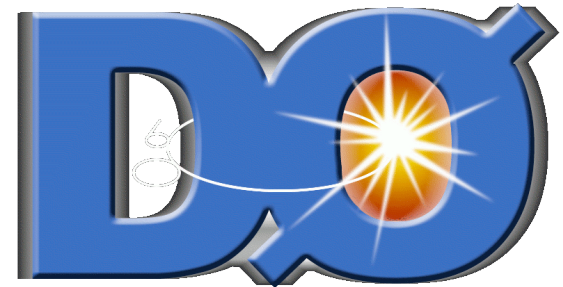
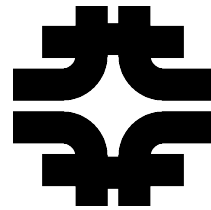


# Statistical Techniques for Combining Tevatron Higgs Boson Searches



Tom Junk  
Fermilab



- Exchanging Experimental Results
- Bayesian Limits
- $CL_s$  Limits
- Discovery Techniques
- Tevatron Combined Results

<http://tevnphwg.fnal.gov>

<http://www-cdf.fnal.gov/physics/new/hdg/hdg.html>

<http://www-d0.fnal.gov/Run2Physics/WWW/results/higgs.htm>

# Exchanging Experimental Results for Combination

At LEP and at the Tevatron, we exchanged histograms of observed and predicted events.

## Many advantages:

- Crosscheck analyzers' work:  
Signal and background checksum  
Limit/discovery recalculations  
check for "broken" bins
  - $s > 0$  when  $b = 0$
  - any observation or prediction  $< 0$
- Can make control plots
- Can try a great variety of statistical treatments
  - Profile Likelihood, Bayesian,  $CL_s$  and compare each one
- Can make expected limits/LLR distributions without approximations
- Can draw the  $\pm 1\sigma$ ,  $\pm 2\sigma$  bands on expected limits with MC
- Can point to excesses and deficits to explain why limits and p-values are as they are
- Can accommodate new cross sections and branching ratios by scaling
- Can pick and choose signals if more than one expected (e.g., do 4th gen analysis with  $H \rightarrow WW$  without WH, ZH and VBF)
- Pre-binned histograms mean combiners don't have to choose binning, reducing mistakes, inconsistencies
- possibly less work for the analyzer

## Disadvantages:

- Lots of work/CPU!
- Have to share preliminary histos (your competitors may find your mistakes!)

# Exchanging Experimental Results for Combination

## Systematic uncertainties itemized by named source

- Asymmetric Rate errors on each predicted component
- Shape errors supplied as alternate shape histograms  
Bin-by-bin ratios are inconvenient -- example  $m_{jj}$  histogram where the variation is “horizontal” and not vertical.  
Need shape interpolators/extrapolators to use them.  
Typically  $\pm 1\sigma$  shape variations are explored one source at a time by analyzers. Analyzers will ask combiners to extrapolate out to arbitrary  $\pm n\sigma$  shapes (!)  
-- practical difficulty: [How to estimate  \$5\sigma\$  systematics?](#)
- Bin-by-bin independent uncertainties (MC statistics)
- Names used to categorize correlations in a way easy to understand and check
- Give names to exchanged template histograms please!

# Cross Section and Branching Ratio Alignment

TABLE I: The production cross sections and decay branching fractions for the SM Higgs boson assumed for the combination.

$m_H$ (GeV/ $c^2$ )	$\sigma_{gg \rightarrow H}$ (fb)	$\sigma_{WH}$ (fb)	$\sigma_{ZH}$ (fb)	$\sigma_{VBF}$ (fb)	$\sigma_{t\bar{t}H}$ (fb)	$B(H \rightarrow b\bar{b})$ (%)	$B(H \rightarrow c\bar{c})$ (%)	$B(H \rightarrow \tau^+\tau^-)$ (%)	$B(H \rightarrow W^+W^-)$ (%)	$B(H \rightarrow ZZ)$ (%)	$B(H \rightarrow \gamma\gamma)$ (%)
100	1861	291.9	169.8	99.5	8.000	80.33	3.542	7.920	1.052	0.1071	0.1505
105	1618	248.4	145.9	93.3	7.062	78.57	3.463	7.821	2.307	0.2035	0.1689
110	1413	212.0	125.7	87.1	6.233	75.90	3.343	7.622	4.585	0.4160	0.1870
115	1240	181.9	108.9	79.07	5.502	71.95	3.169	7.288	8.268	0.8298	0.2029
120	1093	156.4	94.4	71.65	4.857	66.49	2.927	6.789	13.64	1.527	0.2148
125	967	135.1	82.3	67.37	4.279	59.48	2.617	6.120	20.78	2.549	0.2204
130	858	116.9	71.9	62.5	3.769	51.18	2.252	5.305	29.43	3.858	0.2182
135	764	101.5	63.0	57.65	3.320	42.15	1.854	4.400	39.10	5.319	0.2077
140	682	88.3	55.3	52.59	2.925	33.04	1.453	3.472	49.16	6.715	0.1897
145	611	77.0	48.7	49.15	2.593	24.45	1.075	2.585	59.15	7.771	0.1653
150	548	67.3	42.9	45.67	2.298	16.71	0.7345	1.778	68.91	8.143	0.1357
155	492	58.9	37.9	42.19	2.037	9.88	0.4341	1.057	78.92	7.297	0.09997
160	439	50.8	33.1	38.59	1.806	3.74	0.1646	0.403	90.48	4.185	0.05365
165	389	44.6	30.0	36.09	1.607	1.29	0.05667	0.140	95.91	2.216	0.02330
170	349	40.2	26.6	33.58	1.430	0.854	0.03753	0.093	96.39	2.351	0.01598
175	314	35.6	23.7	31.11	1.272	0.663	0.02910	0.073	95.81	3.204	0.01236
180	283	31.4	21.1	28.57	1.132	0.535	0.02349	0.059	93.25	5.937	0.01024
185	255	28.2	18.9	26.81	1.004	0.415	0.01823	0.046	84.50	14.86	0.008128
190	231	25.1	17.0	24.88	0.890	0.340	0.01490	0.038	78.70	20.77	0.006774
195	210	22.4	15.3	23	0.789	0.292	0.01281	0.033	75.88	23.66	0.005919
200	192	20.0	13.7	21.19	0.700	0.257	0.01128	0.029	74.26	25.33	0.005285

All channels must use consistent predictions

# Tevatron Correlated Systematic Errors I

Total Systematic error count: 129 (not counting bin-by-bin errors)

Note: correlation in errors on backgrounds between experiments helps sensitivity! One experiment is another experiment's control sample.

Luminosity: 3.8% Correlated CDF and DØ  $\sigma_{\text{inel}}$ (ppbar)  
4.4% detector-specific

Diboson Cross Sections:  
Defined for  
 $75 < m_H < 105$  GeV

$$\sigma_{W+W^-} = 11.34_{-0.49}^{+0.56} \text{ (scale)} \quad +0.35_{-0.28} \text{ (PDF) pb}$$
$$\sigma_{W\pm Z} = 3.22_{-0.17}^{+0.20} \text{ (scale)} \quad +0.11_{-0.08} \text{ (PDF) pb}$$
$$\sigma_{ZZ} = 1.20_{-0.04}^{+0.05} \text{ (scale)} \quad +0.04_{-0.03} \text{ (PDF) pb}$$

$t\bar{t}$  Cross Section: Moch and Uwer, evaluated at  
 $m_t = 173 \pm 1.2$  GeV is (using MSTW2008 PDFs)

$$\sigma_{t\bar{t}} = 7.04_{-0.36}^{+0.24} \text{ (scale)} \pm 0.14 \text{ (PDF)} \pm 0.30 \text{ (mass)}$$

# Tevatron Correlated Systematic Errors II

Signal Cross Section uncertainties (using MSTW PDFs)

WH, ZH:  $\pm 5\%$

gg $\rightarrow$ H:  $\pm 17.5\%$  (weighted scale over jet samples)  
 $\pm 12.5\%$  (weighted PDF). Pt spectra reweighted to

NNLO+NNLL predictions

Errors taken from Anastasiou, Dissertori, Stockli, and Webber

VBF:  $\pm 10\%$

Theory errors applied to SM interpretations, but taken off for cross-section times branching ratio limits.

*CDF-D0 Uncorrelated errors:*

K-factors (data driven)

trigger efficiency

b-tag efficiency and mistags

jet energy scale

lepton ID, fakes and conversions

MET modeling



Correlated *within*  
CDF and D0 where  
appropriate

## Steps Required for Combination

- Histograms and named rate *and shape* errors exchanged
- Check stacked histograms and systematic tables with analysis documentation total counts:
  - data, signal, background
    - look for bins with  $b=0$  and have data events (bad!)
- Repeat individual channel limits -- compare against approved results.
- Assess correlations on systematics

CDF and D0 teams each do three combinations, using Bayesian and  $CL_s$  techniques.

CDF

D0

Tevatron

Consistency at the better than 10% level required for all combinations at all test masses. Quote Bayesian limits (historical)

# Mini-Review: Bayesian Limits

$$L(r, \theta) = \prod_{\text{channels}} \prod_{\text{bins}} P_{\text{Poiss}}(\text{data} | r, \theta)$$

Where  $r$  is an overall signal scale factor, and  $\theta$  represents all nuisance parameters.

$$P_{\text{Poiss}}(\text{data} | r, \theta) = \frac{(rs_i(\theta) + b_i(\theta))^{n_i} e^{-(rs_i(\theta) + b_i(\theta))}}{n_i!}$$

where  $n_i$  is observed in each bin  $i$ ,  $s_i$  is the predicted signal for a fiducial model (SM), and  $b_i$  is the predicted background. Dependence of  $s_i$  and  $b_i$  on  $\theta$  includes rate, shape, and bin-by-bin independent uncertainties.



# Mini-Review: Bayesian Limits

Including uncertainties on nuisance parameters  $\theta$

$$L'(data | r) = \int L(data | r, \theta) \pi(\theta) d\theta$$

where  $\pi(\theta)$  encodes our prior belief in the values of the uncertain parameters. Usually Gaussian centered on the best estimate and with a width given by the systematic. The integral is high-dimensional. Markov Chain MC integration is quite useful!

Useful for a variety of results:

Limits: 
$$0.95 = \int_0^{r_{lim}} L'(data | r) \pi(r) dr$$

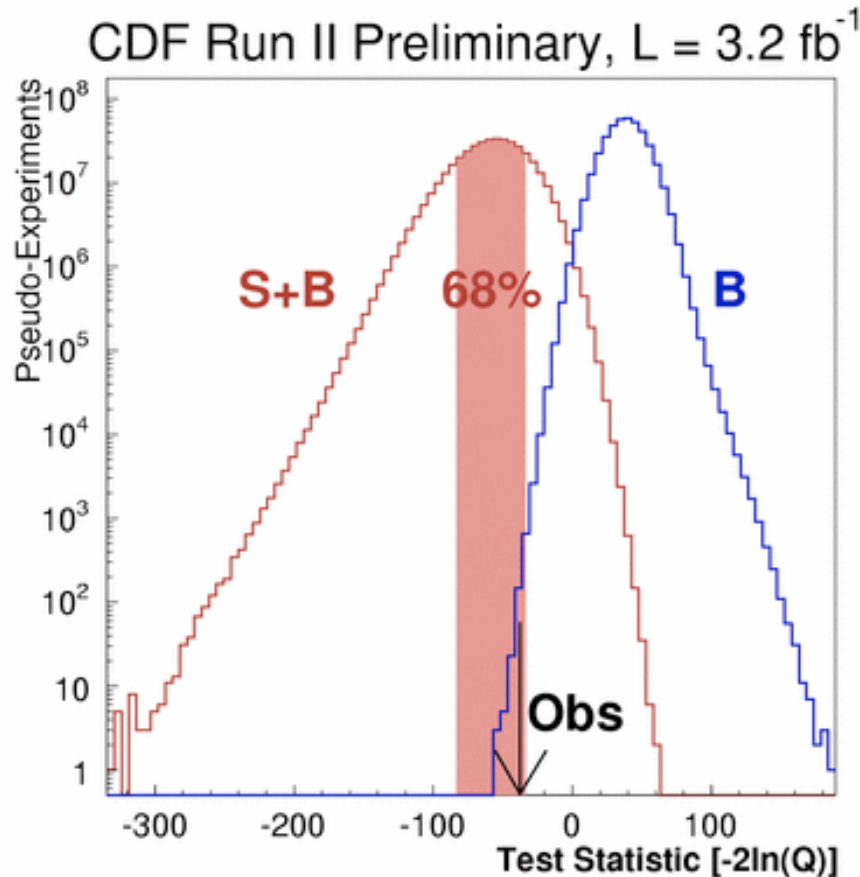
Typically  $\pi(r)$  is constant  
Other options possible.  
Sensitivity to priors a concern.

Measure  $r$ : 
$$0.68 = \int_{r_{low}}^{r_{high}} L'(data | r) \pi(r) dr$$

$$r = r_{max} \begin{matrix} + (r_{high} - r_{max}) \\ - (r_{max} - r_{low}) \end{matrix}$$

Usually: shortest interval containing 68% of the posterior  
(other choices possible)

# Discovery with p-Values



Example: CDF single top.

$$-2\ln Q \equiv LLR \equiv -2\ln \left( \frac{L(\text{data} | s + b, \hat{\theta})}{L(\text{data} | b, \hat{\theta})} \right)$$

100 M s+b and b-only pseudoexperiments, each with fluctuated nuisance parameters, and fit twice.

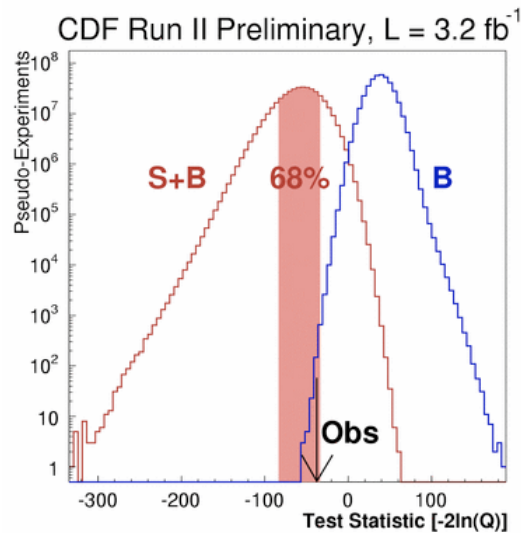
$5\sigma$ : p-value of  $2.77 \times 10^{-7}$  or less.

$3\sigma$ : p-value of  $1.35 \times 10^{-3}$  or less

$2\sigma$ : p-value of 2.28% or less

Buzzword: “Prior Predictive ensemble”

# Fitting and Fluctuating



$$-2\ln Q \equiv LLR \equiv -2\ln \left( \frac{L(\text{data} | s + b, \hat{\theta})}{L(\text{data} | b, \hat{\theta})} \right)$$

- Monte Carlo pseudoexperiments are used to get p-values.
- Test statistic  $-2\ln Q$  is not uncertain for the data.
- Distribution from which  $-2\ln Q$  is drawn is uncertain!

- Nuisance parameter fits in numerator and denominator of  $-2\ln Q$  **do not incorporate systematics into the result.**  
Example -- 1-bin search; all test statistics are equivalent to the event count, fit or no fit.
- Instead, we fluctuate the probabilities of getting each outcome since those are what we do not know. Each pseudoexperiment gets random values of nuisance parameters.
- Can also try values of nuisance parameters that maximize the p-value, but that's very conservative (called the supremum p-value, still needs choices of parameter ranges).
- Why fit at all? It's an optimization. Fitting reduces sensitivity to the uncertain true values and the fluctuated values. For stability and speed, you can choose to fit a subset of nuisance parameters (the ones that are constrained by the data). Or do constrained or unconstrained fits, it's your choice.
- If not using pseudoexperiments but using Wilk's theorem, then the fits are important for correctness, not just optimality.

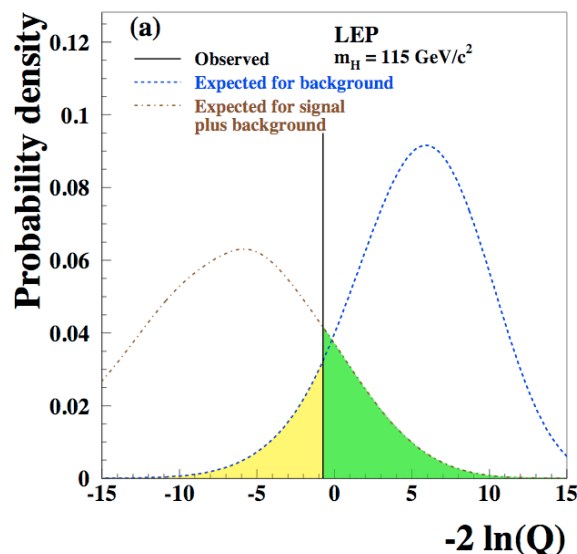
## Mini-Review: $CL_s$ Limits

- Based on p-values using the log likelihood ratio as the test statistic. Neyman-Pearson lemma says LLR is the uniformly most powerful test statistic, although the Neyman-Pearson one fits for the parameter of interest, not just the nuisance parameters, making the null hypothesis a subset of the test hypothesis

$$-2\ln Q \equiv LLR \equiv -2\ln\left(\frac{L(\text{data} | s + b, \hat{\theta})}{L(\text{data} | b, \hat{\theta})}\right)$$

Glen Cowan (really Pearson's) LLR also fits for  $s$  (actually  $r \times s$ ) in the numerator, while  $r = 0$  in the denominator

# Mini-Review: $CL_s$ Limits



p-values:

Yellow area =  $1 - CL_b = 1 - P(-2\ln Q > -2\ln Q_{\text{obs}} \mid b \text{ only})$

Green area =  $CL_{s+b} = P(-2\ln Q > -2\ln Q_{\text{obs}} \mid s+b)$

$$CL_s \equiv CL_{s+b} / CL_b \geq CL_{s+b}$$

Exclude if  $CL_s < 0.05$

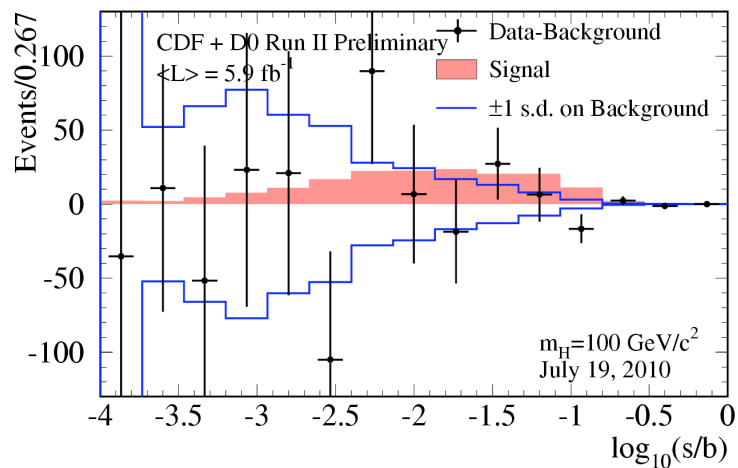
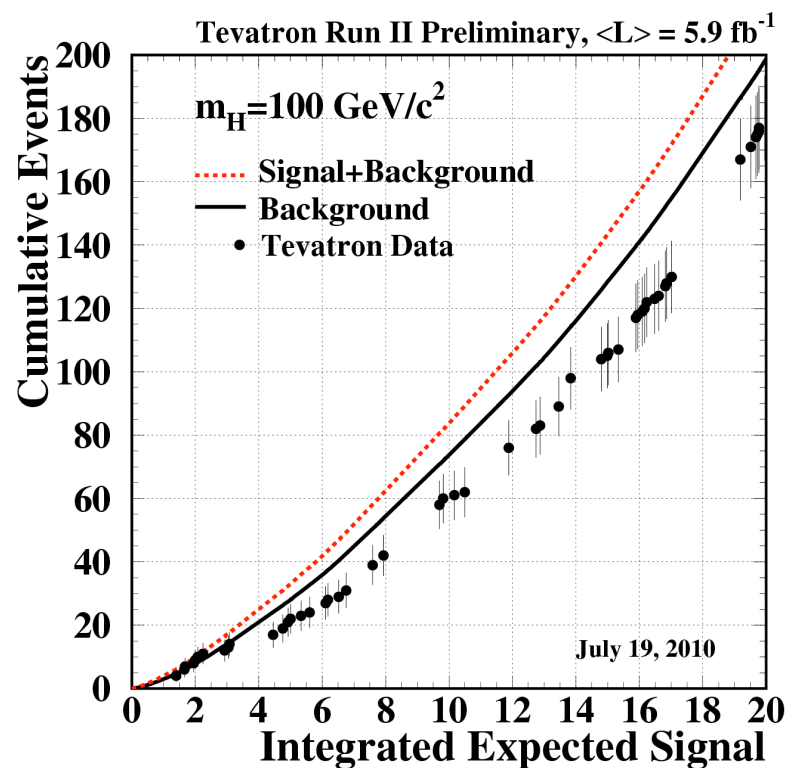
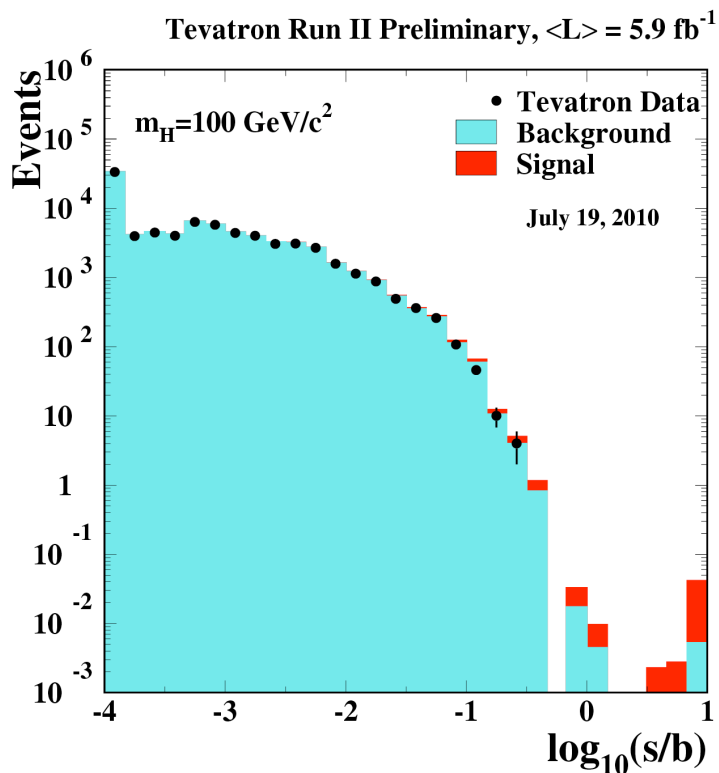
Vary  $r$  until  $CL_s = 0.05$  to get  $r_{\text{lim}}$

- Advantages:

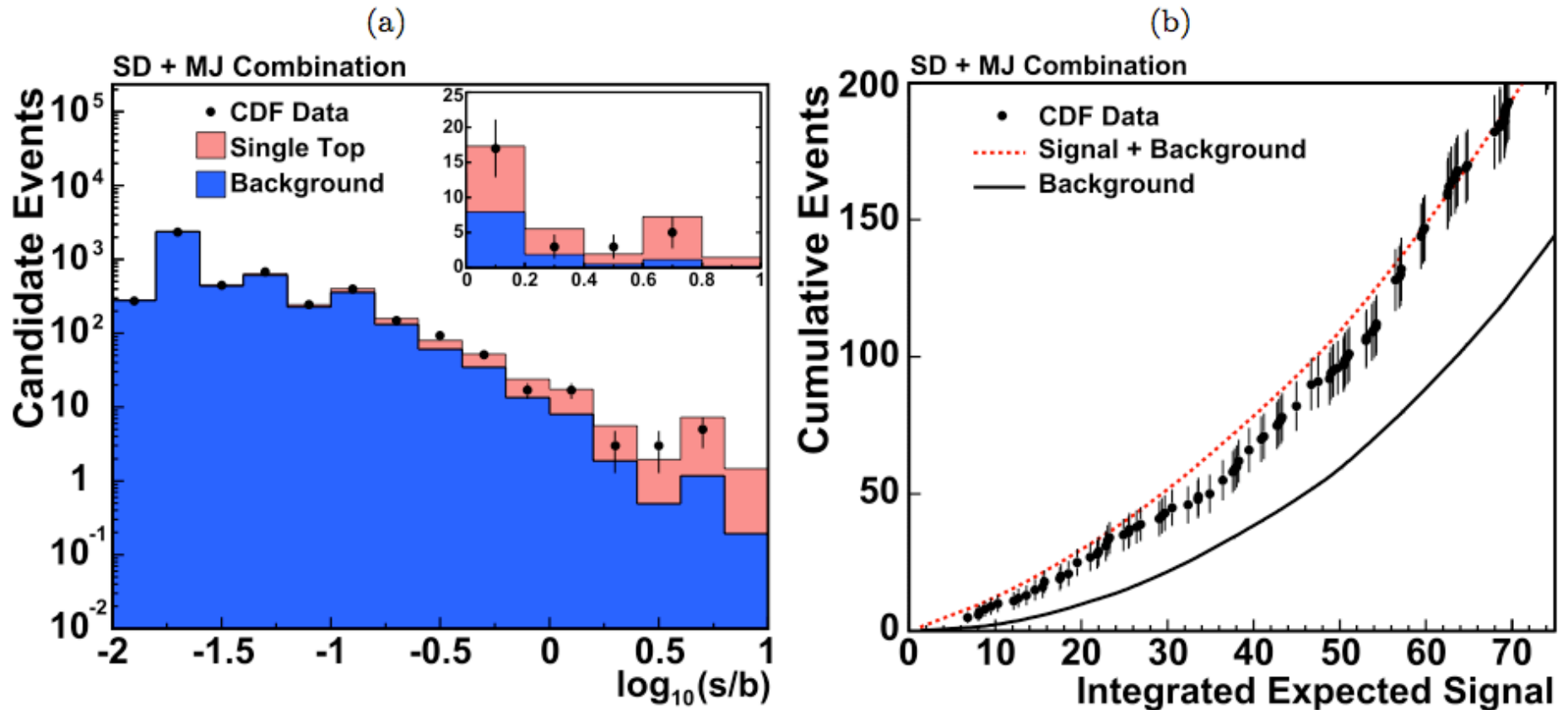
- Exclusion and Discovery p-values are consistent.

Example -- a  $2\sigma$  upward fluctuation of the data with respect to the background prediction appears both in the limit and the p-value as such

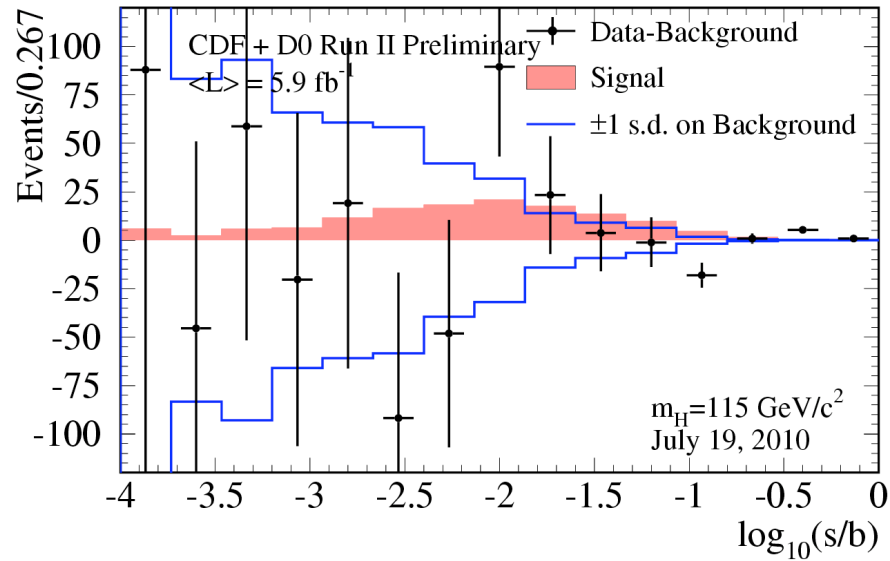
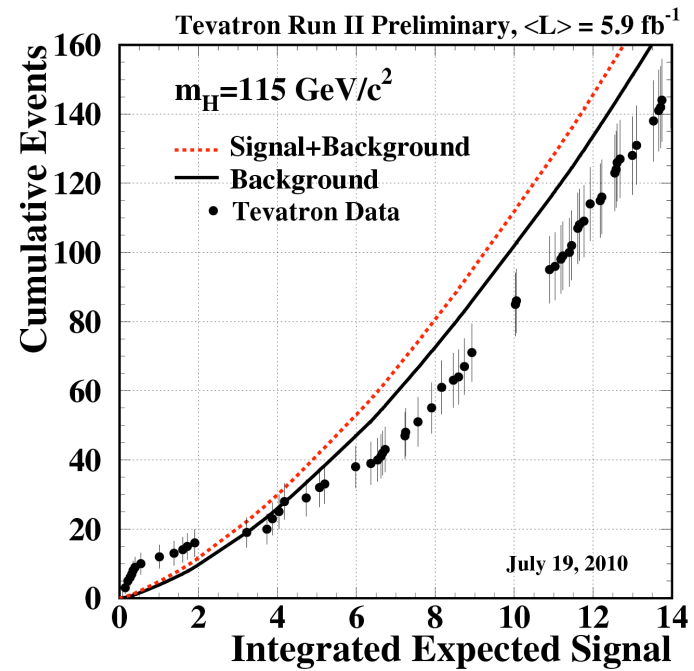
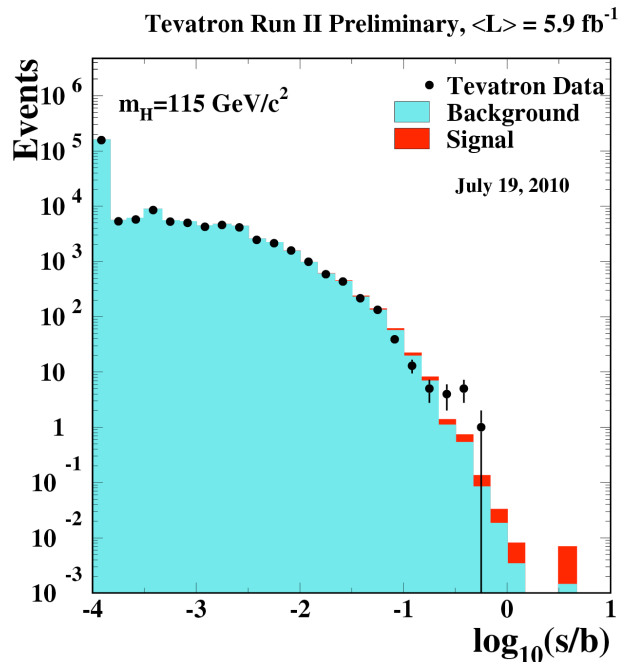
- Does not exclude where there is no sensitivity (big enough search region with small enough resolution and you get a 5% dusting of random exclusions with  $CL_{s+b}$ )



# What These Look Like for a $5.0\sigma$ Observation



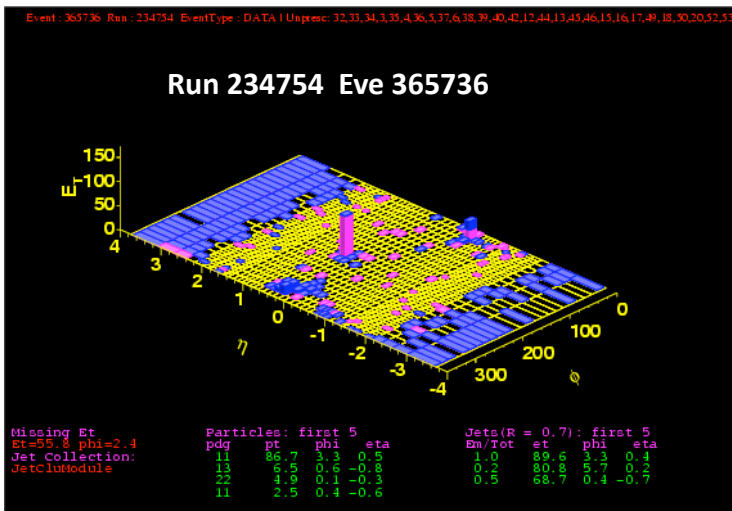
CDF Single Top,  $3.2 \text{ fb}^{-1}$



1. J. J. Stat. Tools for Higgs Combination

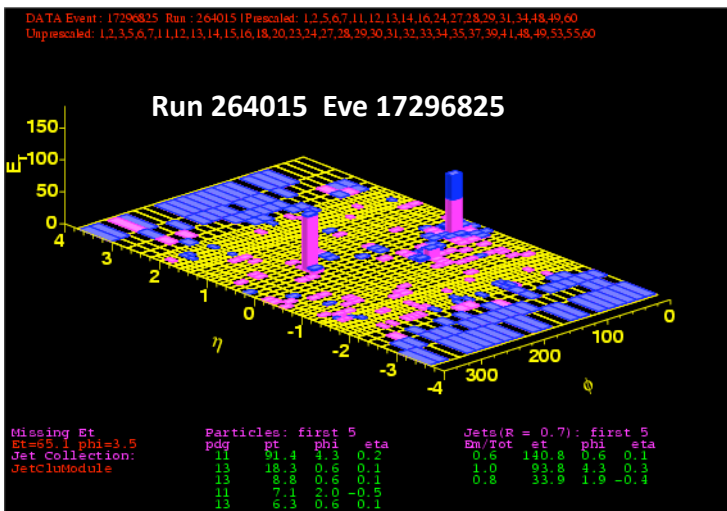


# The Top Five Events at $m_H=115$ GeV



This was found as the most Higgs like event at 1.9/fb analysis (event display was blessed).

Our BNN re-find this event!!



Event : 365736 Run : 234754 EventType : DATA | Unpresc: 32,33,34,3,35,4,36,5,37,6,38,39,40,42,12,44,13,45,46,15,16,17,49,18,30,20,32,35

```
Missing Et
Et=55.8 phi=2.4

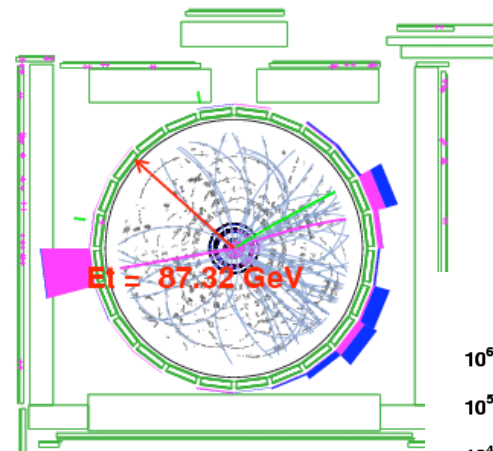
List of Tracks
id pt phi eta

CDF Tracks: first 5
372 -85.7 -3.0 0.5
387 -10.8 -0.5 0.3
351 10.1 -0.9 0.0
352 8.3 -0.5 0.2
403 7.1 -0.5 1.7

To select track type
SelectCdTrack(id)

Dvt Tracks: first 5
4 -180.8 3.3
7 -10.5 5.7
0 -5.7 0.5
3 5.2 0.4
8 5.0 5.8

To select track type
SelectDvtTrack(id)
```

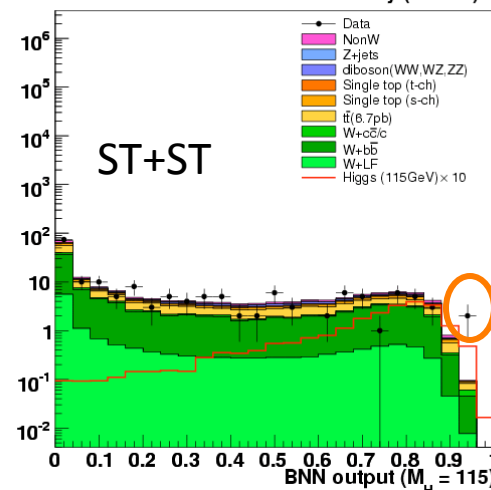


```
Particle
pdg
11
13
22
4.9 0.1 -0.3
11 2.5 0.4 -0.6
To list all particles
ListCdParticles()

Jets(R = 0.7): first
Em/Tot et phi et
1.0 85.5 3.3 0.
0.2 80.8 5.7 0.
... ..
```

Two in WH NN,  
Summer 2009

CDF Run II Preliminary (4.3 fb<sup>-1</sup>)



DATA Event : 17296825 Run : 264015 | Prescaled: 1,2,5,6,7,11,12,13,14,16,24,27,28,29;  
Unprescaled: 1,2,3,5,6,7,11,12,13,14,15,16,18,20,23,24,27,28,29,30,31,32,33,34,35,37,39,

```
Missing Et
Et=65.1 phi=3.5

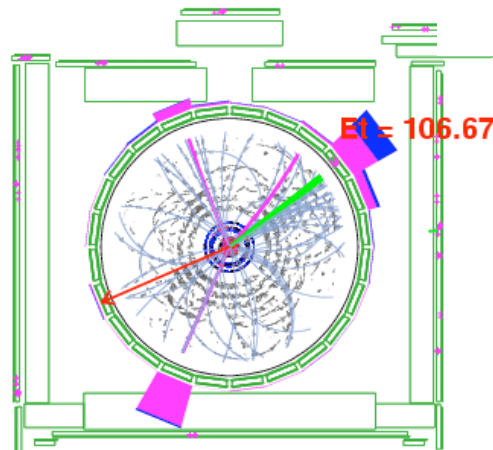
List of Tracks
id pt phi eta

CDF Tracks: first 5
451 84.2 -2.0 0.2
454 84.2 -2.0 0.2
455 -93.7 -2.0 0.2
452 -93.7 -2.0 0.2
386 81.4 -2.0 0.2

To select track type
SelectCdTrack(id)

Dvt Tracks: first 5
0 80.4 4.3
3 30.1 0.5
4 20.1 0.5
5 18.1 0.5
1 7.5 1.8

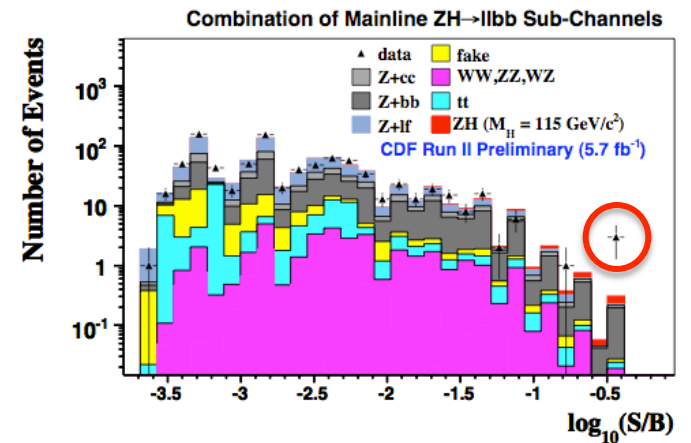
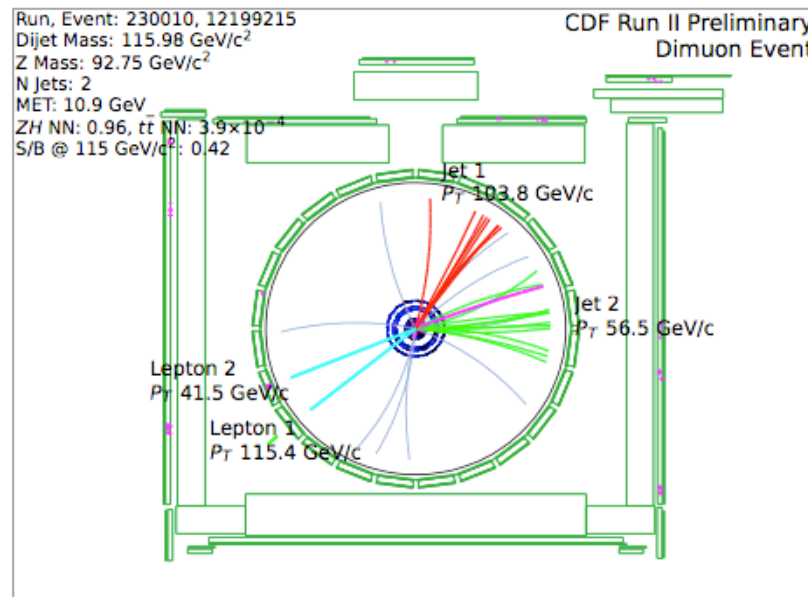
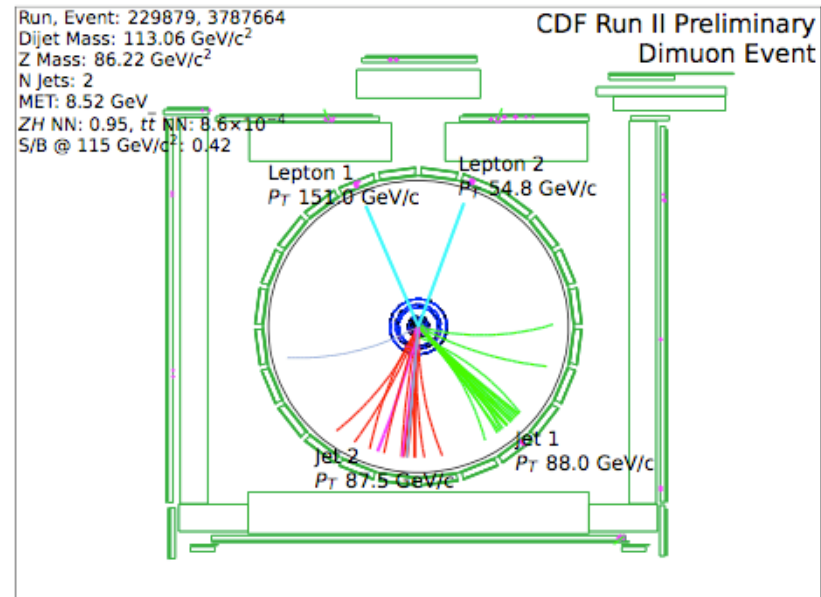
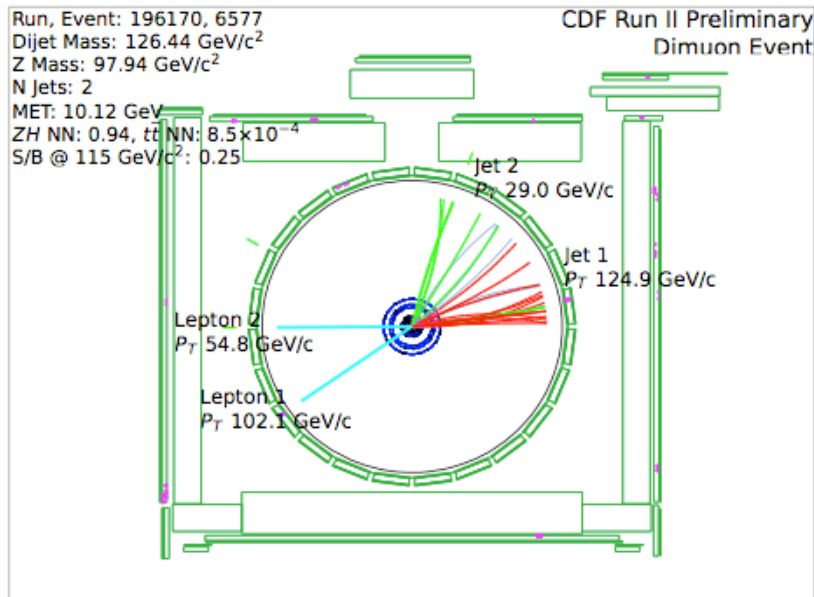
To select track type
SelectDvtTrack(id)
```



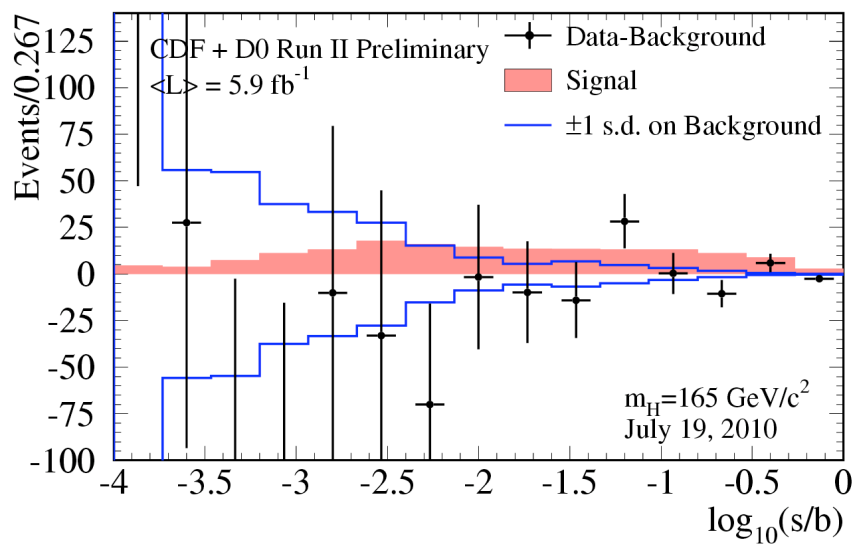
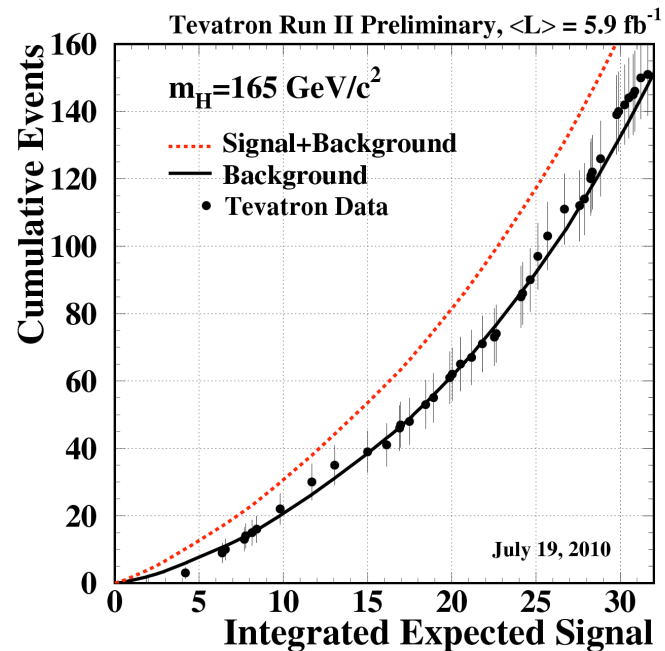
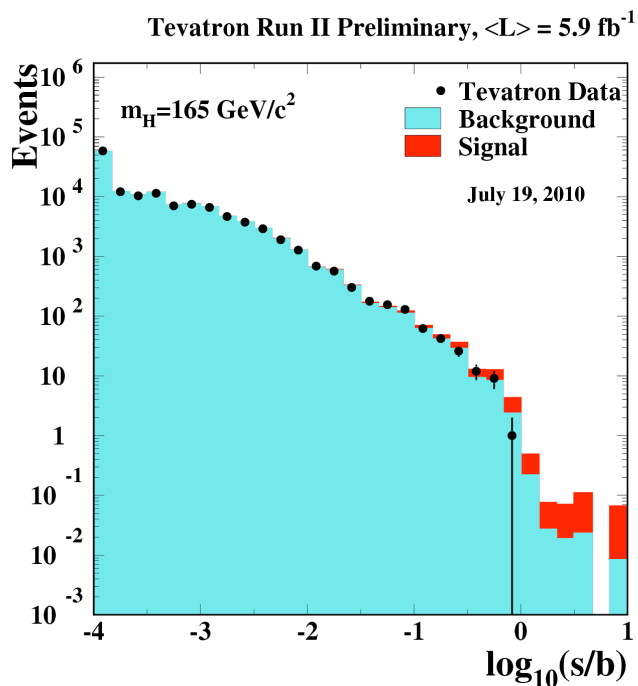
```
Jet
pdg pt pna wca
11 91.4 4.3 0.2
13 18.3 0.6 0.1
11 8.8 0.6 0.1
11 7.1 2.0 -0.5
13 6.3 0.6 0.1
To list all particles
ListCdParticles()

Jets(R = 0.7): first
Em/Tot et phi et
0.5 140.8 0.6 0.
1.0 93.8 4.3 0.
0.8 33.9 1.9 -0.4
To list all jets
ListCdJets()
```

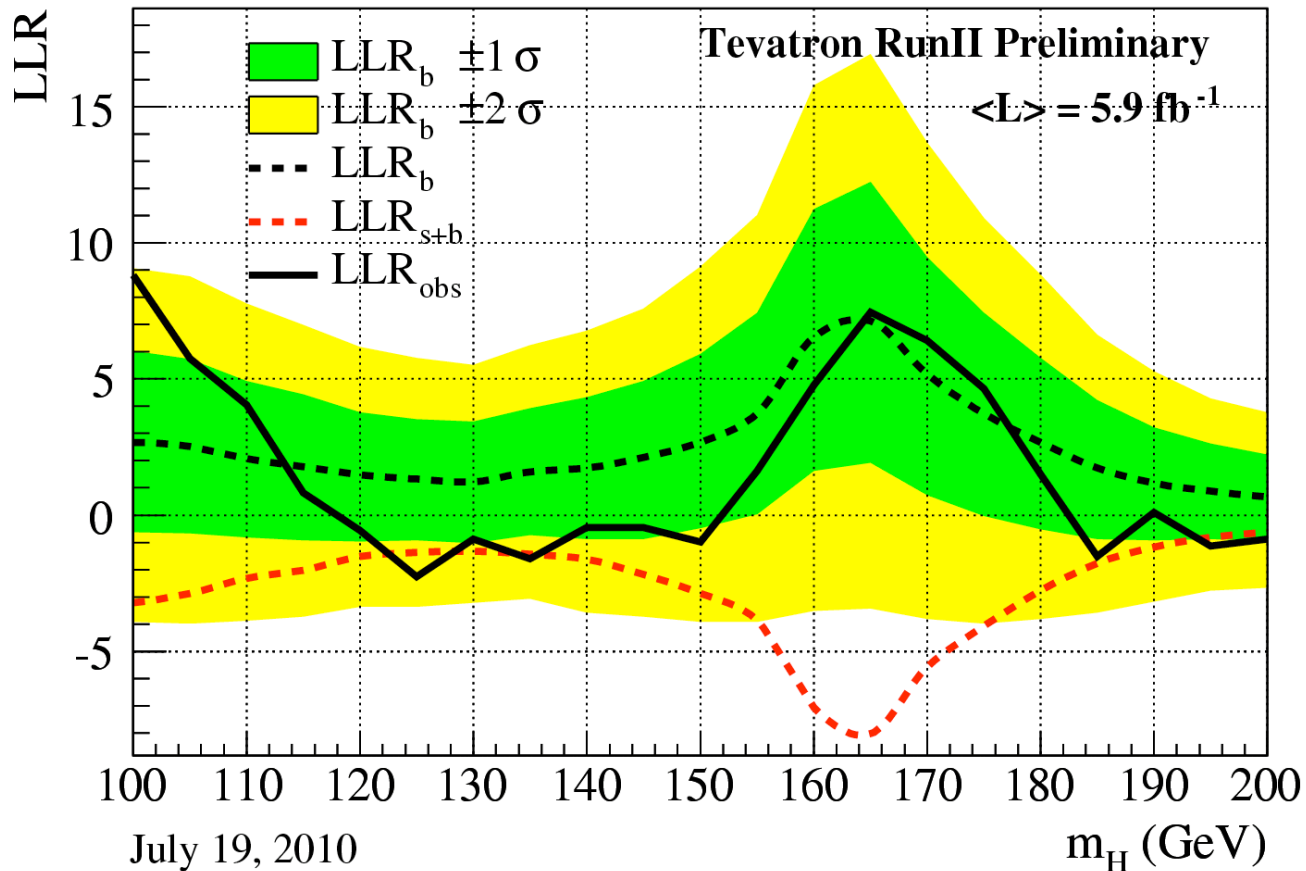
# Three High s/b Candidates in the llbb Channels



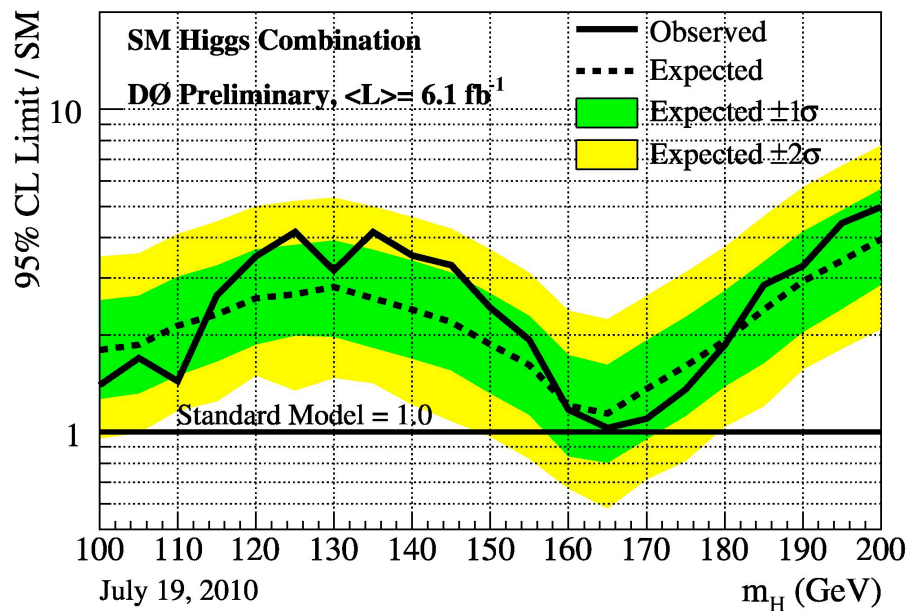
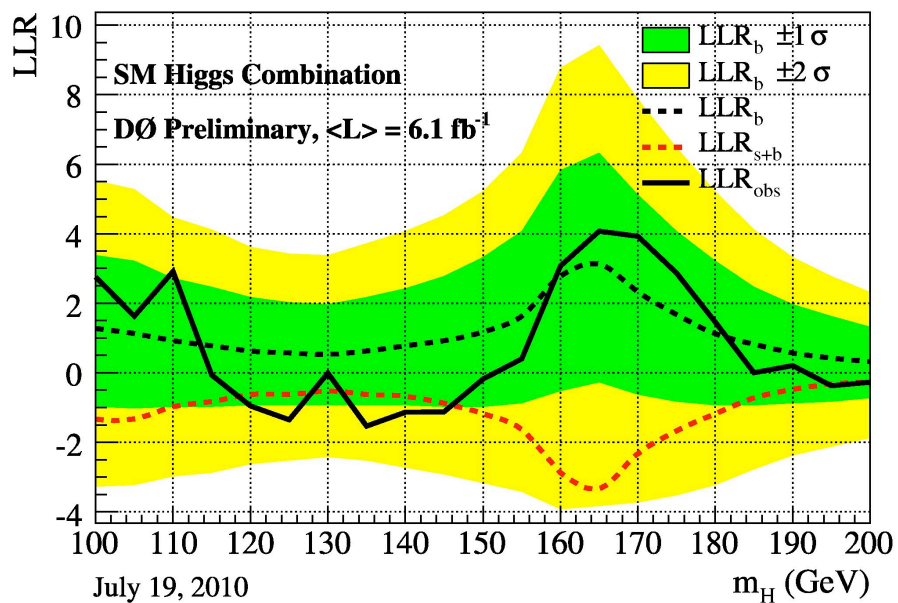
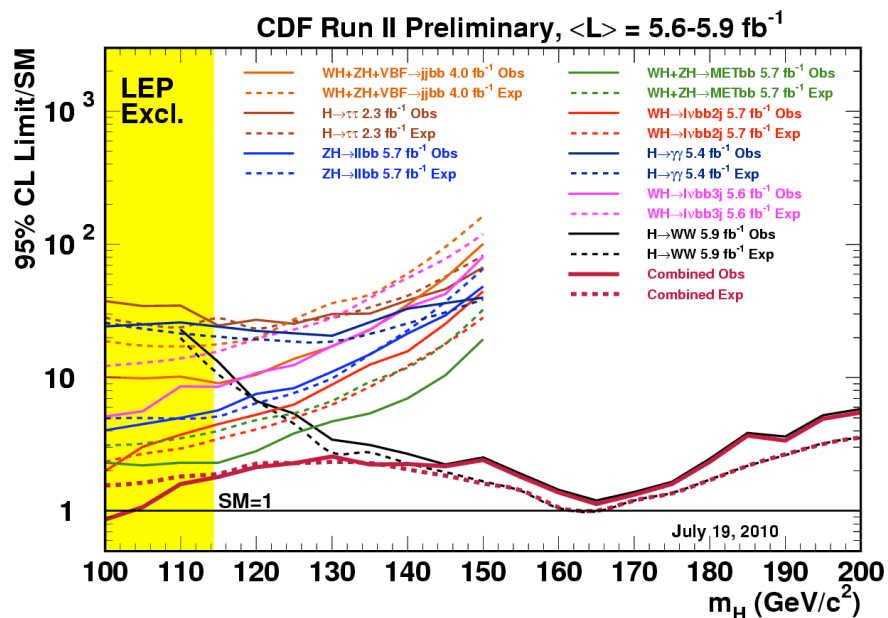
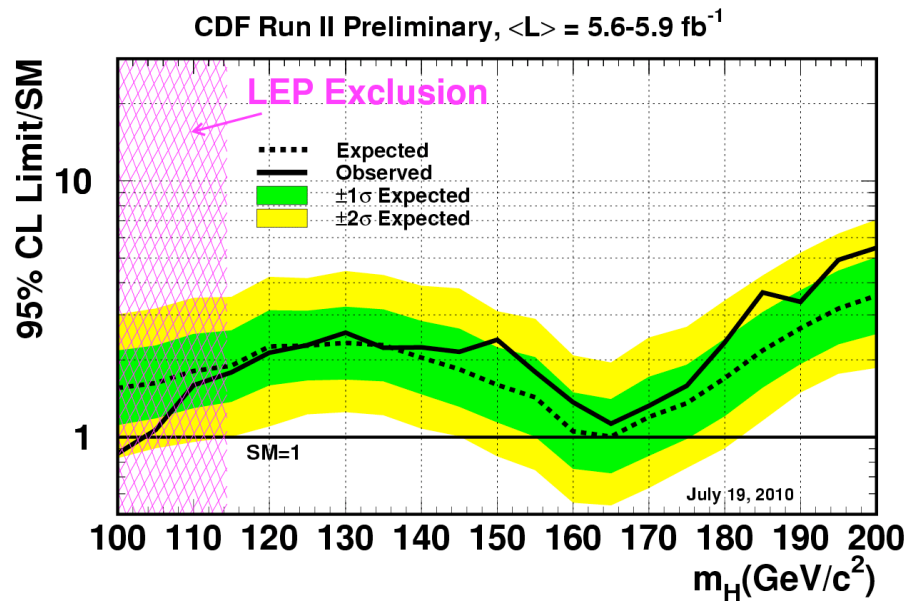
s/b of each of these events is 0.3 to 0.4

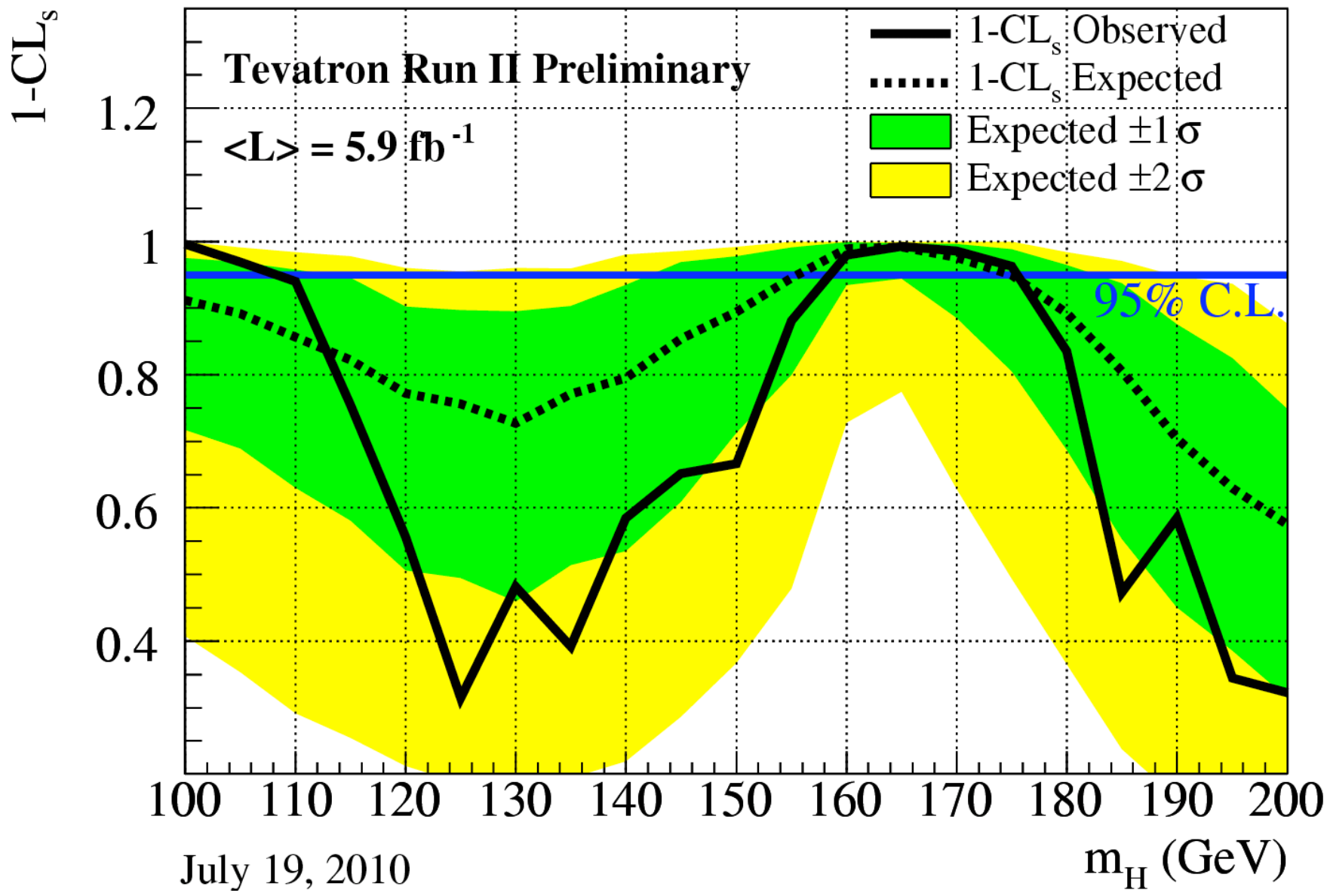


# Looking for a Signal



$$-2 \ln Q \equiv LLR \equiv -2 \ln \left( \frac{L(\text{data} | s + b, \hat{\theta})}{L(\text{data} | b, \hat{\theta})} \right)$$

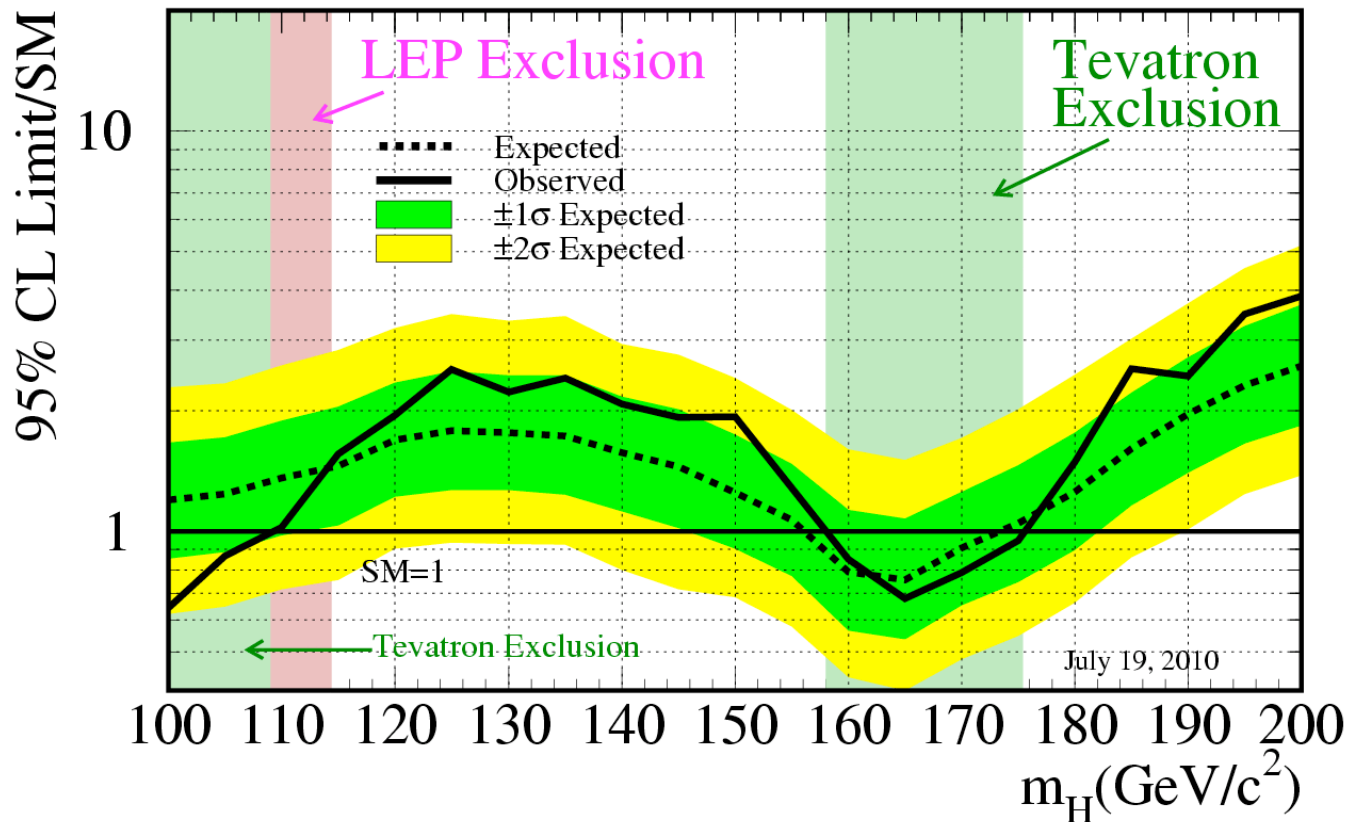




# Tevatron Observed and Expected Limits

Bayesian

Tevatron Run II Preliminary,  $L \leq 6.7 \text{ fb}^{-1}$



Excluded regions:

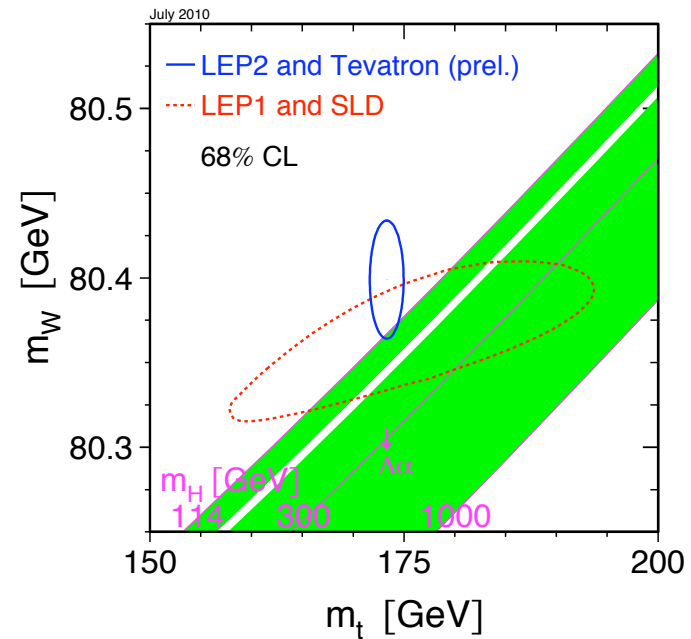
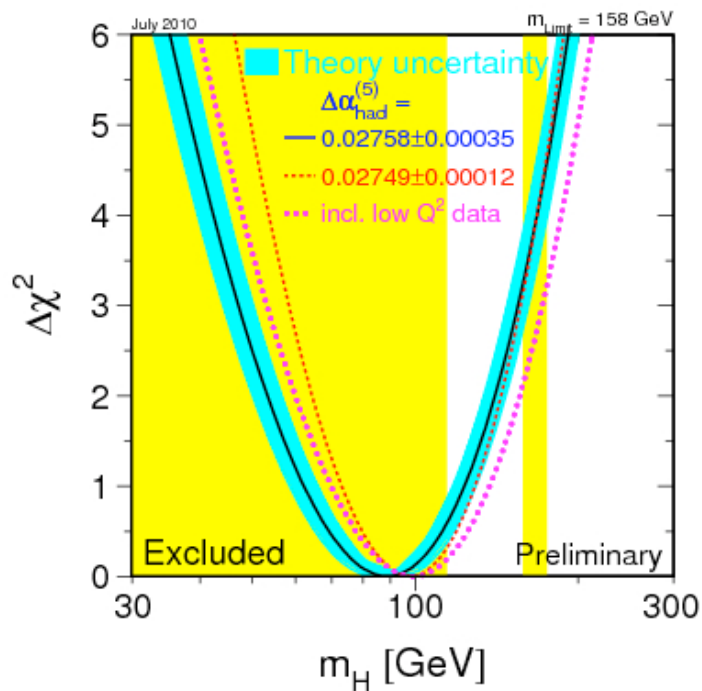
$$158 < m_H < 175 \text{ GeV}$$

$$100 < m_H < 109 \text{ GeV}$$

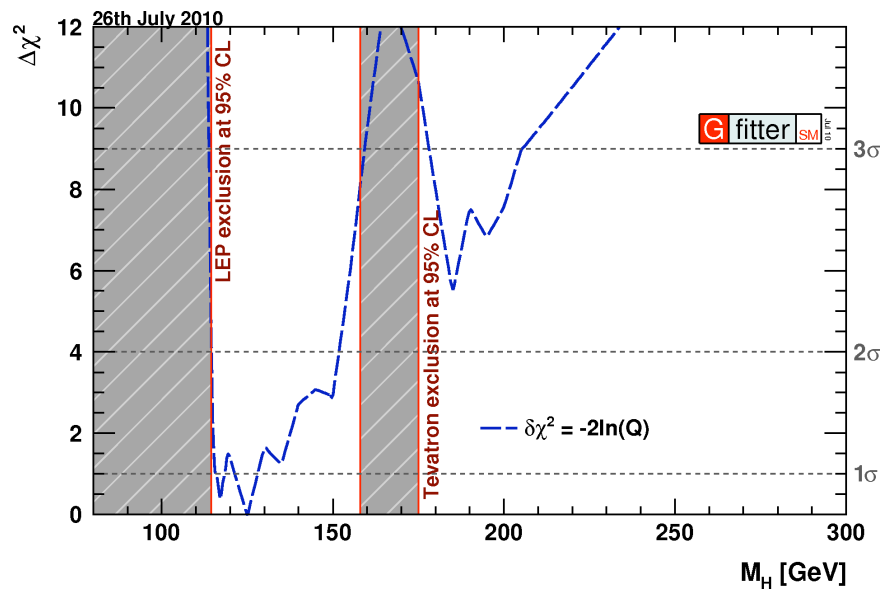
Expected Exclusion  
(if no signal is present):

$$156 < m_H < 173 \text{ GeV}$$

# We have a large impact on the Global Fits

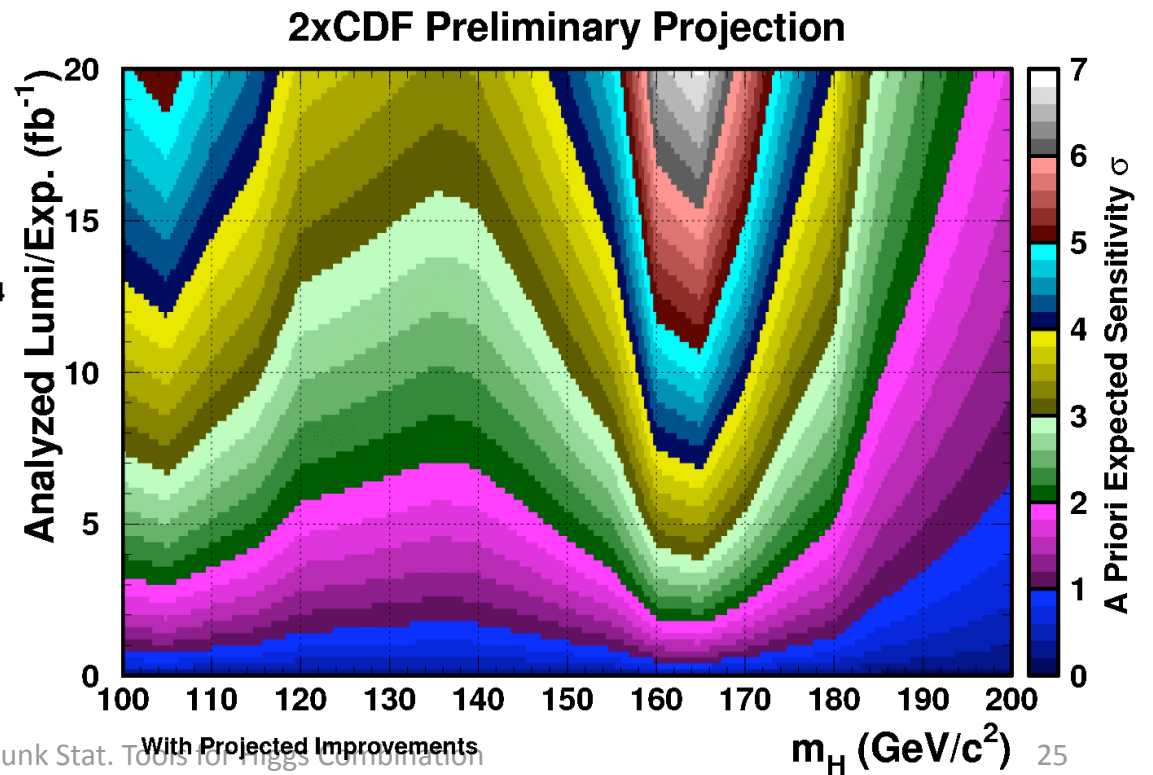
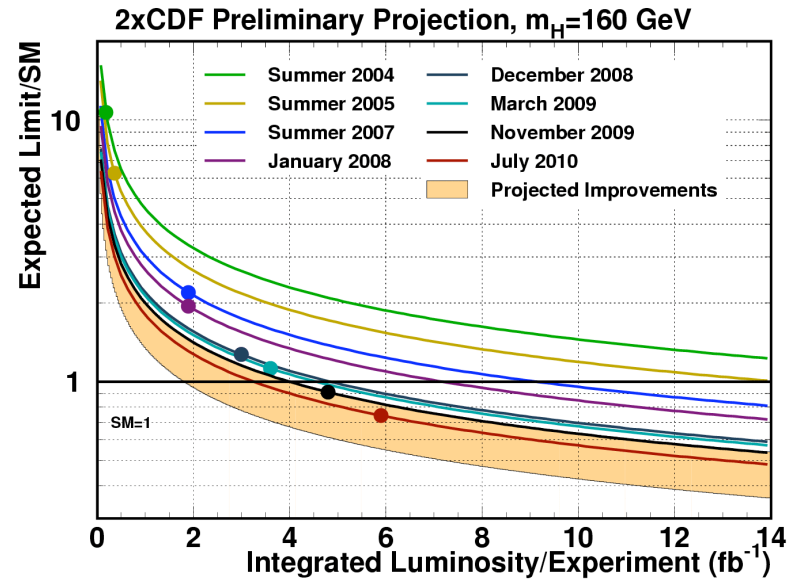
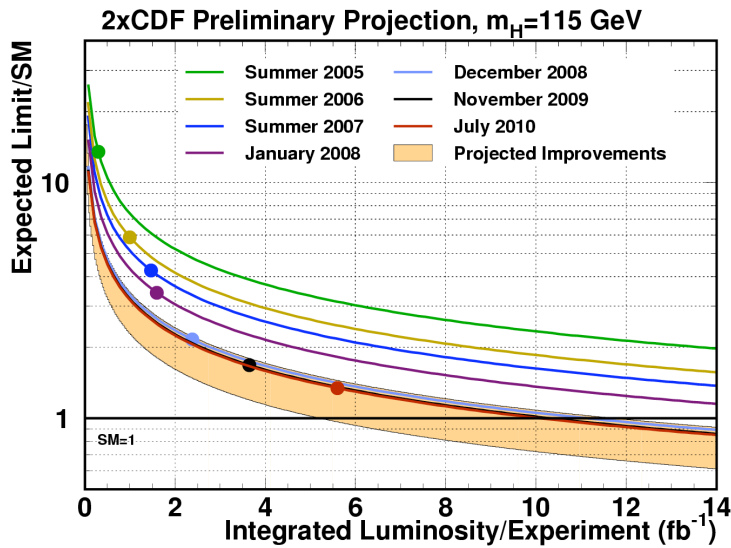


Gfitter Collaboration:



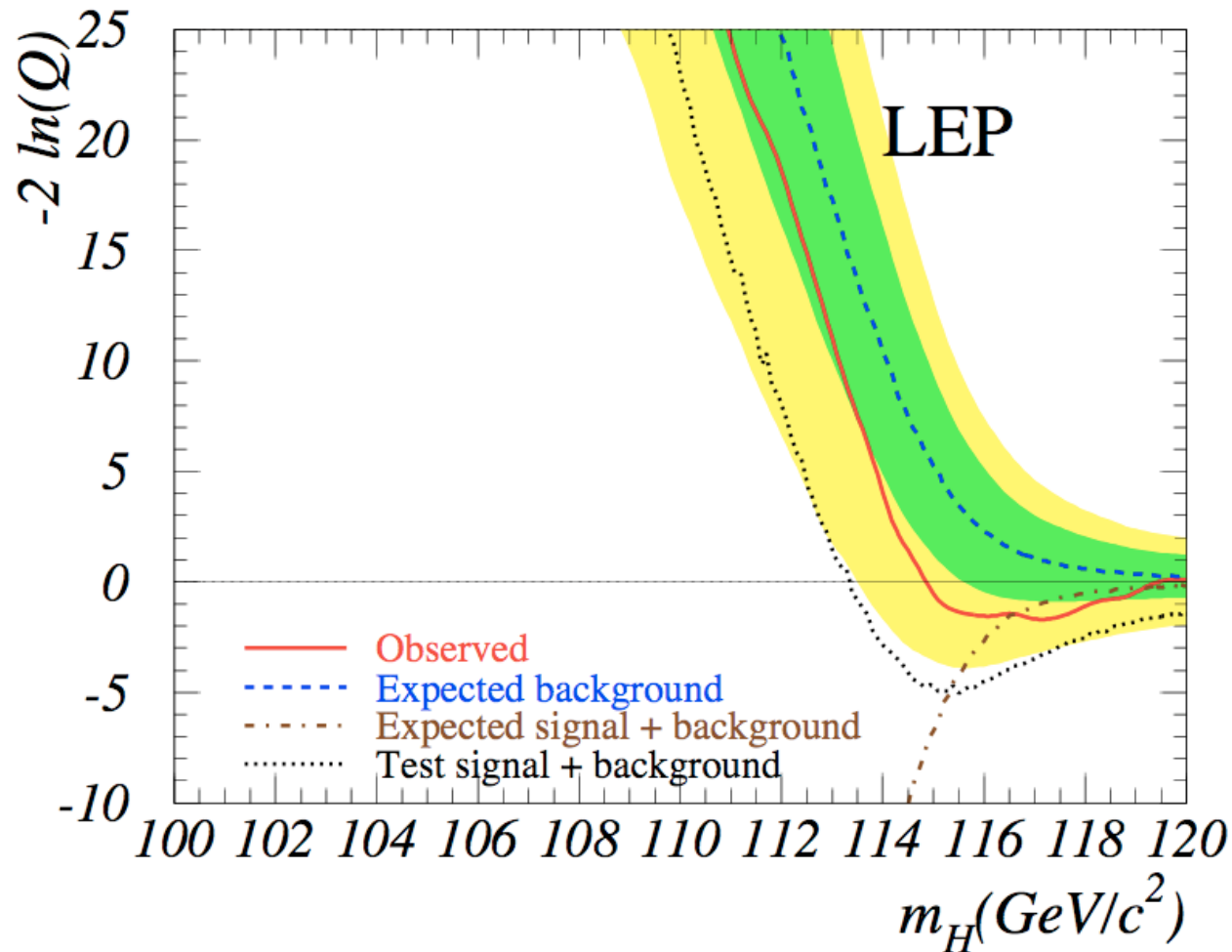


# Projections for Future Performance



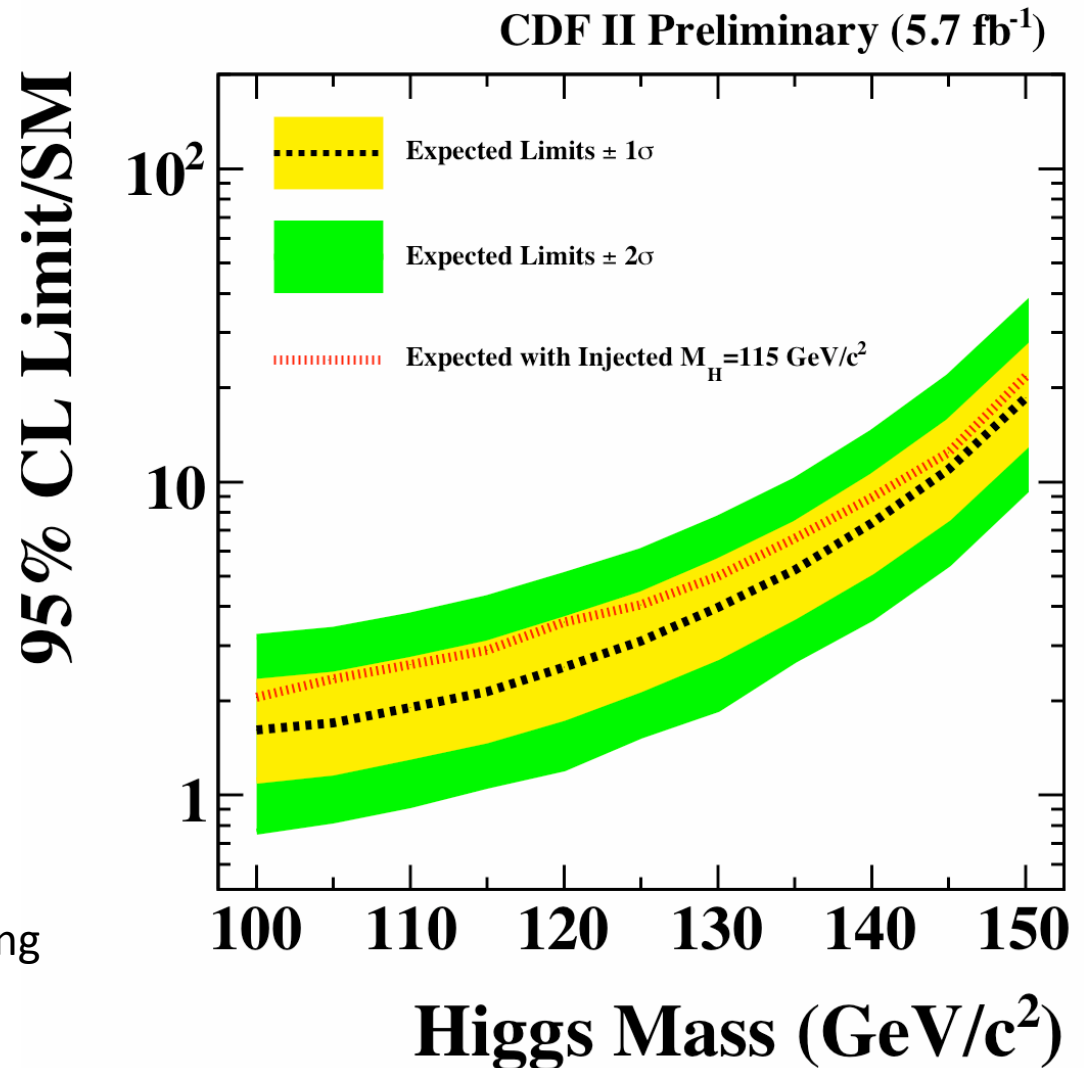
# Look-Aside Histograms

What does a signal with  $m_H = m_1$  look like when seeking  $m_H = m_2$ ?  
So far, not done at Tevatron. Not needed to study the trials factor, but needed to make this plot:



# Studies of Injecting a Signal at $m_H=115$ GeV

- $lvbb$ ,  $METbb$ , and  $llbb$  channels included
- Inject  $SM \cdot 1.0$  signal at  $m_H=115$  GeV on top of SM backgrounds, and generate pseudoexperiments with that.
- Analyze 115 signal+background pseudoexperiments at other test masses – 100 GeV to 150 GeV
- Find the median expected limit assuming signal is there (compute it just as you would without the signal) and compare with the distribution of limits assuming the signal is completely absent.



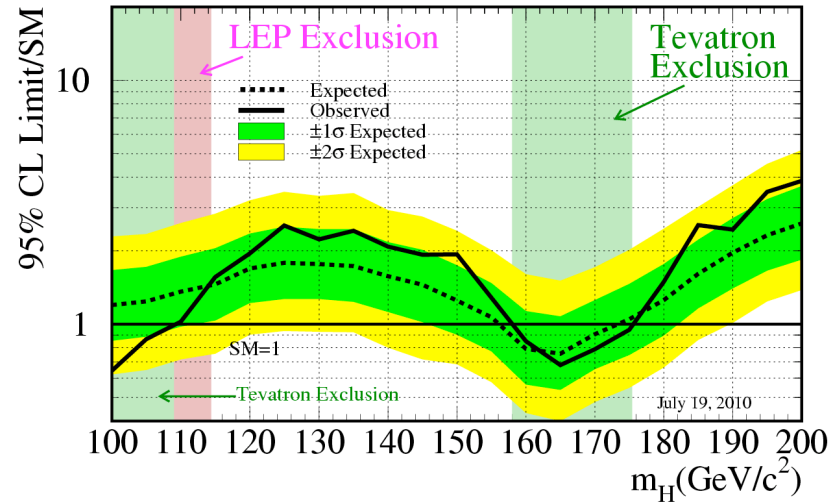
# Summary

The Tevatron is doing very well!

We will run at least through 2011

We are asking to run another 3 years.

Tevatron Run II Preliminary,  $L \leq 6.7 \text{ fb}^{-1}$



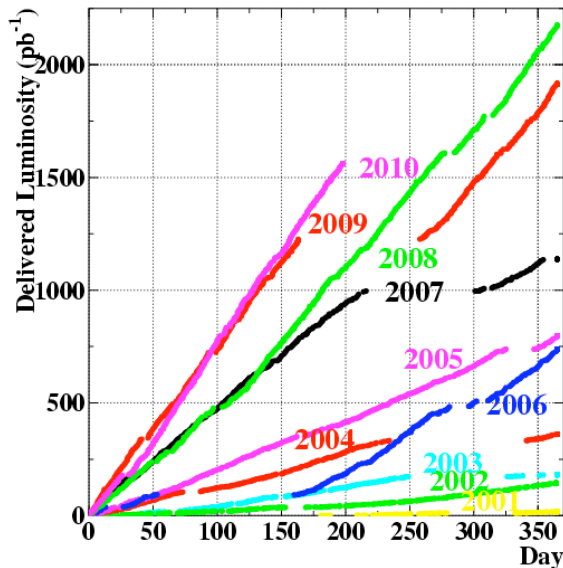
Excluded regions:

$$158 < m_H < 175 \text{ GeV}$$

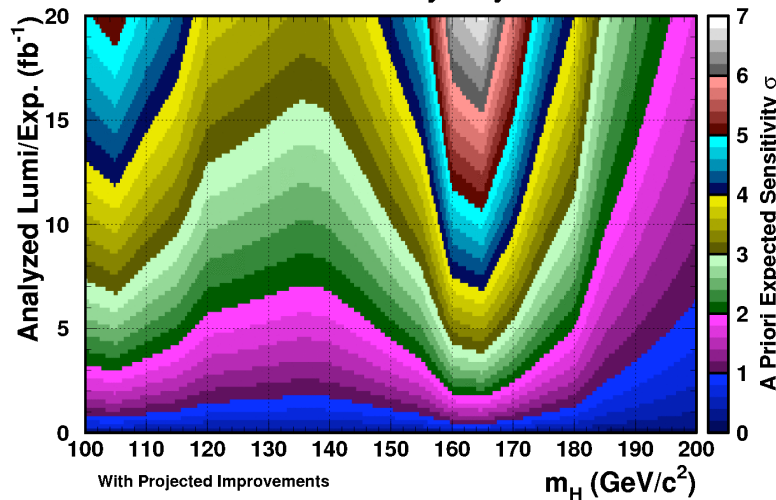
$$100 < m_H < 109 \text{ GeV}$$

Expected Exclusion (if no signal is present):

$$156 < m_H < 173 \text{ GeV}$$



2xCDF Preliminary Projection



# Backup Material

## Commonly Used Tools for Setting Limits and Discovering New Processes in use at the Tevatron

- Bayesian limits -- common at CDF
  - genlimit code by Joel Heinrich, added to mclimit code by Tom Junk
  - Implements posterior integrated over systematic uncertainties with a flat prior in cross section in 1D
  - Method described in PDG statistics review
  - Extra feature -- “correlated prior”
- $CL_s$  limits -- common at D0, but used at CDF as well.
  - Collie code by Wade Fisher in use at D0
  - Method described in PDG statistics review
  - mclimit was originally designed to do  $CL_s$  and still does.
  - TLimit in ROOT is out of date -- no fits for nuisance parameters, no shape errors or bin-by-bin errors

# Measurement and Discovery are Very Different

Buzzwords:

- Measurement = “Point Estimation”
- Discovery = “Hypothesis Testing”

You can have a discovery and a poor measurement!

Example: Expected  $b=2 \times 10^{-7}$  events, expected signal=1 event, observe 1 event, no systematics.

p-value  $\sim 2 \times 10^{-7}$  is a discovery! (hard to explain that event with just the background model). But have  $\pm 100\%$  uncertainty on the measured cross section!

In a one-bin search, all test statistics are equivalent. But add in a second bin, and the measured cross section becomes a poorer test statistic than the ratio of profile likelihoods.

In all practicality, discriminant distributions have a wide spectrum of  $s/b$ , even in the same histogram. But some good bins with  $b < 1$  event

# Sociological Issues

- Discovery is conventionally  $5\sigma$ . In a Gaussian asymptotic case, that would correspond to a  $\pm 20\%$  measurement.
- Less precise measurements are called “measurements” all the time
- We are used to measuring undiscovered particles and processes. In the case of a background-dominated search, it can take years to climb up the sensitivity curve and get an observation, while evidence, measurements, etc. proceed.
- Referees can be confused.



# The Trials Factor

- Also called the “Look Elsewhere Effect”
- Bump-hunters are familiar with it.

What is the probability of an upward fluctuation as big as the one I saw *anywhere* in my histogram?

- Lots of bins → Lots of chances at a false discovery
- Approximation: Multiply smallest p-value by the number of “independent” models sought (not histogram bins!).

Bump hunters: roughly (histogram width)/(mass resolution)

Criticisms:

Adjusted p-value can now exceed unity!

What if histogram bins are empty?

What if we seek things that have been ruled out already?

It's not bins, but the number of independent hypotheses being tested that matters!

Low mass resolution at the Tevatron – not very many independent excesses possible.

# The Trials Factor

More seriously, what to do if the p-value comes from a big combination of many channels each optimized at each  $m_H$  sought?

- Channels have different resolutions (or is resolution even the right word for a multivariate discriminant?)
- Channels vary their weight in the combination as cross sections and branching ratios change with  $m_H$

Proper treatment -- want a p-value of p-values!

(use the p-value as a test statistic)

Run pseudoexperiments and analyze each one at each  $m_H$  studied. Look for the distribution of smallest p-values.

Difficult but possible.