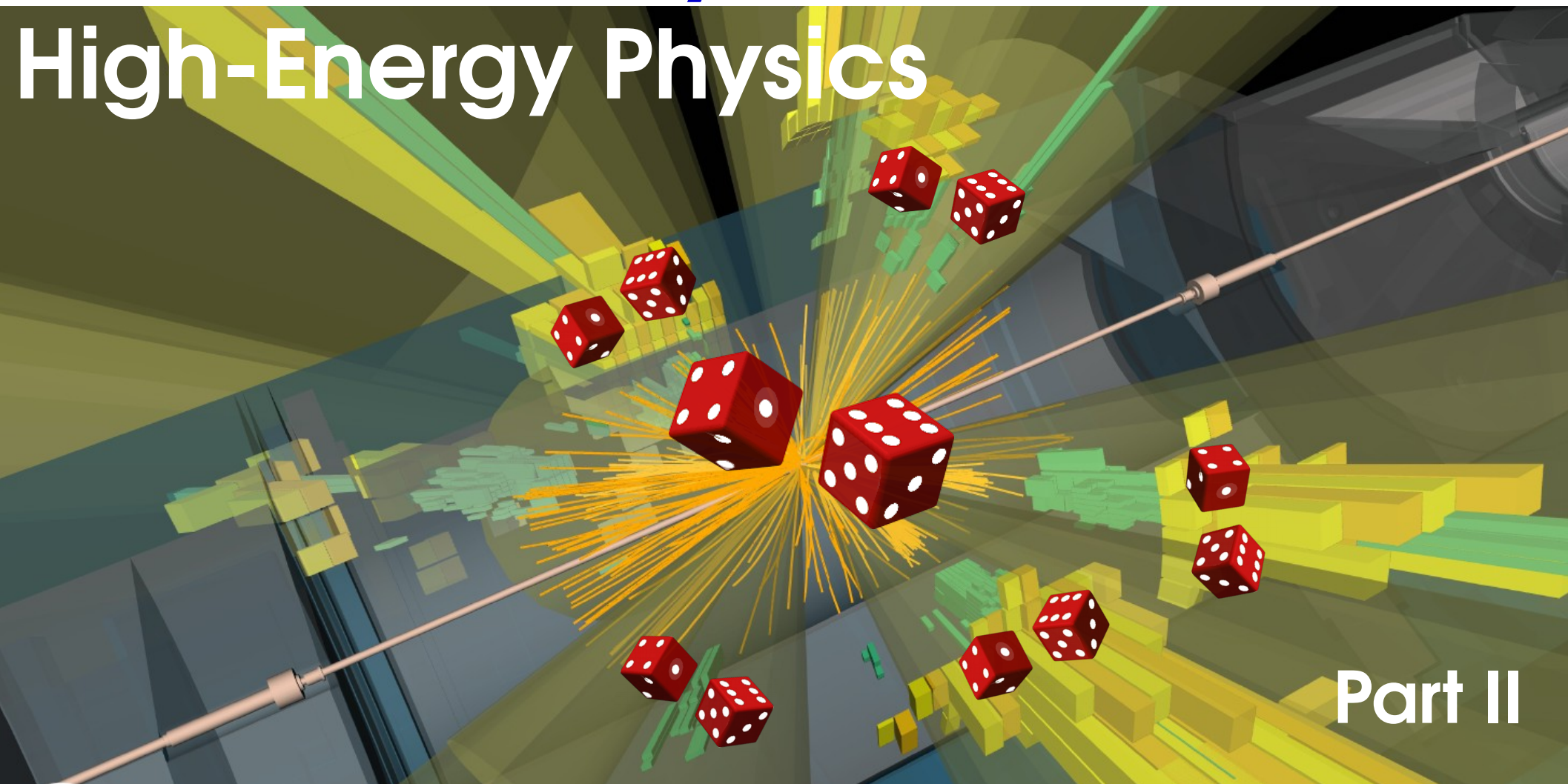


---

# Statistical analysis methods in High-Energy Physics



Part II

---

# Computing Statistical Results

## III. Discovery

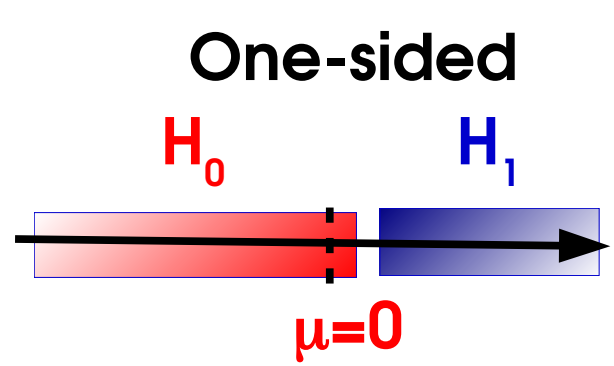
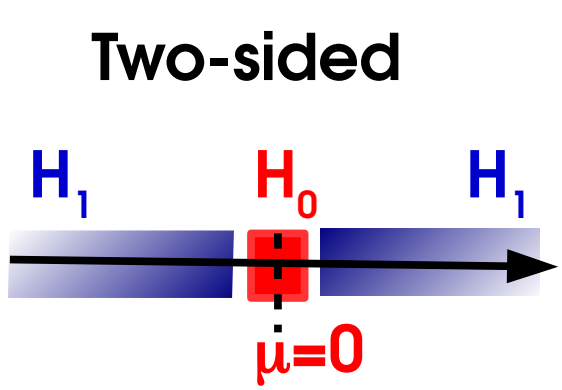
(Continued from yesterday)

# One-sided vs. Two-Sided

If  $\hat{S} < 0$ , is it a *discovery*? (does reject the  $S=0$  hypothesis...)

Usual assumption : only  $\hat{S} > 0$  is a *bona fide* signal

⇒ Change statistic so that  $\hat{S} < 0 \Rightarrow t_0 = 0$  (perfect agreement with  $H_0$ , as for  $\hat{S} = 0$ )



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

$$q_0 = \begin{cases} -2 \log \frac{L(S=0)}{L(\hat{S})} & \hat{S} \geq 0 \\ 0 & \hat{S} < 0 \end{cases}$$

**Test  
Statistic**

$$Z = \Phi^{-1}\left(1 - \frac{p_0}{2}\right)$$

$$Z = \Phi^{-1}(1 - p_0)$$

$p_0$	Z	$p_0$
0.32	1	0.16
0.003	3	0.0015
$6 \times 10^{-7}$	5	$3 \times 10^{-7}$

By convention, factor 2  
in p-values for a given Z

⇒ Same Z in both cases  
for a given signal S

# One-Sided Asymptotics

→ One-sided test:

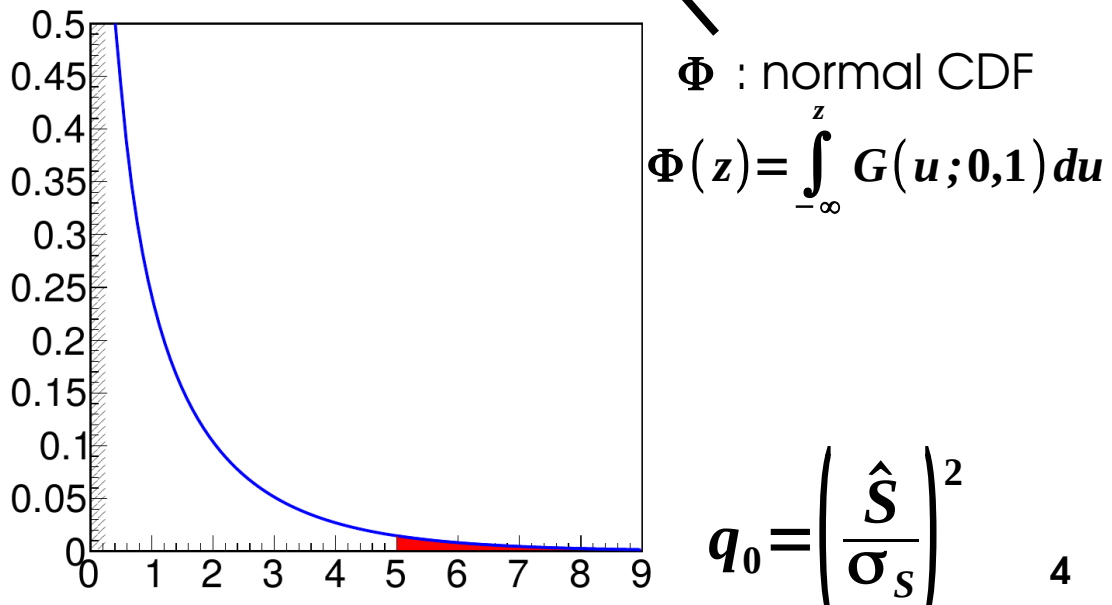
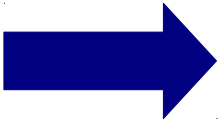
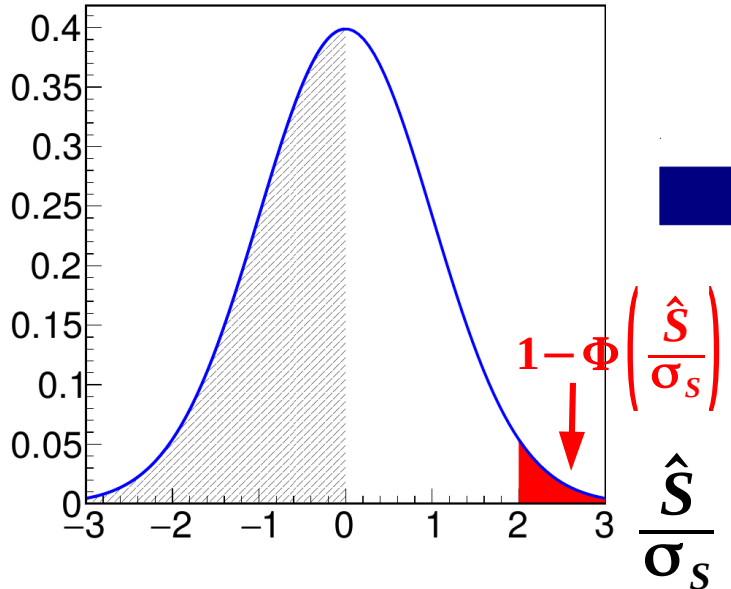


$$q_0 = \begin{cases} -2 \log \frac{L(S=0)}{L(\hat{S})} & \hat{S} \geq 0 \\ 0 & \hat{S} < 0 \end{cases}$$

Asymptotics: "half- $\chi^2$ " distribution:

$$f(q_0 | S=0) = \frac{1}{2} \delta(q_0) + \frac{1}{2} f_{\chi^2(n_{dof}=1)}(q_0)$$

Discovery p-value:  $p_0 = 1 - \Phi(\sqrt{q_0})$     Significance:  $Z = \Phi^{-1}(1 - p_0) = \sqrt{q_0}$



# Example: Gaussian Counting

Count number of events  $n$  in data

→ assume  $n$  large enough so process is Gaussian

→ assume  $B$  is known, measure  $S$

Likelihood :  $L(S) = e^{-\frac{1}{2} \left( \frac{n - (S+B)}{\sqrt{S+B}} \right)^2}$

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

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

Test statistic: assume  $\hat{S} > 0$ ,

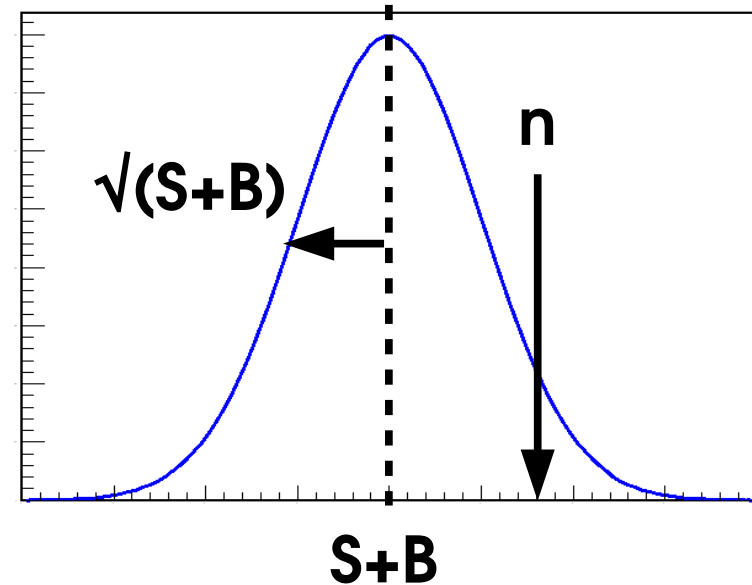
$$q_0 = -2 \log \frac{L(S=0)}{L(\hat{S})} = \lambda(S=0) - \lambda(\hat{S}) = \left( \frac{n-B}{\sqrt{B}} \right)^2 = \left( \frac{\hat{S}}{\sqrt{B}} \right)^2$$

Finally:

$$Z = \sqrt{q_0} = \frac{\hat{S}}{\sqrt{B}}$$

Known formula!

→ Strictly speaking only valid in Gaussian regime



# Example: Poisson Counting

Same problem but now not assuming Gaussianity

$$L(S) = e^{-(S+B)} (S+B)^n \quad \lambda(S) = 2(S+B) - 2n \log(S+B)$$

MLE:  $\hat{S} = n - B$ , same as Gaussian

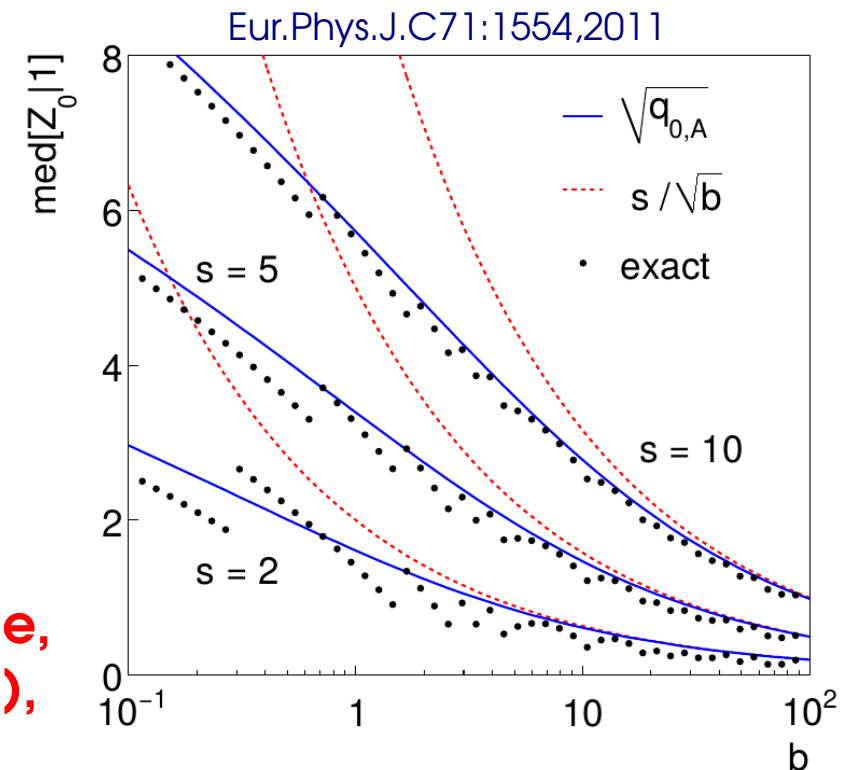
Test statistic (for  $\hat{S} > 0$ ):  $q_0 = \lambda(S=0) - \lambda(\hat{S}) = -2\hat{S} - 2(\hat{S}+B) \log \frac{B}{\hat{S}+B}$

Assuming asymptotic distribution for  $q_0$ ,

$$Z = \sqrt{2 \left[ (\hat{S}+B) \log \left( 1 + \frac{\hat{S}}{B} \right) - \hat{S} \right]}$$

Exact result can be obtained using pseudo-experiments  $\rightarrow$  close to  $\sqrt{q_0}$  result

**Asymptotic formulas justified by Gaussian regime, but remain valid even for small values of  $S+B$  (5!), when  $S$  itself is not Gaussian**



# Example: Multi-bin counting

Likelihood : 
$$L(S) = \prod_{i=1}^N \text{Pois}(n_i; S f_i + B_i)$$

Assume Gaussianity:

$$\lambda(S) = \sum_{i=1}^N \left( \frac{n_i - (S f_i + B_i)}{\sqrt{S f_i + B_i}} \right)^2$$
$$\hat{S} = \frac{\sum_{i=1}^N \frac{f_i}{B_i} \hat{S}_i}{\sum_{i=1}^N \frac{f_i^2}{B_i}}$$

$\hat{S}_i = n_i - B_i$   
↓

**Test statistic:** assuming  $\hat{S} > 0$ ,

$$q_0 = \lambda(S=0) - \lambda(\hat{S}) = \left( \hat{S} \sqrt{\sum_{i=1}^N \frac{f_i^2}{B_i}} \right)^2$$

**Asymptotics:**

$$Z = \sqrt{q_0} = \frac{\hat{S}}{\left( \sum_{i=1}^N \frac{f_i^2}{B_i} \right)^{-1/2}}$$

→

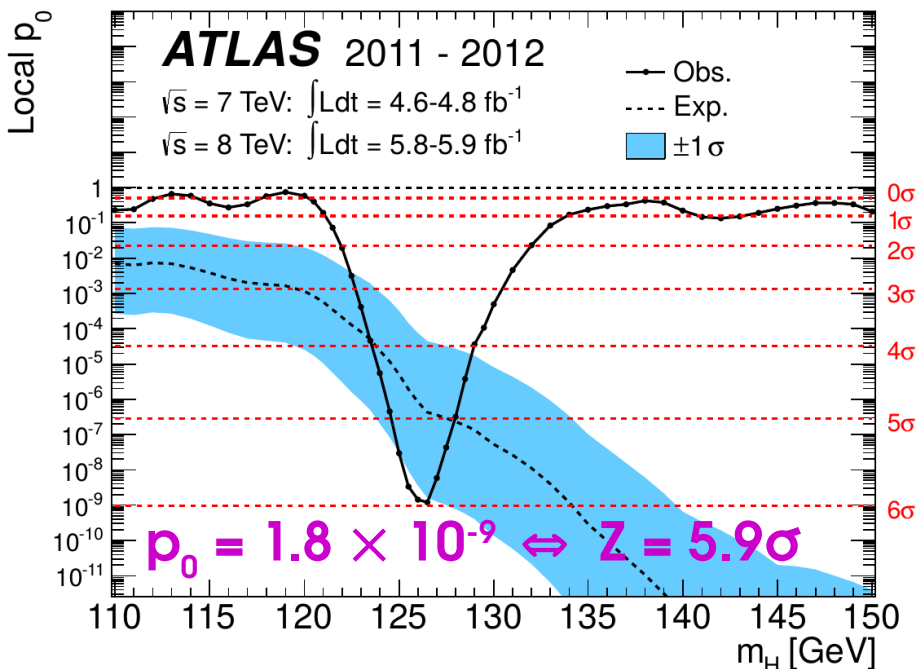
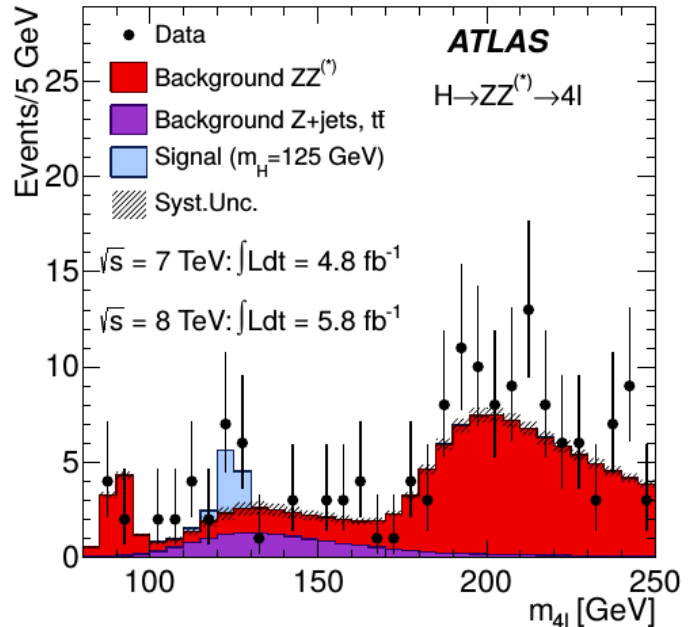
Combined uncertainty  
on  $\hat{S}$  from all the bins

**Always better than**

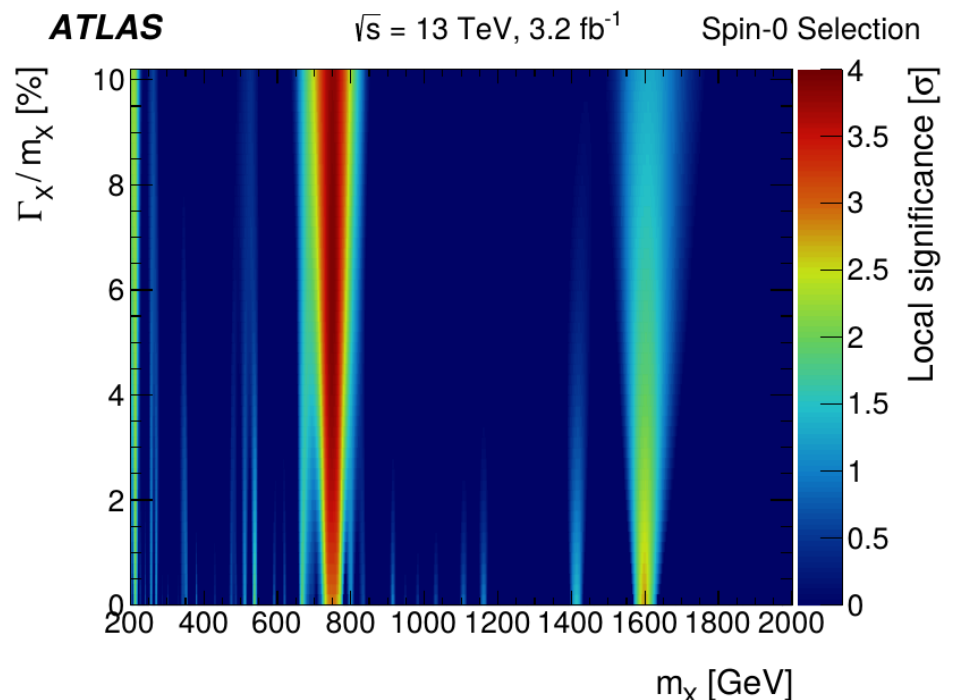
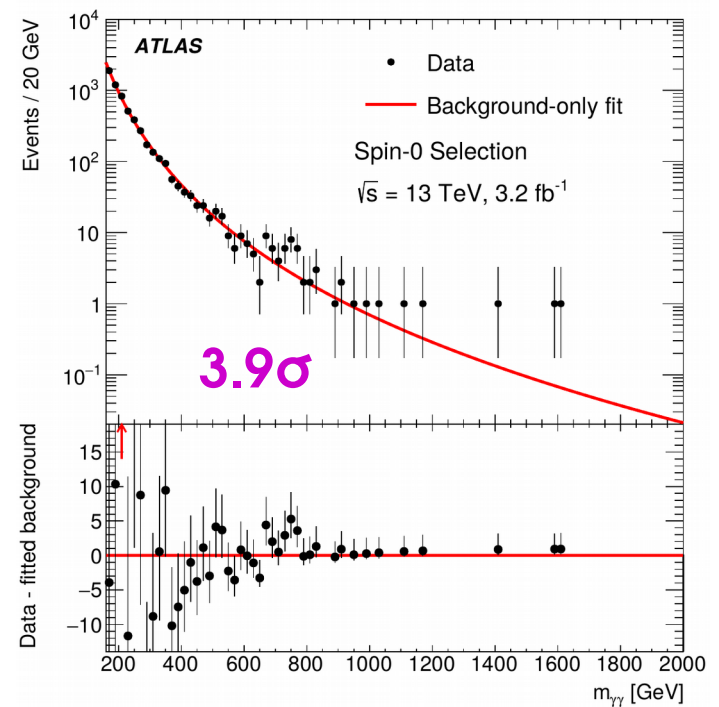
- Any bin by itself (for same  $\hat{S}$ )
- All bins merged together

# Some Examples

## Higgs Discovery: Phys. Lett. B 716 (2012) 1-29



$$Z = \Phi^{-1}(1 - p_0)$$





# Takeaways

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  $\frac{L(H_0)}{L(H_1)}$

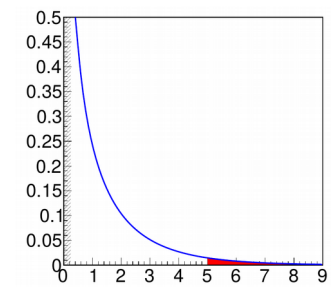
To test for **discovery**, use  $q_0 = \begin{cases} -2 \log \frac{L(S=0)}{L(\hat{S})} & \hat{S} \geq 0 \\ 0 & \hat{S} < 0 \end{cases}$

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

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

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

# What was the question ?



Definition of the p-value:

$$\text{p-value} = \frac{\text{number of signal-like outcomes with only background present}}{\text{all outcomes with only background present}}$$

So  $5\sigma$  significance ( $p_0 \sim 10^{-7}$ )  $\Leftrightarrow$  *Occurs once in  $10^7$  if only background present*

However this is **NOT** "~~One chance in  $10^7$  to be a fluctuation~~"

The first statement is about **data probabilities** –  $P(\text{data}; H_0)$

The second is on  $P(H_0)$  itself – not addressed in the framework described so far  
→ makes sense in a **Bayesian** context, more on this later in these lectures.

It's also a different statement (although they sometimes get confused)

→ If a signal outcome is also very unlikely, **we may not want to reject  $H_0$ , even with  $p_0 \sim 10^{-7}$ .**

# What was the question ?

e.g. Faster-than-light neutrino anomaly

$$(v-c)/c = (2.37 \pm 0.32 \text{ (stat.) } ^{+0.34}_{-0.24} \text{ (sys.)}) \times 10^{-5} \quad \mathbf{6.2\sigma \text{ above } c}$$

*“despite the large significance of the measurement reported here and the stability of the analysis, the potentially great impact of the result motivates the continuation of our studies in order to investigate possible still unknown systematic effects that could explain the observed anomaly.”*

⇒ Very unlikely to be a background fluctuation, but hard to believe **since alternative ( $v > c$ ) is far-fetched**

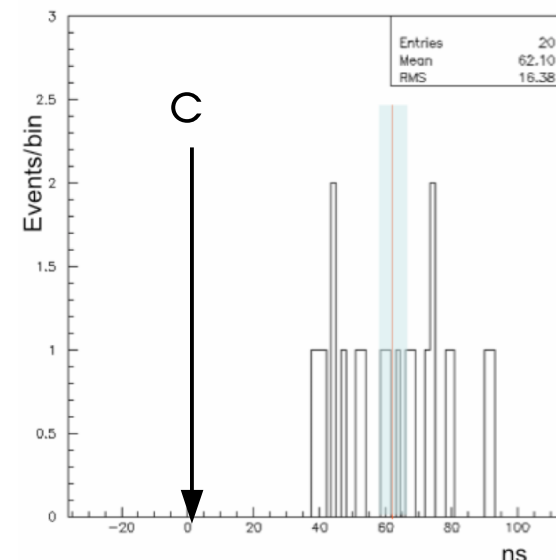
**Alternative:**  $P(\text{fluctuation}) = \frac{\text{number of signal-like outcomes with only B present}}{\text{number of signal-like outcomes from any source (S or B)}}$

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

→ Needs **a priori P(S) and P(B)** → Bayesian methods, discussed later

→ In frequentist context, only have  $p_0 = P(\text{deviation} | B)$

⇒ **However usually same conclusion, assuming P(S) is not  $\ll p_0$ ...**



**“Extraordinary claims require extraordinary evidence”**

# Outline

---

Yesterday:

Statistics basics for HEP

Describing HEP measurements

Computing statistics results:

Discovery

**Today:**

**Computing statistics results:**

Limits

Confidence intervals

**Profiling**

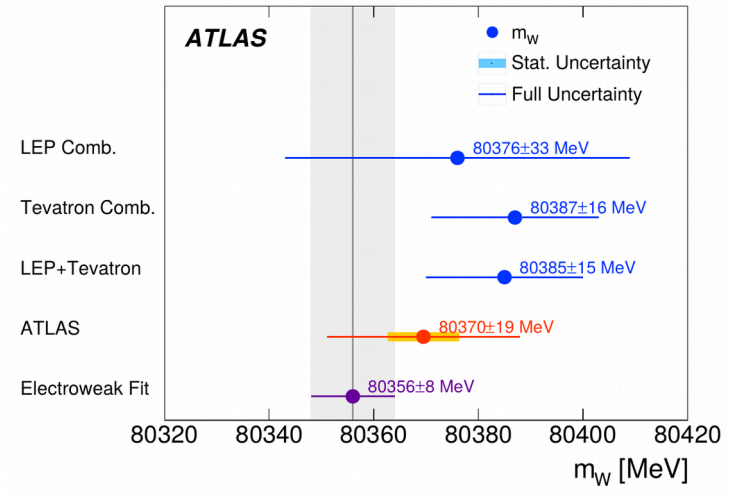
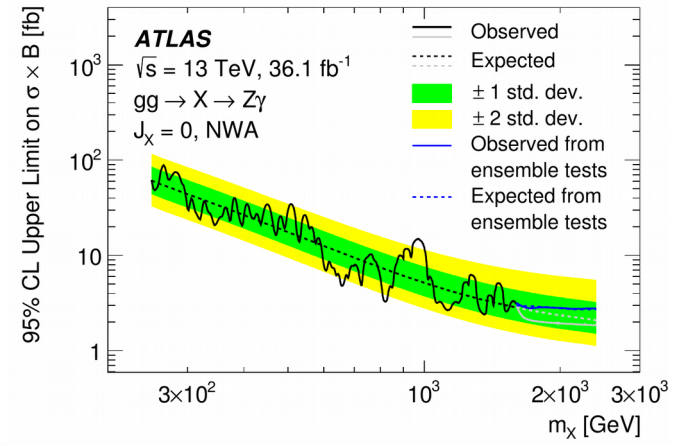
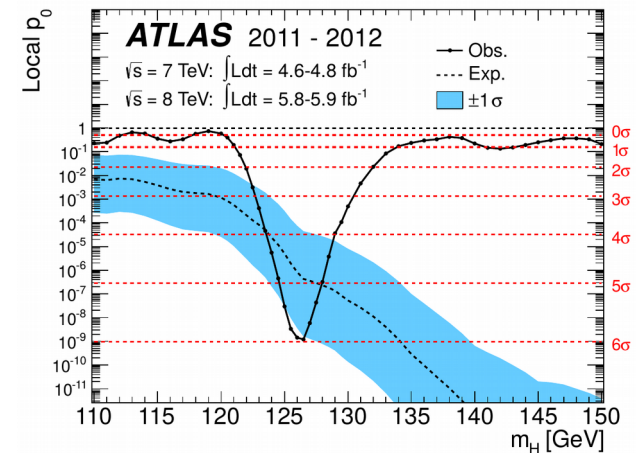
**Look-Elsewhere Effect**

**Bayesian methods**

**Tomorrow:** Practical modeling, Unfolding

# Usual Statistical 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 ?



---

# Upper Limits

# Hypothesis tests for Limits

If no signal in data, testing for discovery not very relevant (report  $0.2\sigma$  excess ?)

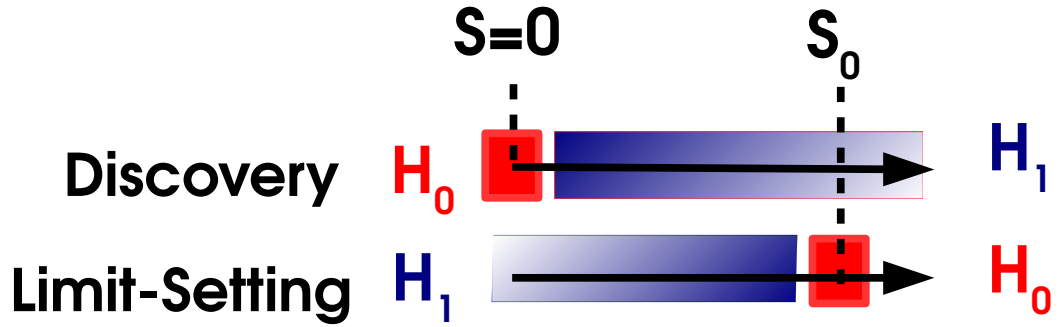
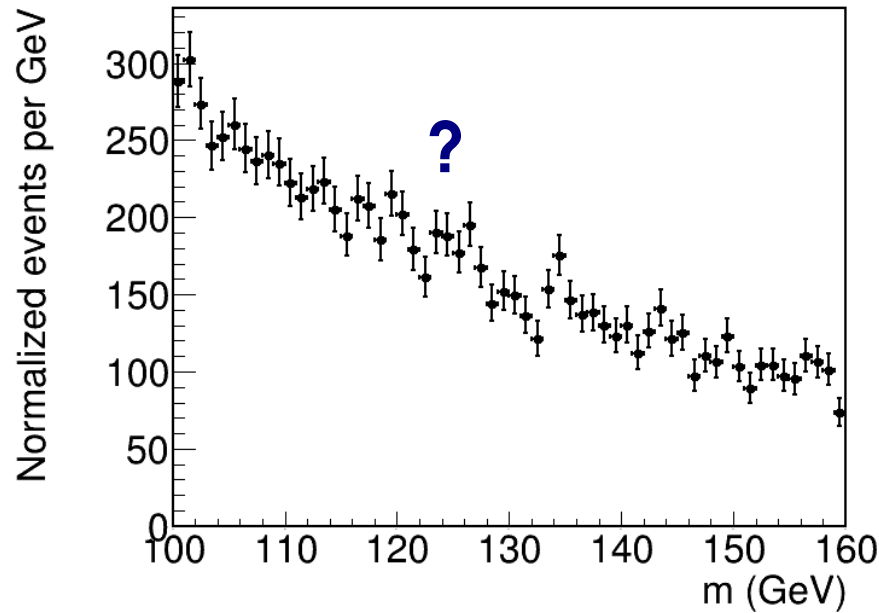
→ More interesting to **exclude large signals** → **Upper limits on signal yield**

For **discovery**

- Try to exclude  $H_0 : S=0$
- Alternative :  $H_1 : S > 0$
- Report p-value for the test (or Z)

For **limit-setting**:

- Try to exclude  $H_0 : S=S_0$
- Alternative :  $H_1 : S < S_0$
- Usually, **adjust  $S_0$  to get a predefined p-value** (typically 5%)
- **Confidence Levels**:  $CL = 1 - p$  ( $p = 5\% \Leftrightarrow 95\% CL$ )



# Hypothesis tests for Limits

If no signal in data, testing for discovery not very relevant (report  $0.2\sigma$  excess ?)

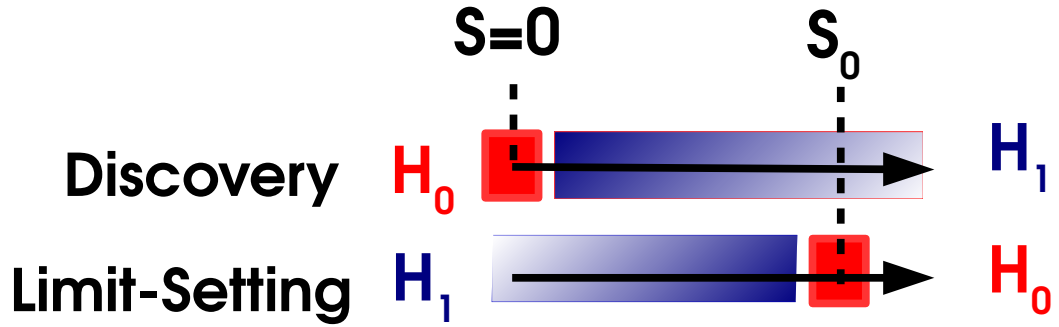
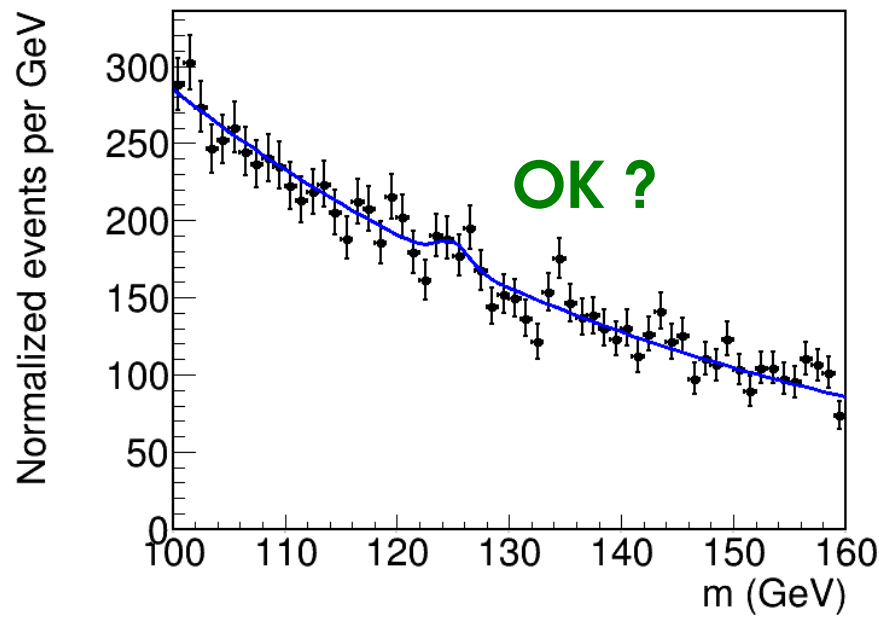
→ More interesting to **exclude large signals** → **Upper limits on signal yield**

For **discovery**

- Try to exclude  $H_0 : S=0$
- Alternative :  $H_1 : S > 0$
- Report p-value for the test (or Z)

For **limit-setting**:

- Try to exclude  $H_0 : S=S_0$
- Alternative :  $H_1 : S < S_0$
- Usually, **adjust  $S_0$  to get a predefined p-value** (typically 5%)
- **Confidence Levels**:  $CL = 1 - p$  ( $p = 5\% \Leftrightarrow 95\% CL$ )





# Hypothesis tests for Limits

If no signal in data, testing for discovery not very relevant (report  $0.2\sigma$  excess ?)

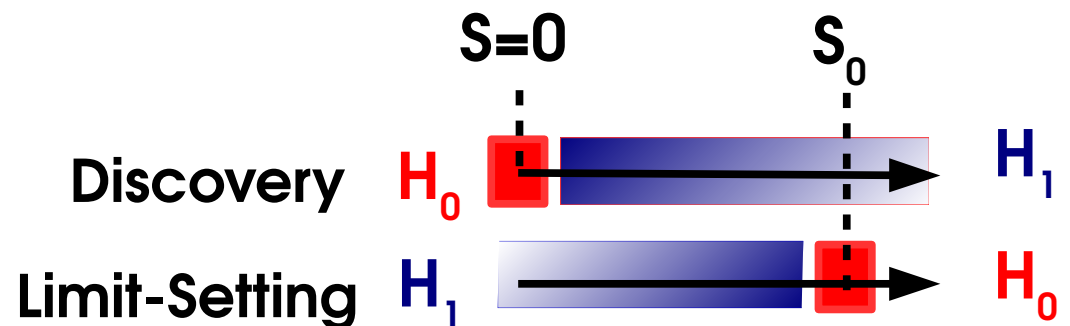
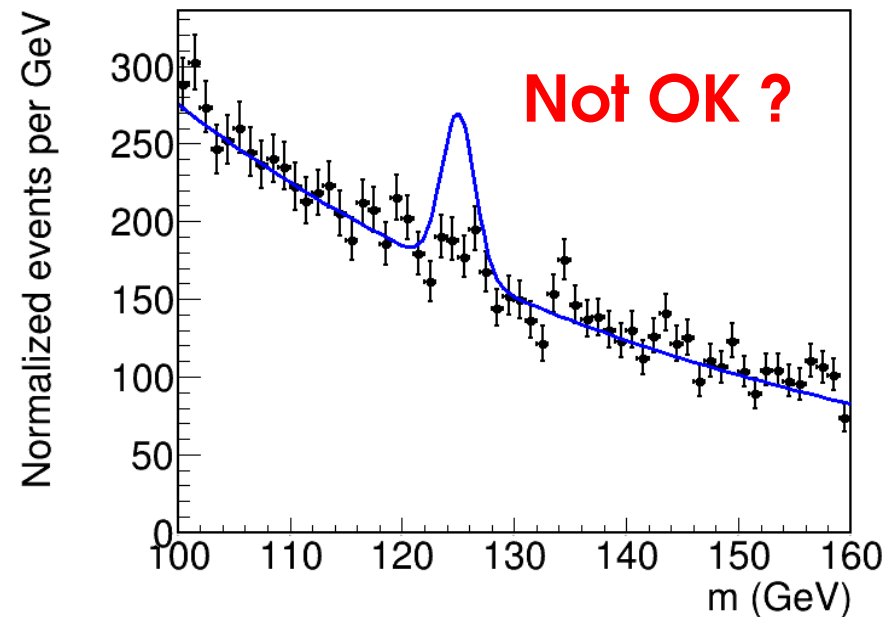
→ More interesting to **exclude large signals** → **Upper limits on signal yield**

For **discovery**

- Try to exclude  $H_0 : S=0$
- Alternative :  $H_1 : S > 0$
- Report p-value for the test (or Z)

For **limit-setting**:

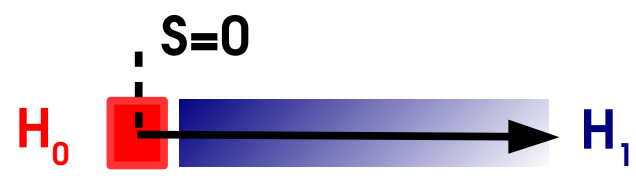
- Try to exclude  $H_0 : S=S_0$
- Alternative :  $H_1 : S < S_0$
- Usually, **adjust  $S_0$  to get a predefined p-value** (typically 5%)
- **Confidence Levels**:  $CL = 1 - p$  ( $p = 5\% \Leftrightarrow 95\% \text{ CL}$ )



# Test Statistic for Limit-Setting

**Discovery :**

- $H_0 : S = 0$
- $H_1 : S > 0$



$$q_0 = -2 \log \frac{L(S=0)}{L(\hat{S})}$$

Compare  
 ← Likelihood of  $H_0$   
 ← Likelihood of  $H_1$

**Limit-setting**

- $H_0 : S = \mu_0$
- $H_1 : S < \mu_0$



$$q_{S_0} = -2 \log \frac{L(S_0)}{L(\hat{S})}$$

Compare  
 ← Likelihood of  $H_0$   
 ← Likelihood of  $H_1$

$\hat{S} \sim S_0$  (no exclusion) :  $q_{S_0} \sim 0$   
 $\hat{S} \ll S_0$  (good exclusion) :  $q_{S_0} \gg 1$

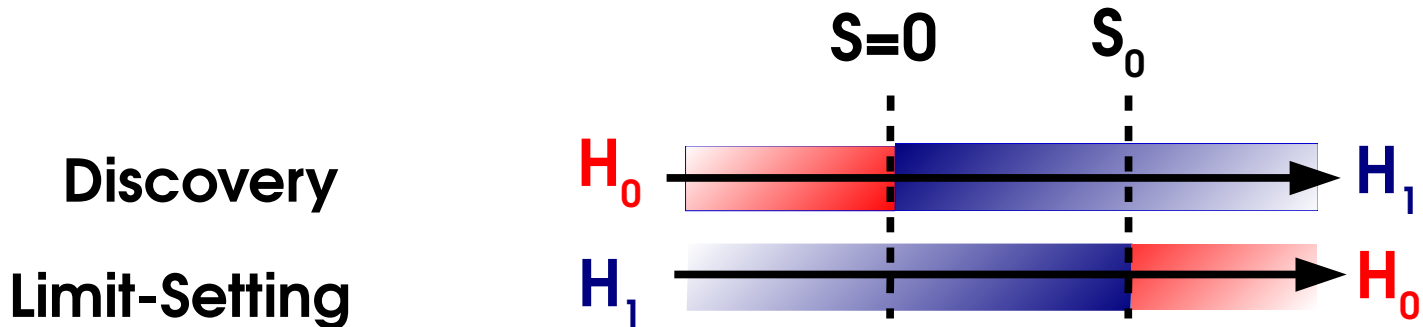
Same as  $q_0$  : large values  
 $\Rightarrow$  good rejection of  $H_0$ .

# One-sided Test Statistic

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  $\hat{S} < 0$ .

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

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

# Inversion : Getting the limit for a given CL

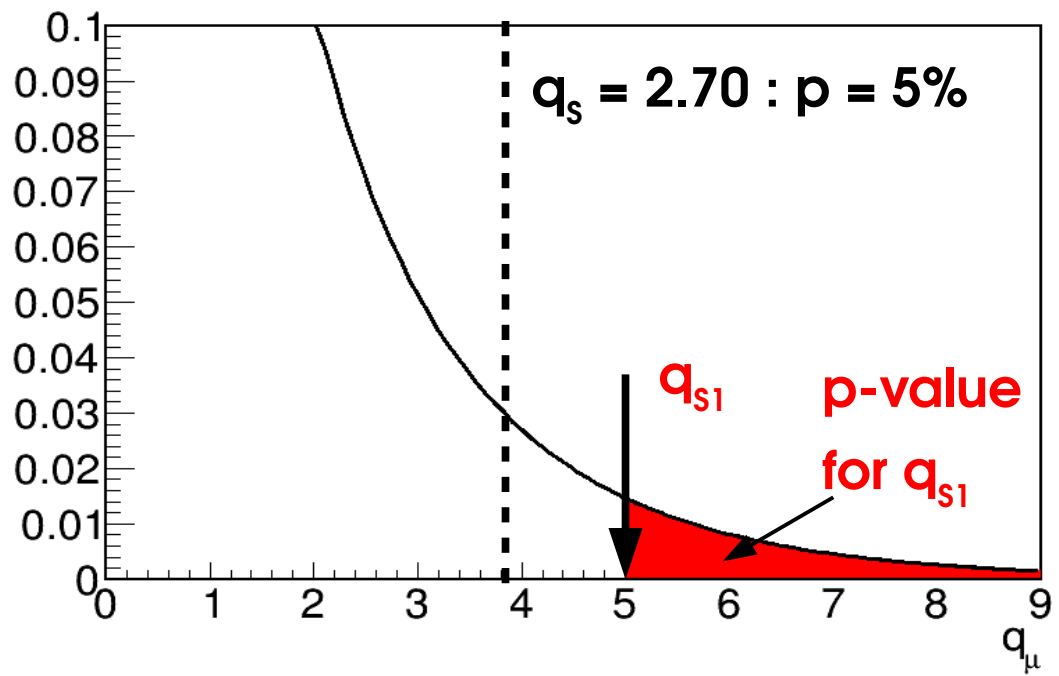
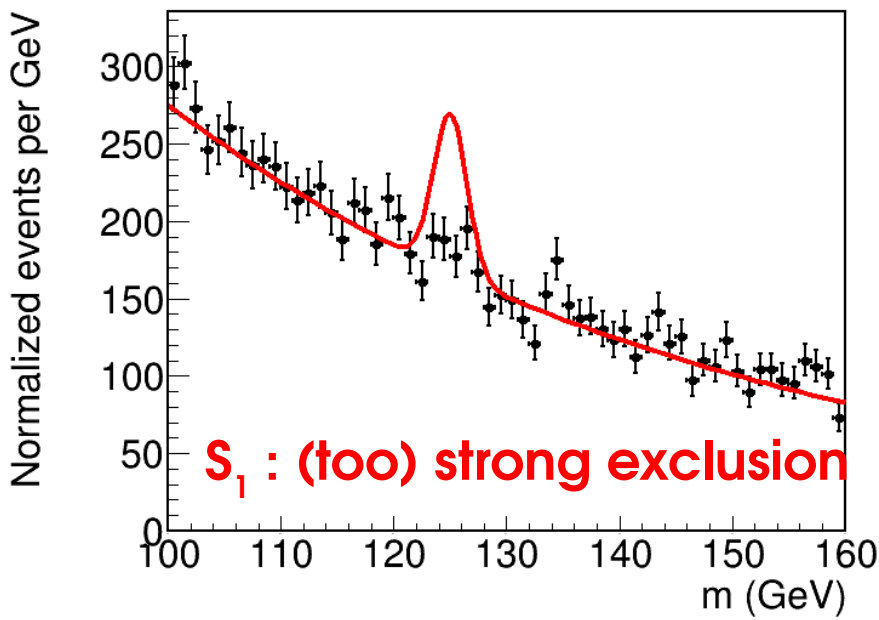
## Procedure

- 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%	$q_s > 5.41$



# Inversion : Getting the limit for a given CL

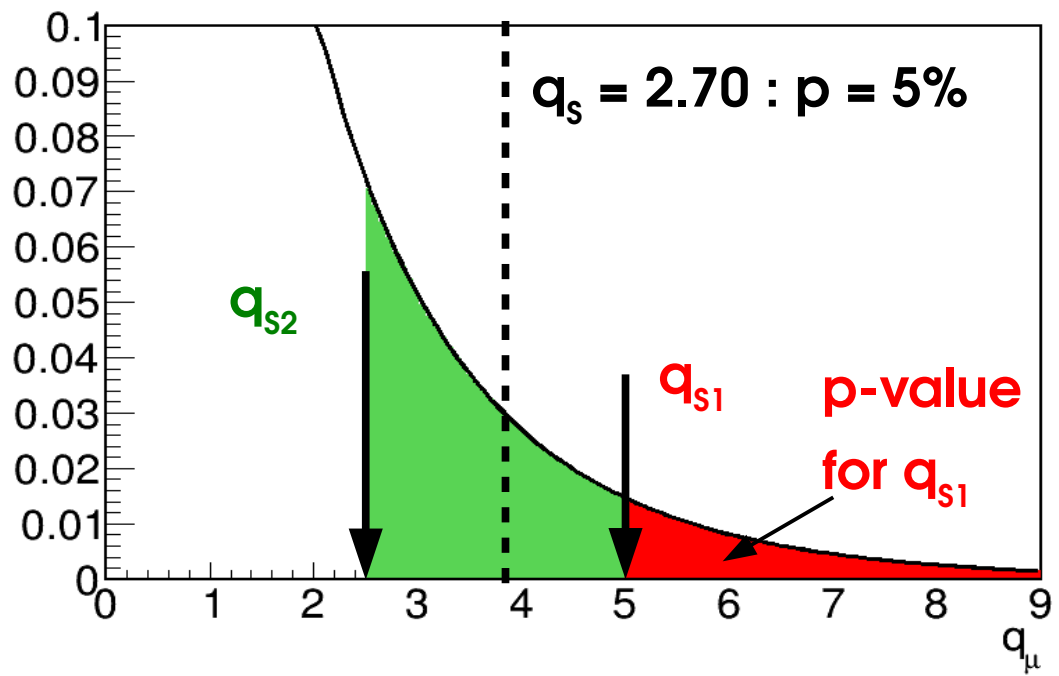
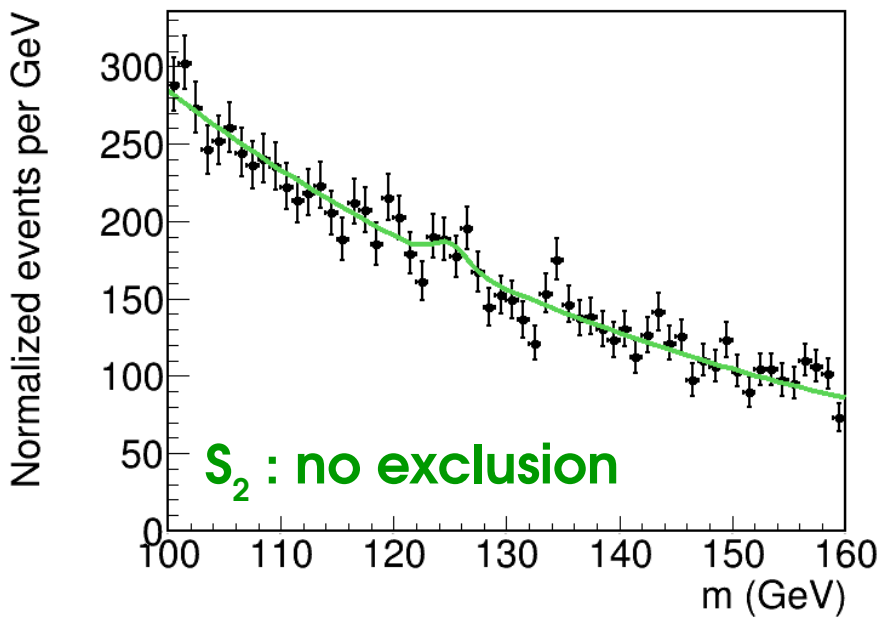
## Procedure

- 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%	$q_s > 5.41$



# Inversion : Getting the limit for a given CL

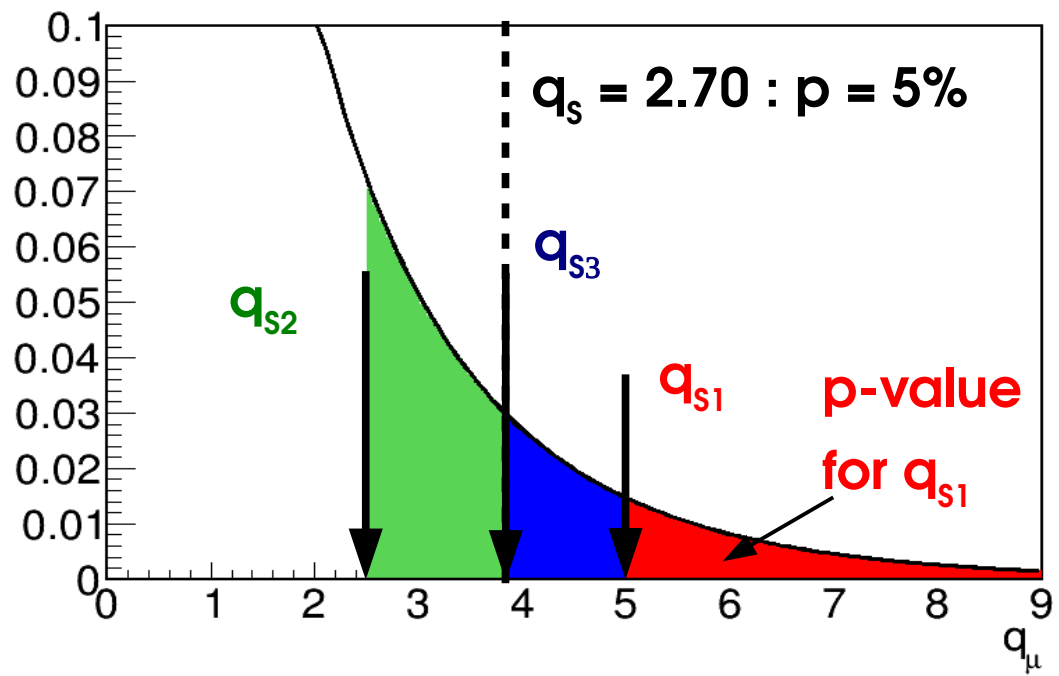
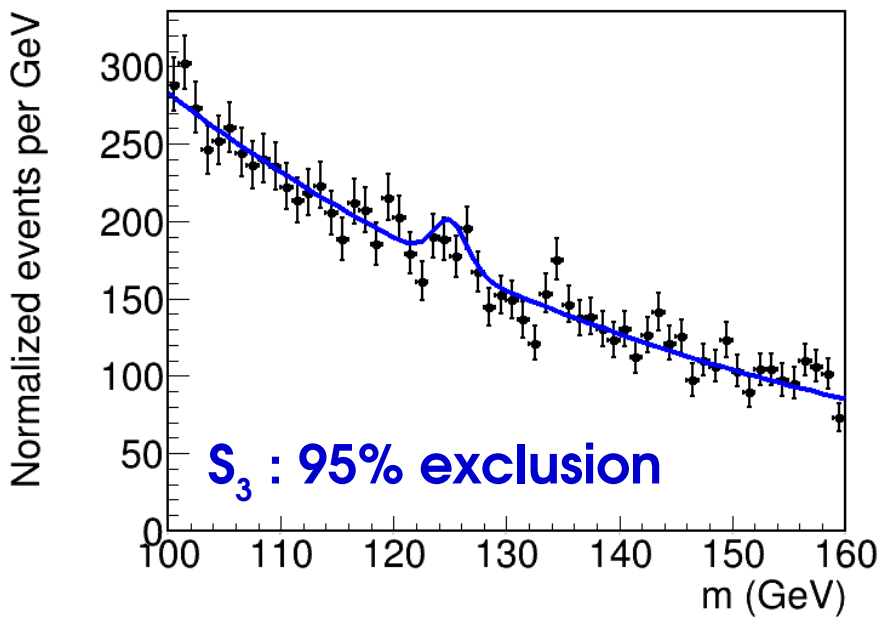
## Procedure

- 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%	$q_S > 5.41$



# Upper Limits: Gaussian Example

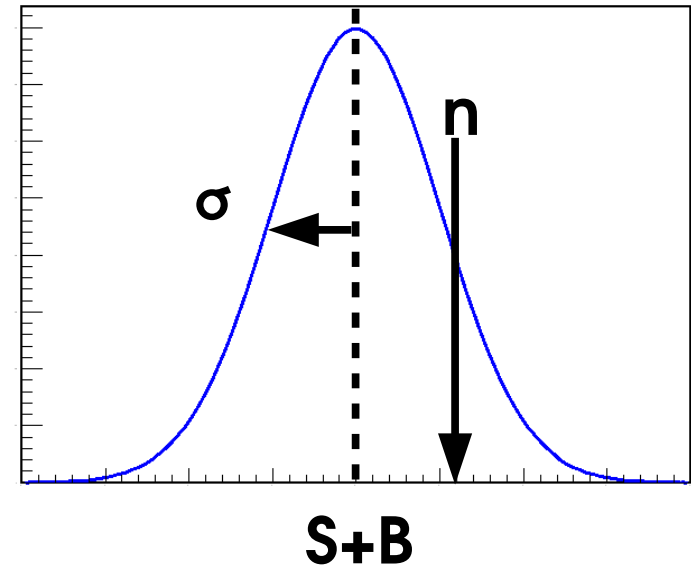
Usual Gaussian counting example with known B:

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

**Reminder:**

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

Significance:  $Z = \hat{S} / \sqrt{B}$



Compute the 95% CL upper limit on S:

$$q_{S_0} = -2 \log \frac{L(S=S_0)}{L(\hat{S})} = \lambda(S_0) - \lambda(\hat{S}) = \left( \frac{n - (S_0 + B)}{\sigma_S} \right)^2 = \left( \frac{S_0 - \hat{S}}{\sigma_S} \right)^2 \quad \text{for } S_0 > \hat{S}$$

so  $q_{S_0} = 2.70$  for  $S_0 = \hat{S} + \sqrt{2.70} \sigma_S$

And finally  $S_{\text{up}} = \hat{S} + 1.64 \sigma_S$  at 95 % CL

# Upper Limit Pathologies

Upper limit:  $S_{up} \sim \hat{S} + 1.64 \sigma_s$ .

**Problem:** for negative  $\hat{S}$ , get **very** good observed limit.

→ For  $\hat{S}$  sufficiently negative, even  $S_{up} < 0$  !

How can this be ?

→ **Background modeling issue ?...** Or:

→ This is a **95%** limit

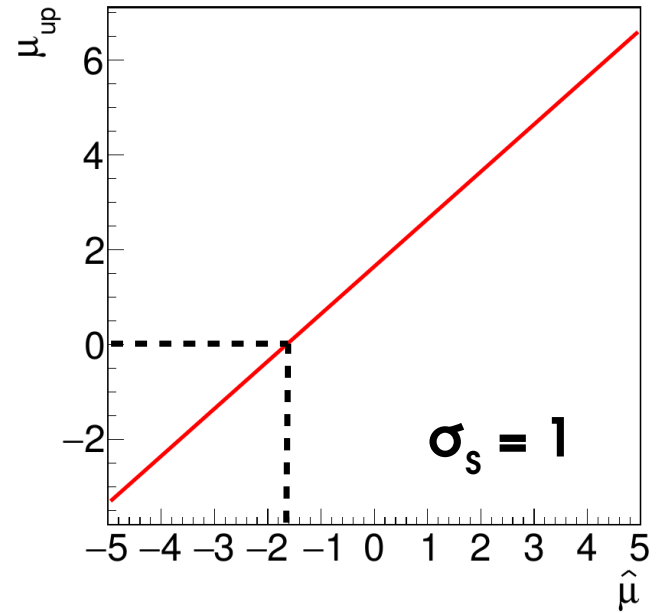
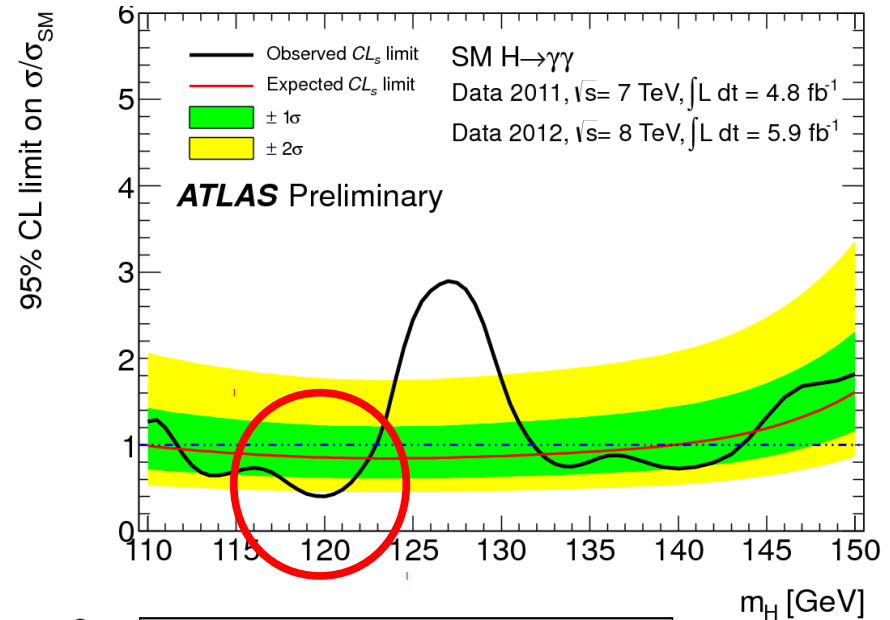
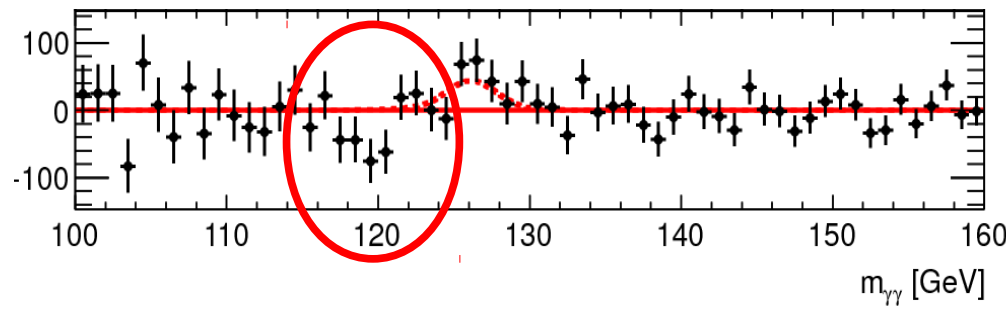
⇒ **5% of the time, the limit wrongly excludes the true value, e.g.  $S^*=0$ .**

But if we assume  $S$  must be  $>0$ , we know a priori this is just a fluctuation.

## Options

→ **live with it:** sometimes report limit  $< 0$

→ **Special procedure to avoid these cases**

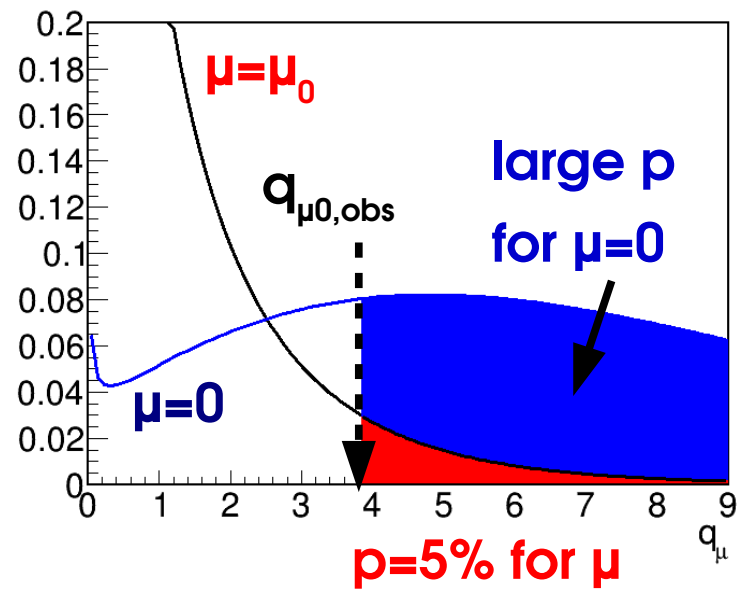
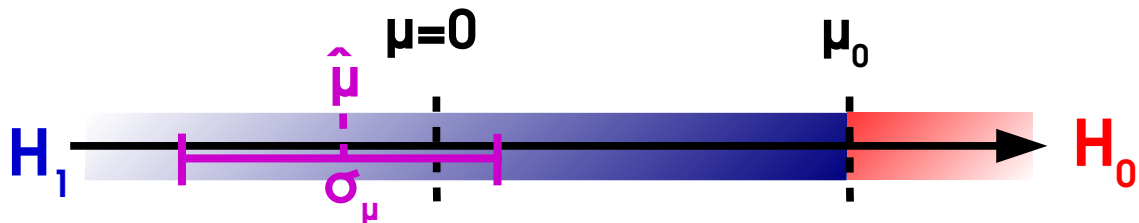




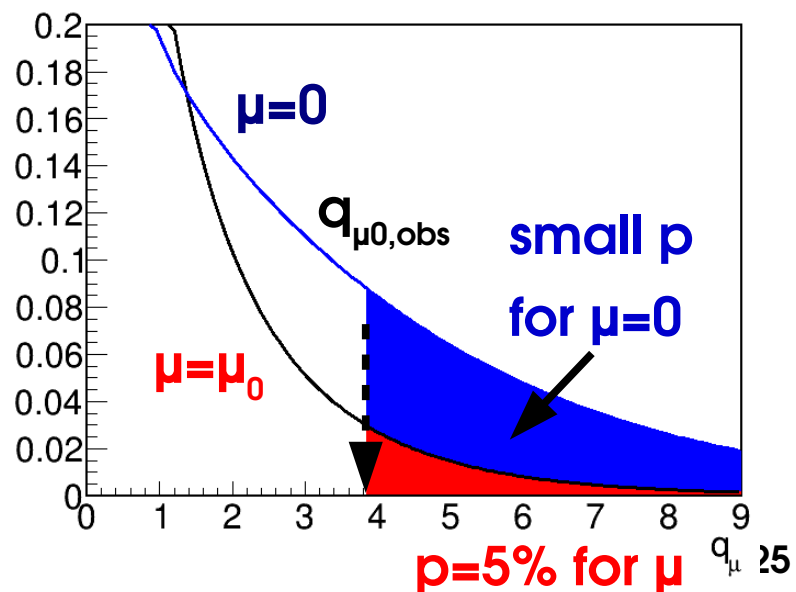
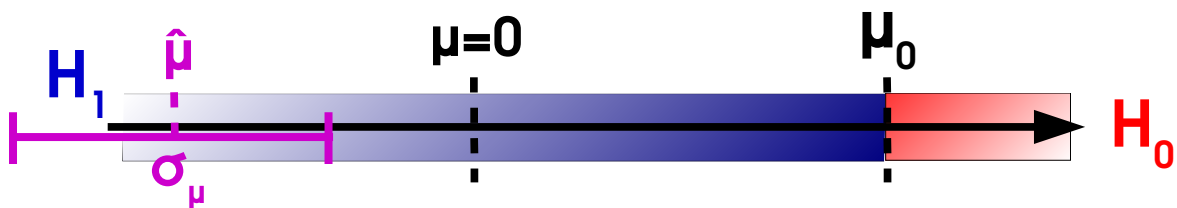
# Upper Limit Pathologies

When setting limits, goal is to exclude large  $\mu$ , to indicate that  $\mu \sim 0$ . What happens at  $\mu=0$  ?

**Normal case:**  $\hat{\mu} \sim 0$ ,  $\mu=0$  not excluded :  
 $\mu_{up} = \hat{\mu} + 1.64 \sigma_\mu > 0$ , large p-value for  $\mu=0$



**Pathological case,** very negative  $\hat{\mu}$ ,  $\mu=0$  also excluded :  
 $\mu_{up} = \hat{\mu} + 1.64 \sigma_\mu < 0$ , p-value for  $\mu=0$  also small



→ However we know a priori that  $\mu \geq 0$   
 ⇒ Inject this information into the procedure

Usual solution in HEP : **CL<sub>s</sub>**.

→ Compute modified p-value

$$p_{CL_s} = \frac{p_{\mu_0}}{p_0}$$

- **p<sub>μ<sub>0</sub></sub>** is the usual p-value (5%)
- **p<sub>0</sub>** is the p-value computed under H(μ=0).

⇒ **Rescale** exclusion at μ<sub>0</sub> by exclusion at μ=0.

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

**Good case** : p<sub>0</sub> ~ O(1)

p<sub>CL<sub>s</sub></sub> ~ p<sub>μ<sub>0</sub></sub> ~ 5%, **no change**.

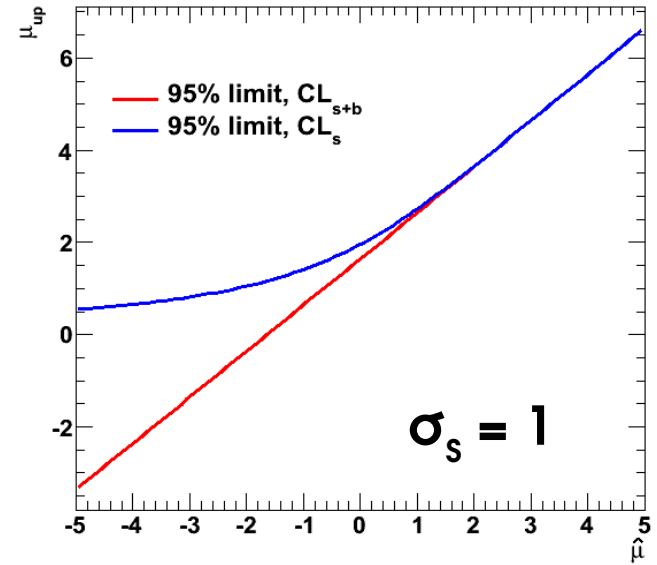
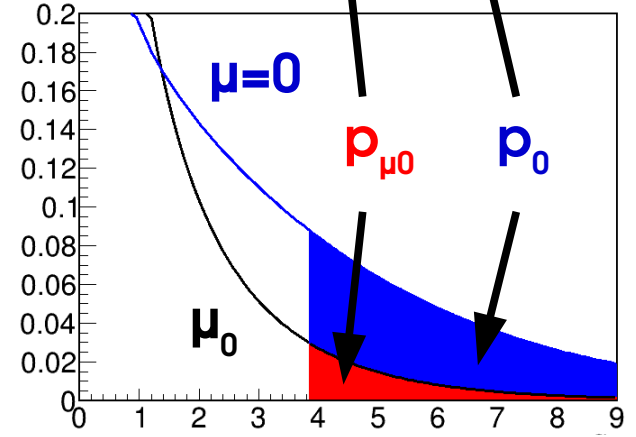
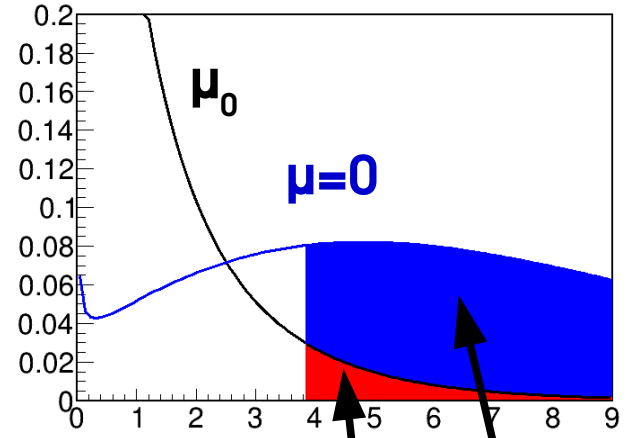
**Pathological case** : p<sub>0</sub> ≪ 1

p<sub>CL<sub>s</sub></sub> ~ p<sub>μ<sub>0</sub></sub>/p<sub>0</sub> ≫ 5%

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

**Drawback: overcoverage**

→ limit is actually >95% CL for small p<sub>0</sub>.



# CL<sub>s</sub> : 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<sub>s+b</sub> limit:  $S_{up} = \hat{S} + 1.64 \sigma_S$  at 95% CL

CL<sub>s</sub> upper limit : still have

so need to solve

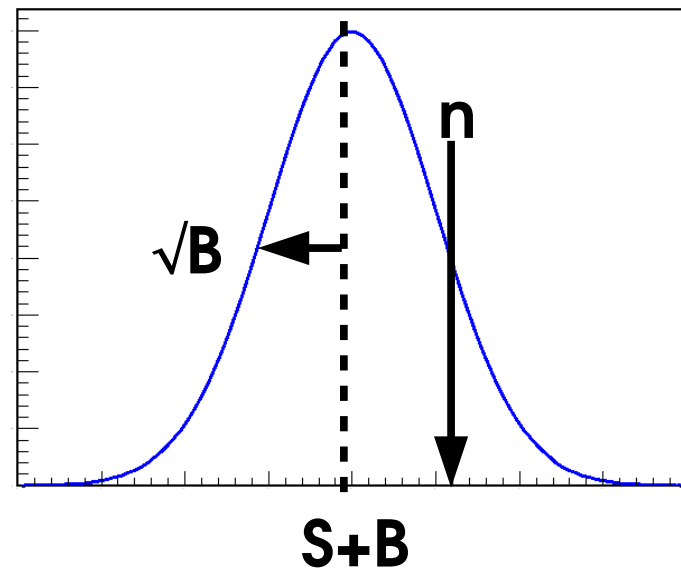
$$q_{S_0} = \left( \frac{S_0 - \hat{S}}{\sigma_S} \right)^2 \quad (\text{for } S_0 > \hat{S})$$

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

for  $\hat{S} = 0$ ,

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

$\Phi(0) = 0.5 \Rightarrow$  at 95% CL, **CL<sub>s</sub> :  $S_{up} = 1.96 \sigma_S$**     **CL<sub>s+b</sub> :  $S_{up} = 1.64 \sigma_S$**



$\hat{S} \sim G(S, \sigma_S)$  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<sub>s</sub>: Poisson Rule of Thumb

Same exercise, for the Poisson case

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

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

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

**⇒ Rule of thumb: when n<sub>obs</sub>=0, the CL<sub>s</sub> 95% CL limit is 3 events (for any B)**

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

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

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

⇒ S<sub>up</sub> ~ 2, exact value depends on B

⇒ Asymptotics not valid in this case – need to use exact results, or toys

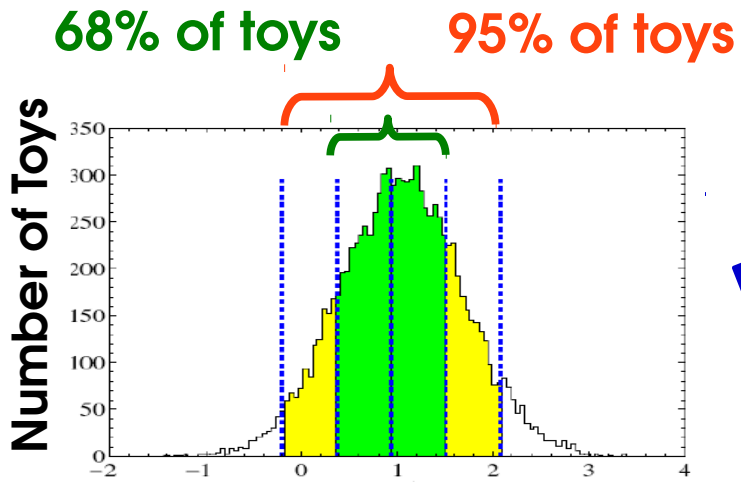
# Expected Limits: Toys

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

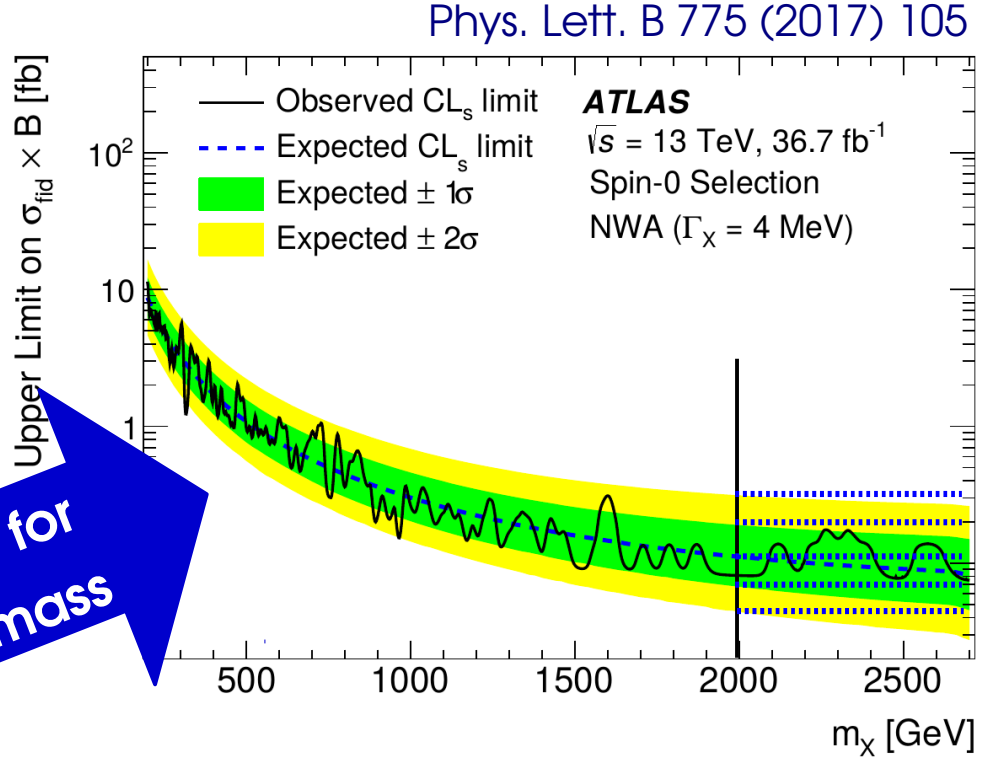
Two main ways to compute:

→ **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



# 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)**

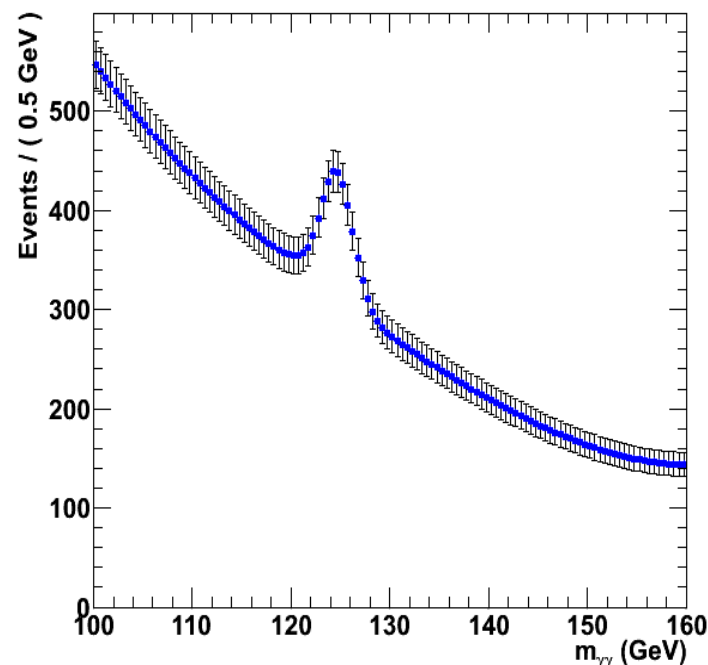
- Get bands from asymptotic formulas:

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<sub>s</sub> : Gaussian Bands

Usual Gaussian counting example with known B:  
95% CL<sub>s</sub> upper limit on S:

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

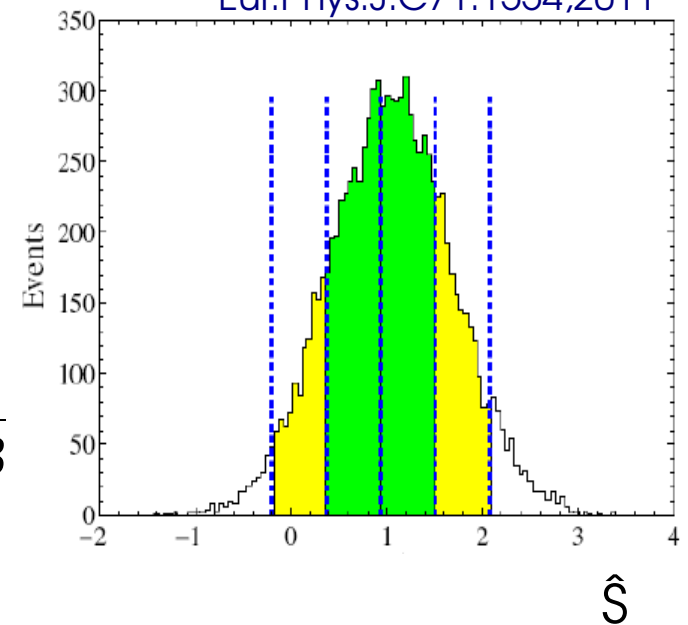
Compute expected bands for S=0:

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

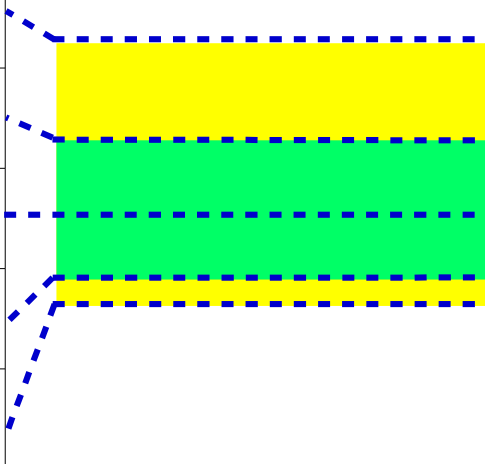
$$S_{\text{up,exp}}^0 = 1.96 \sigma_S$$

→  **$\pm n\sigma$  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



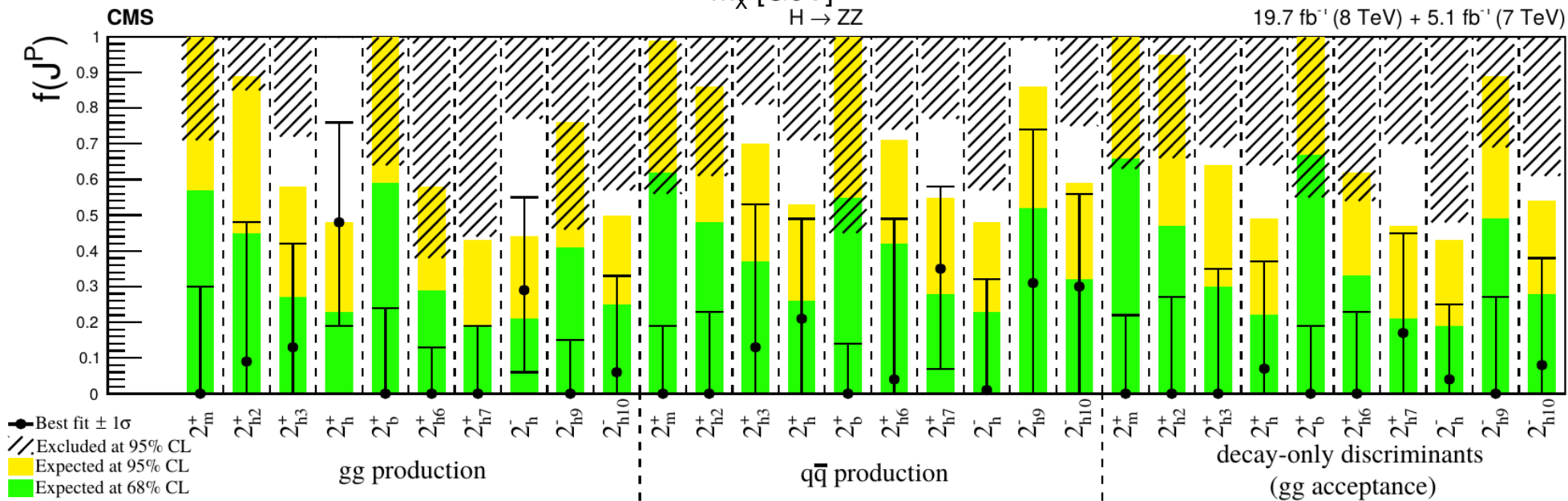
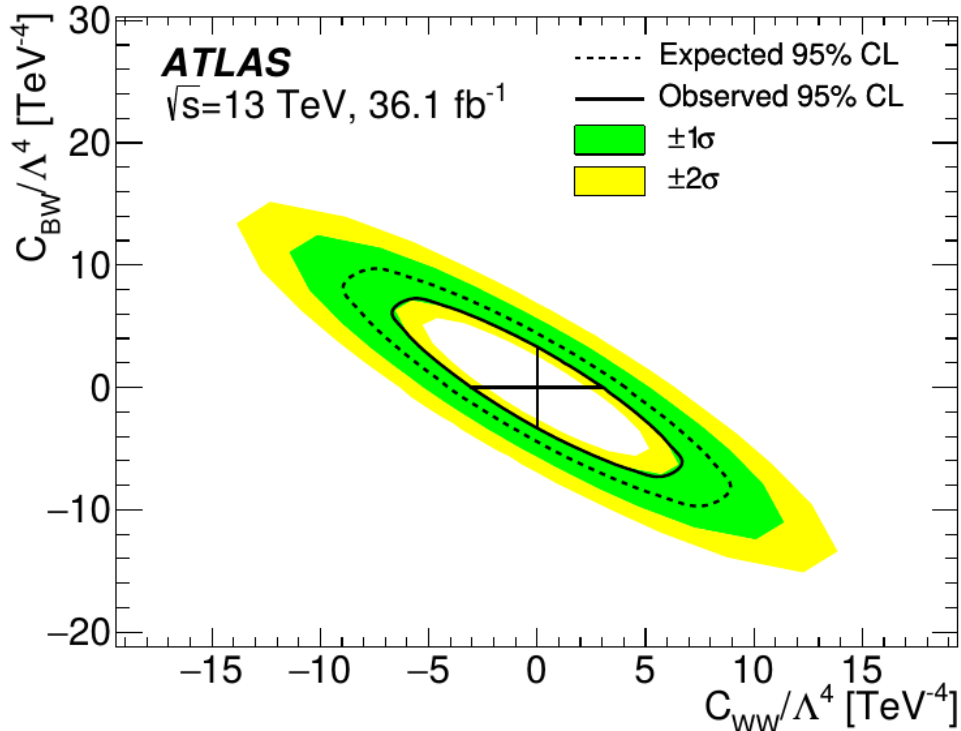
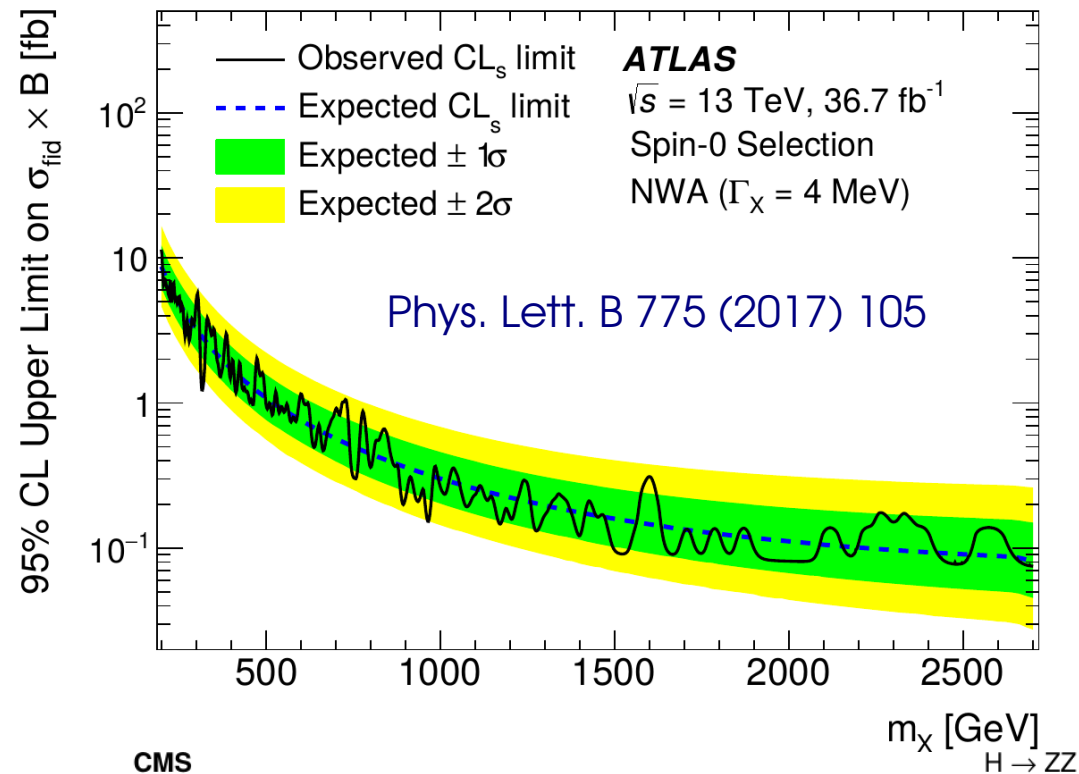
## CLs :

- 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...

# Upper Limit Examples

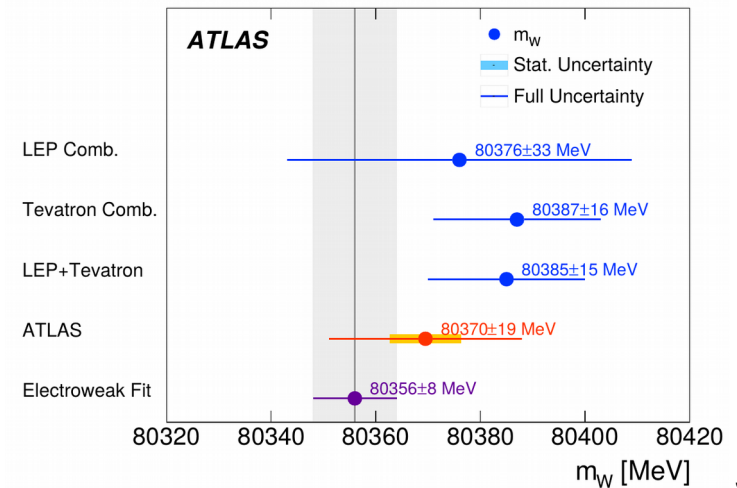
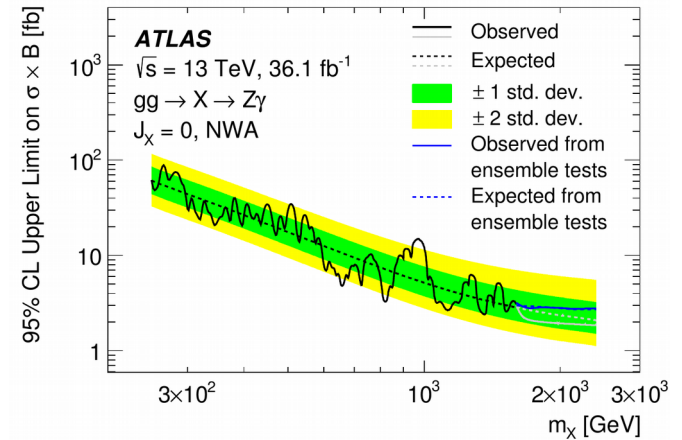
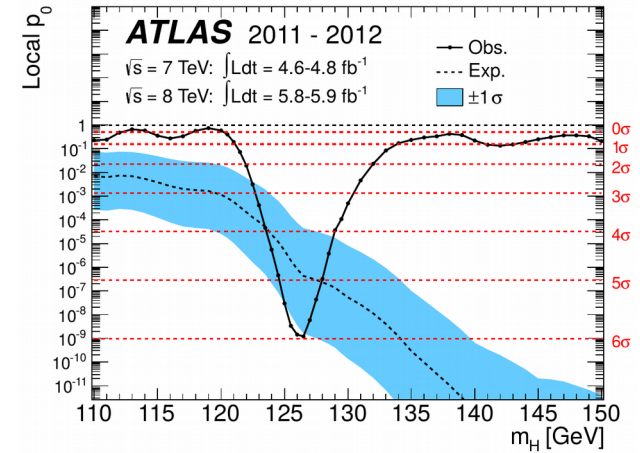
ATLAS 2015-2016 4l aTGC Search





# Usual Statistical 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 ?



# Outline

---

## Computing statistics results:

Limits

**Confidence intervals**

Profiling

Look-Elsewhere Effect

Bayesian methods

---

# Confidence Intervals

# 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  $\mu$ !

→  $\mu^*$  is a **fixed number**, not a random variable

But we can invert the relation:

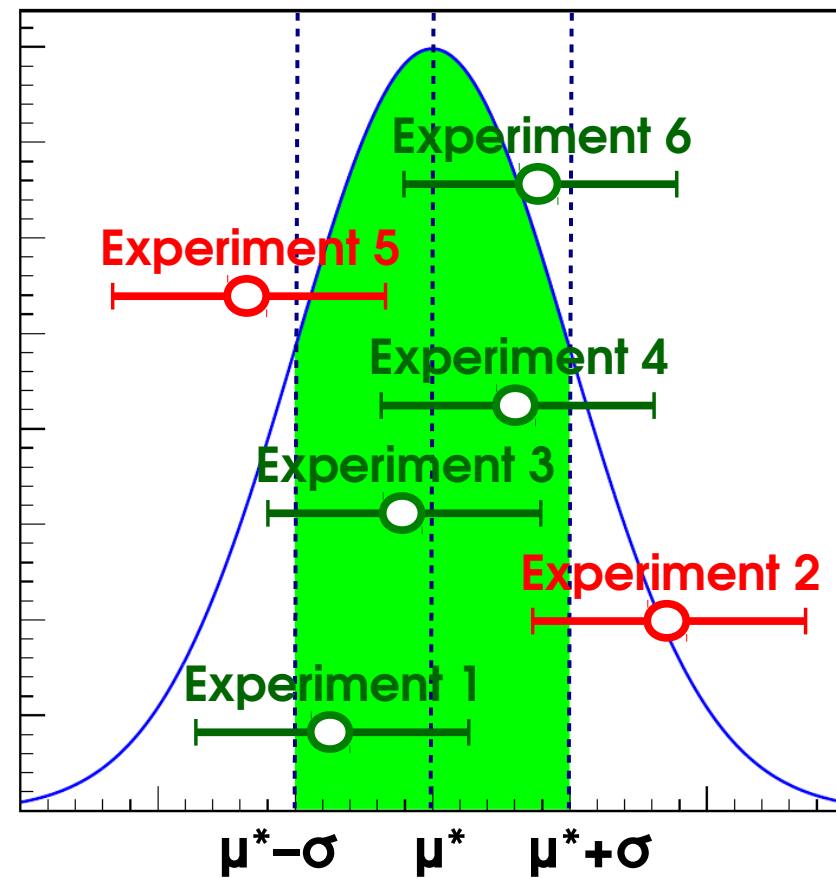
$$P(\mu^* - \sigma < \hat{\mu} < \mu^* + \sigma) = 68\%$$

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

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

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

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

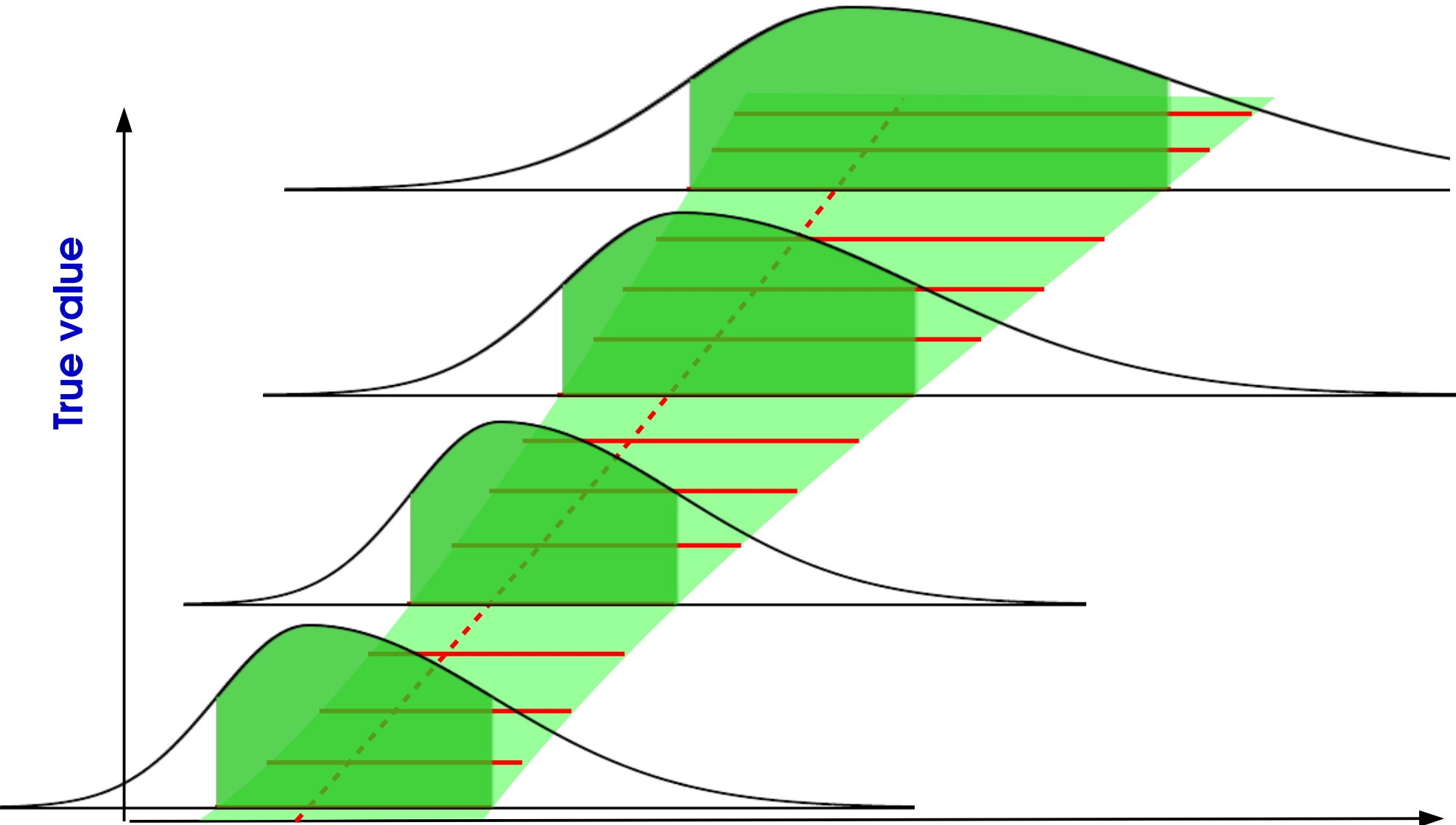


Works in the same way for other interval sizes:  $[\hat{\mu} - Z\sigma, \hat{\mu} + Z\sigma]$  with

Z	1	1.96	2
CL	0.68	0.95	0.955

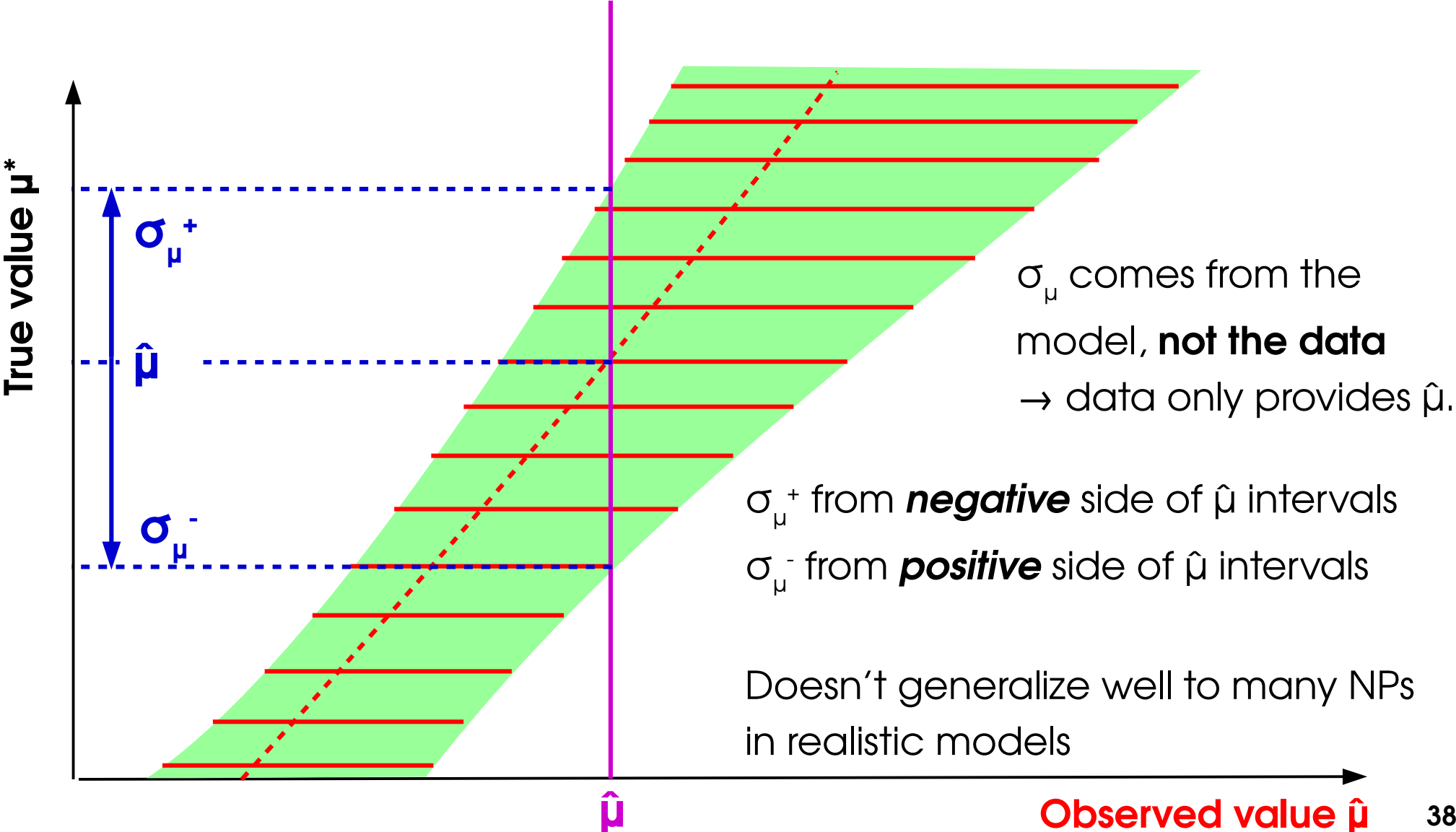
# Neyman Construction

**General case:** Build  $1\sigma$  intervals of observed values for each true value  
⇒ *Confidence belt*

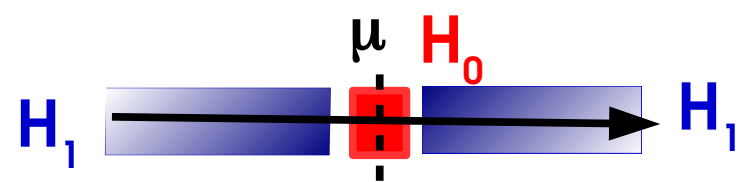


# Inversion using the Confidence Belt

**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  $P(\mu^{\text{obs}} | \mu)$  varies with  $\mu$ .



# 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

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

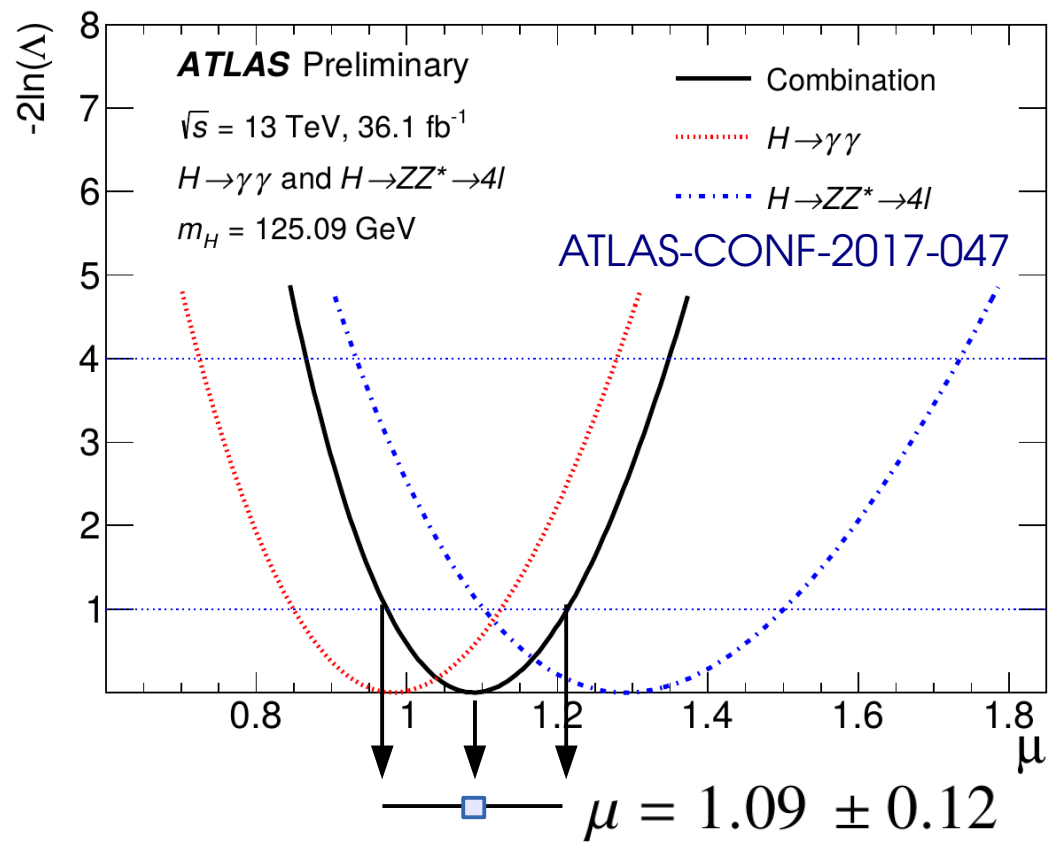
$\mu$  can be several POI!

## Asymptotics:

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

## In practice:

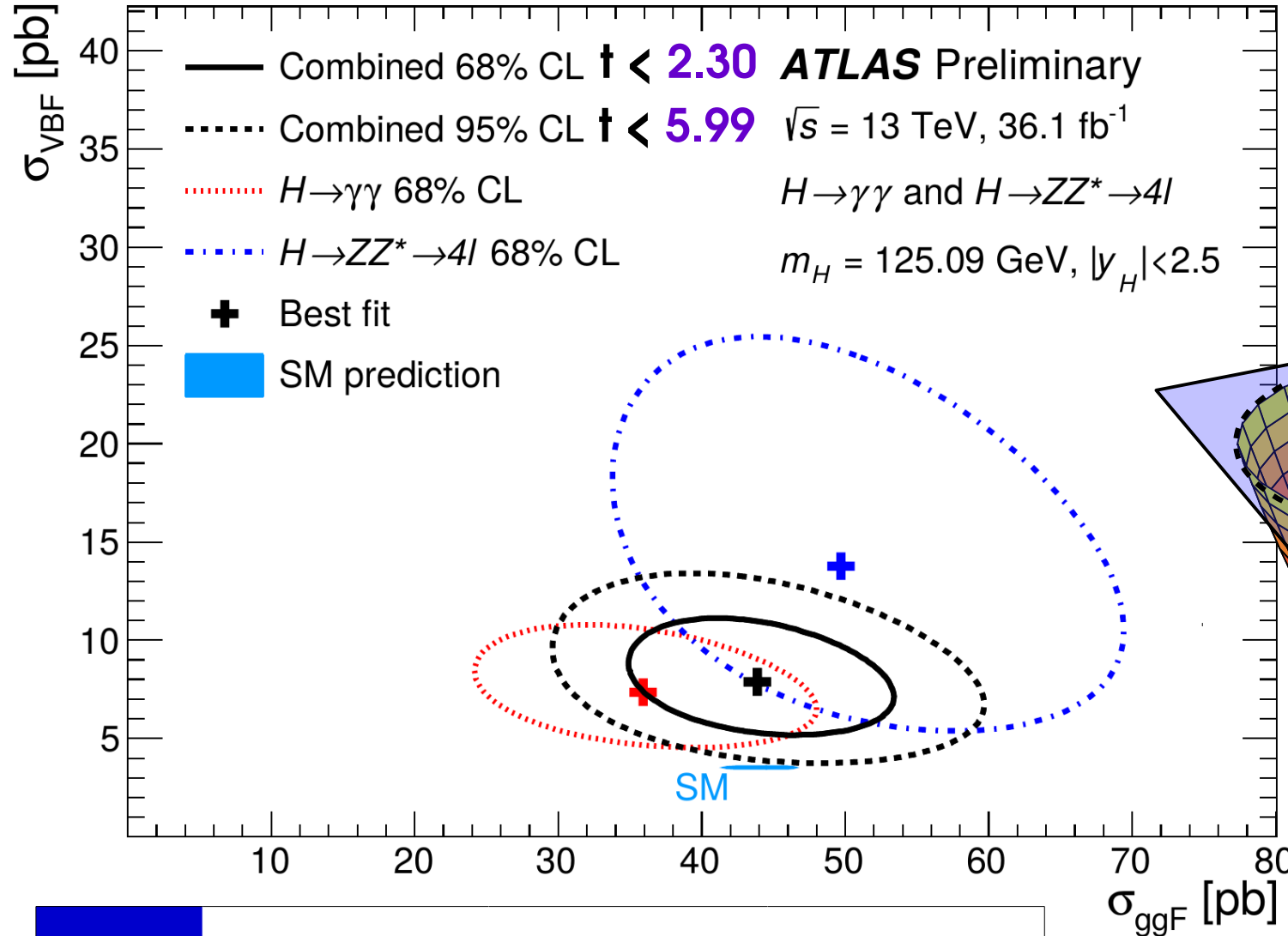
- Plot  $t_{\mu}$  vs.  $\mu$
- The minimum occurs at  $\mu = \hat{\mu}$
- Crossings with  $t_{\mu} = Z^2$  give the  $\pm Z\sigma$  uncertainties (for  $N_{POI}=1$ )



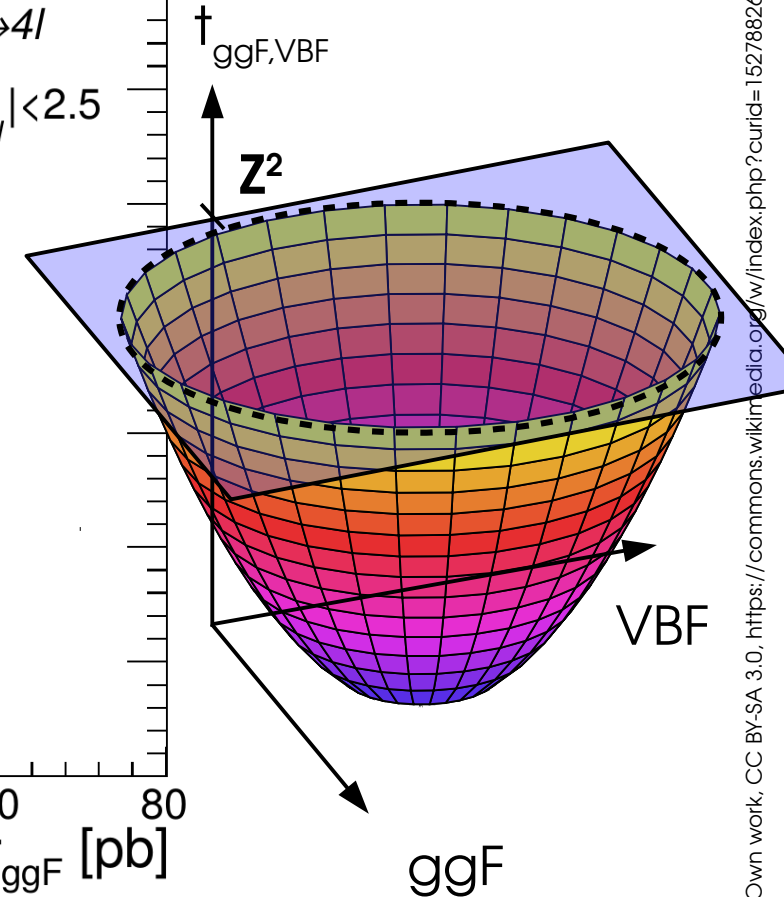
→ **Gaussian case:** parabolic profile,  $t_{\mu} = \left(\frac{\mu - \hat{\mu}}{\sigma}\right)^2 \Rightarrow \mu_{\pm} = \hat{\mu} \pm \sigma$  at  $t_{\mu} = 1$

same result as Neyman construction, also robust against non-Gaussian effects.

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



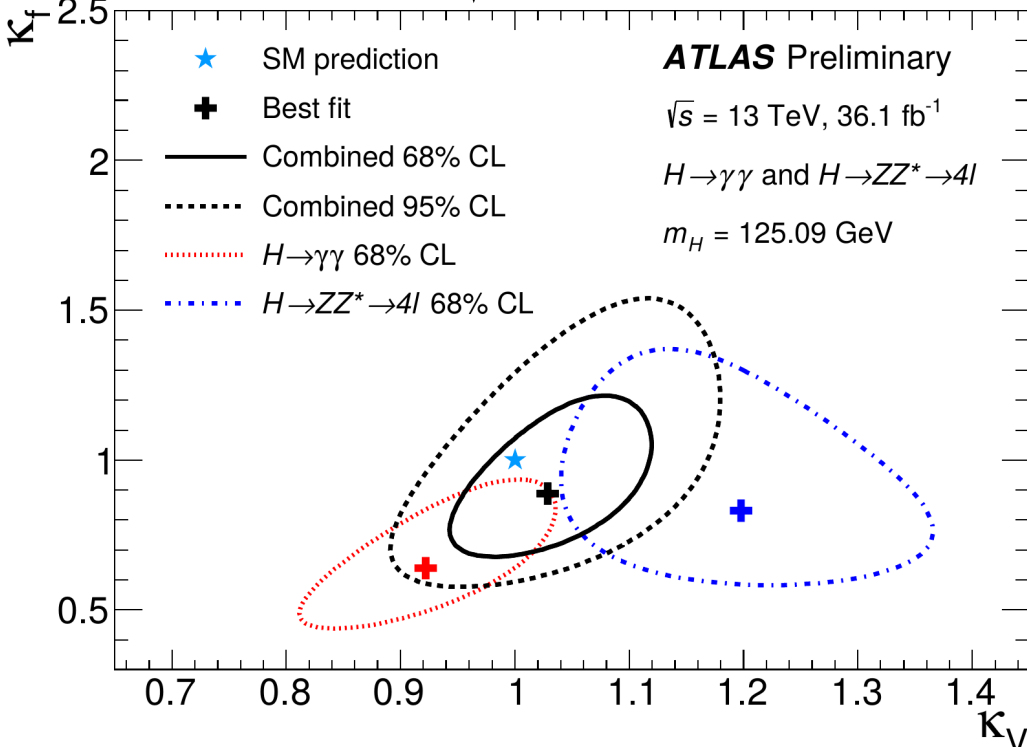
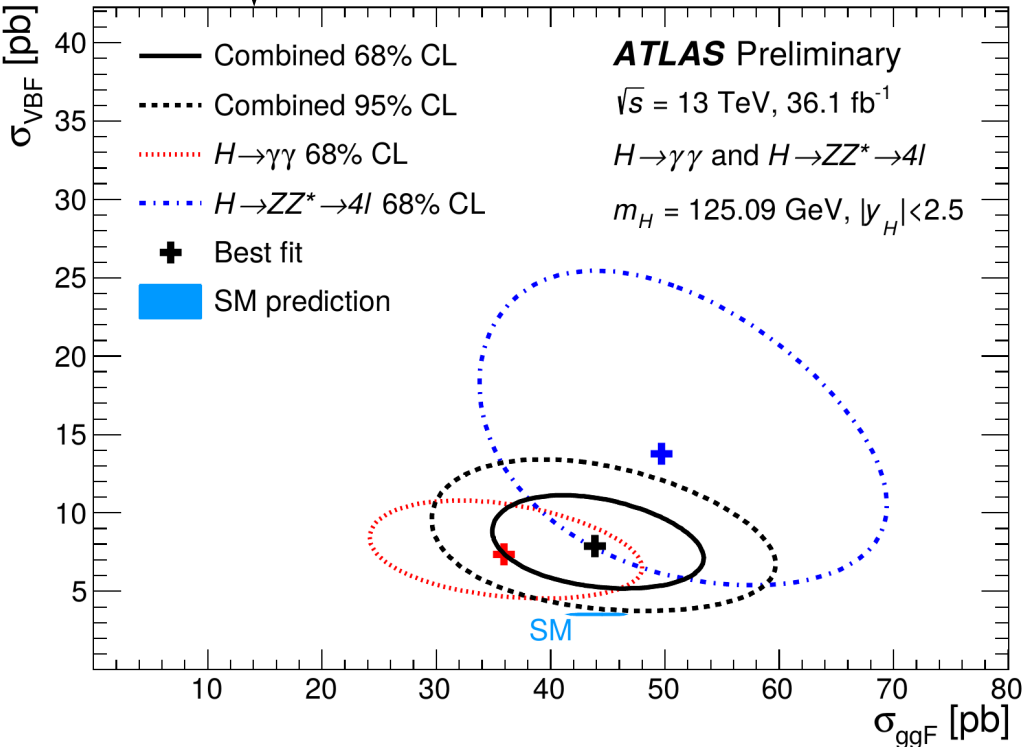
# Reparameterization

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:**

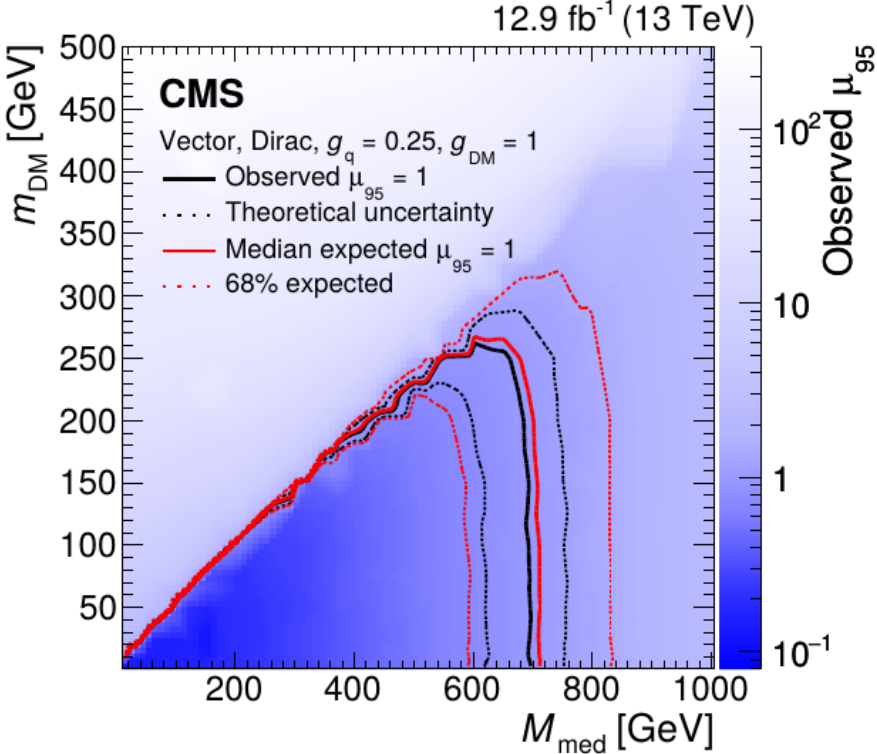
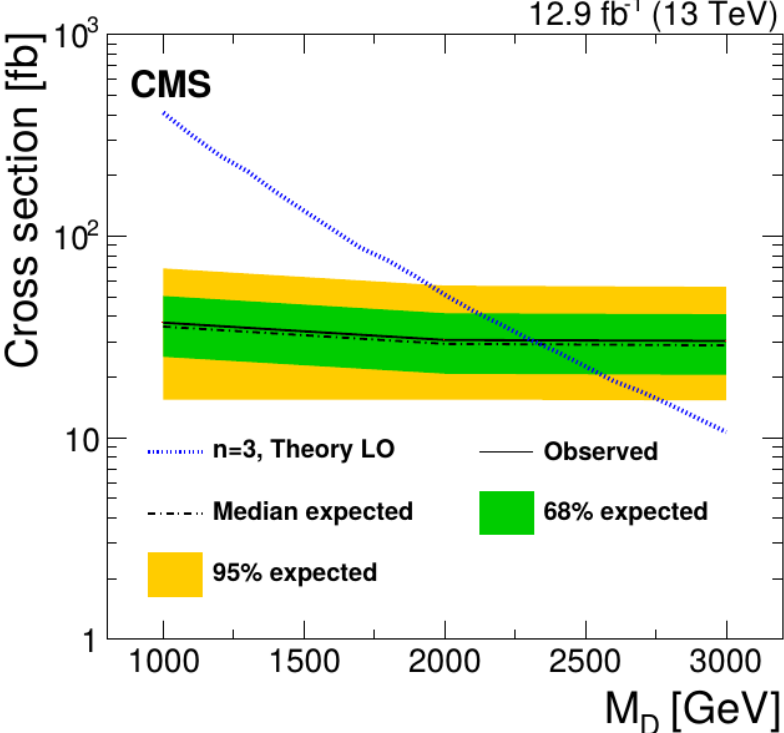
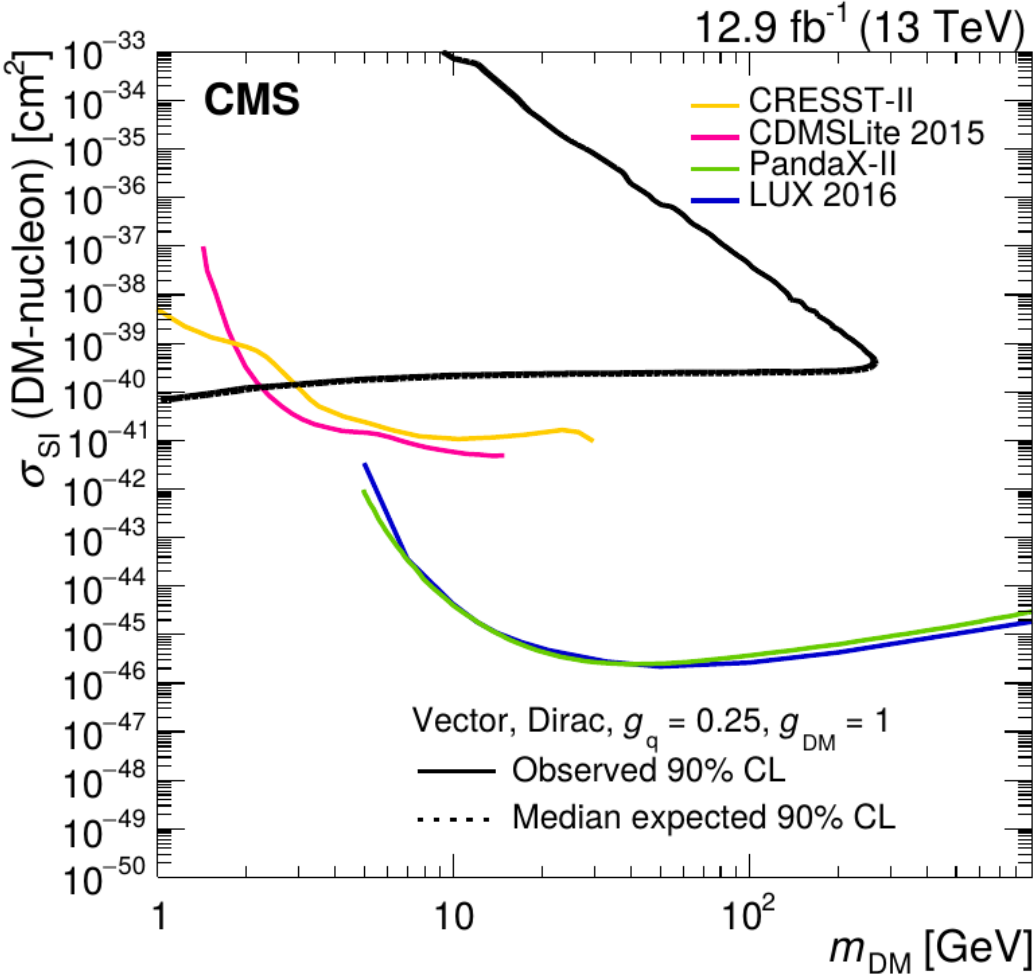
e.g. Higgs couplings:  $\sigma_{ggF}, \sigma_{VBF}$  sensitive to Higgs coupling modifiers  $\kappa_V, \kappa_F$ .

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



# Reparameterization: Limits

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



# Takeaways

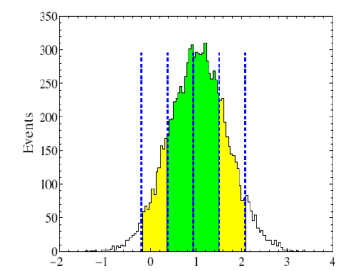
**Limits** : use LR-based test statistic:

→ Use **CL<sub>s</sub> 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**,  $n=0$  :  $S_{up} = 1.96 \sigma_{Gauss}$



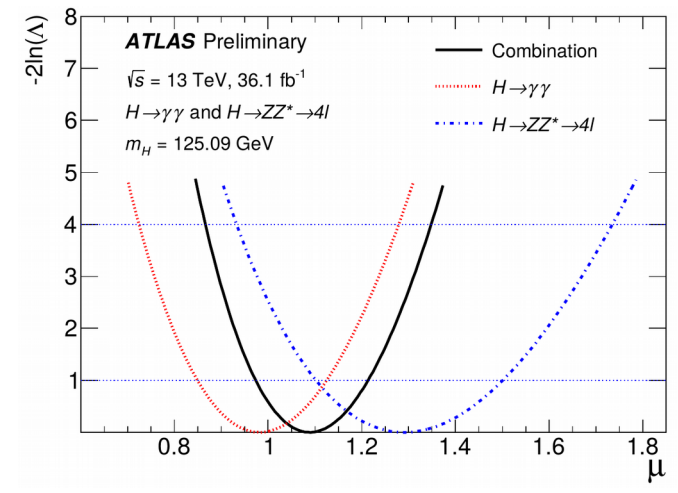
Uncertainty bands: obtain from toys or from Asimov

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

**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**:  $\mu = \hat{\mu} \pm \sigma_{Gauss}$  (1  $\sigma$  interval)

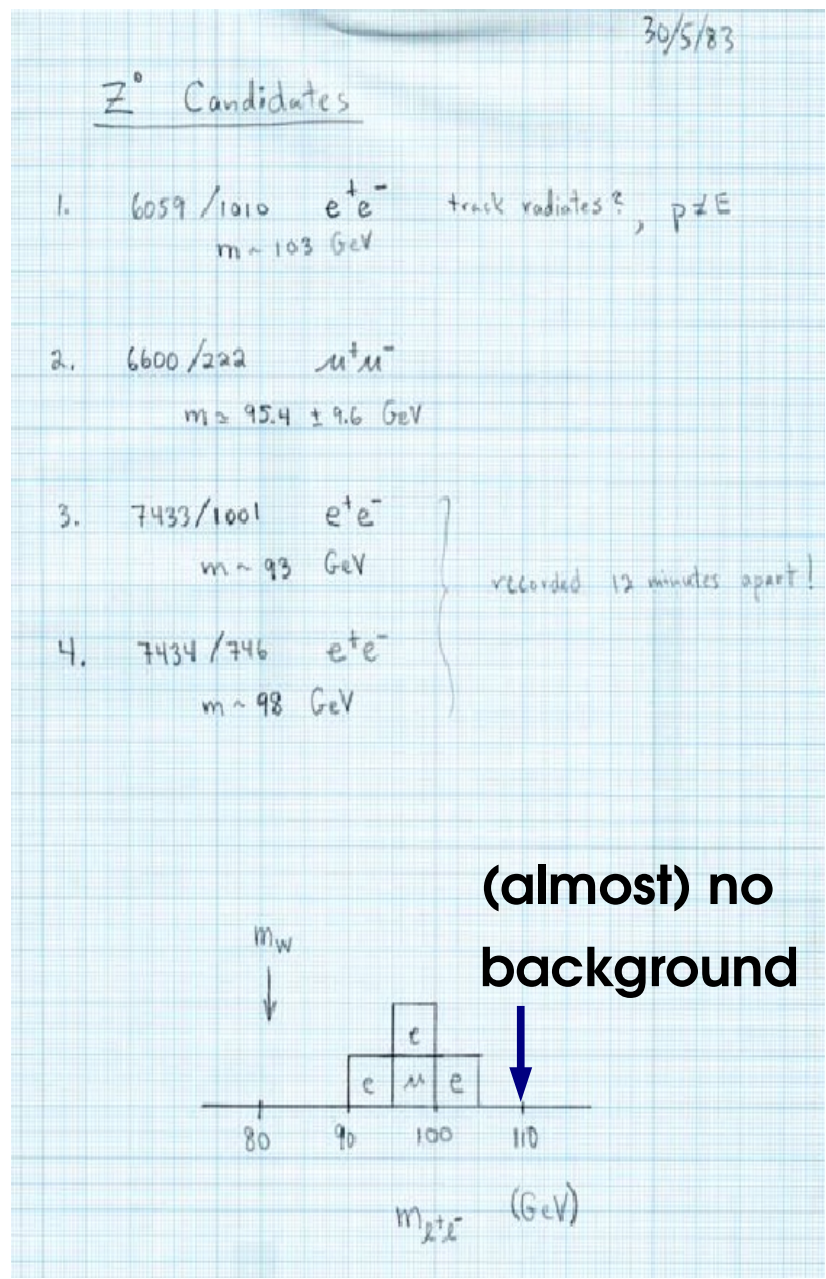
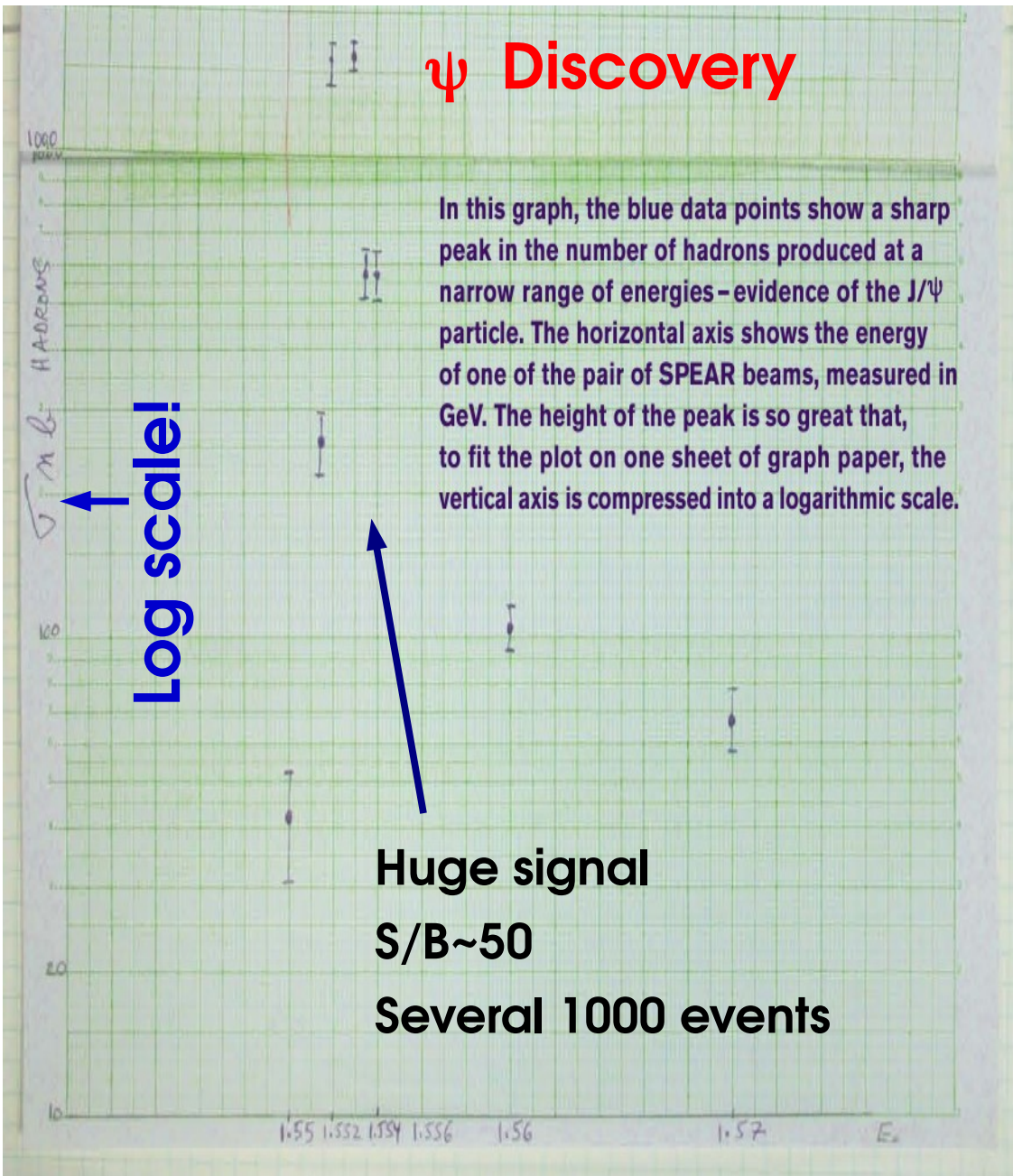


---

# Historical Aside

# Classic Discoveries (1)

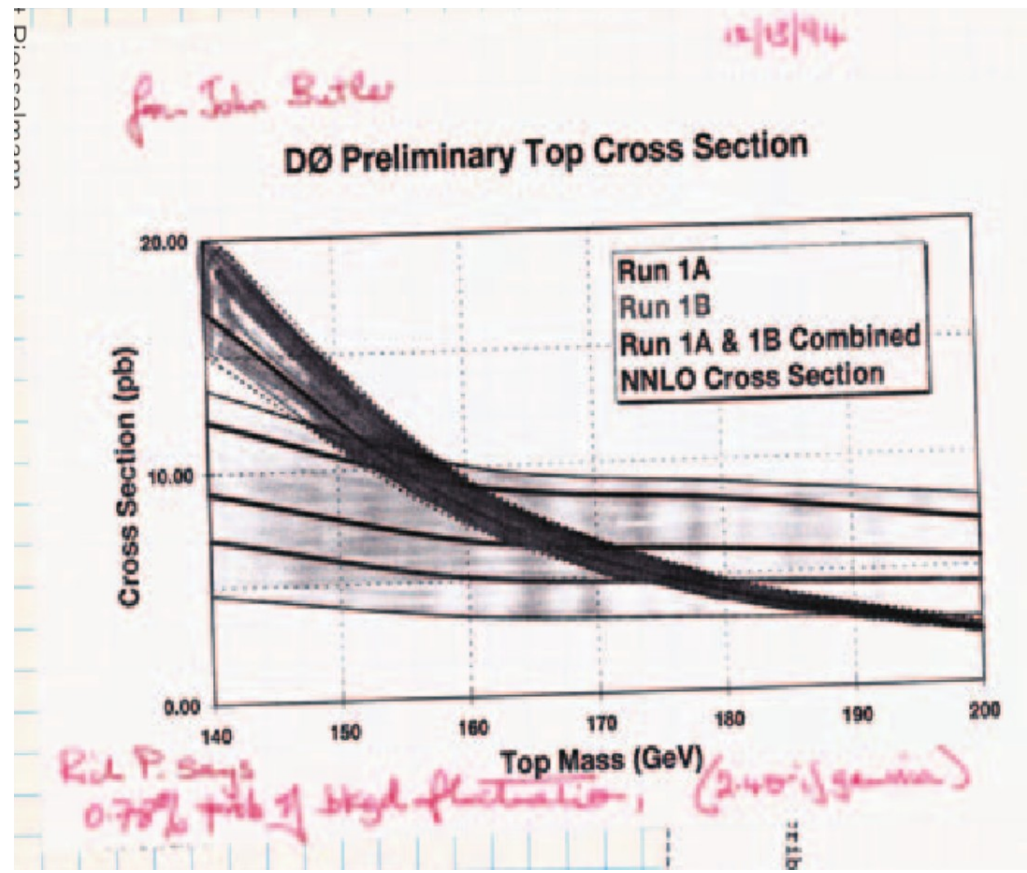
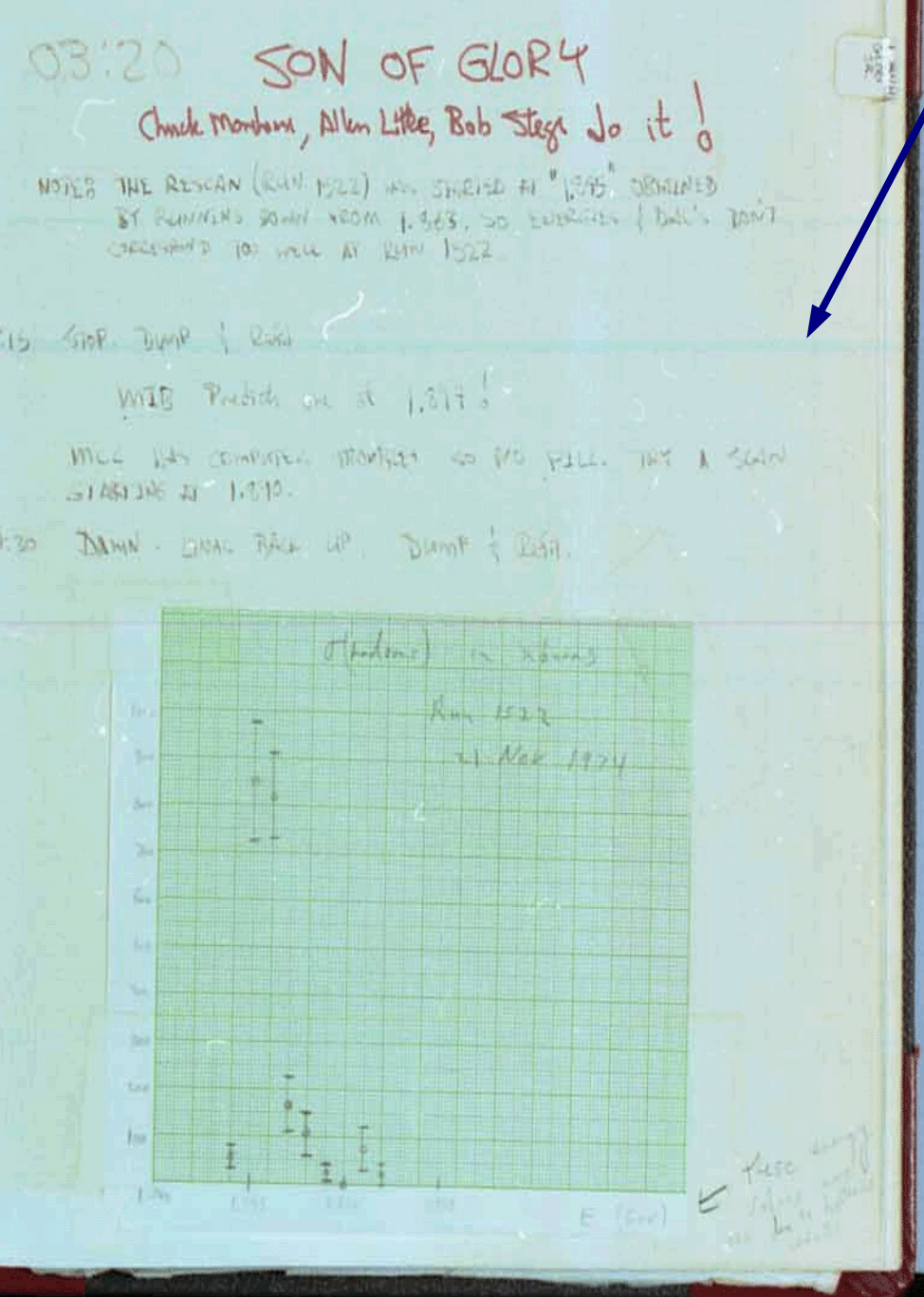
## Z<sup>0</sup> Discovery



Logbook of J. Rohlf, 1983-05-30

# Classic Discoveries (2)

$\psi'$  : discovered online  
by the (lucky) shifters



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

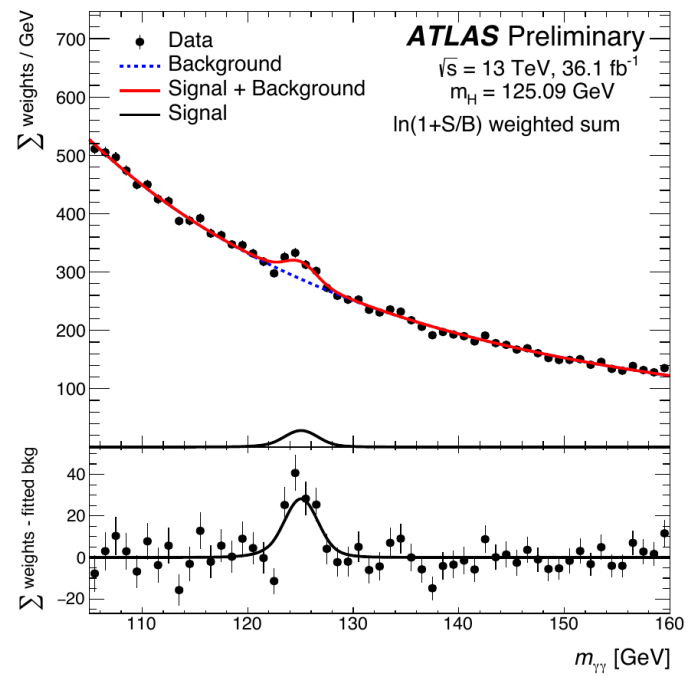
# And now ?

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

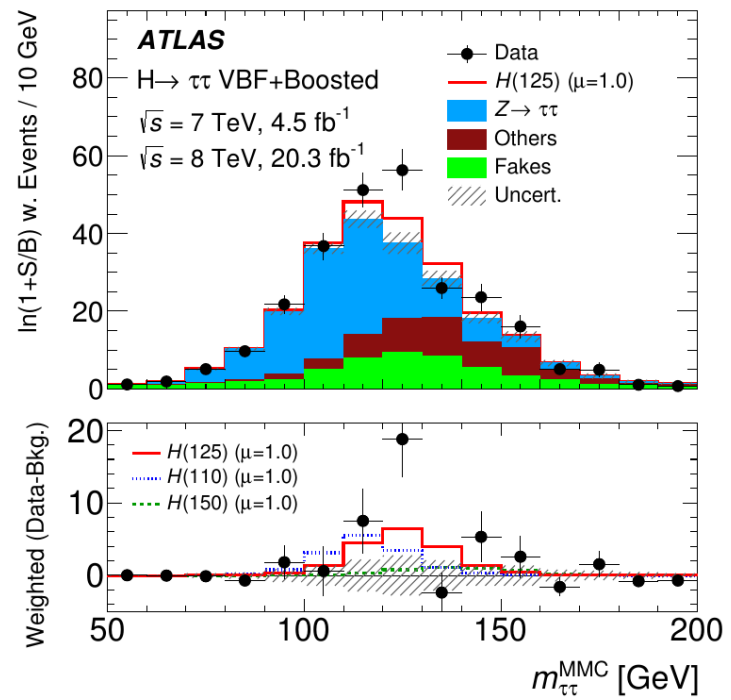
e.g. at LHC:

- **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

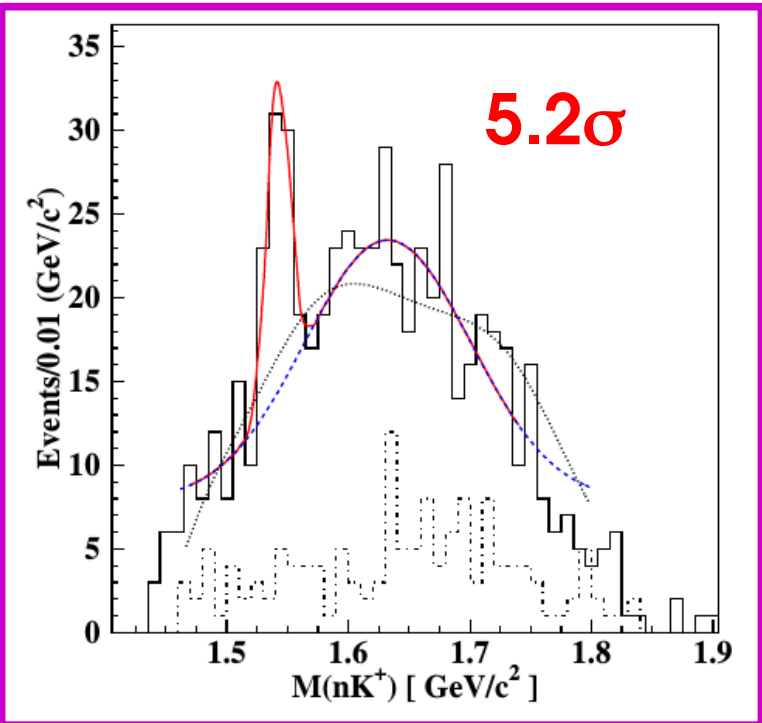
At the present time we can only speculate about the origin of this new effect. The missing transverse energy can be due either to:

- (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]<sup>14</sup> 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/c<sup>2</sup>. 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/c<sup>2</sup>.

## Pentaquarks (2003)



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

## BICEP2 B-mode Polarization (2014)

Selected for a Viewpoint in *Physics*  
 PHYSICAL REVIEW LETTERS

PRL 112, 241101 (2014) week ending  
20 JUNE 2014

---

$\mathcal{G}$

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

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

**Avoid spurious discoveries!**

→ Treatment of modeling uncertainties, systematics in general



# Outline

---

Computing statistics results:

Limits

Confidence intervals

**Profiling**

Look-Elsewhere Effect

Bayesian methods

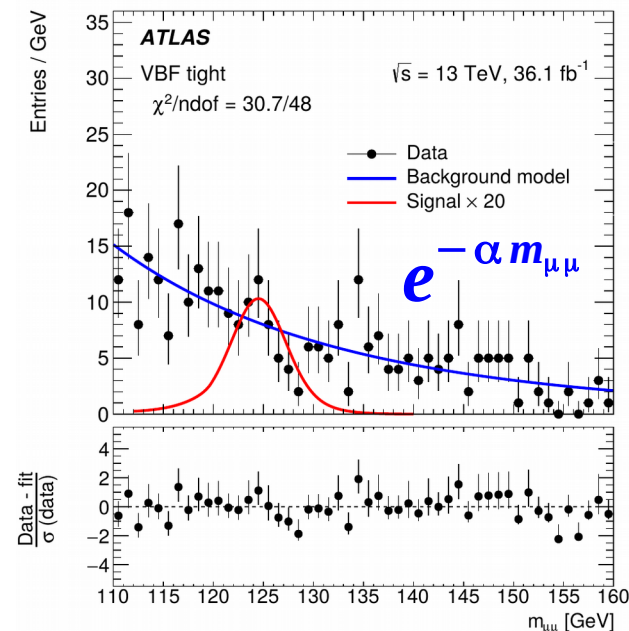
---

# Profiling

# 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 ?*

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

⇒ **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:  $L_{\text{aux}}(\boldsymbol{\theta}; \text{aux. data}) = G(\boldsymbol{\theta}^{\text{obs}}; \boldsymbol{\theta}, \boldsymbol{\sigma}_{\text{syst}})$

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

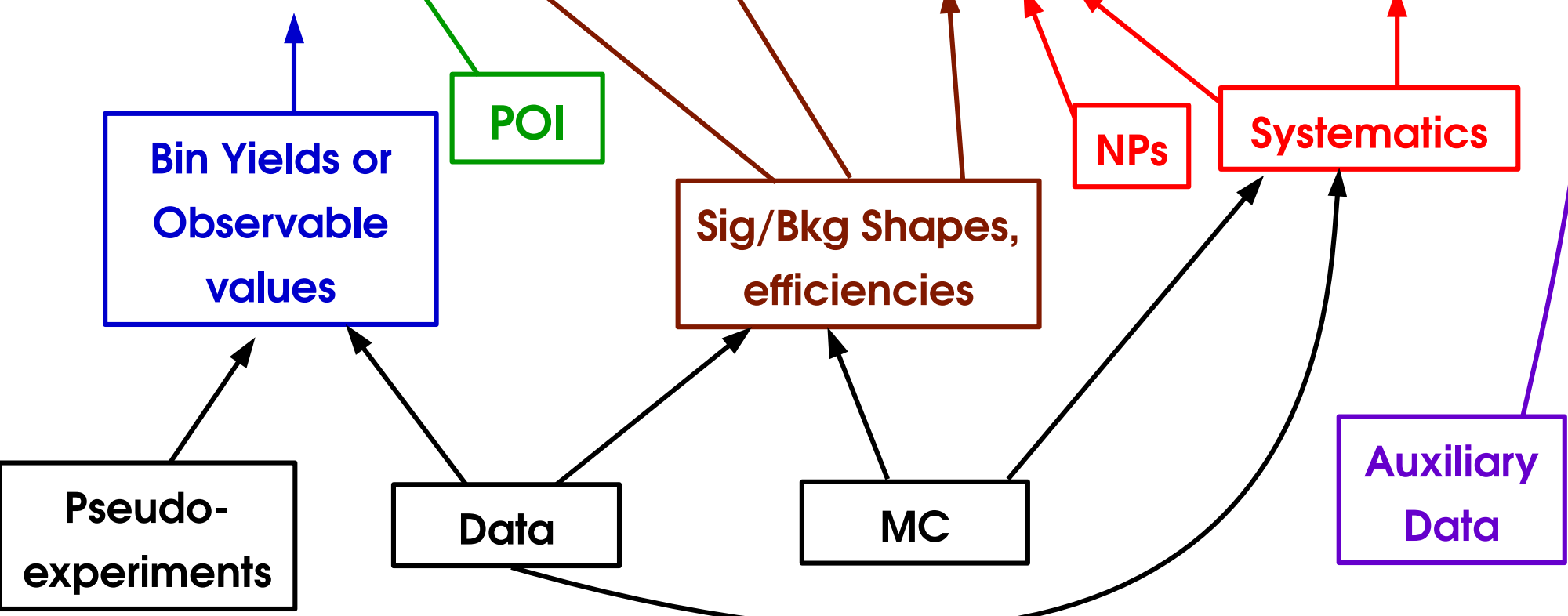
→ **Implemented in the same way nevertheless** (“pseudo-measurement”)

# 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

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



× number of categories!

# Wilks' Theorem

The likelihood usually has NPs:

- **Systematics**
- Parameters fitted in data

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

**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  $\chi^2$  !**  $f(t_{\mu_0} | \mu = \mu_0) = f_{\chi^2(n_{dof}=1)}(t_{\mu_0})$   
**also with NPs present**

→ Profiling “builds in” the effect of the NPs

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

# Effect of Profiling

Systematics still affect the result even after profiling their NPs!

e.g. **Simple counting experiment:**  $N(S, \theta) = S + \theta$ , measure  $N_{\text{obs}}$ , constraint on  $\theta$ .

**1. No NP:**  $N(S) = S$   $t_{S_0} = -2 \log \frac{L(S_0; n_{\text{obs}})}{L(\hat{S}; n_{\text{obs}})}$

→  $\hat{S}$  fit: adjust  $S$  to  $N(\hat{S}) = \hat{S} = n_{\text{obs}}$

→  $S=S_0$  fit:  $S=S_0$  fixed  $\Rightarrow N(S_0) = S_0$ , **cannot adjust**

⇒ **tension** between  $N(S_0)=S_0$  and  $S_{\text{obs}} \Rightarrow$  large  $t_{S_0} \Rightarrow$  **strong exclusion of  $H(S_0)$**

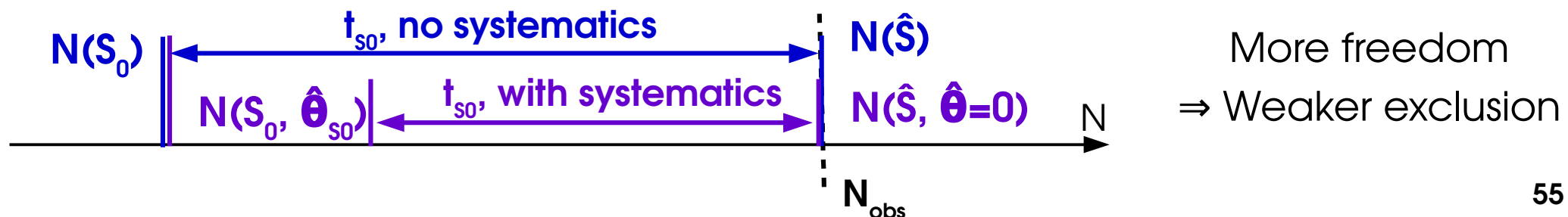
$$t_{S_0} = -2 \log \frac{L(S=S_0, \hat{\theta}_{S_0}; n_{\text{obs}})}{L(\hat{S}, \hat{\theta}; n_{\text{obs}})}$$

**2. With NP:**  $N(\mu, \theta) = S + \theta$

→  $\hat{S}$  fit adjust  $N(\hat{S}, \hat{\theta}) = N(\hat{S}, \hat{\theta}=0) = n_{\text{obs}}$  using  $S$  only (avoid penalty on  $\theta$ )

→  $S=S_0$  fit:  $S=S_0$  fixed, but  $\hat{\theta}_{S_0}$  **can still pull  $N(S_0, \hat{\theta}_{S_0})$  towards  $N_{\text{obs}}$**

⇒ smaller  $t_{S_0} \Rightarrow$  **reduced exclusion of  $H(S_0)$**



# Uncertainty decomposition

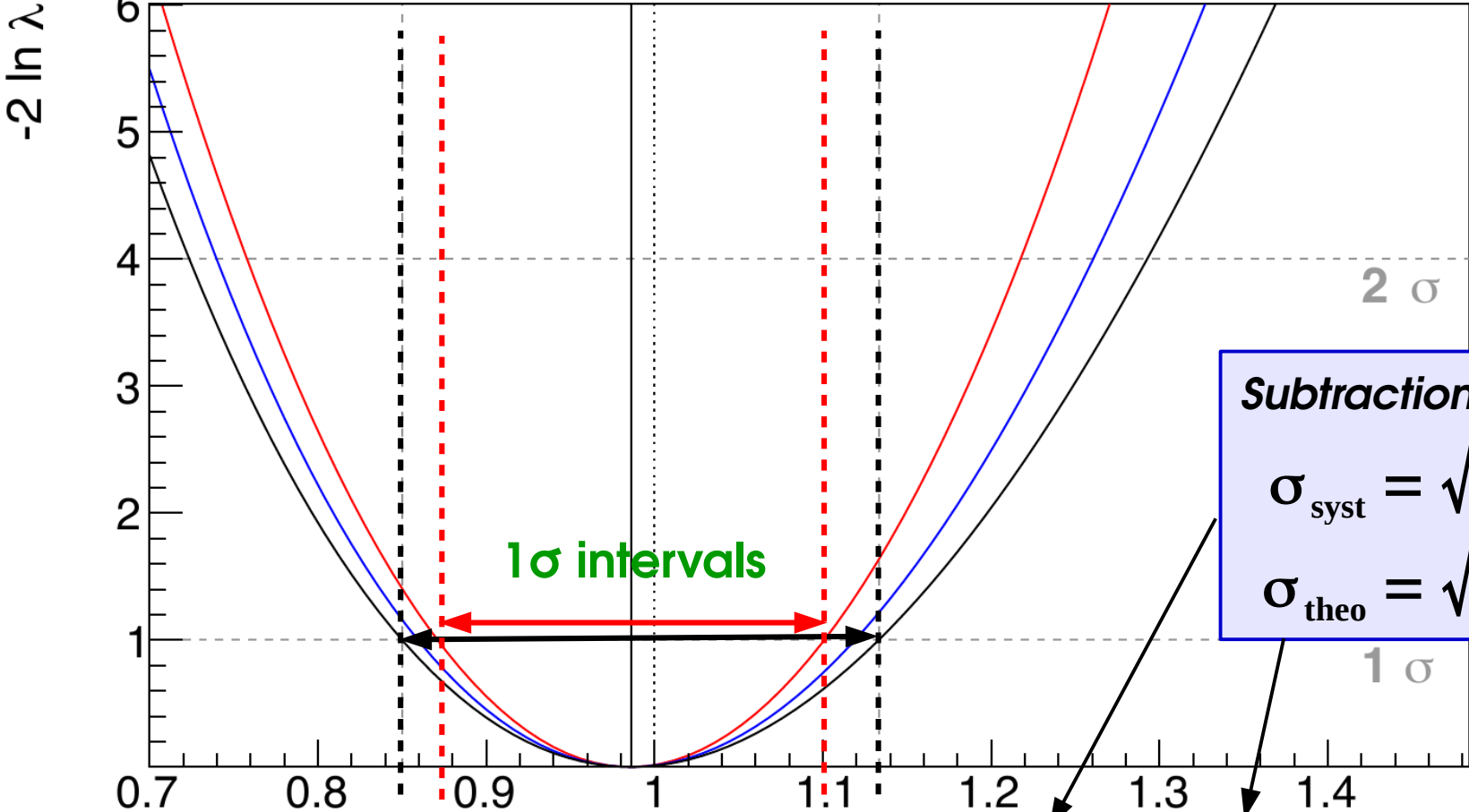
All systematics NPs fixed to 0 : statistical uncertainty only

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

**ATLAS**

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

— Total    — Theory    — Stat



**Subtraction in quadrature**

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

$$\sigma_{\text{theo}} = \sqrt{\sigma_{\text{stat+theo}}^2 - \sigma_{\text{stat}}^2}$$

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



# Gaussian Profiling

Gaussian measurement with 1 POI  $\mu$  and 1 NP  $\theta$ :

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

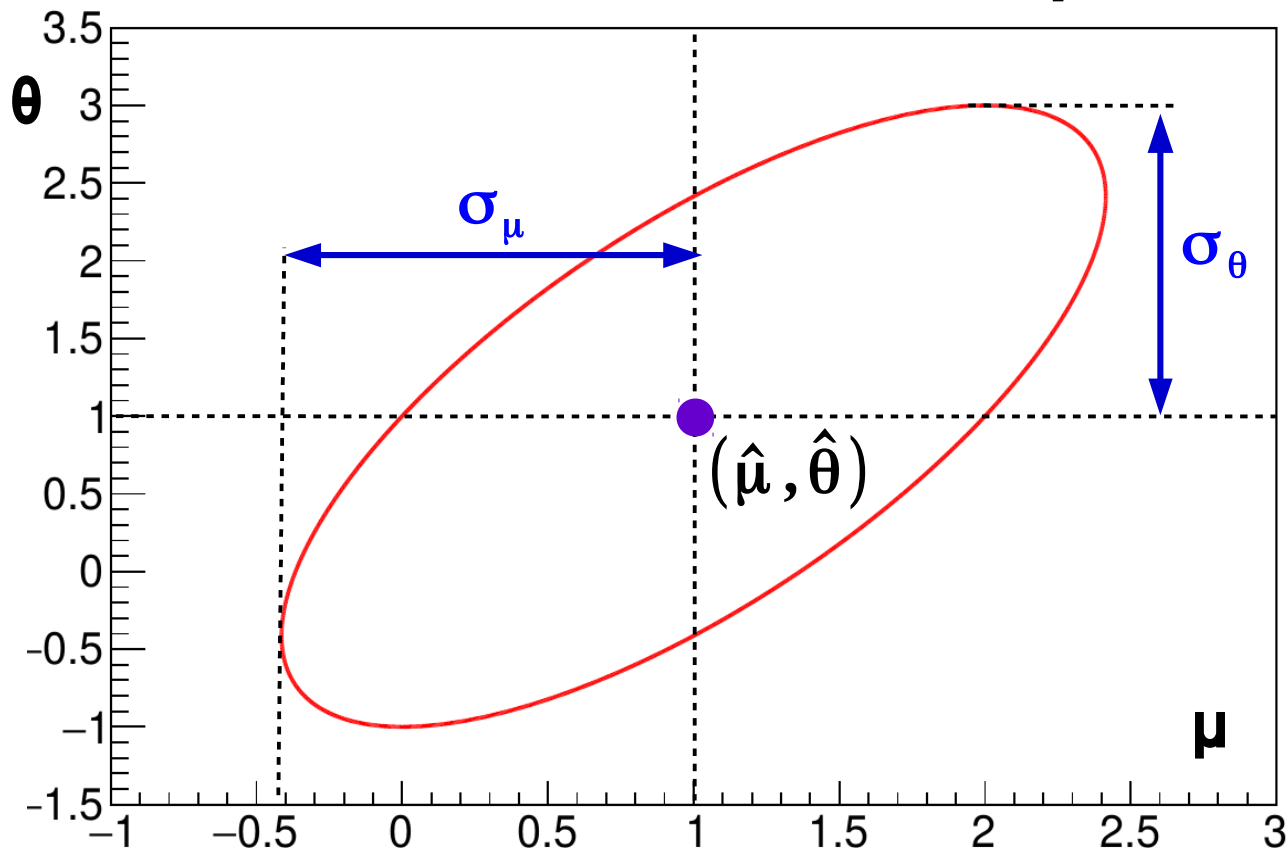
"data"

→  $\lambda(\mu, \theta)$  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 $\mu$ :

- From  $C$ , with  $\theta$  included:  $\sigma_\mu$



# 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  $\theta$  (minimize  $\lambda$  at fixed  $\mu$ ) :

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

Profile likelihood ratio:

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

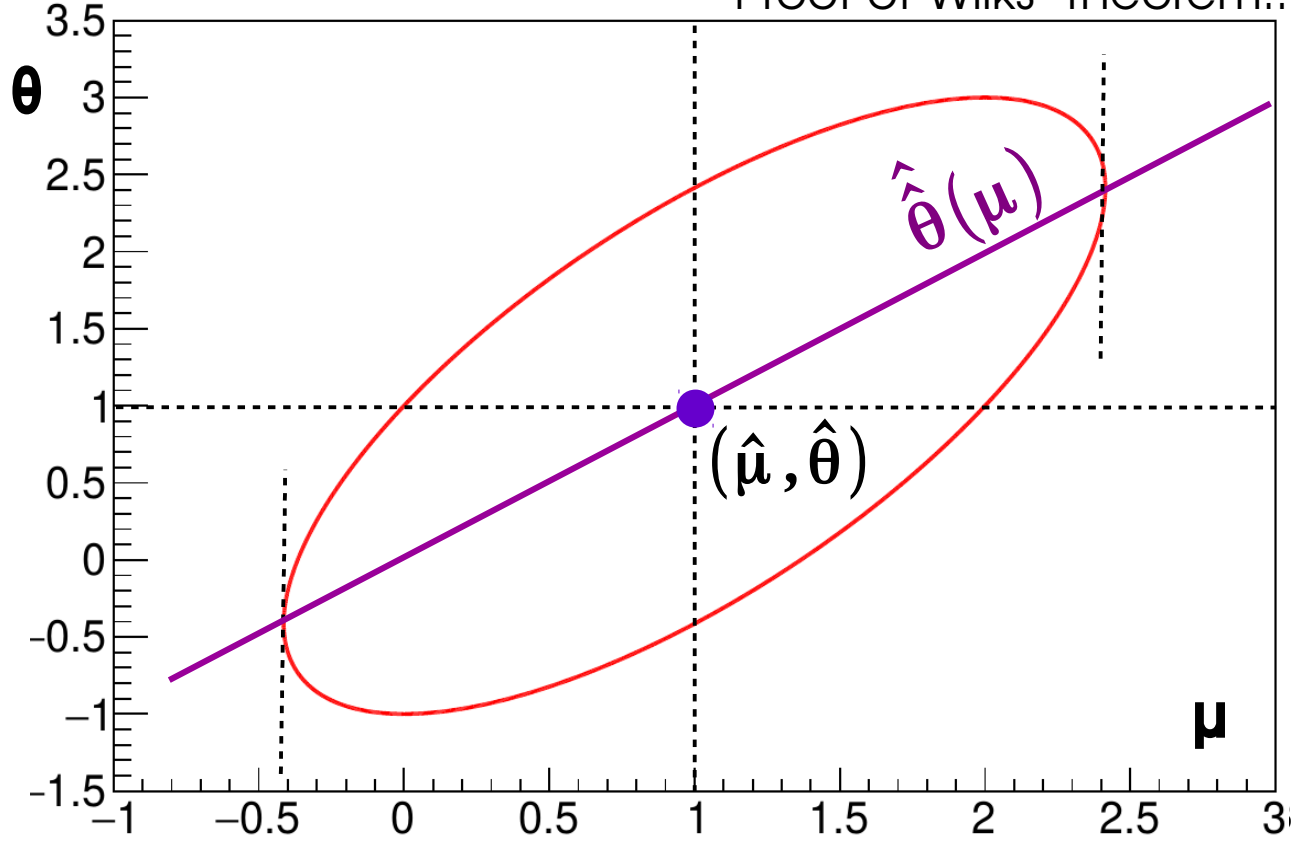
$F_{\mu\mu} \neq C_{\mu\mu}^{-1}$  !!

Proof of Wilks' theorem...

## Uncertainty on $\mu$ :

- From C:  $\sigma_\mu$
- From PLR:  $\sigma_\mu$

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



# Gaussian Profiling

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

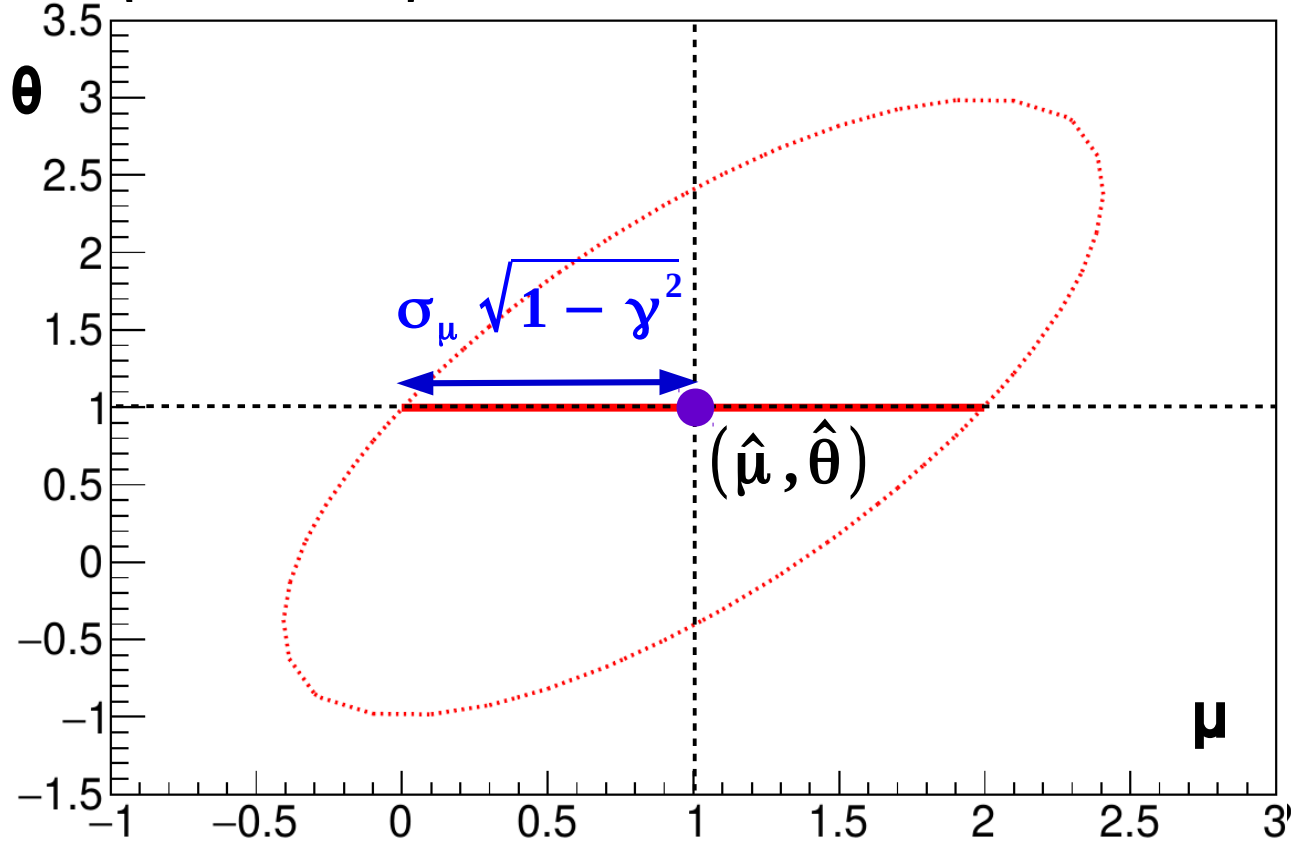
$$F \equiv C^{-1} = \frac{1}{1 - \gamma^2} \begin{bmatrix} \frac{1}{\sigma_\mu^2} & \frac{\gamma}{\sigma_\mu \sigma_\theta} \\ \frac{\gamma}{\sigma_\mu \sigma_\theta} & \frac{1}{\sigma_\theta^2} \end{bmatrix}$$

→ 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 $\mu$ :

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



# Gaussian Profiling

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

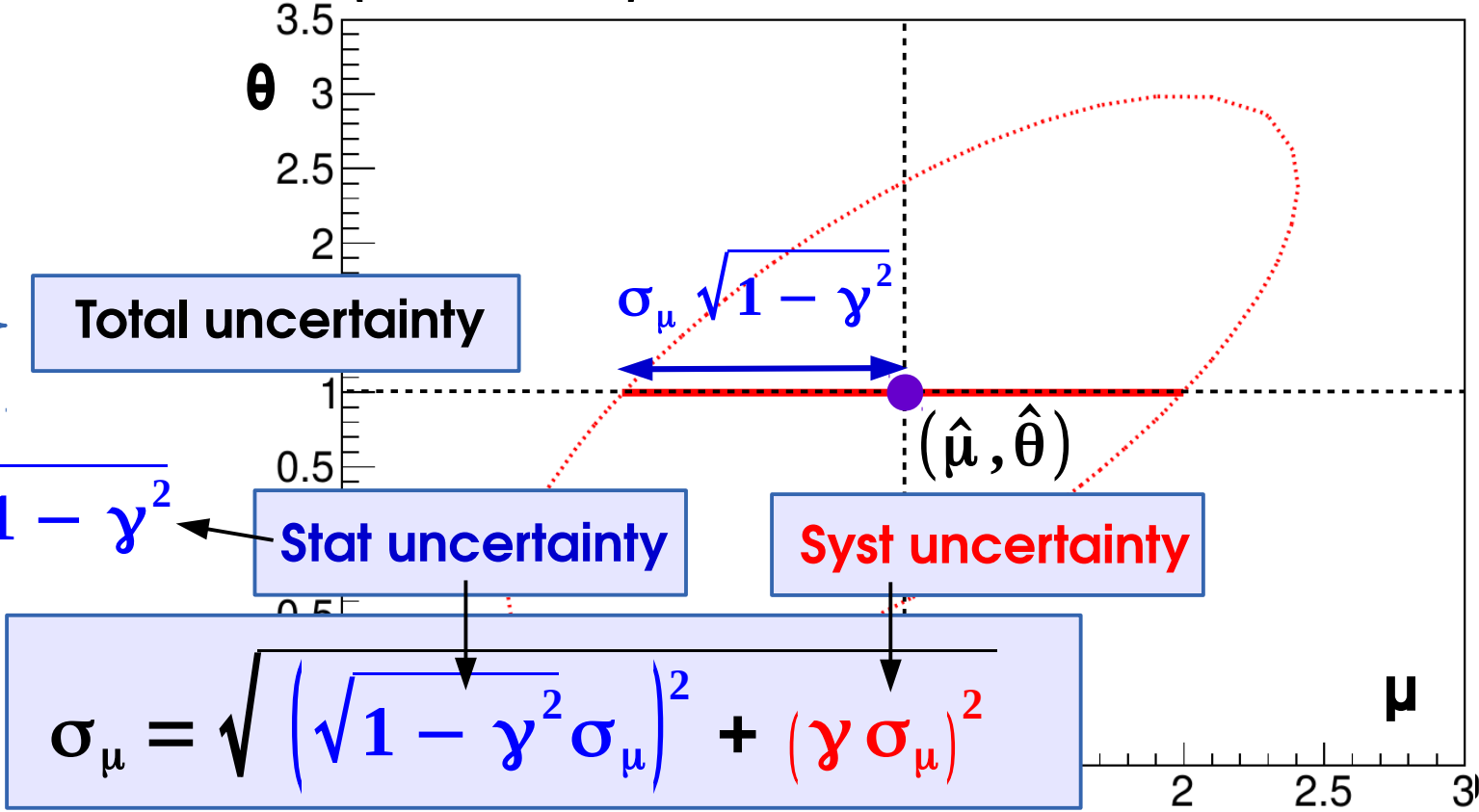
$$F \equiv C^{-1} = \frac{1}{1 - \gamma^2} \begin{bmatrix} \frac{1}{\sigma_\mu^2} & \frac{\gamma}{\sigma_\mu \sigma_\theta} \\ \frac{\gamma}{\sigma_\mu \sigma_\theta} & \frac{1}{\sigma_\theta^2} \end{bmatrix}$$

→ 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 $\mu$ :

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



# Gaussian Profiling

Back to  $\mathbf{N}(\mathbf{S}, \boldsymbol{\theta}) = \mathbf{S} + \boldsymbol{\theta}$  :

→ Measure  $\mathbf{N}_{\text{obs}} \sim \mathbf{G}(\mathbf{N}^*, \sigma_N)$

→ constraint  $\mathbf{G}(\boldsymbol{\theta}, \sigma_\theta)$  on  $\boldsymbol{\theta}$

→ everything still Gaussian:

Then:

$$\left. \begin{array}{ll} \sqrt{1-\gamma^2} \sigma_\mu = \sigma_N & \text{Stat. uncertainty} \\ \gamma \sigma_\mu = \sigma_\theta & \text{Syst. uncertainty} \end{array} \right\} \sigma_\mu = \sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}$$

⇒ Stat uncertainty (on N) and syst (on  $\boldsymbol{\theta}$ ) add in quadrature as expected

## Executive summary:

- **Systematic = NP with an external constraint** (auxiliary measurement)
- Profiling systematics includes their effect into the total uncertainty, as desired
- No special treatment for systematics: treated like any other NP, automatically accounted for through profiling.
- Guaranteed to work only as long as everything is Gaussian, but typically robust against non-Gaussian behavior.

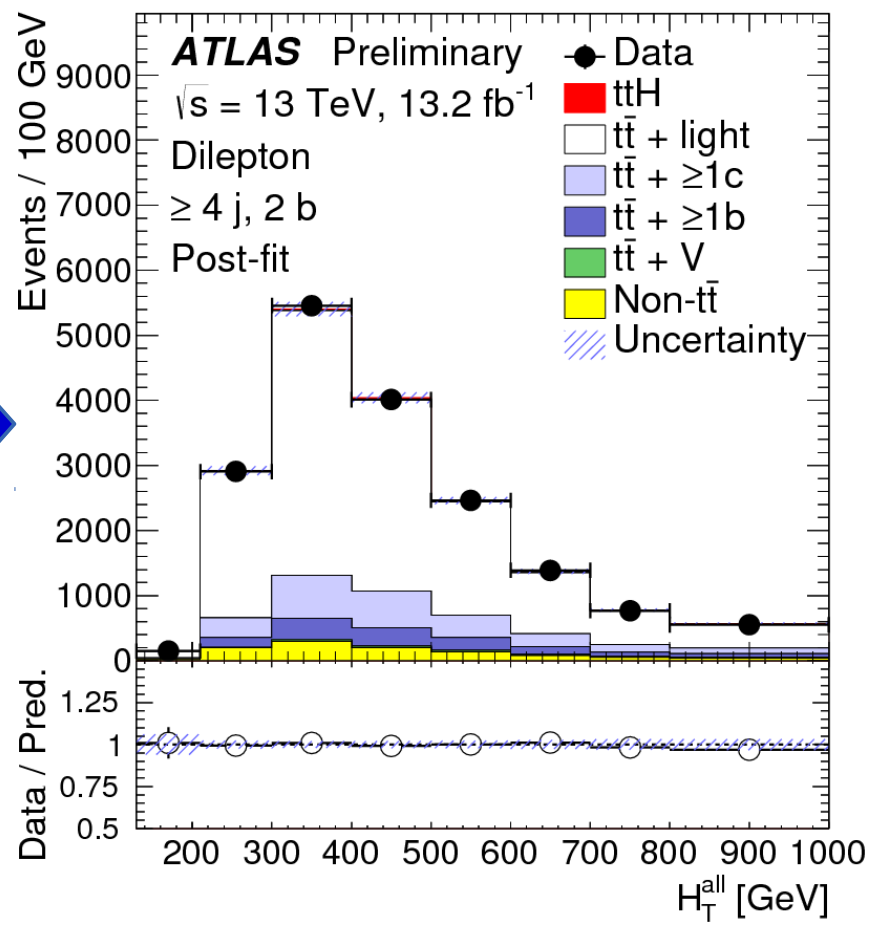
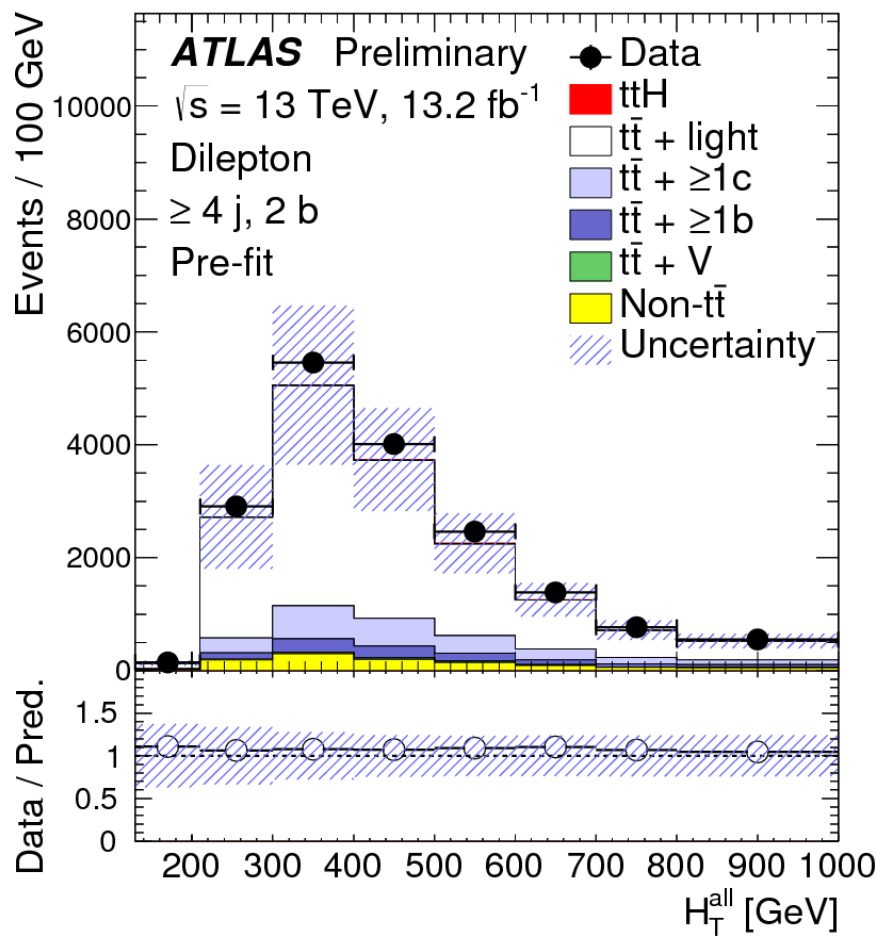
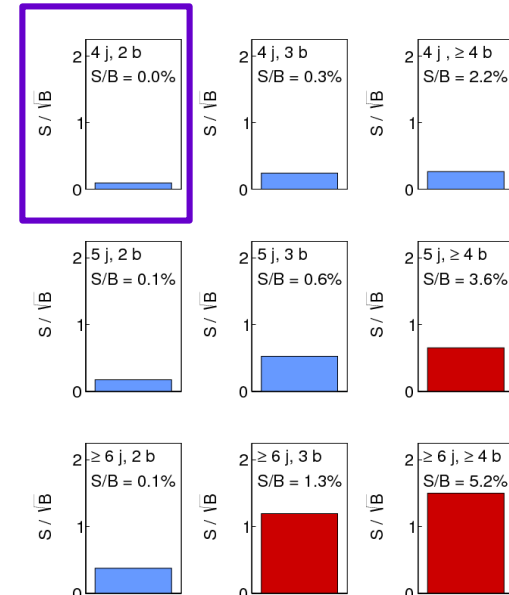
# Profiling Example: $t\bar{t}H \rightarrow bb$

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-2016-080

# Uncertainty decomposition

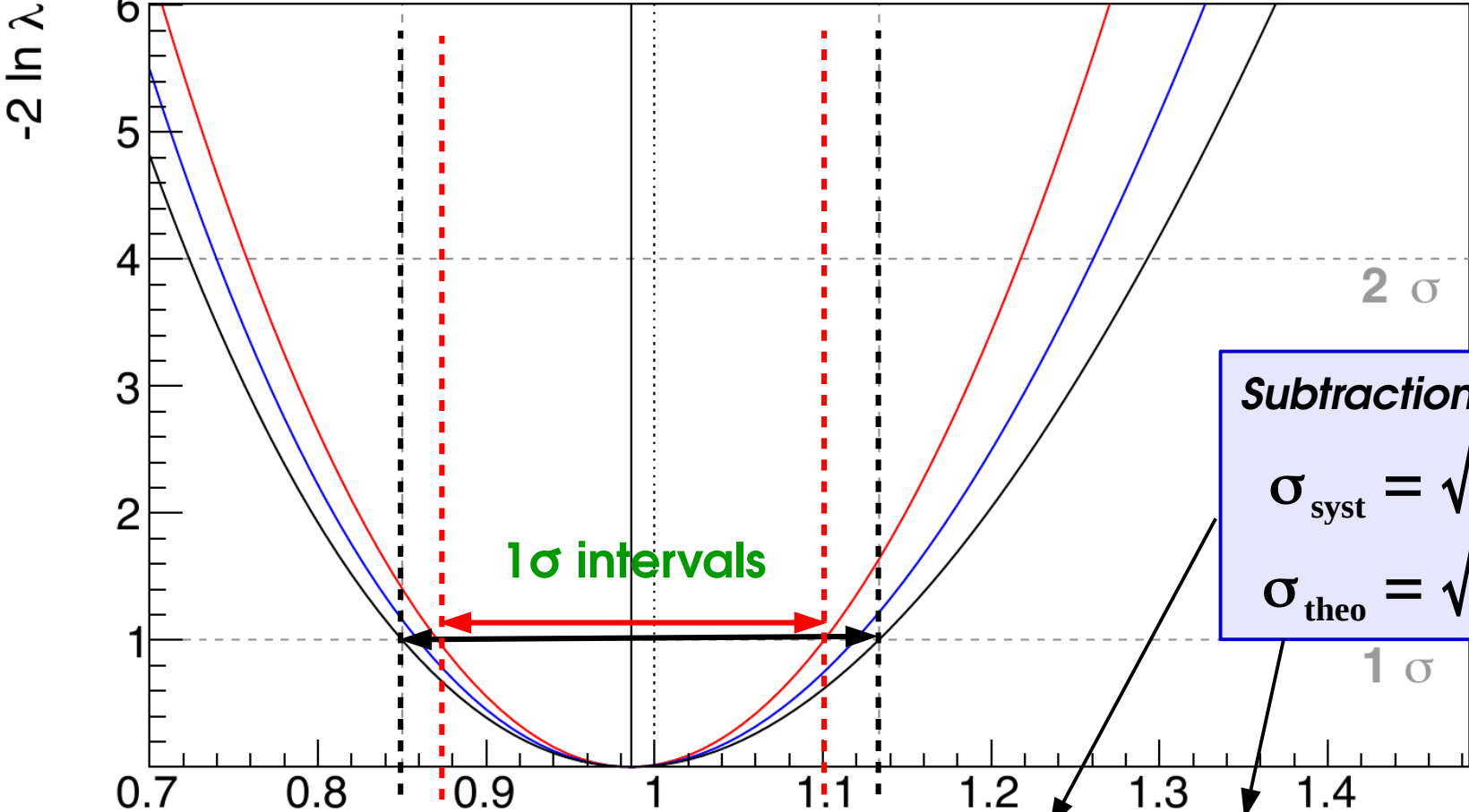
All systematics NPs fixed to 0 : statistical uncertainty only

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

**ATLAS**

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

— Total    — Theory    — Stat



**Subtraction in quadrature**

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

$$\sigma_{\text{theo}} = \sqrt{\sigma_{\text{stat+theo}}^2 - \sigma_{\text{stat}}^2}$$

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

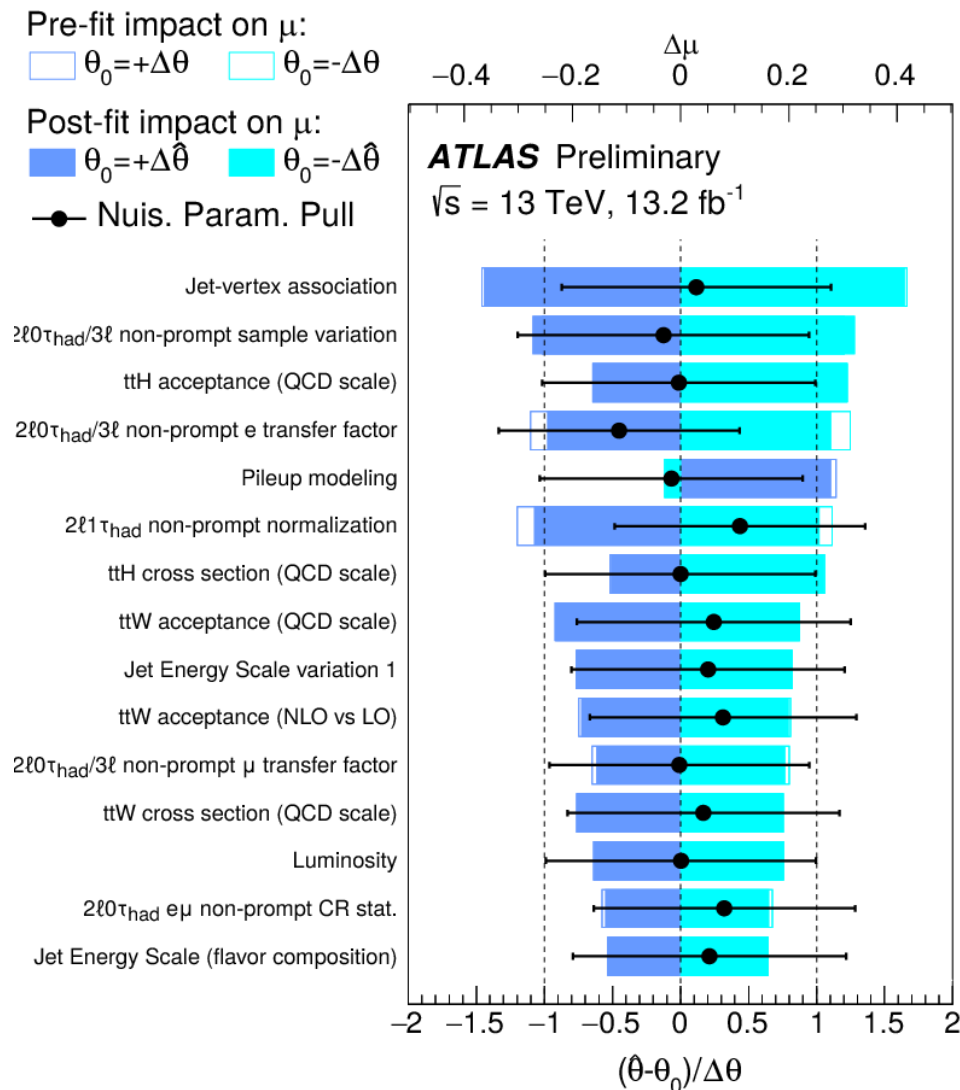
Systematics are described by NPs included in the fit. Nominally:

- **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:

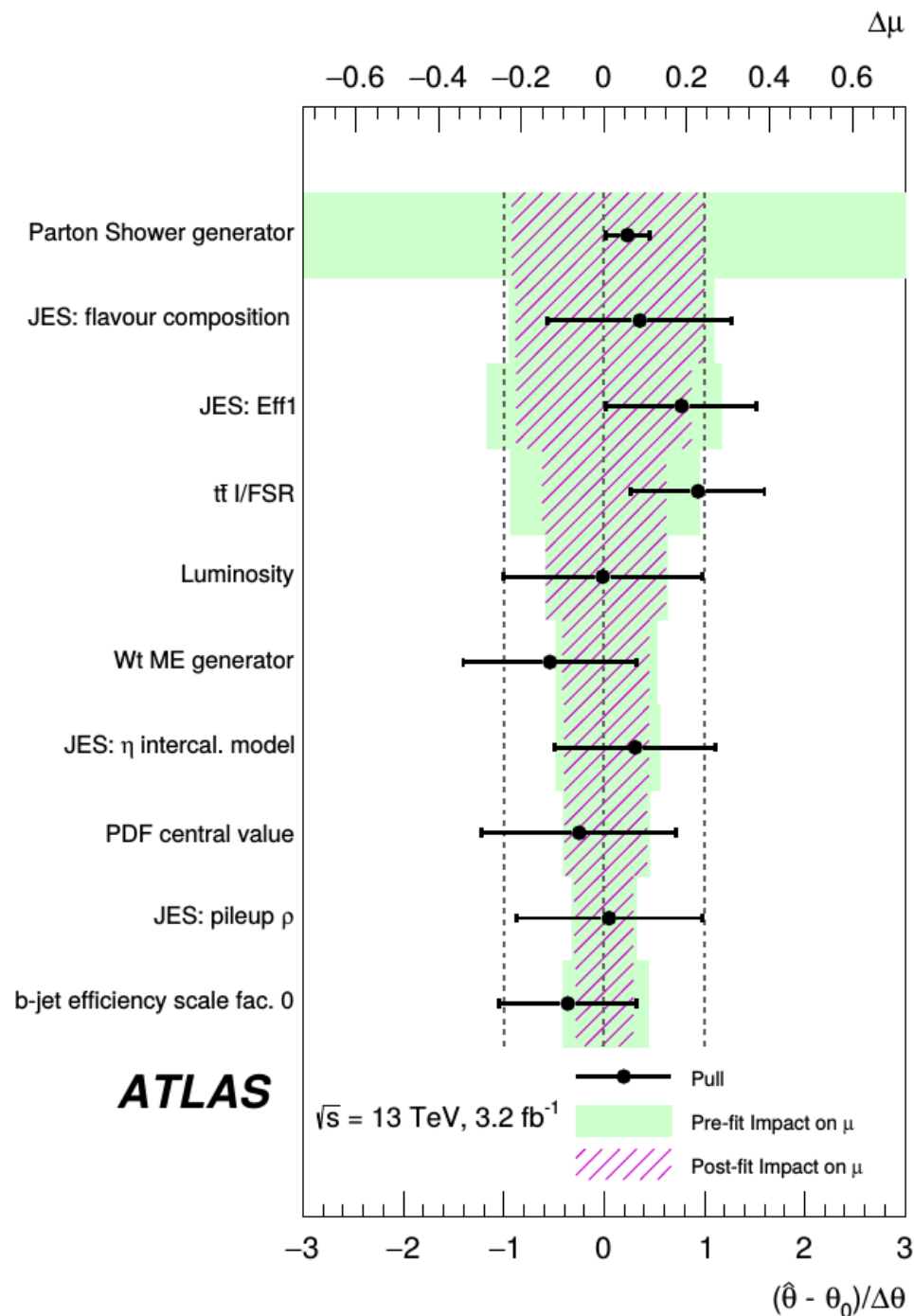
- **NP central value = 0** : corresponds to the pre-fit expectation (usually MC)
- **NP uncertainty = 1** : since NPs normalized to the value of the syst. :

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

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)



# Takeaways

**Systematics:** uncertainties on the **form of the statistical model**

(as opposed to the uncertainties encoded in the model itself)

→ Implemented using additional nuisance parameters in the model

→ Constrained by adding **auxiliary measurements** (sometimes fictitious ones) to the model – usually represented by a single Gaussian for each NP.

$$L(\boldsymbol{\mu}, \boldsymbol{\theta}; \text{data}) = L_{\text{main}}(\boldsymbol{\mu}, \boldsymbol{\theta}; \text{main data}) G(\boldsymbol{\theta}^{\text{obs}}, \boldsymbol{\theta}, 1)$$

⇒ **Systematics treated in the same way as statistical uncertainties**, although we still keep track of **systematics NPs** for bookkeeping purposes

**Profiling:** when testing a hypothesis, use the best-fit values of the nuisance parameters: **profile likelihood ratio**.

$$\frac{L(\boldsymbol{\mu} = \boldsymbol{\mu}_0, \hat{\boldsymbol{\theta}}_{\boldsymbol{\mu}_0})}{L(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\theta}})}$$

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

→ NPs still show up in the PLR as increased uncertainties – Gaussian case:

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

**Profiling can have unintended effects – need to carefully check behavior**

# 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

# Comparison with LEP/TeVatron definitions

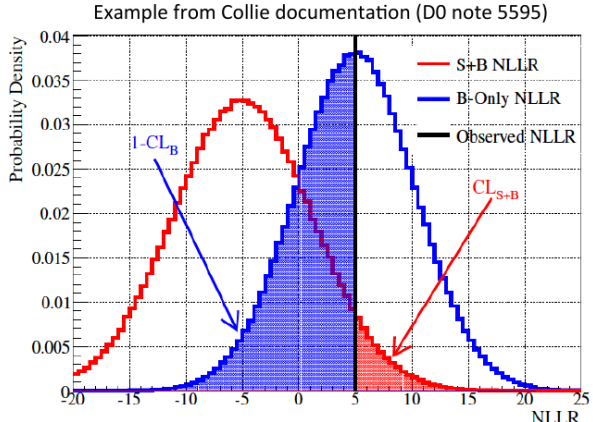
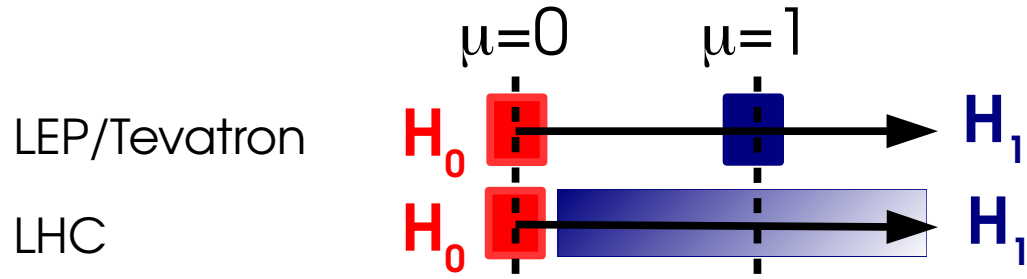
Likelihood ratios are not a new idea:

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

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

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

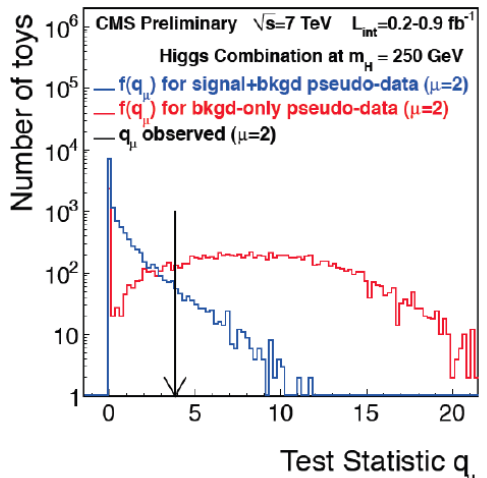
Both compare to  $\mu=1$  instead of best-fit  $\hat{\mu}$



→ Asymptotically:

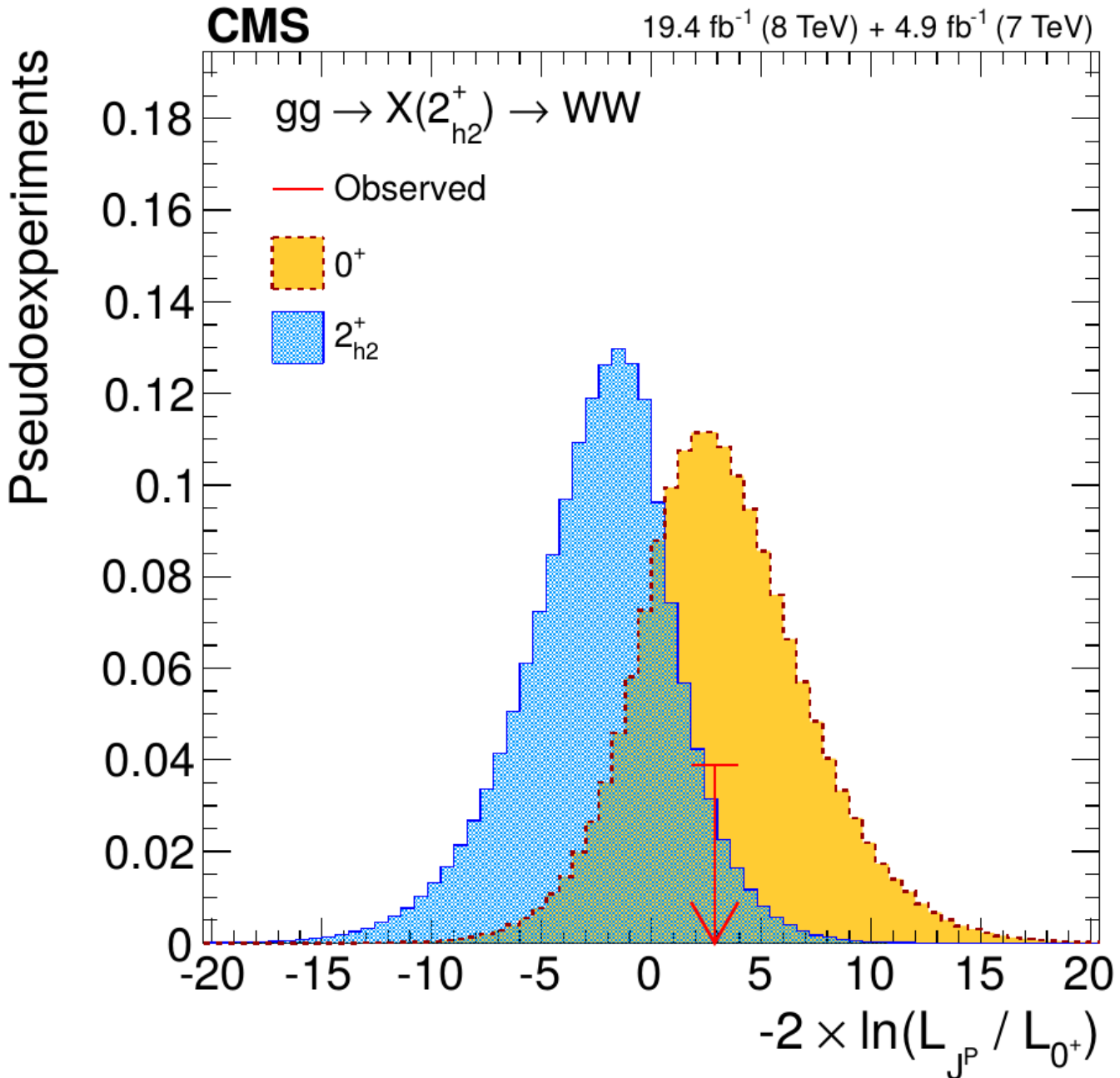
- **LEP/TeVatron**: q linear in  $\mu \Rightarrow \sim \text{Gaussian}$
- **LHC**: q quadratic in  $\mu \Rightarrow \sim \chi^2$

→ Still use TeVatron-style for discrete cases



# Spin/Parity Measurements

Phys. Rev. D 92 (2015) 012004



# 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

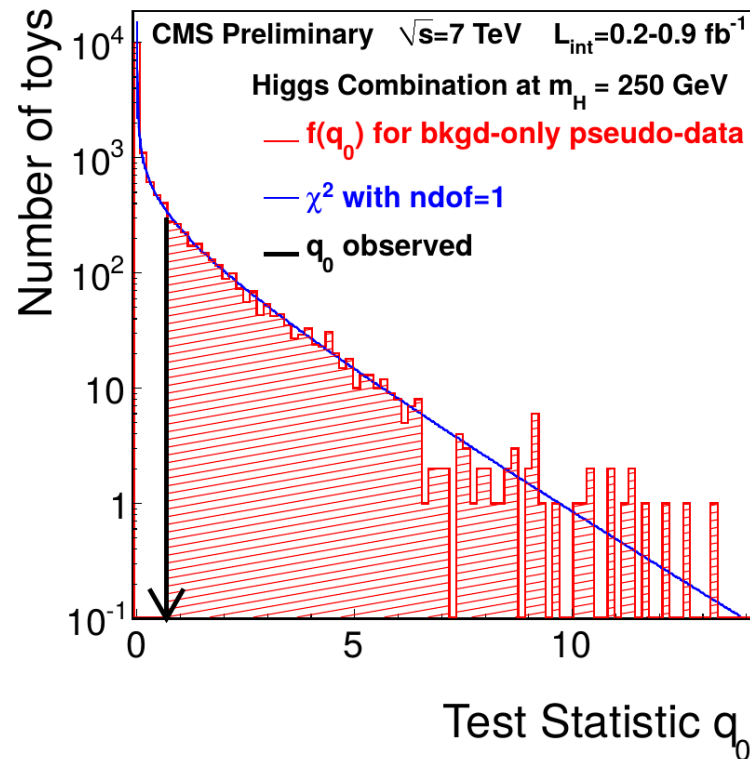
( $\theta^{obs}$ ) of each auxiliary experiment:  $G(\theta^{obs}; \theta, \sigma_{syst})$

→ Samples the true distribution of the PLR

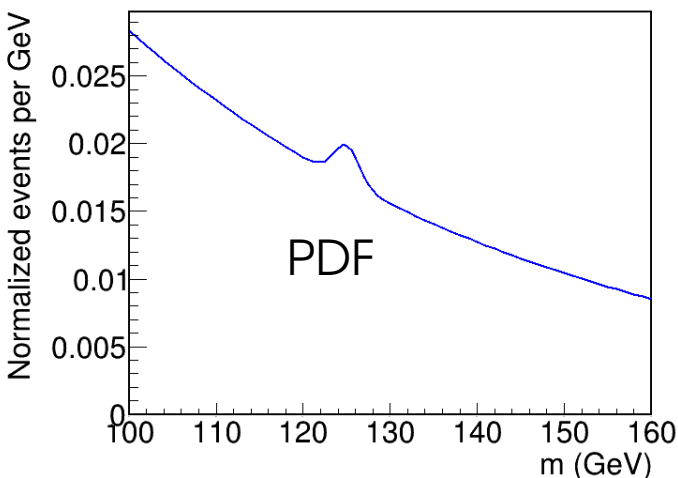
⇒ Integrate above observed PLR to get the p-value

→ Precision limited by number of generated toys,

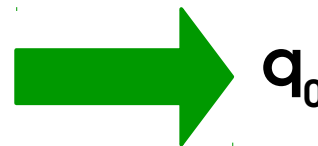
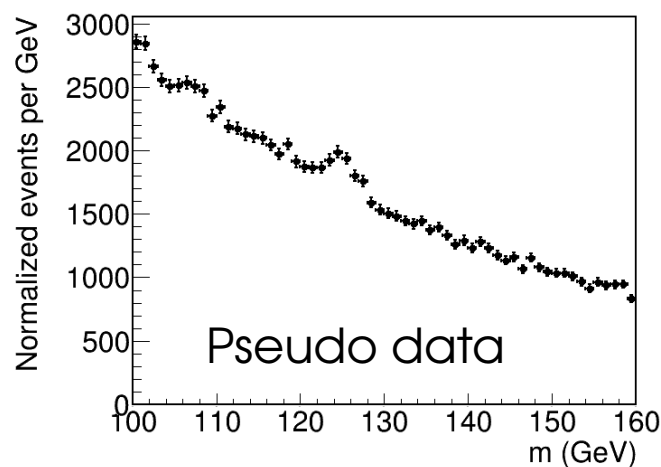
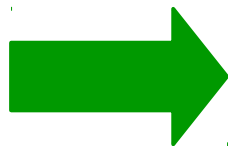
**Small p-values** ( $5\sigma : p \sim 10^{-7}!$ ) ⇒ **large toy samples**



Repeat  $N_{toys}$  times



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

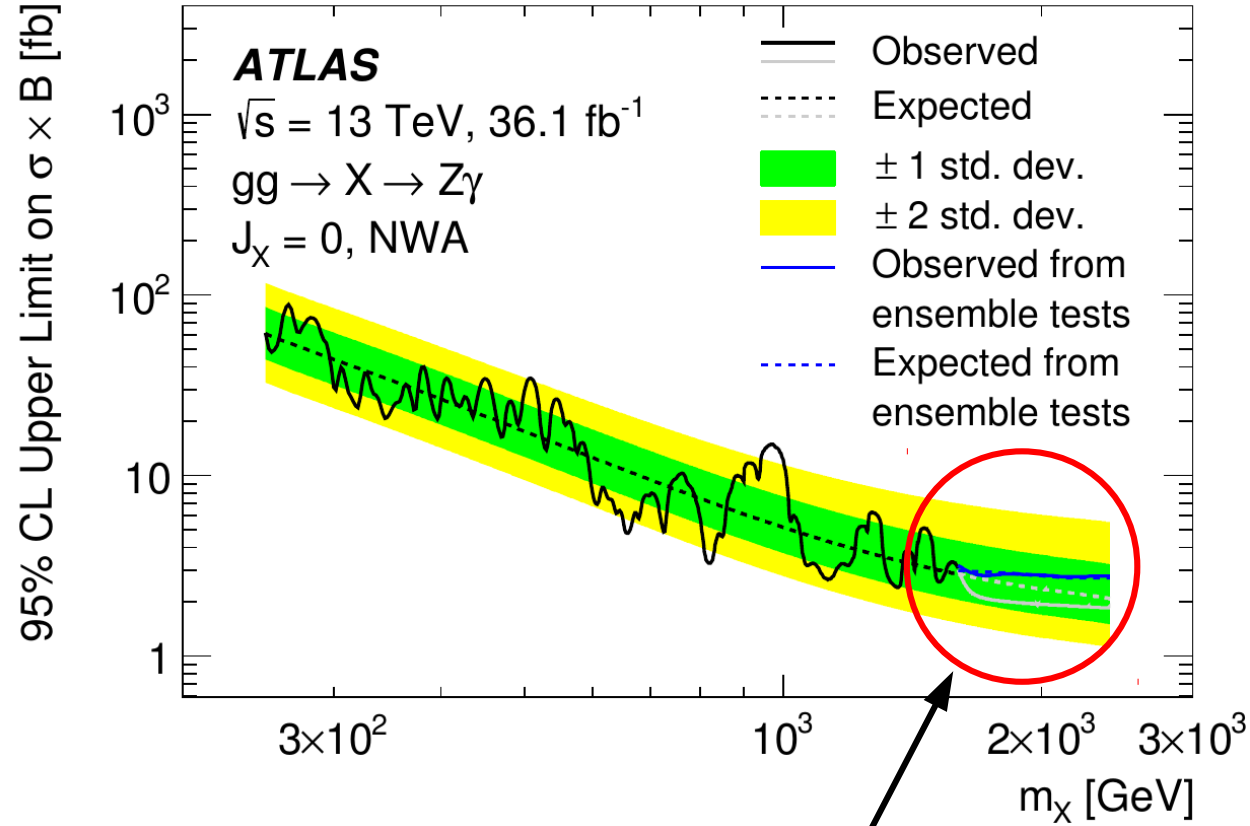
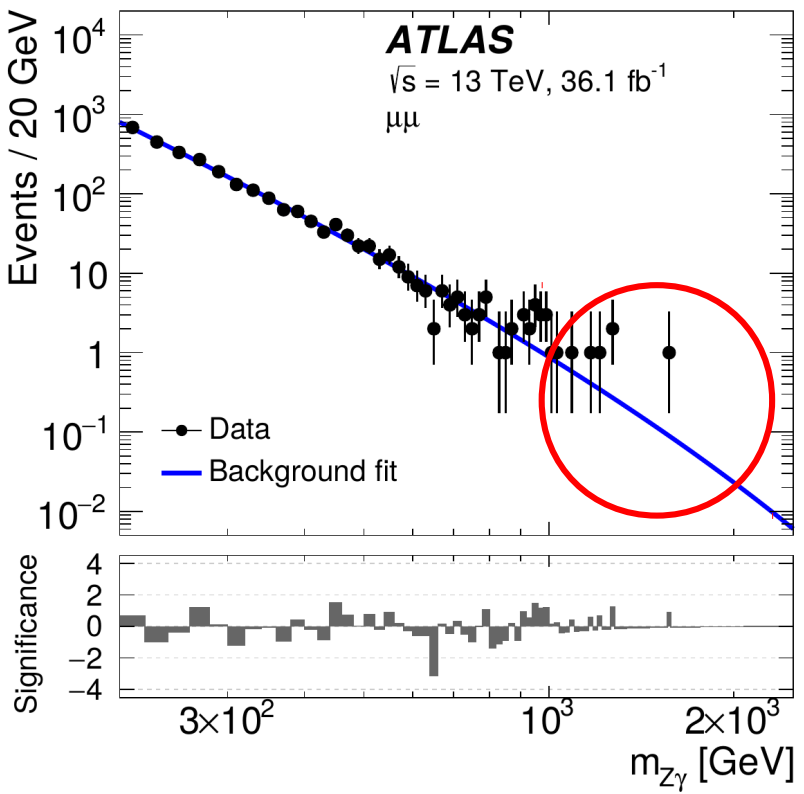


$q_0$

# 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)

# Outline

---

Computing statistics results:

Limits

Confidence intervals

Profiling

**Look-Elsewhere Effect**

Bayesian methods



---

# Look-Elsewhere Effect

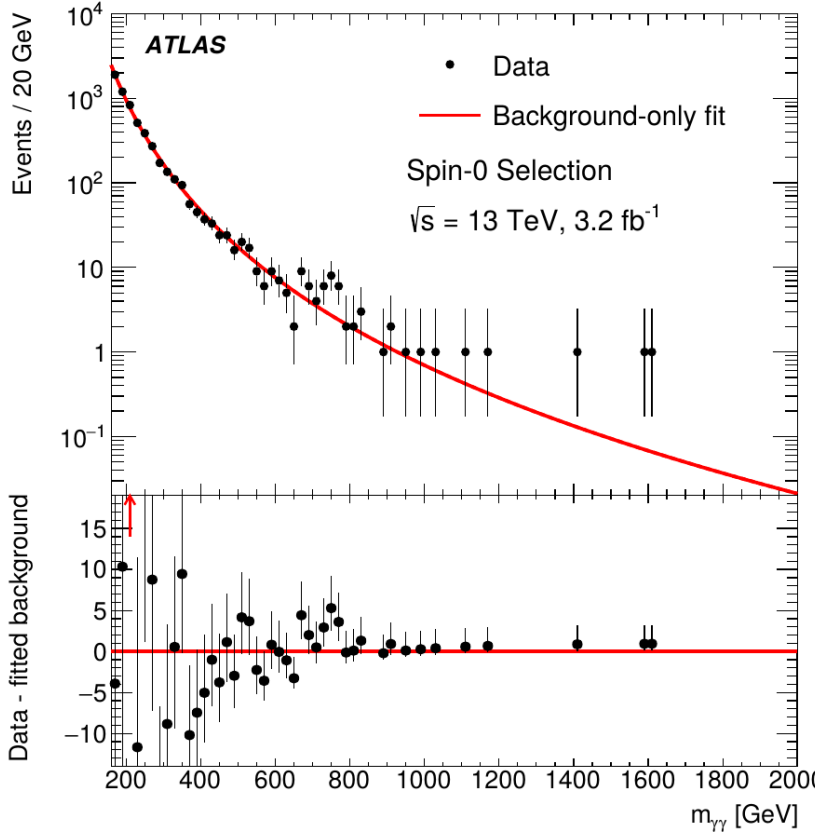
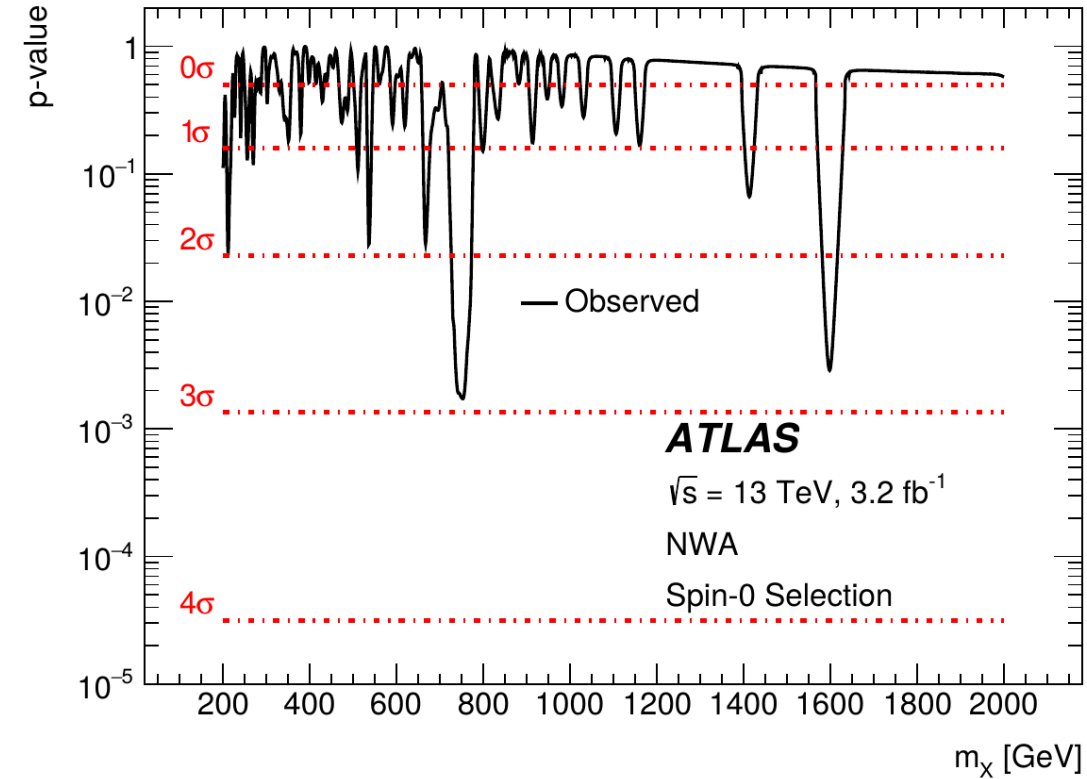
# Look-Elsewhere effect

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**

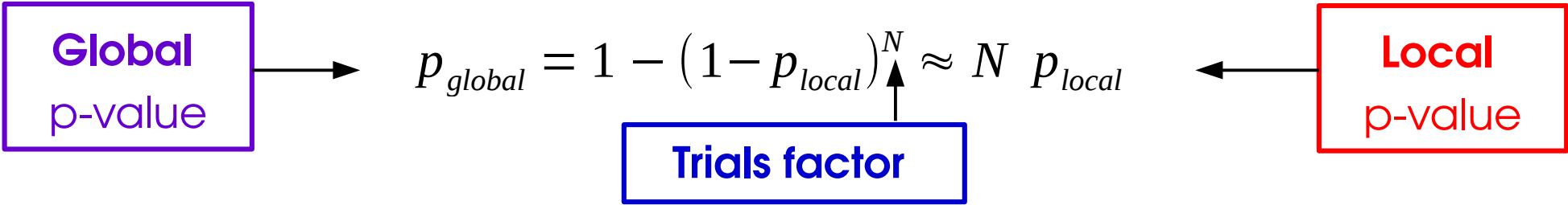


→ More likely to find an excess **anywhere in the range**, rather than in a **predefined** location  
 ⇒ **Look-elsewhere effect** (LEE)

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

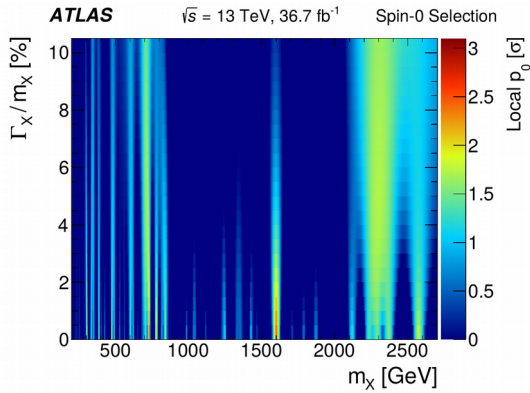


→  $p_{global} > p_{local} \Rightarrow Z_{global} < Z_{local}$  – global fluctuation more likely ⇒ less significant

**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:

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

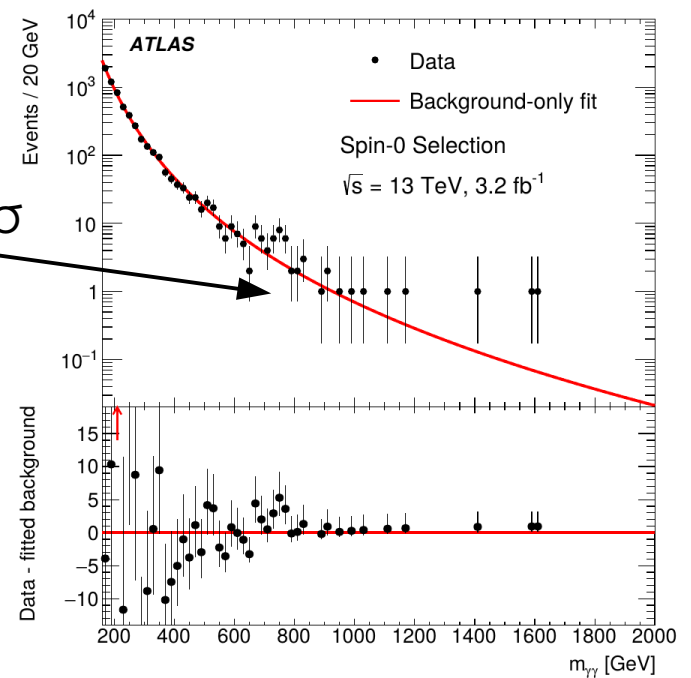
⇒ The frequency at which a given  $Z_0$  is found **is** the global p-value

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...

**Local**  $3.9\sigma$



⊕ **Exact treatment**

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

# 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  $N_{\text{indep}}$** ,  
also depends on  $Z_{\text{local}}$  !

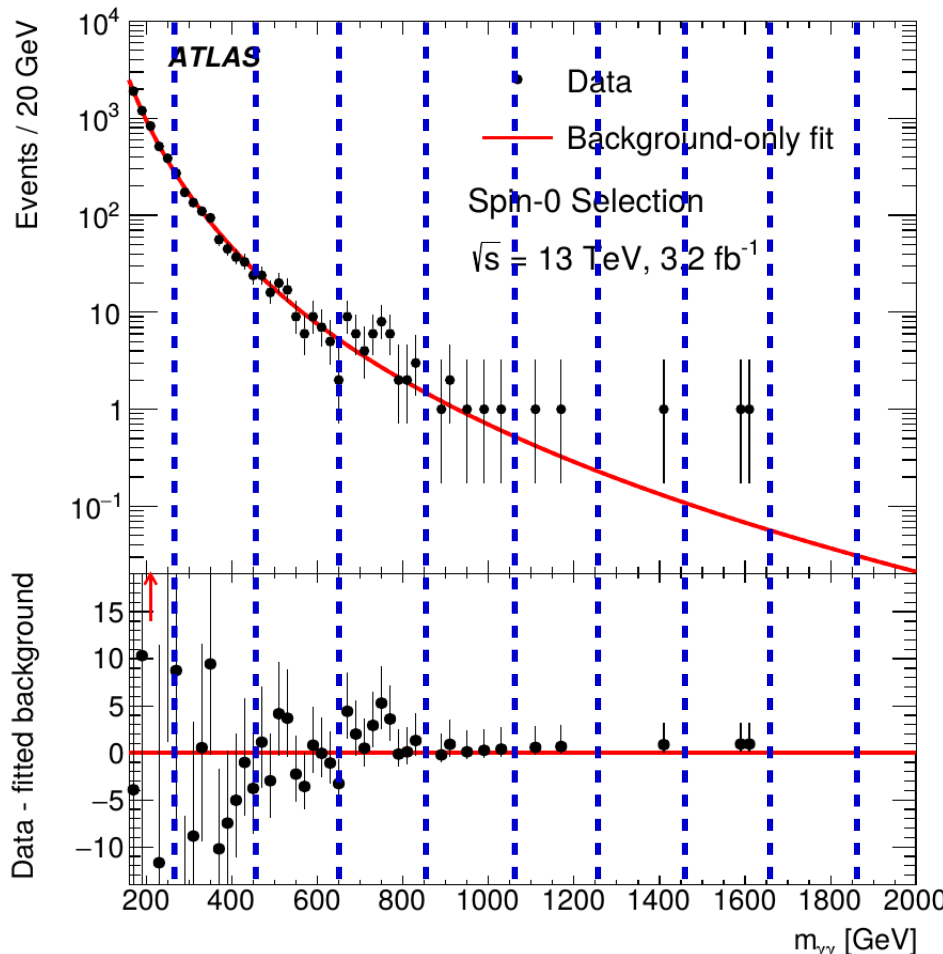
**Why ?**

- slice scan range into  $N_{\text{indep}}$  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  $N_{\text{trials}}$  is  **$> N_{\text{indep}}$**  !



# 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**  $N_{\text{indep}}$ ,  
also depends on  $Z_{\text{local}}$  !

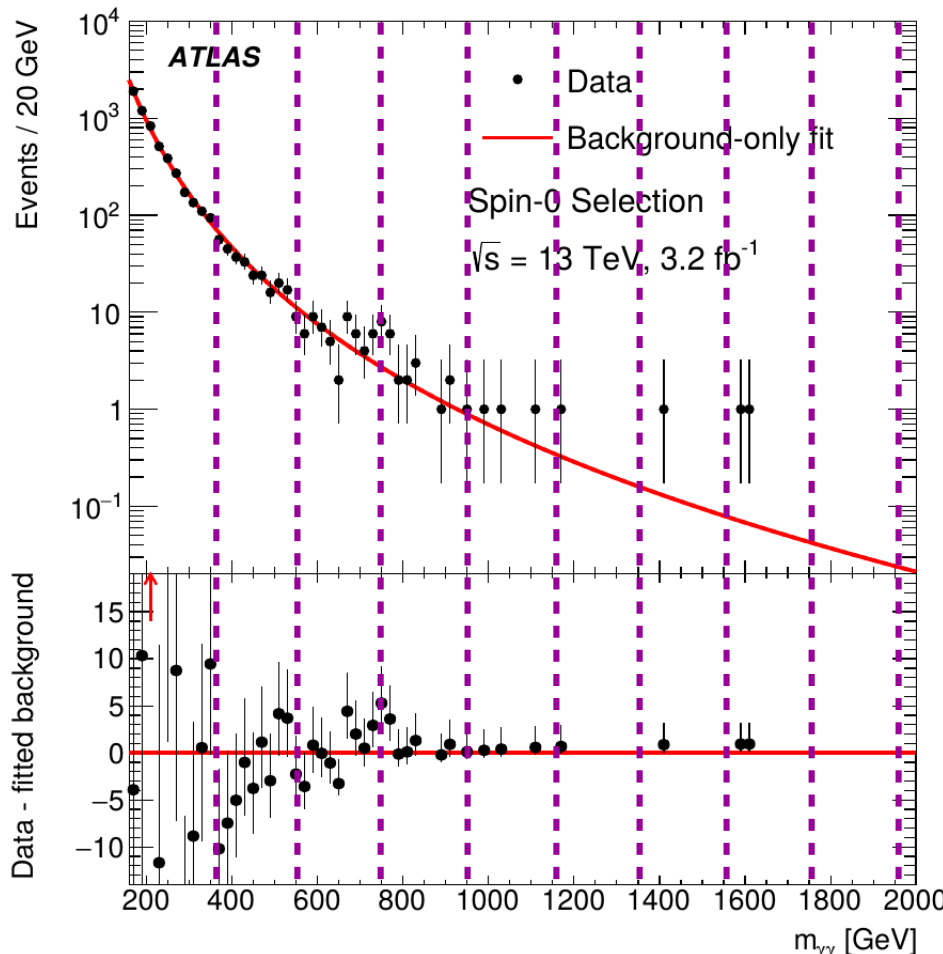
**Why ?**

- slice scan range into  $N_{\text{indep}}$  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  $N_{\text{trials}}$  is **>**  $N_{\text{indep}}$  !



# Illustrative Example

**Test on a simple example:** generate toys with

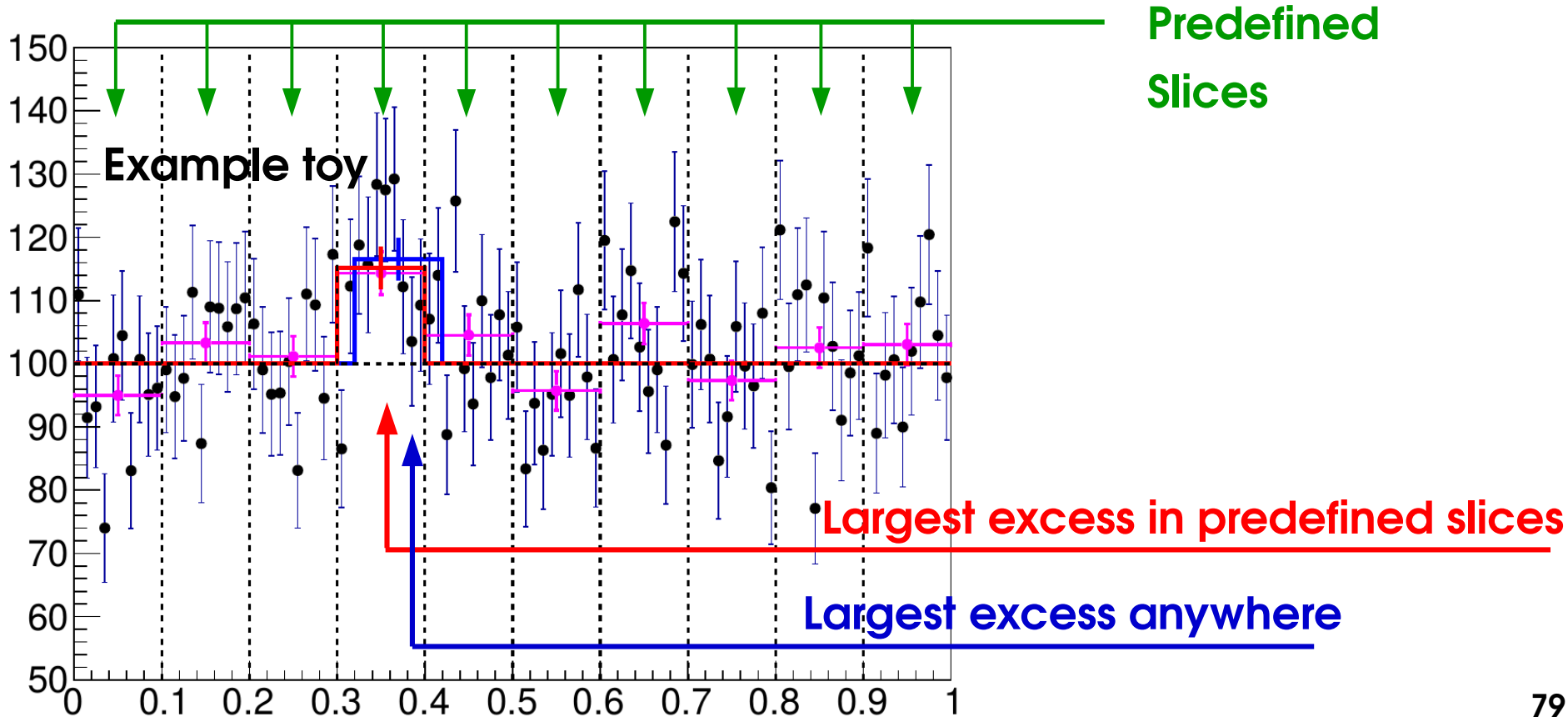
→ flat background (100 events/bin)

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

**Two configurations:**

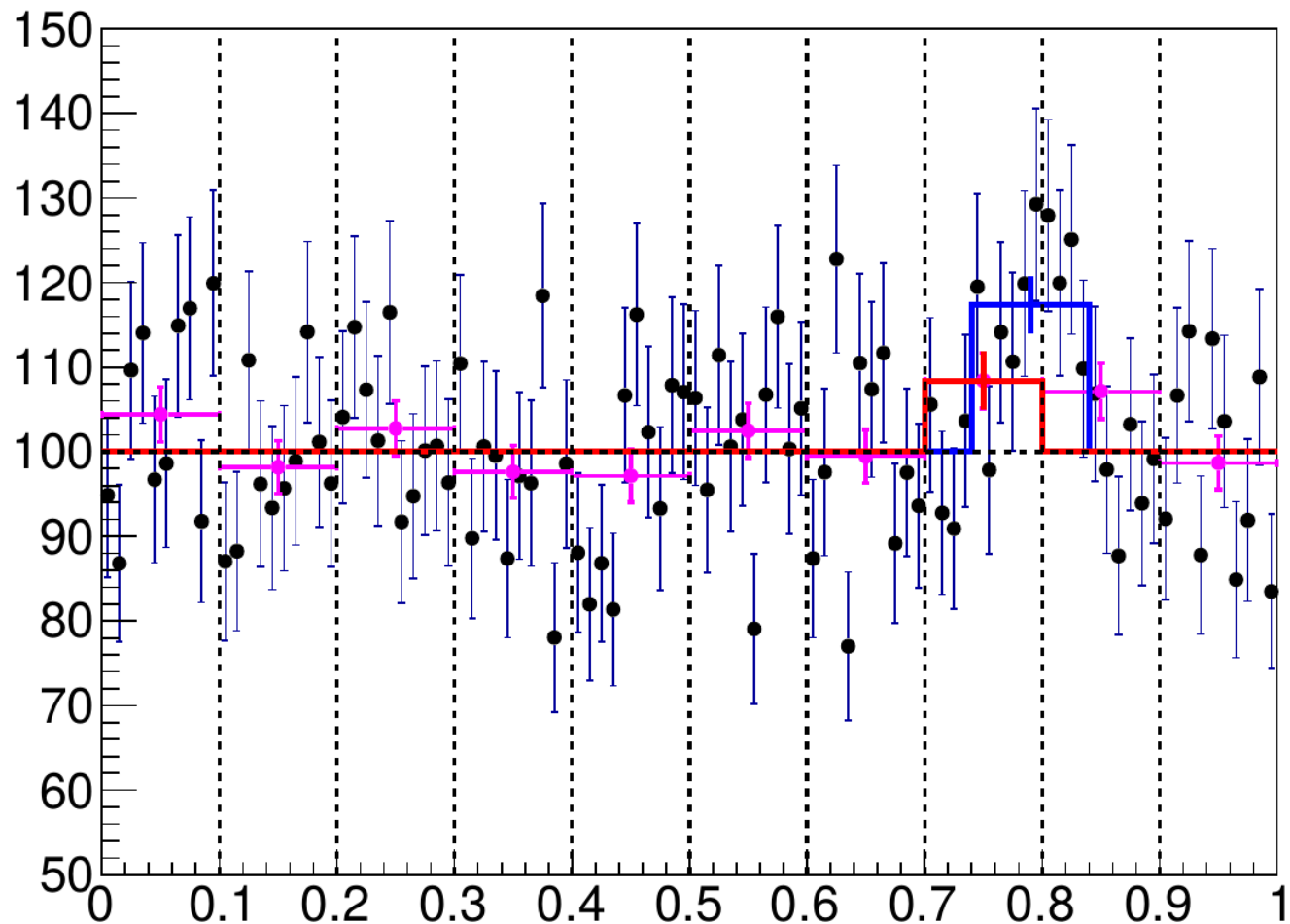
1. Look only in 10 slices of the full spectrum

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



# Illustrative Example (2)

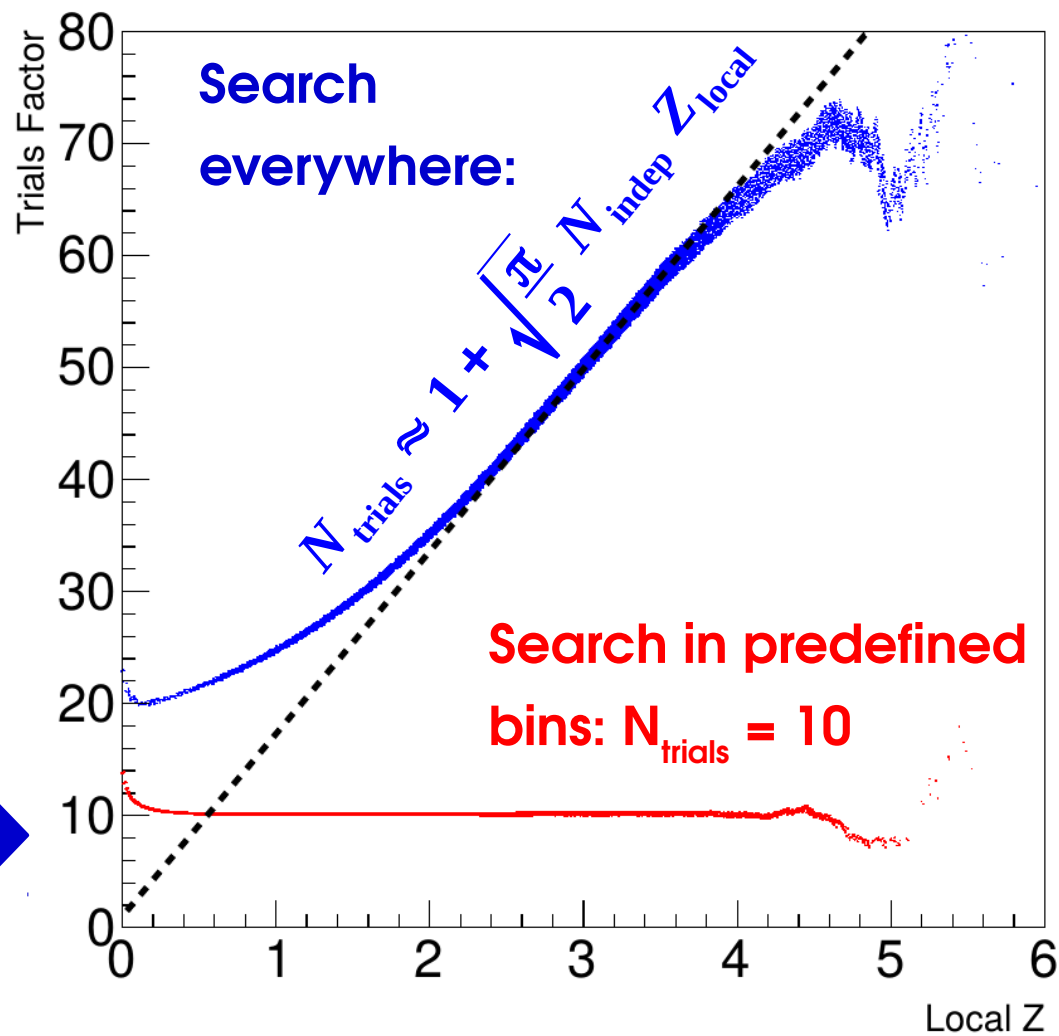
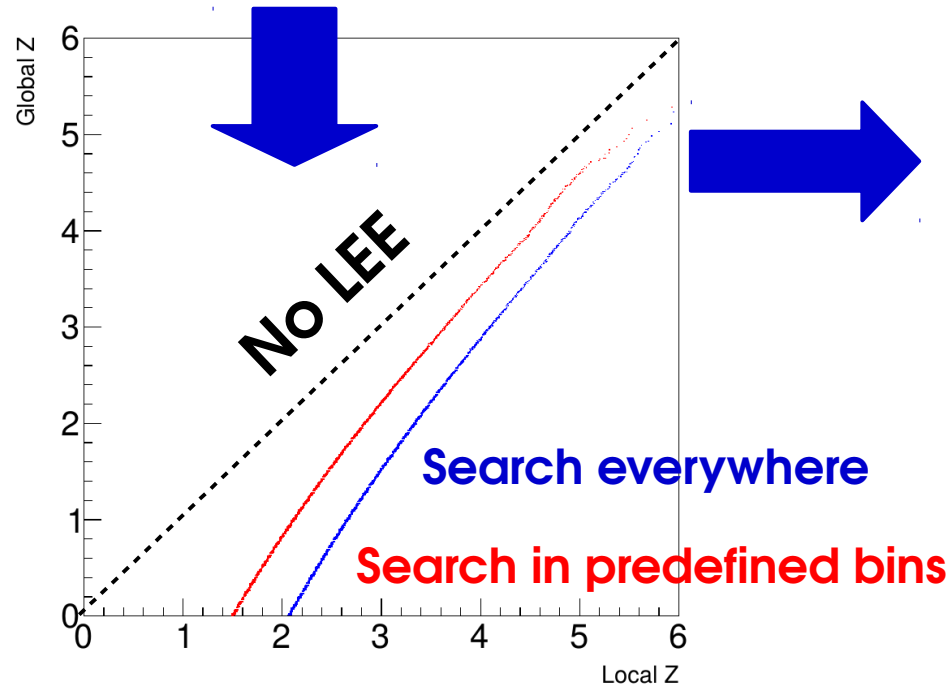
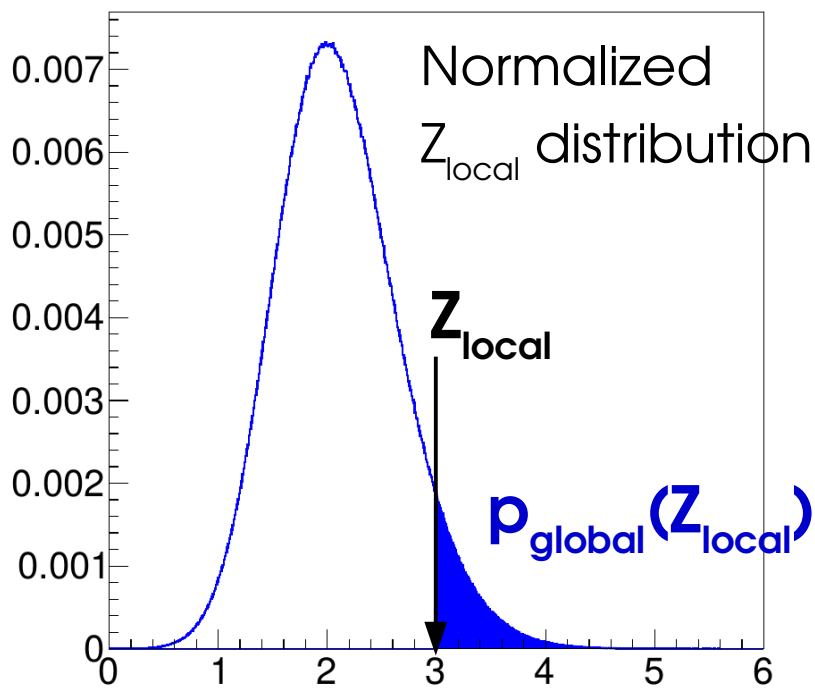
Very different results if the excess is **near a boundary** :



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  $Z_{\text{local}}$  dependence

# $Z_{\text{Global}}$ Asymptotics Extrapolation

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

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

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

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

→ Get  $N_{\text{up}}$  From toys? but high  $Z_{\text{local}} \Rightarrow$  many toys needed

$\Rightarrow$  calibrate for small  $Z_{\text{test}}$ , apply result to higher  $Z_{\text{local}}$ .

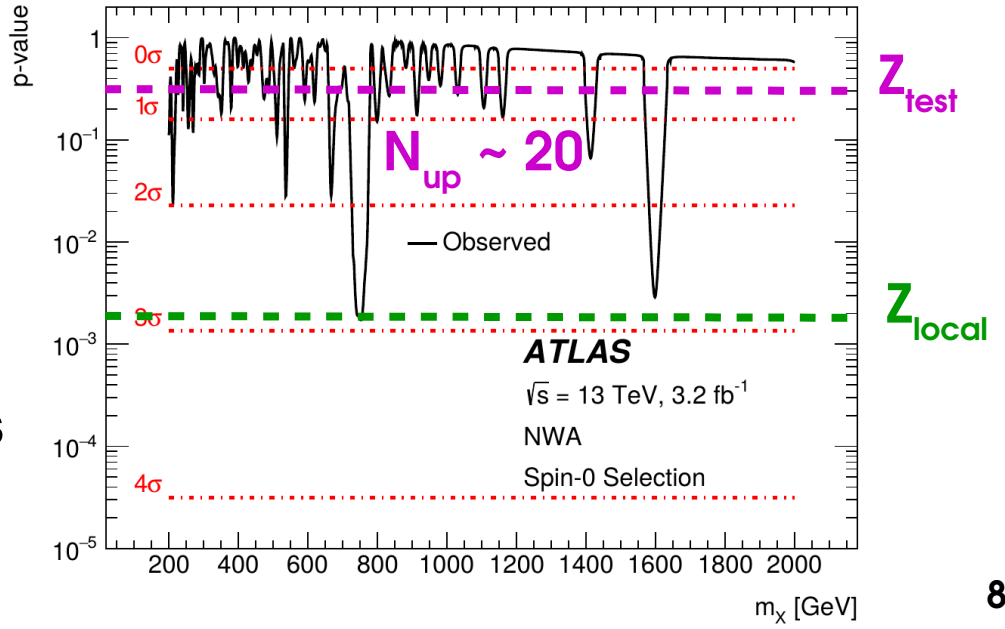
Can choose arbitrarily small  $Z_{\text{test}}$

$\Rightarrow$  many excesses

$\Rightarrow$  can measure  $N_{\text{up}}$  in data (1 "toy")

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

if large stat uncertainty from too few excesses

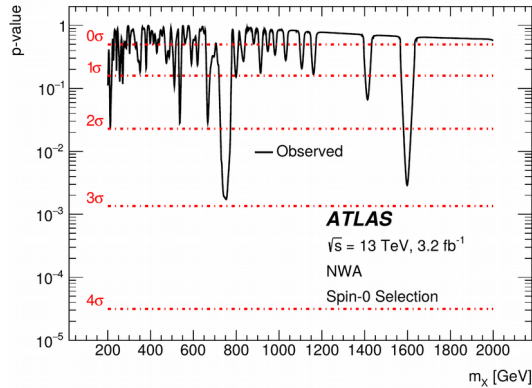


**Generalization to 2D scans:** consider sections at a fixed  $Z_{\text{test}}$ , compute its

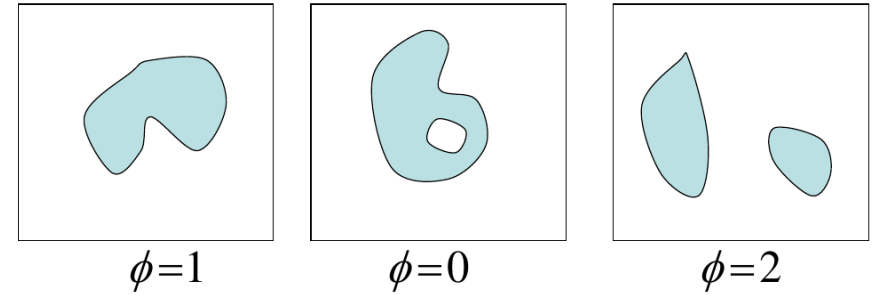
**Euler characteristic**  $\phi$ , and use

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

→ Generalizes 1D bump counting

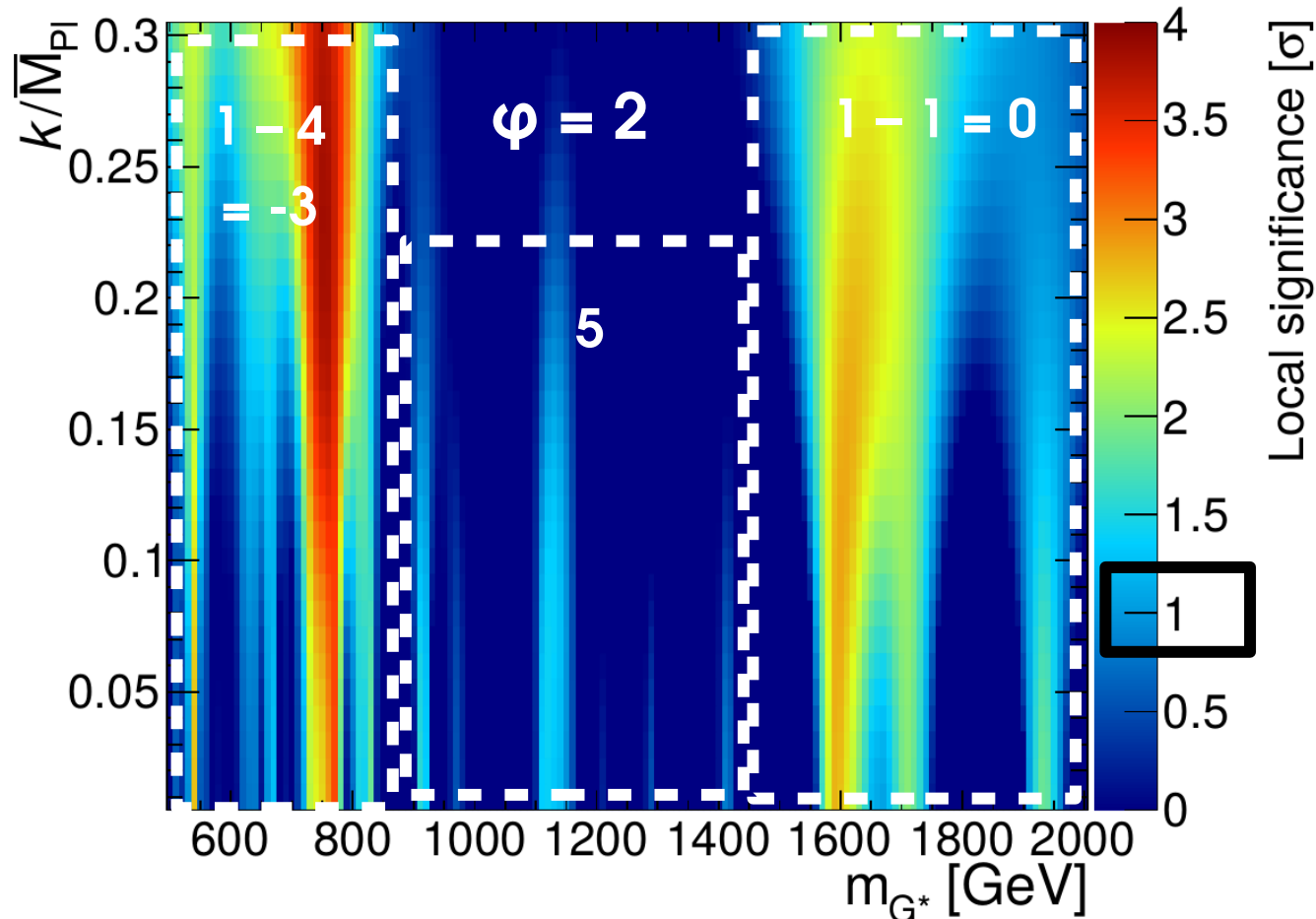


Now need to determine 2 constants  $N_1$  and  $N_2$ , from Euler  $\phi$  measurements at 2 different  $Z_{\text{test}}$  values.



**ATLAS**

$\sqrt{s} = 13 \text{ TeV}, 3.2 \text{ fb}^{-1}$  Spin-2 Selection



# Outline

---

Computing statistics results:

Limits

Confidence intervals

Profiling

Look-Elsewhere Effect

**Bayesian methods**

---

# Bayesian Methods

# 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  $m_H$ .

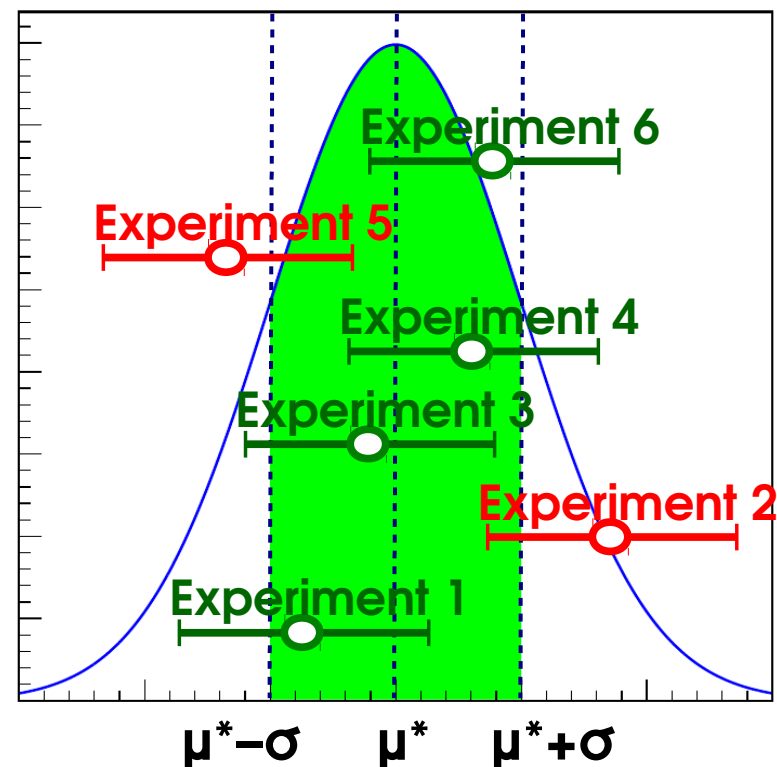
→ “ **$5\sigma$  Higgs discovery**”

- 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  **$P(\text{no Higgs} \mid \text{data})$**  – but all we have is  **$P(\text{data} \mid \text{no Higgs})$**



# 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  $\mu$  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( $\theta$ ) 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(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<sub>s</sub>: Bayesian case

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

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

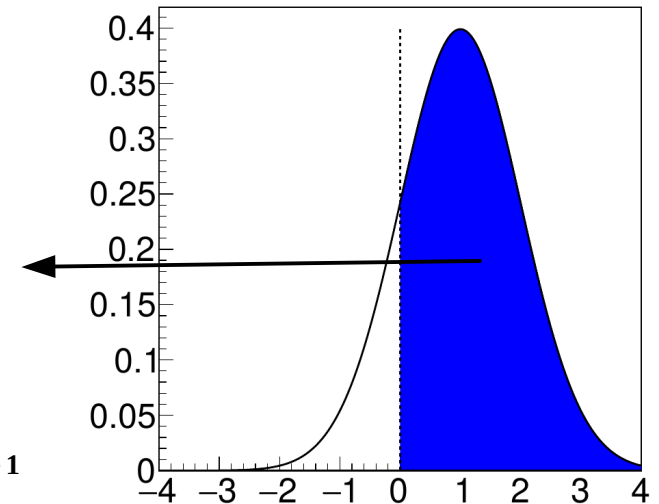
**Bayesian:**  $G(\theta)$  is actually a **prior** on  $\theta \Rightarrow$  perform integral (**marginalization**)

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

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:**

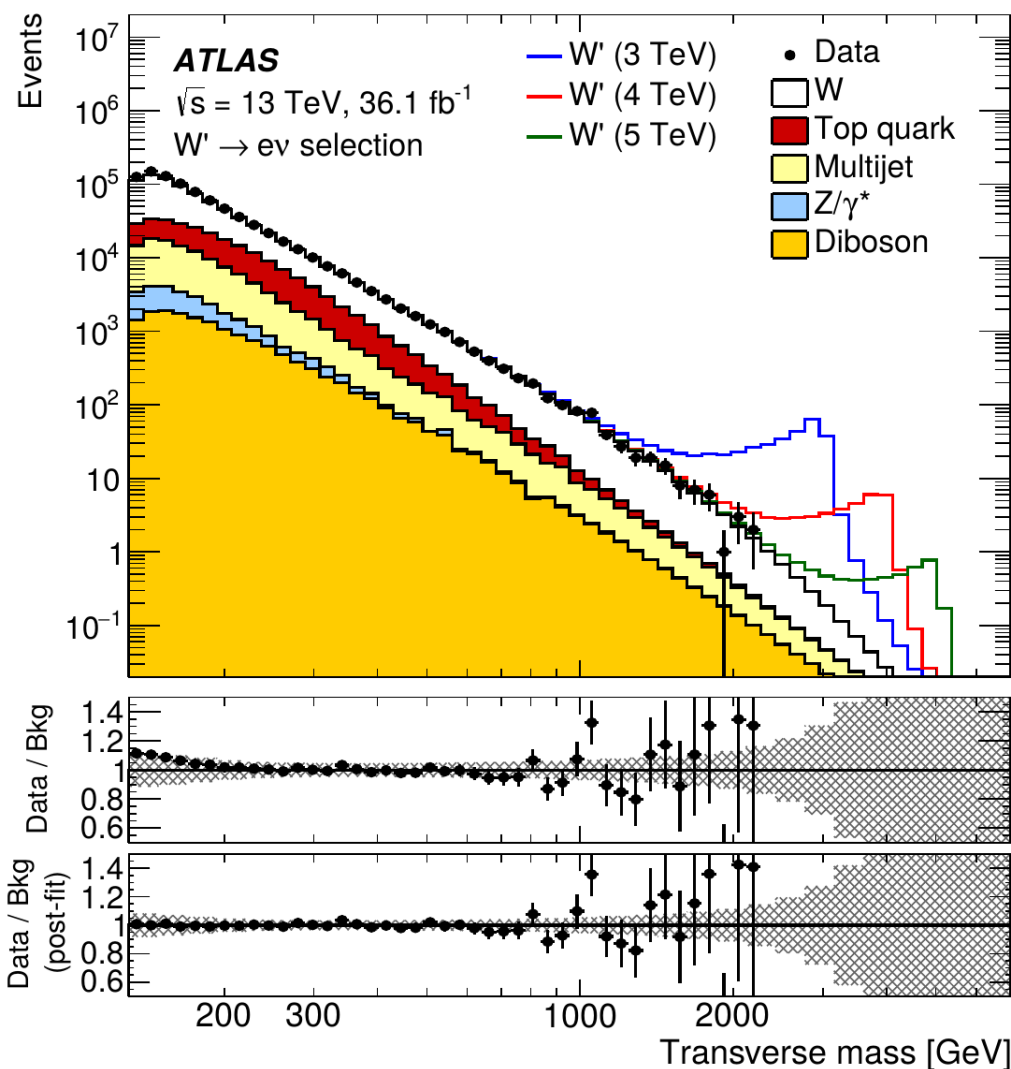
$$\int_{S_{\text{up}}}^{\infty} P(S | n) dS = 5\% = \left[ 1 - \Phi \left( \frac{S_{\text{up}} - (n - B)}{\sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}} \right) \right] \left[ \Phi \left( \frac{n - B}{\sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{syst}}^2}} \right) \right]^{-1}$$

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

same result as CLs!

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

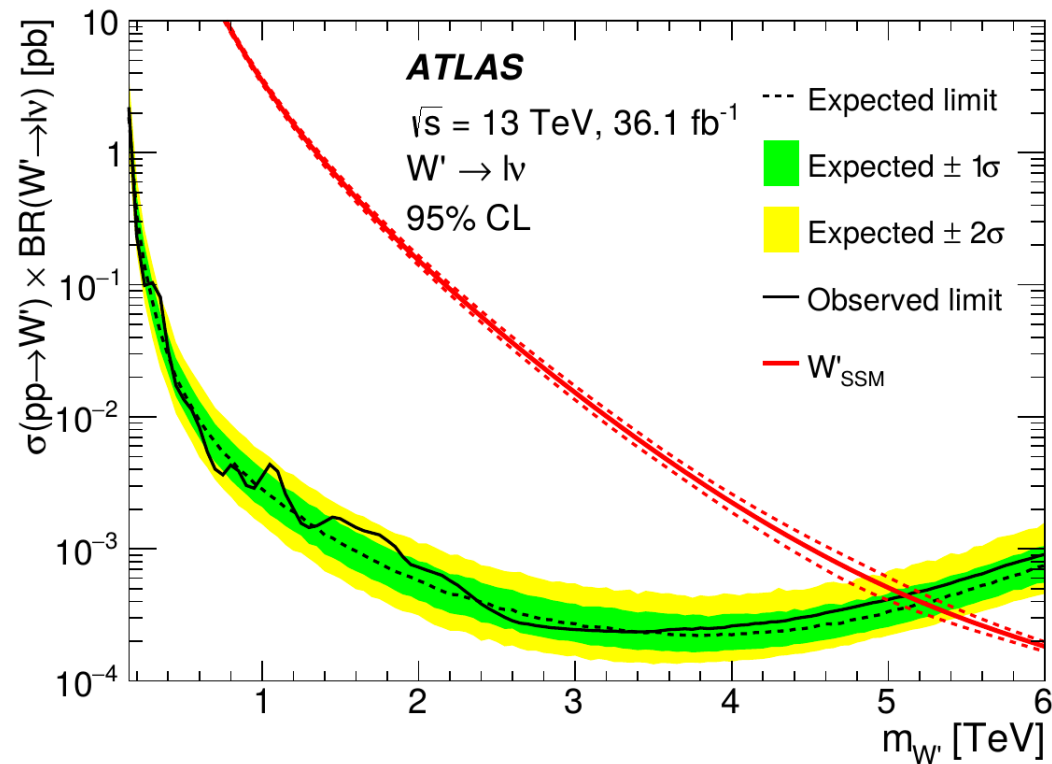
- **POI:**  $W' \sigma \times B \rightarrow$  use flat prior over  $[0, +\infty[$ .
- **NPs:** syst on **signal  $\epsilon$**  (6 NPs), **bkg** (6), **lumi** (1)  $\rightarrow$  integrate over Gaussian priors



Trigger  
 Lepton reconstruction and identification  
 Lepton momentum scale and resolution  
 $E_T^{\text{miss}}$  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 $5\sigma$ ?

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

→ Overly conservative ?

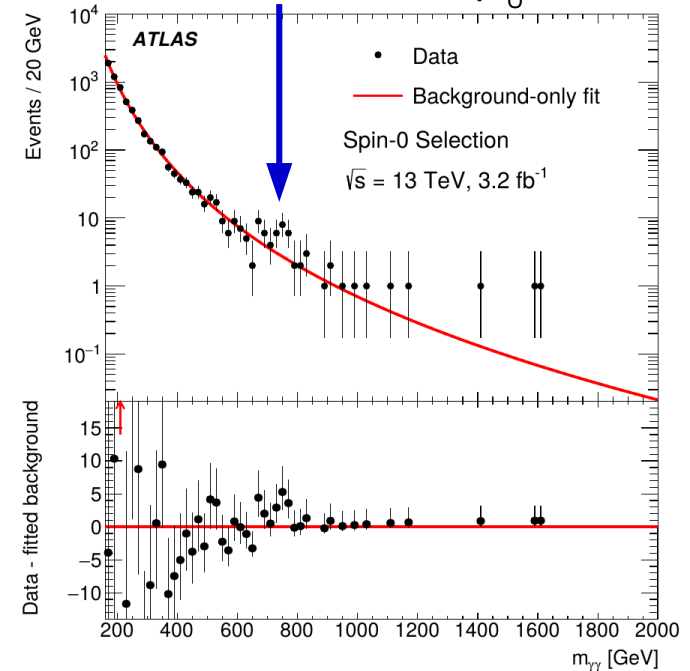
→ Do we even know the sampling distributions so far out ?

**Local**  $3.9\sigma$ ,  $p_0 = 5\text{E-}5$

**Global**  $2.1\sigma$ ,  $p_0 = 2\text{E-}2$

**Reasons for sticking with  $5\sigma$**  (from Louis Lyons):

- **LEE** : searches typically cover multiple independent regions  
 $\Rightarrow$  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  $\Rightarrow$  factor 2 error on Z...
- **History**:  $3\sigma$  and  $4\sigma$  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*

$\Rightarrow$  Stay with  $5\sigma$ ...