } S_0 > \hat{S})$$

$$p_{CL_s} = \frac{p_{S_0}}{p_0} = \frac{1 - \Phi(\sqrt{q_{S_0}})}{1 - \Phi(\sqrt{q_{S_0}} - S_0/\sigma_S)} = 5\%$$

for $\hat{S} = 0$,

$$S_{\text{up}} = \hat{S} + \left[\Phi^{-1} \left(1 - 0.05 \Phi \left(\hat{S}/\sigma_S \right) \right) \right] \sigma_S \text{ at 95\% CL}$$

$\Phi(0) = 0.5 \Rightarrow$ at 95% CL, **CL_s : $S_{\text{up}} = 1.96 \sigma_S$** **CL_{s+b} : $S_{\text{up}} = 1.64 \sigma_S$**



$\hat{S} \sim G(S, \sigma_S)$ so

Under $H_0(S = S_0)$:

$$\sqrt{q_{S_0}} \sim G(0, 1)$$

$$p_{S_0} = 1 - \Phi(\sqrt{q_{S_0}})$$

Under $H_0(S = 0)$:

$$\sqrt{q_{S_0}} \sim G(S_0/\sigma_S, 1)$$

$$p_0 = 1 - \Phi(\sqrt{q_{S_0}} - S_0/\sigma_S)$$

CL_s: Poisson Rule of Thumb

Same exercise, for the Poisson case

Exact computation : sum probabilities of cases “at least as extreme as data” (n)

$$p_{S_0}(n) = \sum_0^n e^{-(S_0+B)} \frac{(S_0+B)^k}{k!} \quad \text{and one should solve } p_{CL_s} = \frac{p_{S_{up}}(n)}{p_0(n)} = 5\% \text{ for } S_{up}$$

$$\text{For } n=0: \quad p_{CL_s} = \frac{p_{S_{up}}(0)}{p_0(0)} = e^{-S_{up}} = 5\% \Rightarrow S_{up} = \log(20) = 2.996 \approx 3$$

⇒ **Rule of thumb: when $n_{obs}=0$, the 95% CL_s limit is 3 events (for any B)**

$$\text{Asymptotics: as before, } q_{S_0} = \lambda(S_0) - \lambda(\hat{S}) = 2(S_0 + B - n) - 2n \log \frac{S_0+B}{n}$$

$$\text{For } n=0, \quad q_{S_0}(n=0) = 2(S_0+B)$$

$$p_{CL_s} = \frac{p_{S_0}}{p_0} = \frac{1 - \Phi(\sqrt{q_{S_0}(n=0)})}{1 - \Phi(\sqrt{q_{S_0}(n=0)} - \sqrt{q_{S_0}(n=B)})} = 5\%$$

⇒ $S_{up} \sim 2$, exact value depends on B

⇒ Asymptotics not valid in this case (n=0) – need to use exact results, or toys

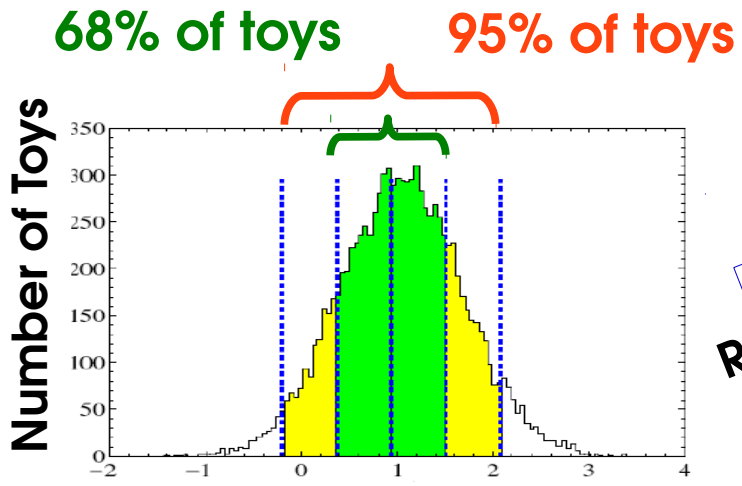
Expected Limits: Toys

Expected results: median outcome under a given hypothesis
 → usually B-only by convention, but other choices possible.

Two main ways to compute:

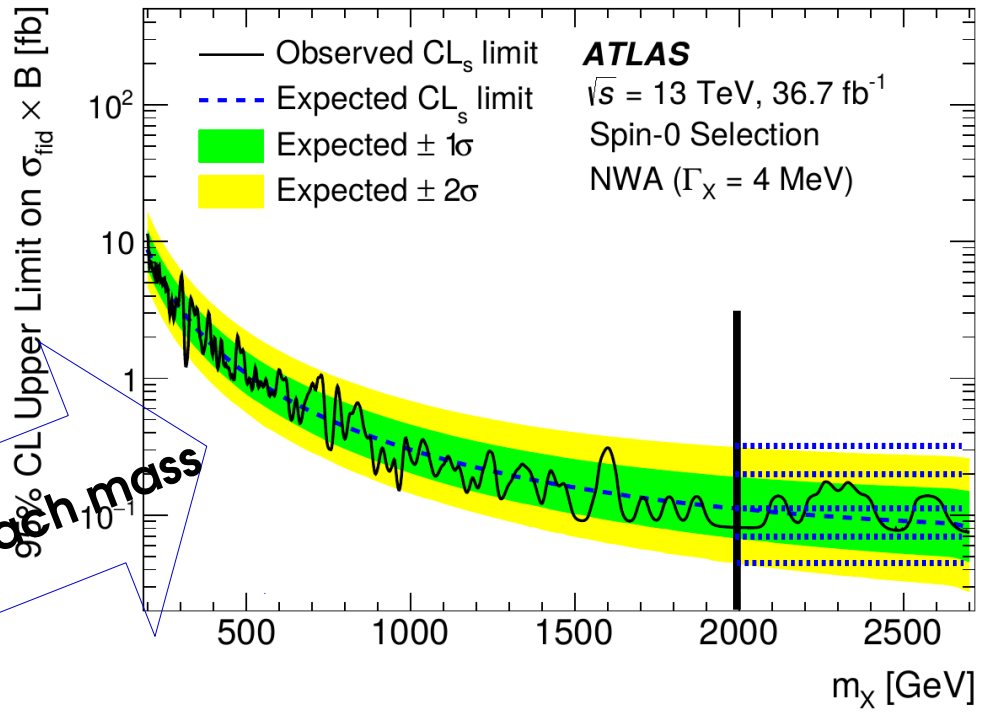
→ **Pseudo-experiments (toys):**

- Generate pseudo-data in B-only hypothesis
- Compute limit
- Repeat and histogram the results
- Central value = median, bands based on quantiles



Repeat for each mass

Phys. Lett. B 775 (2017) 105



Expected Limits: Asimov

Expected results: median outcome under a given hypothesis

→ usually B-only by convention, but other choices possible.

Two main ways to compute:

Strictly speaking, Asimov dataset if
 $\hat{X} = X_0$ for all parameters X ,
where X_0 is the generation value

→ Asimov Datasets

- Generate a “perfect dataset” – e.g. for binned data, set bin contents carefully, no fluctuations.

- Gives the median result immediately:

median(toy results) ↔ result(median dataset)

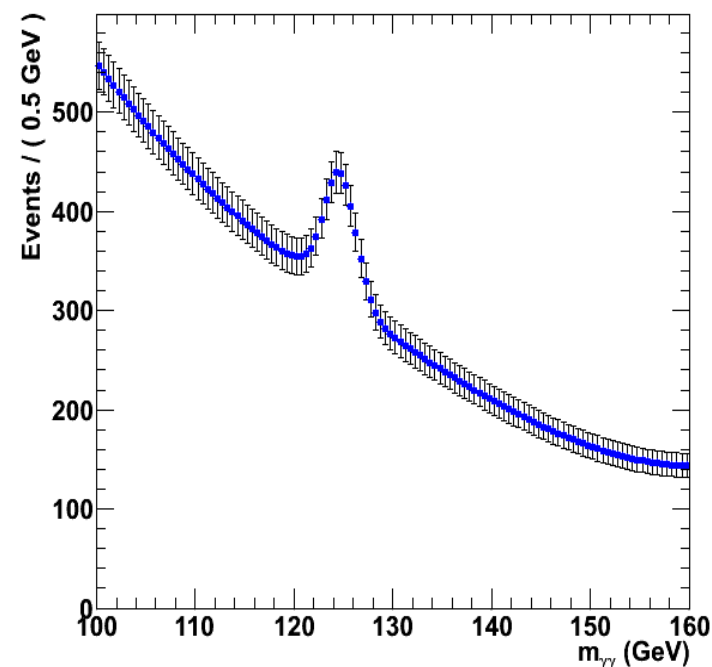
- Get bands from asymptotic formulas:

Band width

$$\sigma_{S_0, A}^2 = \frac{S_0^2}{q_{S_0}(\text{Asimov})}$$

⊕ Much faster (1 “toy”)

⊖ Relies on Gaussian approximation



CL_s : Gaussian Bands

Usual Gaussian counting example with known B:

95% CL_s upper limit on S:

$$S_{\text{up}} = \hat{S} + \left[\Phi^{-1} \left(1 - 0.05 \Phi \left(\hat{S} / \sigma_S \right) \right) \right] \sigma_S \quad \text{with} \quad \sigma_S = \sqrt{B}$$

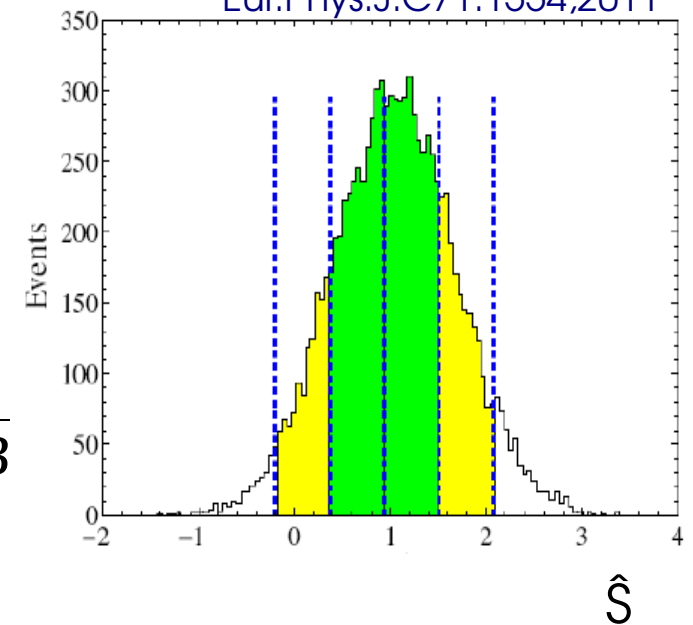
Compute expected bands for S=0:

→ **Asimov dataset** $\Leftrightarrow \hat{S} = 0$:

$$S_{\text{up,exp}}^0 = 1.96 \sigma_S$$

→ **$\pm n\sigma$ bands**:

$$S_{\text{up,exp}}^{\pm n} = \left(\pm n + \left[1 - \Phi^{-1} \left(0.05 \Phi(\mp n) \right) \right] \right) \sigma_S$$



n	$S_{\text{exp}}^{\pm n} / \sqrt{B}$
+2	3.66
+1	2.72
0	1.96
-1	1.41
-2	1.05

CLs :

- Positive bands somewhat reduced,
- Negative ones more so

Band width from $\sigma_{S,A}^2 = \frac{S^2}{q_S(\text{Asimov})}$ depends on S, for non-Gaussian cases, different values for each band...

Outline

Computing Statistical Results

Limits, continued

Confidence Intervals

Profiling

Look-Elsewhere Effect

Bayesian methods

Statistical modeling in practice

BLUE

Gaussian Inversion

If $\hat{\mu} \sim G(\mu^*, \sigma)$, known quantiles :

$$P(\mu^* - \sigma < \hat{\mu} < \mu^* + \sigma) = 68\%$$

This is a probability for $\hat{\mu}$, not μ !

→ μ^* is a **fixed number**, **not a random variable**

But we can invert the relation:

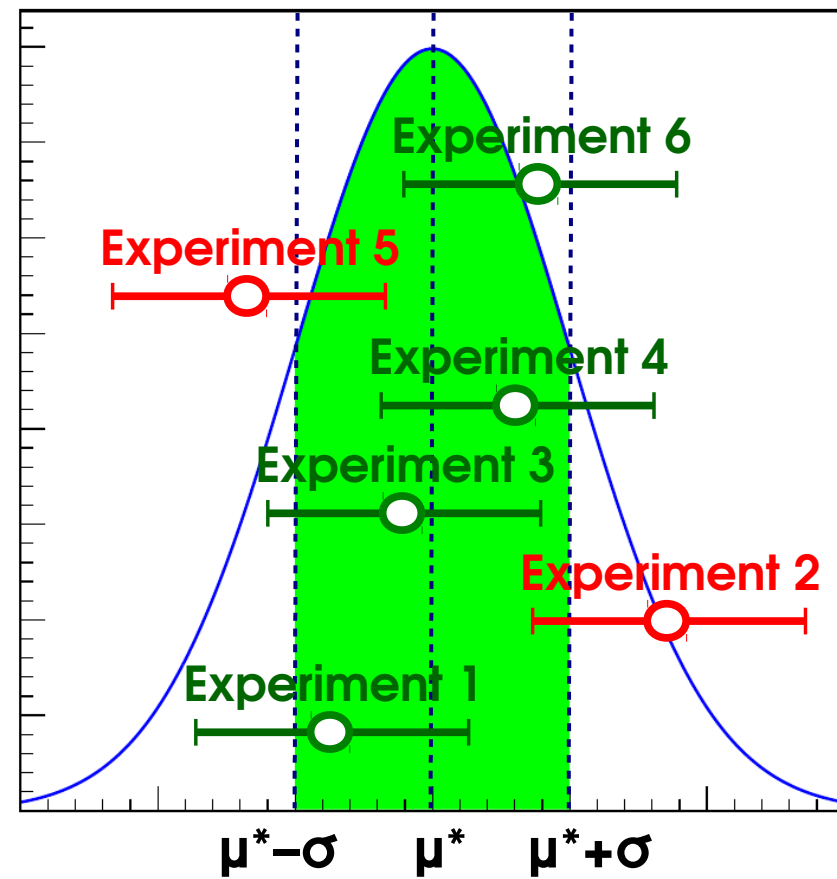
$$P(\mu^* - \sigma < \hat{\mu} < \mu^* + \sigma) = 68\%$$

$$\Rightarrow P(|\hat{\mu} - \mu^*| < \sigma) = 68\%$$

$$\Rightarrow P(\hat{\mu} - \sigma < \mu^* < \hat{\mu} + \sigma) = 68\%$$

→ This gives the desired statement on μ^* : *if we repeat the experiment many times, $[\hat{\mu} - \sigma, \hat{\mu} + \sigma]$ will contain the true value 68% of the time: $\hat{\mu} = \mu^* \pm \sigma$*

This is a statement **on the interval $[\hat{\mu} - \sigma, \hat{\mu} + \sigma]$** obtained for each experiment



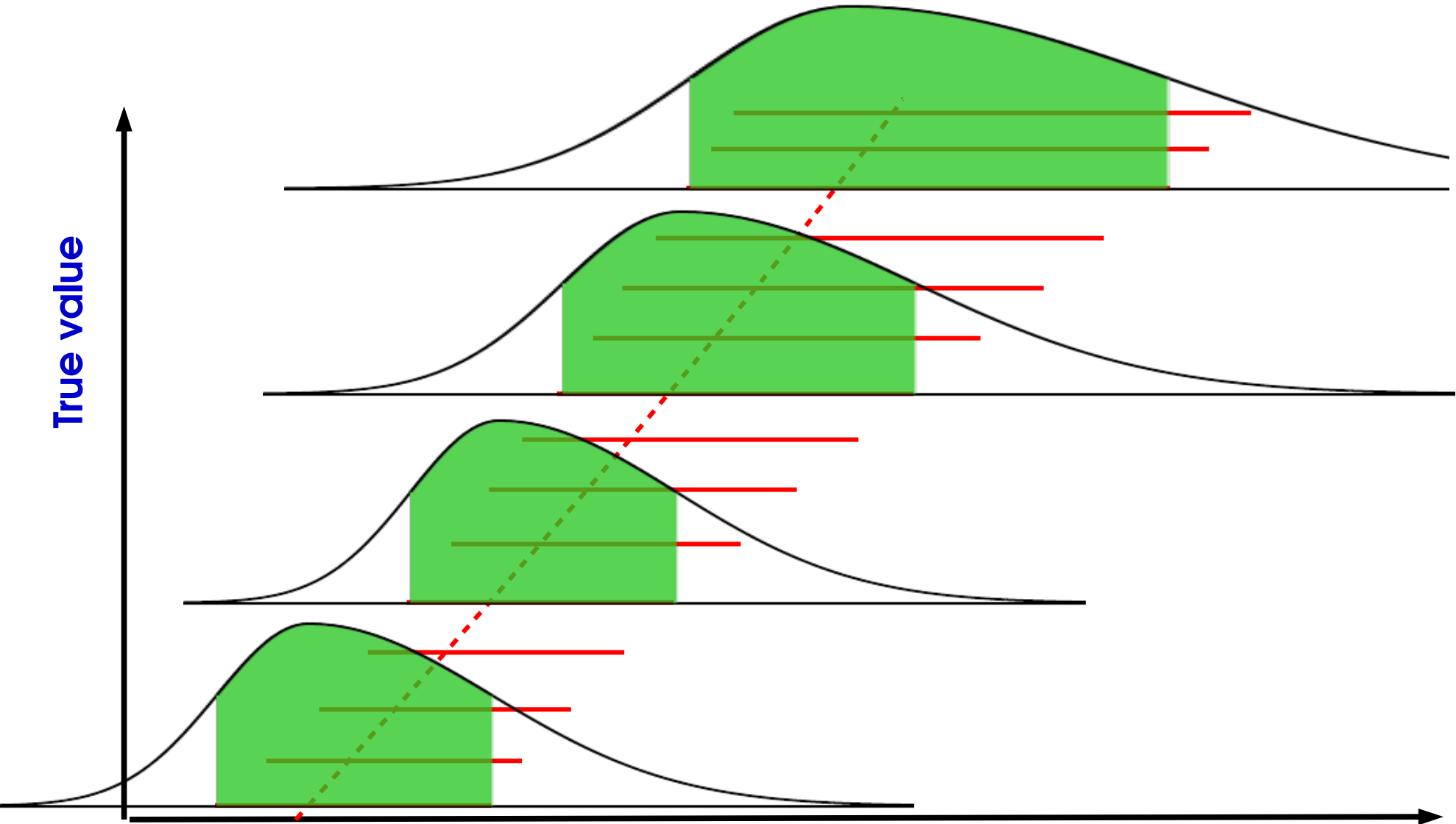
Works in the same way for other interval sizes: $[\hat{\mu} - Z\sigma, \hat{\mu} + Z\sigma]$ with

Z	1	1.96	2
CL	0.68	0.95	0.955

Neyman Construction

General case: Build 1σ intervals of observed values for each true value

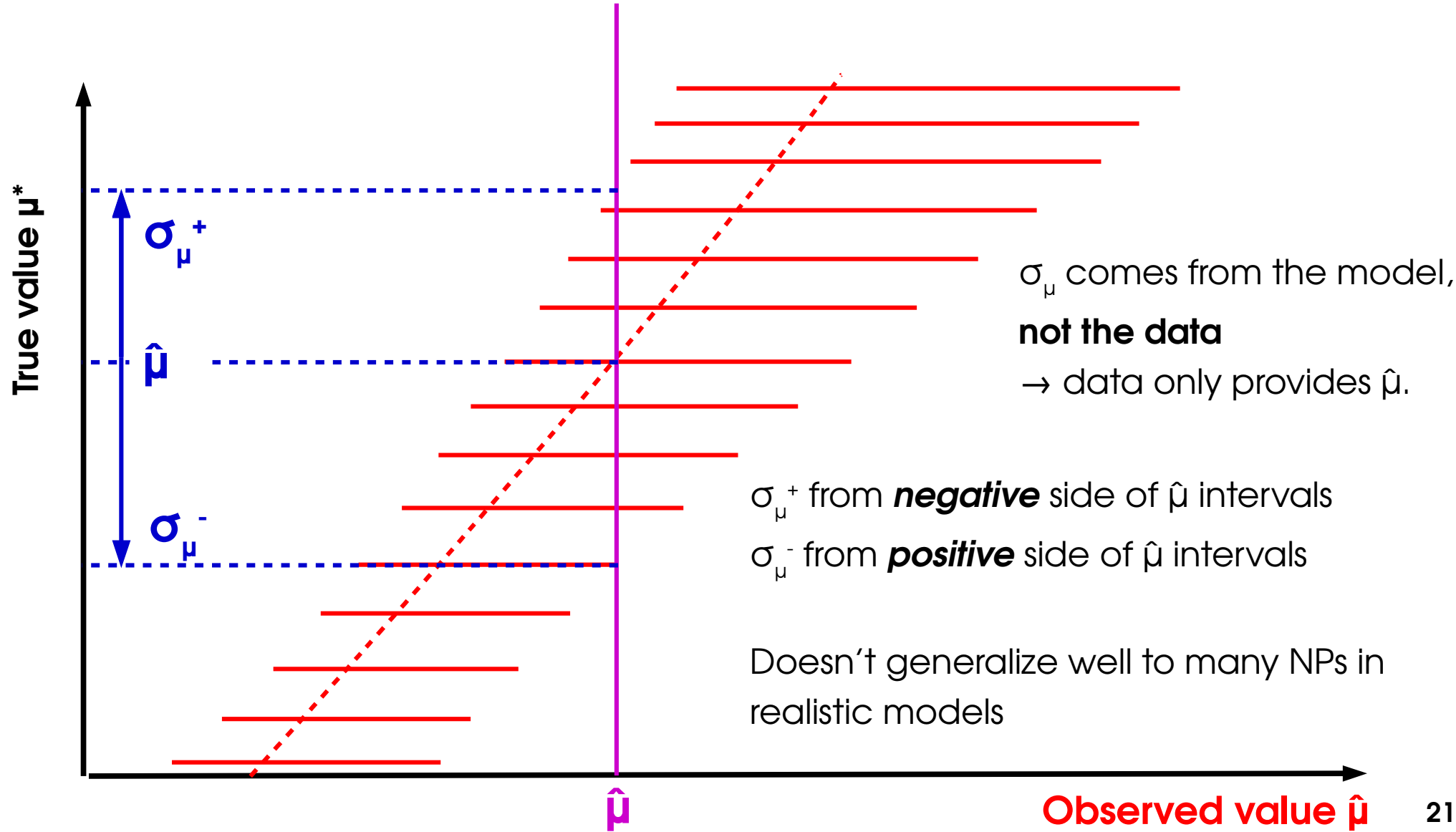
⇒ *Confidence belt*



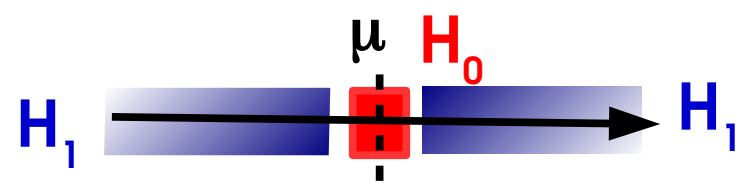
Observed value

Inversion using the Confidence Belt

General case: Intersect belt with given $\hat{\mu}$, get $P(\hat{\mu} - \sigma_{\mu}^{-} < \mu^* < \hat{\mu} + \sigma_{\mu}^{+}) = 68\%$
 → Same as before for Gaussian, works also when $P(\mu^{obs} | \mu)$ varies with μ .



Likelihood Intervals



Confidence intervals from L:

- Test $H(\mu_0)$ against alternative using
- Two-sided test since true value can be higher or lower than observed

$$t_{\mu_0} = -2 \log \frac{L(\mu = \mu_0)}{L(\hat{\mu})}$$

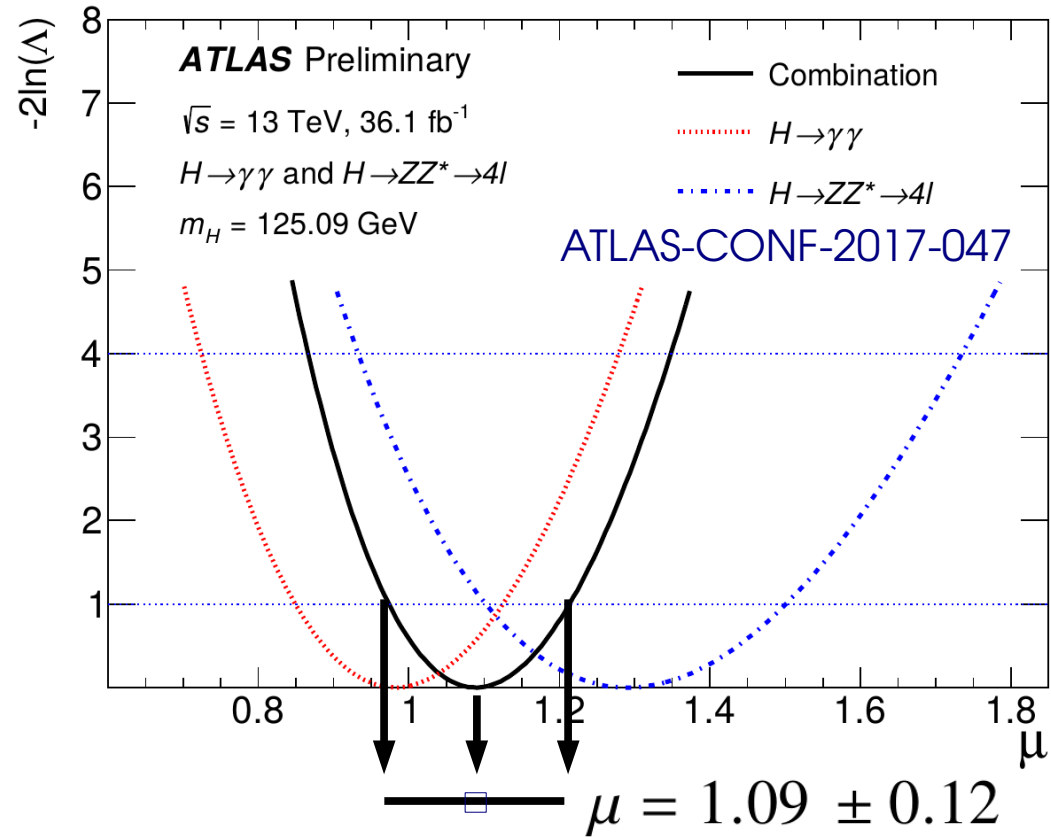
μ can be several POI!

Asymptotics:

- $t_{\mu} \sim \chi^2(N_{\text{POI}})$ under $H(\mu_0)$
- $\sqrt{t_{\mu}} \sim \mathcal{G}(0, 1)$ (Gaussian with $d=N_{\text{POI}}$)

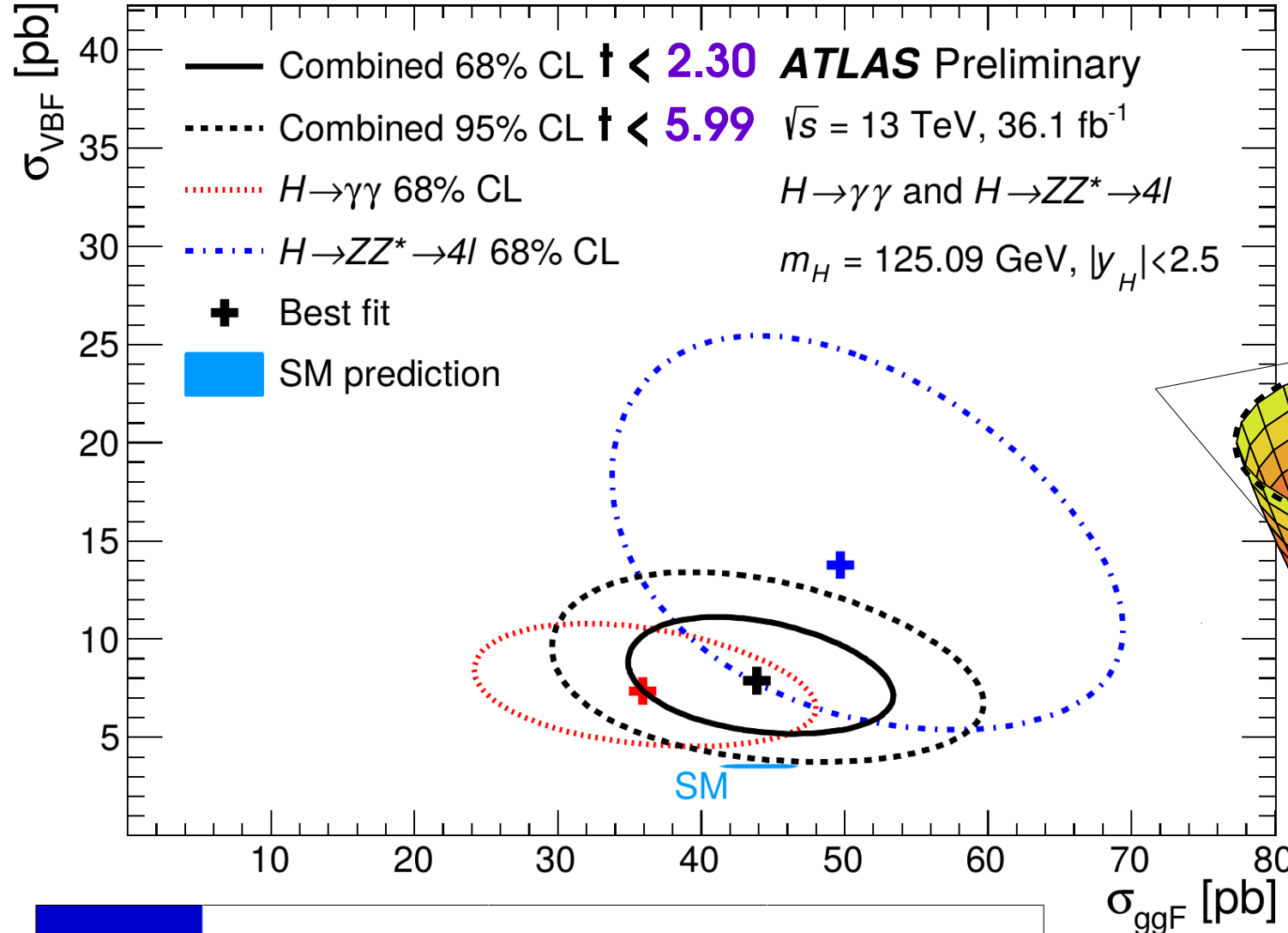
In practice:

- Plot t_{μ} vs. μ
- The minimum occurs at $\mu = \hat{\mu}$
- Crossings with $t_{\mu} = Z^2$ give the $\pm Z\sigma$ uncertainties (for $N_{\text{POI}}=1$)



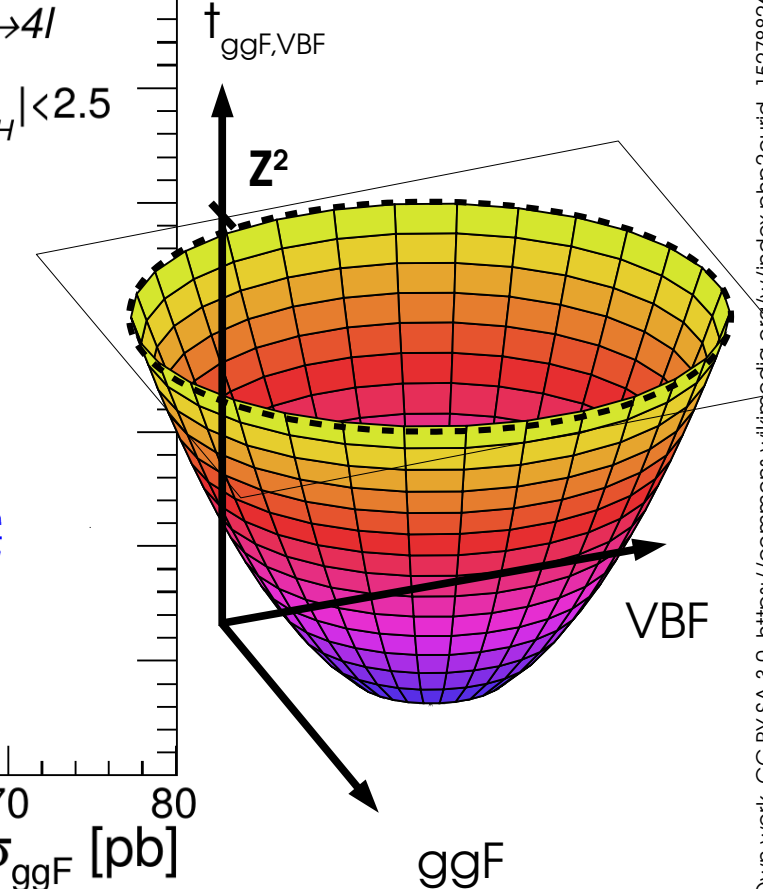
→ **Gaussian case:** parabolic profile, $t_{\mu} = \left(\frac{\mu - \hat{\mu}}{\sigma}\right)^2 \Rightarrow \mu_{\pm} = \hat{\mu} \pm \sigma$ at $t_{\mu} = 1$
 same result as Neyman construction, also robust against non-Gaussian effects.

2D Example: Higgs σ_{VBF} vs. σ_{ggF}



$$t = -2 \log \frac{L(X_0, Y_0)}{L(\hat{X}, \hat{Y})}$$

$$\sim \chi^2(N_{\text{dof}} = 2)$$



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

Gaussian case: elliptic paraboloid surface

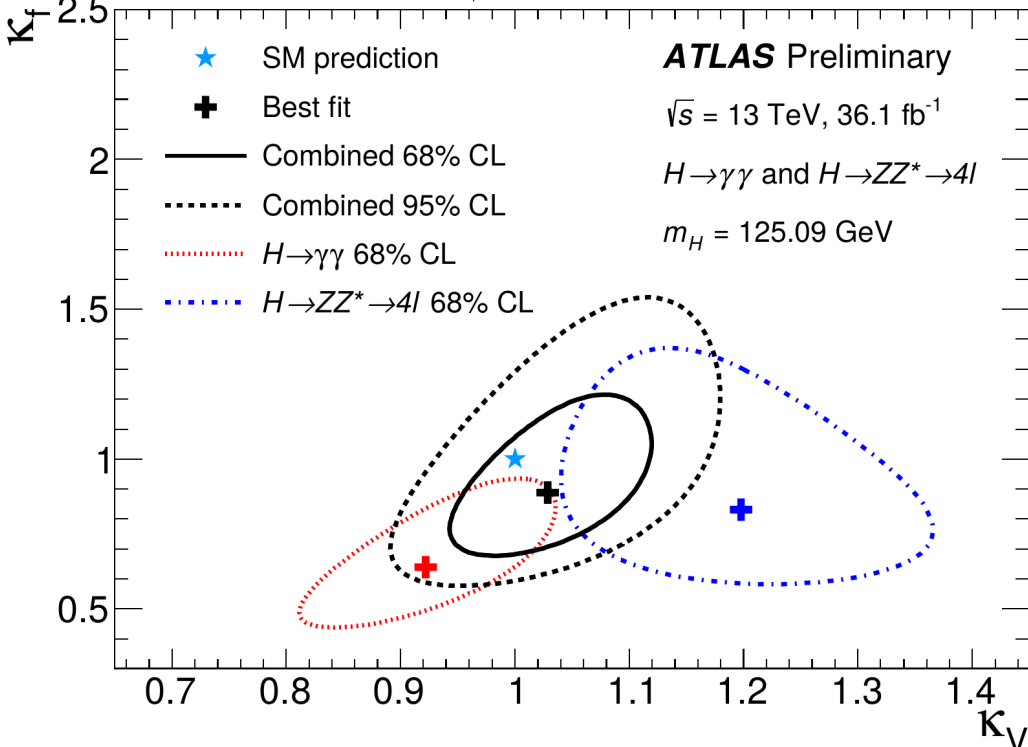
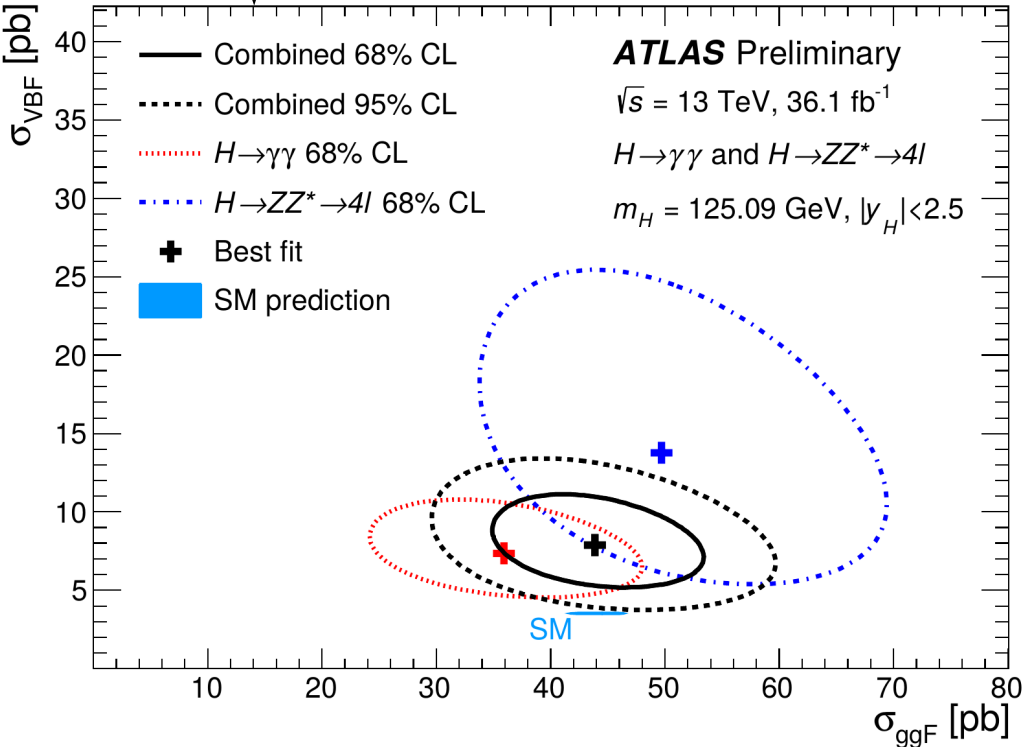
Reparameterization

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

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

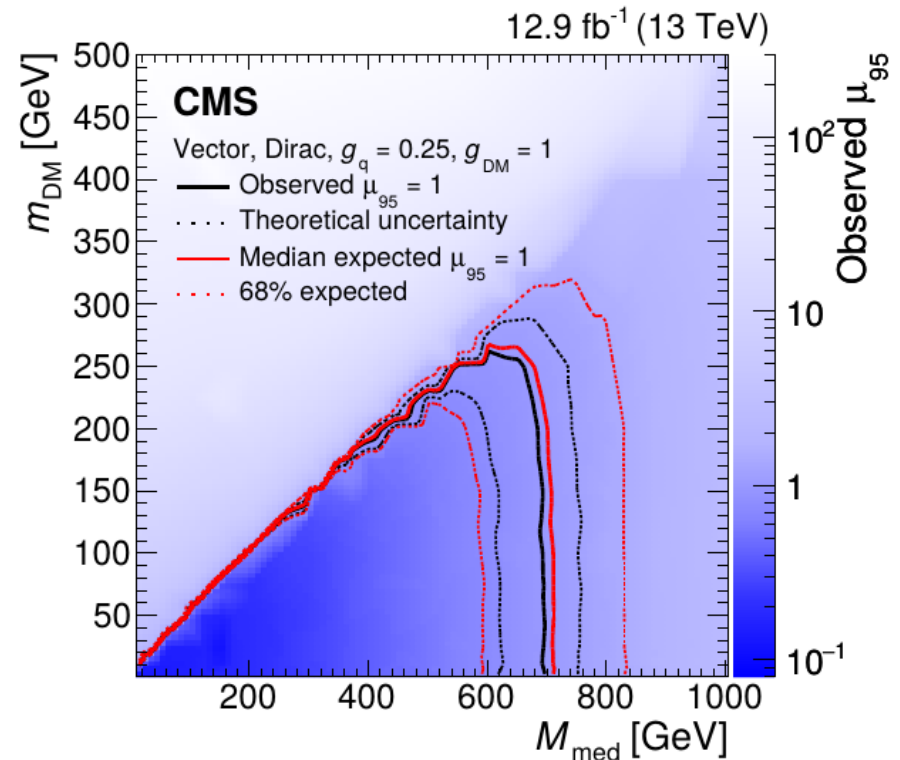
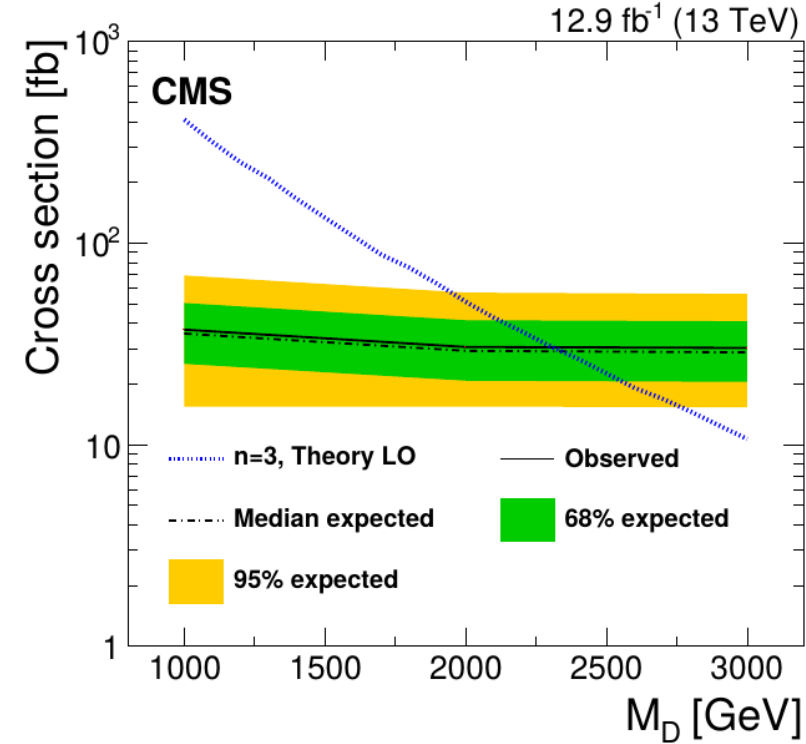
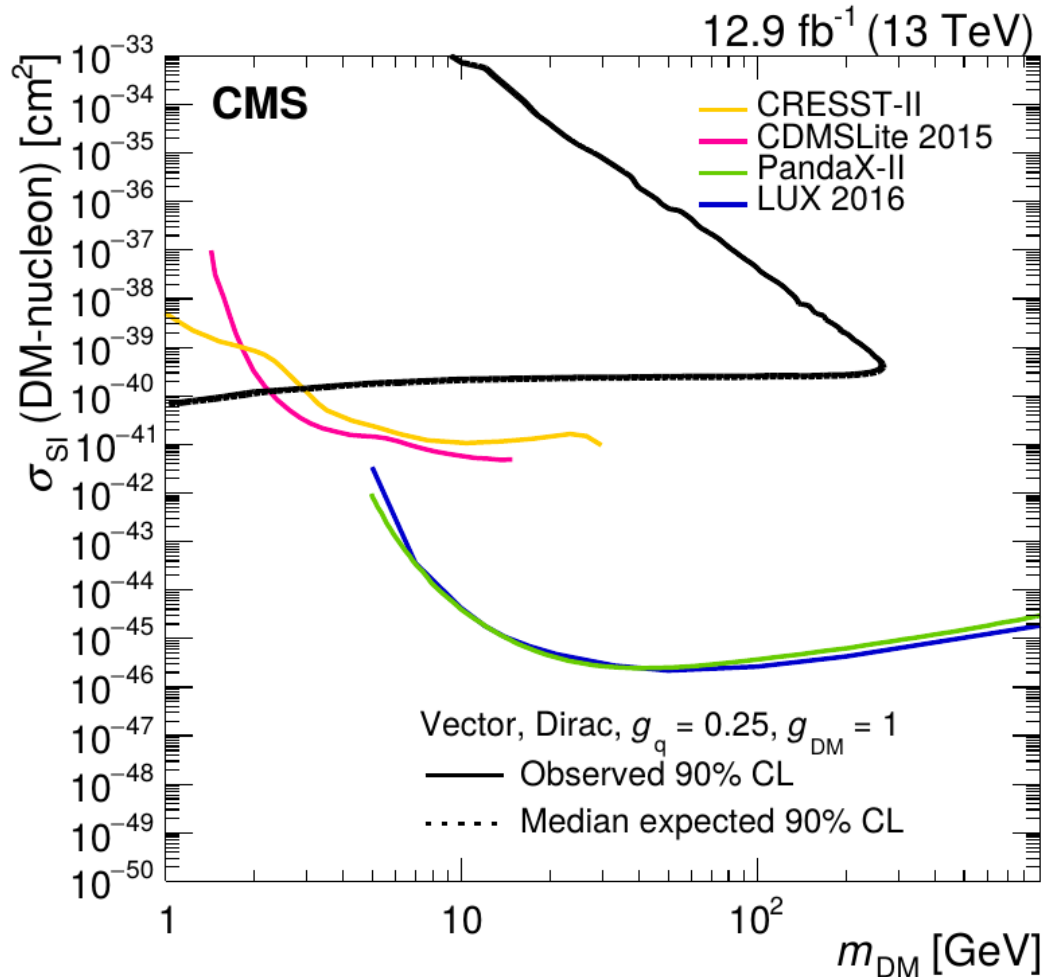
e.g. Higgs couplings: $\sigma_{ggF}, \sigma_{VBF}$ sensitive to Higgs coupling modifiers κ_V, κ_F .

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



Reparameterization: Limits

CMS Run 2 Monophoton Search: measured N_s in a counting experiment reparameterized according to various DM models



Takeaways

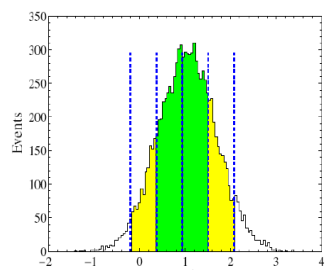
Limits : use LR-based test statistic:

→ Use **CL_s procedure** to avoid negative limits

$$\tilde{q}_{\mu_0} = \begin{cases} 0 & \hat{\mu} \geq \mu_0 \\ -2 \log \frac{L(\mu = \mu_0)}{L(\hat{\mu})} & 0 \leq \hat{\mu} \leq \mu_0 \\ -2 \log \frac{L(\mu = \mu_0)}{L(\mu = 0)} & \hat{\mu} < 0 \end{cases}$$

Poisson regime, $n=0$: $S_{up} = 3$ events

Gaussian regime, $n=0$: $S_{up} = 1.96 \sigma_{Gauss}$



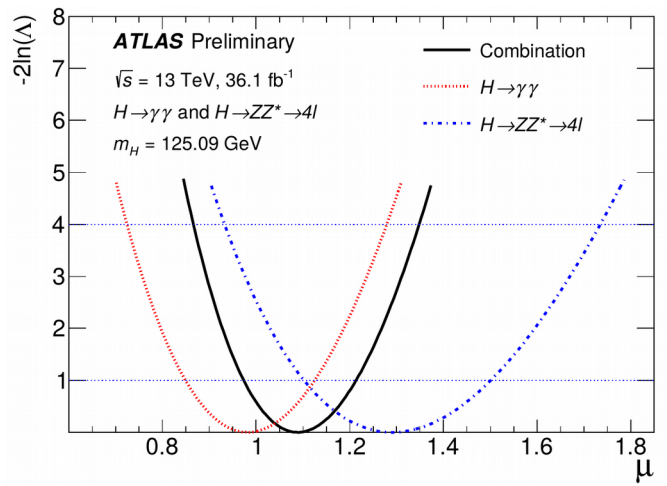
Uncertainty bands: obtain from toys or from Asimov

$$\sigma_{S,A}^2 = \frac{S^2}{q_S(\text{Asimov})}$$

Confidence intervals: use $t_{\mu_0} = -2 \log \frac{L(\mu = \mu_0)}{L(\hat{\mu})}$

→ 1D: crossings with $t_{\mu_0} = Z^2$ for $\pm Z\sigma$ intervals

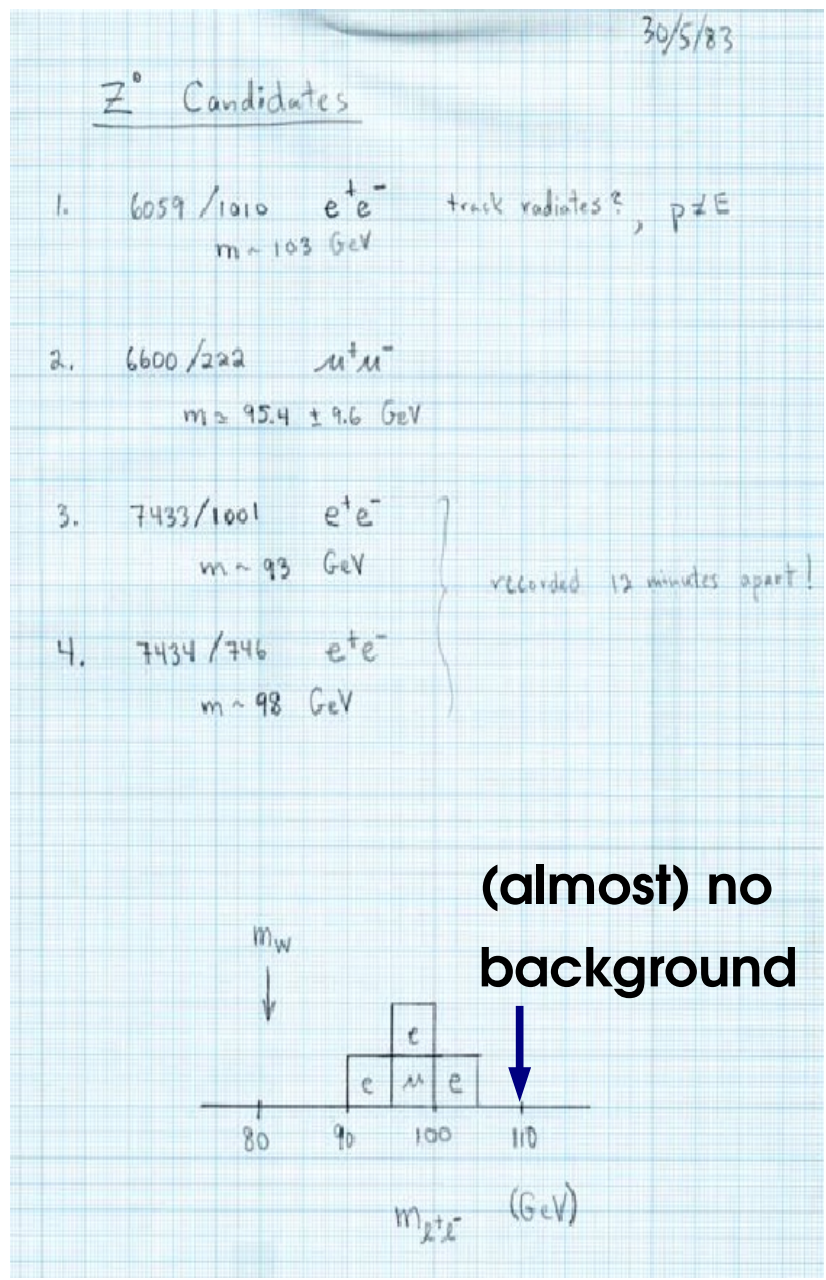
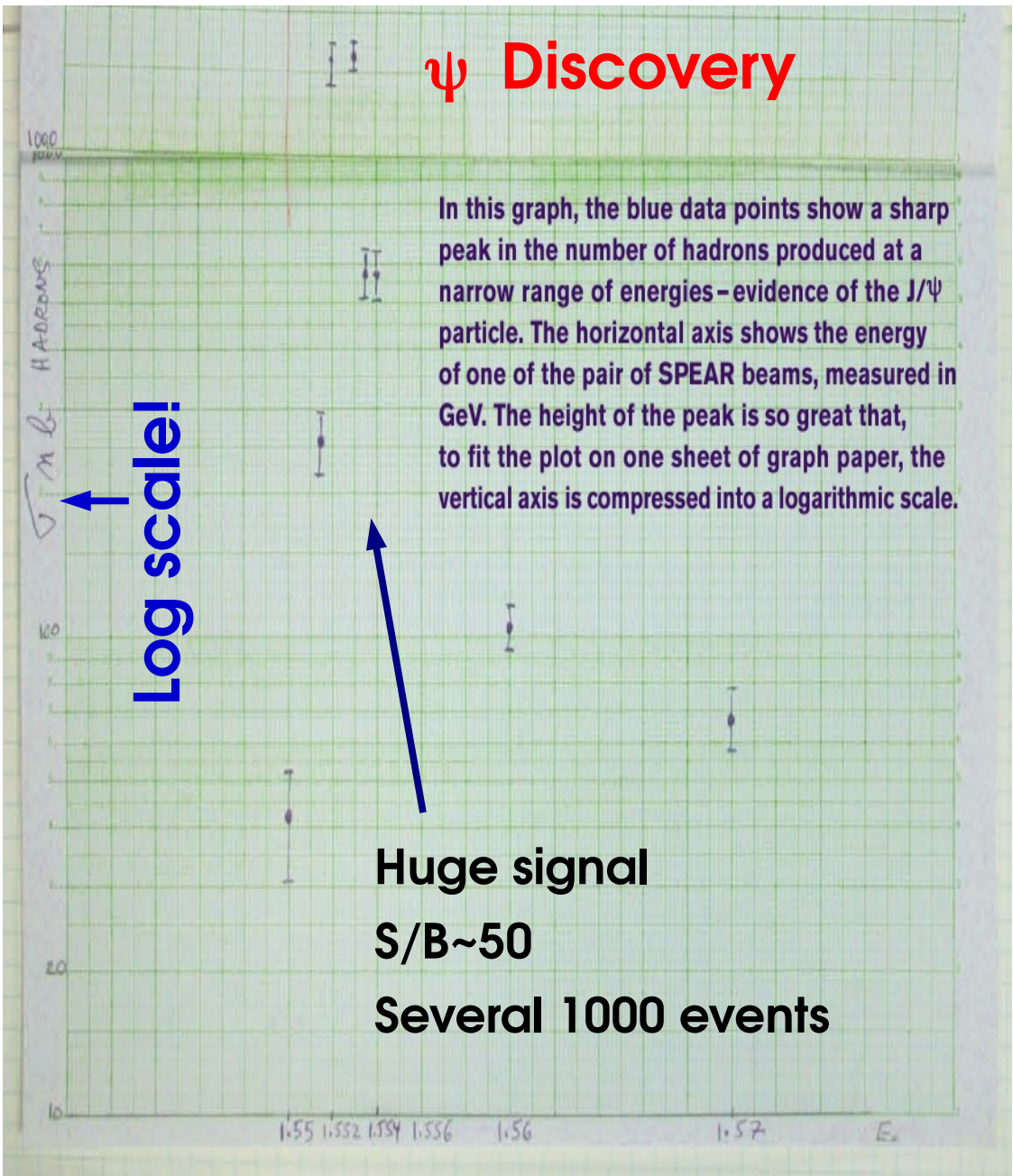
Gaussian regime: $\mu = \hat{\mu} \pm \sigma_{Gauss}$ (1σ interval)



Historical Aside

Classic Discoveries (1)

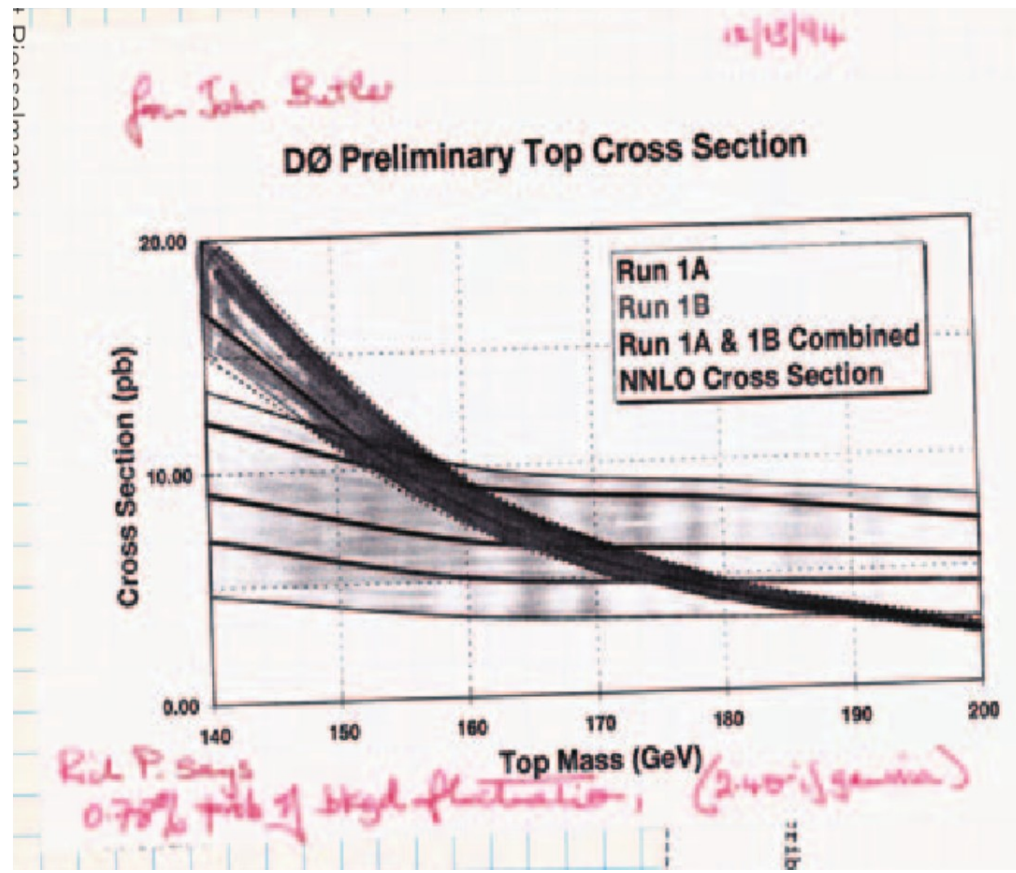
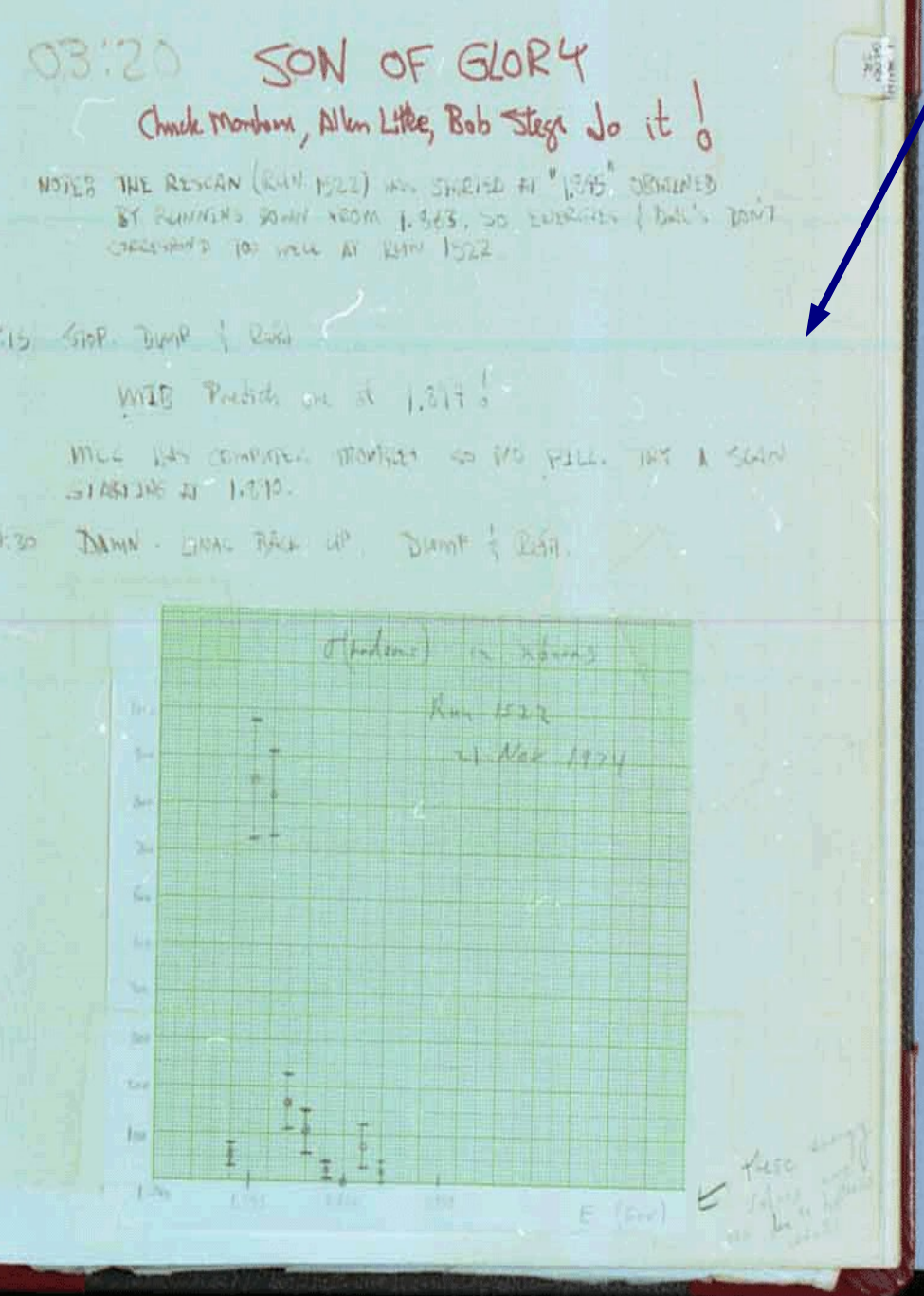
Z⁰ Discovery



Logbook of J. Rohlf, 1983-05-30

Classic Discoveries (2)

ψ' : discovered online
by the (lucky) shifters



First hints of top at DØ:
O(10) signal events,
a few bkg events, 2.4 σ

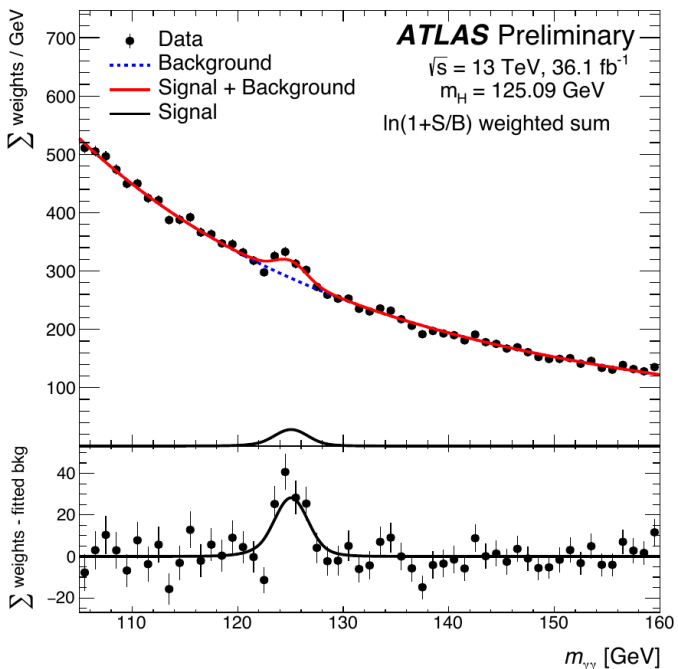
And now ?

Short answer: The high-signal, low-background experiments have been done already (although a surprise would be welcome...)

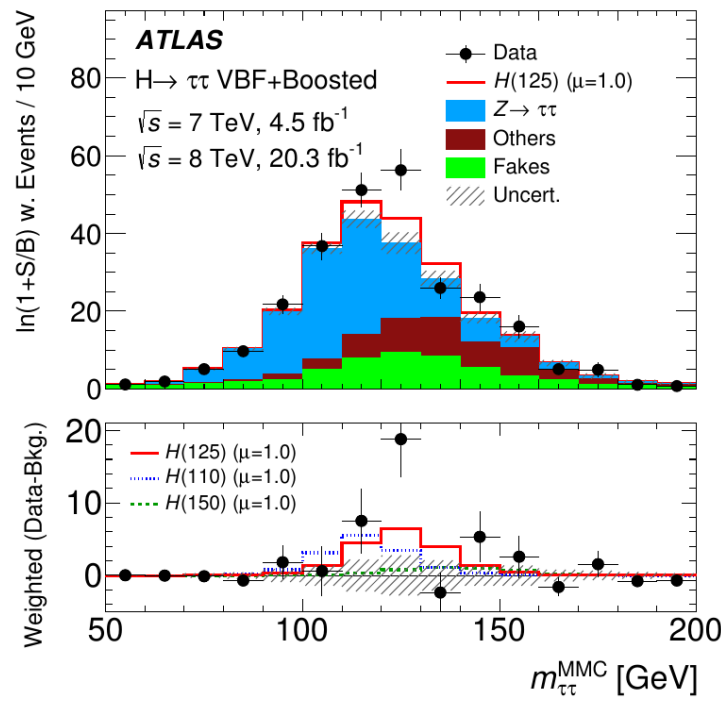
e.g. at LHC:

- **High background levels**, need precise modeling
- **Large systematics**, need to be described accurately
- **Small signals**: need optimal use of available information :
 - **Shape analyses** instead of counting
 - **Categories** to isolated signal-enriched regions

ATLAS-CONF-2017-045



JHEP 12 (2017) 024



Discoveries that weren't

UA1 Monojets (1984)

Volume 139B, number 1,2

PHYSICS LETTERS

3 May 1984

EXPERIMENTAL OBSERVATION OF EVENTS WITH LARGE MISSING TRANSVERSE ENERGY ACCOMPANIED BY A JET OR A PHOTON (S) IN $p\bar{p}$ COLLISIONS AT $\sqrt{s} = 540$ GeV

UA1 Collaboration, CERN, Geneva, Switzerland

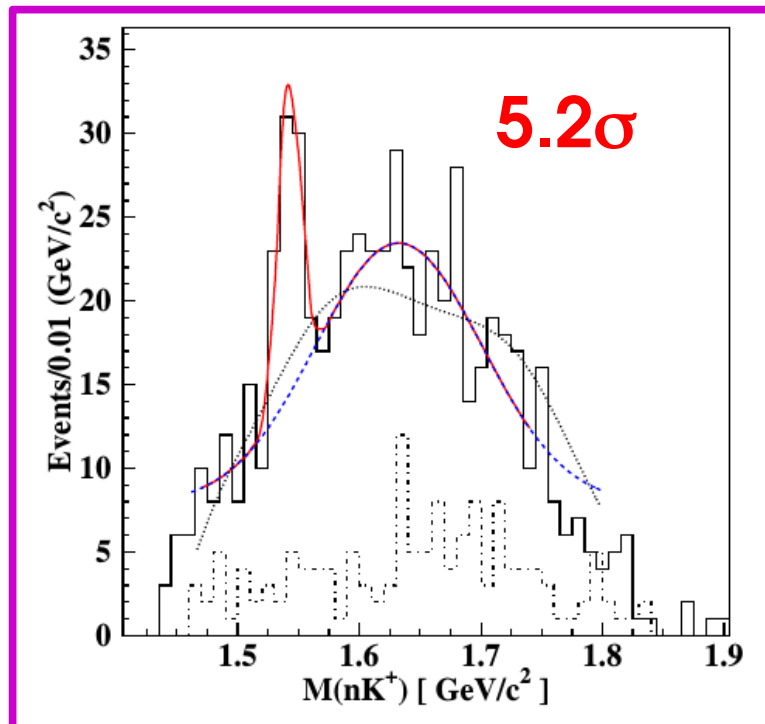
At the present time we can only speculate about the origin of this new effect. The missing transverse energy can be due either to:

(i) One or more prompt neutrinos.
 (ii) Any invisible Z^0 , such as $Z^0 \rightarrow \nu\bar{\nu}$ decay, which is expected to have a large (18%) branching ratio. Note that the corresponding decays into charged lepton pairs $Z^0 \rightarrow e^+e^-$, $Z^0 \rightarrow \mu^+\mu^-$ have lower branching ratios ($\sim 3\%$) and may not have yet been produced within the present statistics.

(iii) New, non-interacting neutral particles.
 The jets appear somewhat narrower and with lower multiplicities than the corresponding QCD jets, although it might be premature to draw conclusions on such limited statistics.


A number of theoretical speculations [9] may be relevant to these results. We mention briefly the possibilities of excited quarks or leptons and of composite or coloured or supersymmetric W's and Higgs. A recent calculation [10]¹⁸ has been made in the context of the present collider experiment, on the rate of events with large missing transverse energy from gluino pair production with each gluino decaying into a quark, antiquark, and photino. The non-interacting photinos may produce large apparent missing energy. For instance, the calculation gives an expectation of about 100 single-jet events with $\Delta E_M > 20$ GeV for a gluino mass of 20 GeV/c². Taking our excess of 5 events above background as an upper limit for such a process, we deduce that the gluino mass must be greater than about 40 GeV/c².


Pentaquarks (2003)



Phys. Rev. Lett. 91, 252001 (2003)

BICEP2 B-mode Polarization (2014)

 Selected for a Viewpoint in *Physics*
 PRL 112, 241101 (2014) PHYSICAL REVIEW LETTERS week ending 20 JUNE 2014


 Detection of *B*-Mode Polarization at Degree Angular Scales by BICEP2

$$r = 0.20_{-0.05}^{+0.07}, \text{ with } r = 0 \text{ disfavored at } 7.0\sigma.$$

Avoid spurious discoveries!

→ Treatment of modeling uncertainties, systematics in general

Outline

Computing Statistical Results

Limits, continued

Confidence Intervals

Profiling

Look-Elsewhere Effect

Bayesian methods

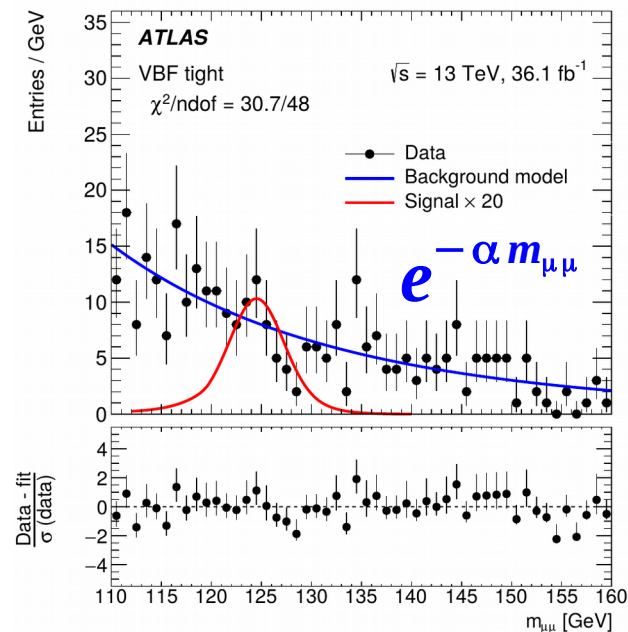
Statistical modeling in practice

BLUE

Nuisances and Systematics

Likelihood typically includes

- **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
 - e.g. shape of $H \rightarrow \mu\mu$ continuum bkg



What about systematics ?

= what we don't know about the random process

⇒ **Parameterize using additional NPs**

→ By definition, **not constrained by the data**

⇒ Cannot be free, or would spoil the measurement (lumi free ⇒ no $\sigma \times B$ measurement!)

⇒ **Introduce a constraint in the likelihood:**

"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 ?*

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

⇒ penalty for $\theta \neq \theta^{\text{nominal}}$

Frequentist Constraints

Prototype: NP measured in a separate *auxiliary experiment*

e.g. luminosity measurement

→ Build the combined likelihood of the main+auxiliary measurements

$$L(\boldsymbol{\mu}, \boldsymbol{\theta}; \text{data}) = L_{\text{main}}(\boldsymbol{\mu}, \boldsymbol{\theta}; \text{main data}) L_{\text{aux}}(\boldsymbol{\theta}; \text{aux. data})$$

Independent
measurements:
⇒ just a product

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

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

→ now same as other NPs: **all uncertainties statistical in nature**

→ Often no clear setup for auxiliary measurements

e.g. theory uncertainties on missing HO terms from scale variations

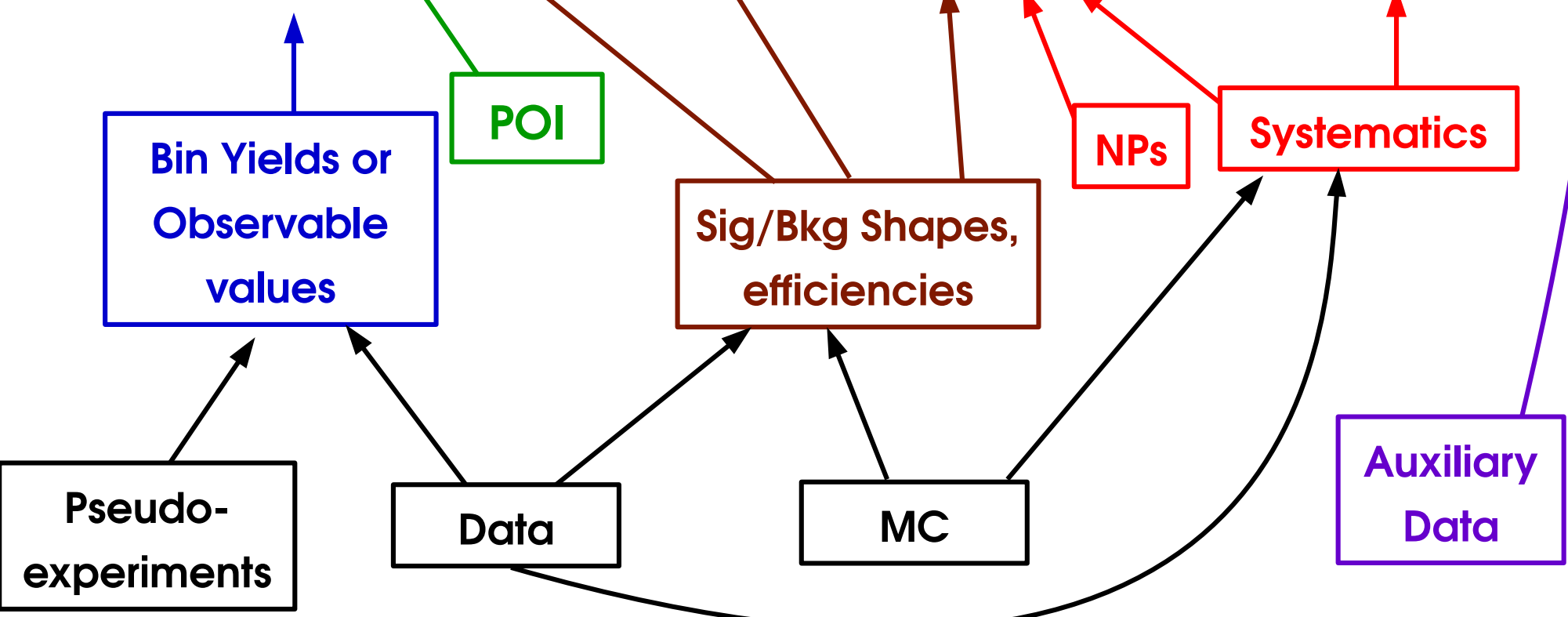
→ **Implemented in the same way nevertheless** (“pseudo-measurement”)

Likelihood, the full version (binned case)

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

Expected bin yield

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



* number of categories!

Reminder: Wilks' Theorem

→ Assume **Gaussian regime** for \hat{S} (e.g. large n_{evts})

Cowan, Cranmer, Gross & Vitells
Eur.Phys.J.C71:1554,2011

⇒ Central-limit theorem :

t_0 is distributed as a χ^2 under the hypothesis H_0

$$t_0 = -2 \log \frac{L(S=0)}{L(\hat{S})}$$

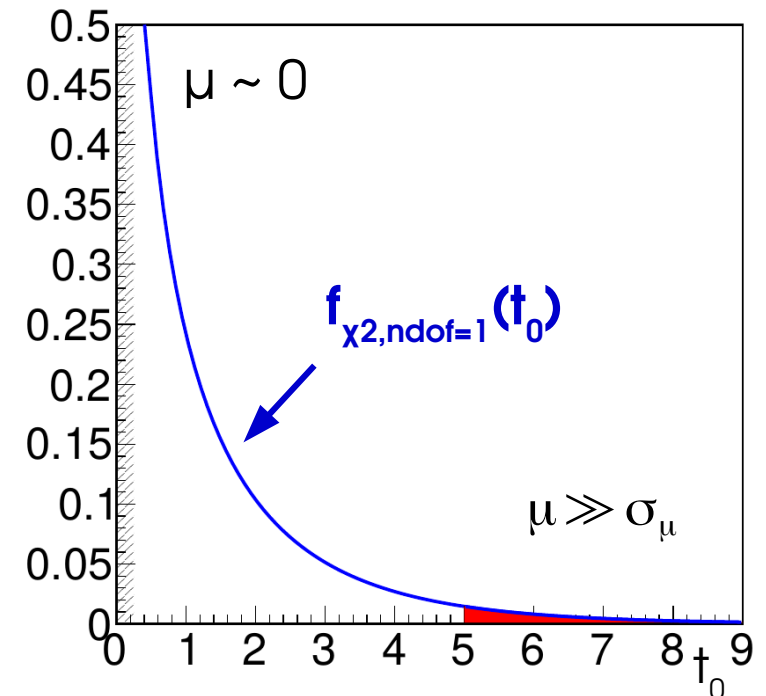
$$f(t_0 | H_0) = f_{\chi^2(n_{\text{dof}}=1)}(t_0)$$

In particular, significance:

$$Z = \sqrt{t_0}$$

By definition,
 $t_0 \sim \chi^2 \Rightarrow \sqrt{t_0} \sim G(0,1)$

Typically works well for event counts $O(5)$
and above (5 already "large" ...)



The 1-line "proof" : asymptotically L and S are Gaussian, so

$$L(S) = \exp \left[-\frac{1}{2} \left(\frac{S - \hat{S}}{\sigma} \right)^2 \right] \Rightarrow t_0 = \left(\frac{\hat{S}}{\sigma} \right)^2 \Rightarrow t_0 \sim \chi^2(n_{\text{dof}}=1) \text{ since } \hat{S} \sim G(0, \sigma)$$

Wilks' Theorem, the Full Version

The likelihood usually has NPs:

- **Systematics**
- Parameters fitted in data

→ What values to use when defining the hypotheses ? → $H(S=0, \theta=?)$

Answer: let the data choose ⇒ use the best-fit values (*Profiling*)

⇒ **Profile Likelihood Ratio** (PLR)

$$t_{\mu_0} = -2 \log \frac{L(\mu = \mu_0, \hat{\theta}_{\mu_0})}{L(\hat{\mu}, \hat{\theta})}$$

$\hat{\theta}_{\mu_0}$ best-fit value for $\mu = \mu_0$ (conditional MLE)
 $\hat{\theta}$ overall best-fit value (unconditional MLE)

Wilks' Theorem: PLR also follows a χ^2 ! $f(t_{\mu_0} | \mu = \mu_0) = f_{\chi^2(n_{dof}=1)}(t_{\mu_0})$
also with NPs present

→ Profiling “builds in” the effect of the NPs

⇒ Can treat the PLR as a **function of the POI only**

Gaussian Profiling

Recall: Gaussian counting, no syst: $t_{S_0} = \left(\frac{S_0 - \hat{S}}{\sigma_S} \right)^2$

Counting exp. with background uncertainty: $\mathbf{n} = \mathbf{S} + \boldsymbol{\theta}$:

$$\left. \begin{array}{l} \rightarrow \text{Main measurement: } \mathbf{n} \sim \mathbf{G}(\mathbf{S} + \boldsymbol{\theta}, \sigma_{\text{stat}}) \\ \rightarrow \text{Aux. measurement: } \boldsymbol{\theta}^{\text{obs}} \sim \mathbf{G}(\boldsymbol{\theta}, \sigma_{\text{syst}}) \end{array} \right\} L(S, \theta) = G(n; S + \theta, \sigma_{\text{stat}}) G(\theta^{\text{obs}}; \theta, \sigma_{\text{syst}})$$

Then: $\lambda(S, \theta) = \left(\frac{n - (S + \theta)}{\sigma_{\text{stat}}} \right)^2 + \left(\frac{\theta^{\text{obs}} - \theta}{\sigma_{\text{syst}}} \right)^2$

For $S = \hat{S}$, matches MLE as it should

MLEs: $\hat{S} = n - \theta^{\text{obs}}$ **Conditional MLE:** $\hat{\theta}(S) = \theta^{\text{obs}} + \frac{\sigma_{\text{syst}}^2}{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2} (\hat{S} - S)$
 $\hat{\theta} = \theta^{\text{obs}}$

PLR: $t_{S_0} = -2 \log \frac{L(S=S_0, \hat{\theta}_{S_0})}{L(\hat{S}, \hat{\theta})}$

$$= \lambda(S_0, \hat{\theta}(S_0)) - \lambda(\hat{S}, \hat{\theta}) = \frac{(S_0 - \hat{S})^2}{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2} \quad \sigma_S = \sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}$$

Stat uncertainty (on n) and syst (on θ) add in quadrature as expected

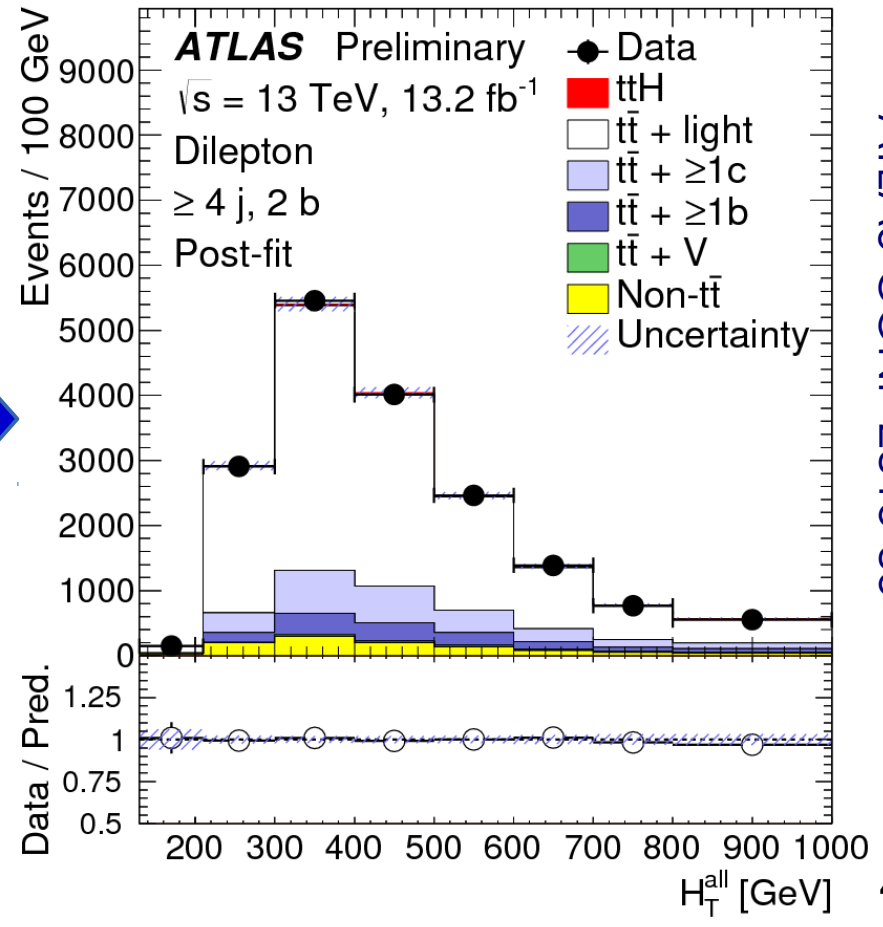
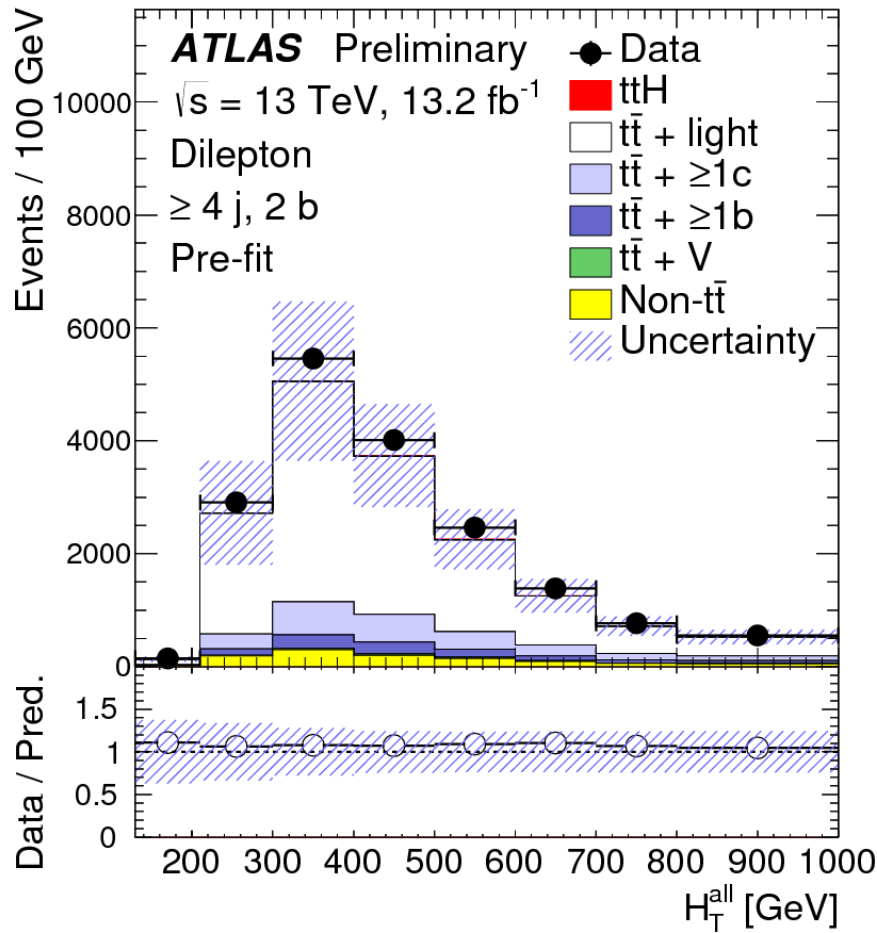
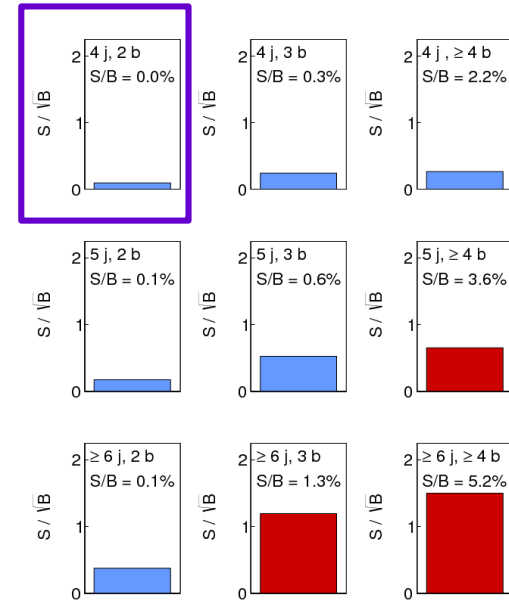
Profiling Example: $t\bar{t}H \rightarrow b\bar{b}$

Analysis uses low-S/B categories to constrain backgrounds.

→ **Reduction in large uncertainties on $t\bar{t}$ bkg**

→ **Propagates to the high-S/B categories** through the statistical modeling

⇒ **Care needed in the propagation** (e.g. different kinematic regimes)



ATLAS-CONF-2016-08

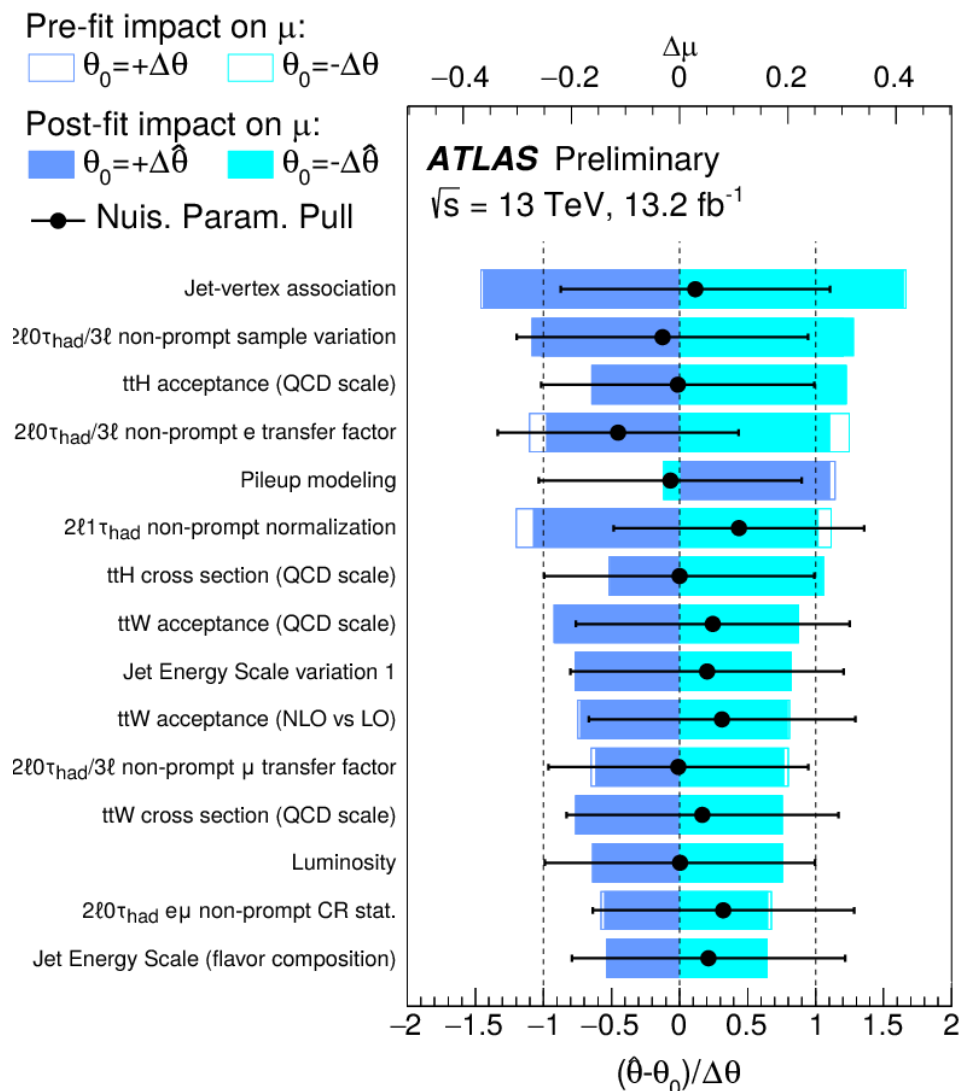
Systematics are described by NPs included in the fit. Nominally:

- **NP central value = 0** : corresponds to the pre-fit expectation (usually MC)
- **NP uncertainty = 1** : since NPs normalized to the value of the syst. :

$$N = N_0 (1 + \sigma_{\text{syst}} \theta), \theta \sim G(0, 1)$$

Fit results provide information on impact of the systematic on the result:

- **If central value $\neq 0$** : some data feature absorbed by nonzero value \Rightarrow Need investigation if large pull
- **If uncertainty < 1** : systematic is constrained by the data \Rightarrow Needs checking if this legitimate or a modeling issue
- **Impact on result** of $\pm 1\sigma$ shift of NP



Pull/Impact plots

Systematics are described by NPs included in the fit. Nominally:

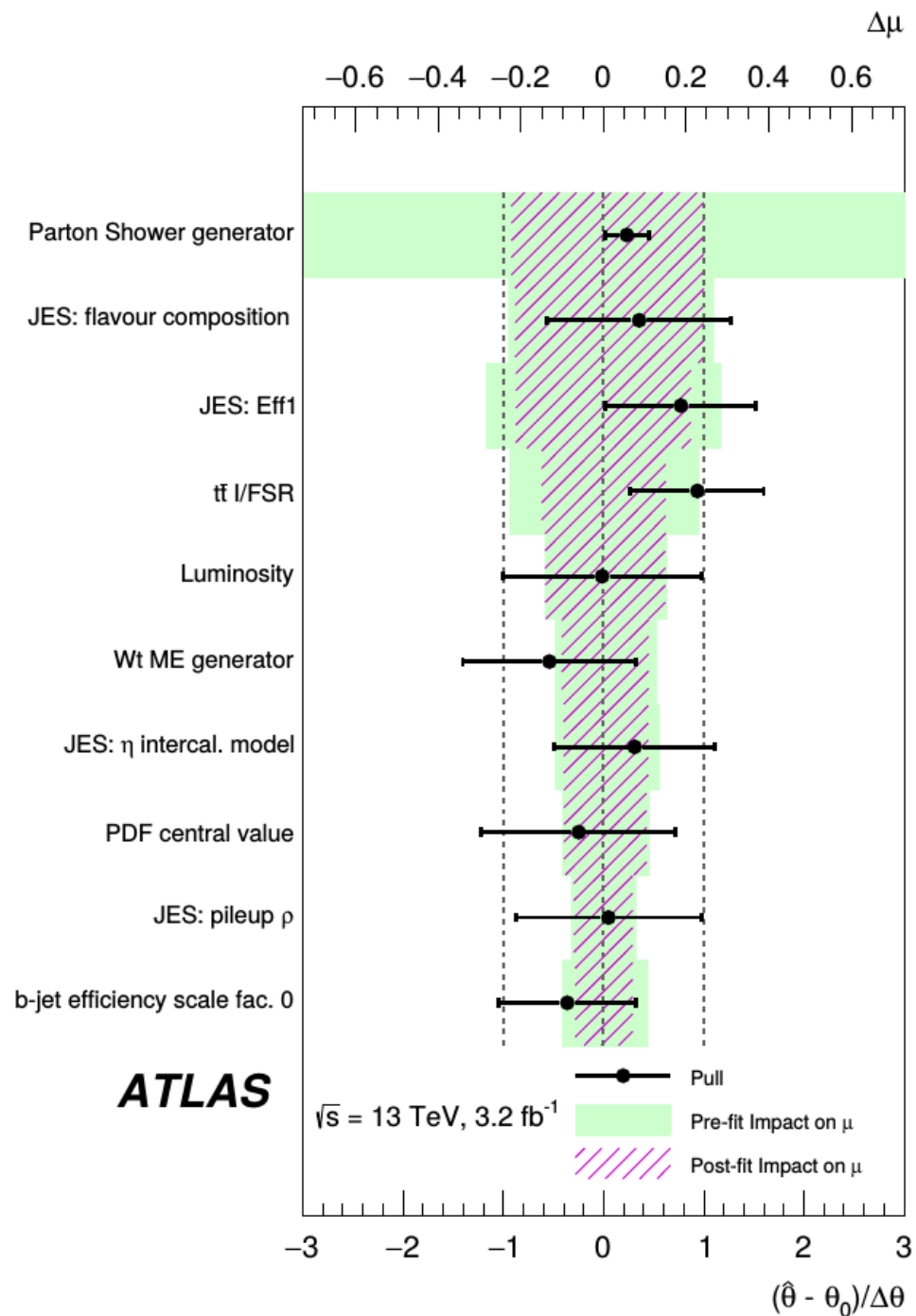
- **NP central value = 0** : corresponds to the pre-fit expectation (usually MC)
- **NP uncertainty = 1** : since NPs normalized to the value of the syst. :

$$N = N_0 (1 + \sigma_{\text{syst}} \theta), \theta \sim G(0, 1)$$

Fit results provide information on impact of the systematic on the result:

- **If central value $\neq 0$** : some data feature absorbed by nonzero value \Rightarrow Need investigation if large pull
- **If uncertainty < 1** : systematic is constrained by the data \Rightarrow Needs checking if this legitimate or a modeling issue
- **Impact on result** of $\pm 1\sigma$ shift of NP

13 TeV single- t XS (arXiv:1612.07231)



Profiling Takeaways

Systematic = NP with an external constraint (auxiliary measurement).

→ No special treatment, treated like any other NP: statistical and systematic uncertainties represented in the same way.

When testing a hypothesis, use the best-fit values of the nuisance parameters: *Profile Likelihood Ratio*.

$$\frac{L(\mu = \mu_0, \hat{\theta}_{\mu_0})}{L(\hat{\mu}, \hat{\theta})}$$

Wilks' Theorem: the PLR has the same asymptotic properties as the LR without systematics: can profile out NPs and just deal with POIs.

Profiling systematics includes their effect into the total uncertainty. Gaussian:

$$\sigma_{\text{total}} = \sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}$$

Guaranteed to work only as long as everything is Gaussian, but typically robust against non-Gaussian behavior.

Profiling can have unintended effects – need to carefully check behavior

Beyond Asymptotics: Toys

Asymptotics usually work well, but break down in some cases – e.g. **small event counts**.

Solution: generate **pseudo data (toys)** using the PDF, under the tested hypothesis

→ Also randomize the observable

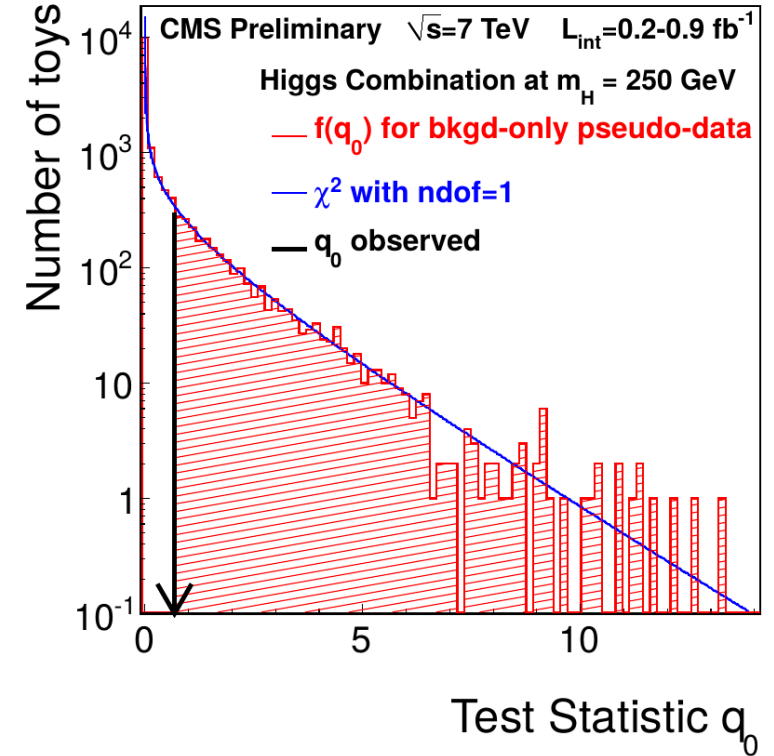
(θ^{obs}) of each auxiliary experiment: $G(\theta^{obs}; \theta, \sigma_{syst})$

→ Samples the true distribution of the PLR

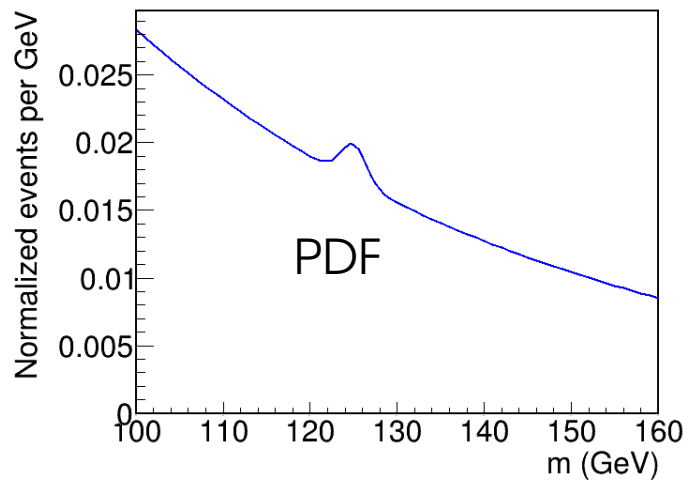
⇒ Integrate above observed PLR to get the p-value

→ Precision limited by number of generated toys,

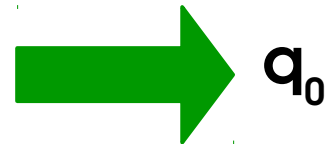
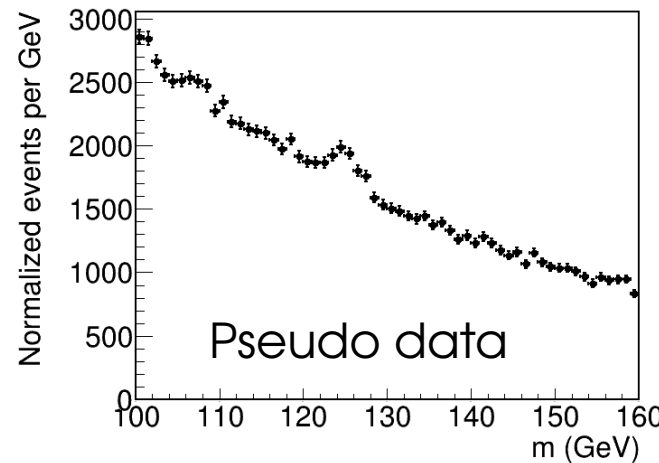
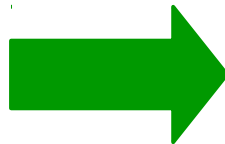
Small p-values ($5\sigma : p \sim 10^{-7}!$) ⇒ **large toy samples**



Repeat N_{toys} times



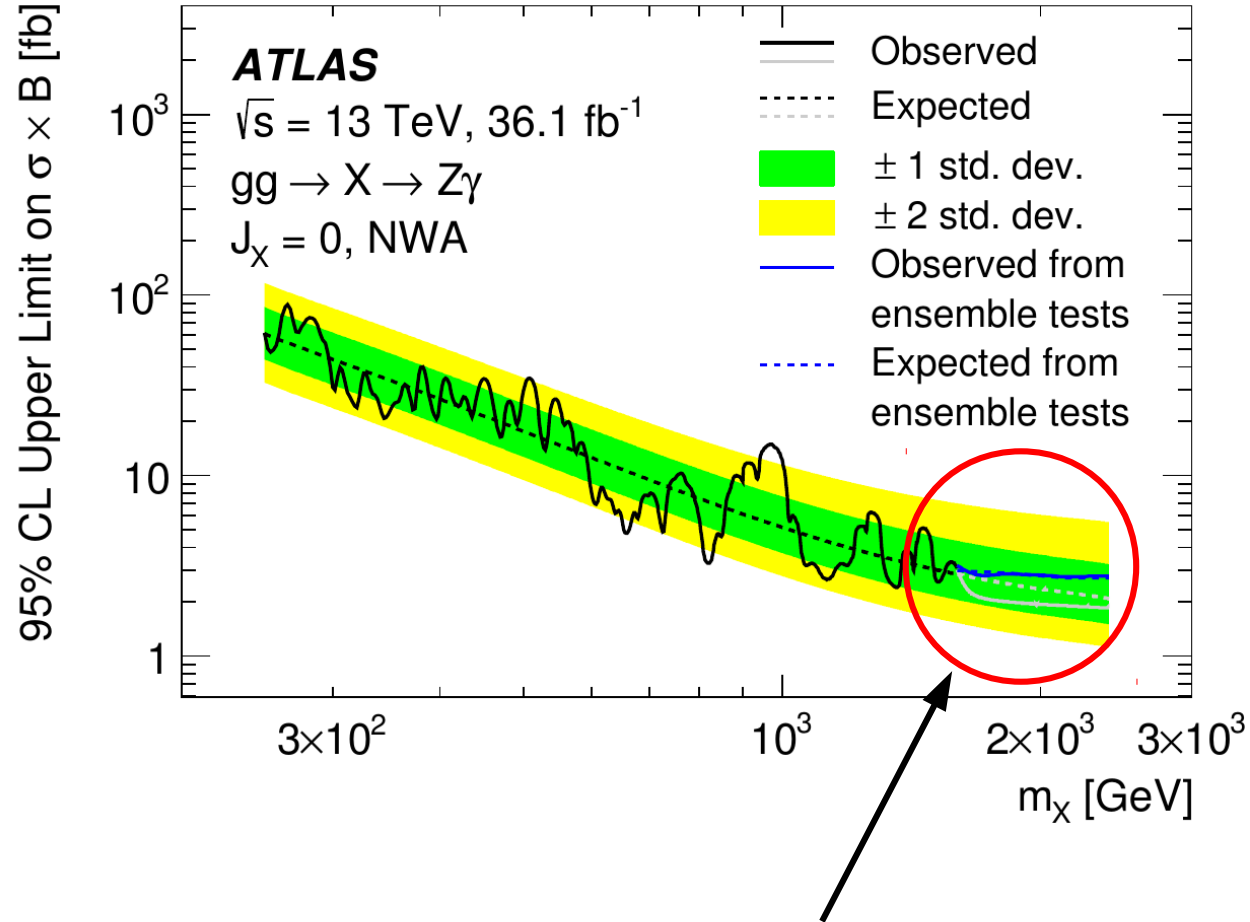
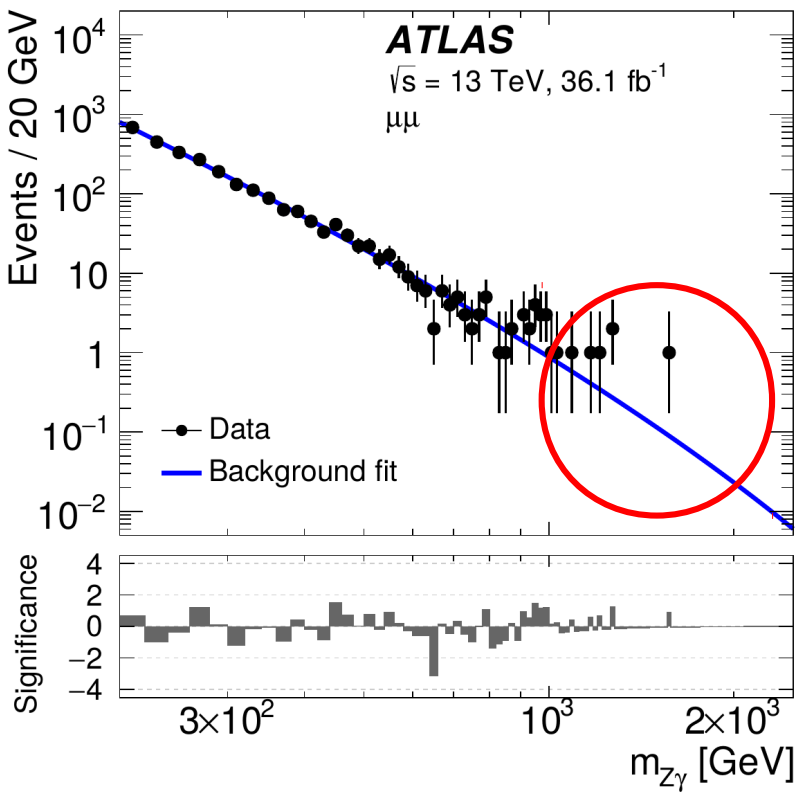
$p(\text{data} | x)$



Toys: Example

ATLAS $X \rightarrow Z\gamma$ Search: covers $200 \text{ GeV} < m_X < 2.5 \text{ TeV}$

\rightarrow for $m_X > 1.6 \text{ TeV}$, low event counts \Rightarrow derive results from toys



Asymptotic results (in gray) give optimistic result compared to toys (in blue)

Summary of Statistical Results Computation

Methods provide:

→ **Optimal use of information from the data under general hypotheses**

→ **Arbitrarily complex/realistic models (up to computing constraints...)**

→ **No Gaussian assumptions in the measurements**

Still often assume Gaussian behavior of PLR – but weaker assumption and can be lifted with toys

Systematics treated as auxiliary measurements – modeling can be tailored as needed

→ **Single PLR-based framework for all usual classes of measurements**

Discovery testing

Upper limits on signal yields

Parameter estimation

Outline

Computing Statistical Results

Limits, continued

Confidence Intervals

Profiling

Look-Elsewhere Effect

Bayesian methods

Statistical modeling in practice

BLUE

Look-Elsewhere Effect

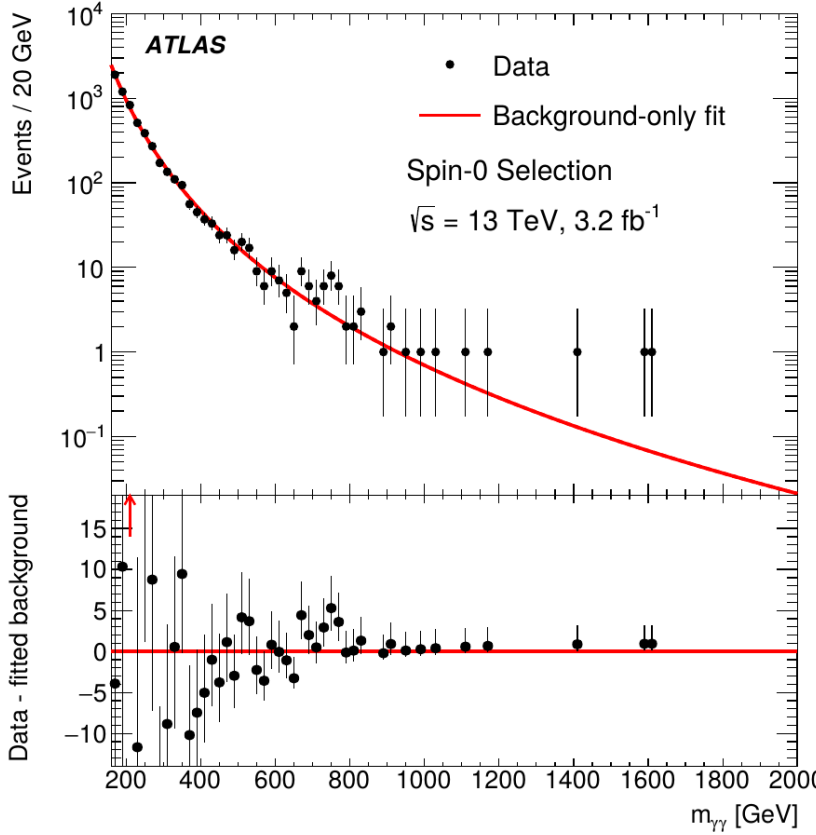
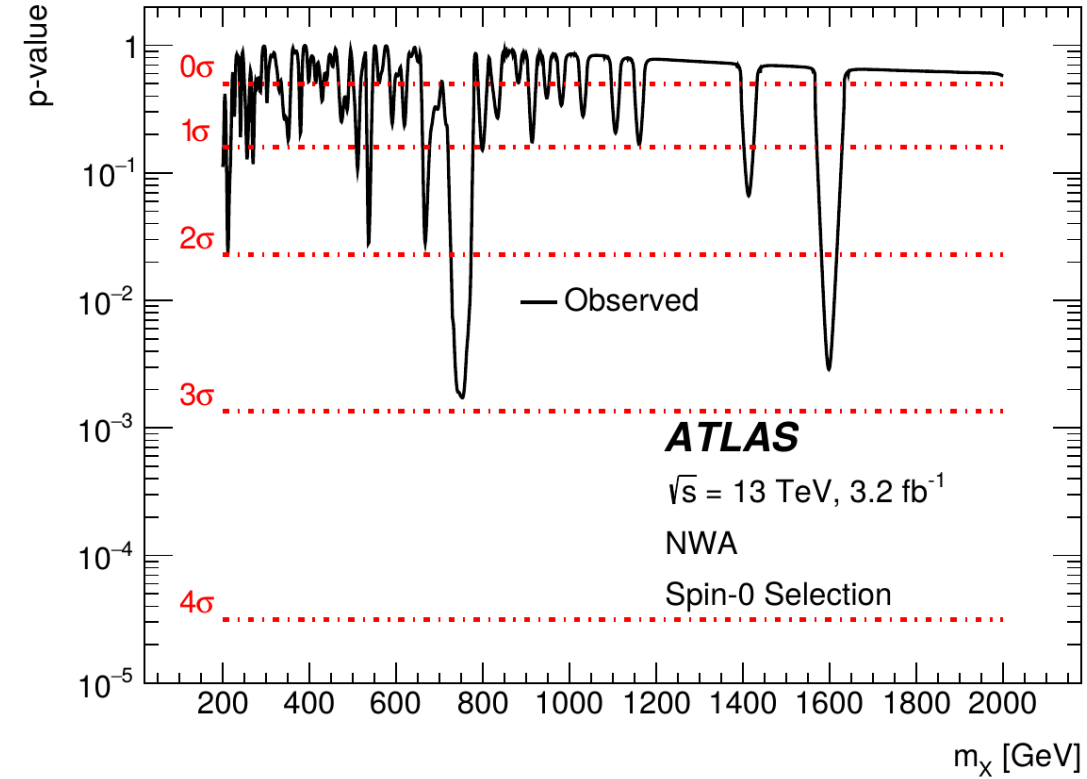
Look-Elsewhere effect

Sometimes, unknown parameters in signal model

e.g. p-values as a function of m_x

⇒ Effectively performing **multiple, simultaneous searches**

→ If e.g. small resolution and large scan range, **many independent experiments**

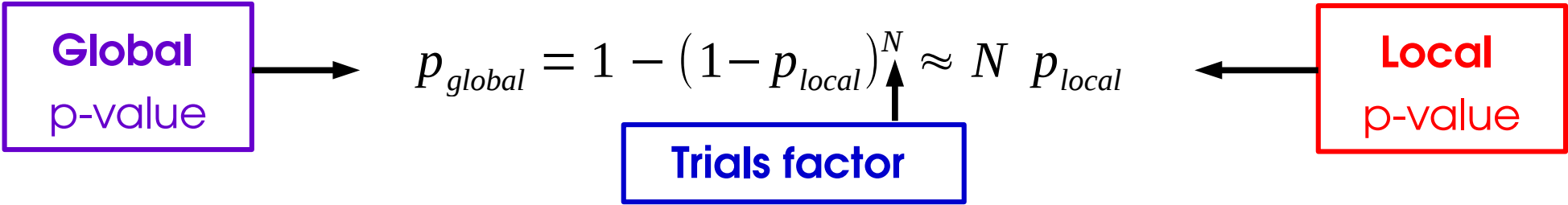


→ More likely to find an excess **anywhere in the range**, rather than in a **predefined** location
 ⇒ **Look-elsewhere effect** (LEE)

Testing the same H_0 , but against different alternatives
 ⇒ different p-values

Global Significance

Probability for a fluctuation **anywhere** in the range → **Global** p-value.
 at a given location → **Local** p-value

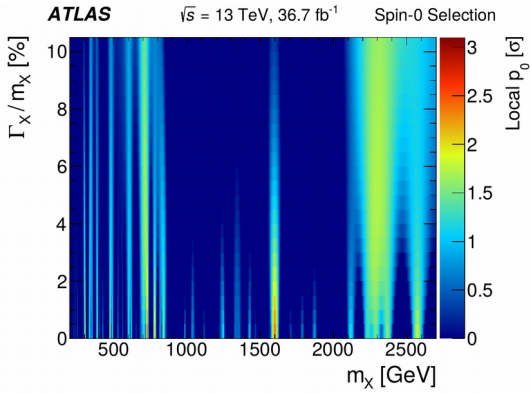


→ $p_{global} > p_{local} \Rightarrow Z_{global} < Z_{local}$ – global fluctuation more likely ⇒ less significant

Trials factor: **naively** = # of independent intervals: $N_{trials} = N_{indep} = \frac{\text{scan range}}{\text{peak width}}$
 However this is usually **wrong** – more on this later

For searches over a parameter range, p_{global} is the relevant p-value

→ Depends on the scanned parameter ranges
e.g. $X \rightarrow \gamma\gamma$: $200 < m_X < 2000 \text{ GeV}$, $0 < \Gamma_X < 10\% m_X$.
 → However what comes out of the usual asymptotic formulas is p_{local} .



How to compute p_{global} ? → **Toys** (brute force) or **asymptotic formulas**.

Global Significance from Toys

Principle: repeat the analysis in toy data:

- generate pseudo-dataset
- perform the search, scanning over parameters as in the data
- report the largest significance found
- repeat many times

⇒ The frequency at which a given Z_0 is found *is* the global p-value

e.g. **$X \rightarrow \gamma\gamma$ Search:** $Z_{\text{local}} = 3.9\sigma$ ($\Rightarrow p_{\text{local}} \sim 5 \cdot 10^{-5}$),
scanning $200 < m_X < 2000$ GeV and $0 < \Gamma_X < 10\% m_X$

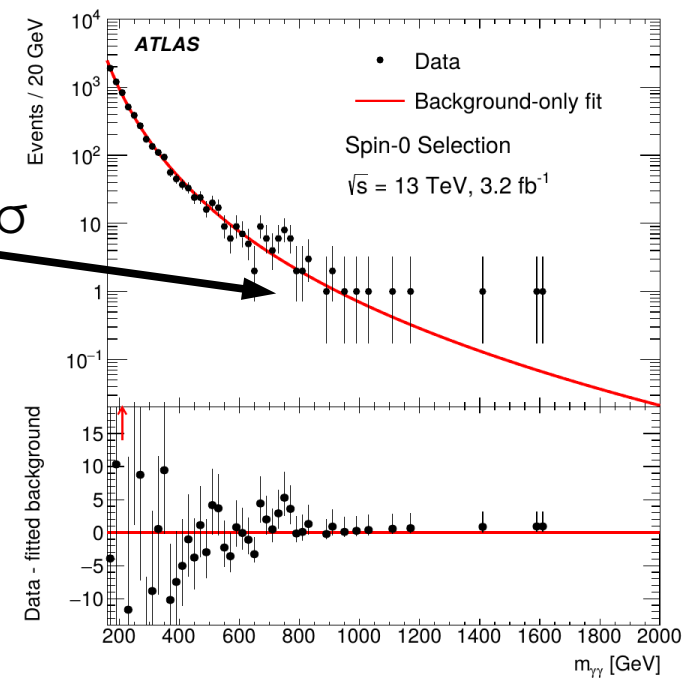
→ In toys, find such an excess 2% of the time

⇒ $p_{\text{global}} \sim 2 \cdot 10^{-2}$, $Z_{\text{global}} = 2.1\sigma$ Less exciting...

⊕ **Exact treatment**

⊖ **CPU-intensive** especially for large Z (need $\sim O(100)/p_{\text{global}}$ toys)

Local 3.9σ



Global Significance from Asymptotics

Principle: approximate the global p-value in the asymptotic limit

→ reference paper: **Gross & Vitells, EPJ.C70:525-530,2010**

$$N_{indep} = \frac{\text{scan range}}{\text{peak width}}$$

Asymptotic trials factor (1 POI):

$$N_{trials} = 1 + \sqrt{\frac{\pi}{2}} N_{indep} Z_{local}$$

→ Trials factor is **not just** N_{indep} ,
also depends on Z_{local} !

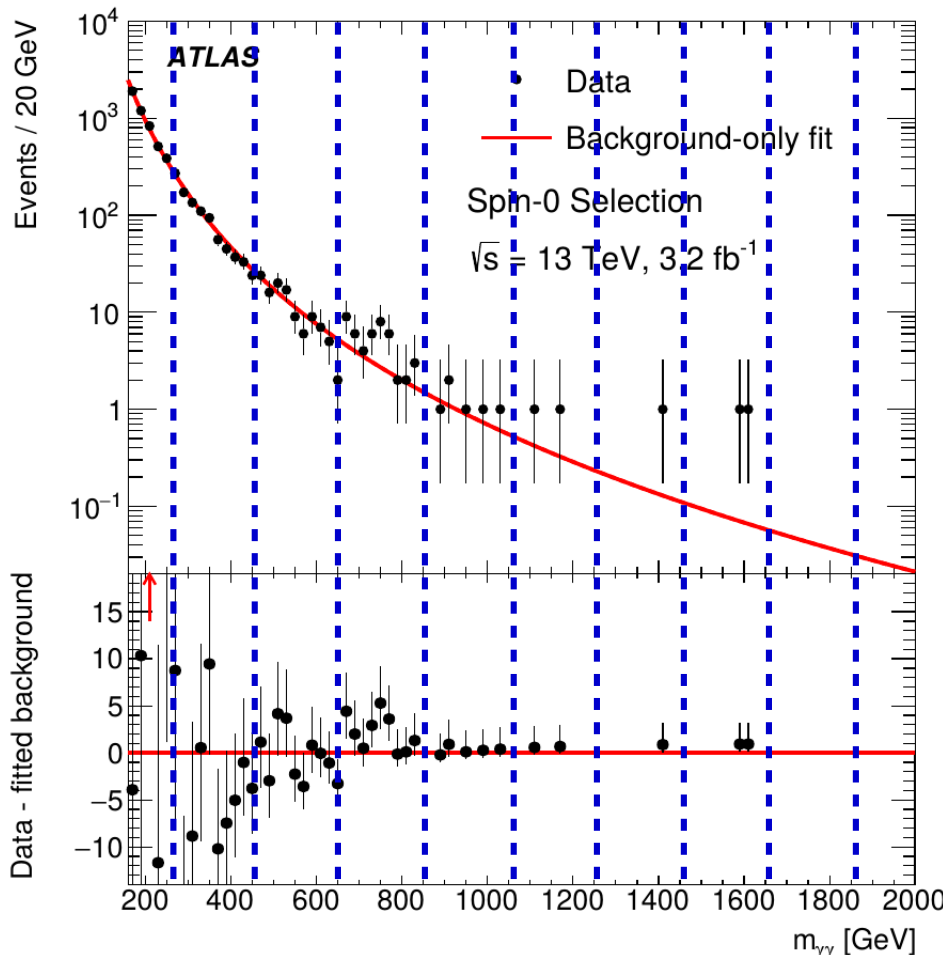
Why ?

- slice scan range into N_{indep} regions of size \sim peak width
- search for a peak in each region

⇒ Indeed gives $N_{trials} = N_{indep}$

However this misses peaks sitting on **edges between regions**

⇒ true N_{trials} is **>** N_{indep} !



Global Significance from Asymptotics

Principle: approximate the global p-value in the asymptotic limit

→ reference paper: **Gross & Vitells, EPJ.C70:525-530,2010**

$$N_{\text{indep}} = \frac{\text{scan range}}{\text{peak width}}$$

Asymptotic trials factor (1 POI):

$$N_{\text{trials}} = 1 + \sqrt{\frac{\pi}{2}} N_{\text{indep}} Z_{\text{local}}$$

→ Trials factor is **not just** N_{indep} ,
also depends on Z_{local} !

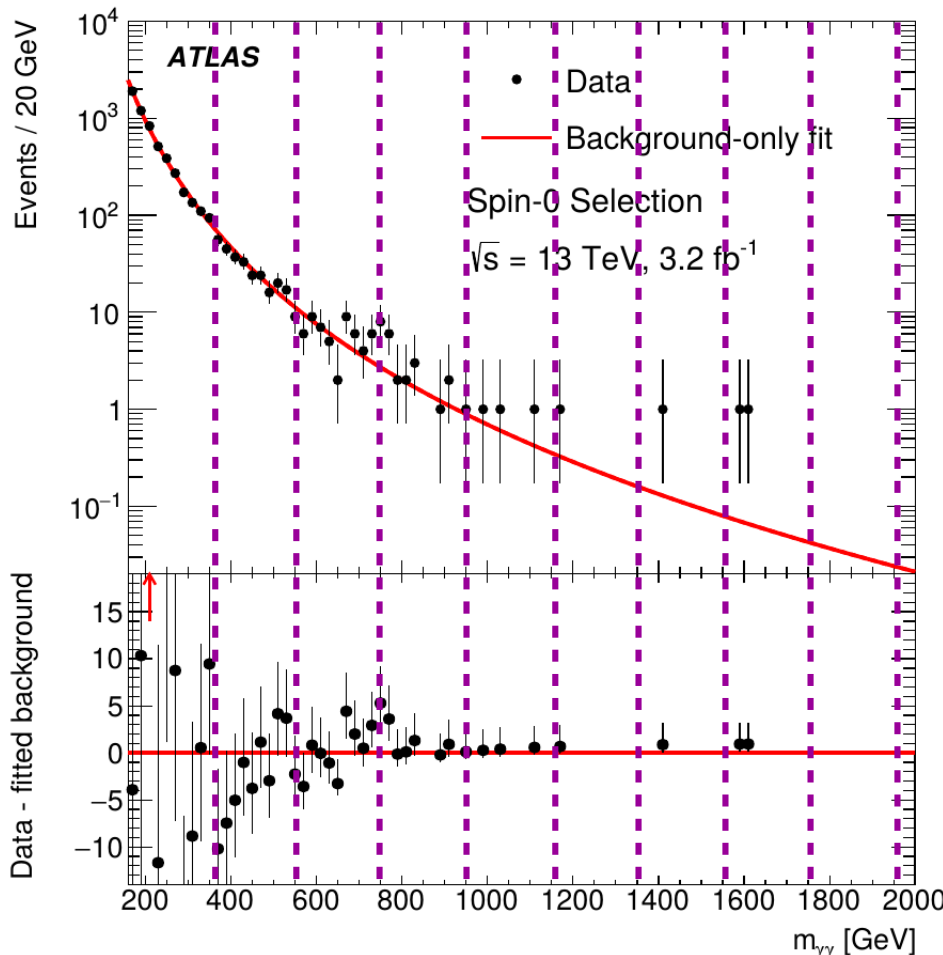
Why ?

- slice scan range into N_{indep} regions of size \sim peak width
- search for a peak in each region

⇒ Indeed gives $N_{\text{trials}} = N_{\text{indep}}$

However this misses peaks sitting on **edges between regions**

⇒ true N_{trials} is **>** N_{indep} !



Illustrative Example

Test on a simple example: generate toys with

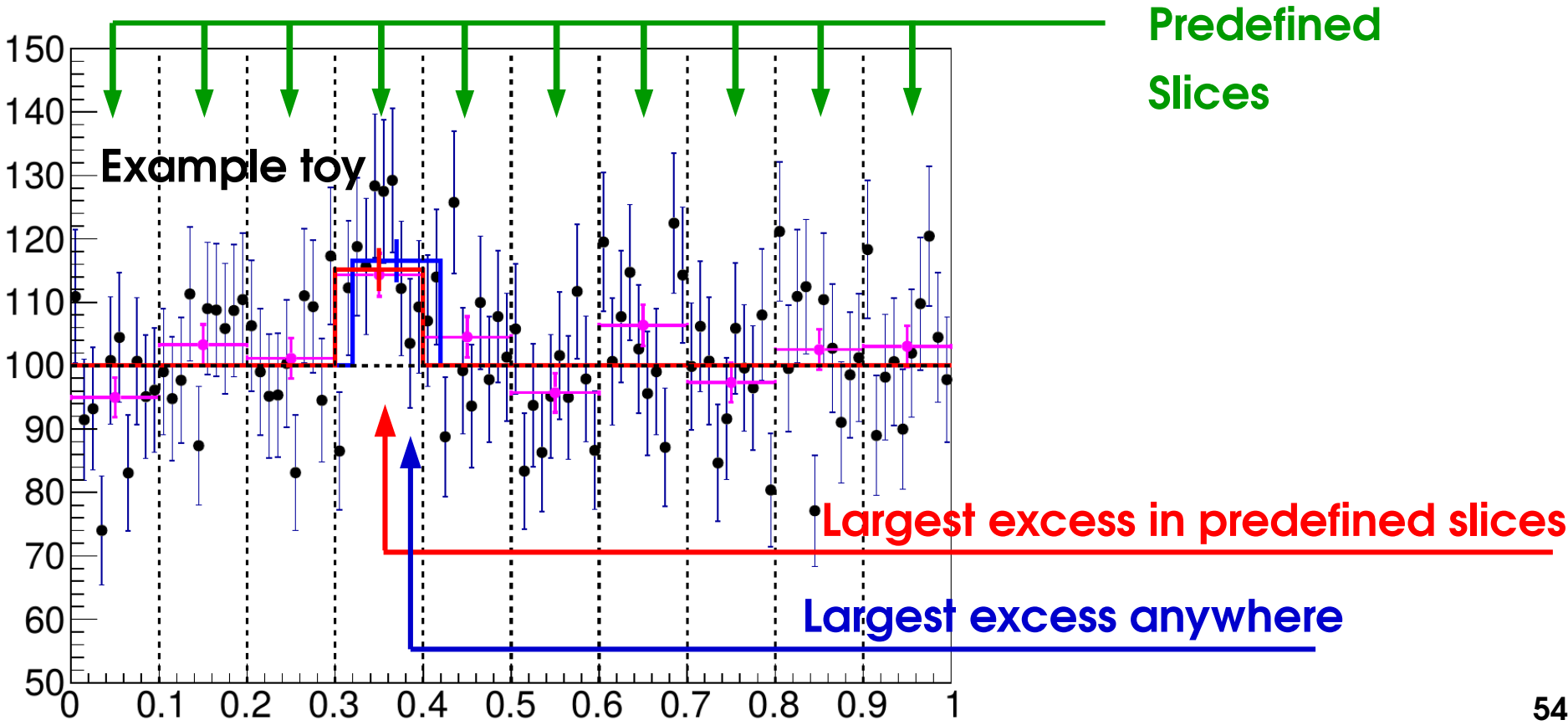
→ flat background (100 events/bin)

→ count events in a fixed-size sliding window, look for excesses

Two configurations:

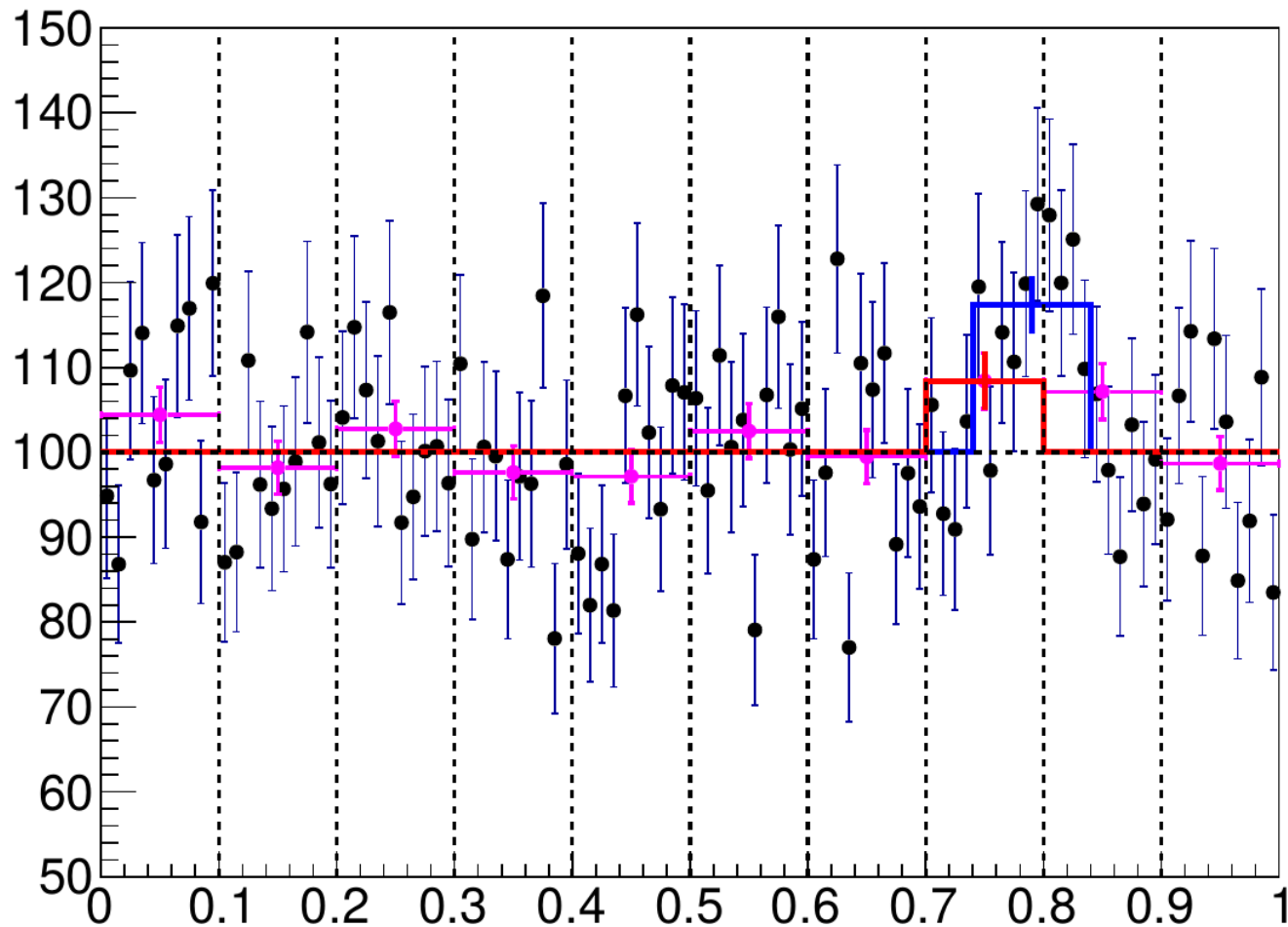
1. Look only in 10 slices of the full spectrum

2. Look in any window of same size as above, anywhere in the spectrum



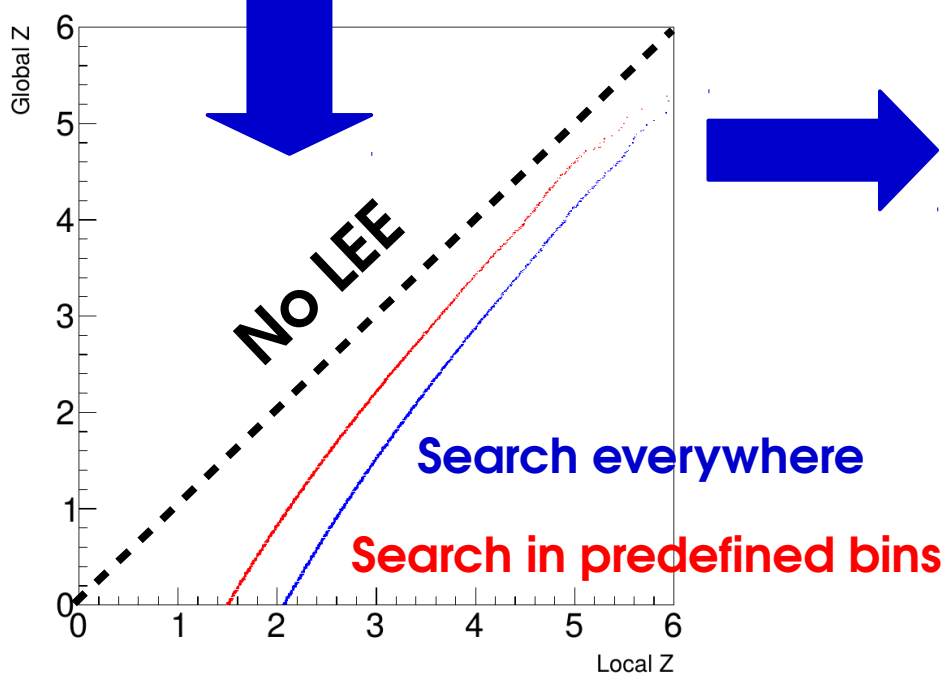
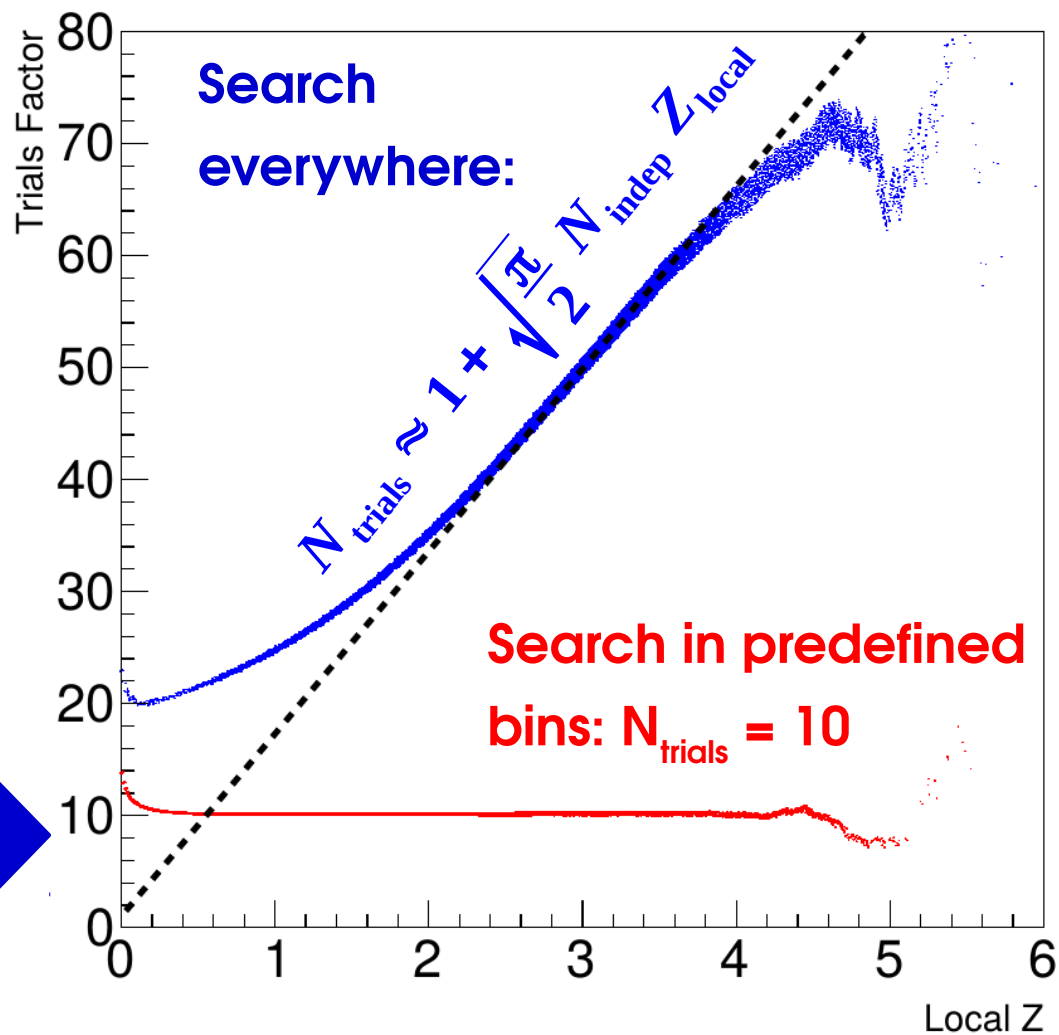
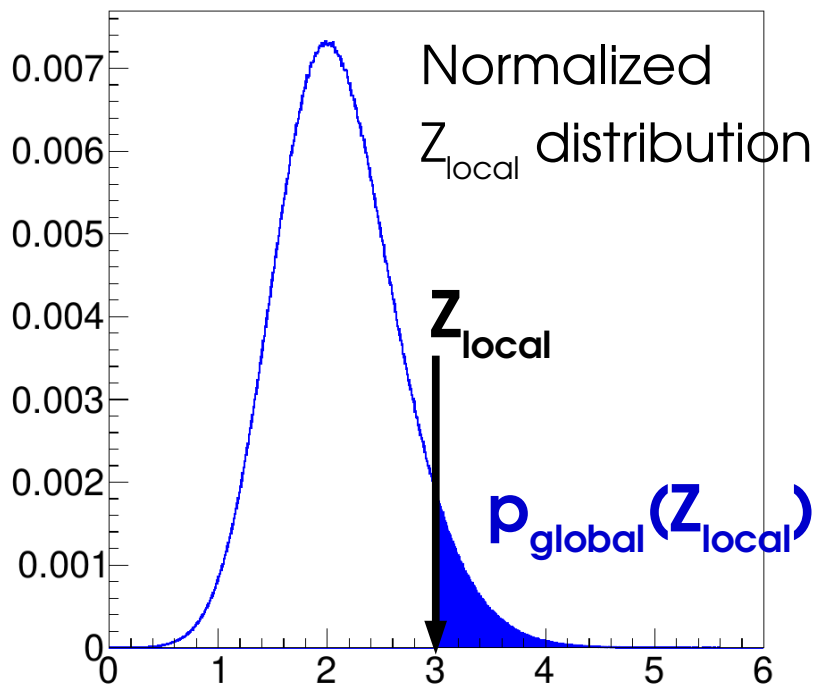
Illustrative Example (2)

Very different results if the excess is **near a boundary** :



1. Look only in 10 slices of the full spectrum
2. Look in any window of same size as above, anywhere in the spectrum

Illustrative Example (3)



Searching everywhere gives the extra Z_{local} dependence

Z_{Global} Asymptotics Extrapolation

Asymptotic trials factor (1 POI):
$$N_{\text{trials}} = 1 + \sqrt{\frac{\pi}{2}} N_{\text{indep}} Z_{\text{local}}$$

How to get N_{indep} ? Usually work with a slightly different formula:

$$N_{\text{trials}} = 1 + \frac{1}{p_{\text{local}}} \langle N_{\text{up}}(Z_{\text{test}}) \rangle e^{\frac{Z_{\text{test}}^2 - Z_{\text{local}}^2}{2}}$$

Number of excesses with $Z > Z_{\text{test}}$

⇒ calibrate for small Z_{test} , apply result to higher Z_{local}

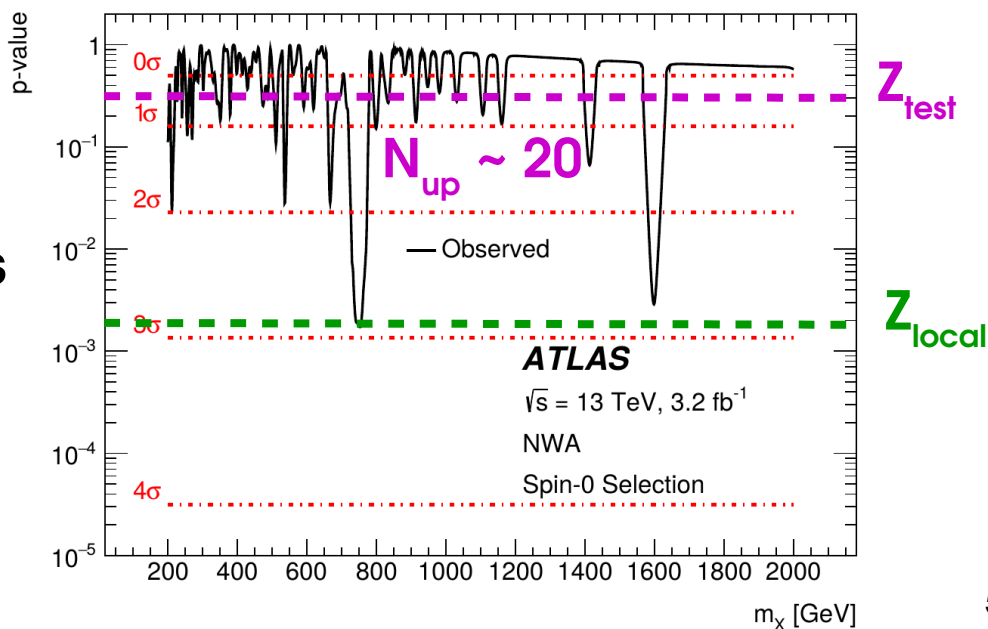
Can choose arbitrarily small Z_{test}

⇒ many excesses

⇒ can measure N_{up} in data (1 “toy”)

Can also measure $\langle N_{\text{up}} \rangle$ in multiple toys

if large stat uncertainty from too few excesses

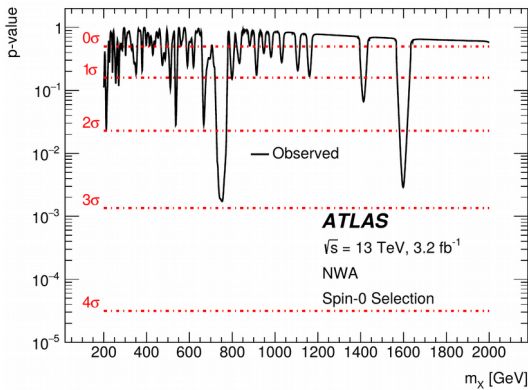


Generalization to 2D scans: consider sections at a fixed Z_{test} , compute its

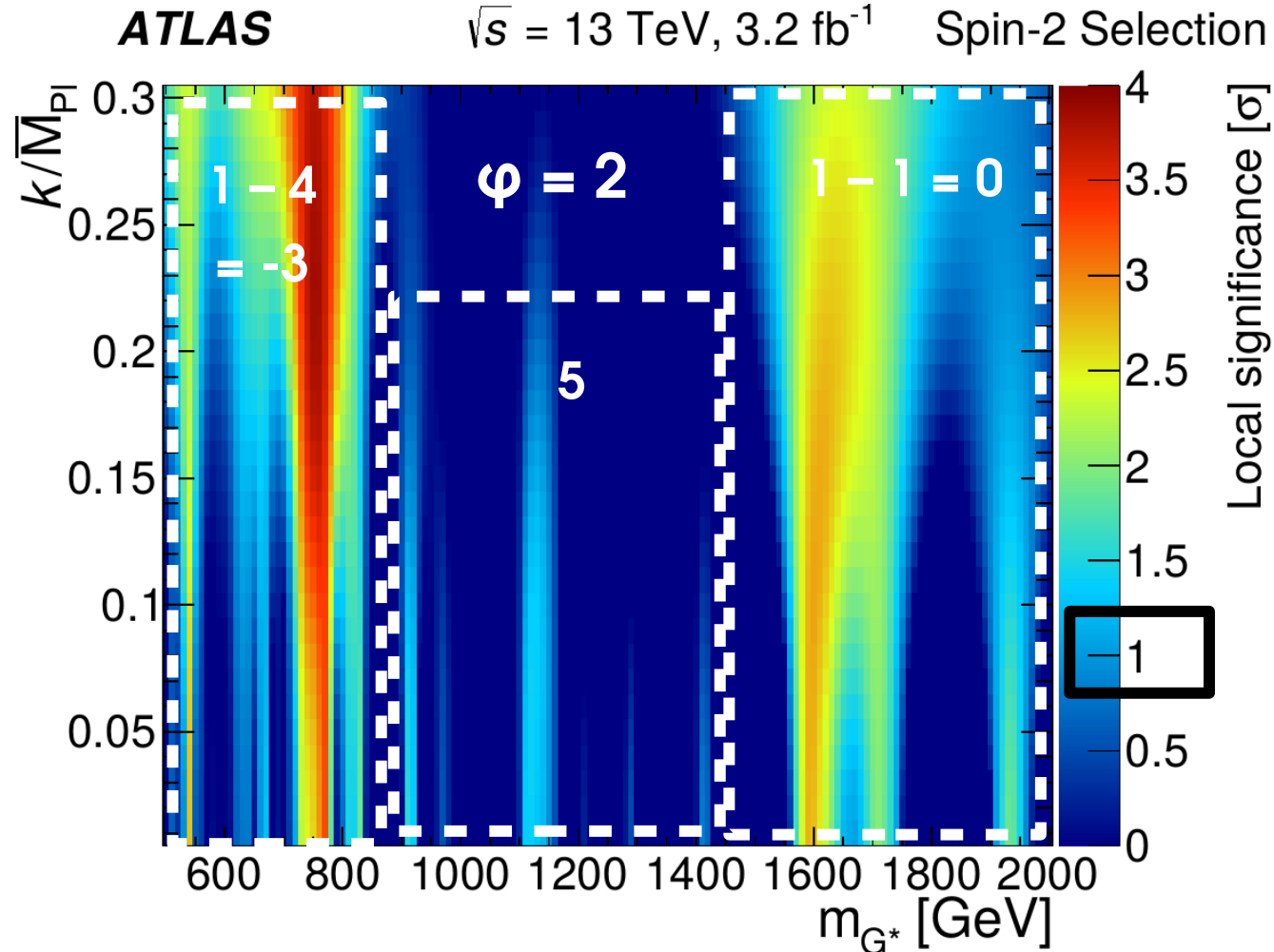
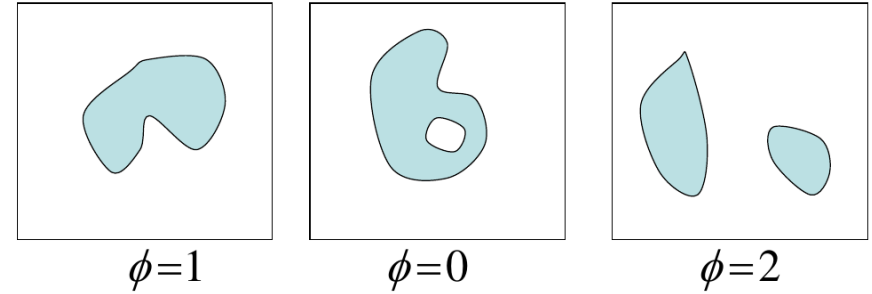
Euler characteristic ϕ , and use

$$p_{\text{global}} \approx E[\phi(A_u)] = p_{\text{local}} + e^{-u/2}(N_1 + \sqrt{u}N_2)$$

→ Generalizes 1D bump counting



Now need to determine 2 constants N_1 and N_2 , from Euler ϕ measurements at 2 different Z_{test} values.



Outline

Computing Statistical Results

Limits, continued

Confidence Intervals

Profiling

Look-Elsewhere Effect

Bayesian methods

Statistical modeling in practice

BLUE

Frequentist vs. Bayesian

All methods described so far are **frequentist**

- Probabilities (p-values) refer to outcomes if the experiment were **repeated identically many times**
- Parameters value are **fixed but unknown**
- Probabilities apply to measurements:

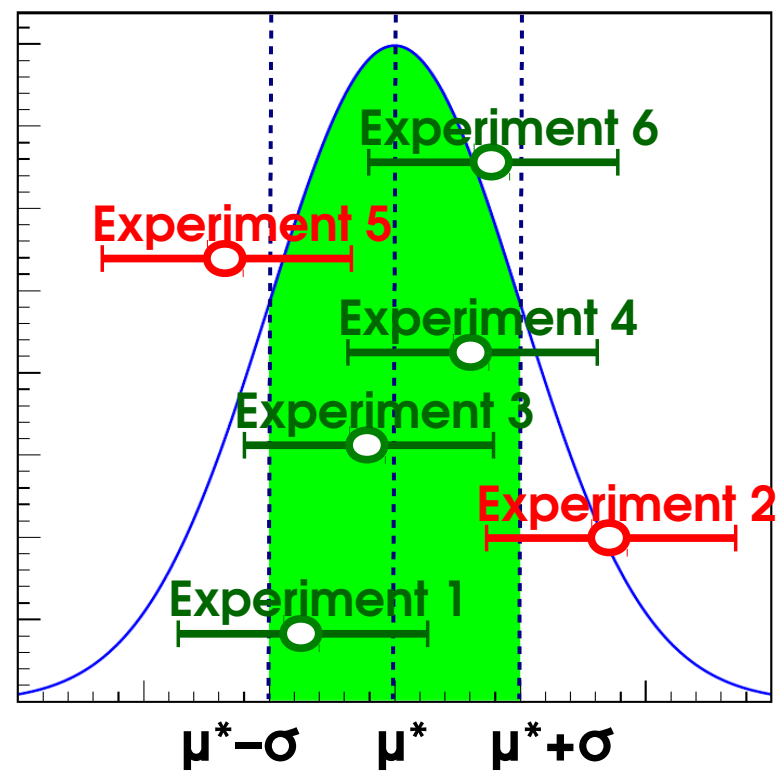
→ “ $m_H = 125.09 \pm 0.24 \text{ GeV}$ ” :

→ i.e. $[125.09 - 0.24 ; 125.09 + 0.24] \text{ GeV}$ has $p=68\%$ to contain **the** true m_H .

→ if we repeated the experiment many times, we would get different intervals, 68% of which would contain the true m_H .

→ “ **5σ Higgs discovery**”

- if there is really no Higgs, such fluctuations observed in $3 \cdot 10^{-7}$ of experiments



Not exactly the crucial question – what we would really like to know is

What is the probability that the excess we see is a fluctuation

→ we want **P(no Higgs | data)** – but all we have is **P(data | no Higgs)**

Frequentist vs. Bayesian

Can use **Bayes' theorem** to address this:

$$P(\mu | data) = \frac{P(data | \mu)}{P(data)} P(\mu)$$

same as in the frequentist formalism (=likelihood)

Prior Probability

irrelevant normalization factor

Can compute $P(\mu | data)$, **if we provide $P(\mu)$**

→ Implicitly, we have now made μ into a random variable

- Is m_H , or the presence of H(125), randomly chosen ?
- In fact, different definition of p: **degree of belief**, not from frequencies.
- $P(\mu)$ **Prior degree of belief** – critical ingredient in the computation

Compared to frequentist PLR:

- ⊕ answers the “right” question
- ⊖ answer depends on the prior

“Bayesians address the questions everyone is interested in by using assumptions that no one believes. Frequentist use impeccable logic to deal with an issue that is of no interest to anyone.” - **Louis Lyons**

Bayesian methods

Probability distribution (= likelihood) : same form as frequentist case, but

P(θ) constraints now **priors for the systematics NPs**, $P(\theta)$

not auxiliary measurements $P(\theta^{\text{mes}}; \theta)$

⊕ Simply integrate them out, no need for profiling: $P(\mu) = \int P(\mu, \theta) d\theta$

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

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

⊖ No simple way to test for discovery

⊖ Integration over NPs can be CPU-intensive

Priors : most analyses still using flat priors in the analysis variable(s)

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

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

Frequentist-Bayesian Hybrid methods (“Cousins-Highland”)

- Integrate out NPs as in Bayesian measurements

- Once only POIs left, Use $P(\text{data} | \mu)$ in a frequentist way

→ “Bayesian NPs, frequentist POIs”

- Some use in Run 1, now phased out in favor of frequentist PLR.

Bayesian methods and CL_s : CL_s computation

Gaussian counting with systematic on background: $n = S + B + \sigma_{\text{syst}} \theta$

$$L(n; S, \theta) = G(n; S + B + \sigma_{\text{syst}} \theta, \sigma_{\text{stat}}) G(\theta_{\text{obs}} = 0; \theta, 1)$$

$$\text{MLE: } \hat{S} = n - B$$

$$\text{Conditional MLE: } \hat{\theta}(\mu) = \frac{\sigma_{\text{syst}}}{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2} (n - S - B) \quad \left. \vphantom{\hat{\theta}(\mu)} \right\} \text{PLR: } \lambda(\mu) = \left(\frac{S + B - n}{\sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}} \right)^2$$

Gaussian \Rightarrow from previous studies, CL_s limit is

$$CL_s: \quad S_{\text{up}}^{CL_s} = n - B + \left[\Phi^{-1} \left(1 - 0.05 \Phi \left(\frac{n - B}{\sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}} \right) \right) \right] \sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}$$

Bayesian methods and CL_s: Bayesian case

Gaussian counting with systematic on background: $n = S + B + \sigma_{\text{syst}} \theta$

$$P(n | S, \theta) = G(n; S + B + \sigma_{\text{syst}} \theta, \sigma_{\text{stat}}) G(\theta | 0, 1)$$

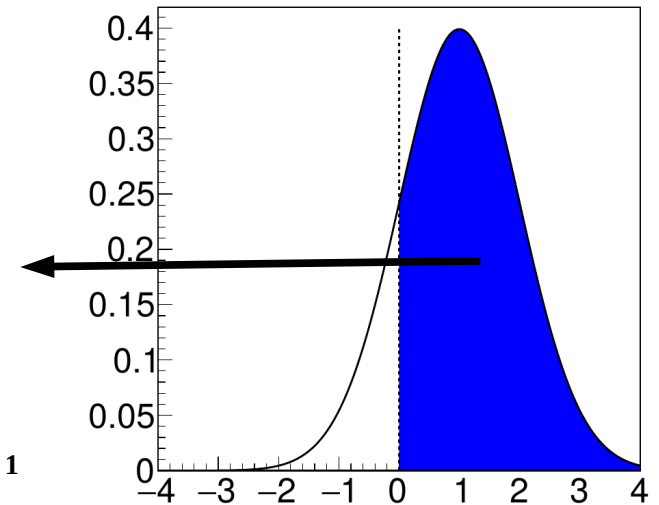
Bayesian: $G(\theta)$ is actually a **prior** on $\theta \Rightarrow$ perform integral (**marginalization**)

$$P(n | S) = G(S; n - B, \sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}) \quad \text{same effect as profiling!}$$

Need $P(S | n) \Rightarrow$ a prior for S – take flat PDF over $S > 0$

\Rightarrow Truncate Gaussian at $S=0$: $P(S | n) = P(n | S) P(S)$

$$P(S | n) = G(S; n - B, \sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}) \left[\Phi \left(\frac{n - B}{\sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}} \right) \right]^{-1}$$



Bayesian Limit:

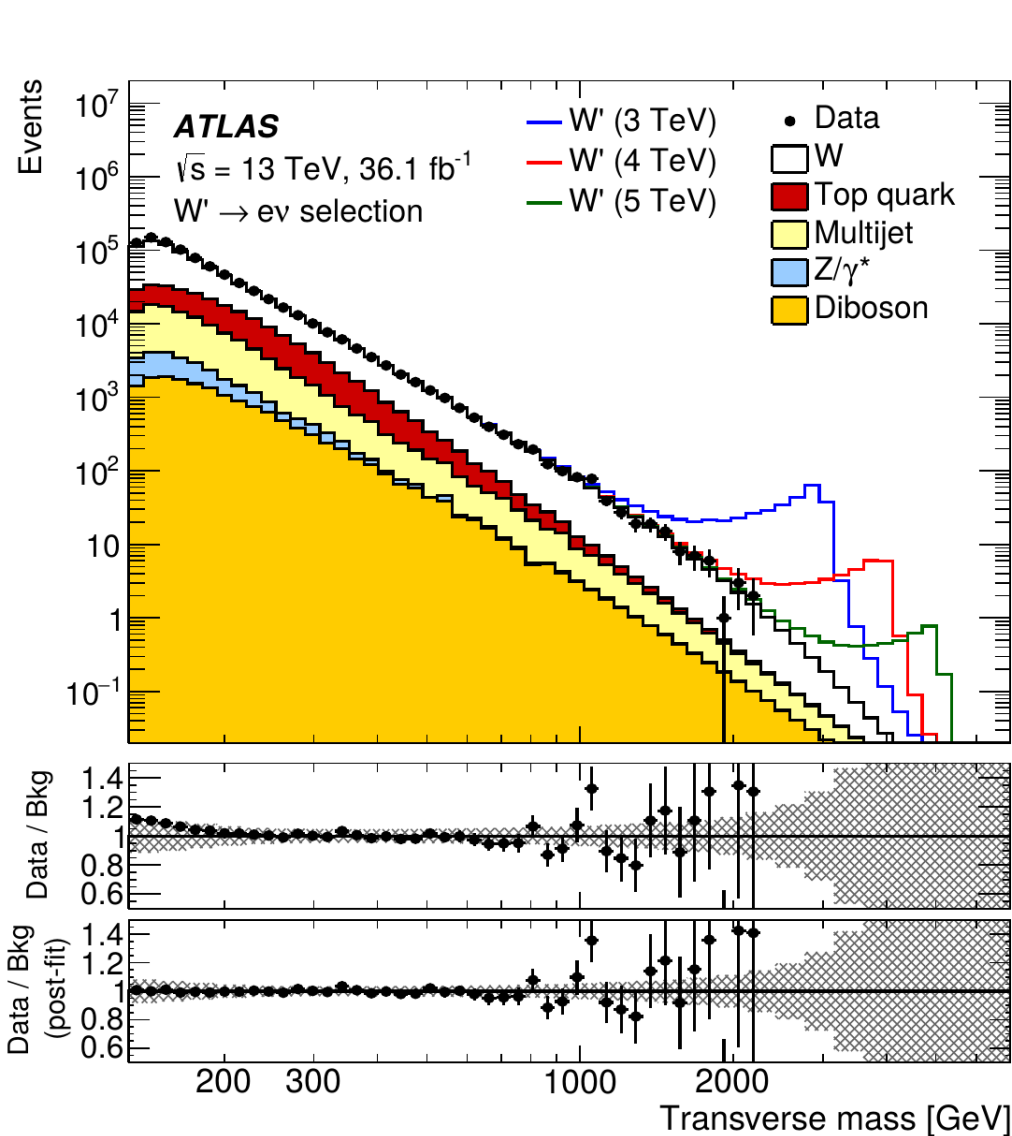
$$\int_{S_{\text{up}}}^{\infty} P(S | n) dS = 5\% = \left[1 - \Phi \left(\frac{S_{\text{up}} - (n - B)}{\sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}} \right) \right] \left[\Phi \left(\frac{n - B}{\sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}} \right) \right]^{-1}$$

$$S_{\text{up}}^{\text{Bayes}} = n - B + \left[\Phi^{-1} \left(1 - 0.05 \Phi \left(\frac{n - B}{\sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}} \right) \right) \right] \sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}$$

same result as CL_s!

Example: $W' \rightarrow l\nu$ Search

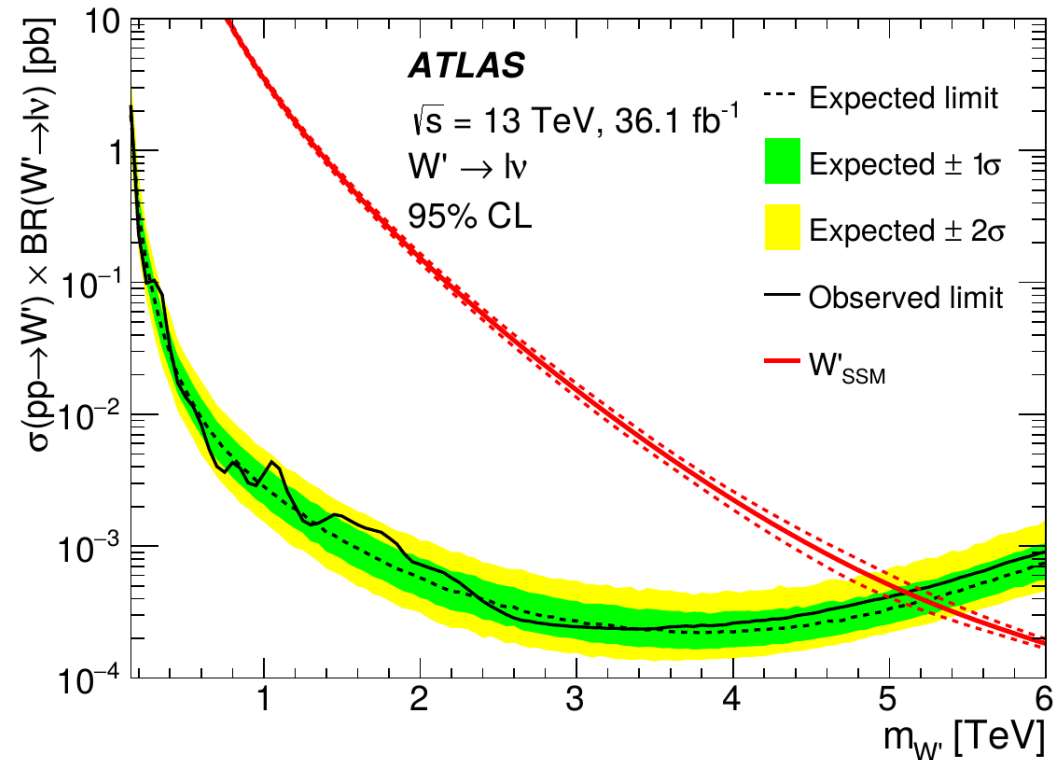
- **POI:** $W' \sigma \times B \rightarrow$ use flat prior over $[0, +\infty[$.
- **NPs:** syst on **signal ϵ** (6 NPs), **bkg** (6), **lumi** (1) \rightarrow integrate over Gaussian priors



Trigger
 Lepton reconstruction and identification
 Lepton momentum scale and resolution
 E_T^{miss} resolution and scale
 Jet energy resolution
 Pile-up

Multijet background
 Top extrapolation
 Diboson extrapolation
 PDF choice for DY
 PDF variation for DY
 EW corrections for DY

Luminosity



Why 5σ ?

One-sided discovery: $5\sigma \Leftrightarrow p_0 = 3 \cdot 10^{-7} \Leftrightarrow 1 \text{ chance in } 3.5\text{M}$

→ Overly conservative ?

→ Do we even know the sampling distributions so far out ?

Local 3.9σ , $p_0 = 5\text{E-}5$

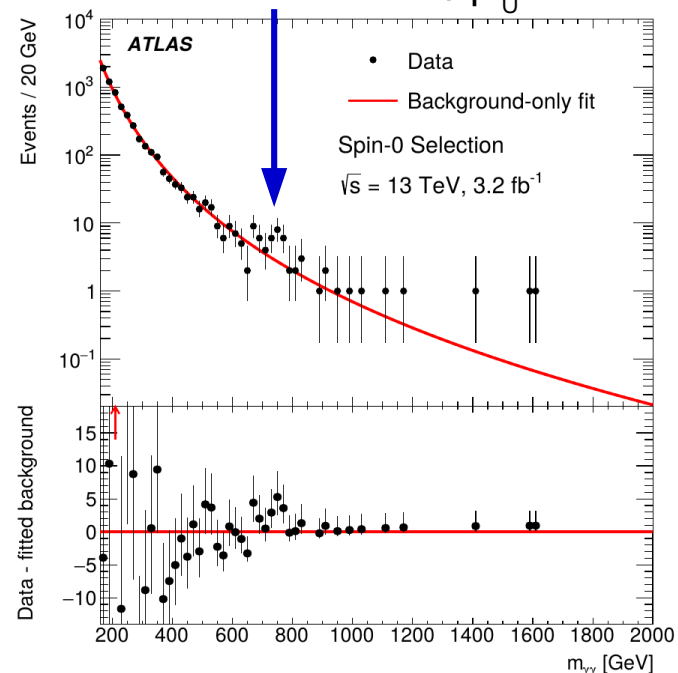
Global 2.1σ , $p_0 = 2\text{E-}2$

Reasons for sticking with 5σ (from Louis Lyons):

- **LEE** : searches typically cover multiple independent regions
 ⇒ Global p-value is the relevant one

$N_{\text{trials}} \sim 1000$: **local $5\sigma \Leftrightarrow O(10^{-4})$** more reasonable

- **Mismodeled systematics**: factor 2 error in syst-dominated analysis ⇒ factor 2 error on Z...
- **History**: 3σ and 4σ excesses do occur regularly, for the reasons above
- **“Subconscious Bayes Factor”** : p-value should be at least as small as the subjective $p(S)$:



$$P(\text{fluct}) = \frac{P(\text{fluct}|B) P(B)}{P(\text{fluct}|S) P(S) + P(\text{fluct}|B) P(B)}$$

Extraordinary claims require extraordinary evidence

⇒ **Stay with 5σ ...**

Outline

Profiling

Look-Elsewhere Effect

Bayesian methods

Statistical modeling in practice

Building binned likelihoods

Choosing PDFs in unbinned likelihoods

Implementing systematics

BLUE

Statistical Modeling: in Practice

Building statistical models

So far focus has been on concepts, but building a statistical model also requires **numerical** inputs:

- **Data PDFs** for all model components
- **Constraint PDFs** for all sources systematics
- **Impact** of each systematic uncertainty on all relevant model parameters

→ Statistical methods are only as accurate (and/or optimal) as the description provided by the model!

Technically, MC simulation provides most of these inputs. However 2 problematic issues:

- **Potential MC/data differences**
- **Limited MC statistics**

Which need to be addressed with (yet more) systematics.

Statistical Modeling:

I. Component PDFs

PDFs : Binned likelihood

Binned case:

→ PDF usually just a normalized histogram, from

- **MC sample** or
- **Data control region (CR)**

⇒ **Statistical uncertainties** on the prediction:

- **Data CR:** counts as **statistical** uncertainty
- **MC sample:** uncertainty can be reduced without collecting more data (just need more CPU!) ⇒ Counted as **systematic**

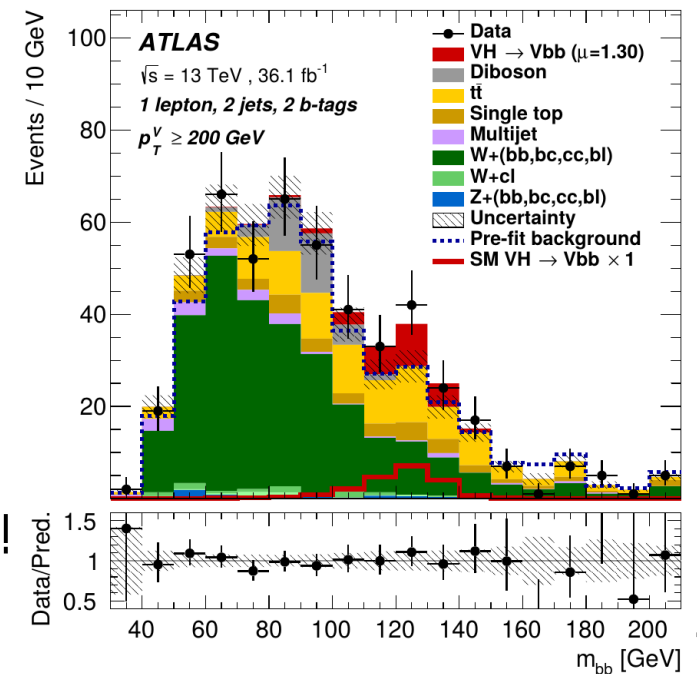
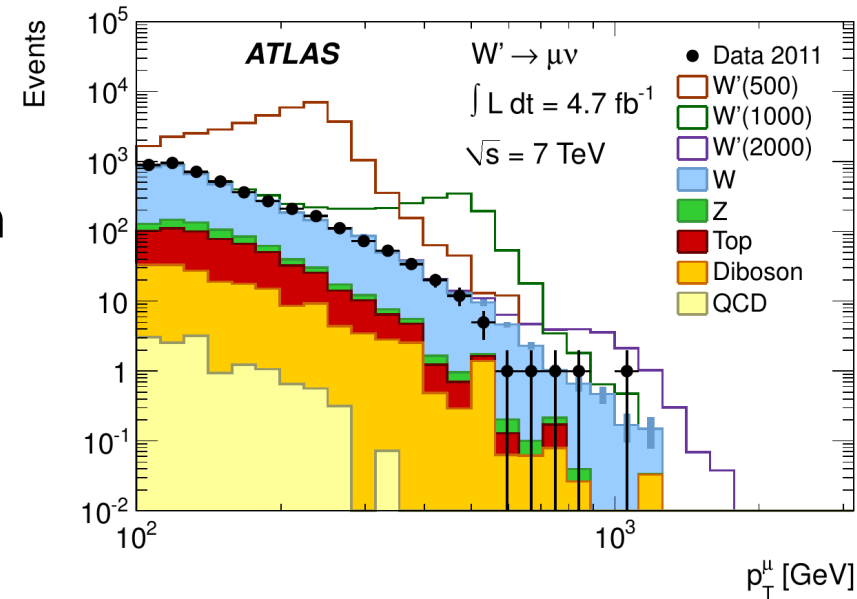
Independent counts in each bin

⇒ a separate **MC statistics NP** in each bin

→ Poisson constraints $\text{Pois}(N_i^{\text{MC}}; N_i^{\text{true}})$

$$\text{Total uncertainty} \sim \sqrt{\sigma_{\text{data stats}}^2 + \sigma_{\text{MC stats}}^2 + \dots}$$

⇒ need enough MC to avoid spoiling the sensitivity!



MC Statistics Requirements

e.g. **Discovery**: Total uncertainty: $\sigma_S^2 \sim \sqrt{\sigma_{\text{data stats}}^2 + \sigma_{\text{MC stats}}^2 + \dots}$

⇒ need $\sigma_{\text{MC stats}} \ll \sigma_{\text{data stats}}$
 $B_{\text{MC}} \gg B_{\text{data}}$

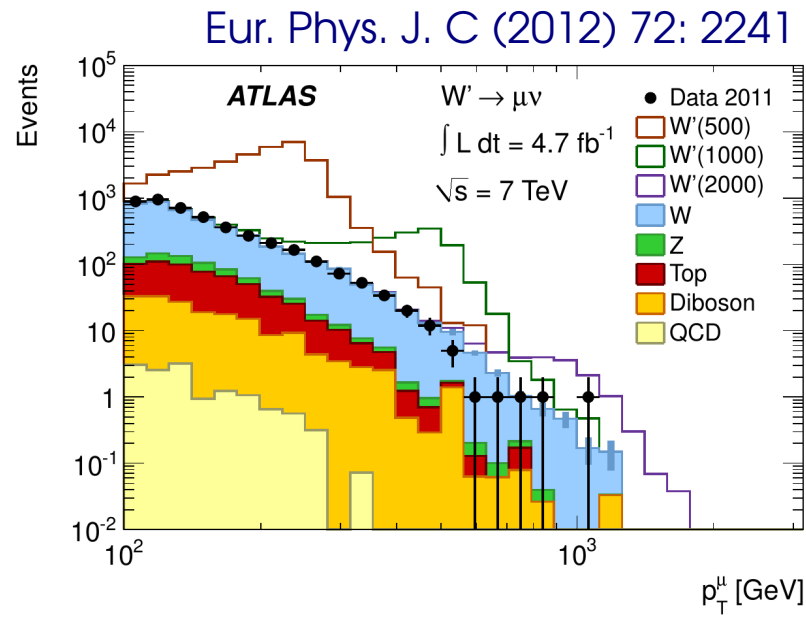
$B_{\text{MC}}/B_{\text{data}}$ (α)	$\sigma_{\text{MC stats}}/\sigma_{\text{data stats}}$ ($1/\sqrt{\alpha}$)	$\sigma_{\text{data+MC stats}}/\sigma_{\text{data stats}}$ ($\sqrt{1+\alpha^{-1}}$)
1	1	1.41
4	0.5	1.12
25	0.2	1.02

By how much ?

In the presence of a signal (e.g. limit-setting, N_{sig} measurement), relevant uncertainty is $\sqrt{(S+B)}$
 ⇒ **S/B** also matters:

$$\frac{\sigma_S}{S} \sim \sqrt{1 + \frac{S}{B} + \frac{B_{\text{data}}}{B_{\text{MC}}} \frac{1}{1+S/B}}$$

- **low S/B** : same problem as for discovery
- **high S/B** : no issue, dominated by uncertainty on signal itself.



PDF shapes: Unbinned likelihood

Smooth backgrounds : Describe distribution using appropriate **function**
 ⇒ **Unbinned likelihood**. Describes sideband + signal region in one fit.

Phys. Rev. Lett. 118 (2017), 191801

Phys. Lett. B241 (1990) 278-282

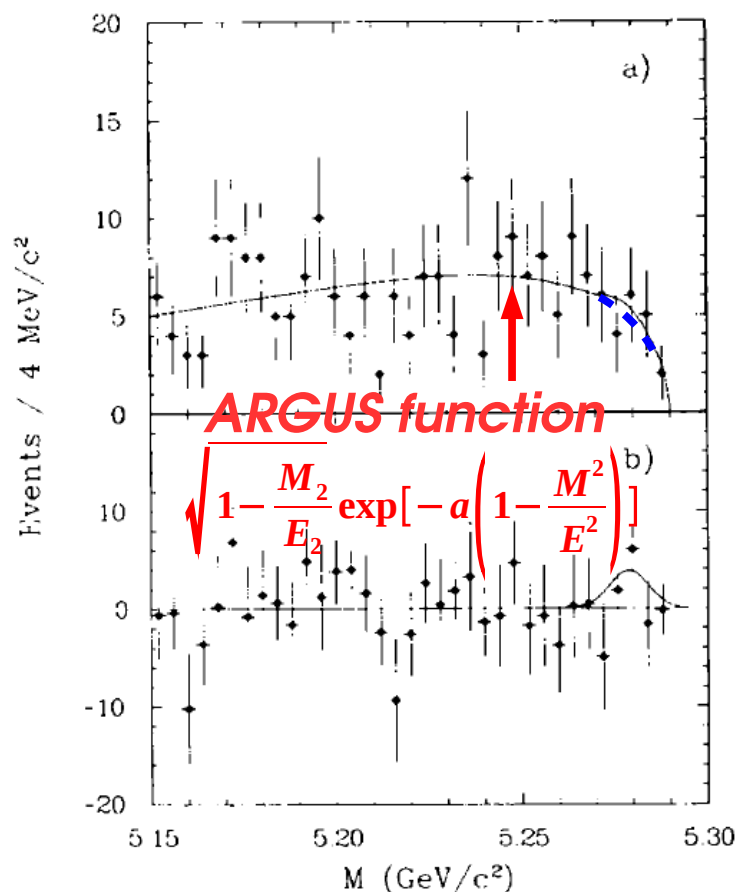
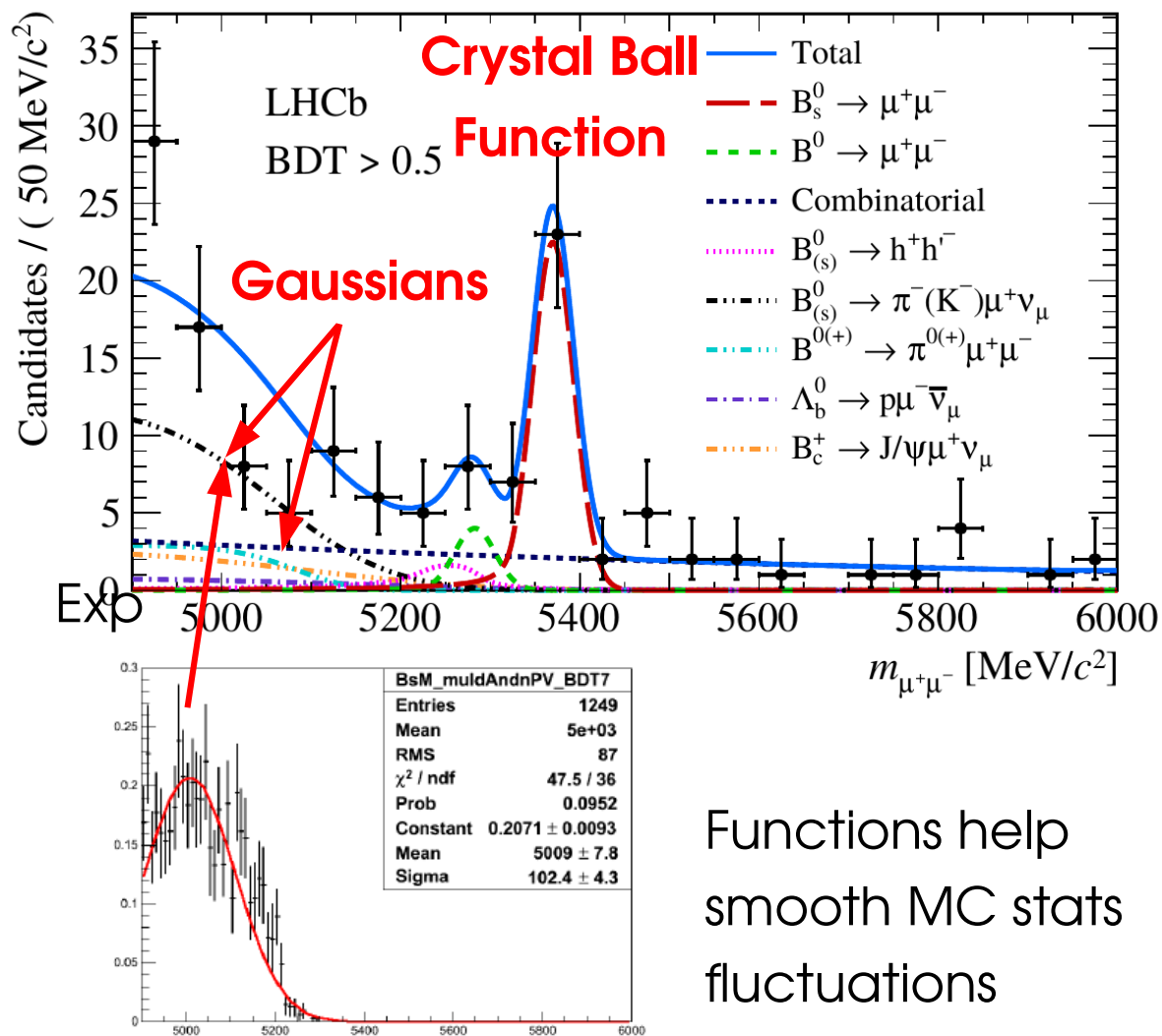


Fig. 1. Invariant mass distribution of the decay $B^+ \rightarrow \pi^+ \pi^0$. (a) At the $\Upsilon(4S)$; the curve shows the result of the maximum likelihood fit described in the text. (b) After subtraction of the continuum contribution. The gaussian curve represents the 90% CL upper limit on the signal from the above fit (see table 1).

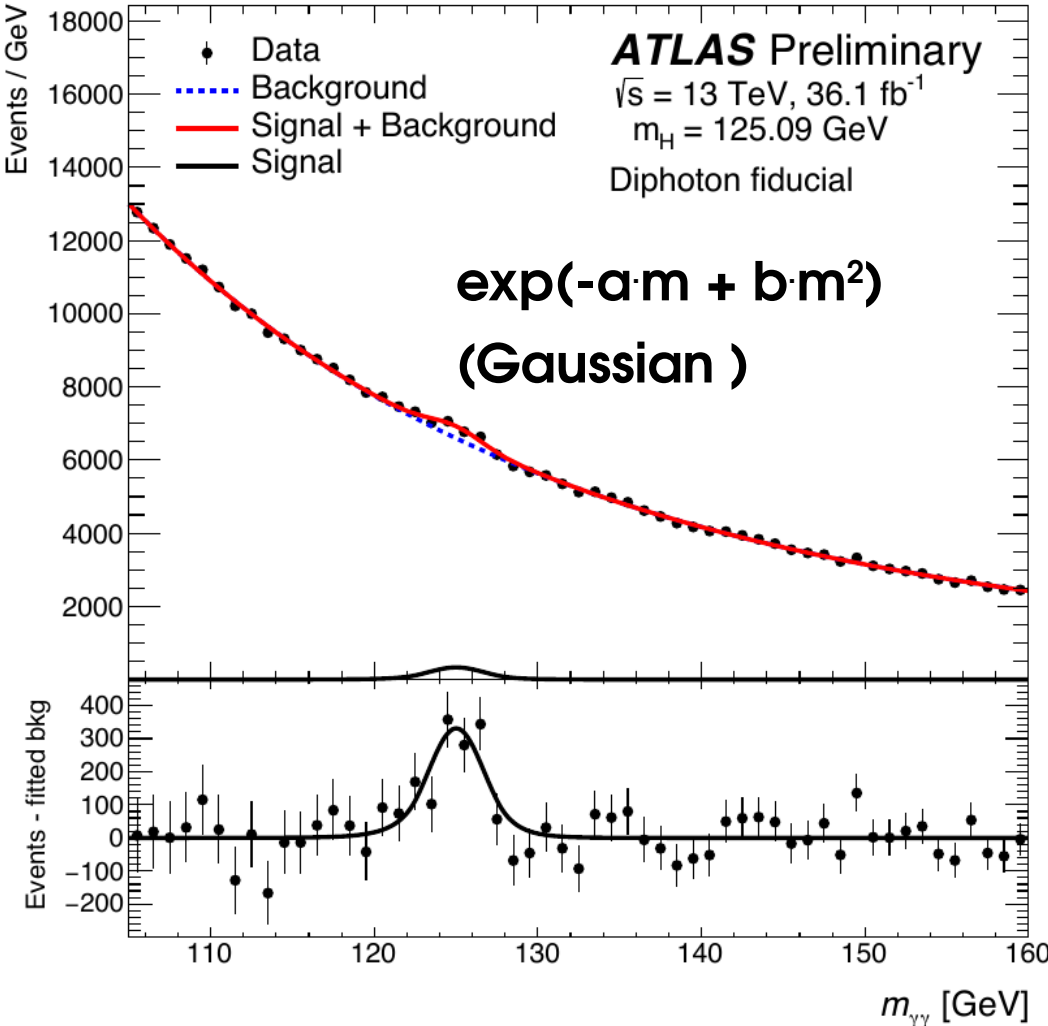


Functions help
smooth MC stats
fluctuations

PDF Shapes: Unbinned likelihood

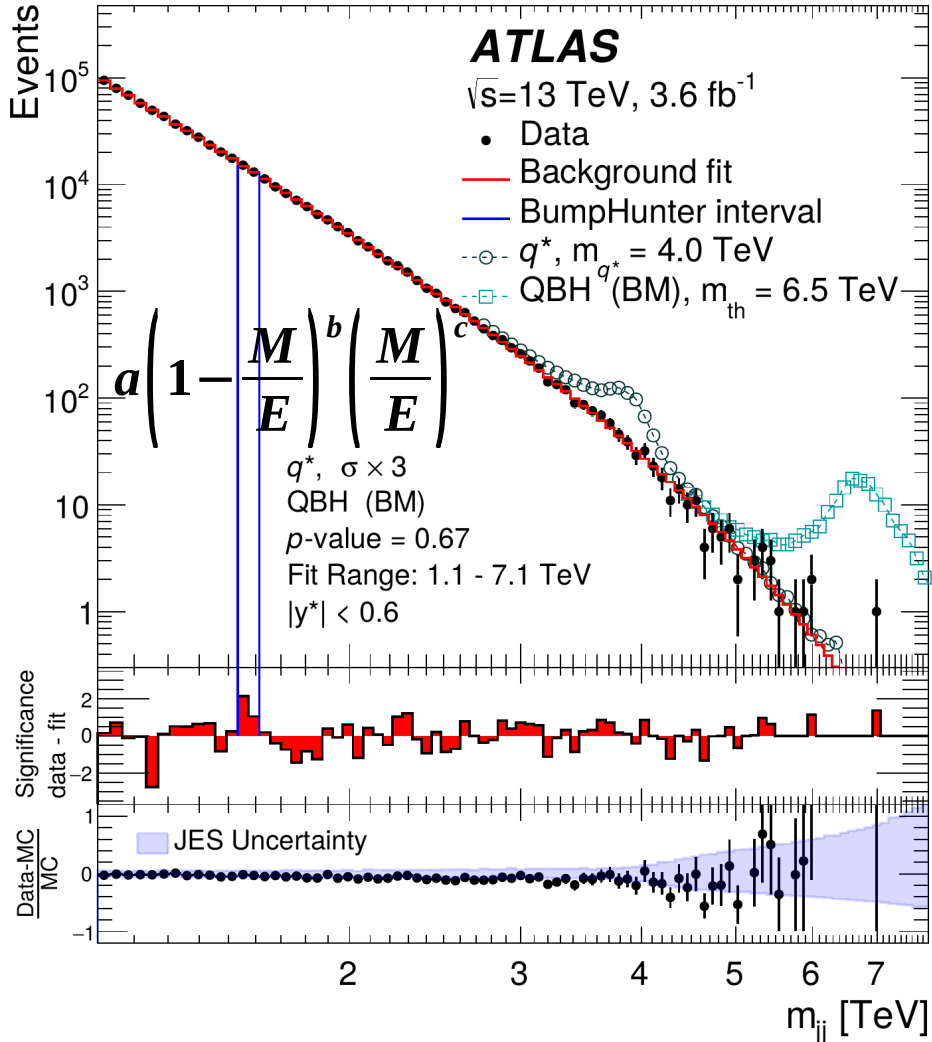
Widely used in HEP for smooth backgrounds (→ no resonances or threshold)

H → γγ Measurements



X → jj Search

Phys.Lett. B754 (2016) 302-322



Signal Bias in Unbinned likelihoods

Function usually ad-hoc (no closed form expression for (theory \otimes detector effects), or usually even theory by itself...)

→ **may not accurately describe the data**

⇒ **Introduce free parameters, fit in sidebands**

→ **Biases may still remain due to functional form itself**

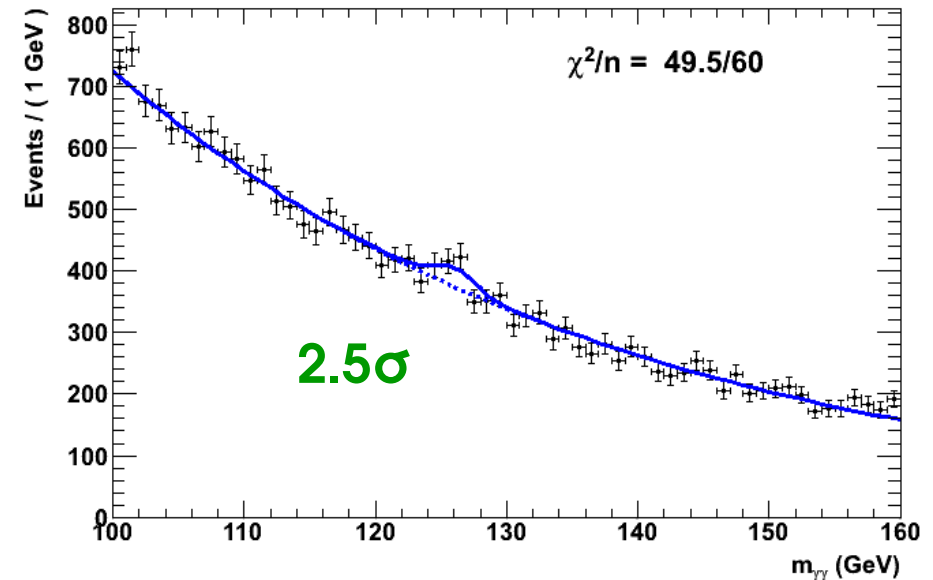
Jan 2012 Higgs search paper
(4.9 fb⁻¹ of 2011 data)

exponential

Problematic especially for **low S/B**

→ small mismodelings of B can be large compared to S.

→ **χ^2 test in sideband may not help**: even a large bias on the scale of S (\ll B) may remain within stat errors in the sideband!



Situation doesn't improve with more luminosity:

→ Reduced statistical uncertainties in sideband, but

→ Also reduced σ_s , in the same proportion

Signal Bias in Unbinned likelihoods

Function usually ad-hoc (no closed form expression for (theory \otimes detector effects), or usually even theory by itself...)

→ **may not accurately describe the data**

⇒ **Introduce free parameters, fit in sidebands**

→ **Biases may still remain due to functional form itself**

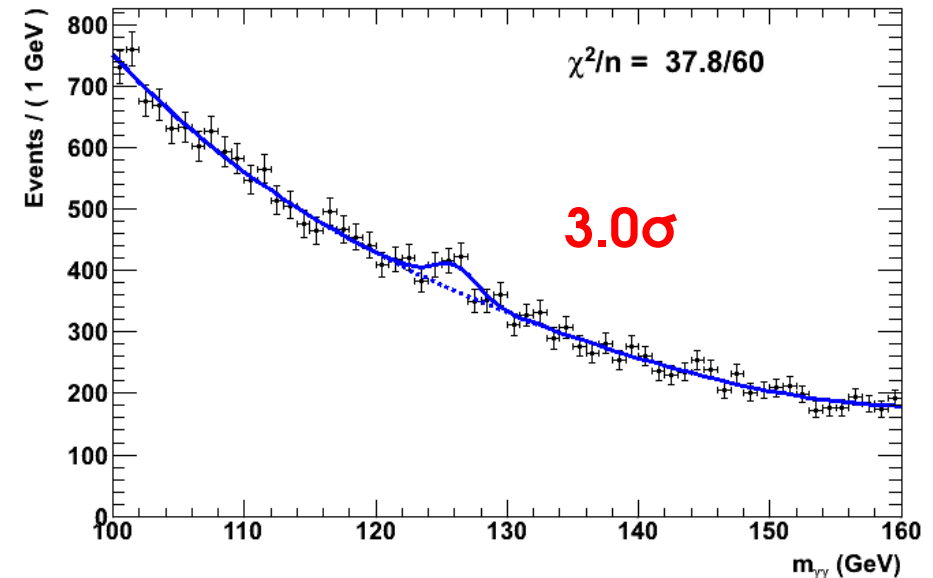
Jan 2012 Higgs search paper
(4.9 fb⁻¹ of 2011 data)

polynomial

Problematic especially for **low S/B**

→ small mismodelings of B can be large compared to S.

→ **χ^2 test in sideband may not help**: even a large bias on the scale of S (\ll B) may remain within stat errors in the sideband!



Situation doesn't improve with more luminosity:

→ Reduced statistical uncertainties in sideband, but

→ Also reduced σ_s , in the same proportion

Signal Bias in Unbinned likelihoods

If data cannot fix B shape, **use MC**
 → Measure signal bias N_{SS} on “credible” shapes taken from MC (**Spurious signal**)
 → take the maximum bias as systematic

Works well if the true distribution is somewhere in the space of MC distributions scanned...

Also **Impose:**

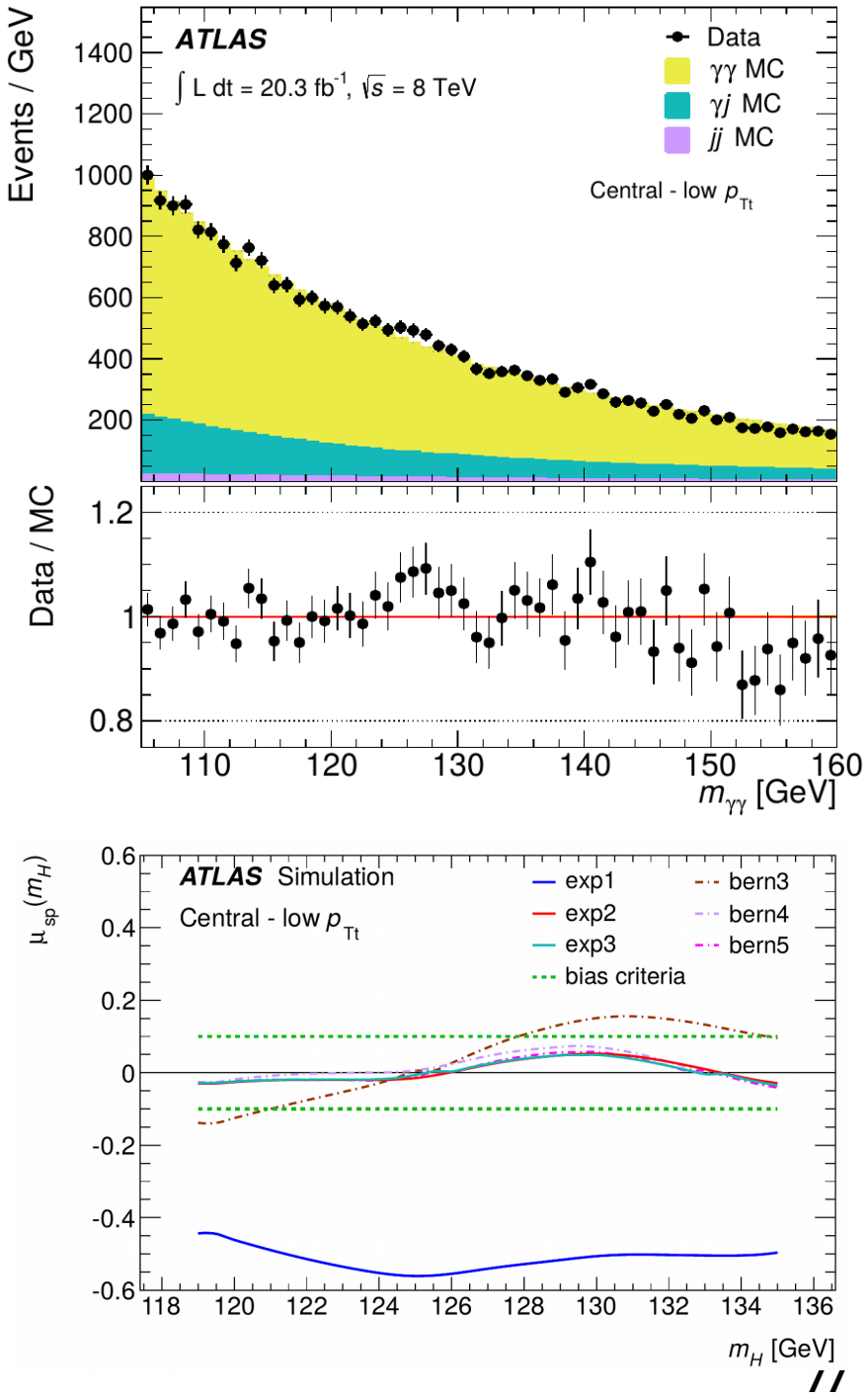
$N_{SS} < 20\% \sigma_{stat}$ (small contribution to σ_{total})

OR

$N_{SS} < 10\% S_{exp}$ (small bias on measured S)

Second criterion more stringent at higher S/\sqrt{B} .

If criteria are not met, move to more complex functions (→ more free parameters)



Signal Bias in Unbinned likelihoods

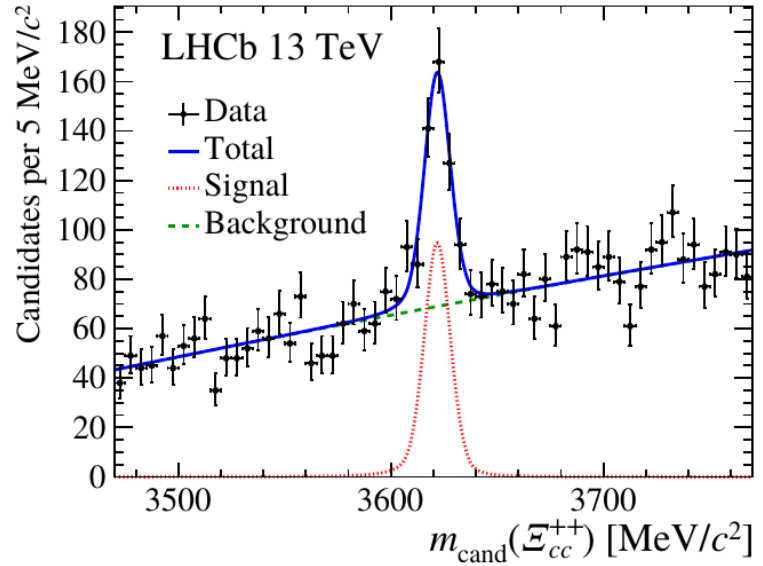
Problem: for small MC stats, measured bias dominated by fluctuations
 → again, need high MC stats ($B_{MC} > 25 B_{data}$) when S/B is low.

B_{MC}/B_{data} (α)	$\sigma_{MC\ stats}/\sigma_{data\ stats}$ ($1/\sqrt{\alpha}$)	$\sigma_{data+MC\ stats}/\sigma_{data\ stats}$ ($\sqrt{1+\alpha^{-1}}$)
1	100%	1.41
4	50%	1.12
25	20%	1.02

← $N_{ss} < 20\% \sigma_{stat}$

→ Can compromise on criterion level
 (50% instead of 20% ?)

→ As before, better situation at at high S/B



Phys. Rev. Lett. 118, 182001 (2017)

Polynomials: various basis choices (Chebyshev, Bernstein,...)

Bernstein basis:

$$B_{k,n}(x) = \binom{n}{k} x^k (1-x)^{n-k} \text{ for } 0 \leq x \leq 1$$

→ **Positive coefficients** ⇒ **positive polynomial everywhere**, useful to avoid numerical issues in $-2 \log(\text{PDF})$ computation

Exponential family : $\exp(\text{polynomial})$

Power laws : $x^\alpha, x^\alpha(1-x)^\beta, \dots$

Gaussians

Crystal Ball Functions

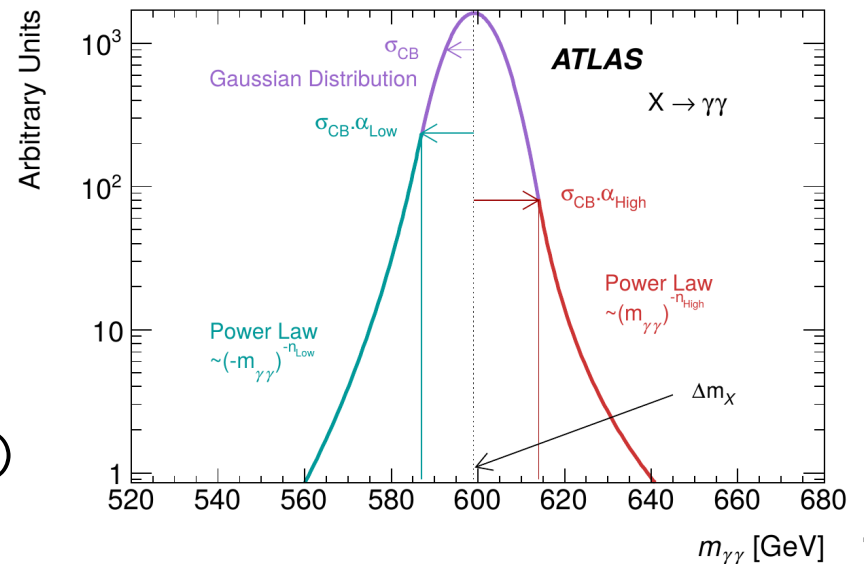
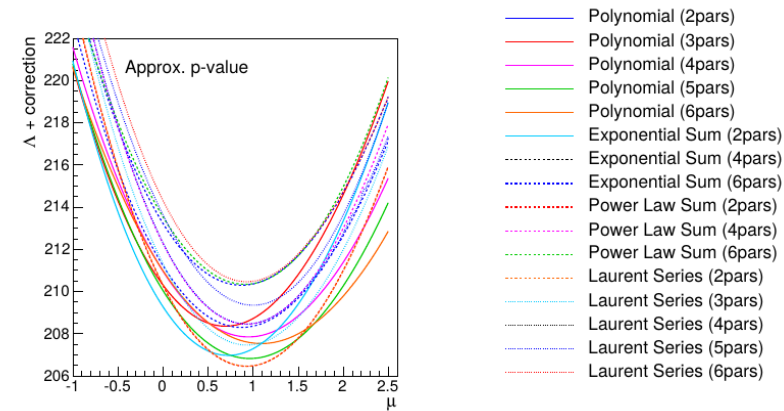
$$t = (m_{\gamma\gamma} - \mu_{\text{CB}}) / \sigma_{\text{CB}}$$

$$N \cdot \begin{cases} e^{-0.5t^2} & \text{if } -\alpha_{\text{low}} \leq t \leq \alpha_{\text{high}} \\ e^{-0.5\alpha_{\text{low}}^2} \left[\frac{\alpha_{\text{low}}}{n_{\text{low}}} \left(\frac{n_{\text{low}}}{\alpha_{\text{low}}} - \alpha_{\text{low}} - t \right) \right]^{-n_{\text{low}}} & \text{if } t < -\alpha_{\text{low}} \\ e^{-0.5\alpha_{\text{high}}^2} \left[\frac{\alpha_{\text{high}}}{n_{\text{high}}} \left(\frac{n_{\text{high}}}{\alpha_{\text{high}}} - \alpha_{\text{high}} + t \right) \right]^{-n_{\text{high}}} & \text{if } t > \alpha_{\text{high}}, \end{cases}$$

→ **Sums of the above**

→ **Convolutions** (resolution \otimes Breit-Wigner, ...)

JINST 10 (2015) no.04, P04015



Discrete Profiling

Idea: treat the **type of function** and **number of parameters** as discrete NPs, profiled in data

→ **Let data choose the best shape**

→ Similar principle as other NPs, except for discrete nature

→ Need a **penalty on N_{pars}** to avoid always choosing functions with high N_{pars}

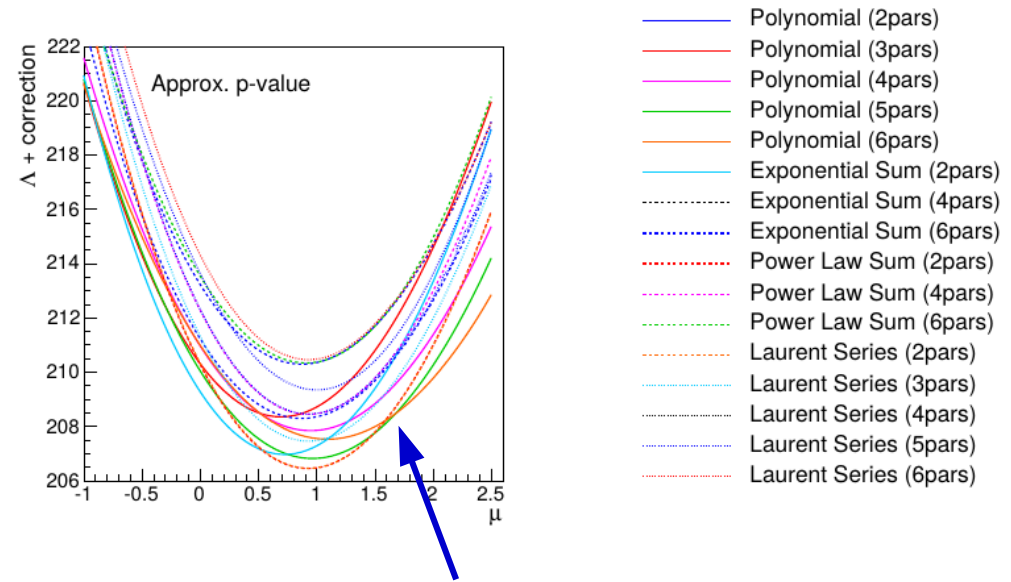
→ Used in the **CMS $H \rightarrow \gamma\gamma$** analysis, works well in this context.

Caveats:

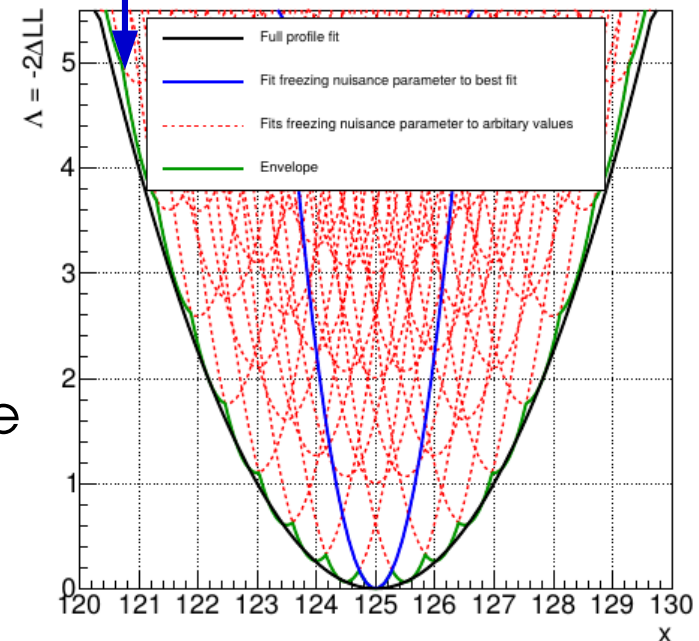
→ for N categories and M functional forms, **M^N possibilities to check** in principle – difficult in practice

→ Need a **well-chosen pool of sensible functions** for the method to work

→ **Large MC samples** for selection and checks



Take lower envelope of all functions when profiling



Gaussian Processes: 1-slide Summary

the data are drawn from one

HUGE^{*}
Gaussian

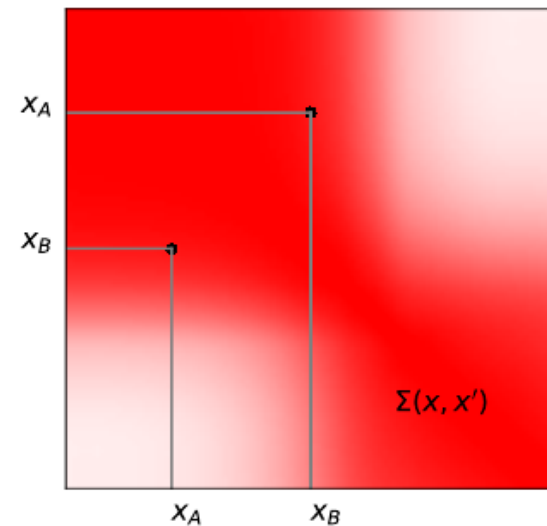
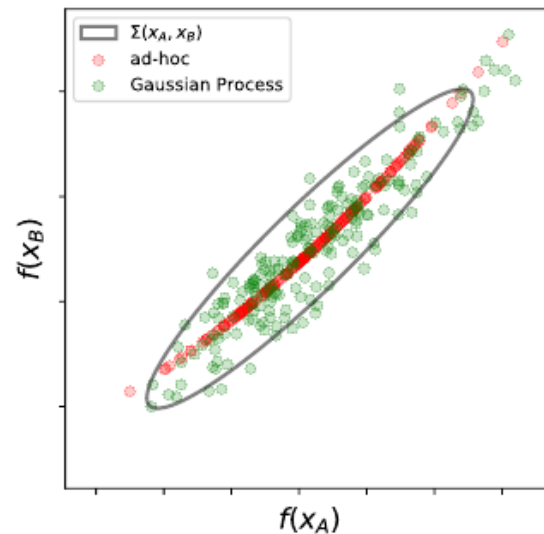
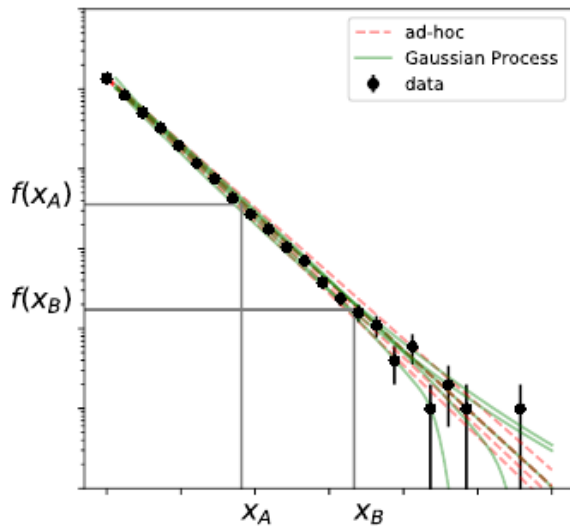
^{*} the dimension is the number of data points.

Image Credits:
K. Cranmer

Gaussian Processes: Longer 1-slide Summary

- Describe background distribution through the **correlations between values at different points**.
- More flexible than a functional form
- **Correlation function (Kernel)** can be
 - Defined using a length scale, to ignore narrow peaks
 - Obtained from first principles (e.g. from known JES/PDF effects)

$$K(x_1, x_2) = \exp\left[-\frac{(x_1 - x_2)^2}{2L^2}\right]$$



arXiv:1709.05681

- ⊕ More flexible than functional form, degrees of freedom less ad-hoc
- ⊖ Still need large MC samples to check for signal bias
- ⊖ Mainly for Gaussian processes, not well-adapted to Poisson regime

Statistical Modeling: II. Systematics

Systematics NPs

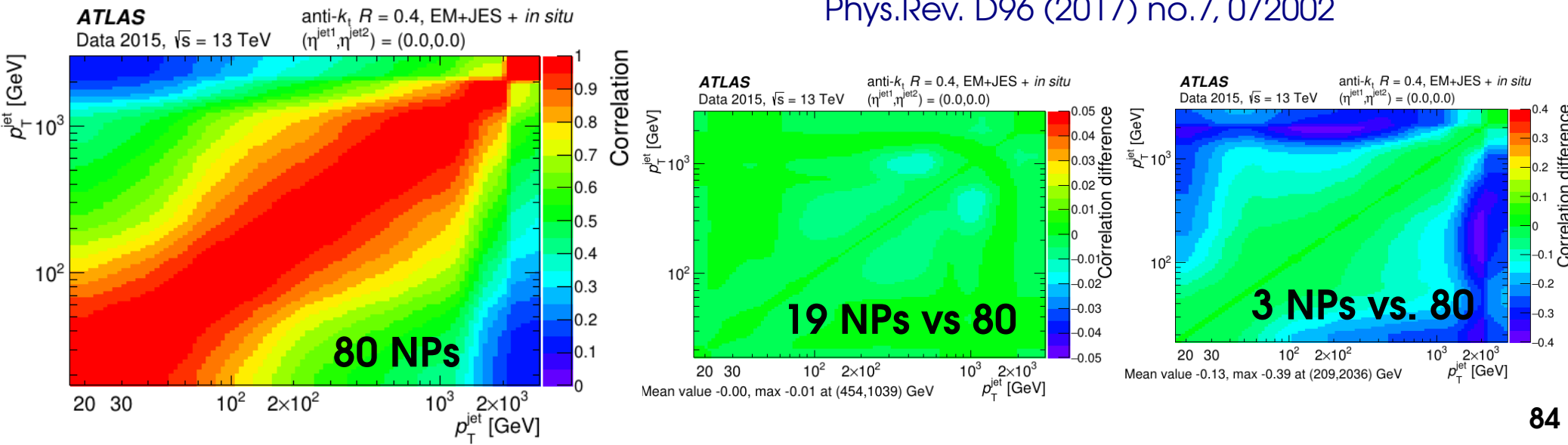
Each systematics NP represent **an independent source of uncertainty**
 ⇒ Usually constrained by a single 1-D PDF (Gaussian, etc.)

Sometimes multiple parameters **conjointly constrained** by an n-dim. PDF.
 → multiple measurements constraining multiple NPs

Assume **n-dim Gaussian** PDF: then can **diagonalize the covariance matrix C**
 and re-express the uncertainties in basis of eigenvector NPs ⇒ **n 1-dim PDFs**

Can also diagonalize to **prune** irrelevant uncertainties: keep NPs with large eigenvalues, combine in quadrature the others

Phys.Rev. D96 (2017) no.7, 072002



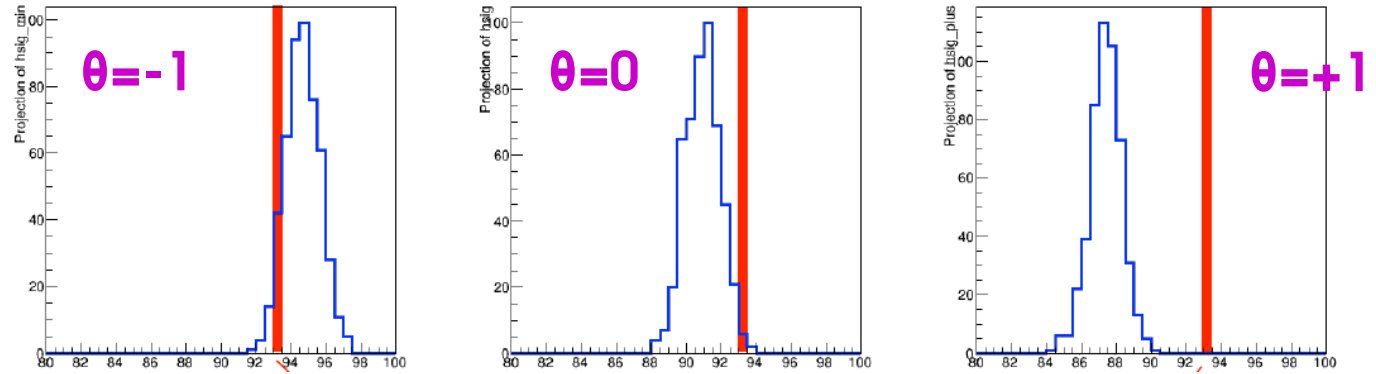
Systematics : Impact on Model

The effect of **each NP** θ_i should be propagated to all the relevant **model parameters** X_j .

→ **Propagation through MC:**

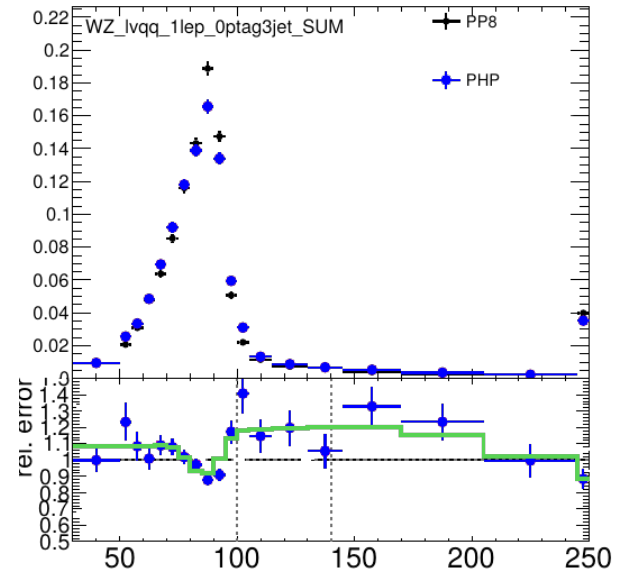
1. Apply $\pm 1\sigma$ systematic variations in MC,
 - ⇒ obtain shifted values $X_j^{\pm} = X_j^0 (1 \pm \Delta_{ij})$.
 - Possibly smooth out MC stats effects

2. Implement systematic in model, e.g. replace or morph shapes:



→ can affect event yields, shapes, etc.

Assumes **Gaussian uncertainties** and **linear impact** on model parameters



Constrained by unit Gaussian

$$X_j \rightarrow X_j^0 (1 + \Delta_{ij} \theta_i)$$

Systematics : Constraints

Ideally, constraint = **likelihood of auxiliary measurement**

⇒ e.g. Poisson for constraint from counting in a low-stat CR.

Sometimes no clear auxiliary measurement

⇒ Semi-arbitrary “pseudo-measurement” motivated by Central Limit Theorem:

- **Gaussian** for additive corrections
- **Log-normal** for multiplicative corrections

Gaussian:

- represent impact as $X_j \rightarrow X_j^0 (1 + \Delta_{ij} \theta_i)$
→ or similar morphing for distributions

Constrained by unit Gaussian

Can include asymmetric variations Δ^+, Δ^- :
$$X_j \rightarrow X_j^0 \left(1 + \begin{cases} \Delta_{ij}^+ \theta_i & \theta_i > 0 \\ \Delta_{ij}^- \theta_i & \theta_i < 0 \end{cases} \right)$$

However discontinuity in derivative at 0, so use smooth interpolation instead, e.g. implementation in `RooStats::HistFactory::FlexibleInterpVar`.

Systematics : Log-normal Constraint

Log-normal: $x \sim \text{log-normal}$ if $\log(x)$ is normal

→ **always** > 0 , useful to avoid numerical issues

→ **PDF:**

$$P(s; X_0, \kappa) = \frac{1}{x \kappa \sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{\log(x) - X_0}{\kappa}\right)^2\right)$$

However usually simpler to implement as :

$$X_j \rightarrow X_j^0 \exp(\kappa_{ij} \theta_i)$$

where θ_i is constrained by a unit Gaussian as usual

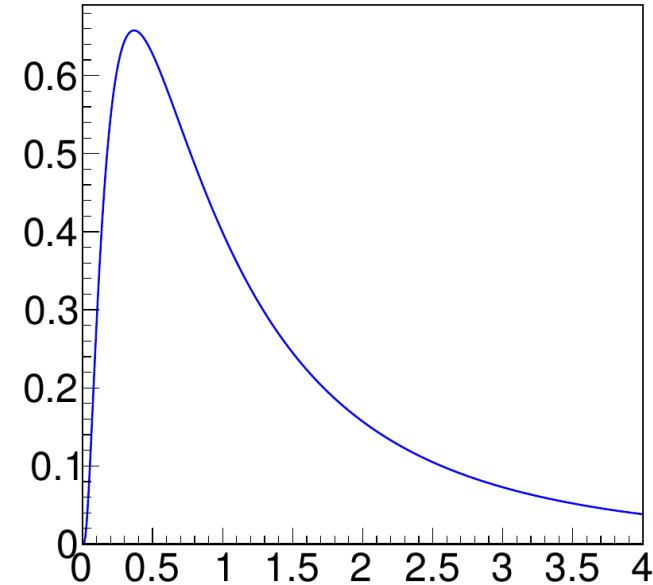
→ Correct form for a multiplicative uncertainty:

$$\log \sqrt[n]{(X_0 k_1)(X_0 k_2) \dots (X_0 k_n)} = \frac{1}{n} \sum_{i=1}^n \log(X_0 k_i) \stackrel{n \rightarrow \infty}{\sim} G(\log X_0, \frac{\text{RMS}(\log(k))}{\sqrt{n}} = \kappa)$$

Similarly to Gaussian → represent $\mathbf{X} = \mathbf{X}_0 \mathbf{e}^{\kappa\theta} \sim G(\log X_0, \kappa)$ if $\theta \sim G(0,1)$

Which κ to use ? $\kappa = \text{RMS}(X)$ only at first order. For larger uncertainties,

e.g. **Match $\pm 1\sigma$ variations:** $X_j(\theta=\pm 1) = X_j^\pm \Rightarrow \kappa_\pm = \pm \log(X_j^\pm / X_j^0)$



Implemented in RooStats::HistFactory::FlexibleInterpVar.

Systematics : Theory Constraints

Missing high-order terms in perturbative calculations: evaluate from scale variations – but no underlying random process. Possible constraint shapes:

- **Gaussians** (ATLAS/CMS Higgs analyses, see [Yellow Report 4, I.4.1.d](#))

- Usually several independent “sources” of uncertainty (QCD/EW/resummation)
- ⇒ overall uncertainty may be rather Gaussian
- Numerically well-behaved
- Uncertainties add in quadrature as usual

- **Flat constraints** : “100% confidence” intervals

- no preference for any value in the range
- Need regularization to avoid numerical issues
- uncertainties add linearly

→ For Higgs cross-sections, rather similar results for both cases

Constraints : Two-point systematics

Sometimes differences between 2 discrete cases → e.g. Pythia vs. Herwig
Solutions:

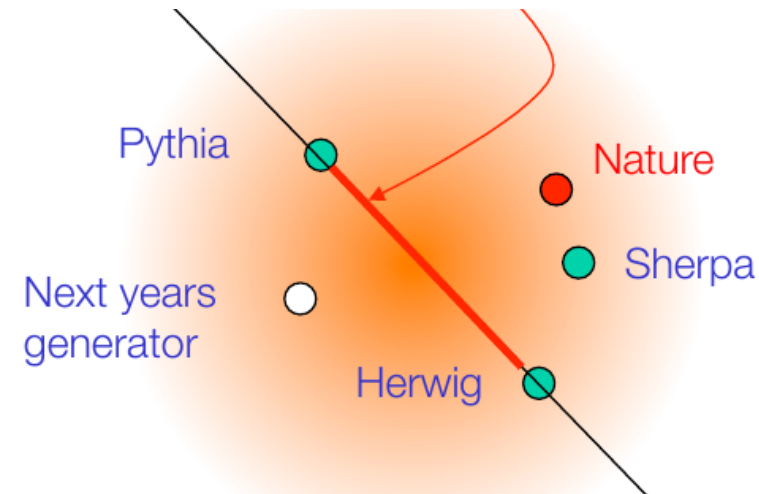
- Results for one case only
- Full results for both cases
- Single result with an uncertainty that covers the difference
→ **Two-point uncertainty**

Usually implemented as **1D linear interpolations** between the two cases
→ However cannot guarantee this covers the space of possible configurations

⇒ ***This is not even a pseudo-measurement...***

Ideally, need to define proper uncertainties within a single model, which would cover the other cases
→ e.g. showering uncertainties within Pythia, covering Herwig

→ **Usually a difficult task**



Profiling Issues

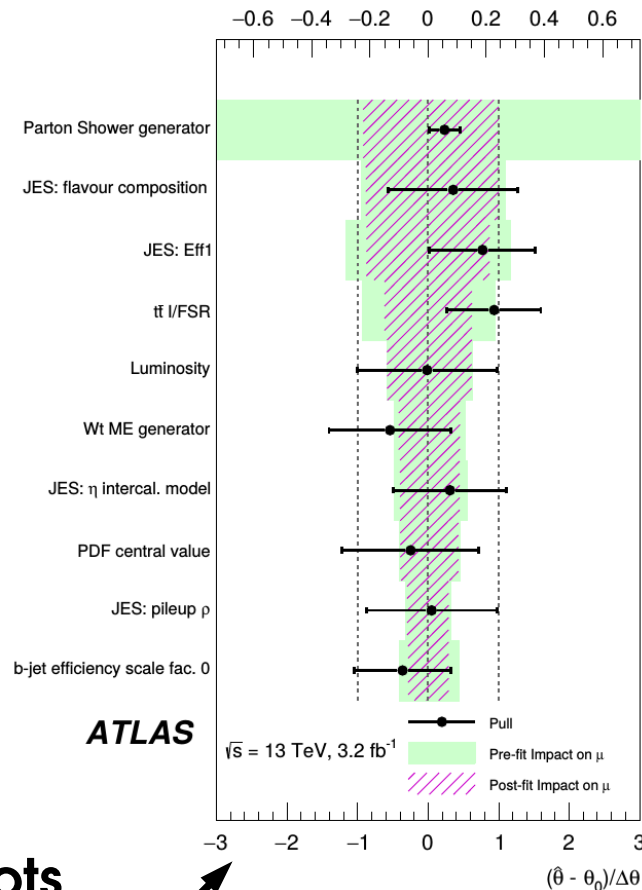
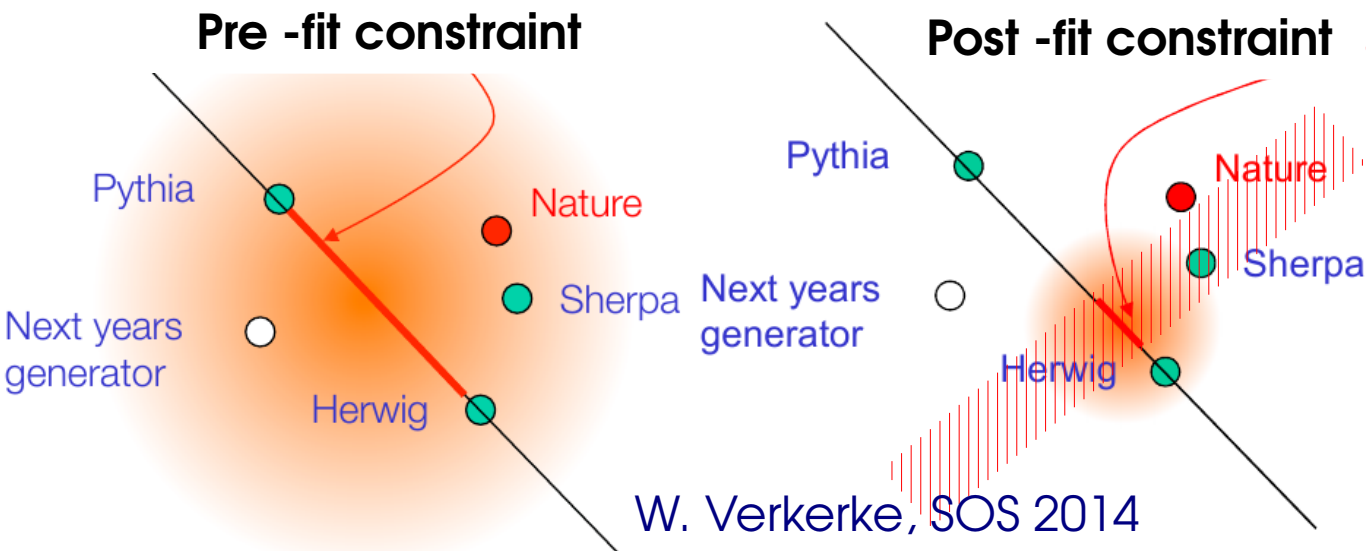
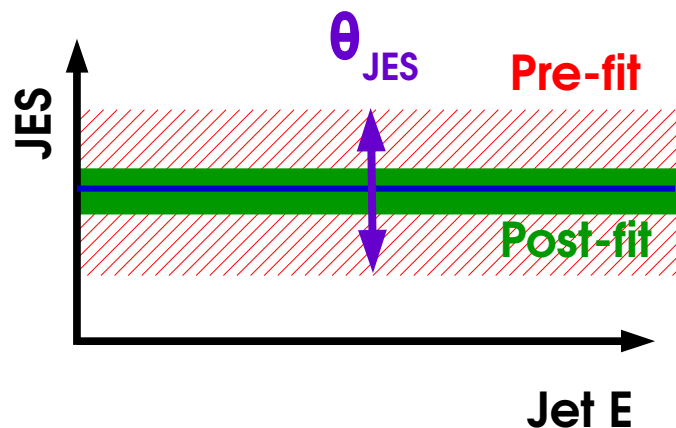
Too simple modeling can have unintended effects

→ 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



NP central values and uncertainties in pull/impact plots provide important “debugging” information for profiling

Outline

Profiling

Look-Elsewhere Effect

Bayesian methods

Statistical modeling in practice

Building binned likelihoods

Choosing PDFs in unbinned likelihoods

Implementing systematics

BLUE

BLUE

BLUE

Commonly-used ansatz for combination of measurements:

1. **Build a χ^2 :** (same as $-2\log L$ for Gaussian L)

$$\chi^2(\mathbf{X}) = \sum_i \left(\mathbf{X}_i^{\text{obs}} - \mathbf{X} \right) \mathbf{C}_{ij}^{-1} \left(\mathbf{X}_j^{\text{obs}} - \mathbf{X} \right)$$

\mathbf{C}_{ij} : covariance matrix of measurements:

$$\mathbf{C} = \begin{bmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 & \cdots \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

ρ : correlation coefficients

2. **Estimate combined X from minimum of $\chi^2(\mathbf{X})$**

- In the Gaussian case, equivalent to ML solution \Rightarrow inherits good properties:
 - **Unbiased** : $\langle \hat{\mathbf{X}} \rangle = \mathbf{X}^*$
 - **Optimal**: minimizes the combined uncertainty
- Solution is a linear combination of the inputs:

$$\boldsymbol{\lambda} = \frac{\mathbf{C}^{-1} \mathbf{J}}{\mathbf{J}^T \mathbf{C}^{-1} \mathbf{J}}, \quad \mathbf{J} = \begin{pmatrix} 1 \\ 1 \\ \vdots \end{pmatrix}$$

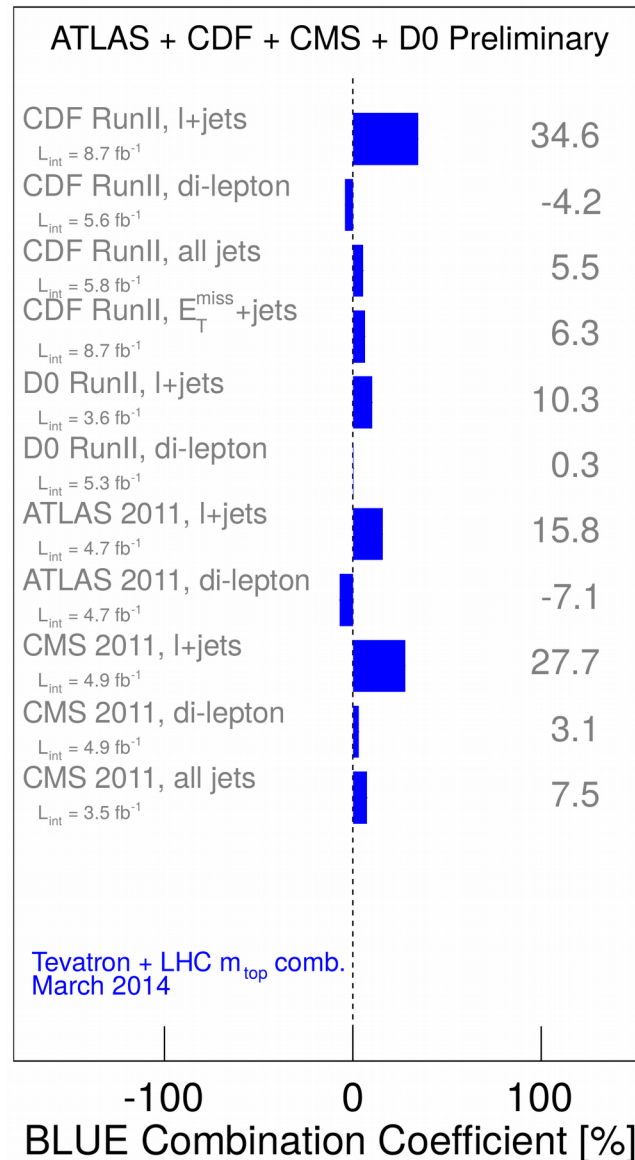
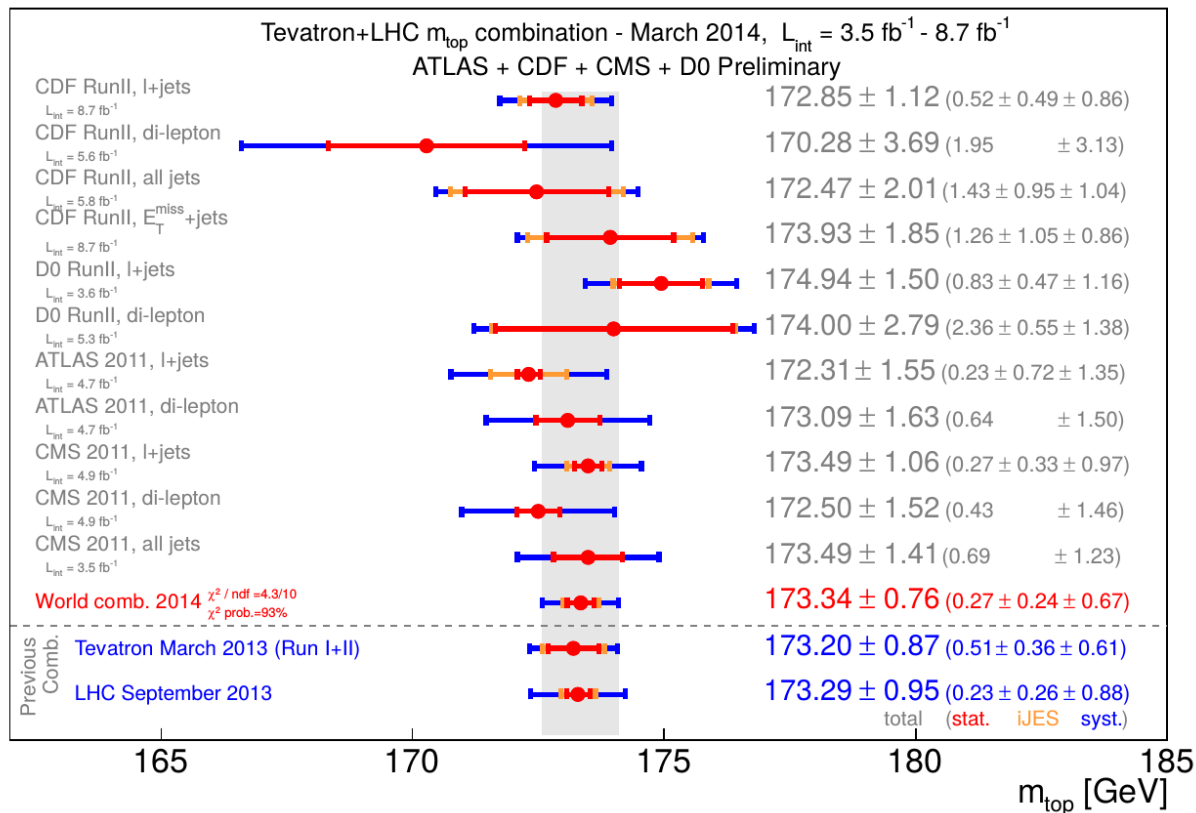
λ_i = combination weight of measurement i

$$\hat{\mathbf{X}} = \sum_i \lambda_i \mathbf{X}^{\text{obs}, i}$$

\Rightarrow “**Best Linear Unbiased Estimator**” (**BLUE**)

BLUE Example

Example: World m_{top} combination



Limitation: relies on **Gaussian assumptions** (satisfied in this case!)

Negative weights possible! (for large correlations, see [Eur. Phys. J. C 74 \(2014\), 2717](#))

BLUE and PLR

PLR Computation: 2 measurements
+ 1 auxiliary measurement

$$\begin{aligned} X_1 &= X + \Delta_1 \theta \sim G(X^*, \sigma_1) \\ X_2 &= X + \Delta_2 \theta \sim G(X^*, \sigma_2) \\ \theta &\sim G(0, 1) \end{aligned}$$

Single measurement: $\lambda(X, \theta) = \frac{1}{\sigma_1^2} (X + \Delta_1 \theta - X_1^{\text{obs}})^2 + (\theta - \theta^{\text{obs}})^2$

MLEs:
$$\begin{cases} \hat{\theta} = \theta^{\text{obs}} \\ \hat{X} = X_1^{\text{obs}} - \Delta_1 \theta^{\text{obs}} \end{cases}$$

PLR:
$$\lambda(X) = \frac{(X - \hat{X})^2}{\sigma_{1, \text{tot}}^2} \quad \sigma_{1, \text{tot}}^2 = \sigma_1^2 + \Delta_1^2$$

Combination:
$$\lambda(X, \theta) = \frac{1}{\sigma_1^2} (X + \Delta_1 \theta - X_1^{\text{obs}})^2 + \frac{1}{\sigma_2^2} (X + \Delta_2 \theta - X_2^{\text{obs}})^2 + (\theta - \theta^{\text{obs}})^2$$

MLE:
$$\hat{X} = \lambda_1 X_1^{\text{obs}} + \lambda_2 X_2^{\text{obs}} + \lambda_\theta \theta^{\text{obs}} \quad \lambda_{1(2)} = \frac{\sigma_{2(1), \text{tot}}^2 - \Delta_1 \Delta_2}{\sigma_{1, \text{tot}}^2 + \sigma_{2, \text{tot}}^2 - 2\Delta_1 \Delta_2}$$

PLR:
$$\lambda(X) = \frac{(X - \hat{X})^2}{\sigma_{X, \text{tot}}^2} \quad \sigma_{X, \text{tot}}^2 = \frac{\sigma_{1, \text{tot}}^2 \sigma_{2, \text{tot}}^2 - \Delta_1^2 \Delta_2^2}{\sigma_{1, \text{tot}}^2 + \sigma_{2, \text{tot}}^2 - 2\Delta_1 \Delta_2}$$

BLUE computation: measurements X_1 and X_2 with uncorrelated statistical uncertainties σ_1 and σ_2 , correlated systematics Δ_1 and Δ_2 .

Single measurement: stat uncertainty σ_1 , systematic Δ_1

- Uncorrelated uncertainties
- Assume everything is Gaussian

⇒ Uncertainties add in quadrature:

$$\sigma_{1, \text{tot}}^2 = \sigma_1^2 + \Delta_1^2$$

Combination:

$$C = \begin{bmatrix} \sigma_{1, \text{tot}}^2 & \rho \sigma_{1, \text{tot}} \sigma_{2, \text{tot}} \\ \rho \sigma_{1, \text{tot}} \sigma_{2, \text{tot}} & \sigma_{2, \text{tot}}^2 \end{bmatrix} \quad \rho = \frac{\Delta_1 \Delta_2}{\sigma_{1, \text{tot}} \sigma_{2, \text{tot}}}$$

BLUE weights

$$\hat{X} = \lambda_1 X_1^{\text{obs}} + \lambda_2 X_2^{\text{obs}}$$

$$\lambda_{1(2)} = \frac{\sigma_{2(1), \text{tot}}^2 - \rho \sigma_{1, \text{tot}} \sigma_{2, \text{tot}}}{\sigma_{1, \text{tot}}^2 + \sigma_{2, \text{tot}}^2 - 2\rho \sigma_{1, \text{tot}} \sigma_{2, \text{tot}}}$$

Propagate uncertainties from C:

$$\sigma_{X, \text{tot}}^2 = \frac{\sigma_{1, \text{tot}}^2 \sigma_{2, \text{tot}}^2 (1 - \rho^2)}{\sigma_{1, \text{tot}}^2 + \sigma_{2, \text{tot}}^2 - 2\rho \sigma_{1, \text{tot}} \sigma_{2, \text{tot}}}$$

Negative BLUE Weights

Occasionally, negative BLUE weights:

Can happen if $\rho \sim 1$:

$$\lambda_2 = \frac{\sigma_{1, \text{tot}} (\sigma_{1, \text{tot}} - \rho \sigma_{2, \text{tot}})}{\sigma_{1, \text{tot}}^2 + \sigma_{2, \text{tot}}^2 - 2\rho \sigma_{1, \text{tot}} \sigma_{2, \text{tot}}} < 0 \text{ for } \rho > \frac{\sigma_{1, \text{tot}}}{\sigma_{2, \text{tot}}}$$

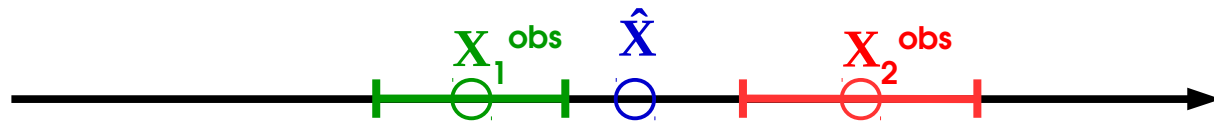
Not intuitive! (Can also have $\lambda_2 = 0$ for $\sigma_{1, \text{tot}} = \rho \sigma_{2, \text{tot}} \dots$)

Can be explained in the PLR picture:

$$X_1 = X + \Delta \theta$$

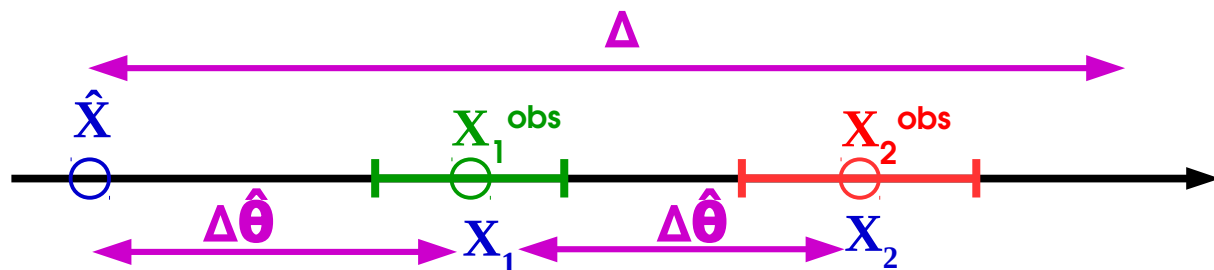
$$X_2 = X + 2\Delta \theta$$

Without correlated systematics ($\Delta = 0$):



$$\lambda_{1(2)} = \frac{\sigma_{2(1)}^2}{\sigma_1^2 + \sigma_2^2} > 0$$

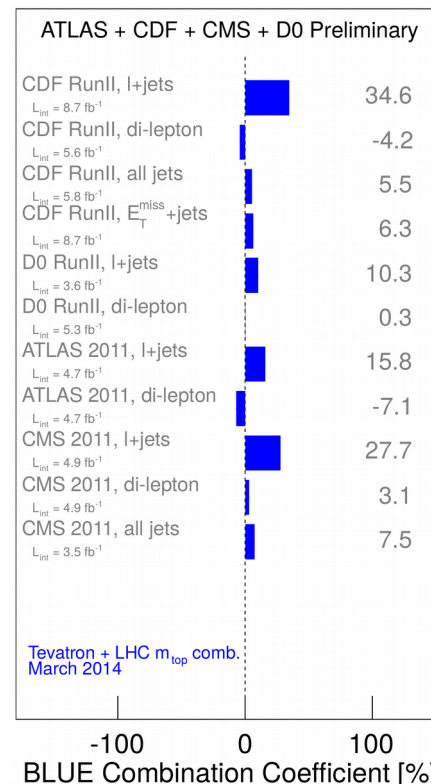
With large correlated systematics ($\Delta \gg \sigma_{1,2}$)



$\hat{\theta}$ value makes X_1 and X_2 match observations, small pull on θ if Δ is large

$$\lambda_1 < 0$$

$\rho \sim 1 \Rightarrow \theta$ measurement is important \Rightarrow possibly very different MLE than $X_1 \oplus X_2 \dots_{97}$



Uncertainty Decomposition

Often useful to break down uncertainties into components (stat + syst, etc.)

PLR approach: perform measurement twice

1. With all uncertainties included
→ **nominal uncertainty** σ_{total}
2. Removing some uncertainties
(e.g. all syst uncertainties) → $\sigma_{\text{no-syst}}$

⇒ Subtract in quadrature:

$$\sigma_{\text{syst}} = \sqrt{\sigma_{\text{total}}^2 - \sigma_{\text{no-syst}}^2}$$

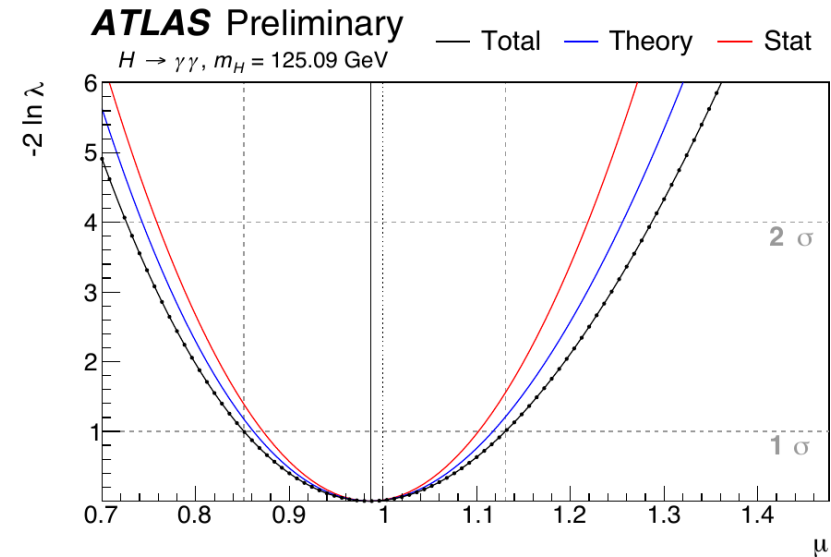
BLUE-based approach:

1. Propagate each source of uncertainty (stat & syst) to the observables
2. Propagate through to the measurement using the BLUE weights

$$\hat{X} = \sum_i \lambda_i X^{\text{obs},i}$$

The two methods are not completely equivalent (recently discovered!)

→ In the BLUE case, weights still computed including systematics effects

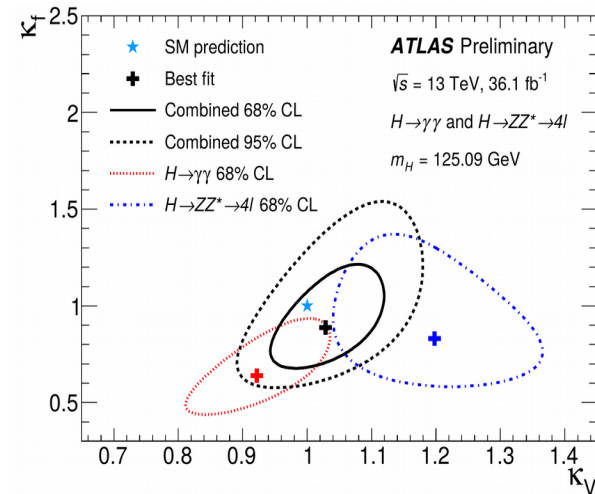
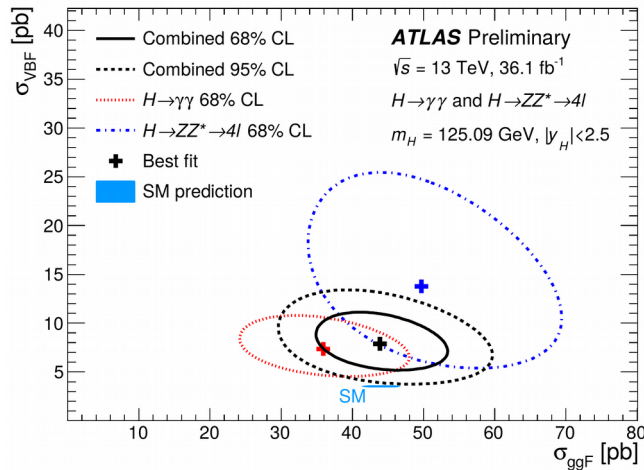


Presentation of Results

Presentation of Results

Measurements often recast to constrain a particular theory model.

→ Ideally, by **reparameterizing the likelihood** and repeating the measurement



⇒ **Done by experiments for selected benchmark models**

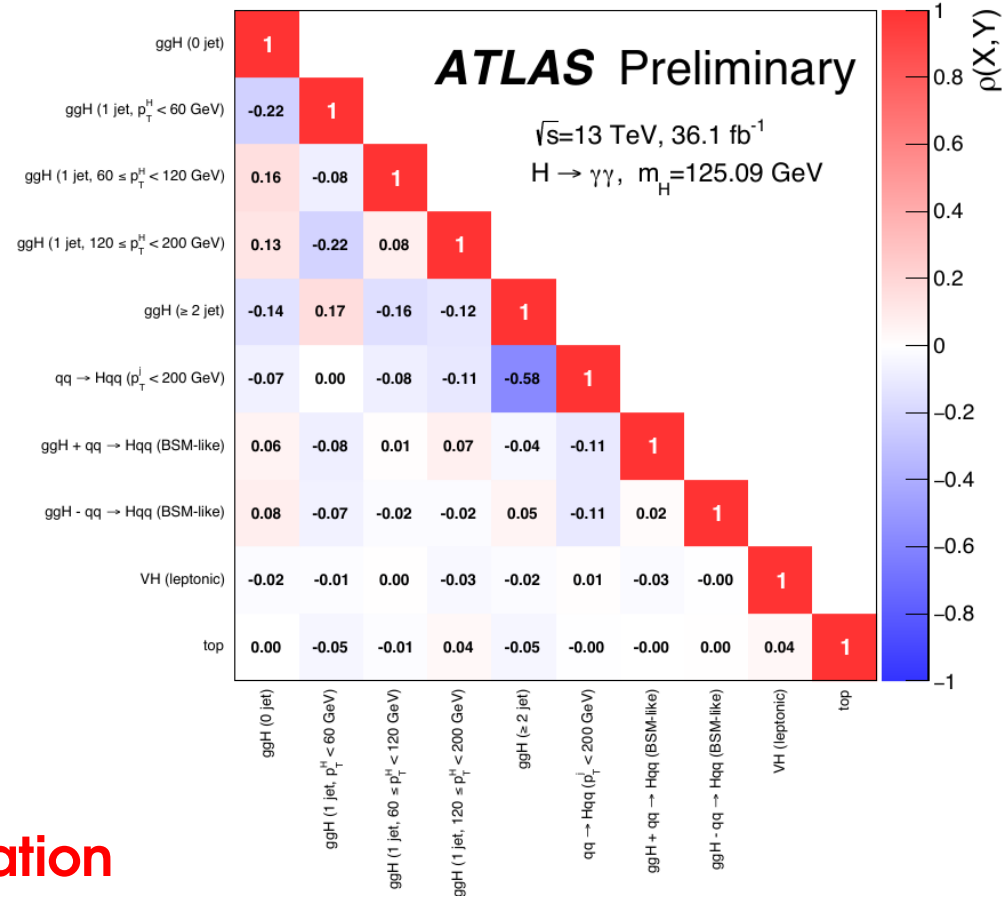
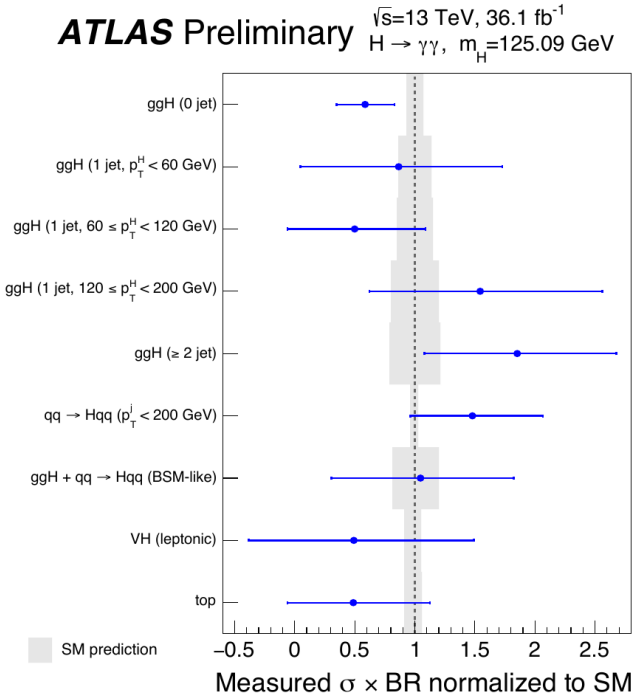
→ However, **often too complex to implement widely:**

- Full likelihood typically not published
- theorists typically do not want to deal with 4000 NPs...

→ **Other approaches:** e.g. reimplementing the analysis in a public fast-simulation framework (e.g. SUSY searches). However clear accuracy limitations

Presentation of Results

→ **Current solution**: publish covariance matrices in **HEPData**, together with the individual measurements



- **Only valid in the Gaussian approximation**
- To go further, need some form of **simplified likelihoods**
- **Profile likelihood** – function of POI only (NPs profiled out)
- **Additional terms** for non-Gaussian effects
- Significantly more complex (many dimensions!)
- Will be needed eventually as measurements become syst-dominated

Conclusion

- Significant evolution in the statistical methods used in HEP
- Variety of methods, adapted to various situations and target results
- Allow to
 - model the statistical process with high precision in difficult situations (large systematics, small signals)
 - make optimal use of available information
- Implemented in standard RooFit/RooStat toolkits within the ROOT framework, as well as other tools (BAT)
- Improvement and uniformization efforts are still ongoing
- Still many open questions and areas that could use improvement
→ e.g. how to present results with all available information to the “outside”

Extra Slides

Uncertainty decomposition

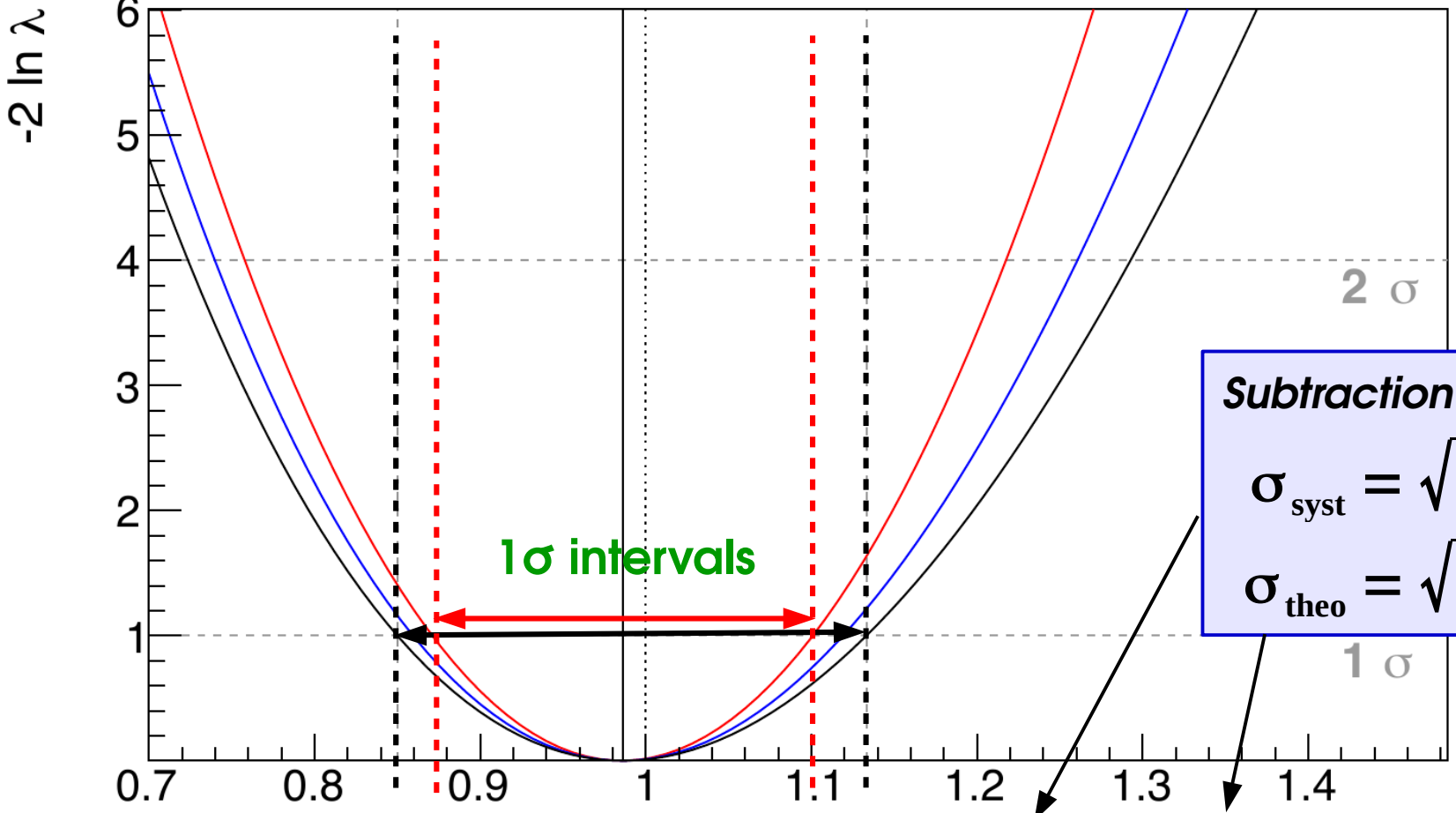
All systematics NPs fixed to 0 : statistical uncertainty only

exp. syst. NPs fixed to 0 : stat+theory uncertainty

ATLAS

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

— Total — Theory — Stat



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

Gaussian Profiling

Gaussian measurement with 1 POI μ and 1 NP θ :

$$L(\mu, \theta; \hat{\mu}, \hat{\theta}) = \exp \left[-\frac{1}{2} \begin{pmatrix} \mu - \hat{\mu} \\ \theta - \hat{\theta} \end{pmatrix}^T C^{-1} \begin{pmatrix} \mu - \hat{\mu} \\ \theta - \hat{\theta} \end{pmatrix} \right] \quad C = \begin{bmatrix} \sigma_\mu^2 & \gamma \sigma_\mu \sigma_\theta \\ \gamma \sigma_\mu \sigma_\theta & \sigma_\theta^2 \end{bmatrix}$$

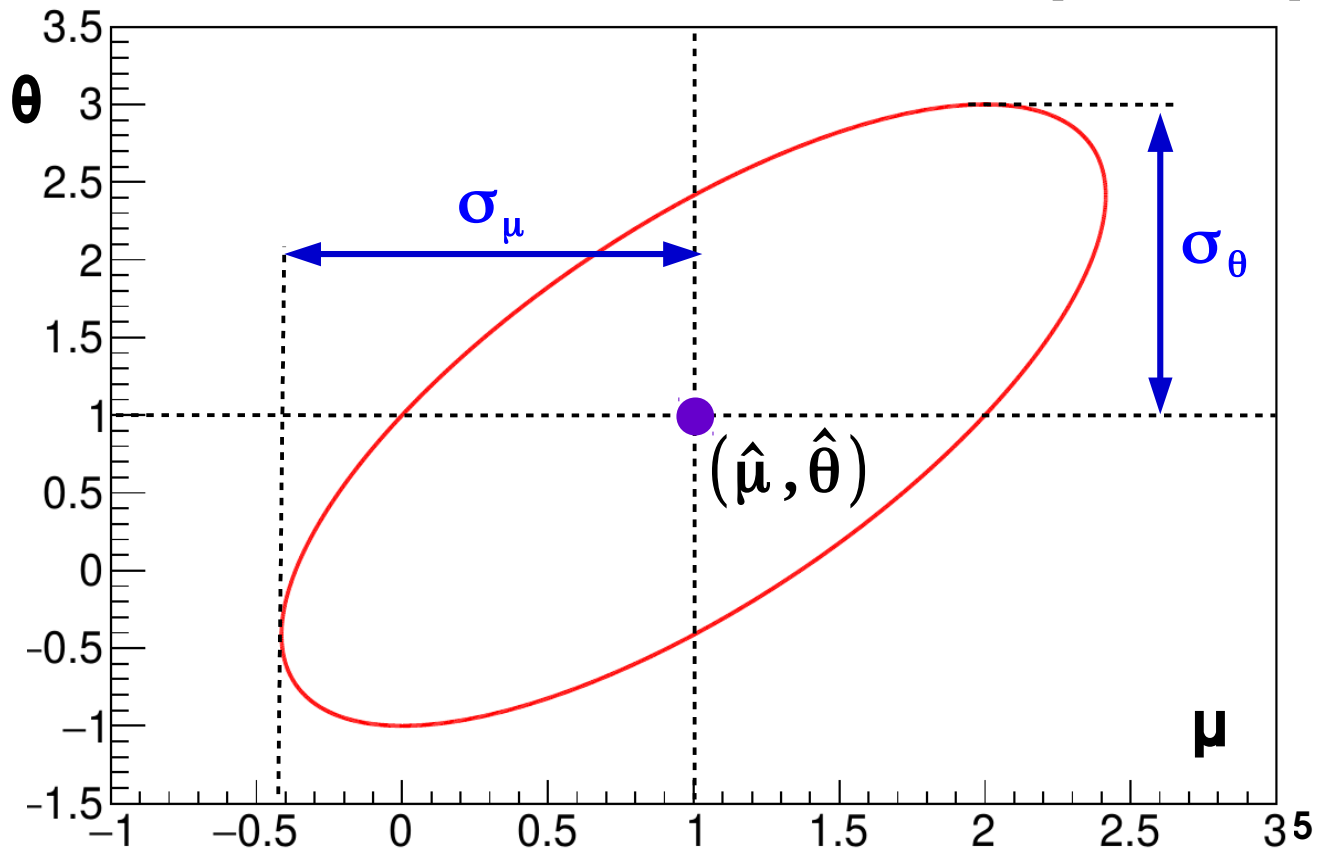
"data"

→ $\lambda(\mu, \theta)$ defines an **ellipse**:

$$\lambda(\mu, \theta; \hat{\mu}, \hat{\theta}) = F_{\mu\mu}(\mu - \hat{\mu})^2 + 2F_{\mu\theta}(\mu - \hat{\mu})(\theta - \hat{\theta}) + F_{\theta\theta}(\theta - \hat{\theta})^2 \quad F \equiv C^{-1} = \begin{bmatrix} F_{\mu\mu} & F_{\mu\theta} \\ F_{\mu\theta} & F_{\theta\theta} \end{bmatrix}$$

Uncertainty on μ :

- From C , with θ included: σ_μ



Gaussian Profiling

$$C = \begin{bmatrix} \sigma_\mu^2 & \gamma \sigma_\mu \sigma_\theta \\ \gamma \sigma_\mu \sigma_\theta & \sigma_\theta^2 \end{bmatrix}$$

$$F = \begin{bmatrix} F_{\mu\mu} & F_{\mu\theta} \\ F_{\mu\theta} & F_{\theta\theta} \end{bmatrix}$$

$$\lambda(\mu, \theta; \hat{\mu}, \hat{\theta}) = F_{\mu\mu}(\mu - \hat{\mu})^2 + 2F_{\mu\theta}(\mu - \hat{\mu})(\theta - \hat{\theta}) + F_{\theta\theta}(\theta - \hat{\theta})^2$$

Profiled θ (minimize λ at fixed μ) :

$$\hat{\theta}(\mu) = \hat{\theta} - F_{\theta\theta}^{-1} F_{\theta\mu}(\mu - \hat{\mu})$$

Profile likelihood ratio:

$$\lambda(\mu, \hat{\theta}(\mu); \hat{\mu}, \hat{\theta}) = (F_{\mu\mu} - F_{\mu\theta} F_{\theta\theta}^{-1} F_{\theta\mu})(\mu - \hat{\mu})^2 = C_{\mu\mu}^{-1}(\mu - \hat{\mu})^2 = \left(\frac{\mu - \hat{\mu}}{\sigma_\mu}\right)^2$$

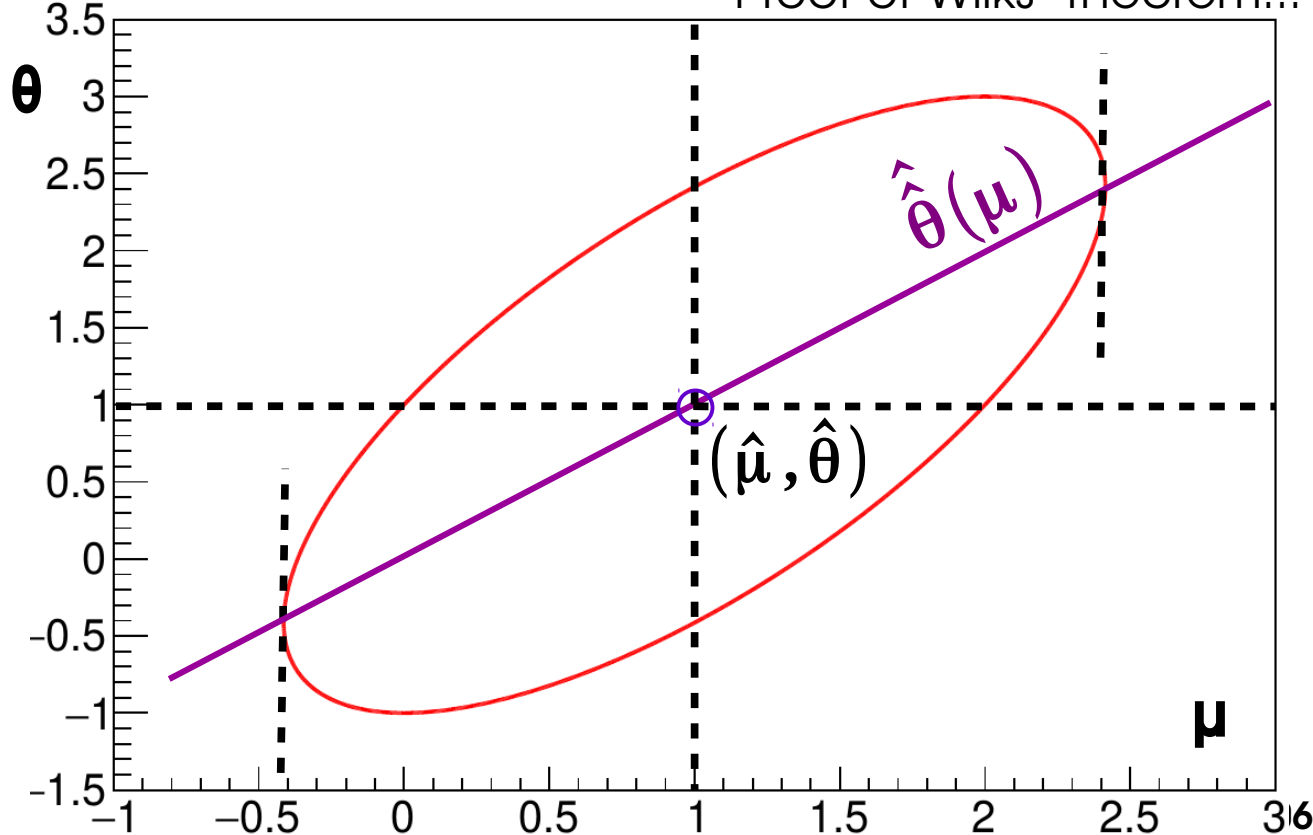
$F_{\mu\mu} \neq C_{\mu\mu}^{-1} !!$

Proof of Wilks' theorem...

Uncertainty on μ :

- From C: σ_μ
- From PLR: σ_μ

Profiled θ **crosses ellipse at vertical tangents** by definition (L is lower at other points on the tangent)



Gaussian Profiling

$$\lambda(\mu, \theta; \hat{\mu}, \hat{\theta}) = F_{\mu\mu}(\mu - \hat{\mu})^2 + 2F_{\mu\theta}(\mu - \hat{\mu})(\theta - \hat{\theta}) + F_{\theta\theta}(\theta - \hat{\theta})^2$$

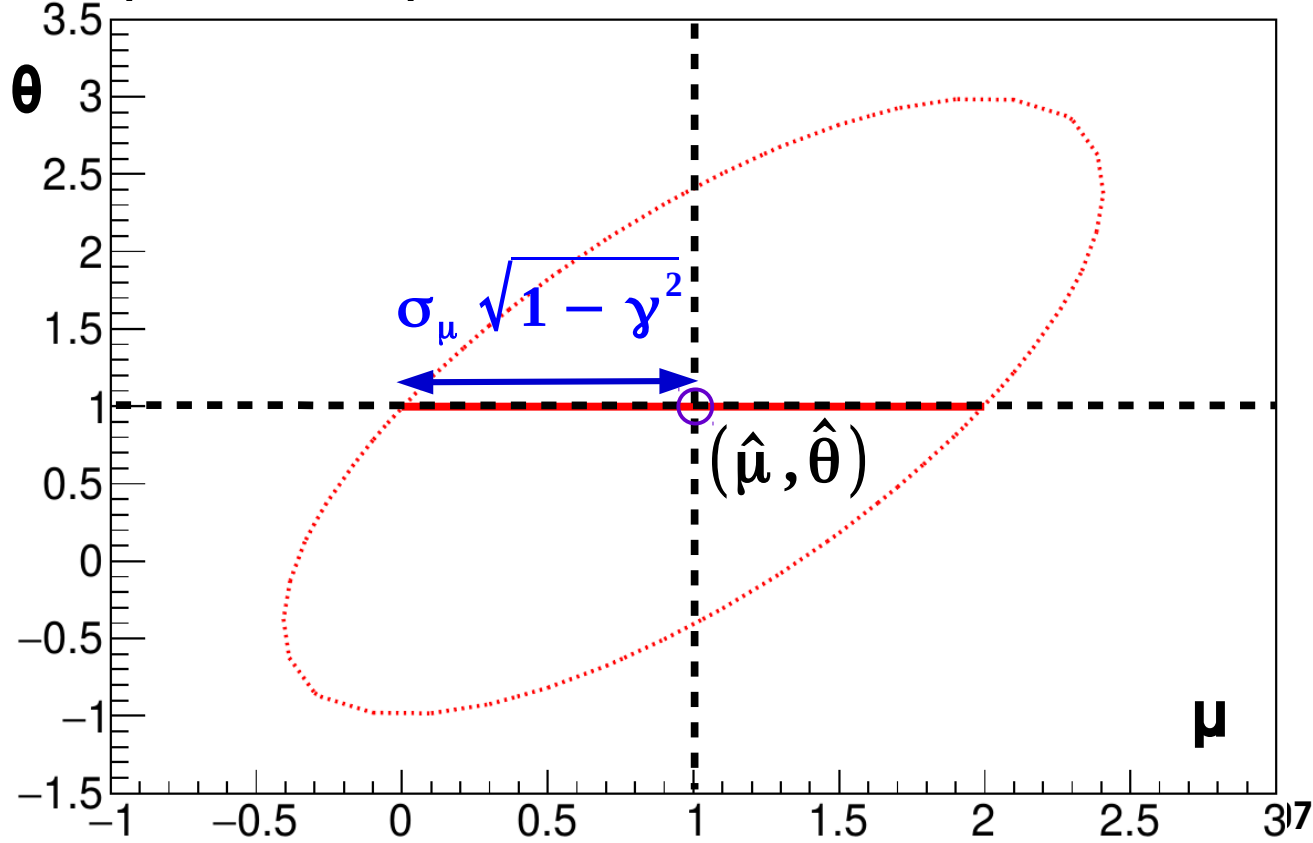
$$F \equiv C^{-1} = \frac{1}{1 - \gamma^2} \begin{bmatrix} \frac{1}{\sigma_\mu^2} & \frac{\gamma}{\sigma_\mu \sigma_\theta} \\ \frac{\gamma}{\sigma_\mu \sigma_\theta} & \frac{1}{\sigma_\theta^2} \end{bmatrix}$$

→ For fixed $\theta = \hat{\theta}$, $\lambda(\mu)$ defines an interval:

$$\lambda(\mu, \theta = \hat{\theta}; \hat{\mu}, \hat{\theta}) = F_{\mu\mu}(\mu - \hat{\mu})^2 = \left(\frac{\mu - \hat{\mu}}{\sigma_\mu \sqrt{1 - \gamma^2}} \right)^2$$

Uncertainty on μ :

- From C: σ_μ
- From PLR: σ_μ
- From $\lambda(\mu)$: $\sigma_\mu \sqrt{1 - \gamma^2}$



Gaussian Profiling

$$\lambda(\mu, \theta; \hat{\mu}, \hat{\theta}) = F_{\mu\mu}(\mu - \hat{\mu})^2 + 2F_{\mu\theta}(\mu - \hat{\mu})(\theta - \hat{\theta}) + F_{\theta\theta}(\theta - \hat{\theta})^2$$

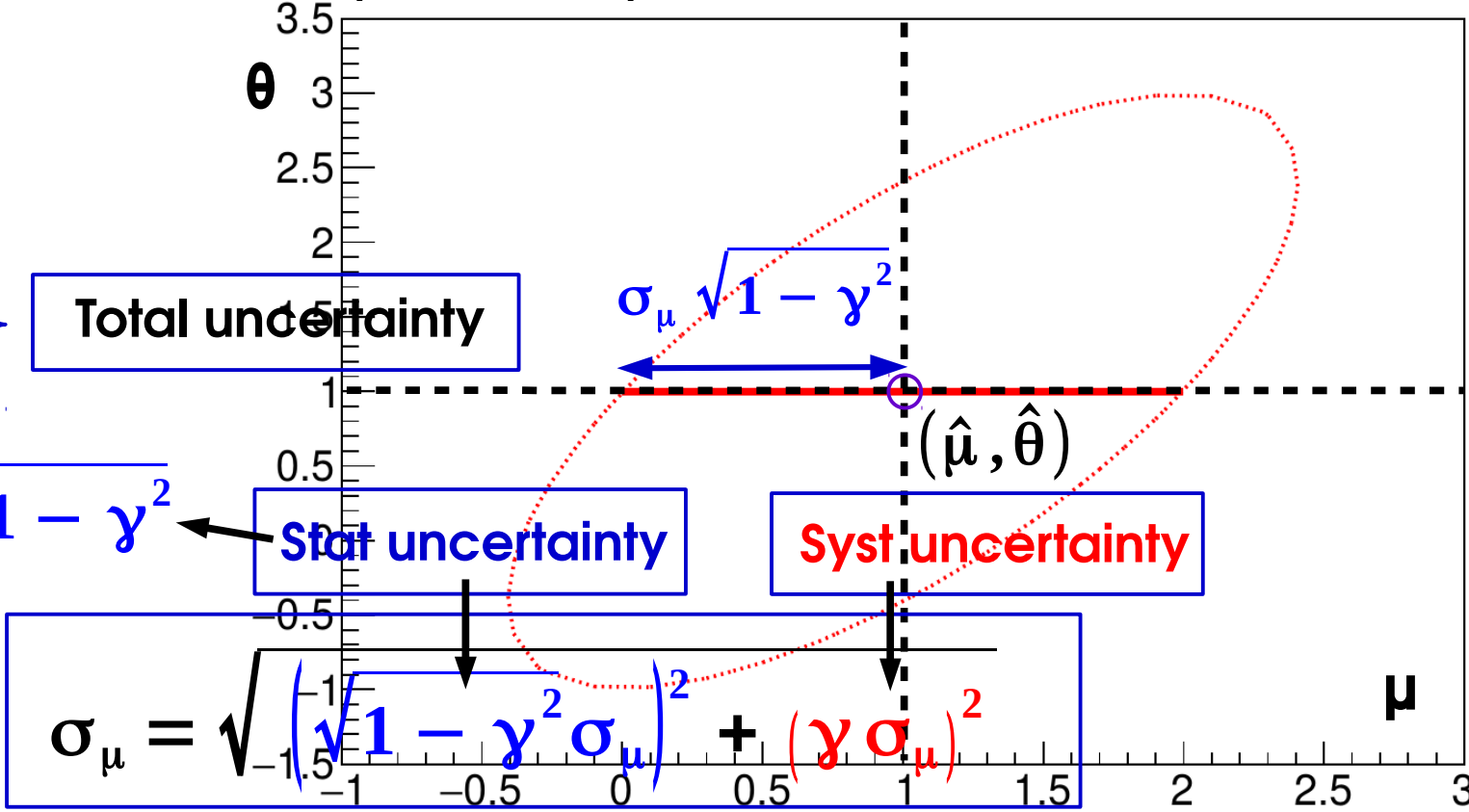
$$F \equiv C^{-1} = \frac{1}{1 - \gamma^2} \begin{bmatrix} \frac{1}{\sigma_\mu^2} & \frac{\gamma}{\sigma_\mu \sigma_\theta} \\ \frac{\gamma}{\sigma_\mu \sigma_\theta} & \frac{1}{\sigma_\theta^2} \end{bmatrix}$$

→ For fixed $\theta = \hat{\theta}$, $\lambda(\mu)$ defines an interval:

$$\lambda(\mu, \theta = \hat{\theta}; \hat{\mu}, \hat{\theta}) = F_{\mu\mu}(\mu - \hat{\mu})^2 = \left(\frac{\mu - \hat{\mu}}{\sigma_\mu \sqrt{1 - \gamma^2}} \right)^2$$

Uncertainty on μ :

- From C: σ_μ
- From PLR: σ_μ
- From $\lambda(\mu)$: $\sigma_\mu \sqrt{1 - \gamma^2}$



Comparison with LEP/TeVatron definitions

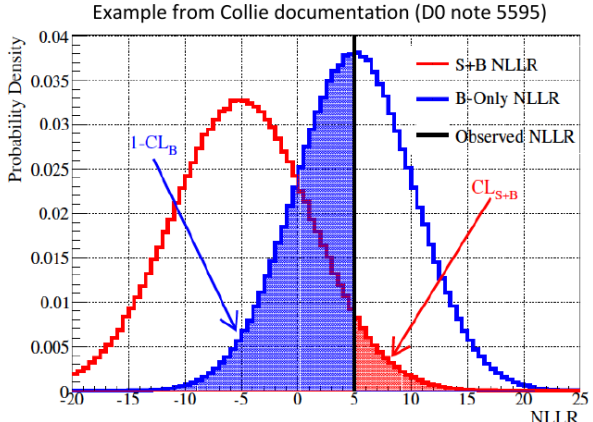
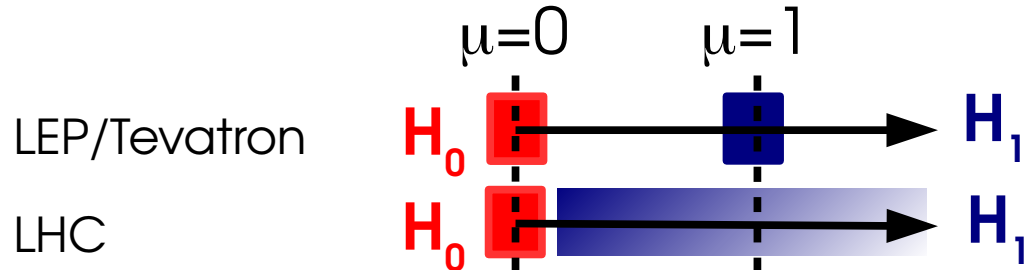
Likelihood ratios are not a new idea:

- **LEP**: Simple LR with NPs from MC
 - Compare $\mu=0$ and $\mu=1$
- **TeVatron**: PLR with profiled NPs

$$q_{LEP} = -2 \log \frac{L(\mu=0, \tilde{\theta})}{L(\mu=1, \tilde{\theta})}$$

$$q_{TeVatron} = -2 \log \frac{L(\mu=0, \hat{\theta}_0)}{L(\mu=1, \hat{\theta}_1)}$$

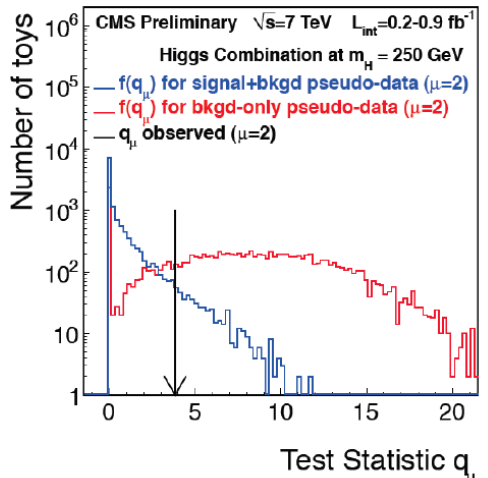
Both compare to $\mu=1$ instead of best-fit $\hat{\mu}$



→ Asymptotically:

- **LEP/TeVatron**: q linear in $\mu \Rightarrow \sim$ Gaussian
- **LHC**: q quadratic in $\mu \Rightarrow \sim \chi^2$

→ Still use TeVatron-style for discrete cases



Spin/Parity Measurements

Phys. Rev. D 92 (2015) 012004

