

# Statistics & Analysis

## Q&A

### PHYS428/PHYS576 **Advanced Techniques in Experimental Particle Physics**

Fred James's lectures

[http://preprints.cern.ch/cgi-bin/setlink?base=AT&categ=Academic\\_Training&id=AT00000799](http://preprints.cern.ch/cgi-bin/setlink?base=AT&categ=Academic_Training&id=AT00000799)

<http://www.desy.de/~acatrain/>

Glen Cowan's lectures

[http://www.pp.rhul.ac.uk/~cowan/stat\\_cern.html](http://www.pp.rhul.ac.uk/~cowan/stat_cern.html)

Louis Lyons

<http://indico.cern.ch/conferenceDisplay.py?confId=a063350>

Bob Cousins gave a CMS lecture, may give it more publicly Gary Feldman "Journeys of an Accidental Statistician"

<http://www.hepl.harvard.edu/~feldman/Journeys.pdf>

<http://histfitter.web.cern.ch/histfitter/>



# Example: Maximum Likelihood Estimator

example:  $\text{PDF}(x) = \text{Gauss}(x, \mu, \sigma) \rightarrow$

$$L(\mu|x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

→ estimator for  $\mu_{true}$  from the data measured in an experiment  $x_1, \dots, x_N$

→ full likelihood

$$L(\mu|x) = \prod_i^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_i-\mu)^2}{2\sigma^2}\right)$$

→ typically:

$$-2\ln(L(\mu|x)) = \sum_i^N \left( \frac{(x_i-\mu)^2}{2\sigma^2} \right) + N \frac{1}{\sqrt{2\pi}\sigma^2}$$

Note: It's a function of  $\mu$ !

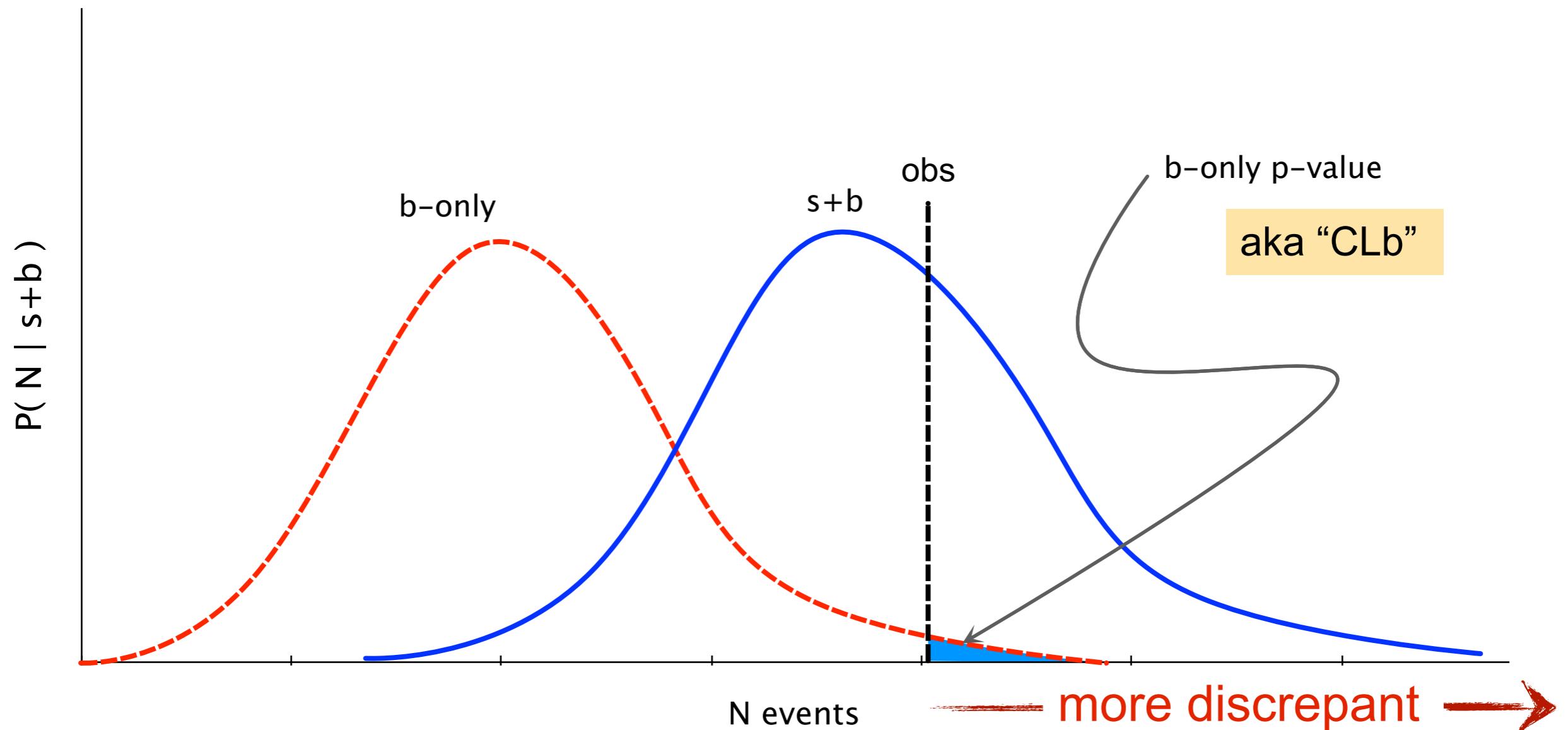
$$-2\Delta\ln(L(\mu)) = \sum_i^N \left( \frac{(x_i-\mu)^2}{2\sigma^2} \right)$$

→  $\chi^2$ , least squares

# Discovery sensitivity: p-value

Discovery: test b-only (null:  $s=0$  vs. alt:  $s>0$ )

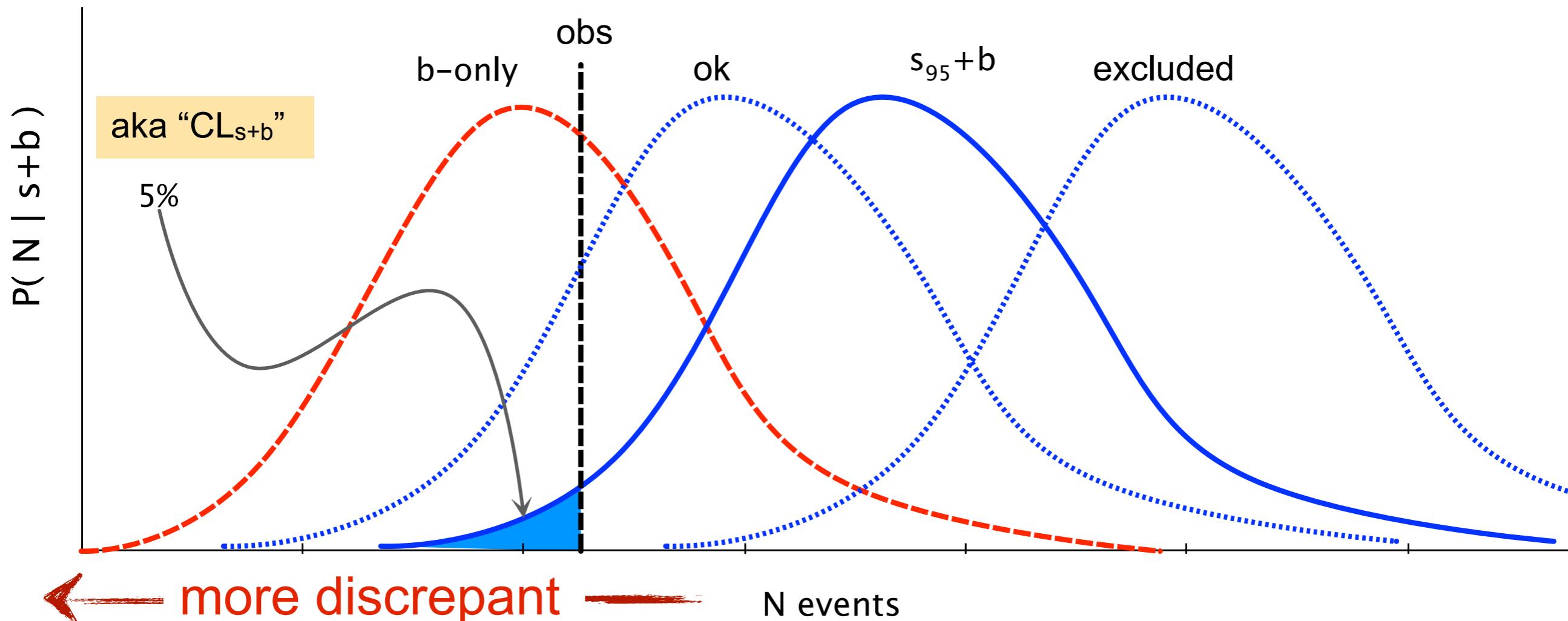
- note, **one-sided** alternative. larger N is “more discrepant”



# Upper Limit: Confidence Level

What is meant by “95% upper limit” ? See the picture below?

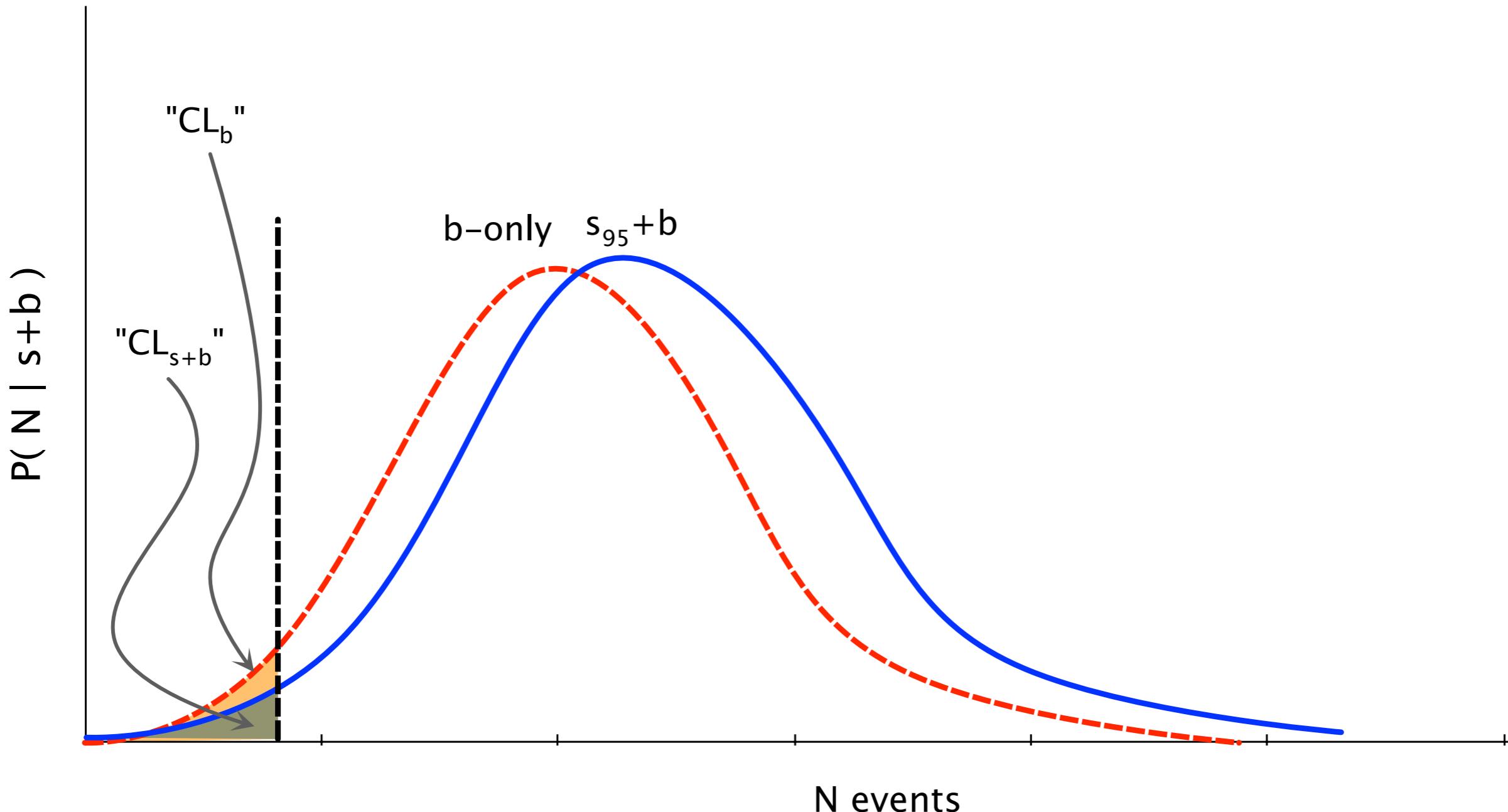
- ie. increase  $s$ , until the probability to have data “more discrepant” is  $< 5\%$



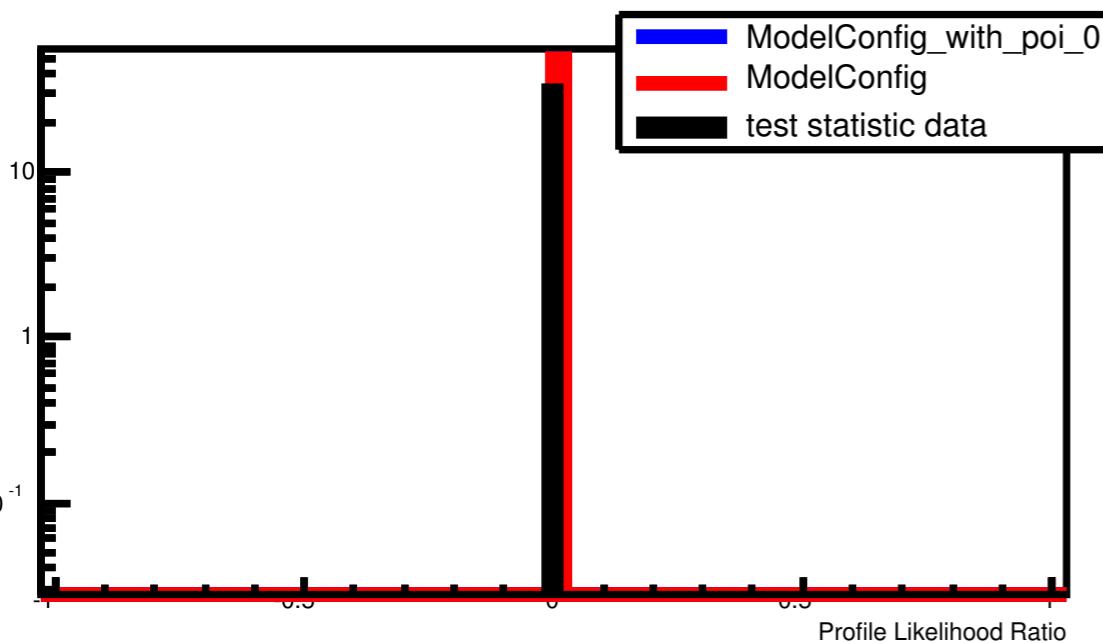
# CL<sub>s</sub> Upper Limit: Confidence Level

CL<sub>s</sub> is known to be “conservative” (over-cover): expected limit covers with 97.5%

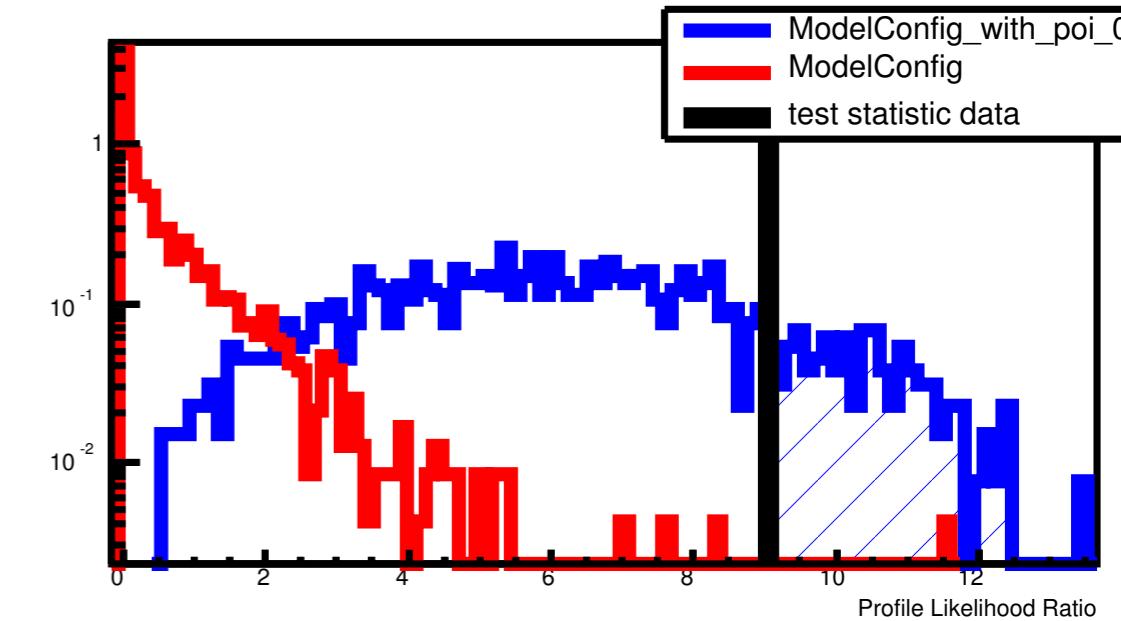
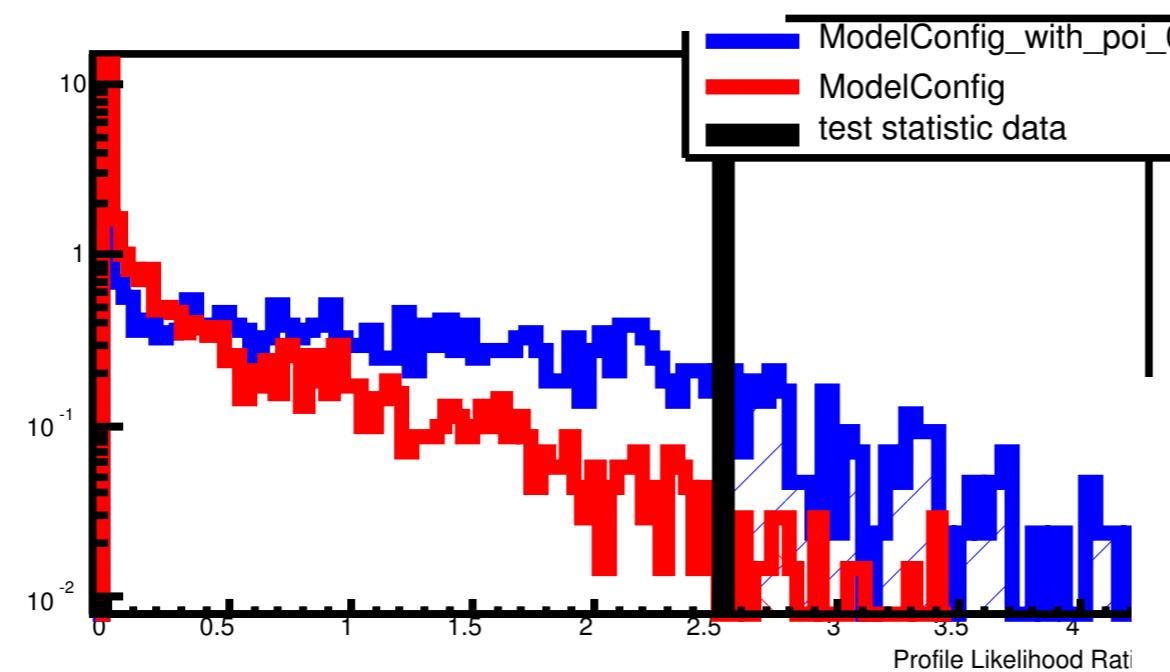
- Note: CL<sub>s</sub> is NOT a probability



# Test Statistics

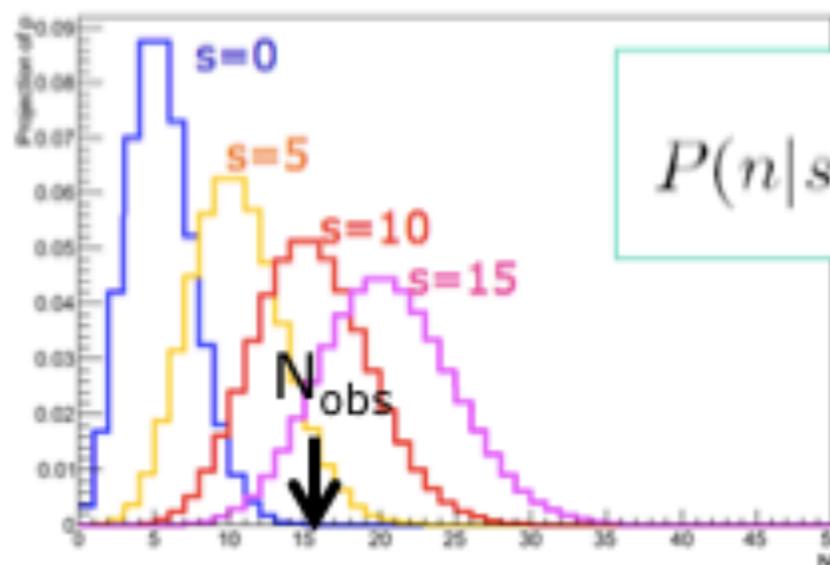


$$q_{\mu_{\text{sig}}} = -2 \log \frac{L(\mu_{\text{sig}}, \hat{\theta})}{L(\hat{\mu}_{\text{sig}}, \hat{\theta})}$$



# Likelihood

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



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

NB:  $b$  is a constant in this example

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

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

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



Frequentist statistics



Bayesian statistics

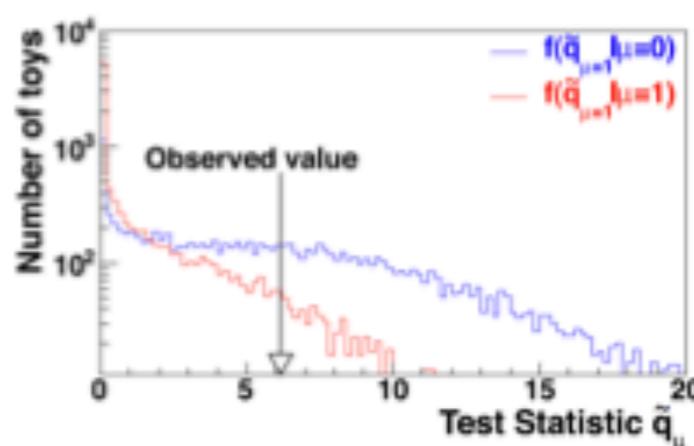


Maximum Likelihood

# Likelihood

Frequentist statistics

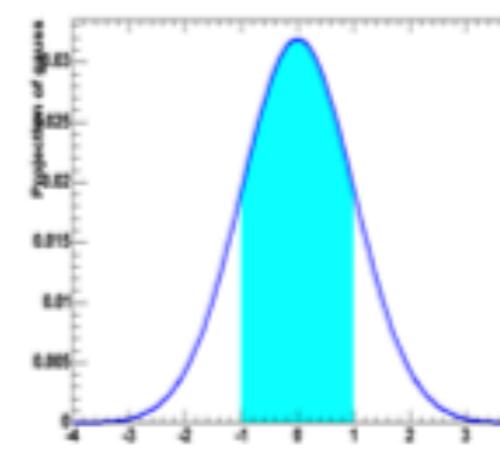
$$\lambda_\mu(\vec{N}_{obs}) = \frac{L(\vec{N} | \mu)}{L(\vec{N} | \hat{\mu})}$$



Confidence interval  
or p-value

Bayesian statistics

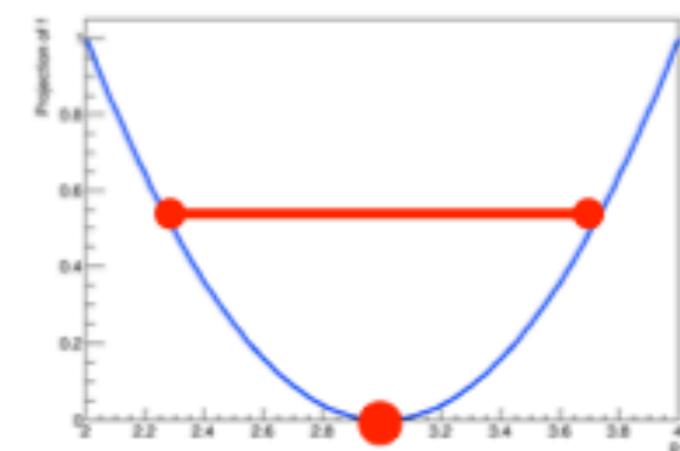
$$P(\mu) \propto L(x | \mu) \cdot \pi(\mu)$$



Posterior on s  
or Bayes factor

Maximum Likelihood

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

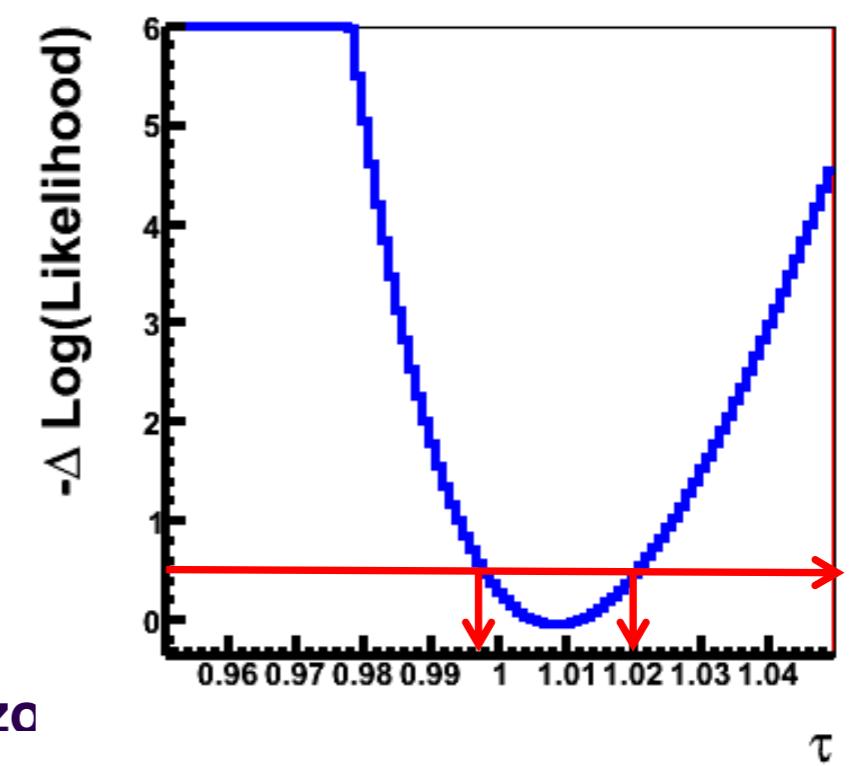
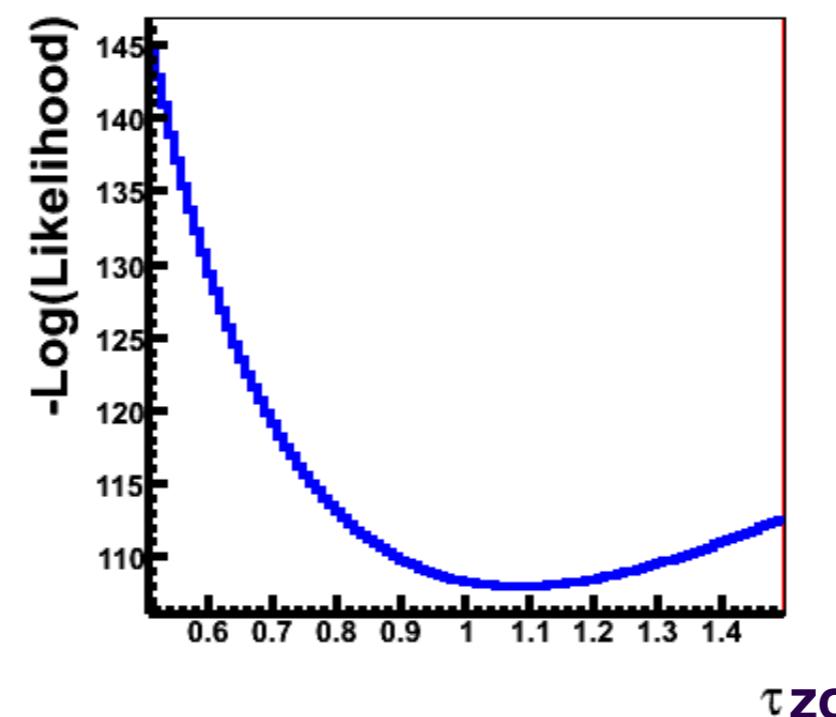
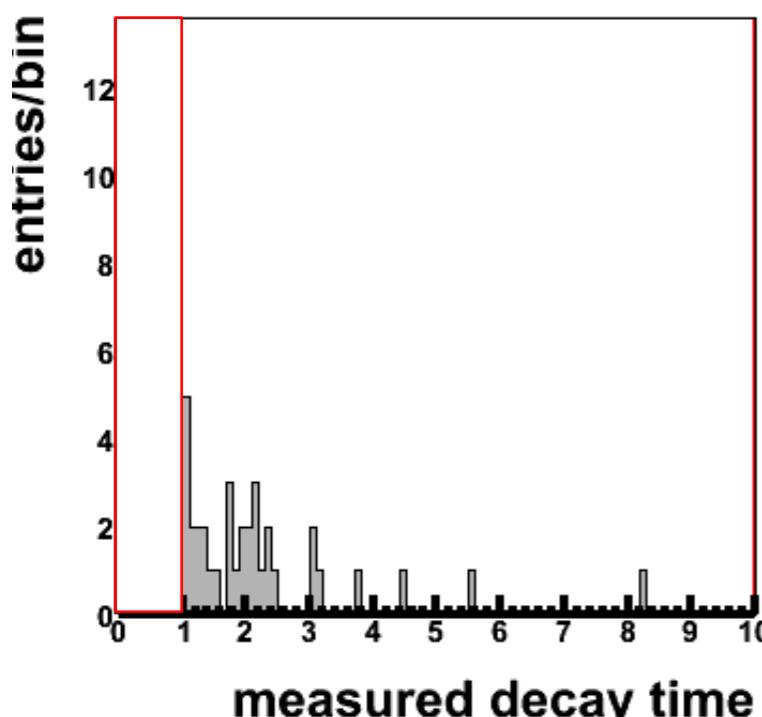


$s = x \pm y$   
Wouter Verkerke, NKHEF

# Maximum Likelihood Estimator

$$-\ln(L(\theta)) \approx -\ln(L(\hat{\theta})) - \underbrace{\left[ \frac{d\ln(L\theta)}{d\theta} \right]_{\hat{\theta}} (\theta - \hat{\theta})}_{\text{minimum}} - \frac{1}{2} \left[ \frac{d^2\ln(L\theta)}{d\theta^2} \right]_{\hat{\theta}} (\theta - \hat{\theta})^2 + ..$$

$$-\ln(L(\theta)) \approx -\ln(L(\hat{\theta})) + \frac{1}{2\sigma^2} (\theta - \hat{\theta})^2$$

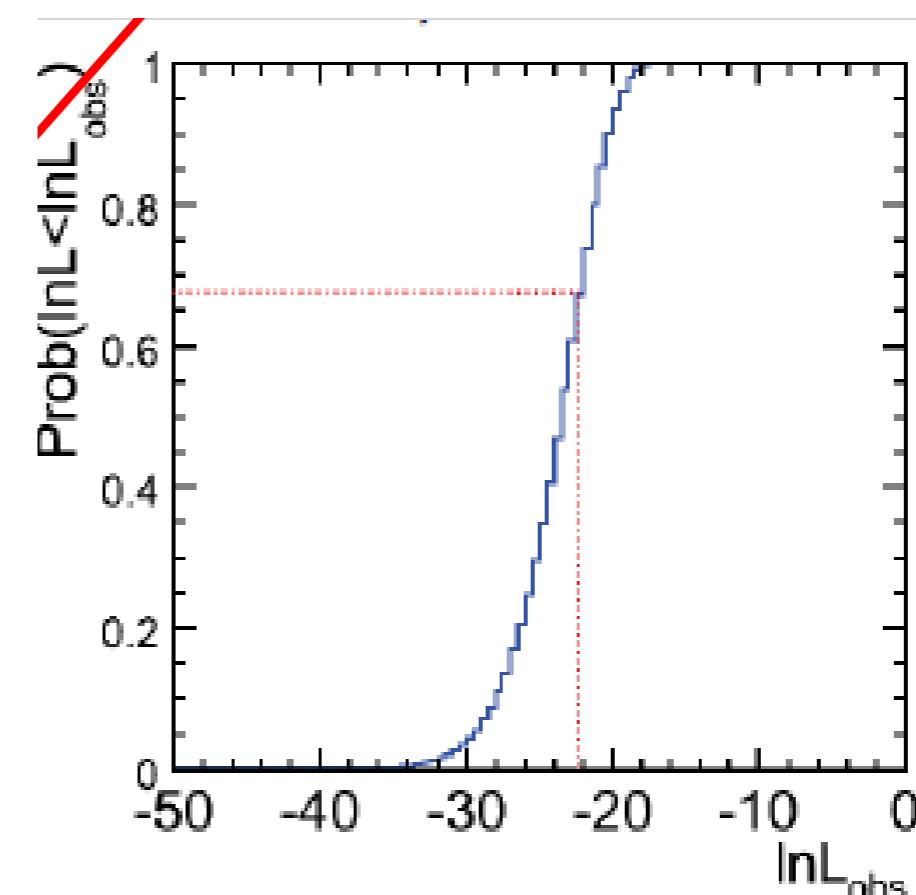
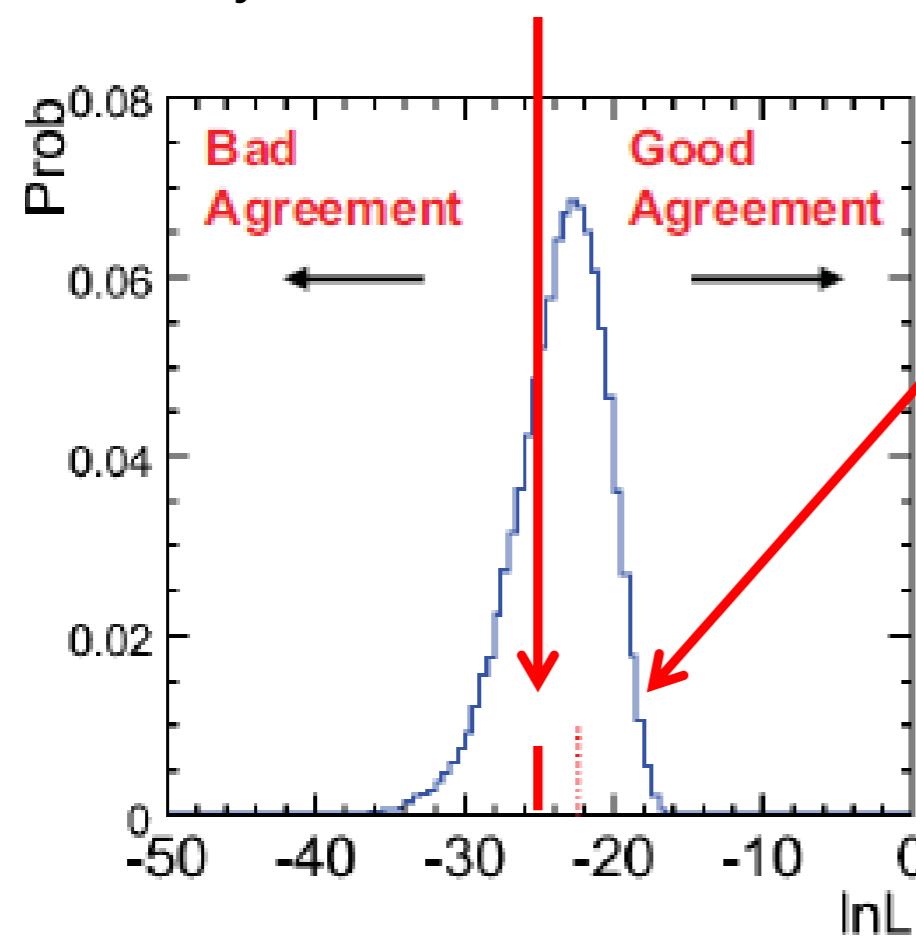


read off parabolic  $\ln(L)$  curve:

$$-\ln(L(\hat{\theta} \pm \sigma_\theta)) = -\ln(L(\hat{\theta})) + \frac{1}{2}$$

# Goodness-of-fit

- So far we know the “uncertainty” on the fitted value of  $\theta$ , but...
- did the fitted model “really” describe the data?
- The value of the  $\ln L$  (log Likelihood) at the minimum does not “mean anything” → **calibrate!**
  - determine the distribution of  $\ln L$  fit results with Monte Carlo toys!
  - check your “data”-fit



# Goodness-of-fit

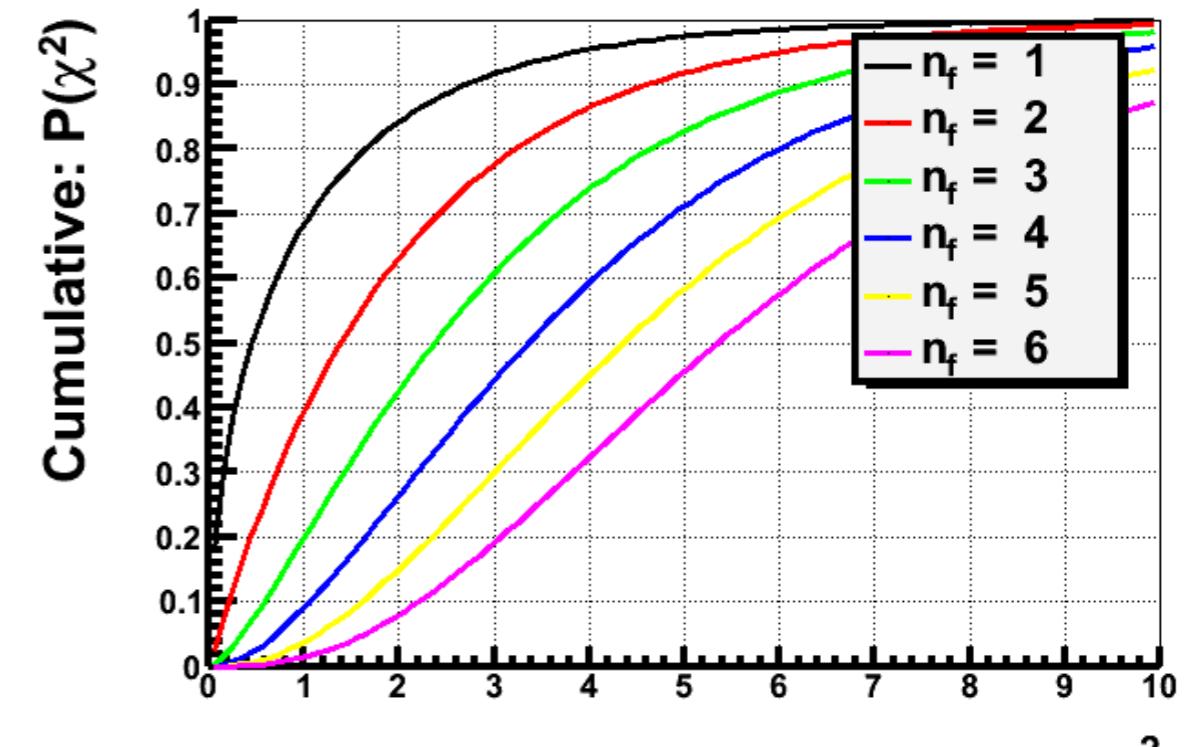
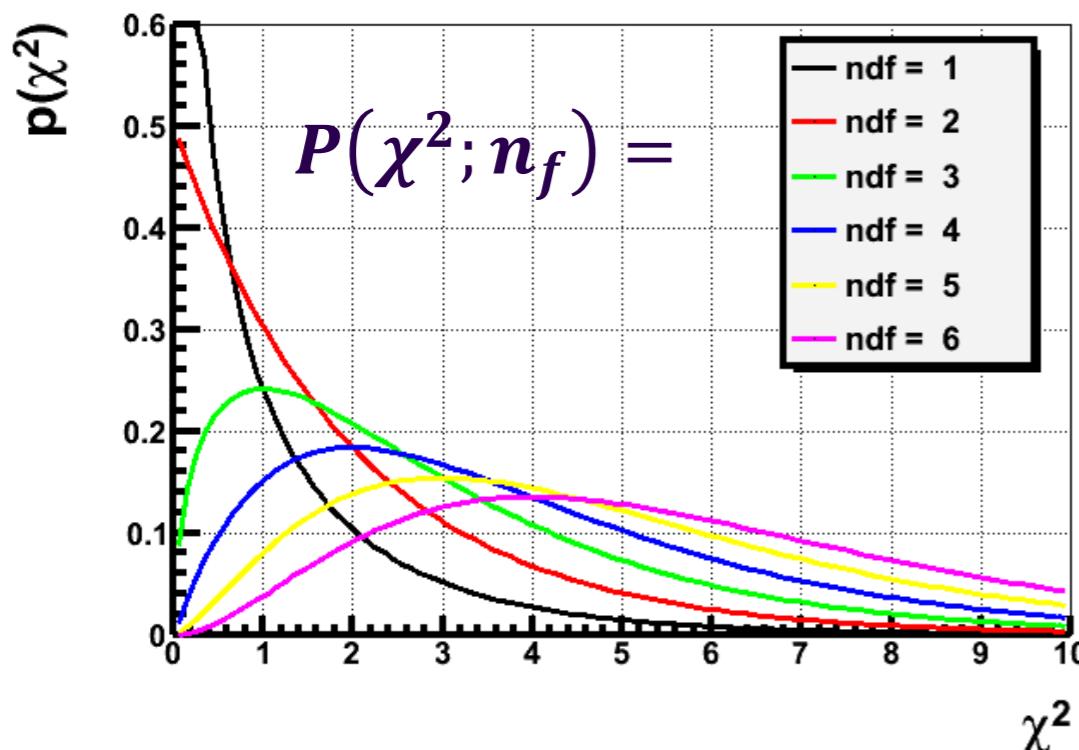
- Easier with Gaussian distributed variables

→ least square fit →  $\chi^2$

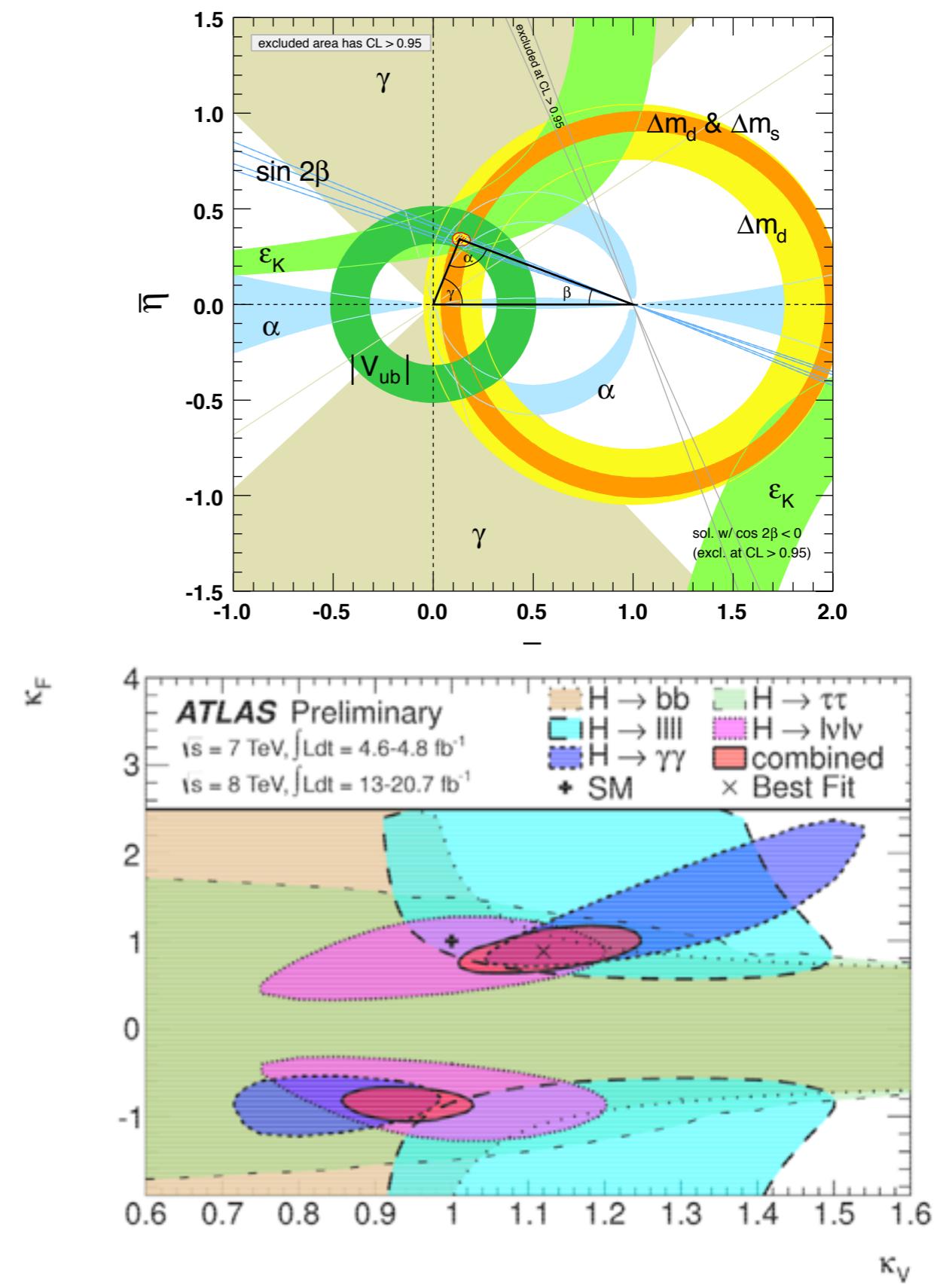
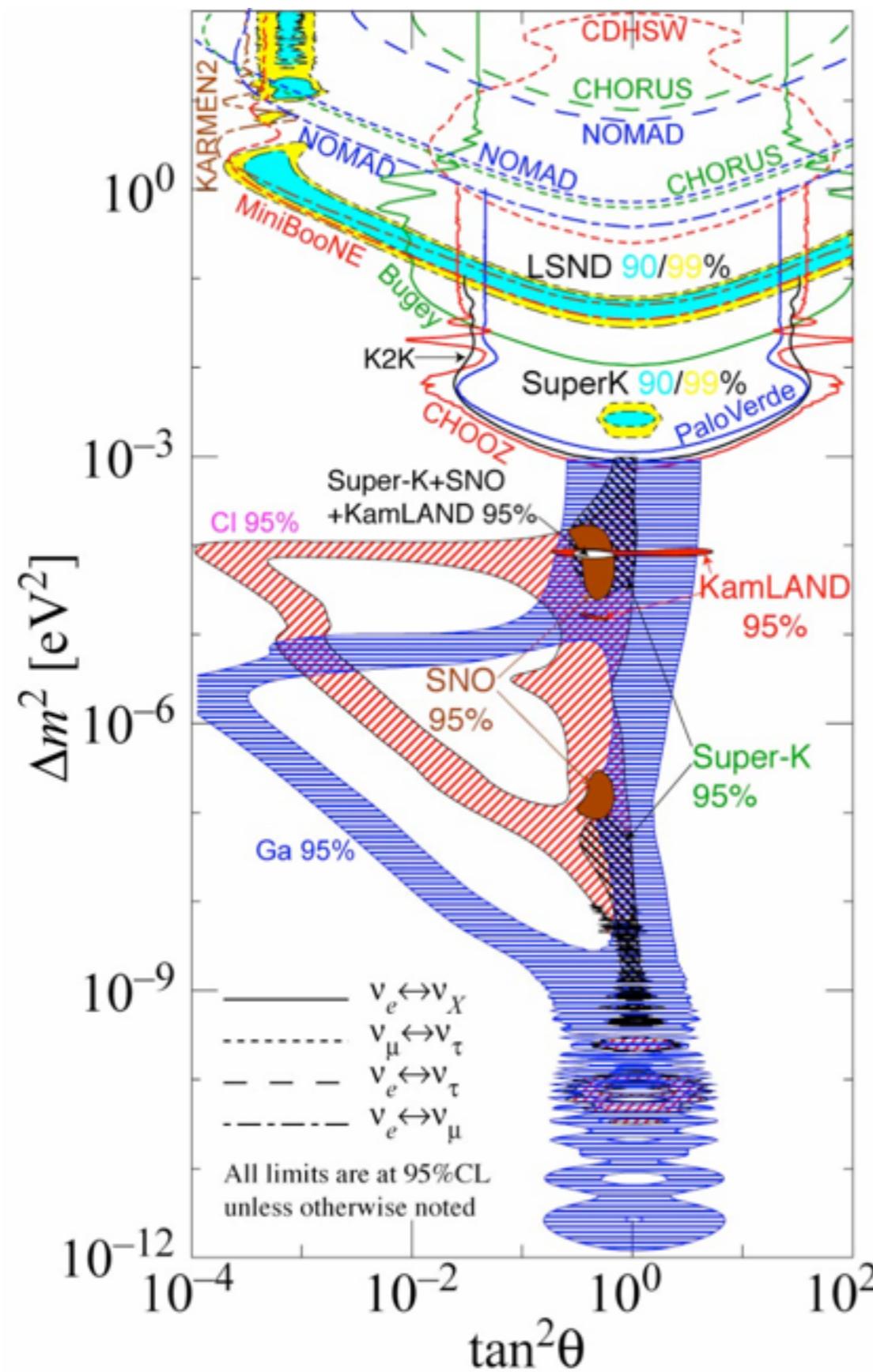
$$\chi^2 = \sum_i^n \left( \frac{(\hat{\mu}_i - \mu_i(\theta))^2}{2\sigma^2} \right)$$

has known distribution:  $E[\chi^2] = n_f$ : #number of “degree of freedom”  
i.e.  $n - \#fitted\ parameters$

Chi2 Probability: The 1-cumulative distr. of  $P(\chi^2, n_f)$  distribution  
how often to expect “worse” fit result (i.e. with larger  $\chi^2$  value at min.)

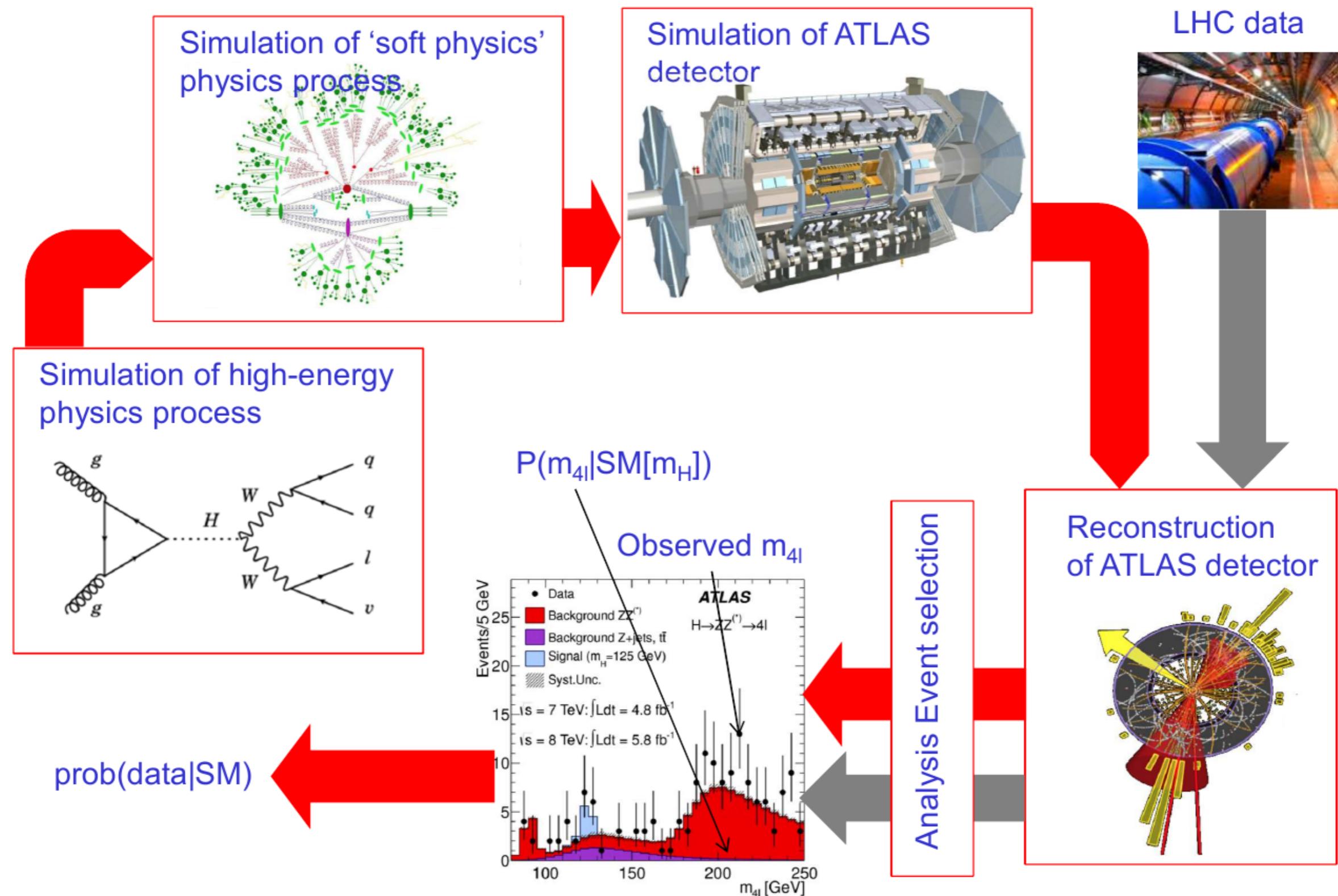


# Examples of Confidence Intervals

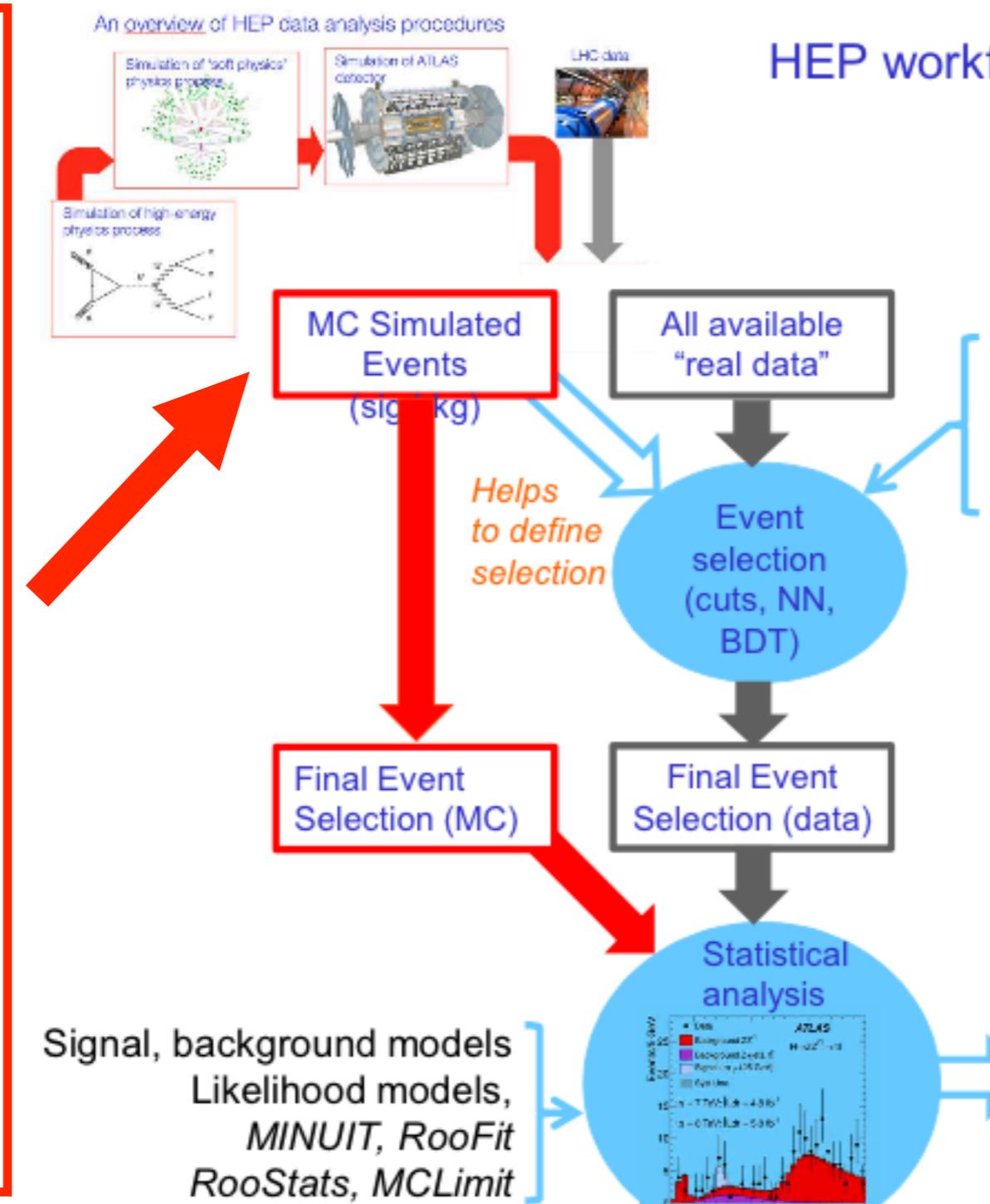
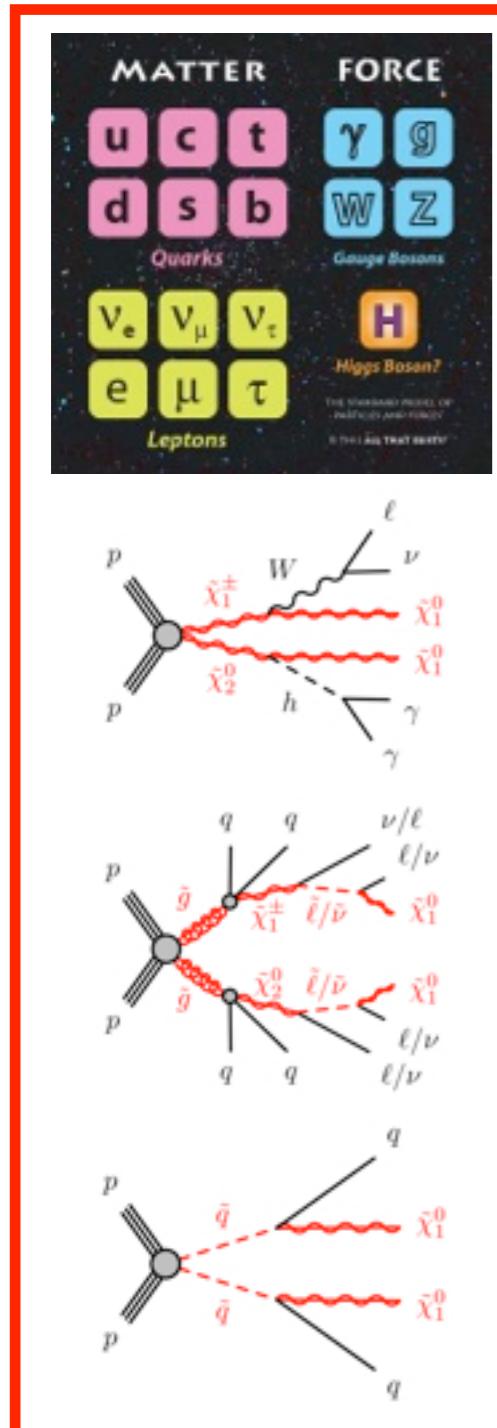


# Physics Analysis

# Particle Physics Workflow



# Analysis Workflow

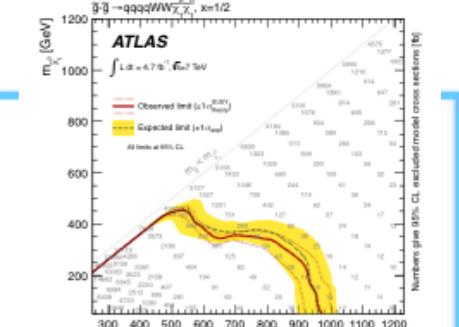


## HEP workflow: analysis view

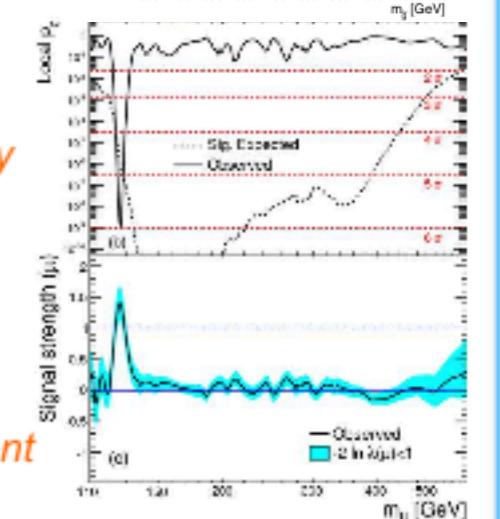
N-tuples  
Cut-flows,  
Multi-variate analysis (NN,BDT)  
ROOT, TMVA, NeuroBayes

### Final Result

*Limit*

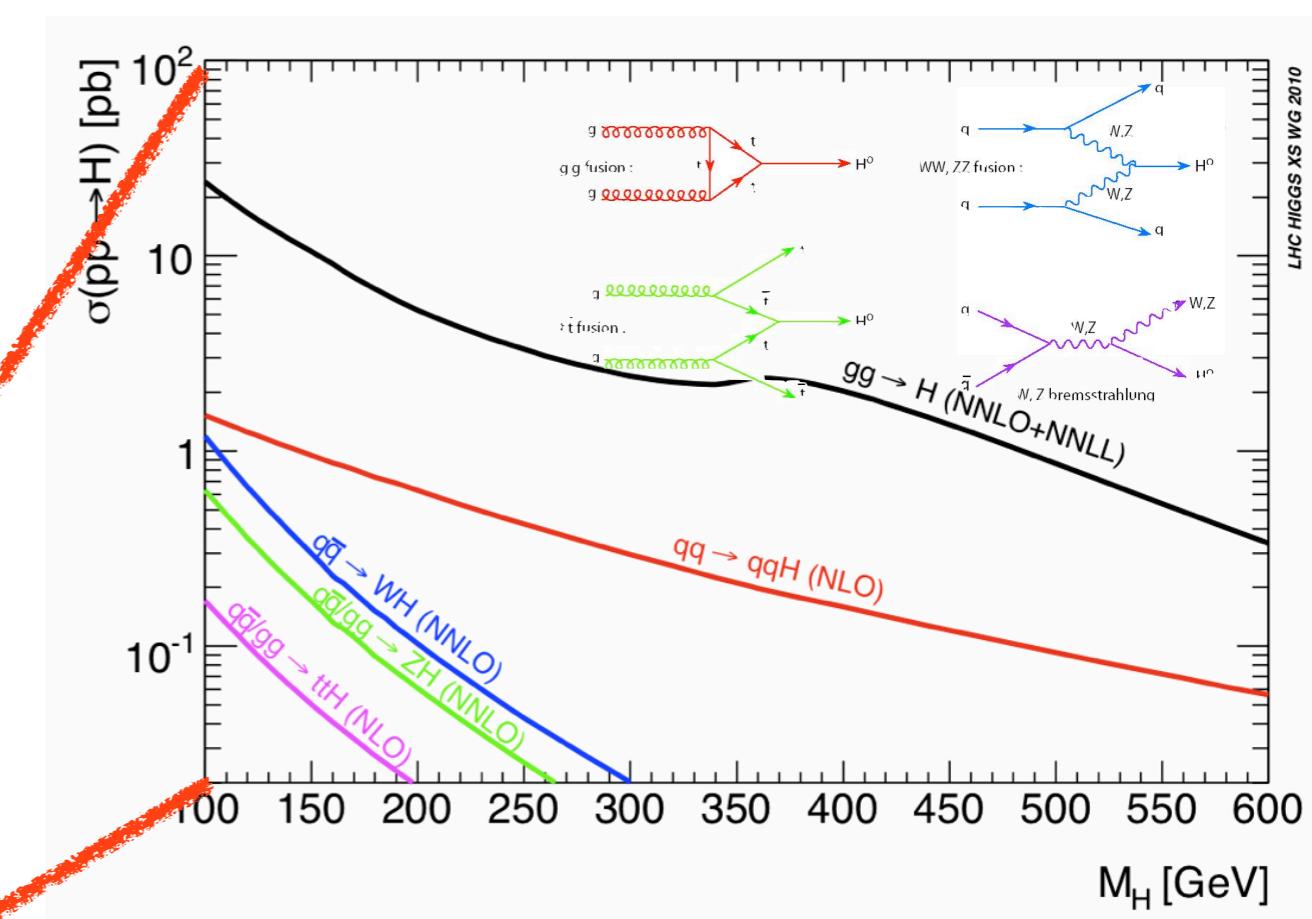
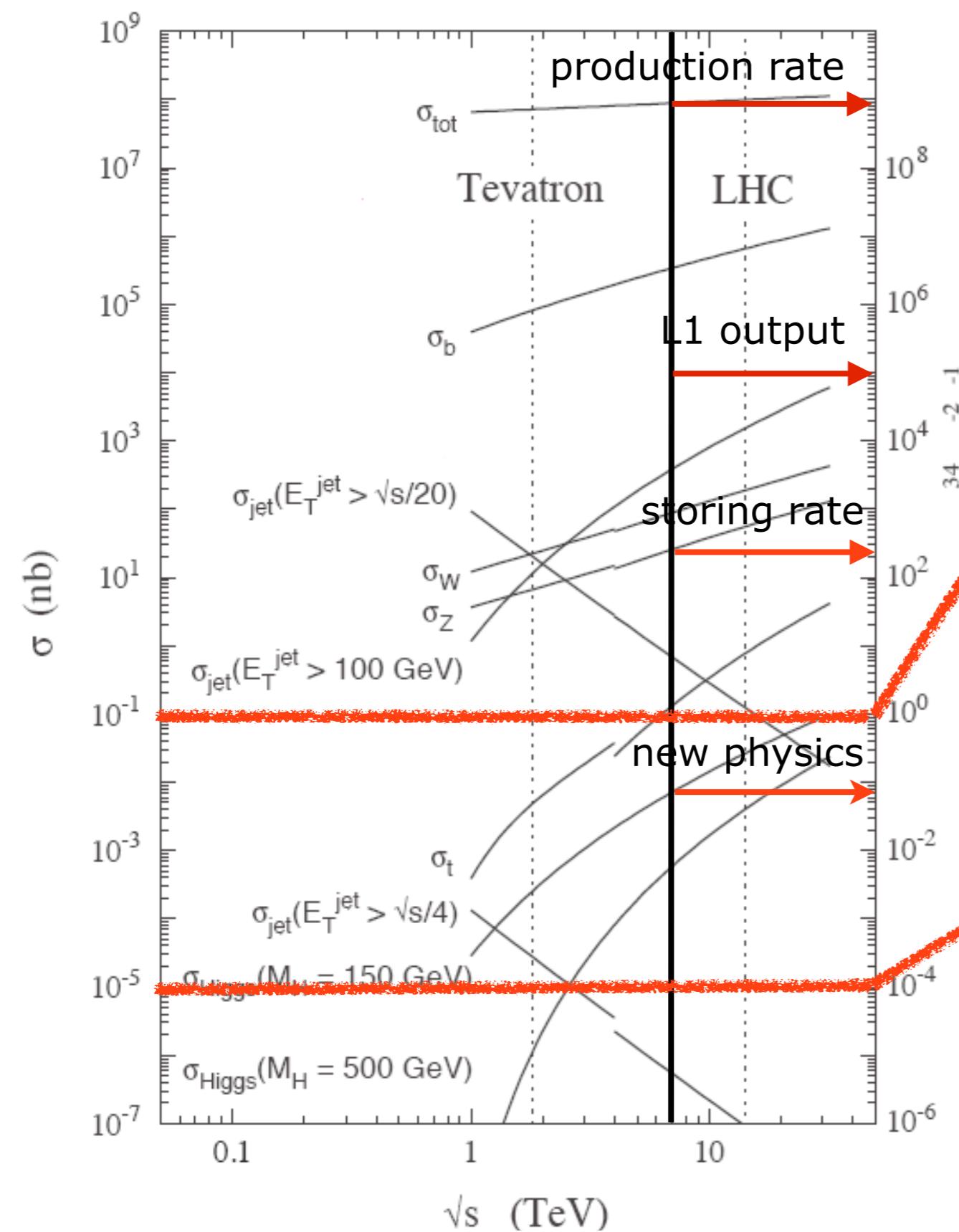


*Discovery*



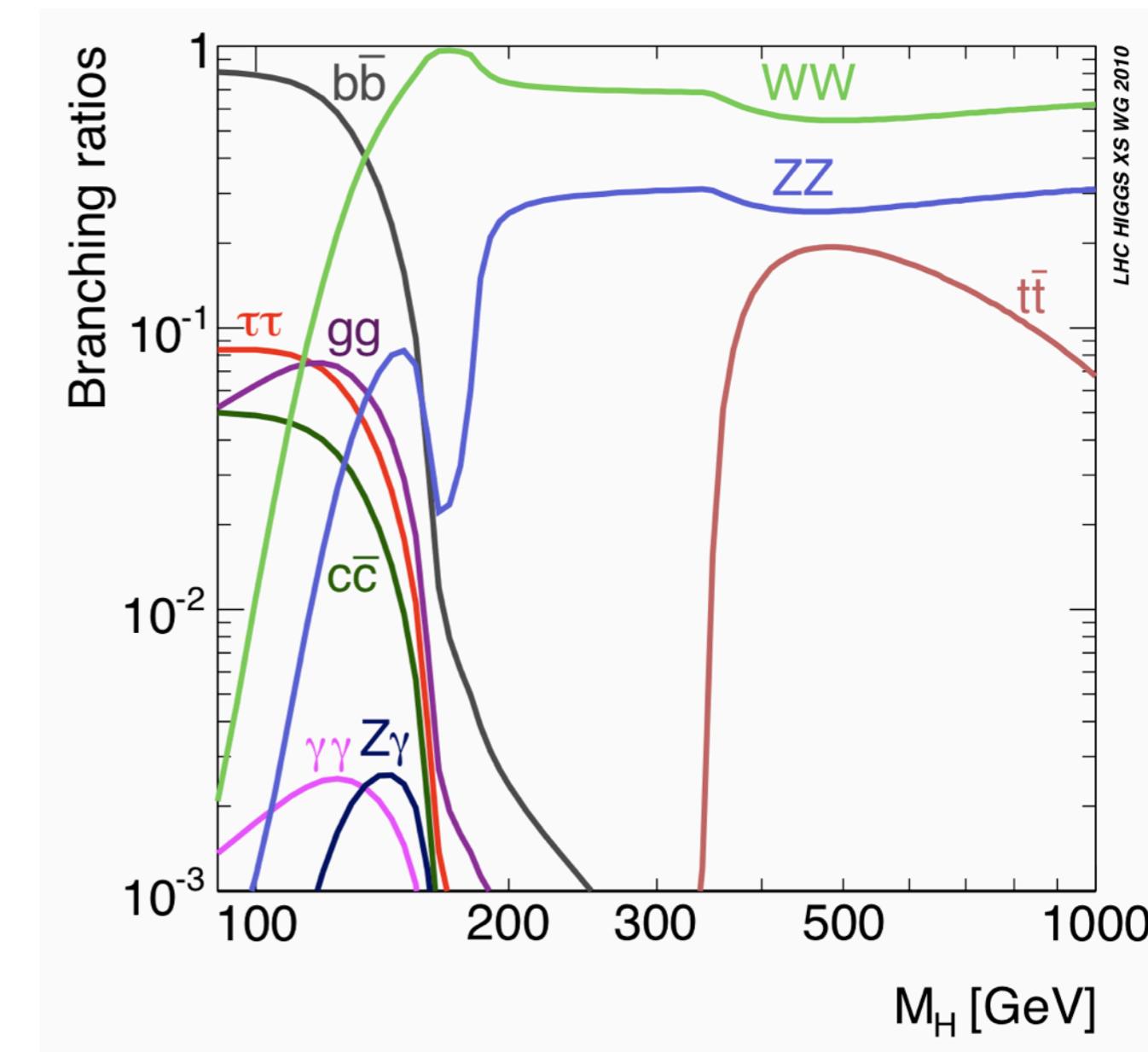
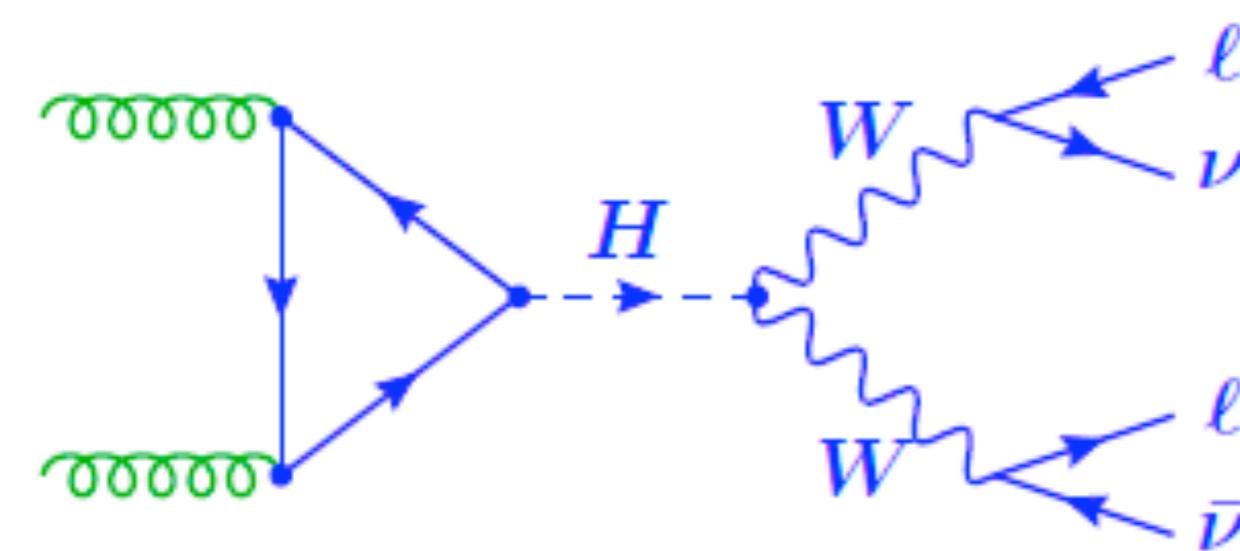
*Measurement*

# Extremely Rare Production processes



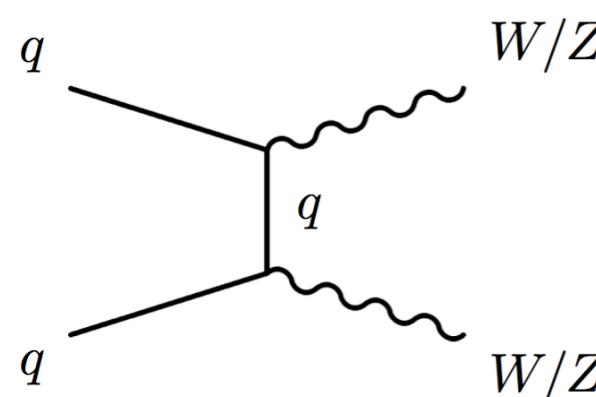
# $H > WW > l\nu\bar{l}\nu$ : Signal

$gg \rightarrow H \rightarrow WW \rightarrow l\bar{\nu}l\nu$

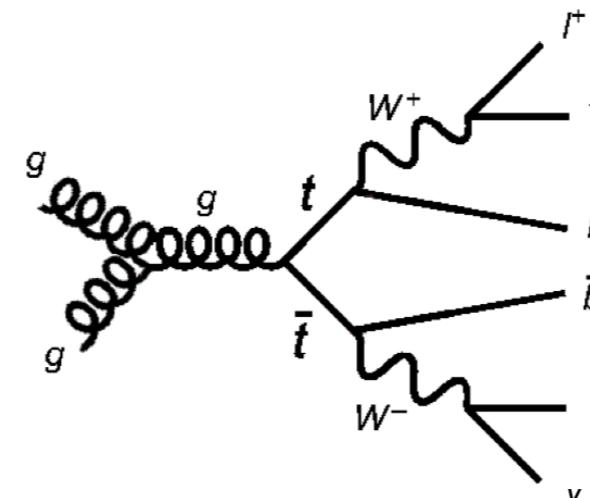


# Backgrounds

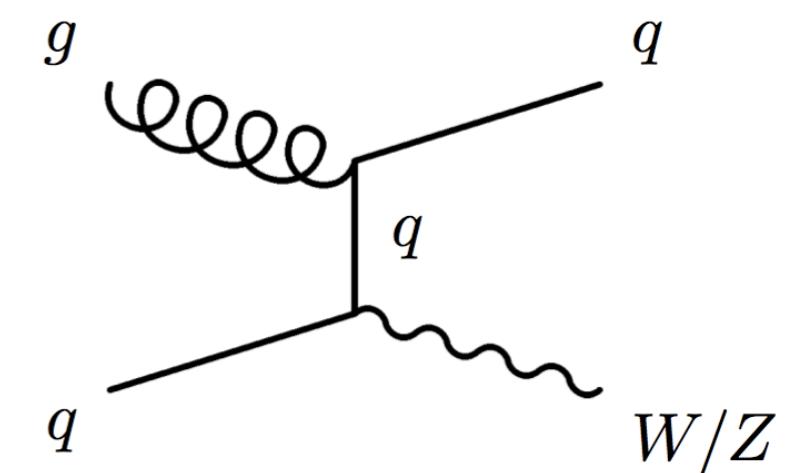
- two identified leptons + missing energy in the final state



$$\sigma = 4.5 \text{ pb}$$



$$\sigma = 15 \text{ pb}$$



$$\begin{aligned} W: \sigma &= 31 \cdot 10^3 \text{ pb} \\ Z: \sigma &= 3.5 \cdot 10^3 \text{ pb} \end{aligned}$$

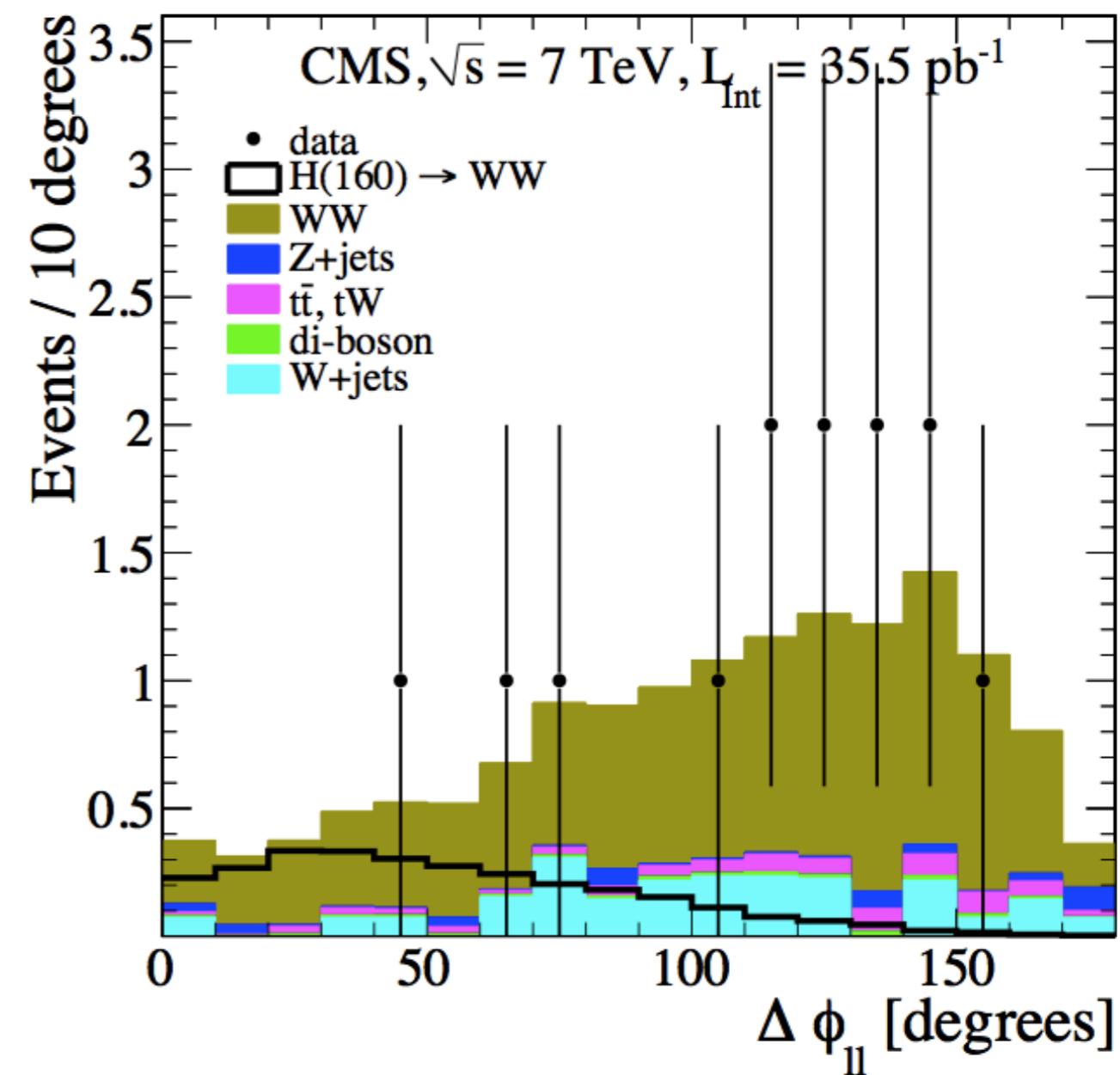
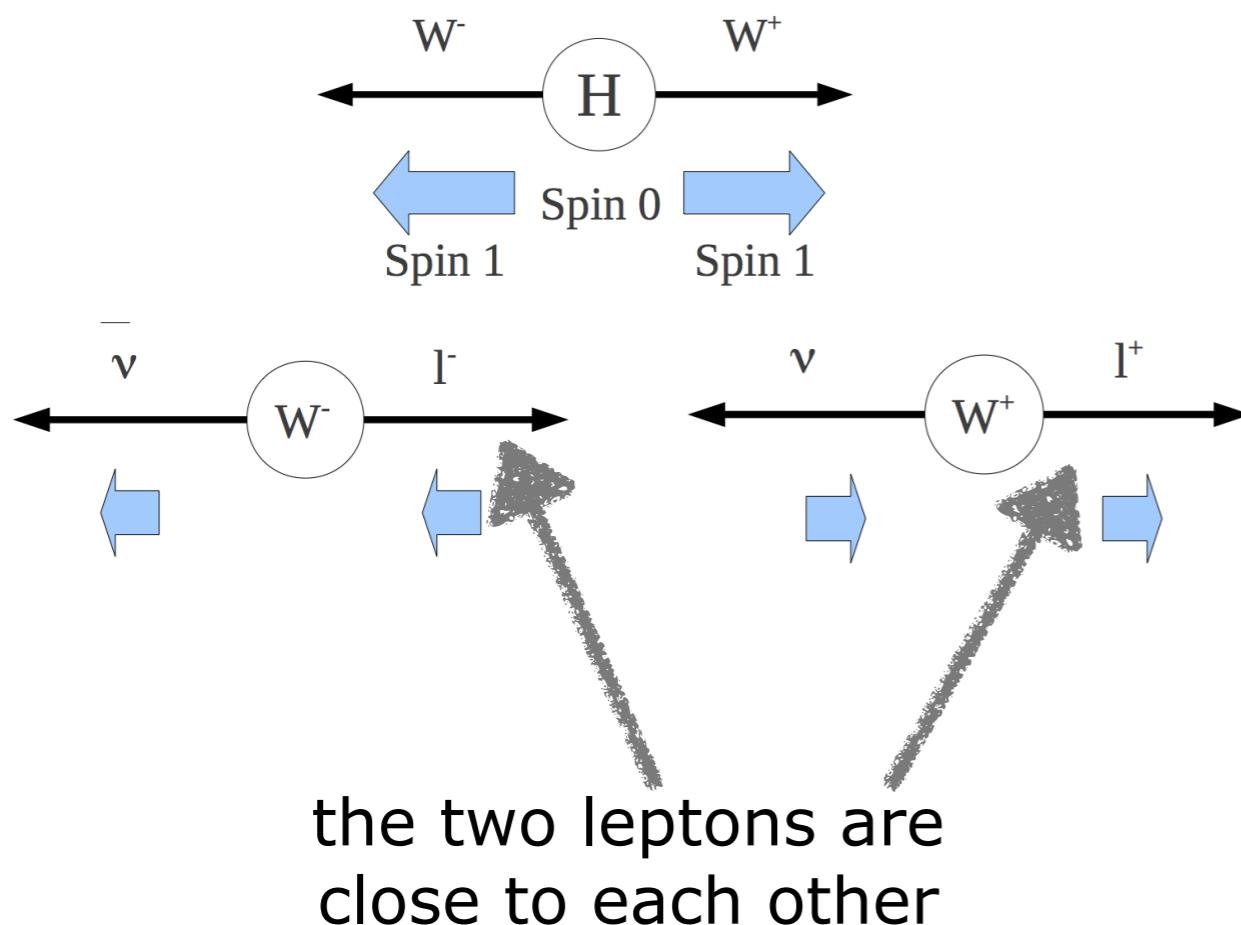
irreducible: same final state of the signal, exploit different kinematics of the production

there are two additional  $b$ -jets in the detector, due to the top decay, veto on jets (or on  $b$ -jets)

jets in the detector can give a lepton-like signature (non prompt leptons, or fake leptons from track+calo deposit): very high cross section

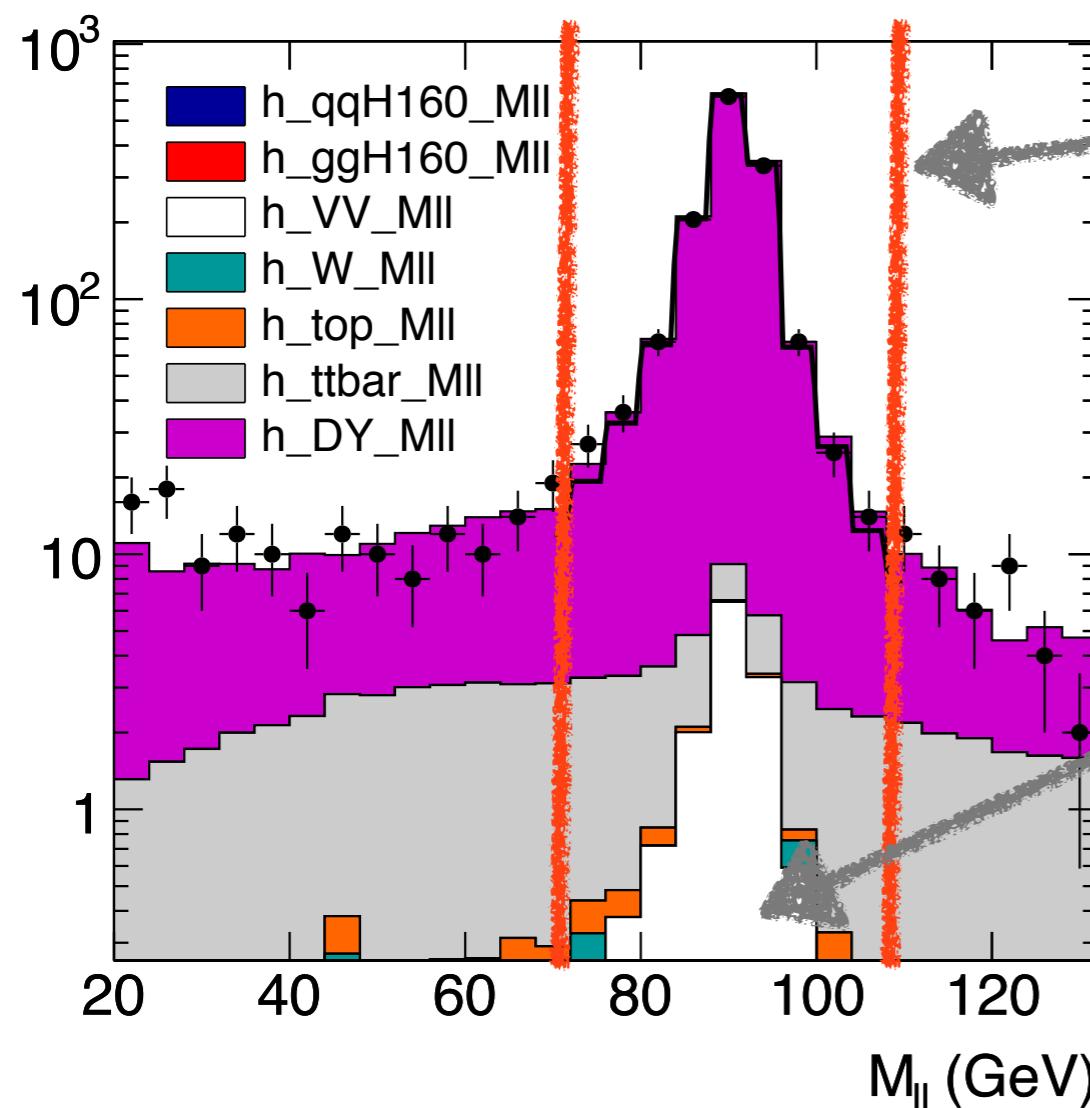
# Signal Features

example of a discriminating variable;



# Cross-section measurement

$$\sigma = \frac{N_{obs} - N_{bkg}}{\varepsilon \cdot \int \mathcal{L} dt}$$



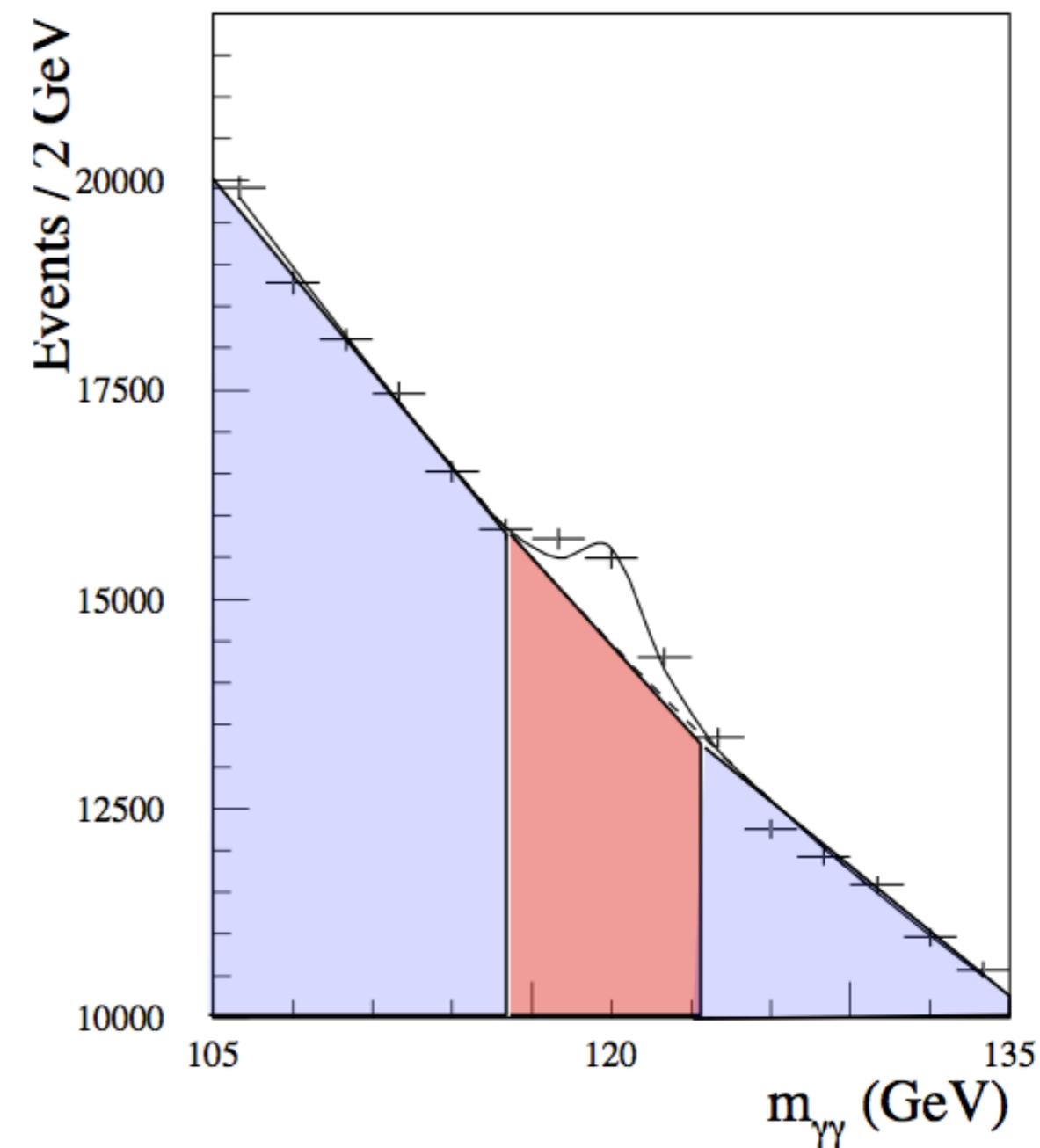
the Z contribution is dominant under the Z resonance peak

the uncertainty is due to the statistics available and the contamination of other samples

evaluate the **impact on the analysis**: probably does not need the same precision as a cross-section measurement

# Side-bands method

- when the background is **expected to behave smoothly**, for example in case of random combinations
- assume a simple shape, and **extrapolate the background under the signal peak** from the sides
- **fit the distribution** with signal shape (a resonance) and background (exponential, linear)



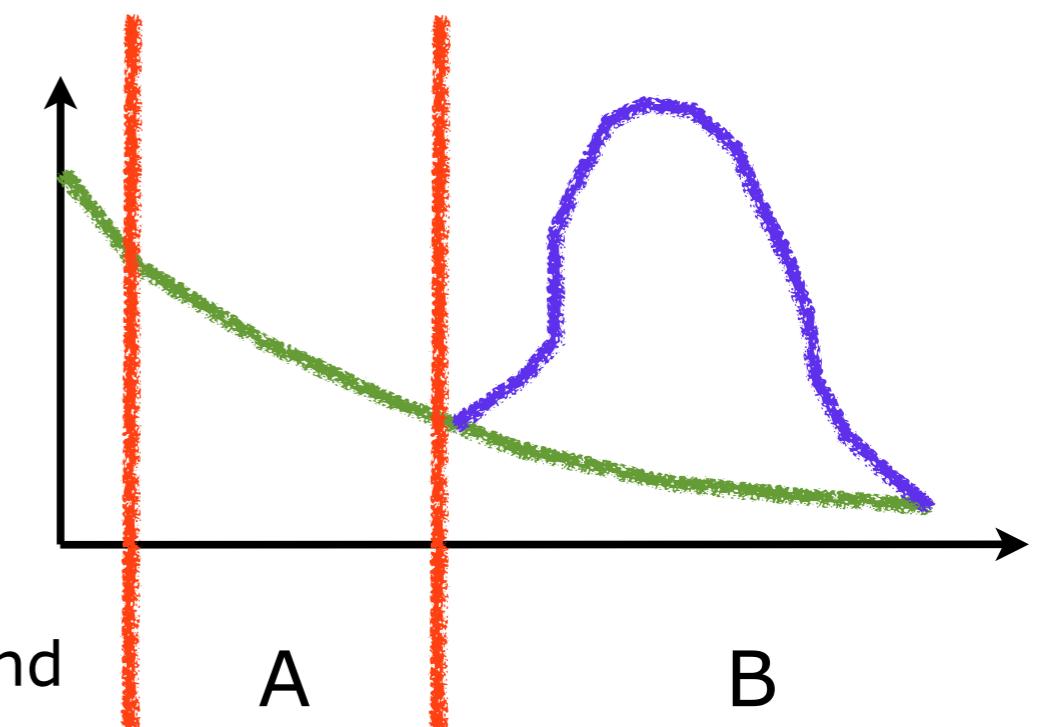
# Control region

- assume the **knowledge of the background shape**, at a certain level of the analysis
- fit the shape to data in a **signal-free region**, where that background is dominant and extrapolate it to the signal region
- in case of low statistics, **count the number of events** and extrapolate

$$N_{\text{inferred}}^{bkg-A} = N_{DATA}^{bkg-B} \left( \frac{N_{MC}^{bkg-A}}{N_{MC}^{bkg-B}} \right)$$

uncertainty due  
to the data  
statistics and  
other systematic  
sources

uncertainty due to:  
background  
model, control region  
contamination, propagation of  
other systematic sources,  
MC stats

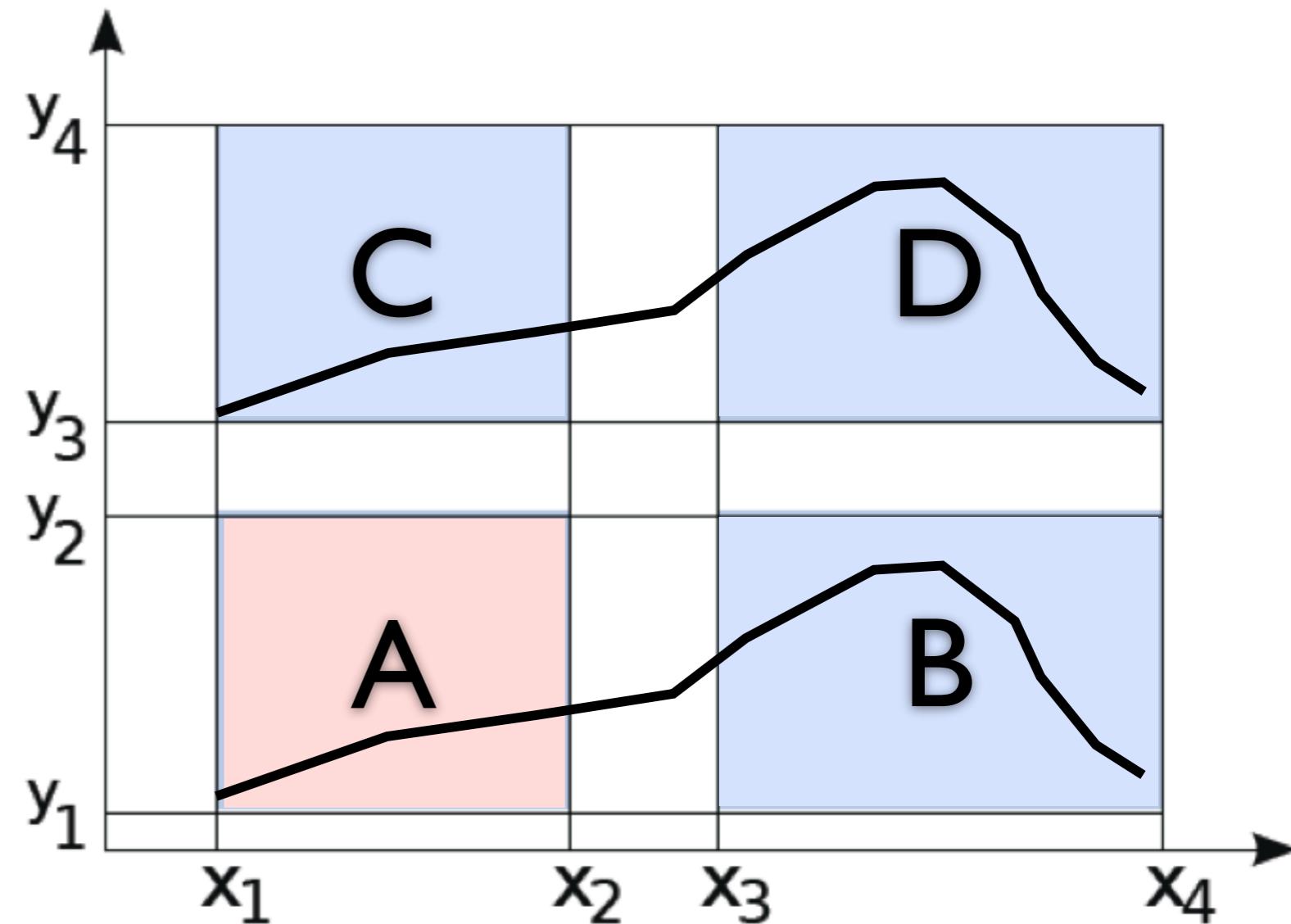


# ABCD method

$$N_A^{bkg} = N_B^{bkg} \frac{N_C^{bkg}}{N_D^{bkg}}$$

uncertainty due to the data statistics and other systematic sources

uncertainty due to: control region contamination, propagation of other systematic sources

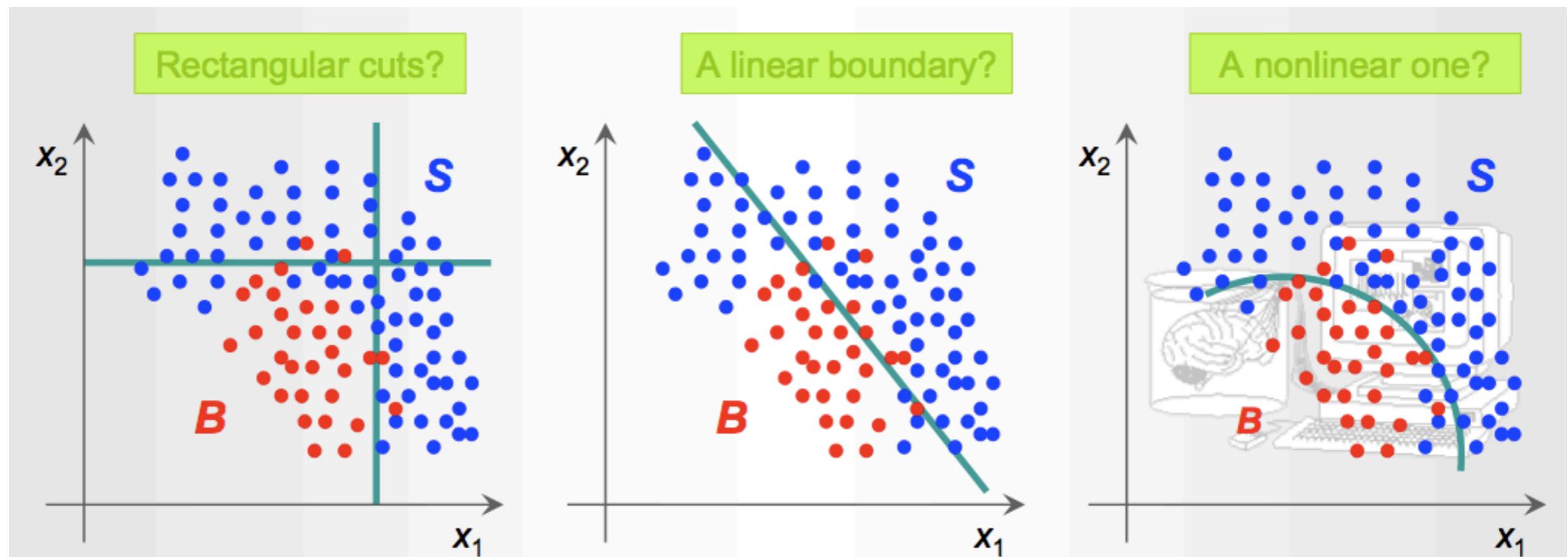


- assume the bkg pdf to be factorized:  $f^{bkg}(x, y) = f_x^{bkg}(x) \cdot f_y^{bkg}(y)$
- the **correlation check** done with simulation is a less stringent requirement than the good description of the shape
- in case of low statistics, **count the number** of events and extrapolate

# Multivariate Techniques

- **Toolkit for Multivariate Data Analysis with ROOT**, <http://tmva.sourceforge.net/>  
H. Voss, **Multivariate Data Analysis and Machine Learning in High Energy Physics**

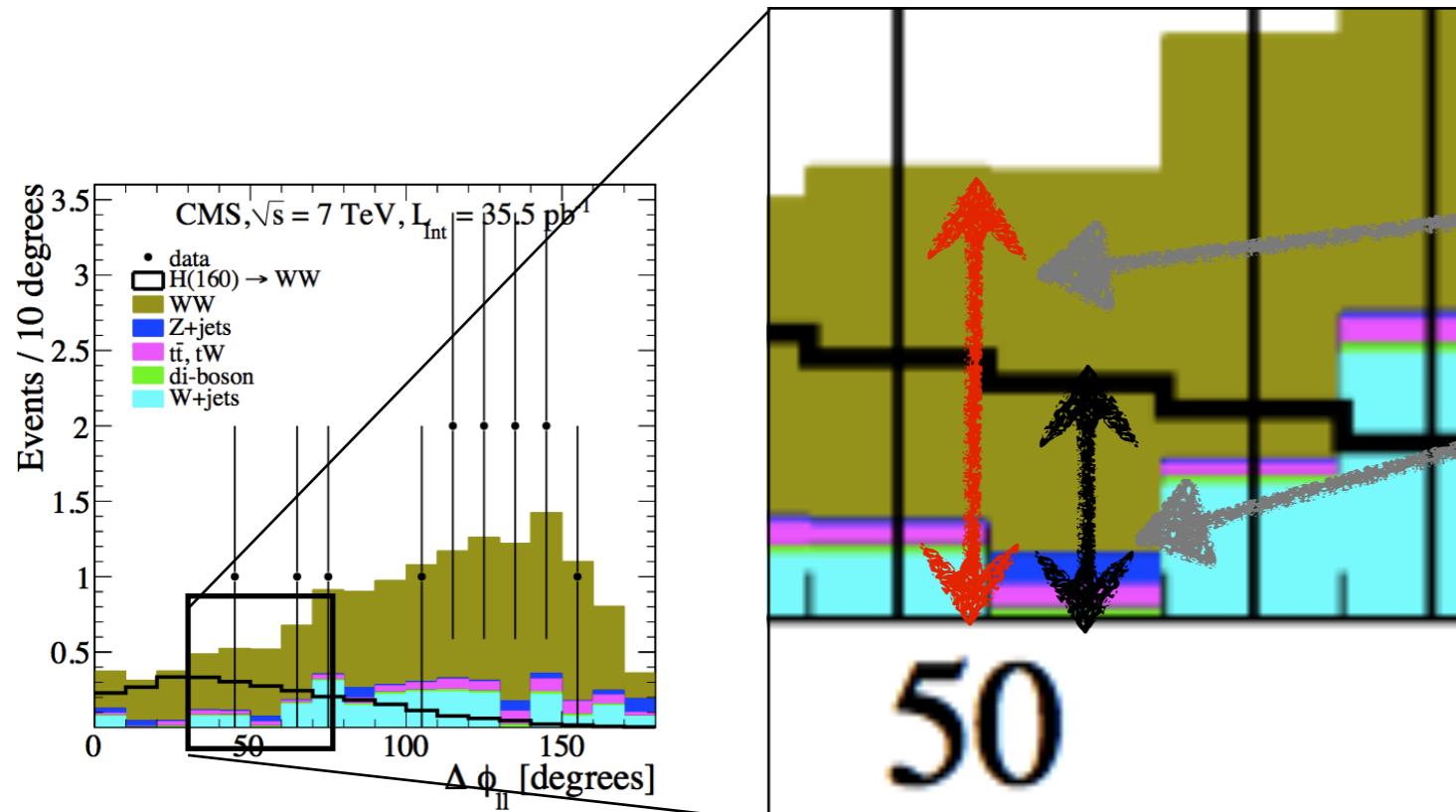
# Signal vs Background



- rectangular selections do not fully exploit the **topology of the events**
- build **more sophisticated discriminants** to separate signal from backgrounds
- need a **good knowledge** of both signal and background
- need high Monte Carlo **statistics**

# Likelihood Discriminant

search for a classification of the events, that maps the set of the analysis variables into a single one



$$y_i = f(\vec{x}_i) : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$R_L(i) = \frac{L_S(i)}{L_B(i) + L_S(i)}$$

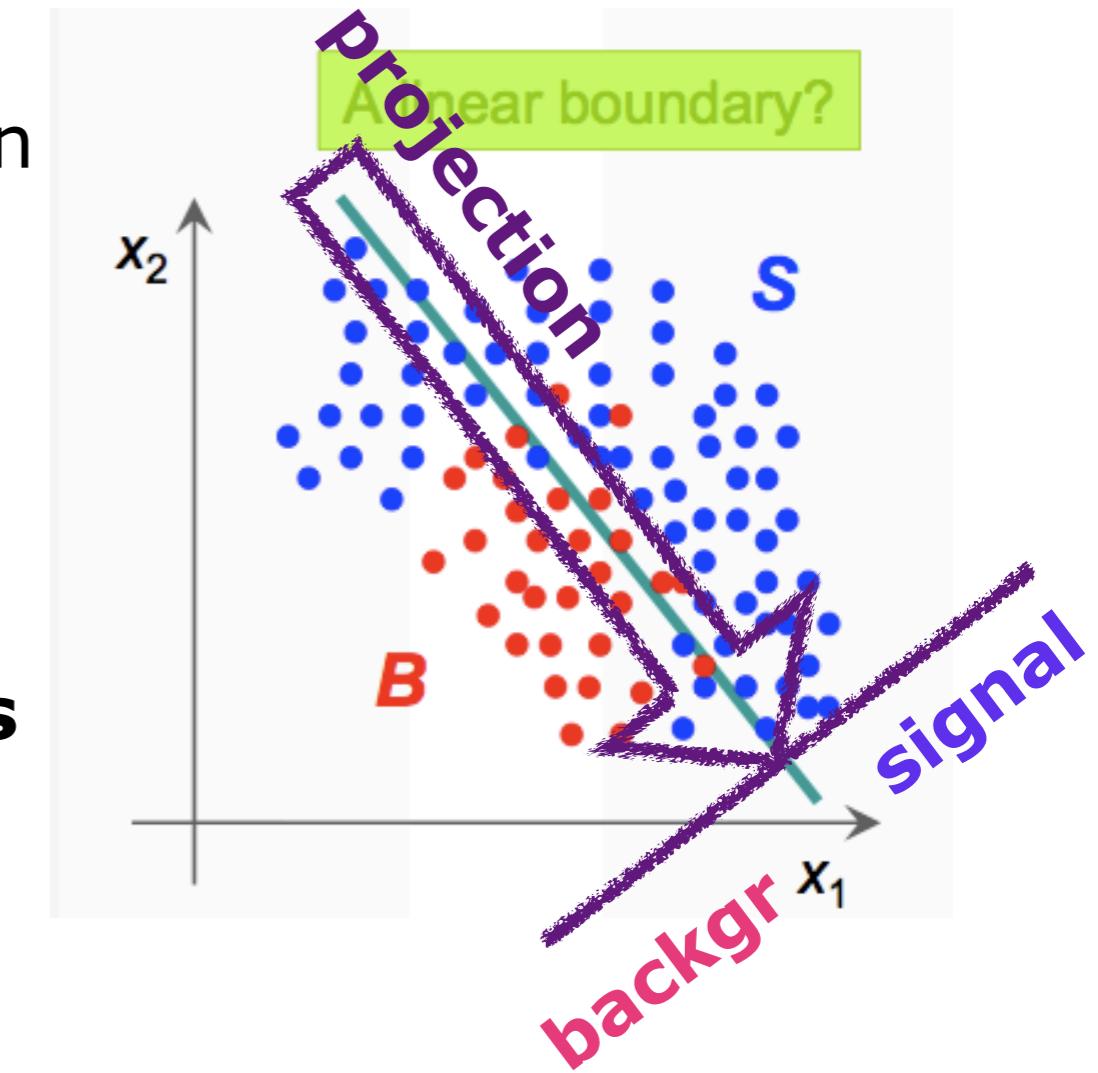
estimates for each event  
the confidence of being  
signal-like

For more, uncorrelated,  
variables: “easily” built  
For linearly correlated variables,  
first decorrelate them

$$L_S(i) = f_S(\vec{x}_i) = \prod_{j \in (\text{vars})} f_{S,j}(x_{ij})$$

# Fisher Discriminant

- project **high-dimensional dataset onto a line** and perform classification in this one-dimensional space
- **optimization**: maximize the distance between means, while minimizing the variance within each class
- very effective with **linear correlations**



build the linear combination:

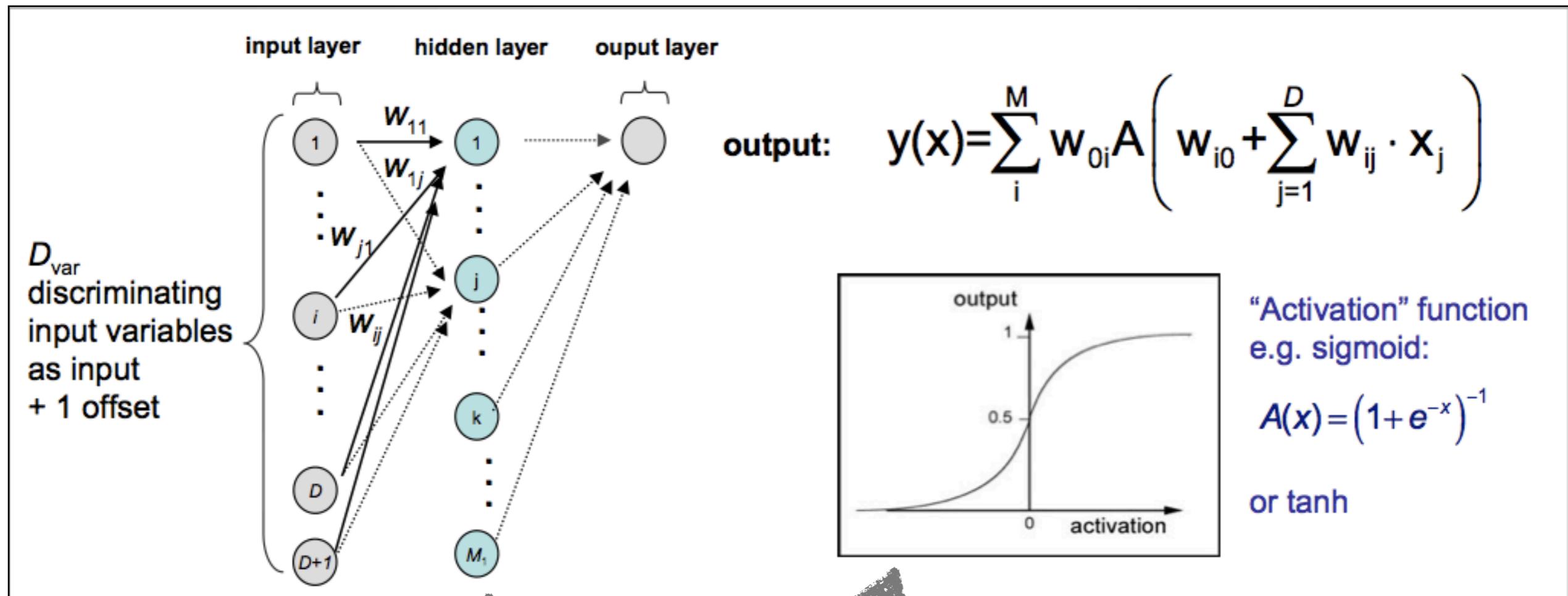
$$y(\vec{x}, \vec{w}) = w_0 + \sum x_i \cdot w_i$$

find the best weight by minimizing the criterion:

$$J(\vec{w}) = \frac{(\langle y_S \rangle - \langle y_B \rangle)^2}{\sigma_{y_S}^2 + \sigma_{y_B}^2}$$

# Neural Network

- to cope with non-linear correlations, try a more sophisticated combination of the inputs



input  
variables,  
none of them  
is a smoking  
gun

factors of the  
“base” in which  
the non linear y is  
decomposed

the non-linear  
base element

need to find the  
weights, i.e. to train the  
neural network

# On the training

**loss function:** how many times I make a mistake in the classification

$$L(w) = \sum_i^{\text{events}} (y(x_i) - y(C))^2$$

predicted event type      true event type

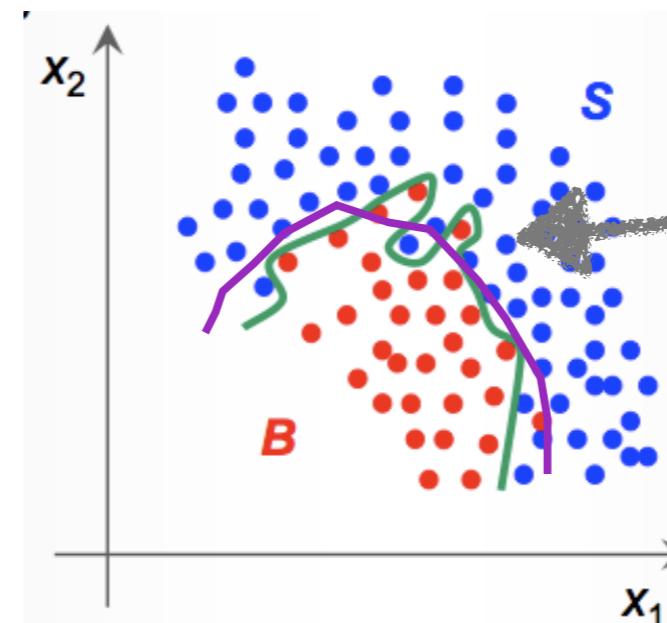
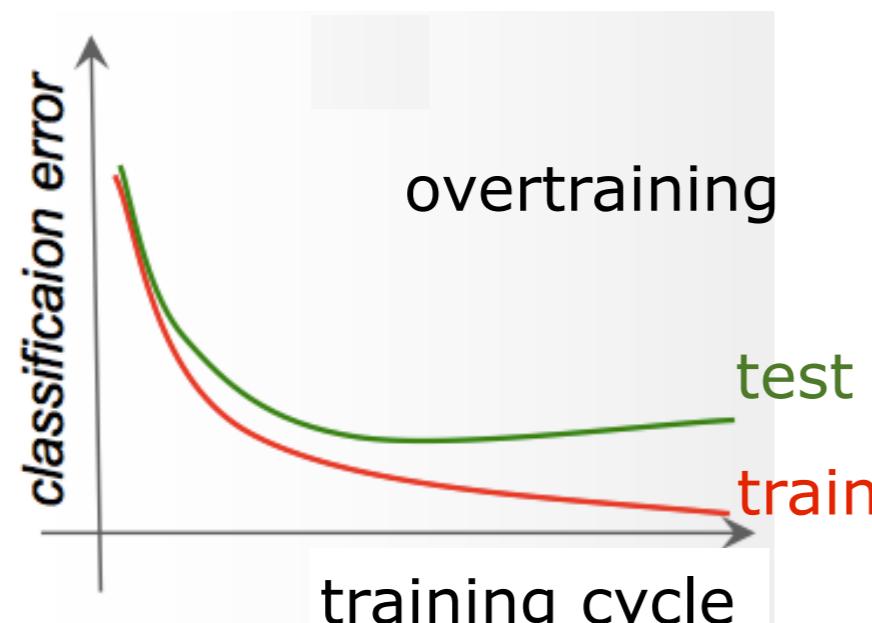
values of  $y$ :  
1 = signal  
0 = background

**minimize** the loss function:

- start from random weights
- change them according to the  $L$  gradient
- loop several times on the training samples

$$w^{n+1} = w^n + \eta \cdot \vec{\nabla}_w L(w)$$

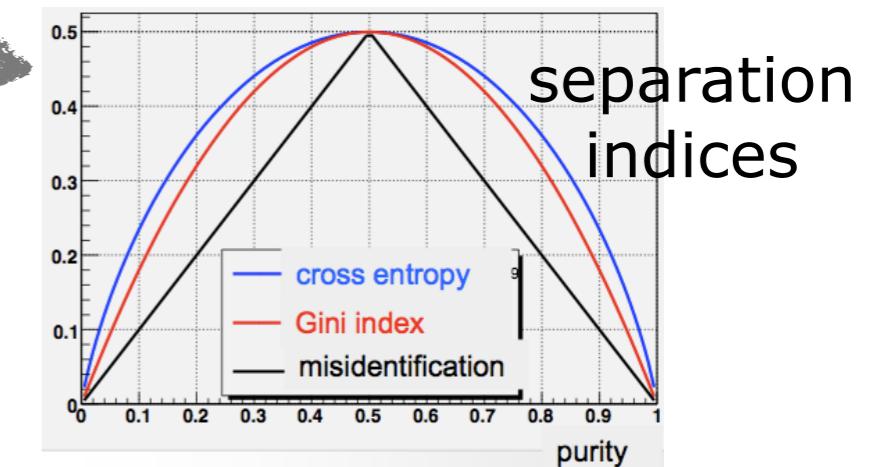
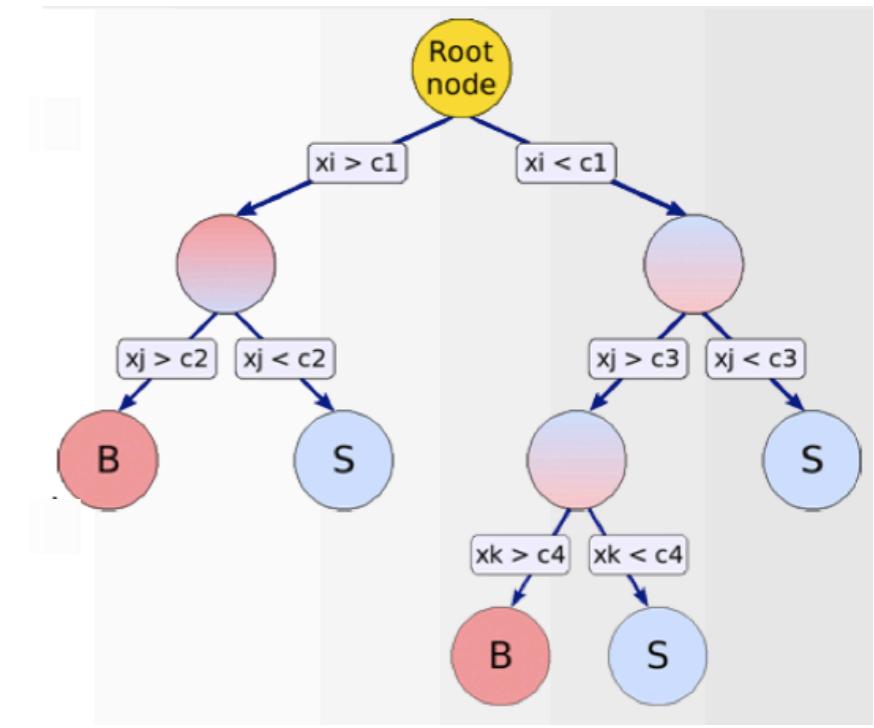
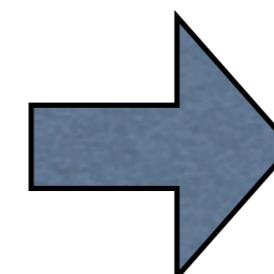
divide the simulated sample into **training** and **testing**,  
continue the training until the performances on the training stabilize,  
stop before the ones on the testing worsen



in overtraining,  
the NN is  
adapting to  
statistical  
fluctuations of  
the training  
sample

# Boosted Decision Trees

- rank the variables in terms of **discriminating power**
- apply **subsequent selections** in each of the variables
  - minimal #events per node
  - maximum number of nodes
  - maximum depth specified
  - a split doesn't give a minimum separation gain
- stop when:
  - maximum number of nodes
  - maximum depth specified
  - a split doesn't give a minimum separation gain
- in each final node (leaf) return S/B discrimination (discrete or continuous)
- independent of monotonous variable transformations
- immune against outliers
- weak variables are ignored
- **very sensitive to statistical fluctuations in training data**

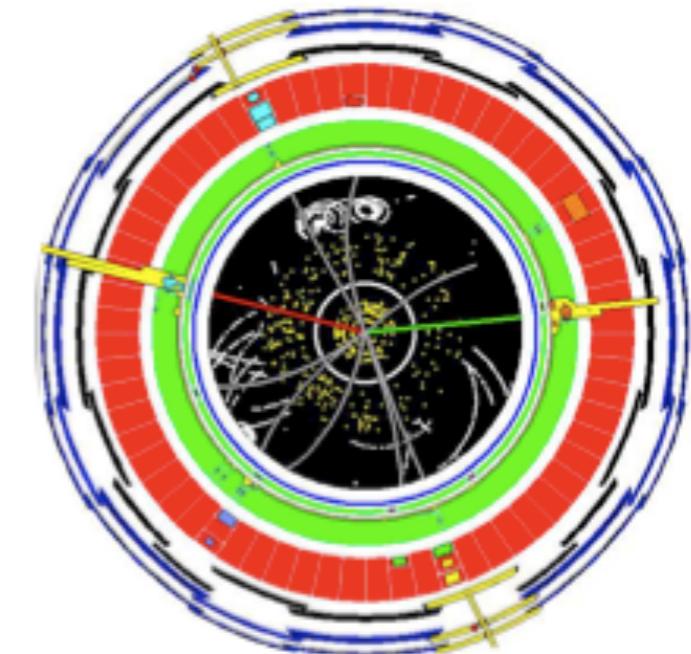
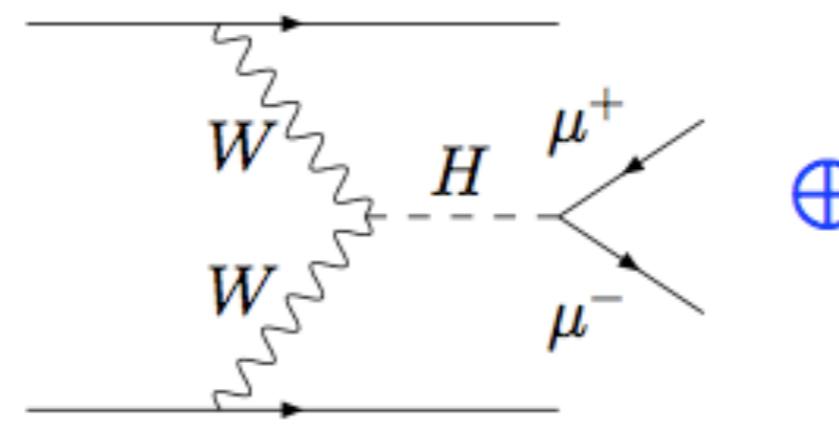


boosting: combine a whole forest of Decision Trees, derived from the same sample, e.g. using different event weights.

# Matrix Element

- the MVA techniques is **describe the final state topology** with a parametrization, built on the simulation
- matrix elements are this description, at the level of the physical process

need to include the effects due to the detector for each physics object considered

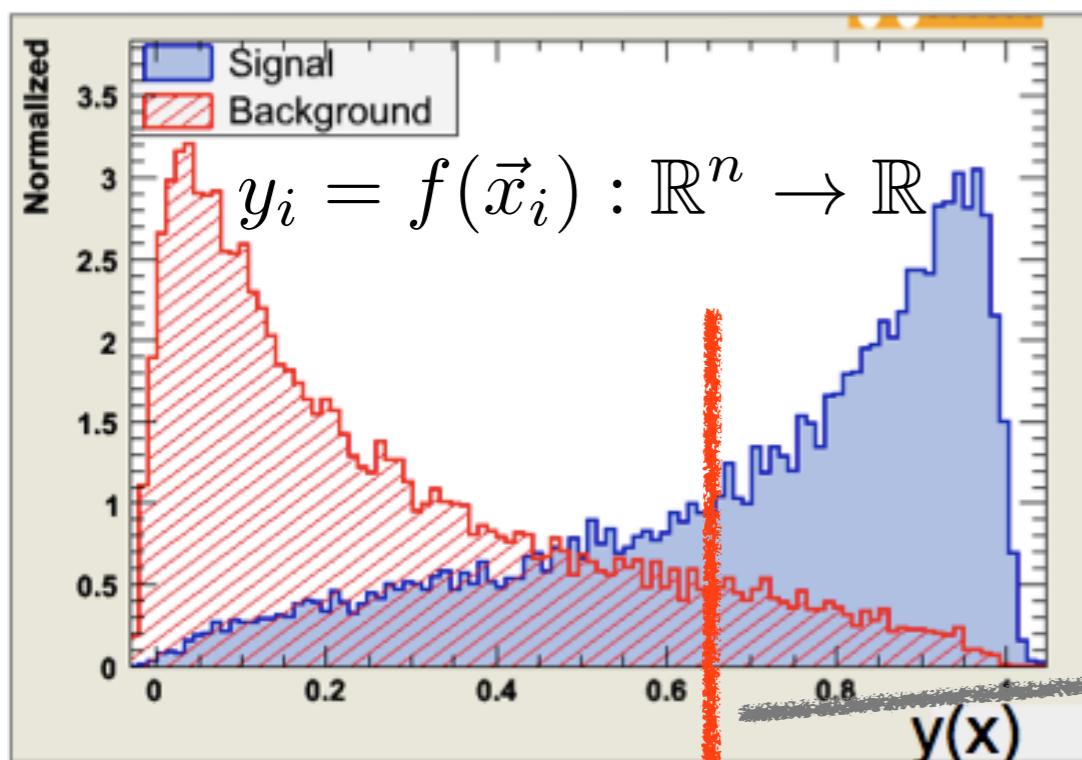


$$P(\mathbf{x}|M_t) = \frac{1}{N} \int d\Phi |\mathcal{M}_{t\bar{t}}(p; M)|^2 \prod_{jets} f(p_i, j_i) f_{PDF}(q_1) f_{PDF}(q_2)$$

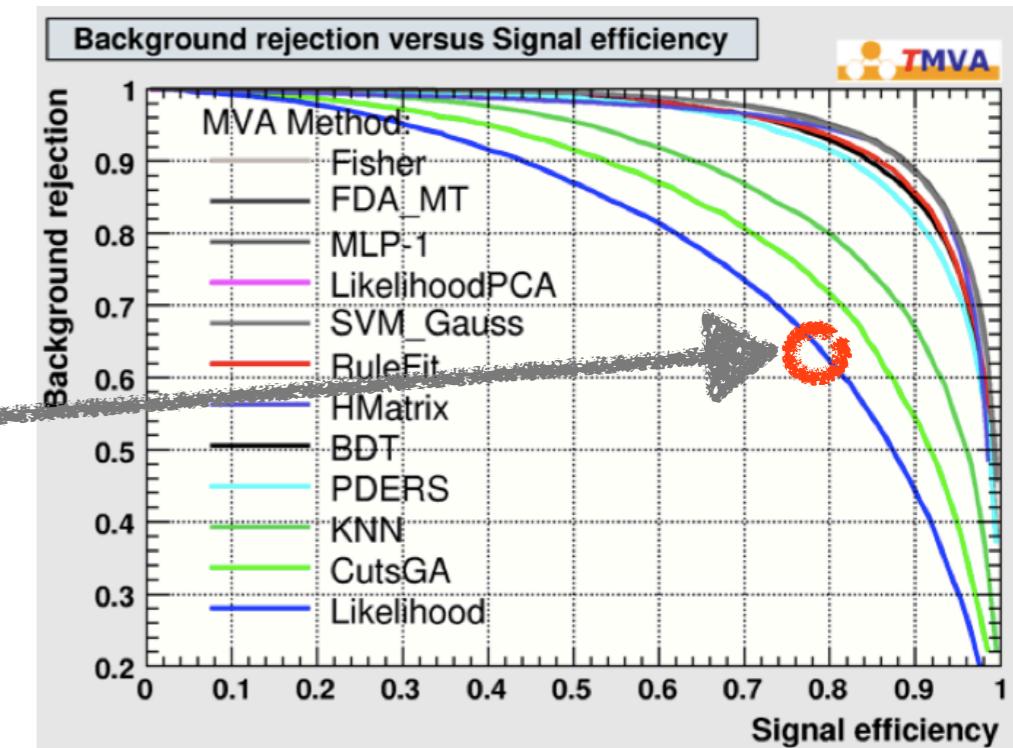
↑  
Phase-space Integral  
↑  
Matrix Element  
↓  
Transfer Functions

calculate the probability for each background and **build a likelihood ratio**

# Event Selection Optimization



for each discriminant, now make the choice:  
what is signal, what is background?



Receiver Operating Characteristics (ROC)  
Curve: how efficiency versus purity

- choose the working point by maximizing a figure of merit:

$$\frac{S}{\sqrt{B}}$$

search:  
sensitivity over  
background  
fluctuations

$$\frac{S}{\sqrt{S + B}}$$

known signal:  
sensitivity over  
fluctuations of the total  
sample

$$\frac{S}{\sqrt{B + (\Delta B)^2}}$$

search:  
sensitivity over  
background fluctuations  
plus systematics

# Method Comparison

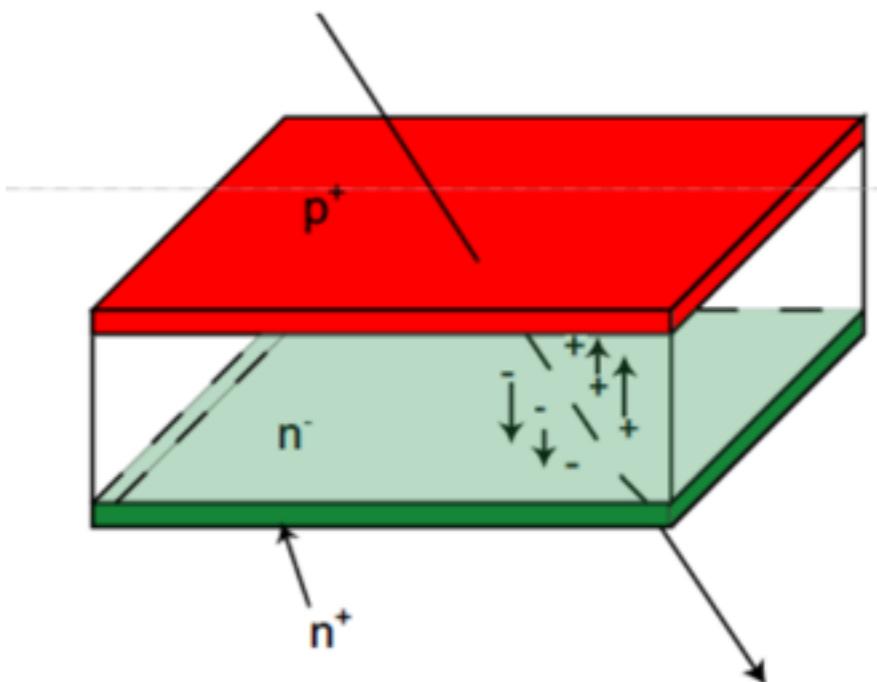
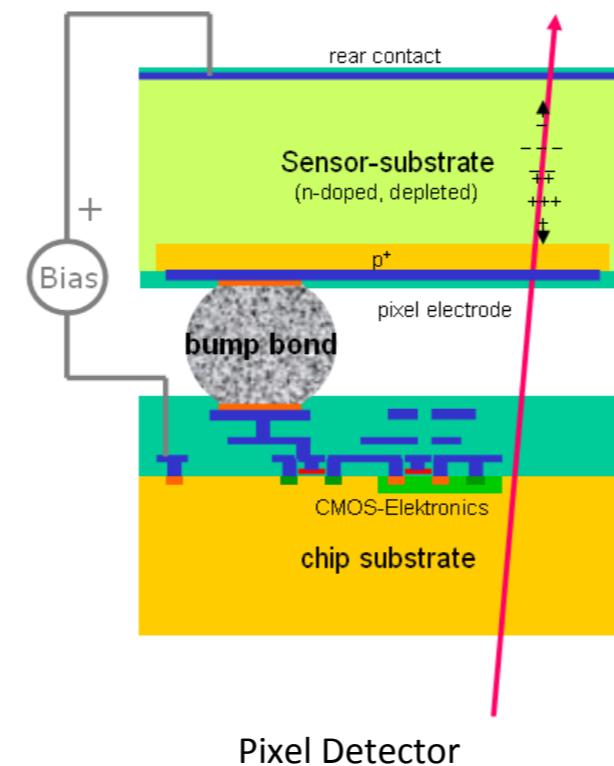
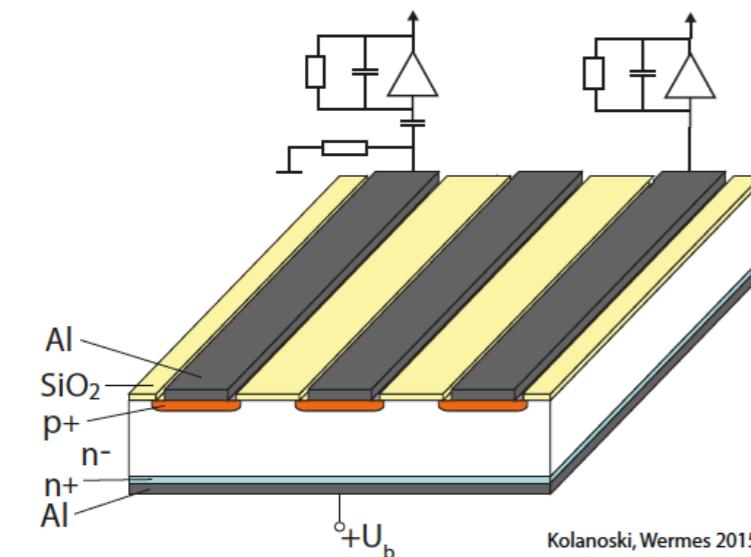
- useful table for the choice of the method to be used, among the ones provided by TMVA

Criteria		Classifiers								
		Cuts	Likeli-hood	PDERS / k-NN	H-Matrix	Fisher	MLP	BDT	RuleFit	SVM
Performance	no / linear correlations	😊	😊	😊	😊	😊	😊	😊	😊	😊
	nonlinear correlations	😊	😢	😊	😢	😢	😊	😊	😊	😊
Speed	Training	😢	😊	😊	😊	😊	😊	😢	😊	😢
	Response	😊	😊	😢/😊	😊	😊	😊	😊	😊	😊
Robust-ness	Overtraining	😊	😊	😊	😊	😊	😢	😢	😊	😊
	Weak input variables	😊	😊	😢	😊	😊	😊	😊	😊	😊
Curse of dimensionality		😢	😊	😢	😊	😊	😊	😊	😊	😊
Transparency		😊	😊	😊	😊	😊	😢	😢	😢	😢

# Capacitance of Silicon Pixels

<https://arxiv.org/abs/hep-ex/0003032.pdf>

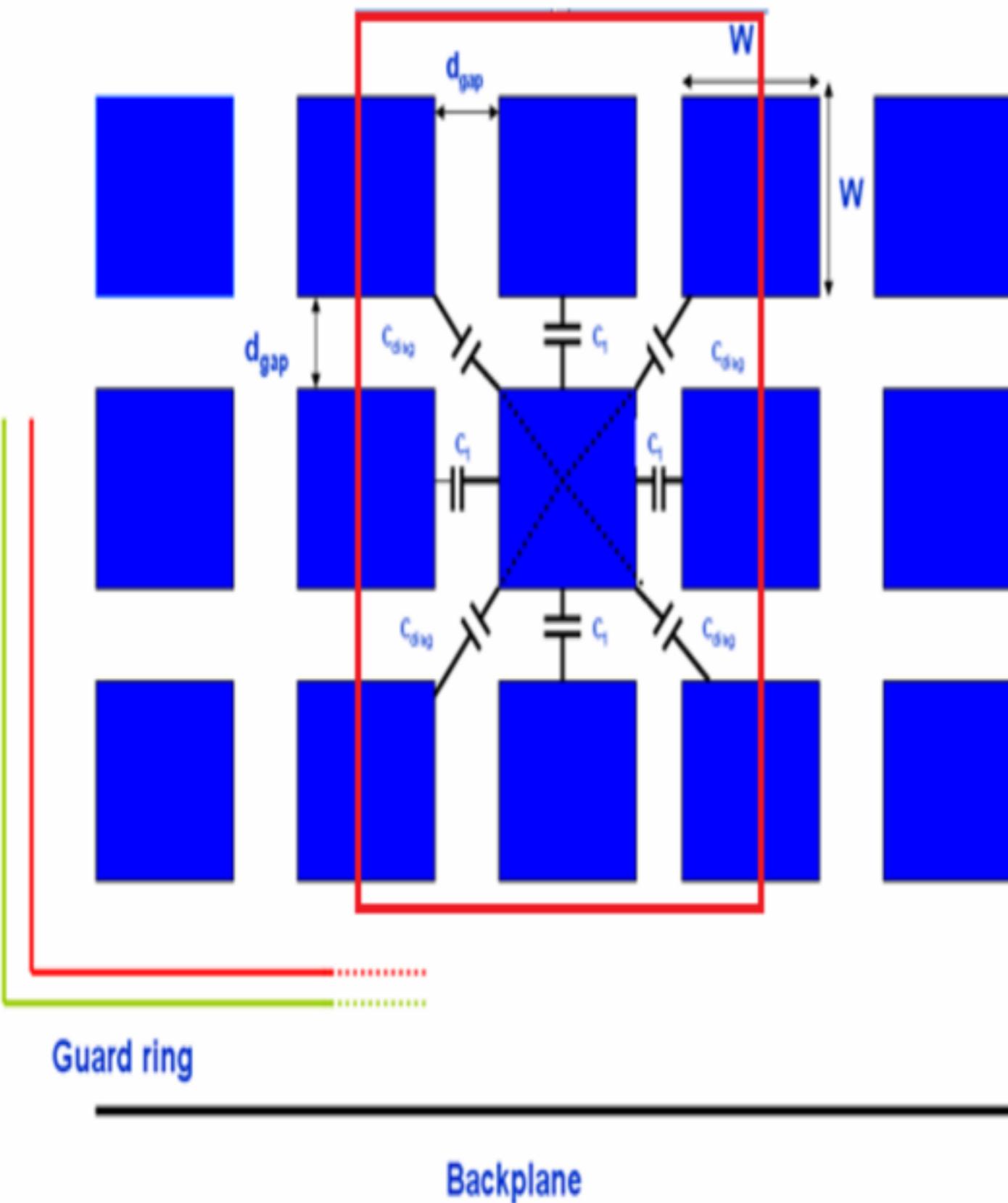
# Semiconductor detectors



**in Si bulk fully depleted**

- $w_i = 3.65 \text{ eV per e/h}$
- a high energy particle  
→  $\sim 80 \text{ e/h per } \mu\text{m}$
- all charge collected
- $\sim 20\,000 \text{ e/h per } 250 \mu\text{m} = 3 \text{ fC}$

# TCAD Simulation: capacitance calculation



Capacitance calculations in p<sup>+</sup>n silicon pixel sensors using three dimensional TCAD simulation approach

Ajay K. Srivastava<sup>a,1</sup> E. Fretwurst<sup>a</sup>, R.Klanner<sup>a</sup>

<sup>a</sup>Institute for Experimental Physics, University of Hamburg, Hamburg 22761, Germany

Internal note (within AGIPD collaboration)

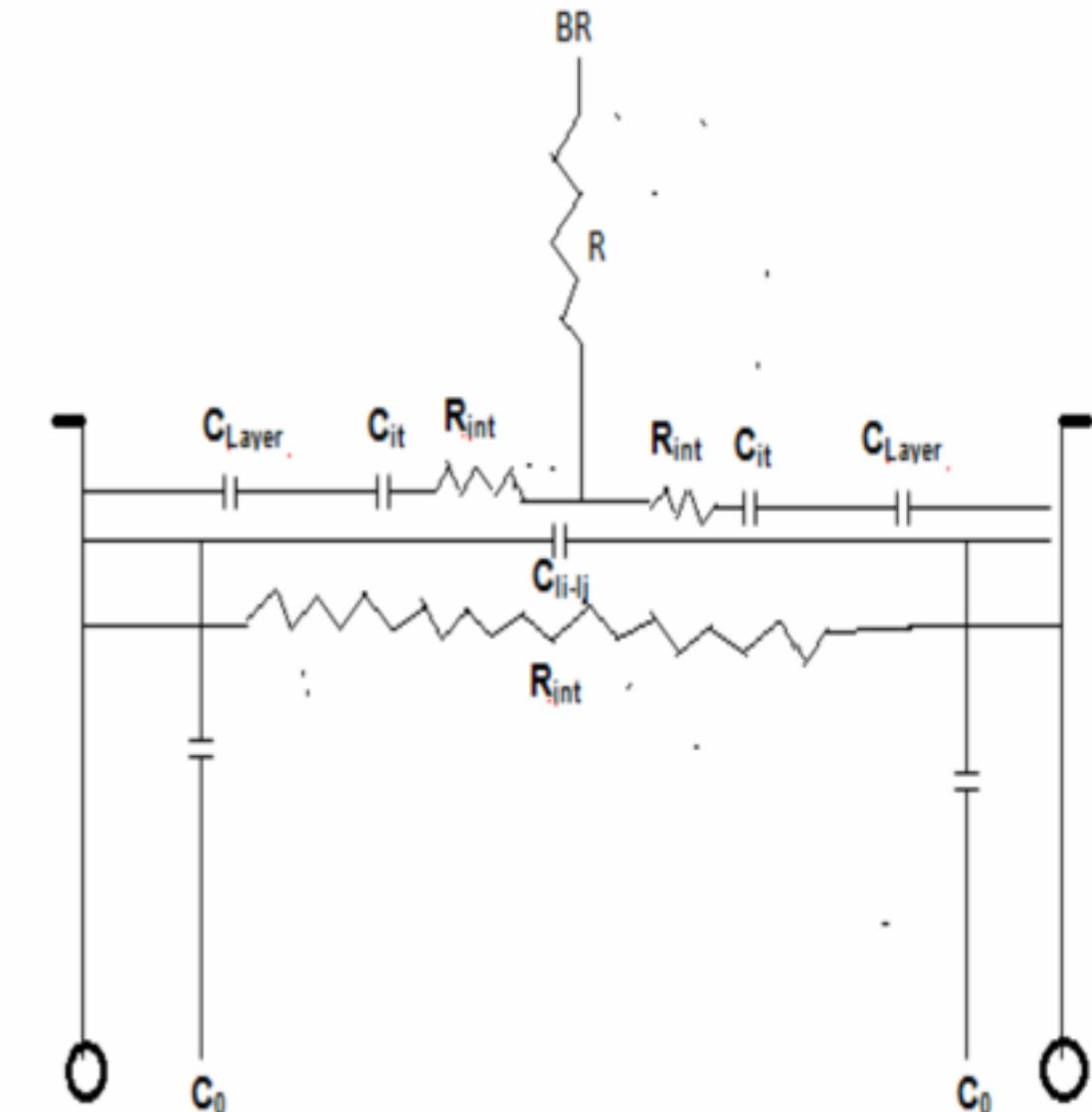


Figure 2: Two lumped model for the C/V curve description of two adjacent strip/pixel sensor.

Figure 1: Layout of a 5x5 sensor pixel array. The simulated portion is marked by a red colour.

# Simualtion Results

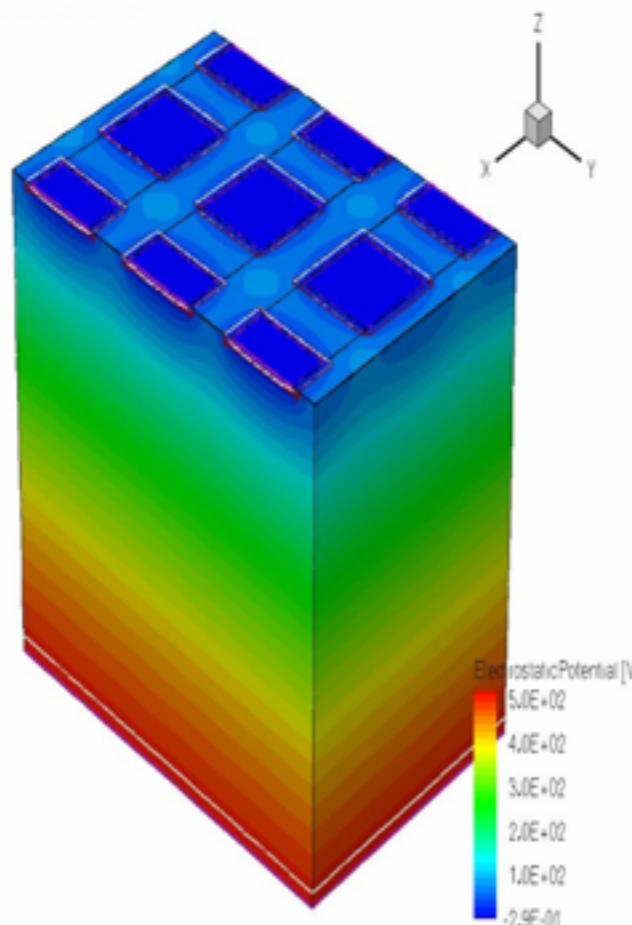


Table.1: Comparison of 3-D simulated capacitances with analytical expressions for  $80 \mu\text{m}$   $\text{p}^+$  pixel gap.

Capacitances	Analytical calculation [4] in fF	Simulation [S] in fF	Error [%]
$C_0$	7.68	8	5
$C_1$	5.099	5.48	22
$C_{\text{diag}}$	1.355	1.94	34

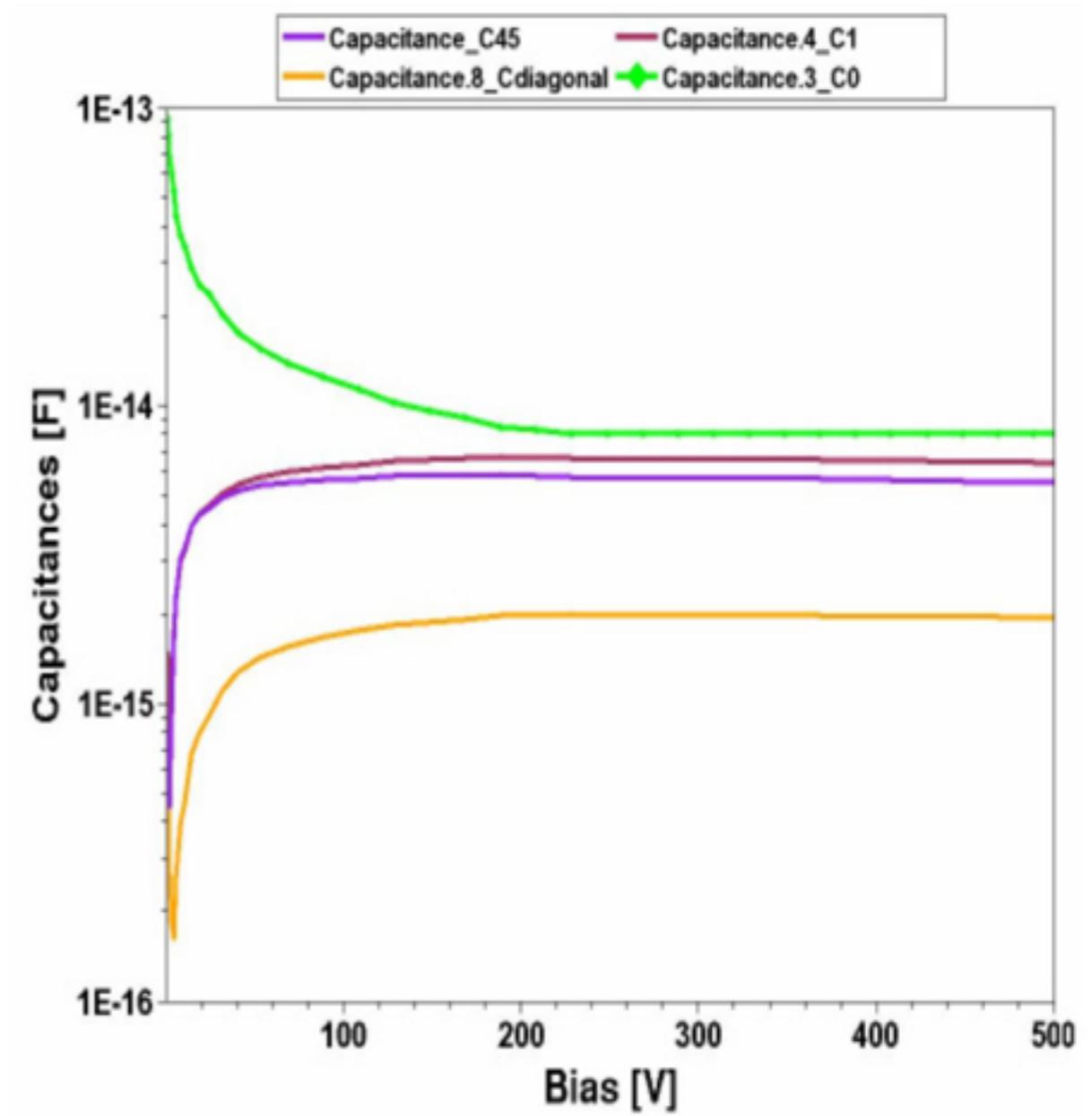


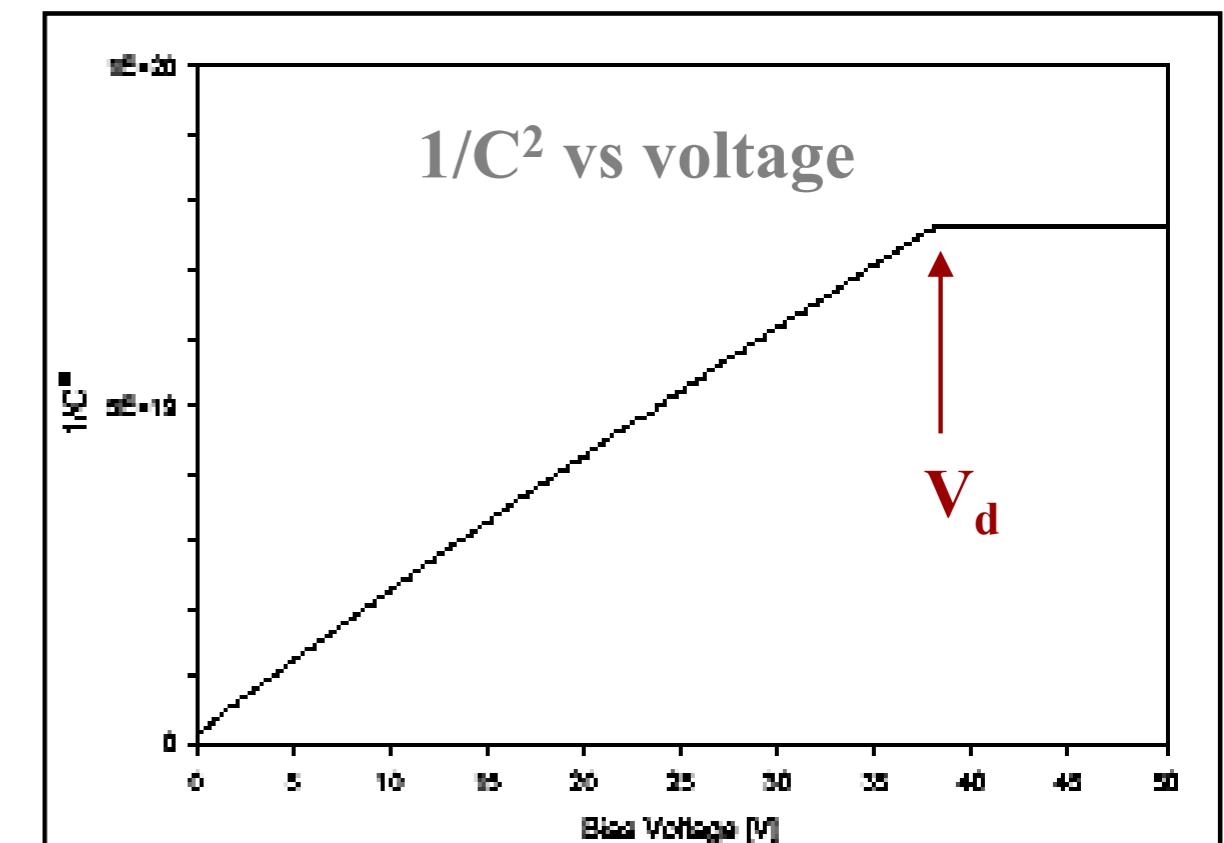
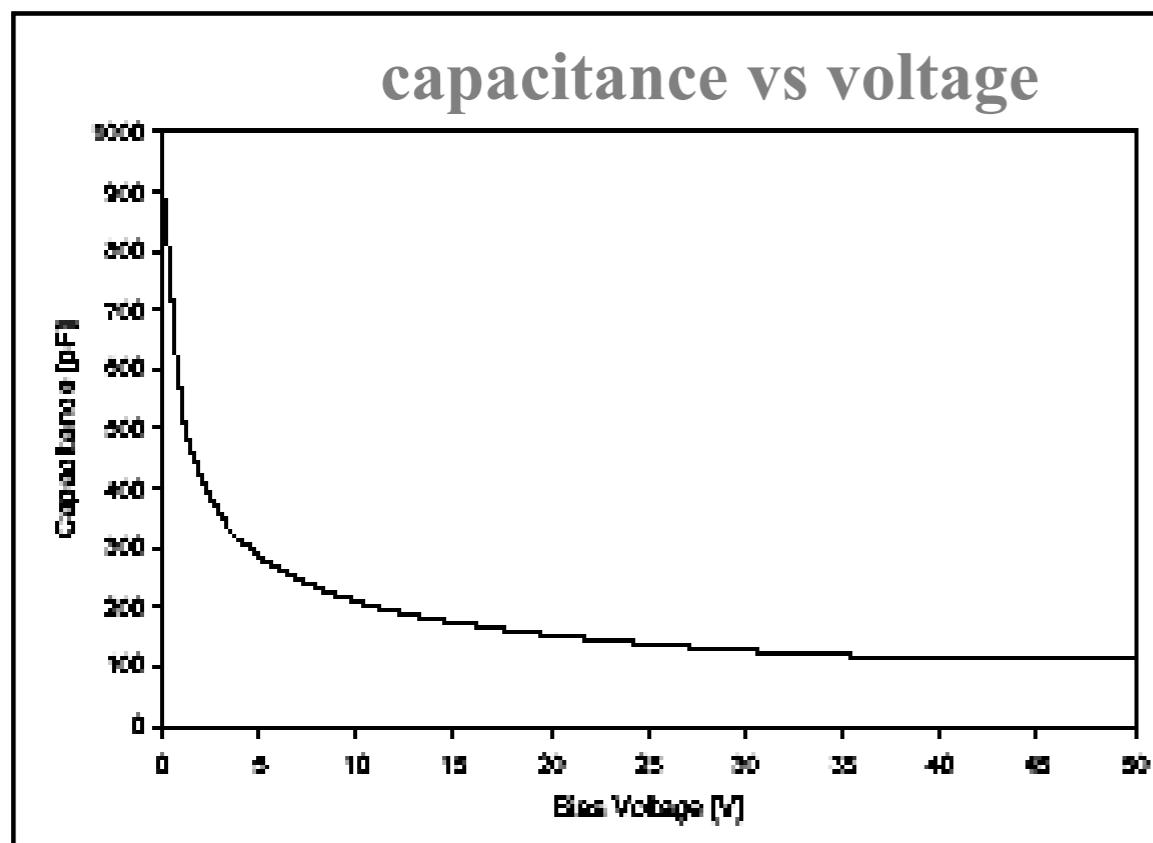
Figure.4: Interpixel capacitance as a function of the applied voltage at 1 MHz.

# Detector Capacitance

- > Capacitance is similar to parallel-plate capacitor
- > Fully depleted detector capacitance defined by geometric capacitance

$$C = \sqrt{\frac{\epsilon_0 \epsilon_r}{2\mu\rho|V|}} \cdot A$$

One normally measures the depletion behaviour (finds the depletion voltage) by measuring the capacitance versus reverse bias voltage.



# Capacitance measurements

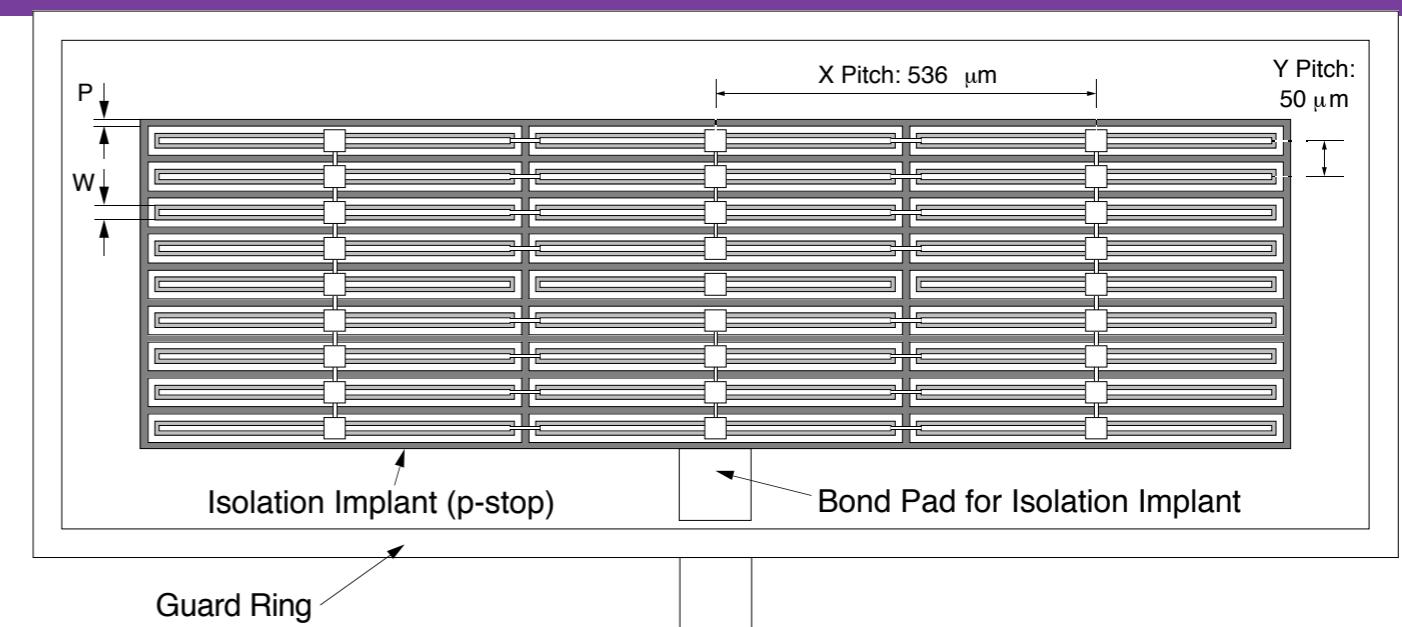
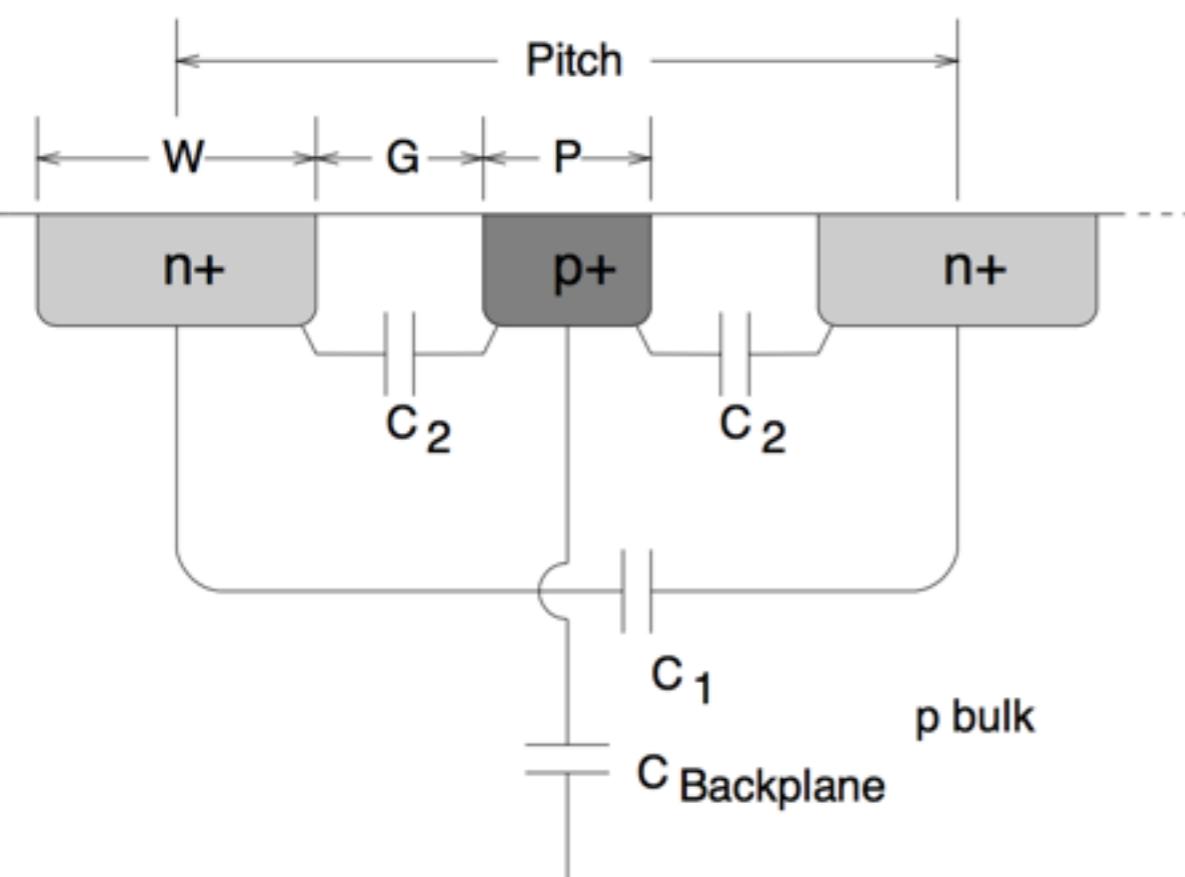
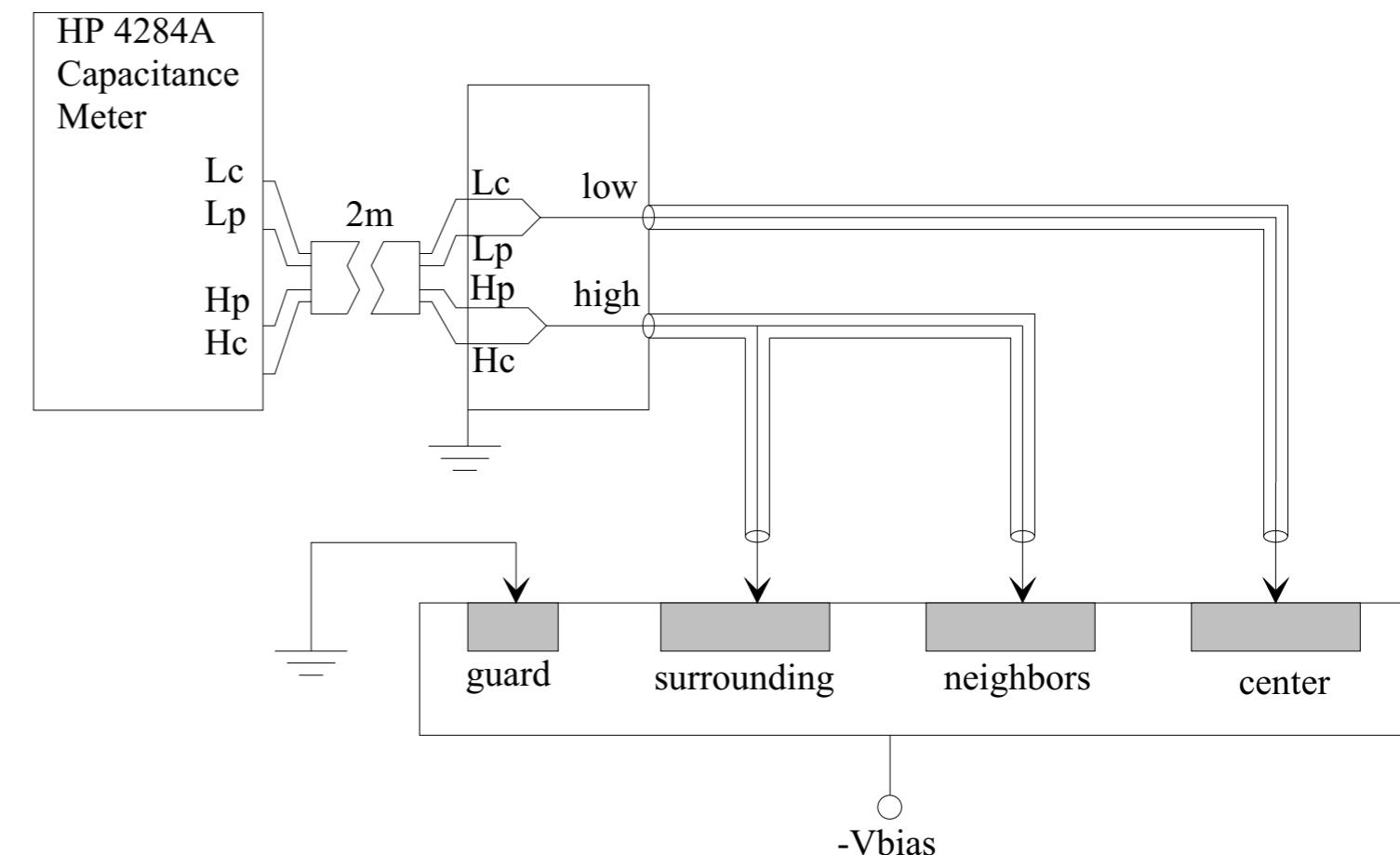


Figure 2: One pixel array in the LBNL *p*-type test structure.



# Results

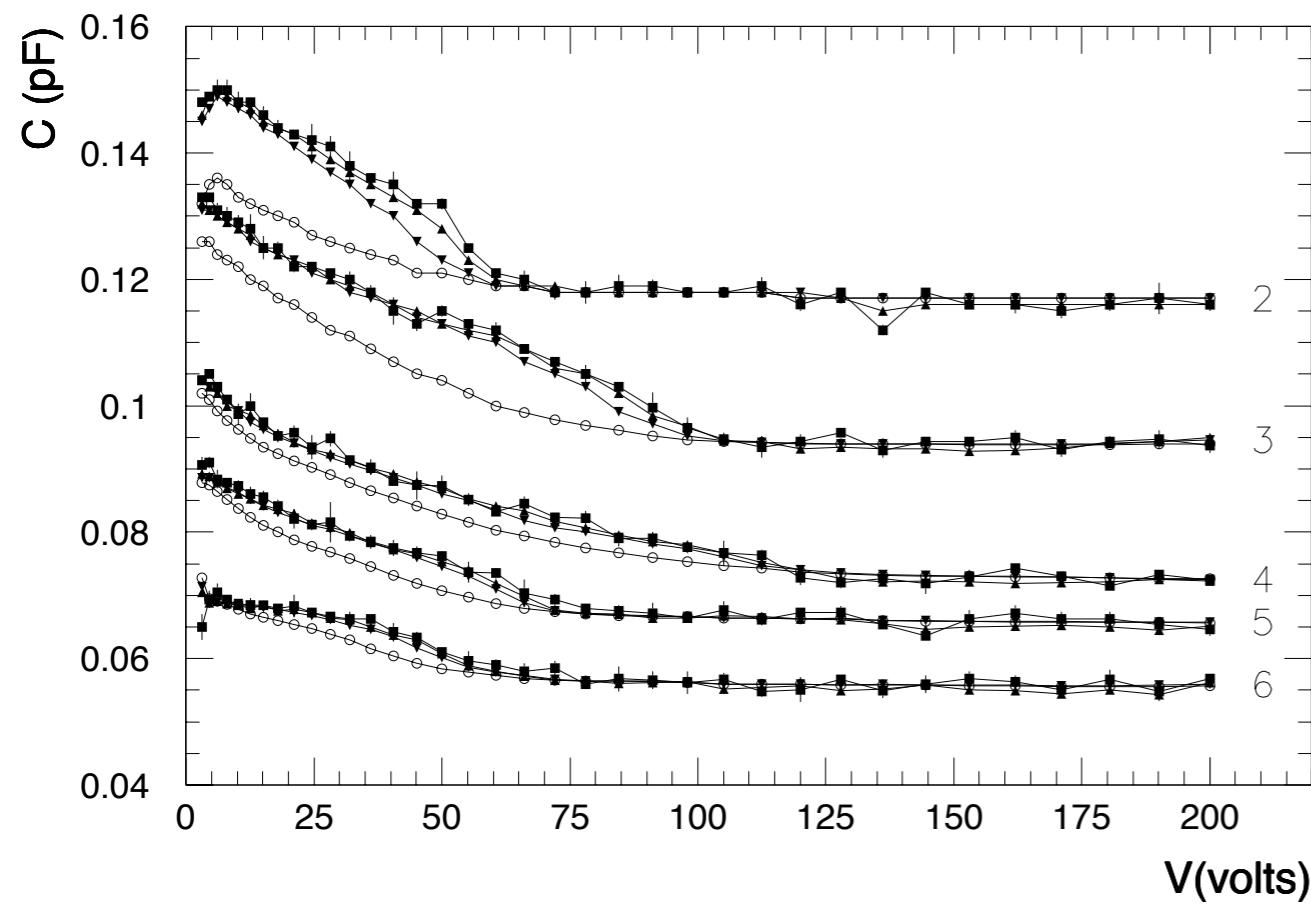


Figure 9: Inter-pixel capacitance of an unirradiated *n*-type LBNL detector. The multiple curves represent measurements made at frequencies 3, 10, and 100 kHz and 1 MHz. The families of curves labelled 2–6 show measurements of the arrays with the corresponding numbers (see Table 1 for their characteristics). The open circles are the 1 MHz data.

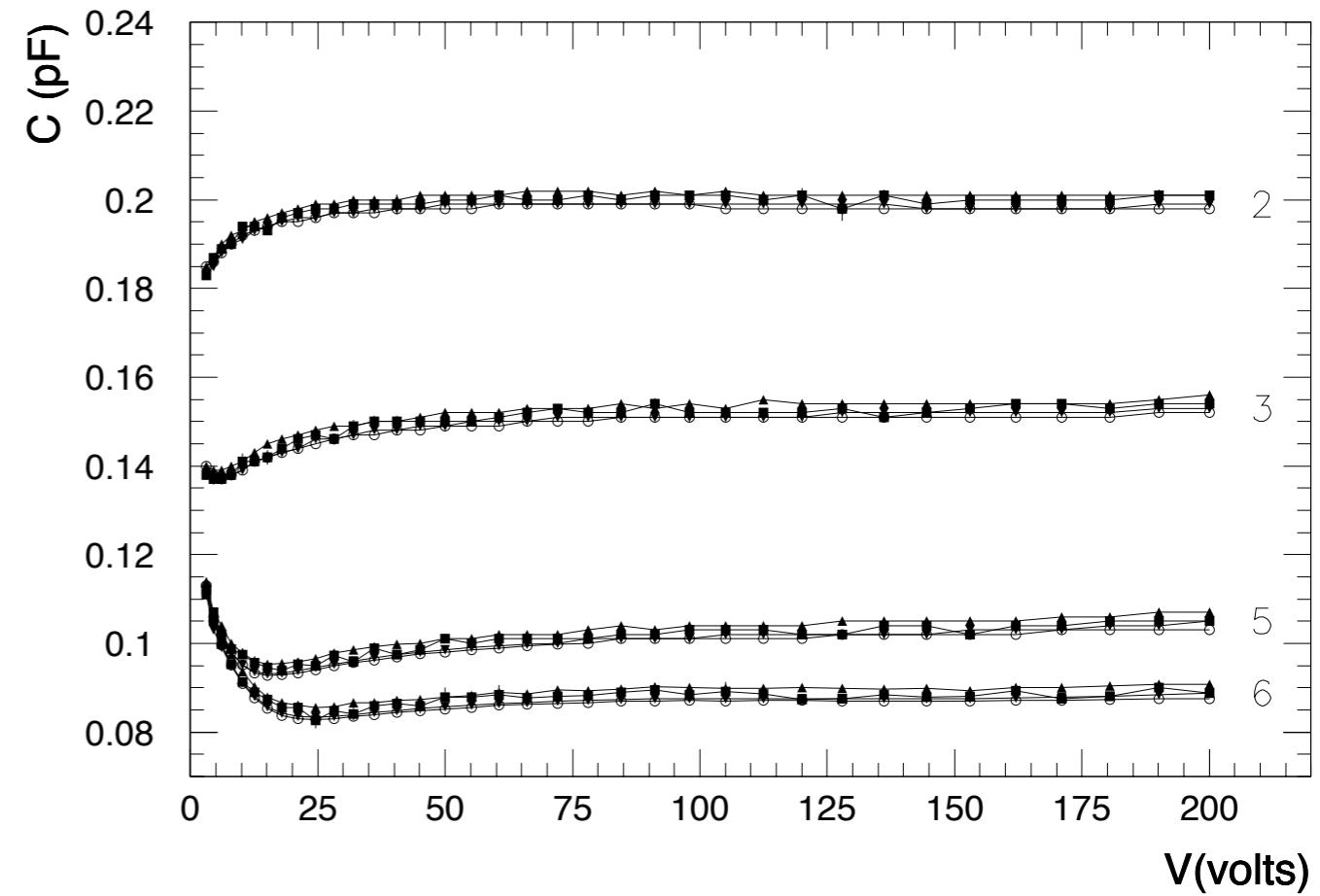
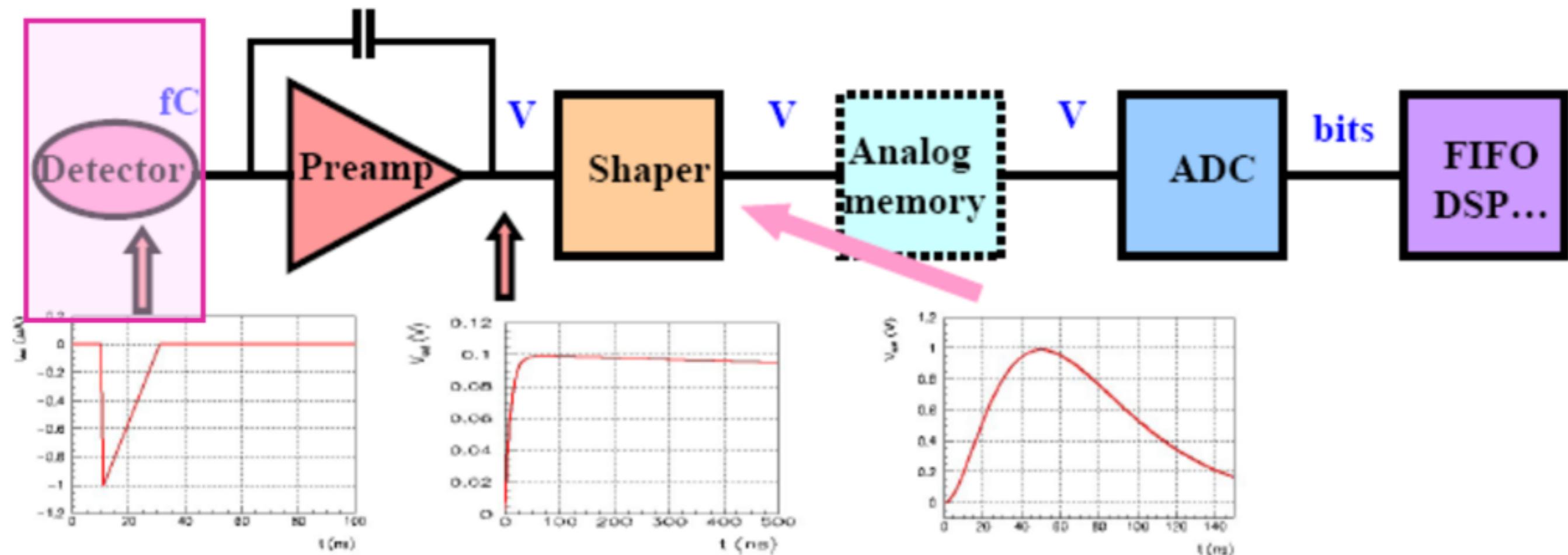


Figure 10: Inter-pixel capacitance of an unirradiated *p*-type LBNL detector. The multiple curves represent measurements made at frequencies 3, 10, and 100 kHz and 1 MHz. The families of curves labelled 2, 3, 5, and 6 show measurements of the arrays with the corresponding numbers (see Table 1 for their characteristics).

# Readout chain

Most front-ends follow a similar architecture



- Very small signals (fC) -> need **amplification**
- Measurement of **amplitude** (ADCs)
- Thousands to millions of channels