

# Practical Statistics for Particle Physics



Daniel Whiteson, UC Irvine  
HCPSS, 2014: Lecture 2

# Outline

- I. Mathematical preliminaries
- II. Fitting
- III. Data models
- IV. Hypothesis testing
- V. Tools and examples

# Models

Full MC

Fast MC

Effective models

Data-driven models

# Uncertainties

We have a recipe for

$f(\text{data} \mid \text{theory})$

But is it right?

# Uncertainties

We have a recipe for

$f(\text{data} | \text{theory})$

But is it right?

Theory has **lots**  
of **nuisance** parameters:  
cross-sections, LO, NLO...  
showering details  
hadronization details  
detector response

There is some point in NP  
space which gives  
the most accurate model  
but we don't know  
where it is!

# Systematics

## The Good

NP can be constrained in some control region. Uncertainty decreases with luminosity.

eg. Background cross-section

B-tagging efficiency

Jet energy scale

## The Ugly

Underlying theoretical approach

eg PYTHIA vs HERWIG

## The Bad

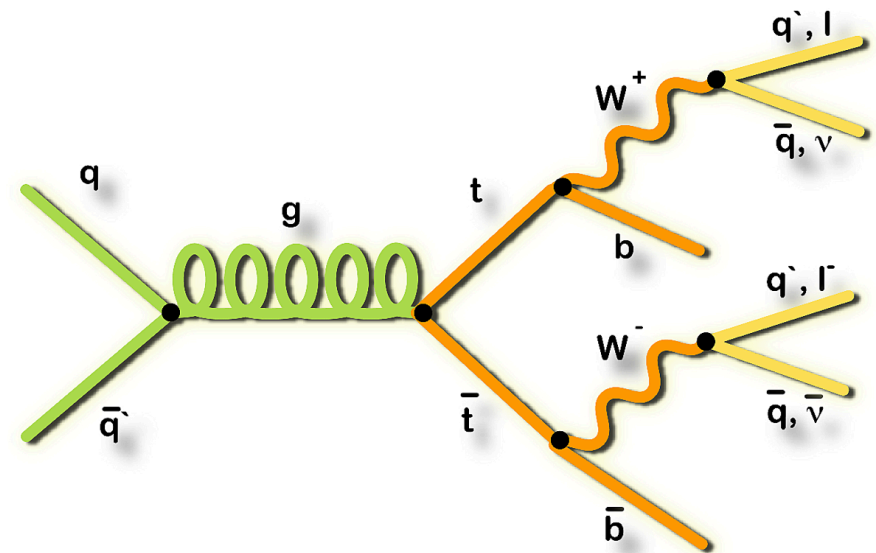
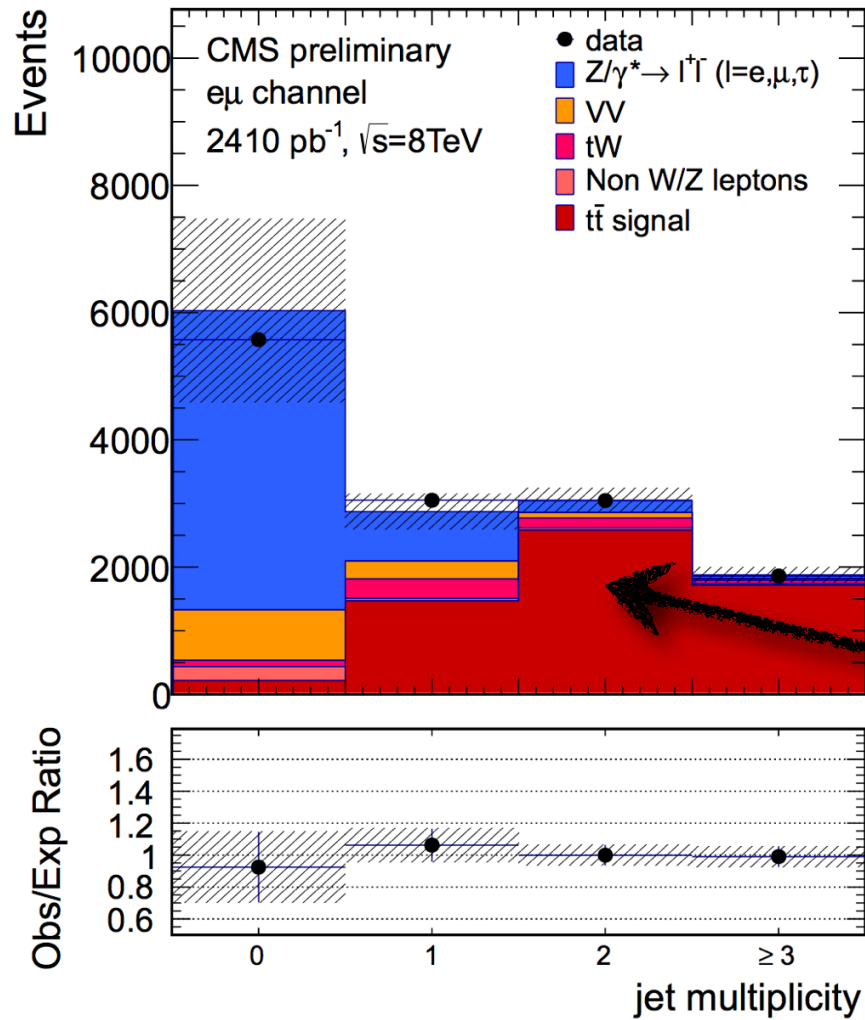
Parameters of underlying heuristic

eg PYTHIA tunes



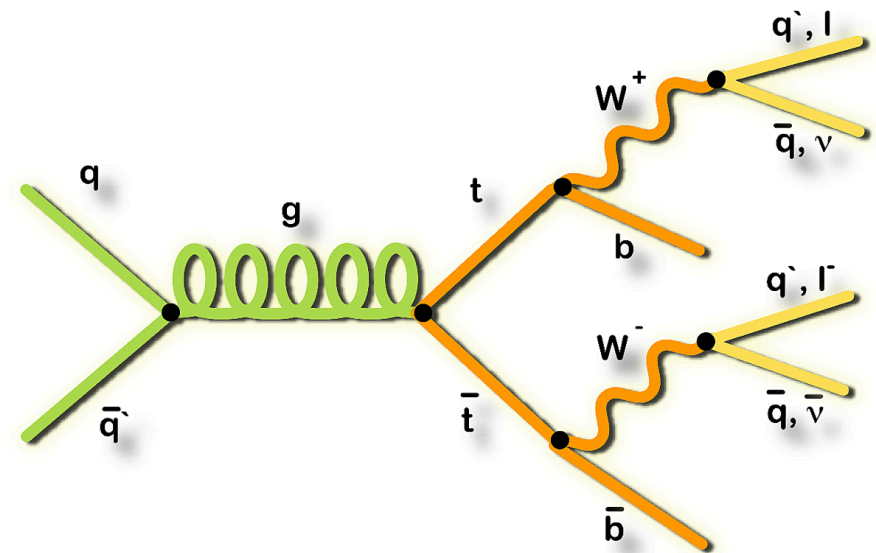
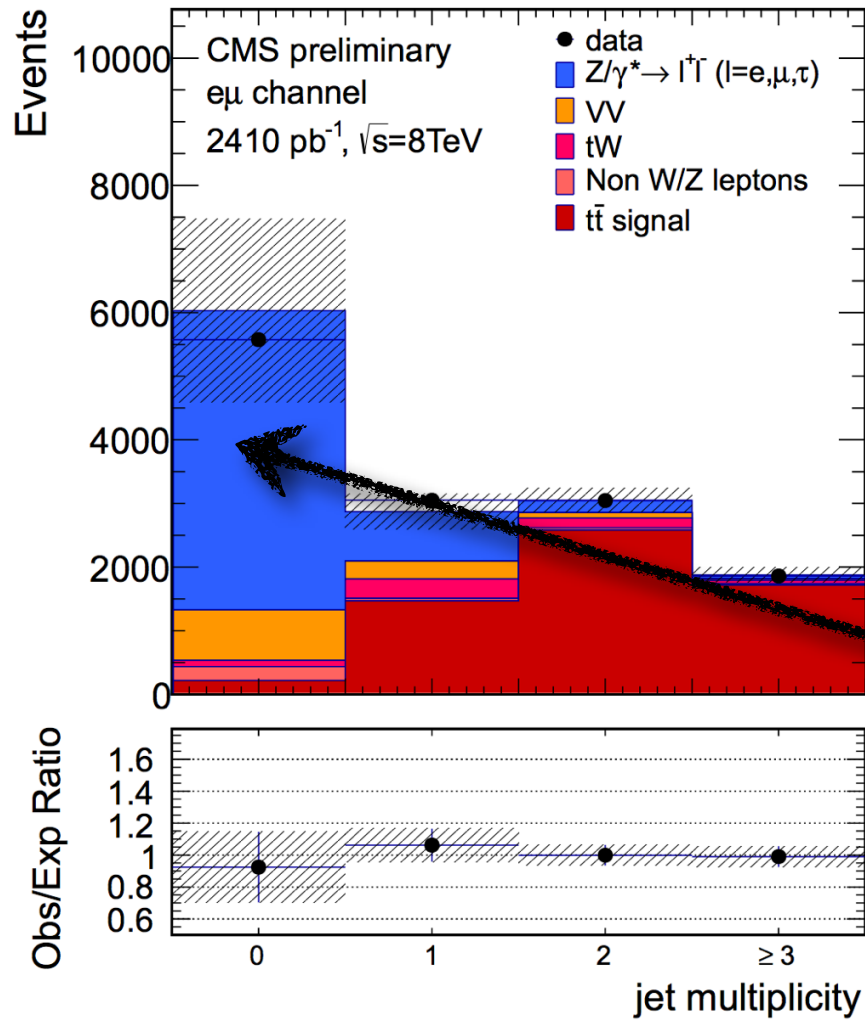
(Pekka Servino)

# Example



Signal region:  
 $\geq 2$  jets

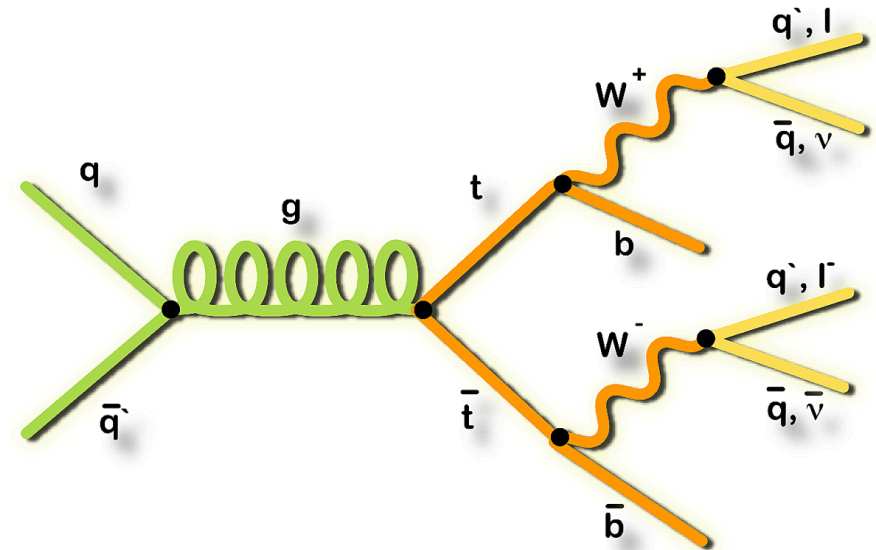
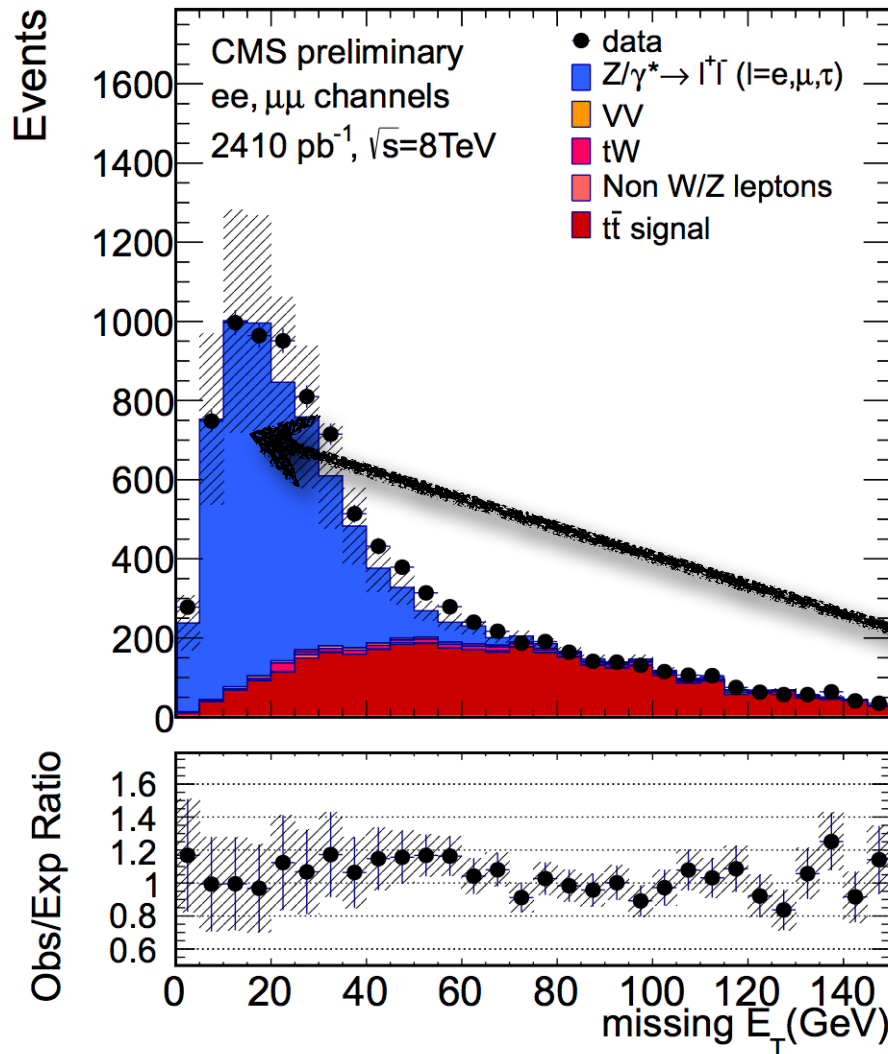
# Example



Z control region  
== 0 jets



# Example

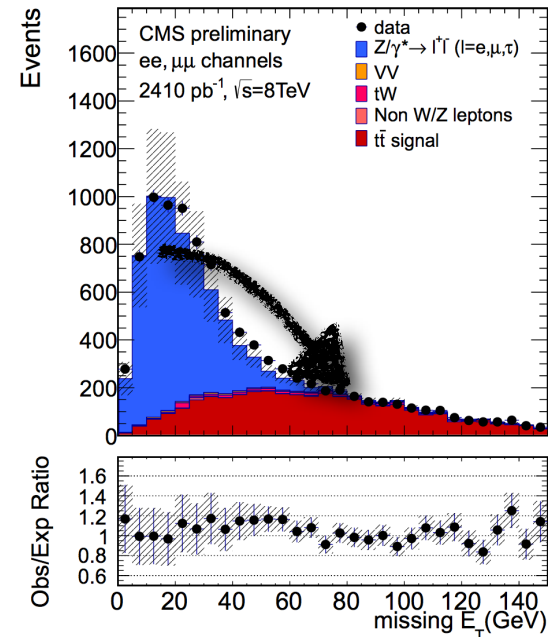
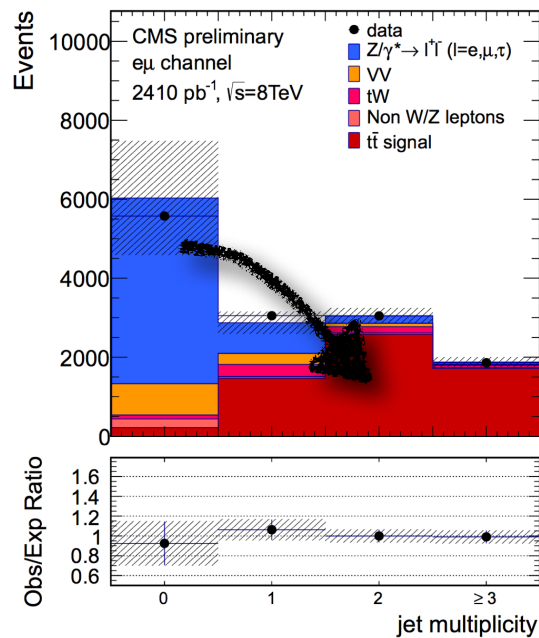


Z control region

Small Missing  $E_T$

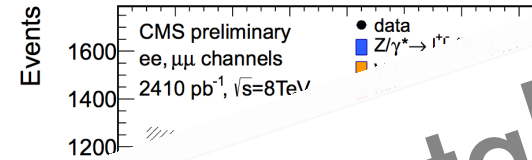
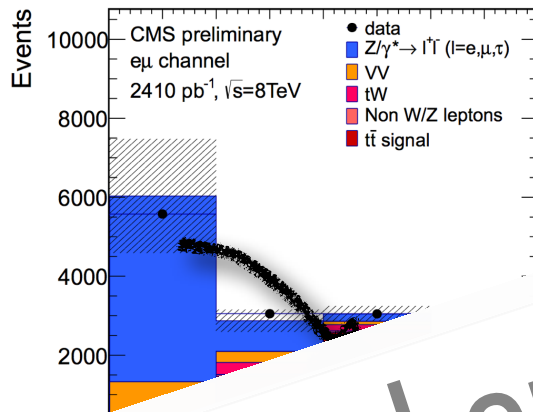


# Extrapolation



Measure background in one region, **extrapolate** knowledge to another. How well do we know rates of jet production, tails of MET?

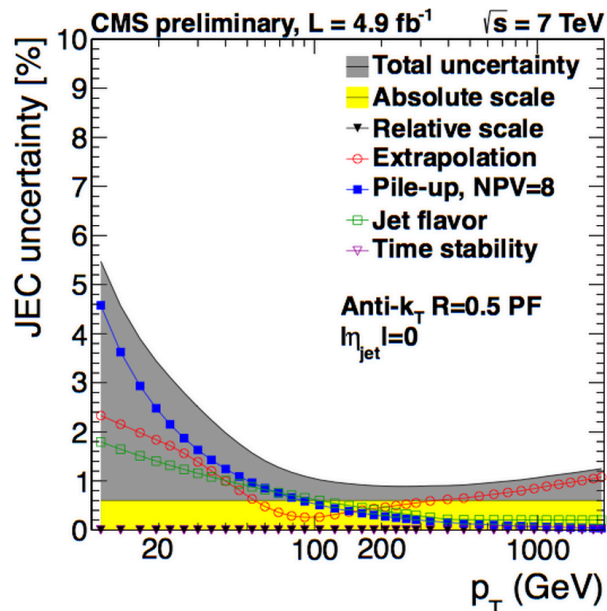
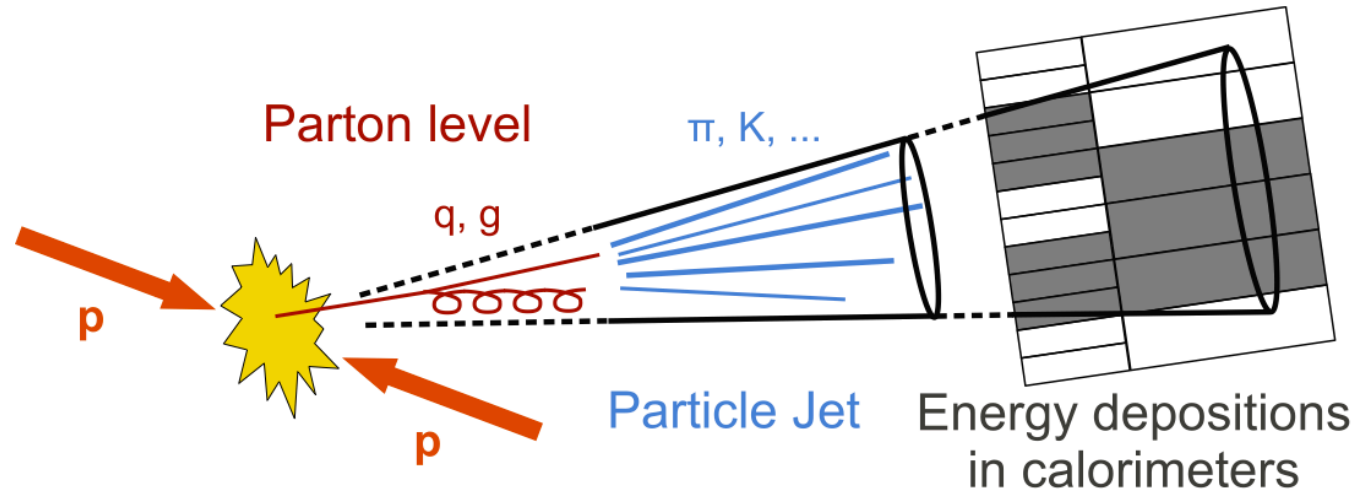
# Extrapolation



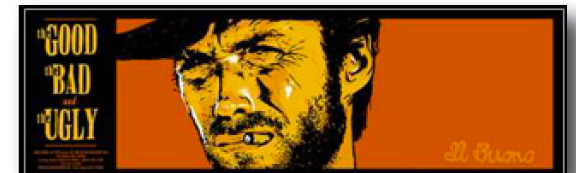
This is where experimental  
cleverness and creativity  
happens!

... and in one  
extrapolate knowledge  
another. How well do we know  
rates of jet production, tails of MET?

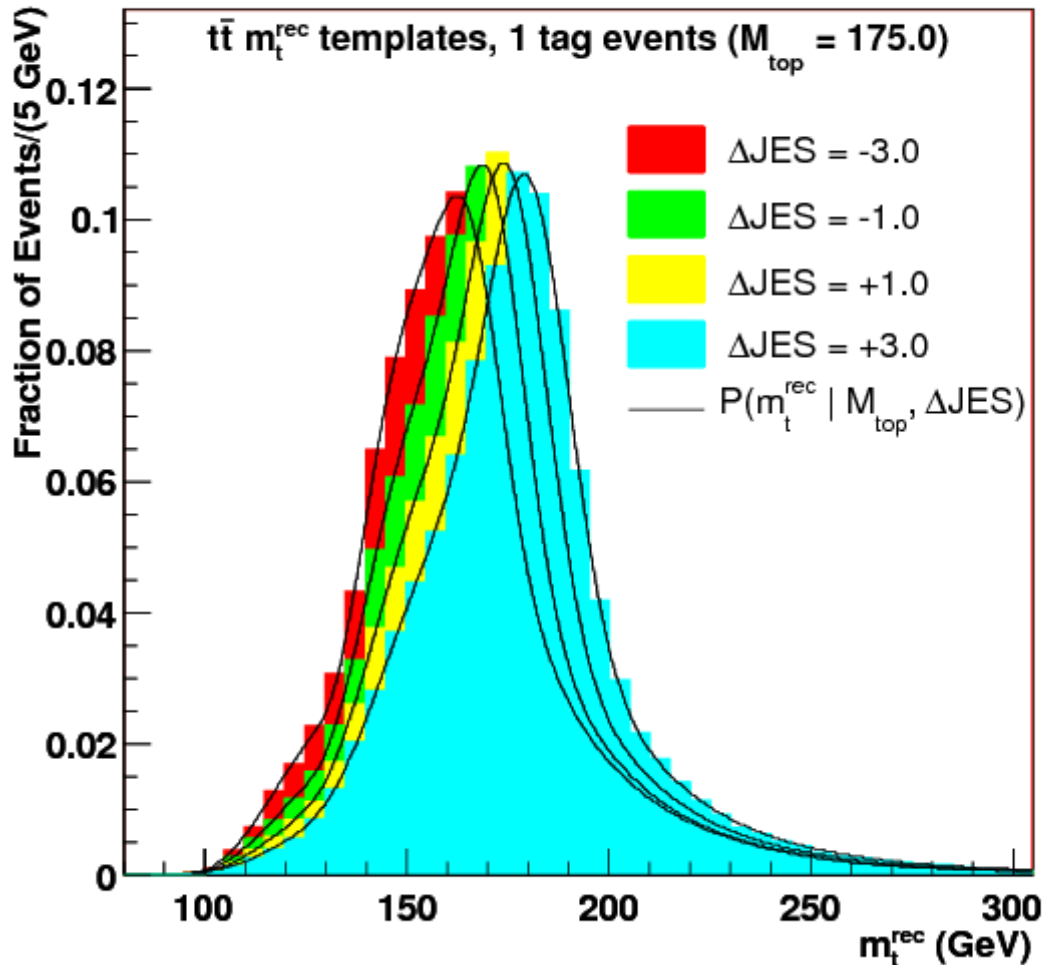
# Jet energy scale



Many steps in jet production  
 Lots of opportunities for mistakes  
 Calibrate in  $jj$ , photon+jet  
 Extrapolate to your dataset



# the shift method

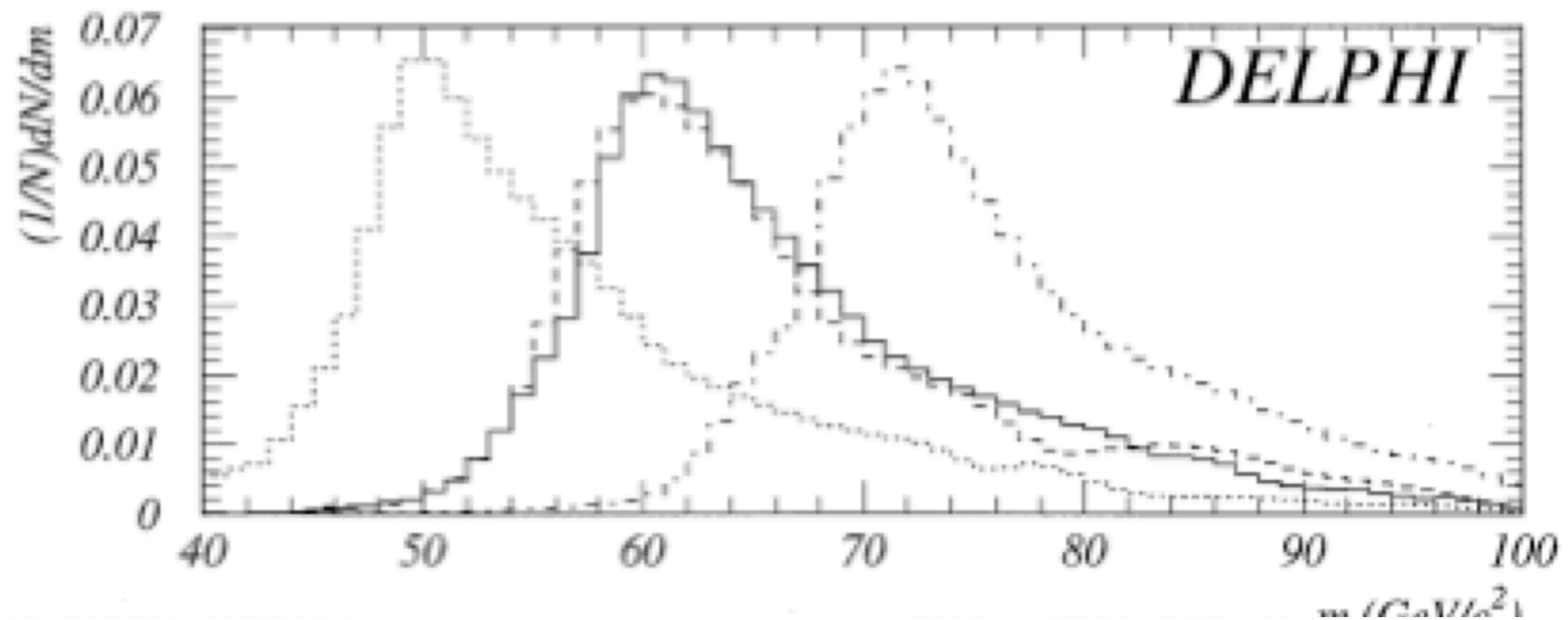


Generating samples at arbitrary values of NP can be expensive!

Often, just generate a few and interpolate.

# Histogram interpolation

*A.L. Read / Nuclear Instruments and Methods in Physics Research A 425 (1999) 357–360*



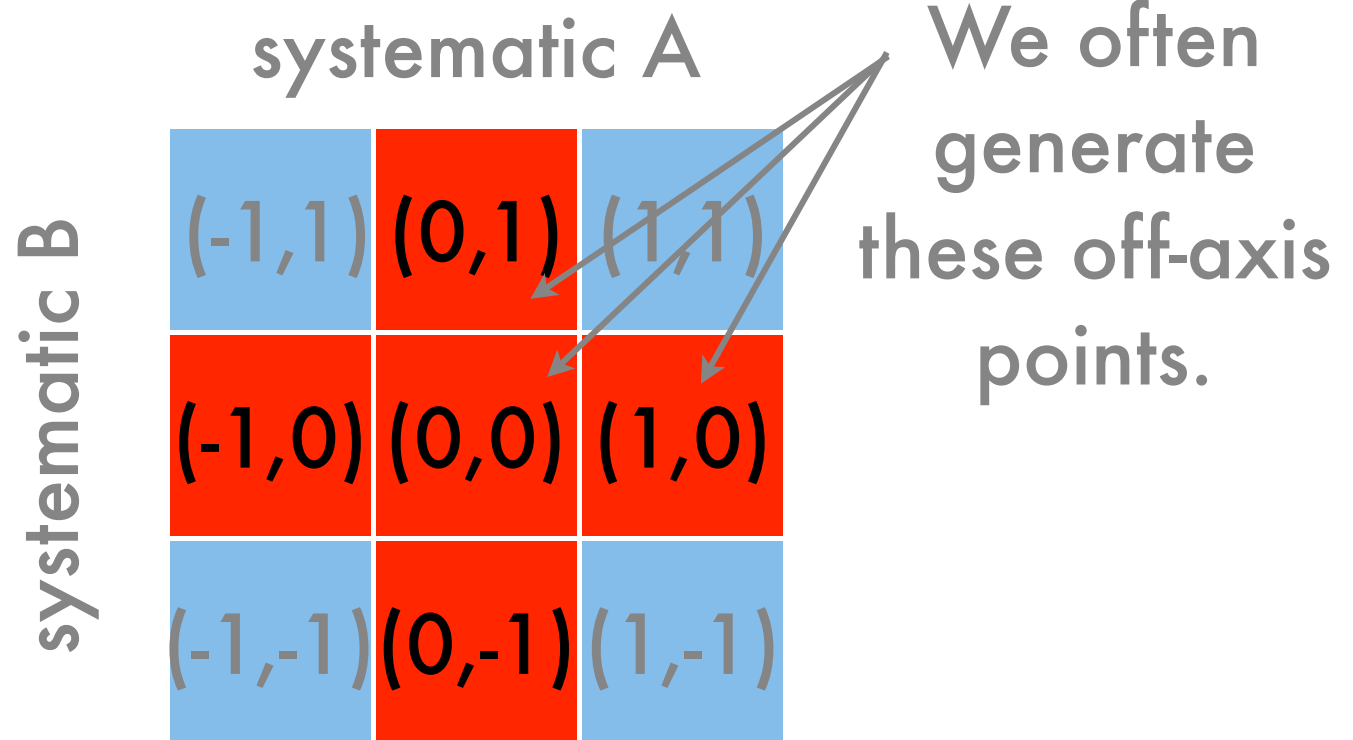
# dimensions

systematic A

systematic B

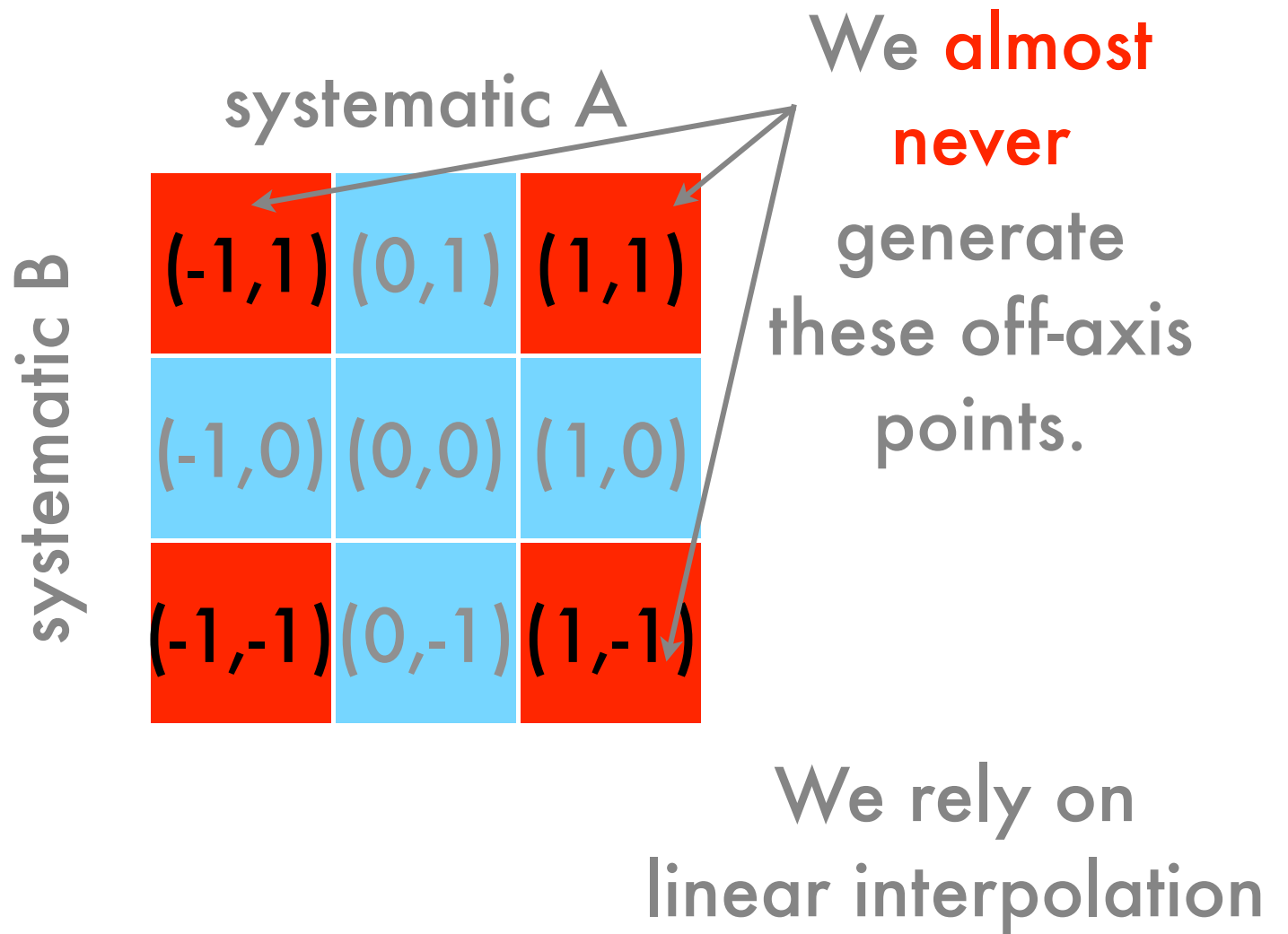
$(-1,1)$	$(0,1)$	$(1,1)$
$(-1,0)$	$(0,0)$	$(1,0)$
$(-1,-1)$	$(0,-1)$	$(1,-1)$

# dimensions

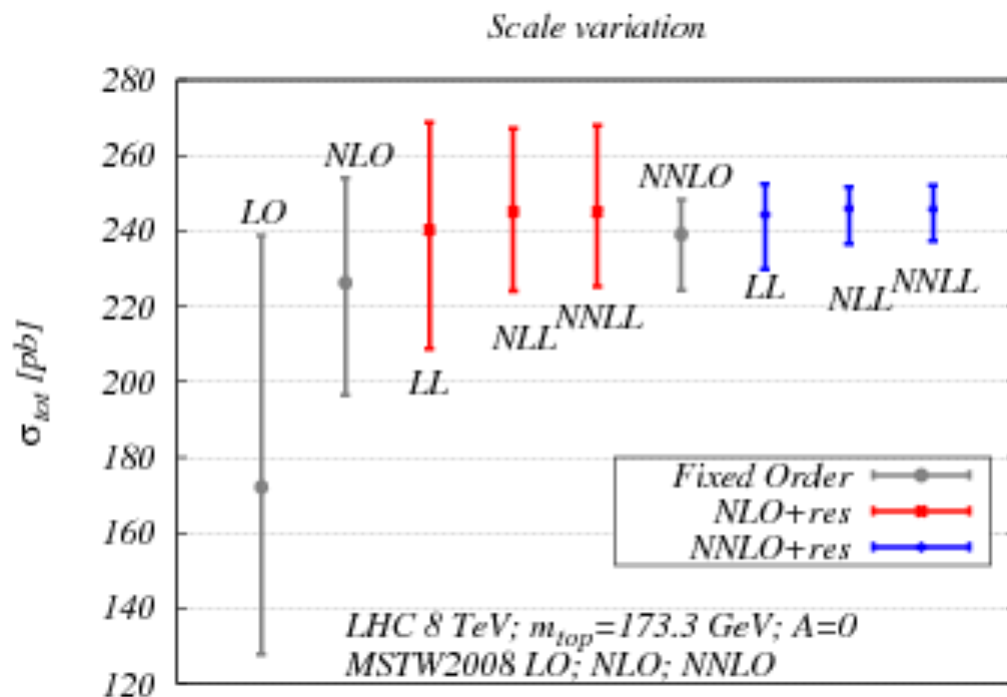




# dimensions



# Uncertainties



Uncertainty:  
shift renormalization,  
factorization scales  
by 2, 1/2  
measure change.

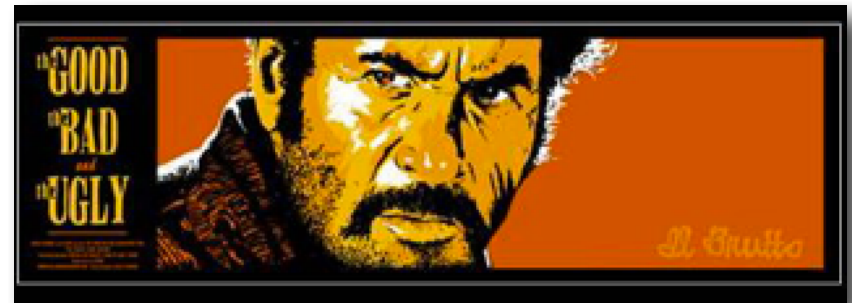
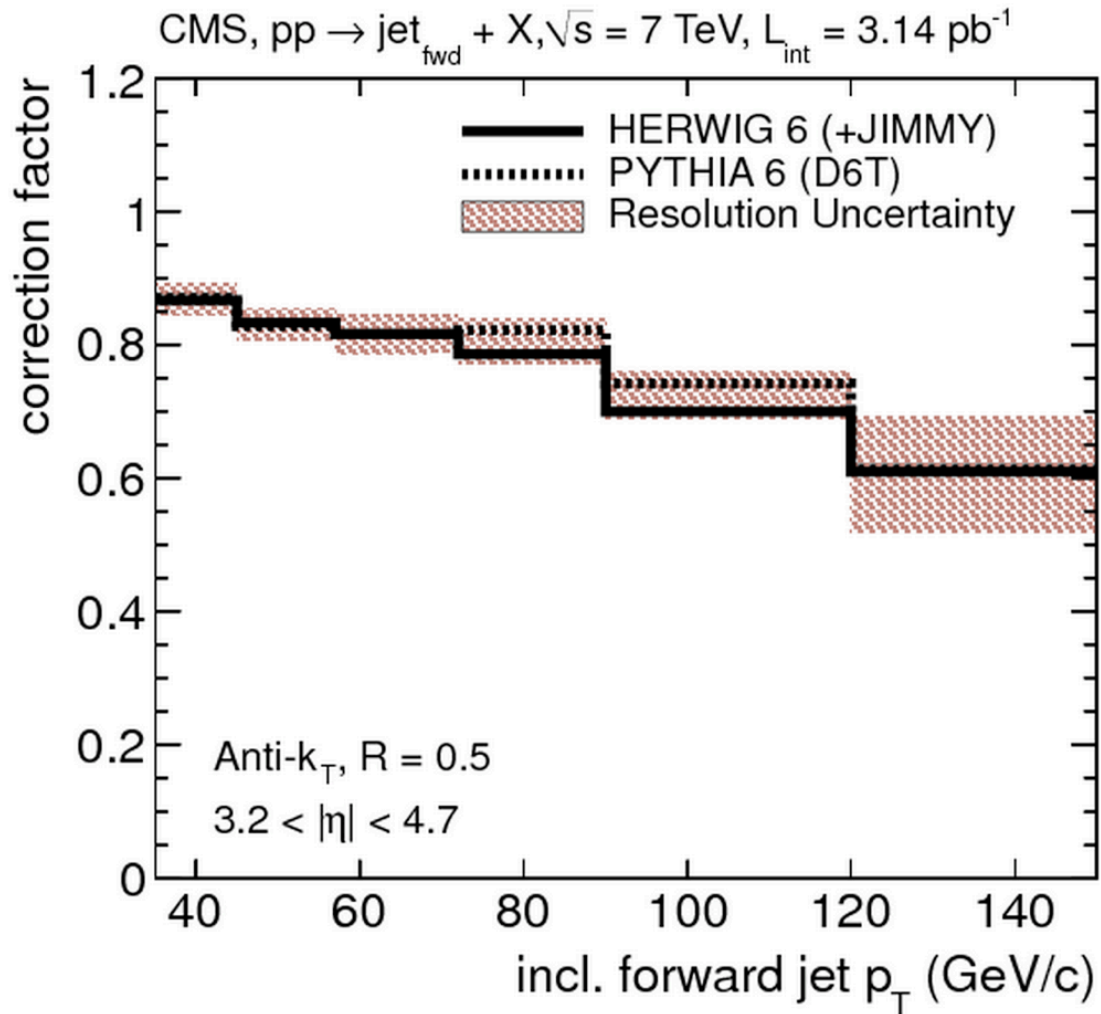
Why 2, 1/2?

Just convention

Not 1 sigma!



# Generators



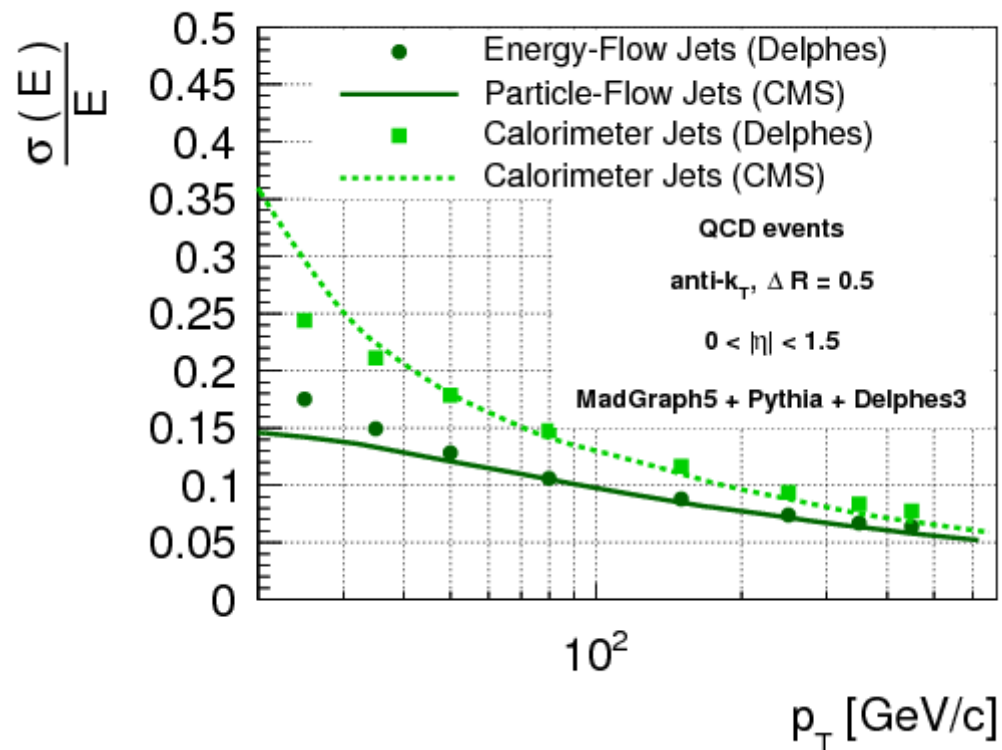
**fast-MC model**

# fast MC model



**DELPHES**  
fast simulation

Begin with generated events  
but rather than simulating microphysics,  
**smear** particles according to **resolution**.

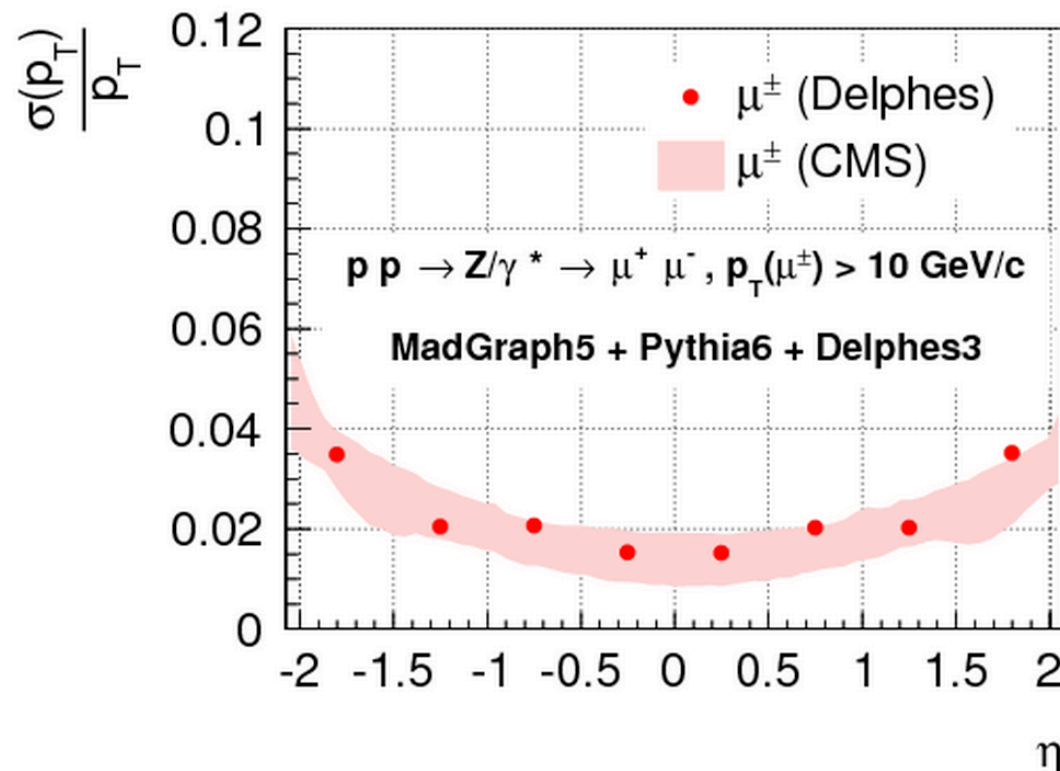


# fast MC model

Less accurate, but same issues as full MC:  
No analytic PDF  
Uncertainties in simulation



**DELPHES**  
fast simulation



# The dream

$f(\text{data} \mid \text{final-state particles } P)$

$\times f(\text{final state particles } P \mid \text{showered particles } S)$

$\times f(\text{showered particles } S \mid \text{hard scatter products } M)$

$\times f(\text{hard scatter products } M \mid \text{theory})$

Sum over all possible intermediate  $P, S, M$

# ME approach

If we have a parametrized detector response,  
can we parameterize

$$f(\text{data} \mid \text{final-state particles } P)$$

$$\times f(\text{final state particles } P \mid \text{showered particles } S)$$

$$\times f(\text{showered particles } S \mid \text{hard scatter products } M)$$

$$\times f(\text{hard scatter products } M \mid \text{theory})$$





# ME approach

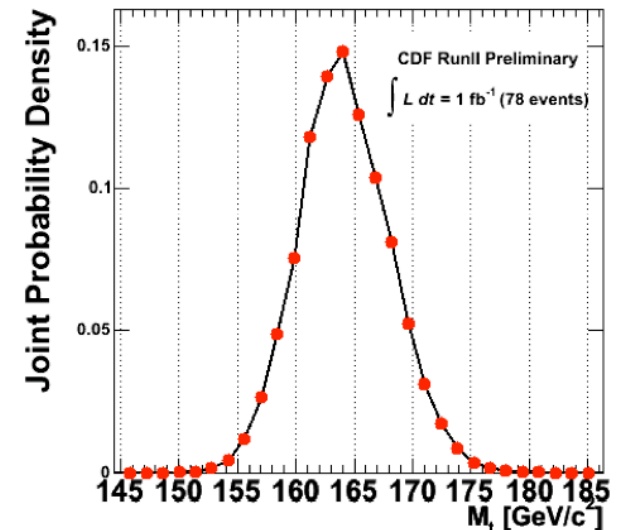
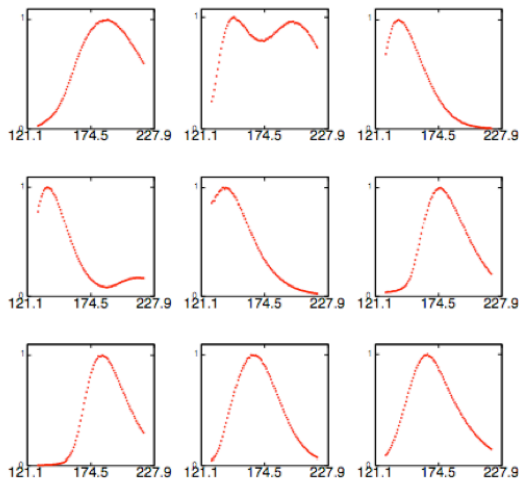
Yes we can!

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

Phase-space  
Integral

Matrix  
Element

Transfer  
Functions



# ME approach

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

Phase-space  
Integral

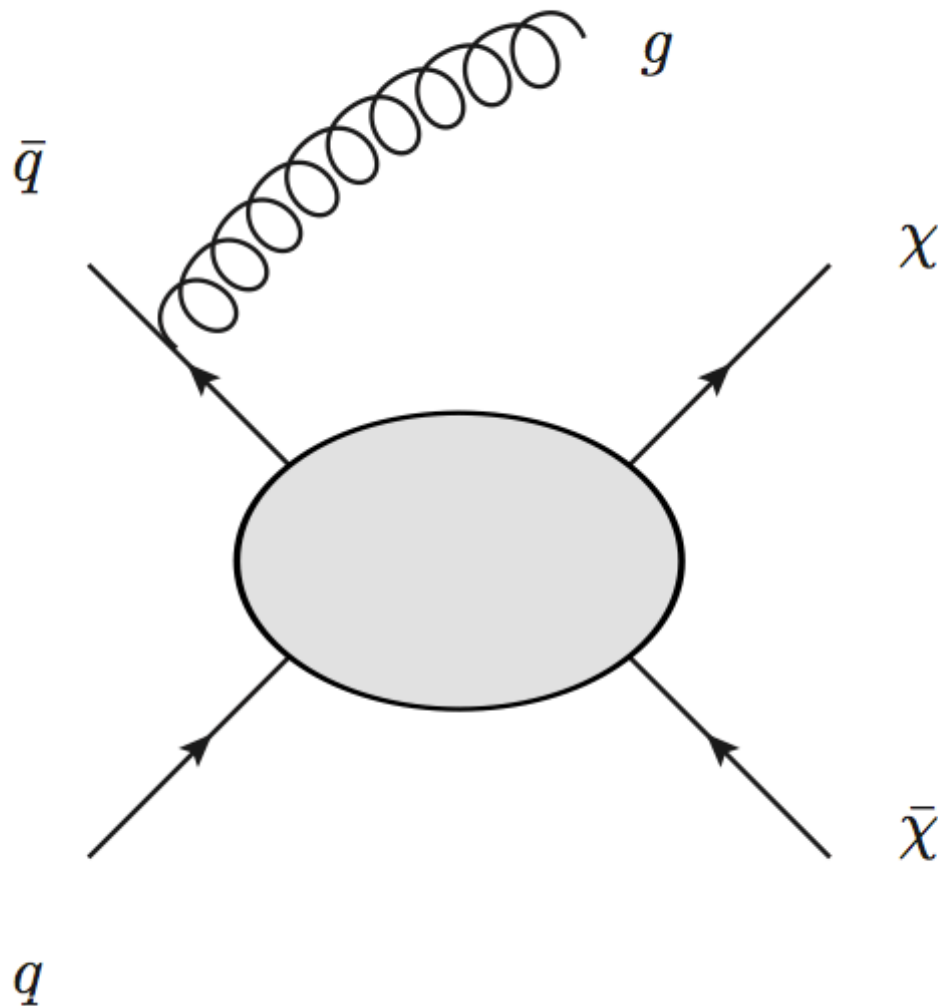
Matrix  
Element

Transfer  
Functions

Transfer functions reflect a very complex process  
By necessity, approximations, and therefore  
uncertainties.

**Data-driven model**

# Example: dark matter



Final state:

Two WIMPs+**jet**

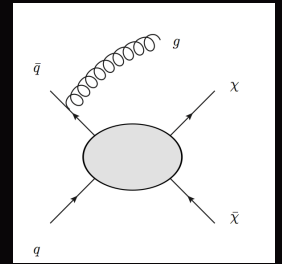
Detector signature

**Jet** + **MET**

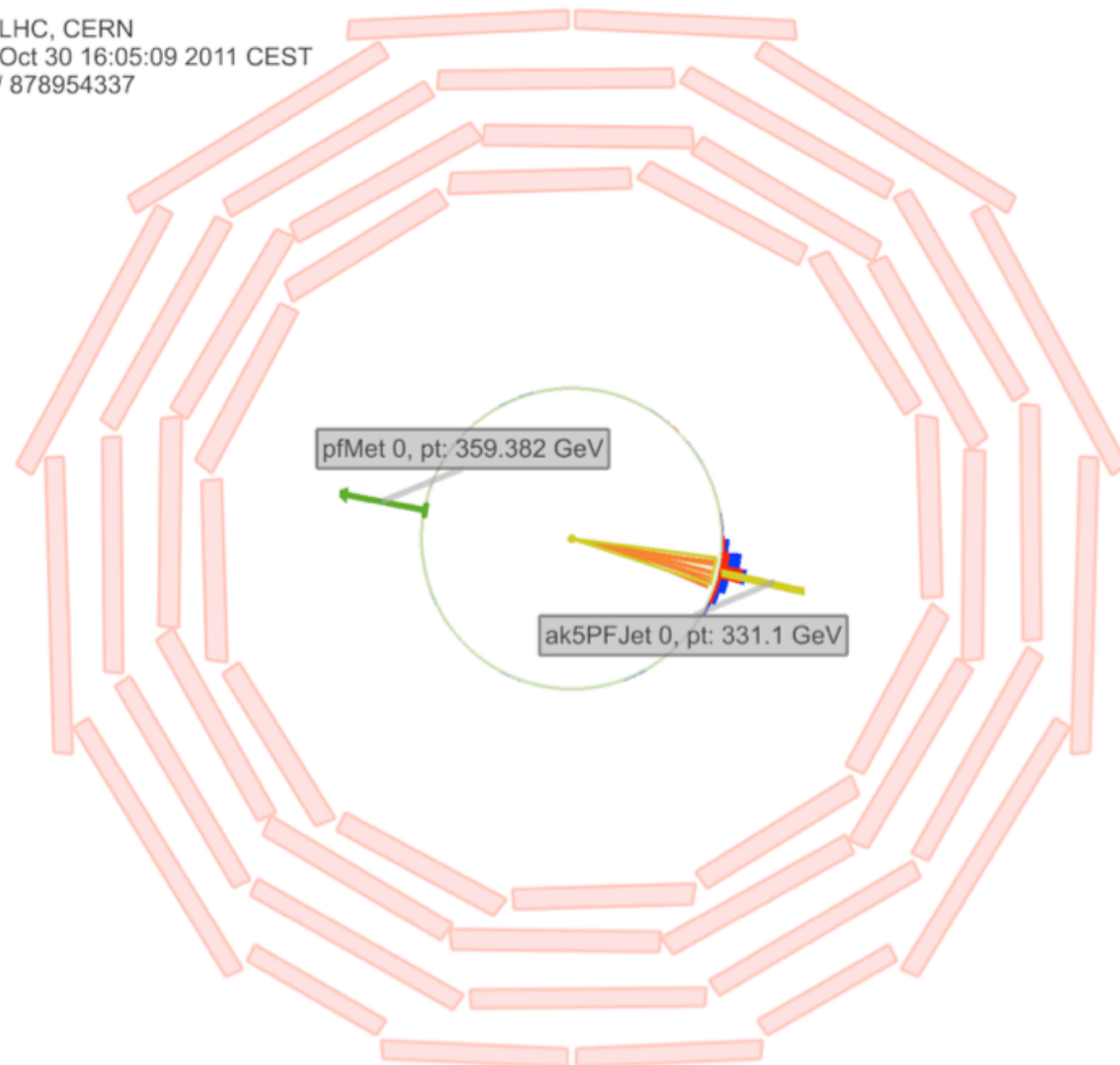
# Mono-jet



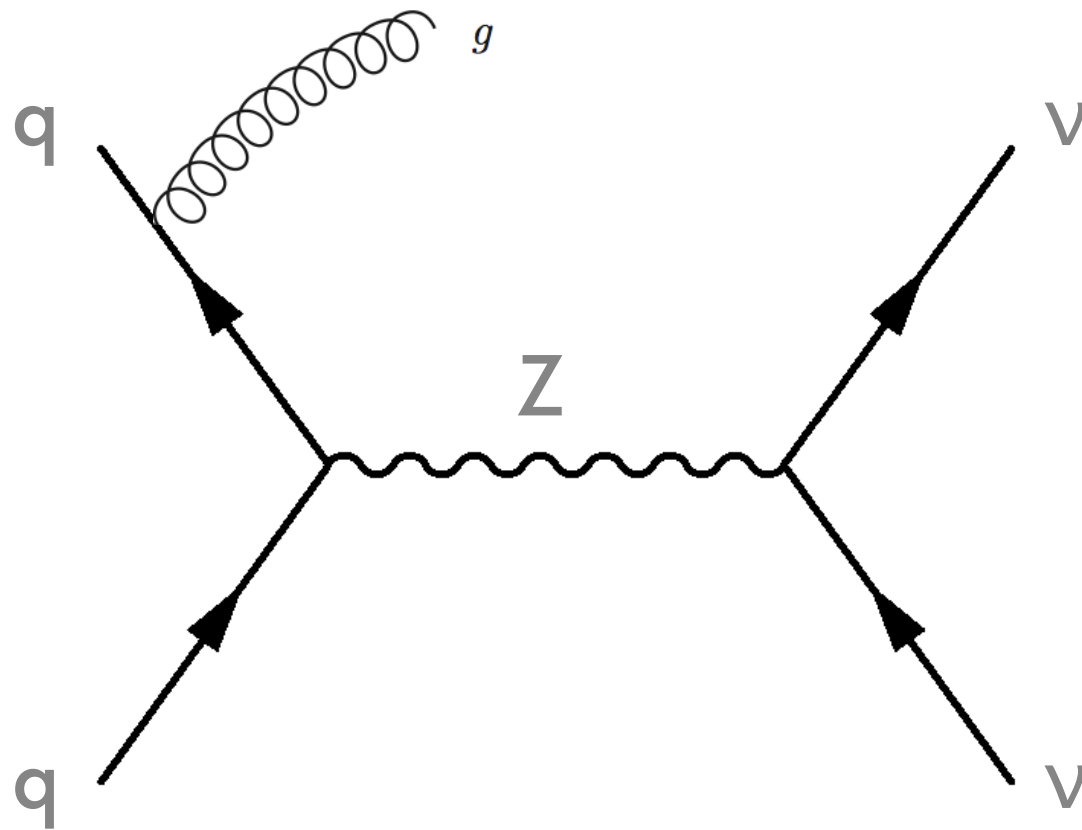
# Event display



CMS Experiment at LHC, CERN  
Data recorded: Sun Oct 30 16:05:09 2011 CEST  
Run/Event: 180250 / 878954337  
Lumi section: 481



# Backgrounds



Final state:

jet + MET

Process:

$Z \rightarrow \nu \nu$ , with jet

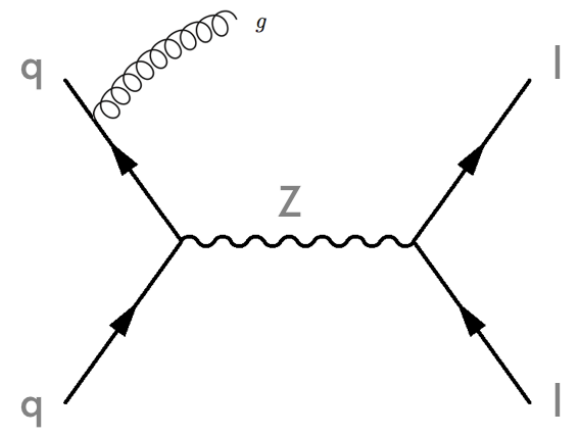
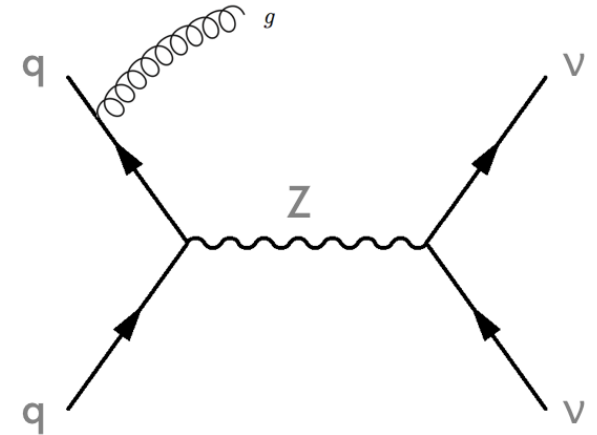
# Backgrounds

How to estimate?

Idea:  $Z \rightarrow \nu \nu$  from  $Z \rightarrow \parallel$

Approach:

- (1) measure  $Z$  to  $\parallel$  + jet
- (2) scale by known branching ratios





# Details

$$N[Z(\mathbf{v}\mathbf{v})] = N[Z(\mathbf{I}\mathbf{I})] \times \text{BF}[Z(\mathbf{v}\mathbf{v})] / \text{BF}[Z(\mathbf{I}\mathbf{I})]$$

# Details

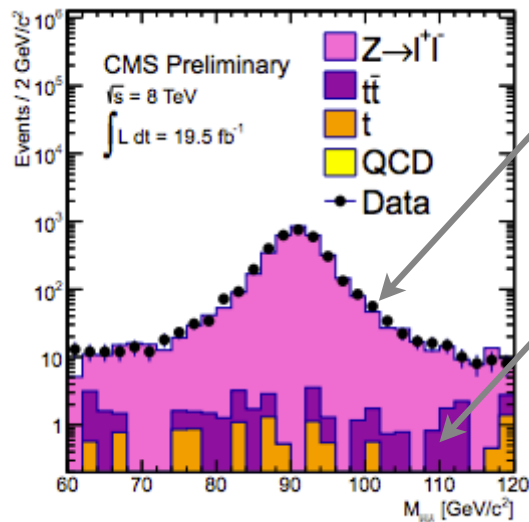
$$N[Z(\mathbf{vv})] = N[Z(\mathbf{ll})] \times \text{BF}[Z(\mathbf{vv})] / \text{BF}[Z(\mathbf{ll})]$$

$$N[Z(\mathbf{ll})] = N(\mathbf{ll}) - N(\text{bg}) / \varepsilon$$

# Details

$$N[Z(\nu\nu)] = N[Z(l\bar{l})] \times \text{BF}[Z(\nu\nu)] / \text{BF}[Z(l\bar{l})]$$

$$N[Z(l\bar{l})] = N(l\bar{l}) - N(\text{bg}) / \epsilon$$

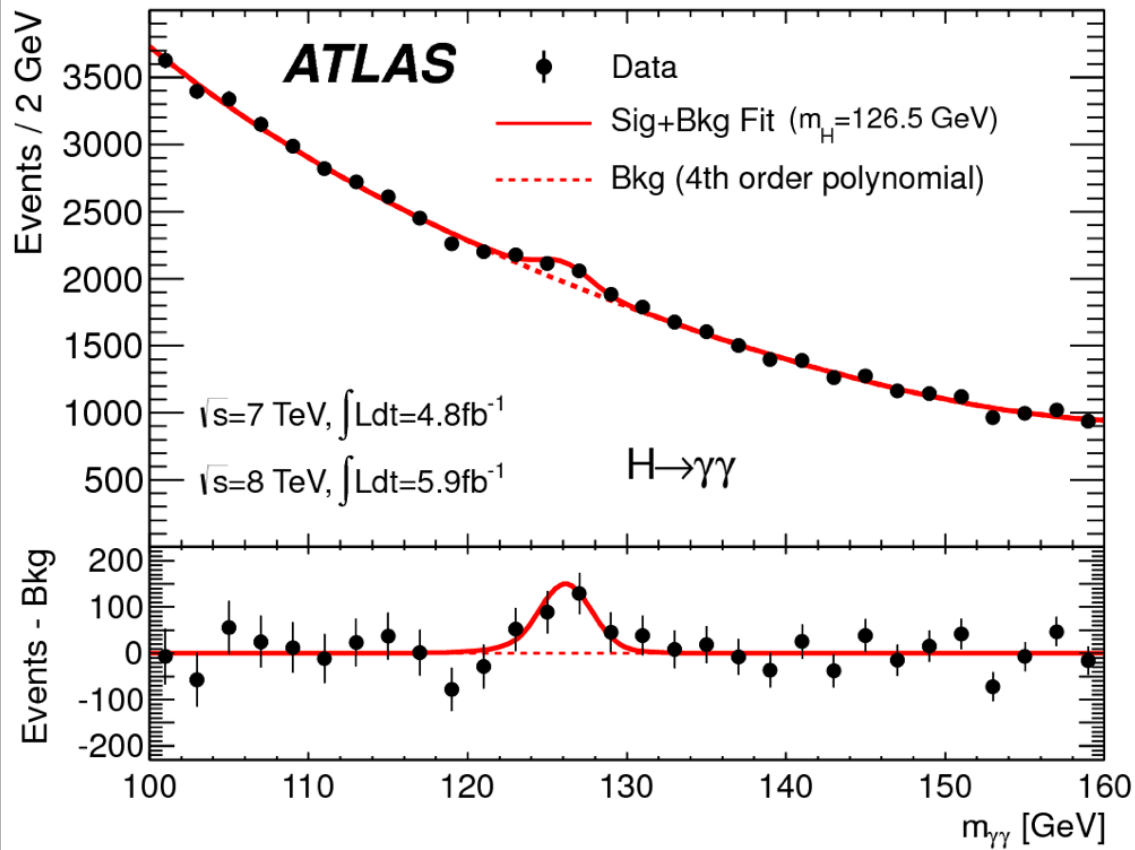


From simulation

CMS PAS EXO-12-048

# Effective Model

# Effective Model

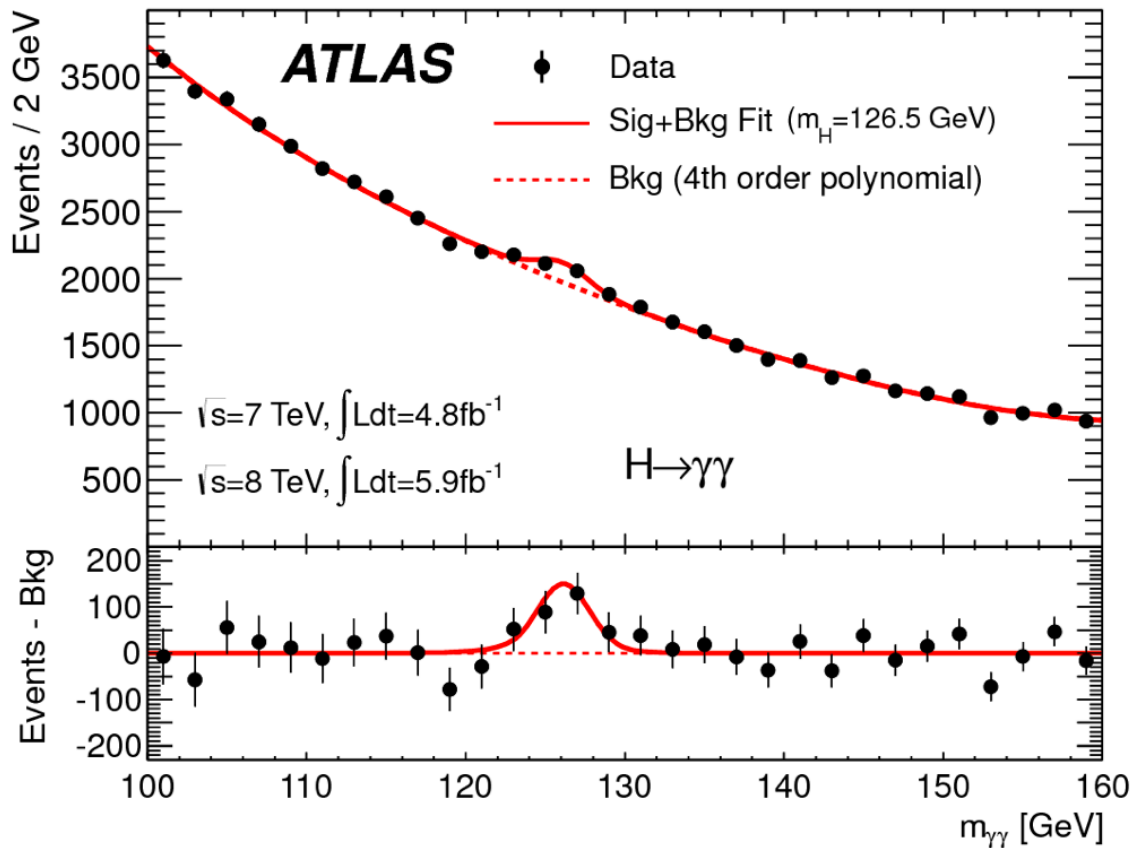


High-level arguments

dependence of background  
eg "smooth background"

shape of signal  
eg "narrow resonance"

# Effective Model



Uncertainties

Prediction under peak  
depends on quality of  
side band fits

Background function  
evaluate in control regions

# Summary of models

## MC simulation

- sample of events from on/off simulation
- estimate PDF from events

## Fast MC simulation

- simpler generation model
- still estimate PDF from events

## Data-driven model

- extrapolate from control regions

## Effective model

- parametrized functional form

# Summary of models

## Pros

## Cons

MC Simulation

detailed descr  
of micro physics

very slow  
must reconstruct PDF

Fast MC

fast

approximate

Data-driven

Calculations by Nature

Extrapolations  
from CR have  
uncertainties

Effective model

fast,  
physical justification

approximate  
no details of  
underlying effects



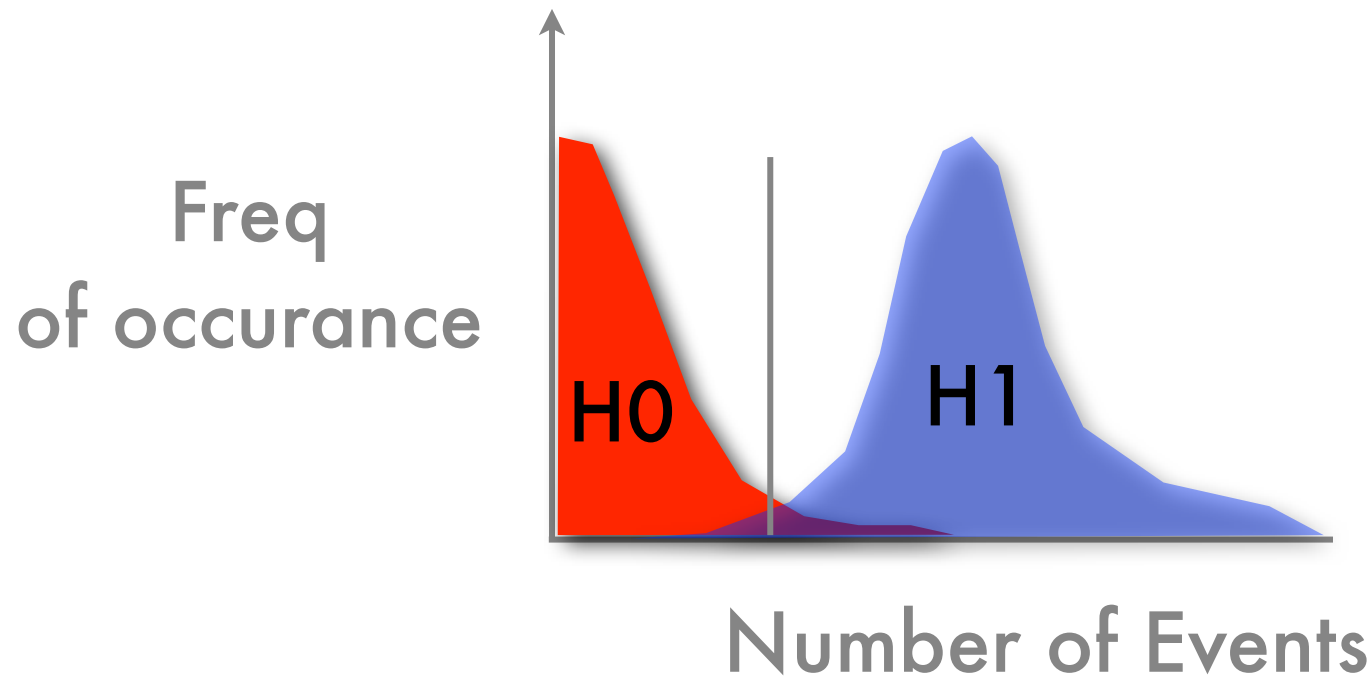
# Hypothesis Testing

# Hypothesis Testing

	BSM Particle is real	BSM Particle is not real
Claim Discovery	True Positive	False Positive Type I error $\alpha$
No Claim of Discovery	False Negative Type II error	True Negative

$\beta$ , power =  $1 - \beta$

# Example

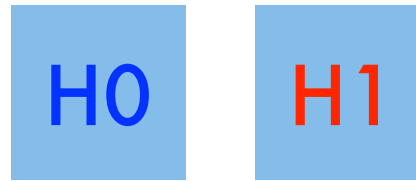


A threshold makes sense.  
Choice of position balances  
Type I/II errors

Typically:  
fix  $\alpha$   
minimize  $\beta$

# Generalize

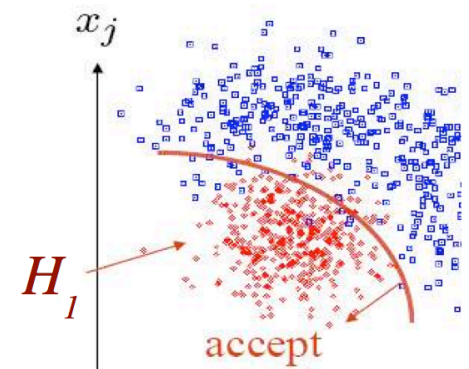
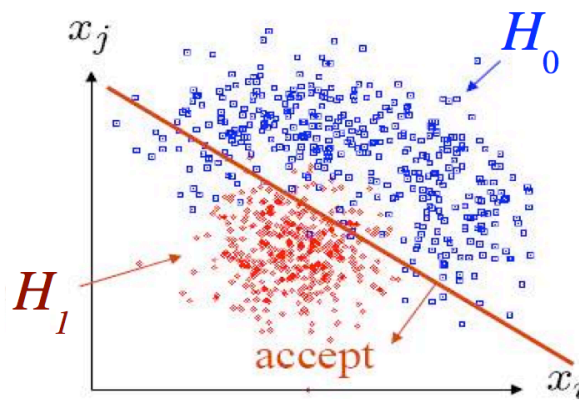
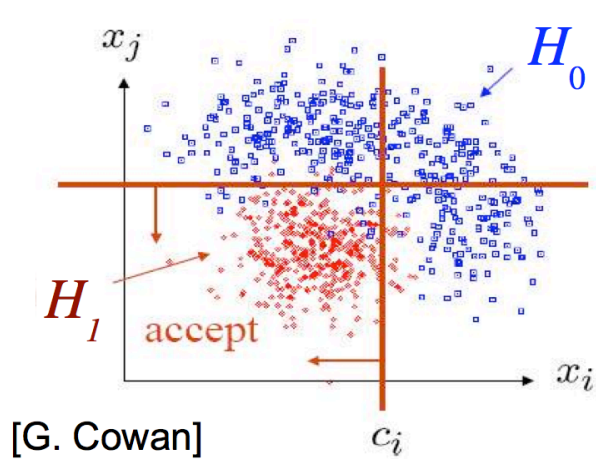
Hypothesis  
Testing



Parameter  
Estimation

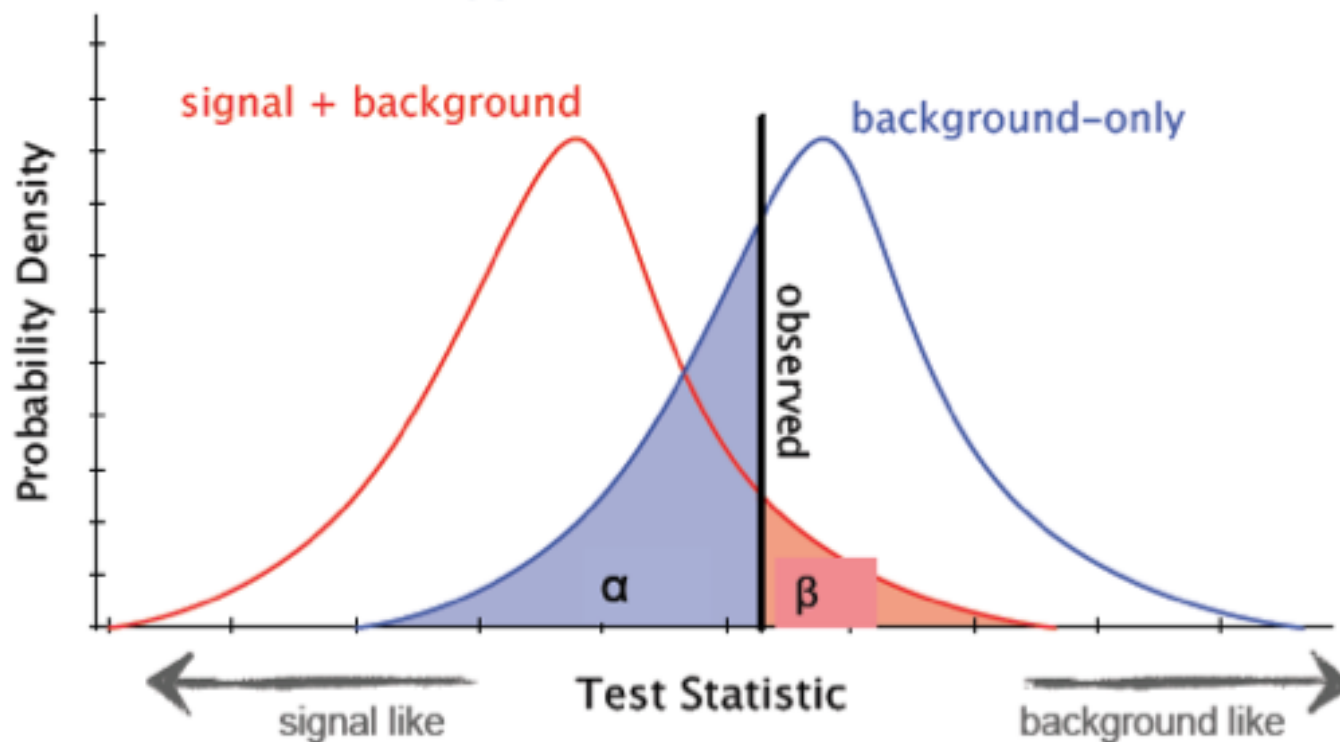


# More complicated



# Test statistic

Reduce vector of observables to 1 number



How to build distribution of TS? (Usually MC)

How to choose TS?

(K. Cranmer)

# Neyman-Pearson

## Statement of the problem:

Given some prob that we wrongly reject the Null hypothesis

$$\alpha = P(x \notin W | H_0)$$

Find the region  $W$  (where we accept  $H_0$ ) such that we minimize the prob

$$\beta = P(x \in W | H_1)$$

	BSM Particle is real	BSM Particle is not real
Claim Discovery	True Positive	False Positive $\alpha$ Type I error
No Claim of Discovery	False Negative Type II error	True Negative
$\beta$ , power=1- $\beta$		

# Neyman-Pearson

NP lemma says that the best test statistic is the likelihood ratio:

$$\frac{P(x|H_1)}{P(x|H_0)} > k_\alpha$$

(Gives smallest  $\beta$  for fixed  $\alpha$ )

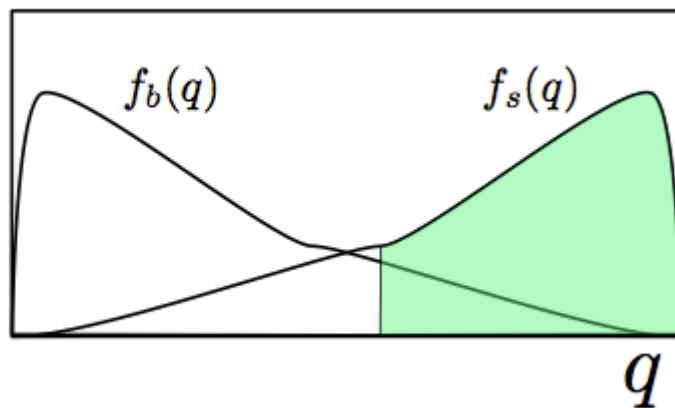
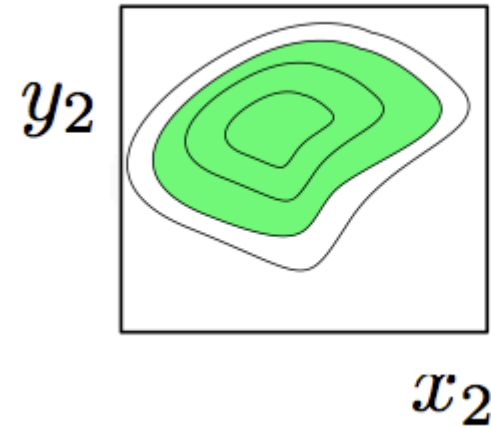
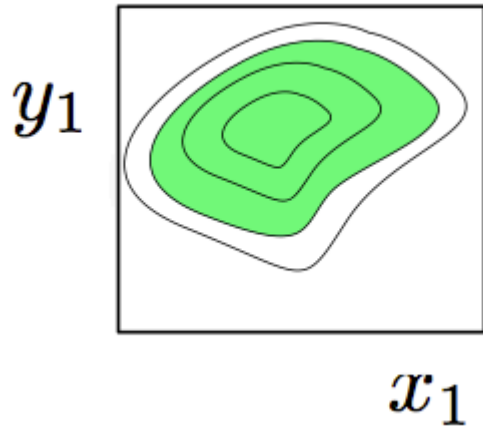
	BSM Particle is real	BSM Particle is not real
Claim Discovery	True Positive	False Positive Type I error $\alpha$
No Claim of Discovery	False Negative Type II error	True Negative

$\beta$ , power=1- $\beta$



# What does the TS do?

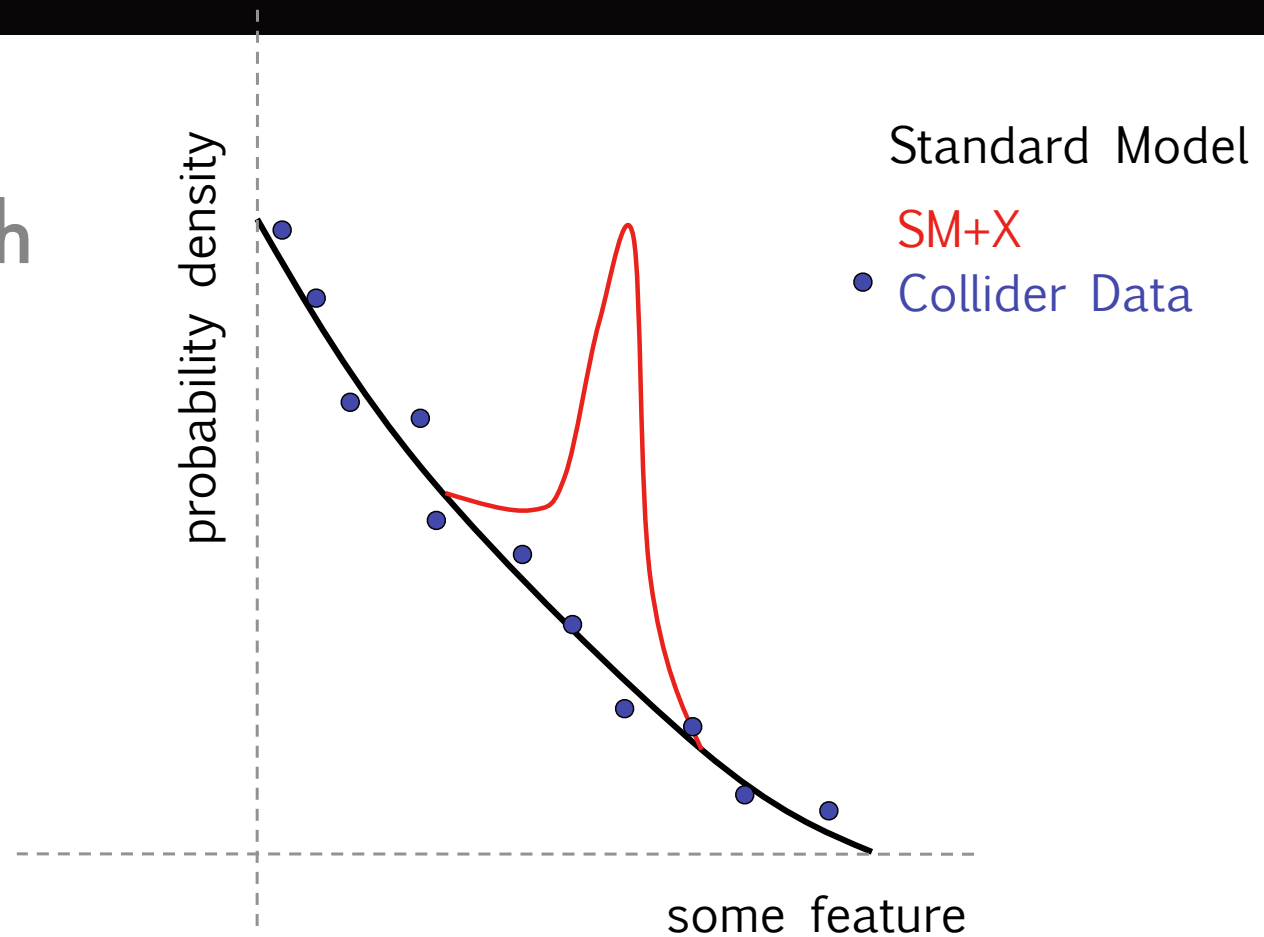
Finds a region in variable space



(K. Cranmer)

# How to find NP

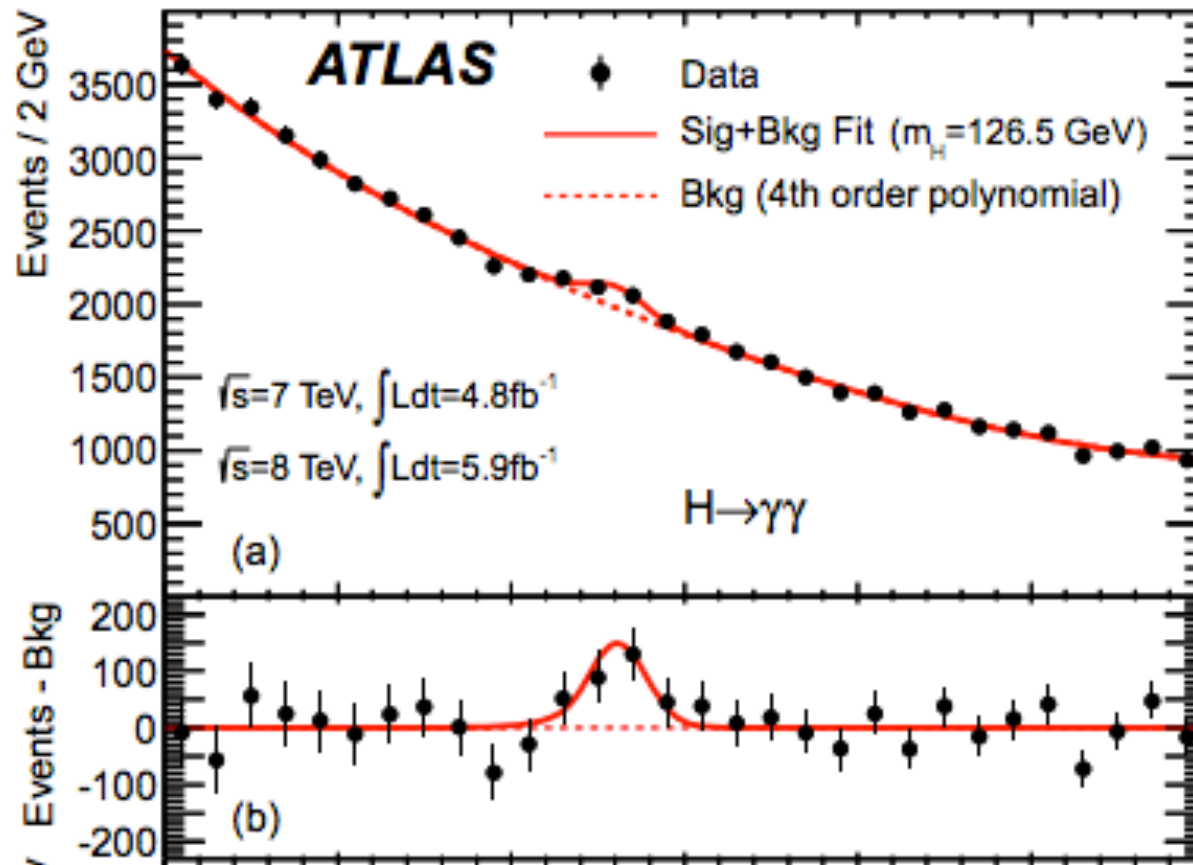
Isolate some feature in which two theories **SM**, **SM+X** can be best distinguished.



The data can tell us which hypothesis is preferred via a likelihood ratio:

$$\frac{L_{SM+X}}{L_{SM}} = \frac{P(\text{data} \mid \text{SM+X})}{P(\text{data} \mid \text{SM})}$$

e.g.

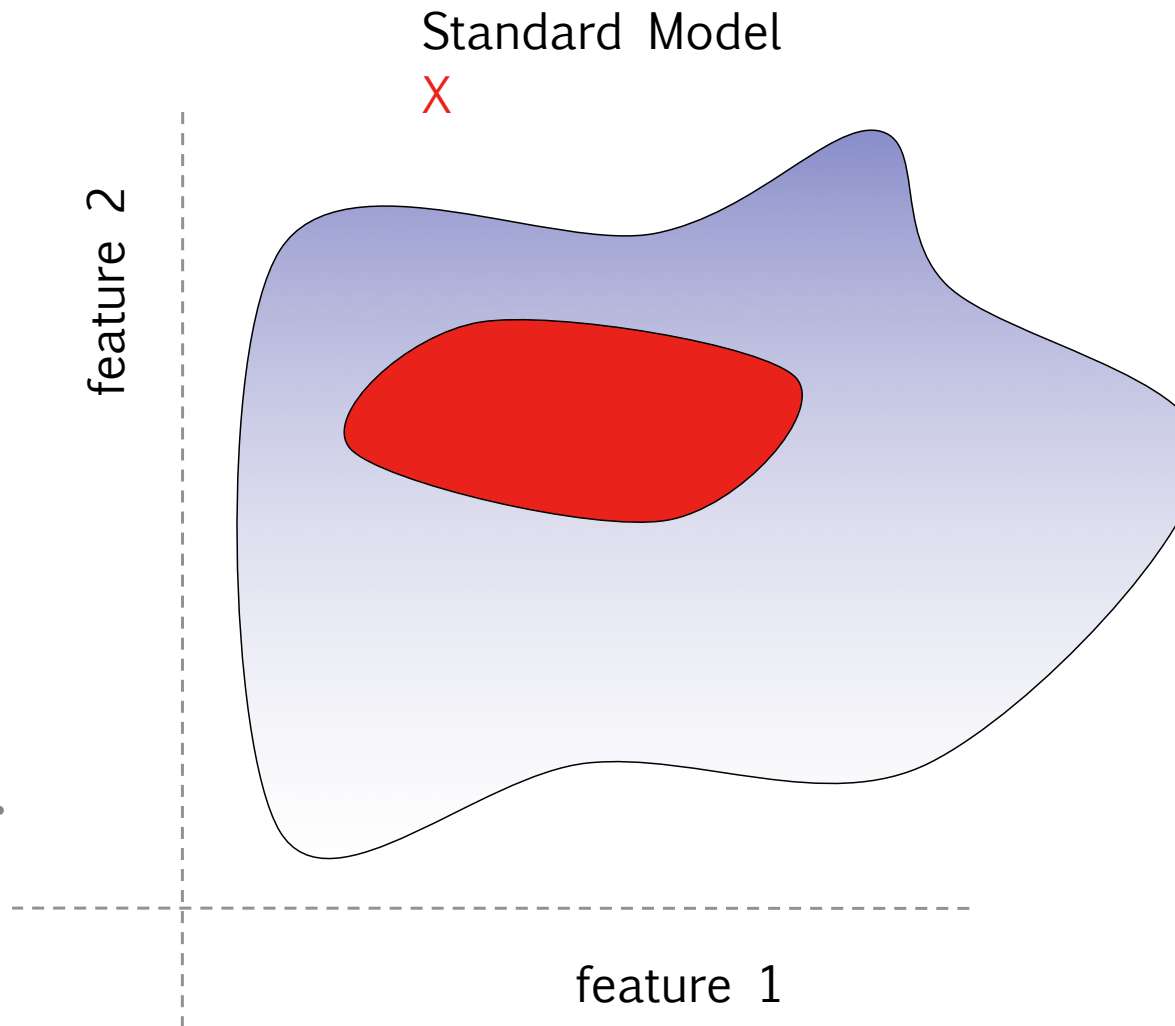


# But...

Reality is more complicated.

The full space can be very high dimensional.

Calculating likelihood in **d**-dimensional space requires  $\sim 100^d$  MC events.

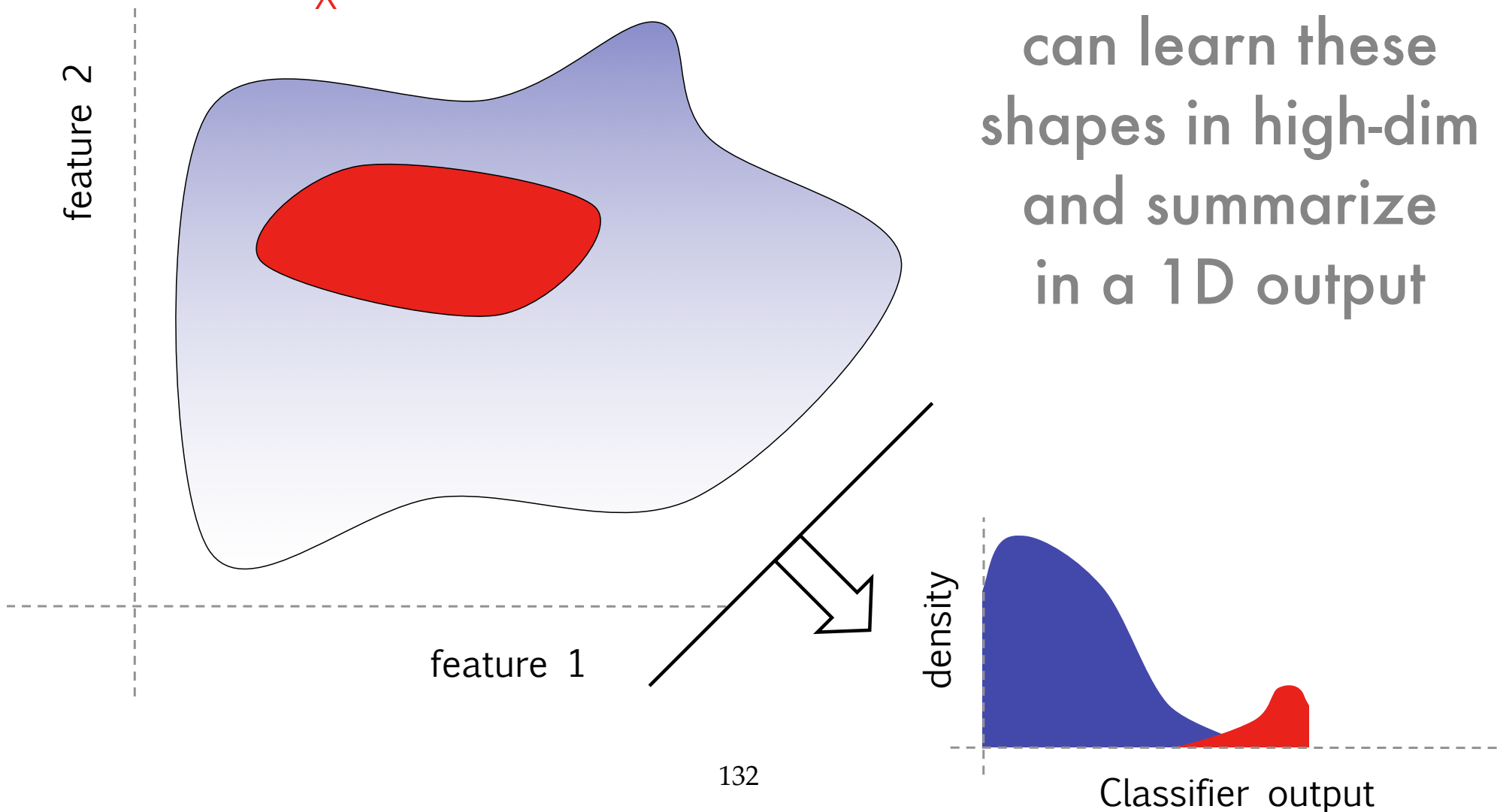


# ML tools

Standard Model

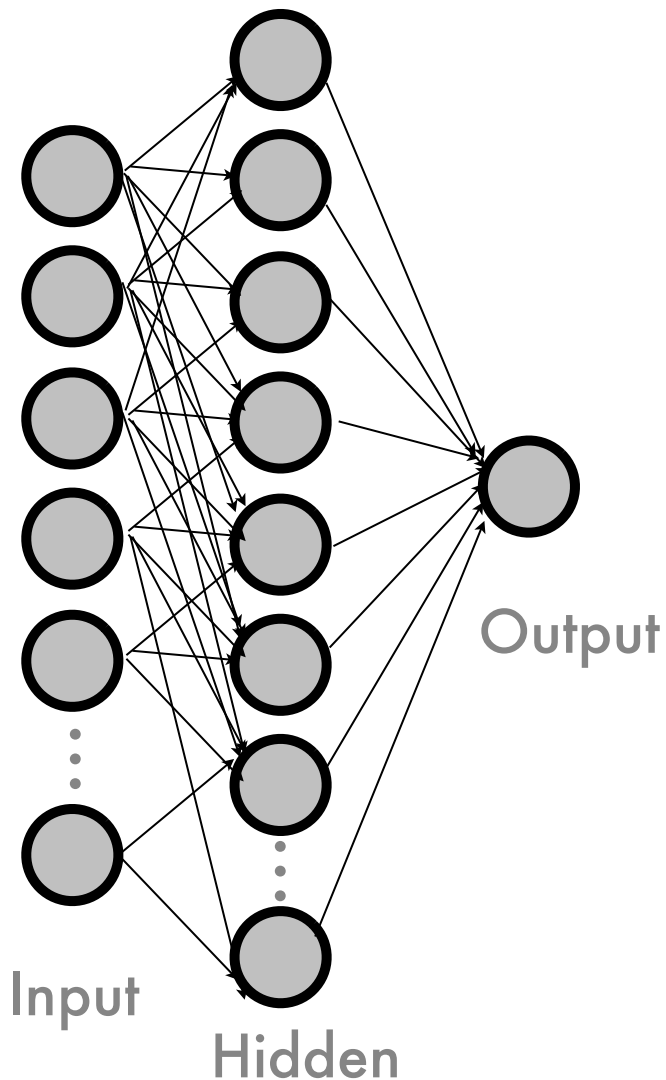
X

Neural networks  
can learn these  
shapes in high-dim  
and summarize  
in a 1D output



# Neural Networks

Essentially a functional fit with many parameters



## Function

Each neuron's output is a function of the weighted sum of inputs.

## Goal

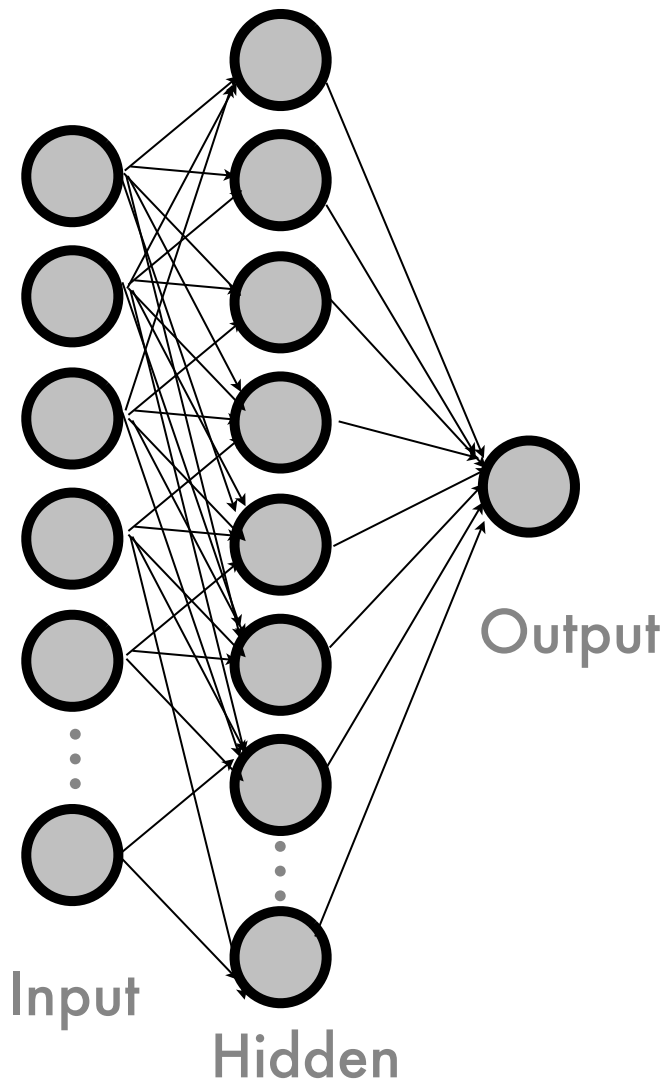
find set of weights which give most useful function

## Learning

give examples, back-propagate error to adjust weights

# Neural Networks

Essentially a functional fit with many parameters



## Problem:

Networks with  $> 1$  layer are very difficult to train.

## Consequence:

Networks are not good at learning non-linear functions.  
**(like invariant masses!)**

## In short:

Can't just throw 4-vectors at NN.

# Search for Input

ATLAS-CONF-2013-108

Can't just use  $4v$

Can't give it too many inputs

Painstaking search through input feature space.

Variable	VBF			Boosted		
	$\tau_{\text{lep}}\tau_{\text{lep}}$	$\tau_{\text{lep}}\tau_{\text{had}}$	$\tau_{\text{had}}\tau_{\text{had}}$	$\tau_{\text{lep}}\tau_{\text{lep}}$	$\tau_{\text{lep}}\tau_{\text{had}}$	$\tau_{\text{had}}\tau_{\text{had}}$
$m_{\tau\tau}^{\text{MMC}}$	•	•	•	•	•	•
$\Delta R(\tau, \tau)$	•	•	•		•	•
$\Delta\eta(j_1, j_2)$	•	•	•			
$m_{j_1, j_2}$	•	•	•			
$\eta_{j_1} \times \eta_{j_2}$		•	•			
$p_{\tau}^{\text{total}}$		•	•			
sum $p_{\tau}$					•	•
$p_{\tau}(\tau_1)/p_{\tau}(\tau_2)$					•	•
$E_{\tau}^{\text{miss}} \phi$ centrality		•	•	•	•	•
$x_{\tau 1}$ and $x_{\tau 2}$						•
$m_{\tau\tau, j_1}$				•		
$m_{\ell_1, \ell_2}$				•		
$\Delta\phi_{\ell_1, \ell_2}$				•		
sphericity				•		
$p_{\tau}^{\ell_1}$				•		
$p_{\tau}^{j_1}$				•		
$E_{\tau}^{\text{miss}}/p_{\tau}^{\ell_2}$				•		
$m_{\tau}$		•			•	
$\min(\Delta\eta_{\ell_1, \ell_2, \text{jets}})$	•					
$j_3$ $\eta$ centrality	•					
$\ell_1 \times \ell_2$ $\eta$ centrality	•					
$\ell$ $\eta$ centrality		•				
$\tau_{1,2}$ $\eta$ centrality			•			

Table 3: Discriminating variables used for each channel and category. The filled circles identify which variables are used in each decay mode. Note that variables such as  $\Delta R(\tau, \tau)$  are defined either between the two leptons, between the lepton and  $\tau_{\text{had}}$ , or between the two  $\tau_{\text{had}}$  candidates, depending on the decay mode.



# Search for Input

ATLAS-CONF-2013-108

Can't just use 4v

Can't give it +  
many inp

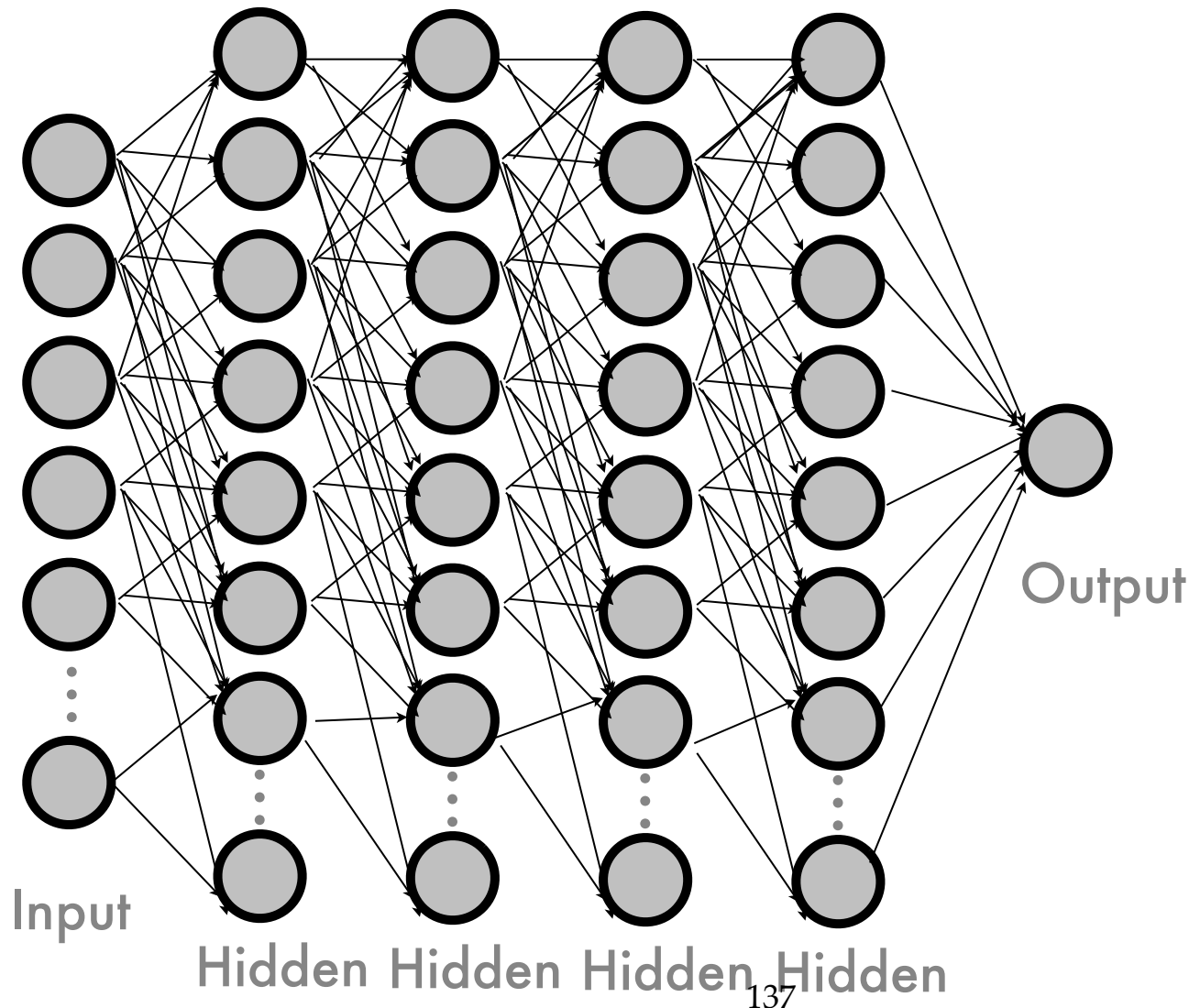
Painstaking  
through inp  
feature space.

**Also true for  
BDTs, SVNs, etc**

Variable	VBF	Boosted		
	$\tau_{lep}\tau_{lep}$	$\tau_{lep}\tau_{had}$	$\tau_{had}\tau_{had}$	
$m_{TT}^{MMC}$		•	•	•
$\Delta P_T$		•	•	
$P_T^{j1}$		•		
$E_T^{miss} / P_T^{\ell 2}$		•		
$m_T$		•		•
$\min(\Delta\eta_{\ell_1, \ell_2, jets})$	•			
$j_3$ centrality	•			
$\ell_1 \times \ell_2$ centrality	•			
$\ell$ centrality		•		
$\tau_{1,2}$ centrality			•	

Table 3: Discriminating variables used for each channel and category. The filled circles identify which variables are used in each decay mode. Note that variables such as  $\Delta R(\tau, \tau)$  are defined either between the two leptons, between the lepton and  $\tau_{had}$ , or between the two  $\tau_{had}$  candidates, depending on the decay mode.

# Deep networks



New tools  
let us  
train  
deep  
networks.

How well  
do they work?

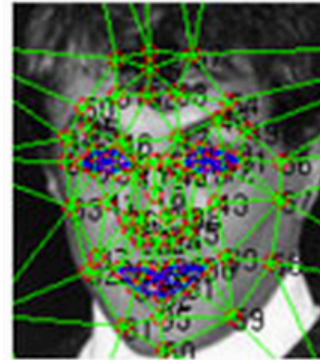
# Real world applications



(a)



(b)



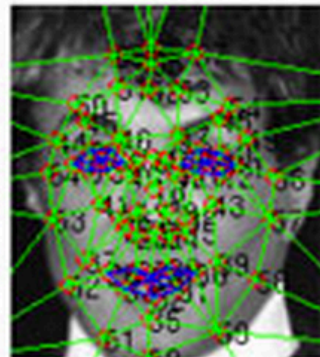
(c)



(d)



(e)



(f)



(g)



(h)

**Head turn:** DeepFace uses a 3-D model to rotate faces, virtually, so that they face the camera. Image (a) shows the original image, and (g) shows the final, corrected version.

# Paper

## Deep Learning in High-Energy Physics: Improving the Search for Exotic Particles

P. Baldi,<sup>1</sup> P. Sadowski,<sup>1</sup> and D. Whiteson<sup>2</sup>

<sup>1</sup>*Dept. of Computer Science, UC Irvine, Irvine, CA 92617*

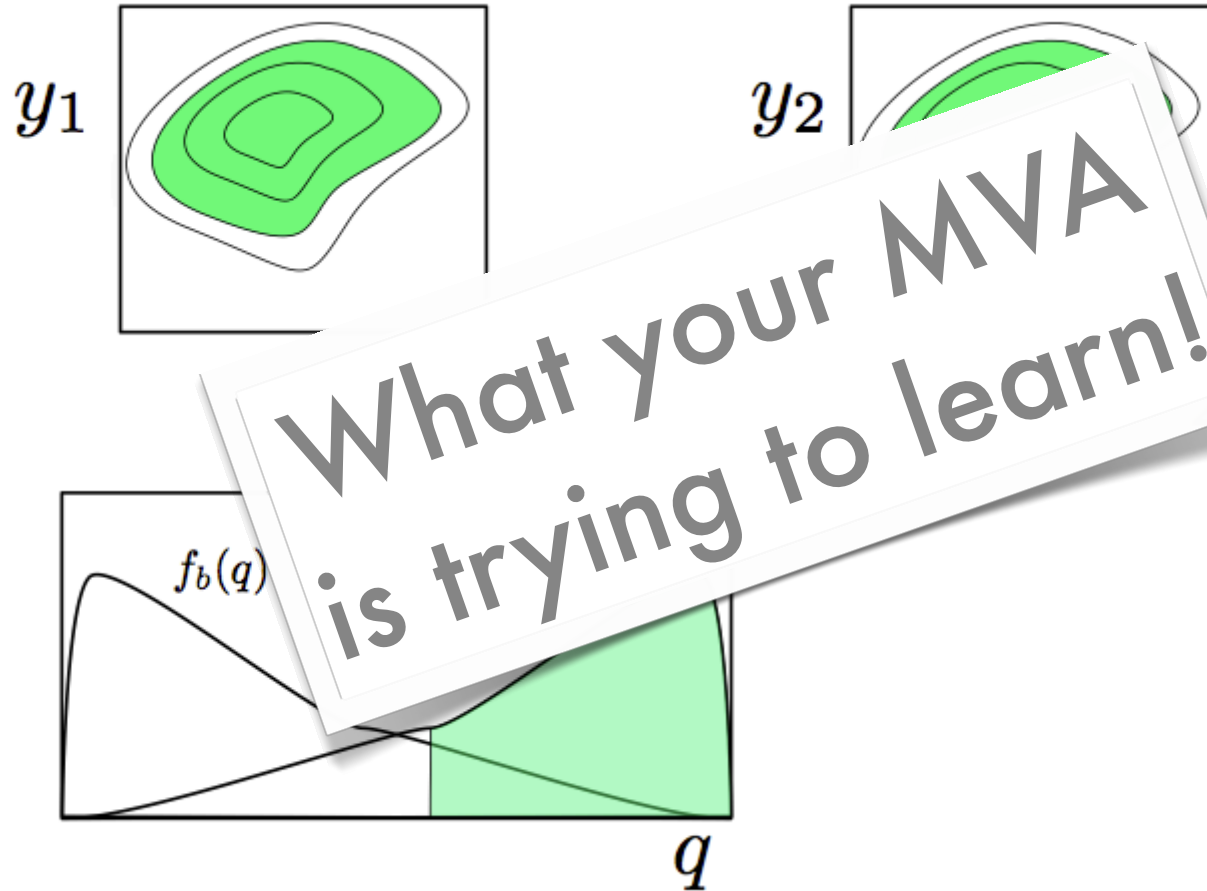
<sup>2</sup>*Dept. of Physics and Astronomy, UC Irvine, Irvine, CA 92617*

arXiv: 1402.4735

Accepted in *Nature Comm.*

# What does the TS do?

Finds a region in variable space



(K. Cranmer)

# Test statistic

Define  $\mu$  to be signal strength,  
 $\mu=0$  is no signal  
 $\mu=1$  is theory prediction

At LEP, this was used:

$$Q_{LEP} = L_{s+b}(\mu = 1) / L_b(\mu = 0)$$

Where the nuisance parameters  
are fixed to their nominal values

# Test statistic

Define  $\mu$  to be signal strength,  
 $\mu=0$  is no signal  
 $\mu=1$  is theory prediction

At LEP, this was used:

$$Q_{LEP} = \frac{L(data|\mu = 1, b, \nu)}{L(data|\mu = 0, b, \nu)}$$

This also means the background estimate doesn't vary.

# Tevatron

Still consider two points (0,1)  
but now float the NPs at those points

$$Q_{TEV} = L_{s+b}(\mu = 1, \hat{\nu}) / L_b(\mu = 0, \hat{\nu}')$$

Ratio of profiled likelihoods:  
the model is adapted to the data  
even in the signal region





# LHC

## Profile likelihood

$$\lambda(\mu = 0) = \frac{L(\text{data} | \mu = 0, \hat{b}(\mu = 0), \hat{v}(\mu = 0))}{L(\text{data} | \hat{\mu}, \hat{b}, \hat{v})}$$

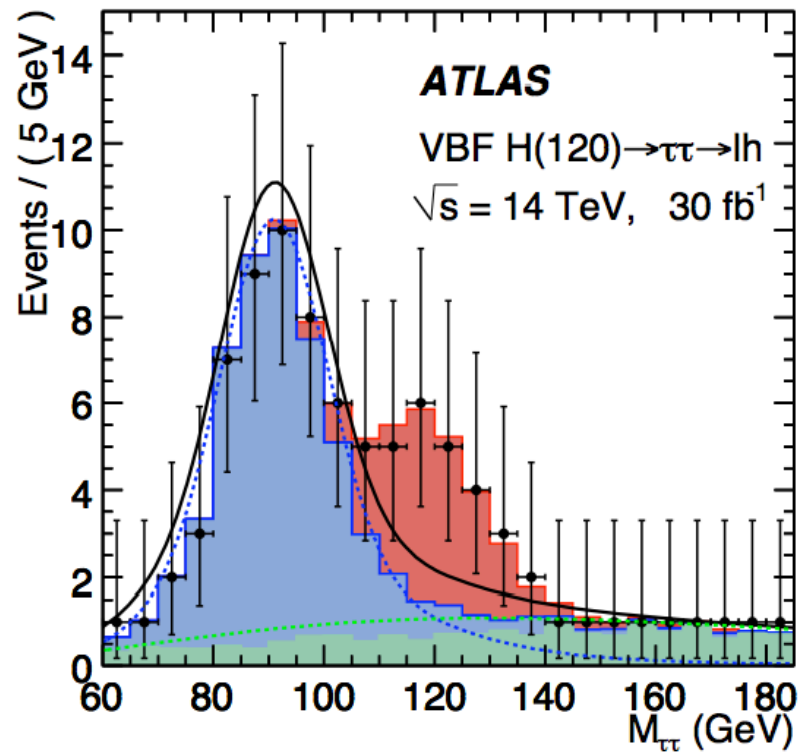
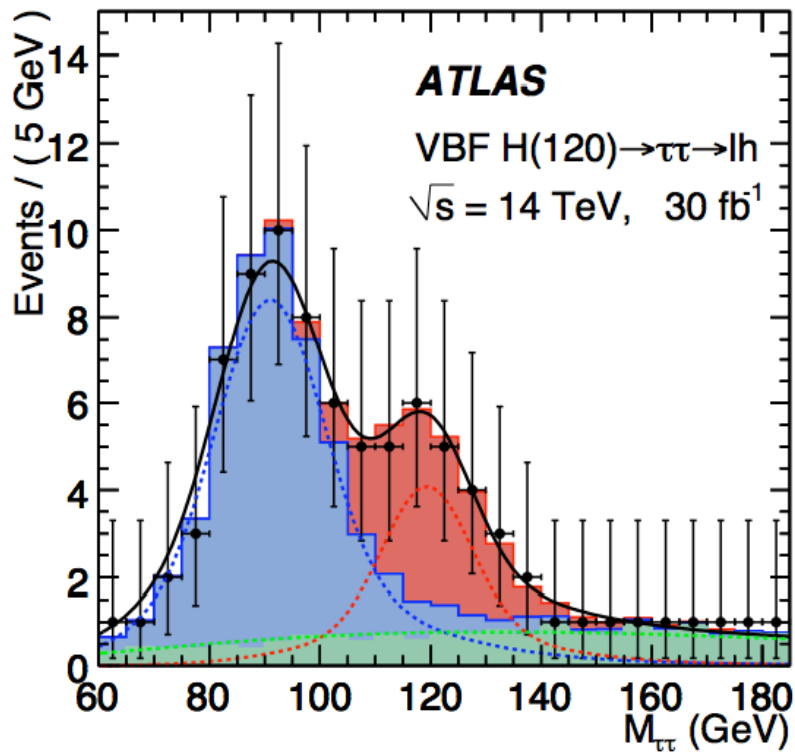
fit best value of NPs at  $\mu=0$

and at **best fit value of  $\mu$**

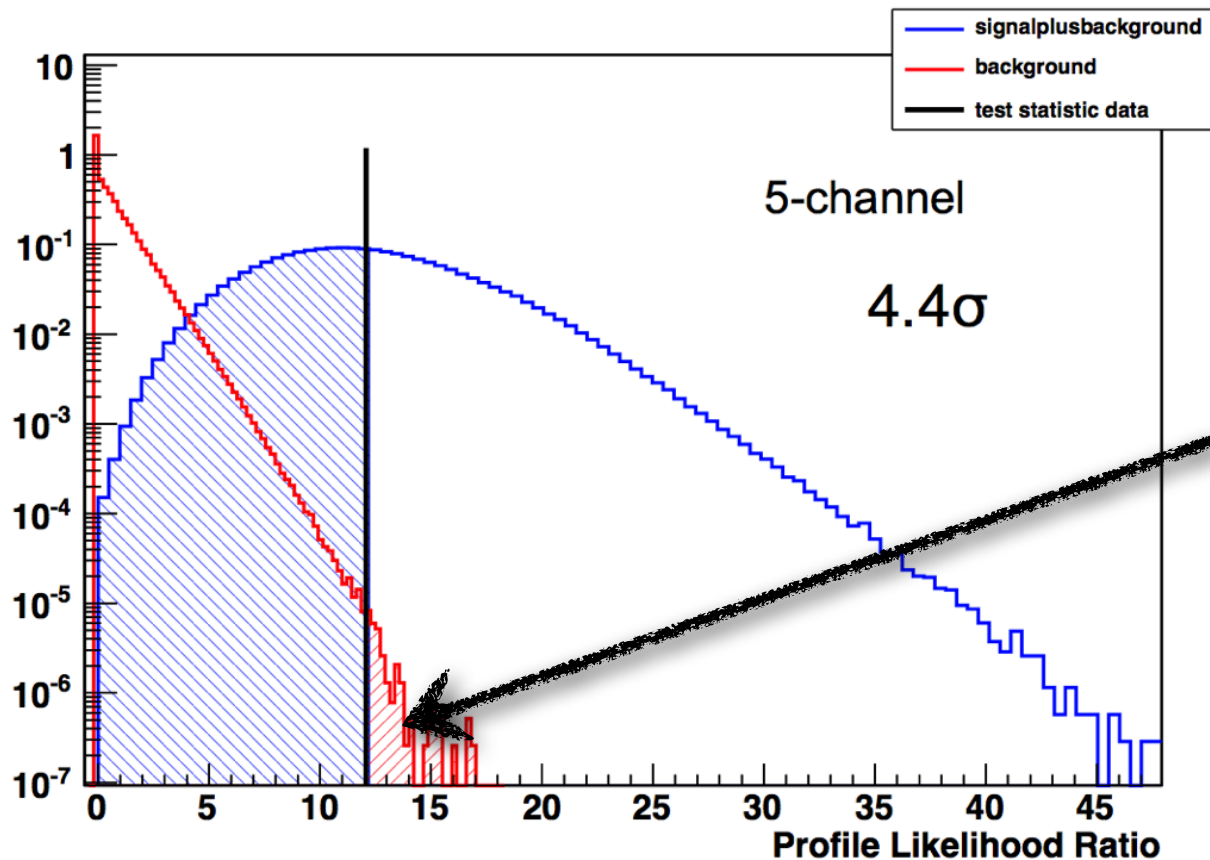
# Two fits to data

$$\lambda(\mu = 0) = \frac{L(\text{data}|\mu = 0, \hat{b}(\mu = 0), \hat{v}(\mu = 0))}{L(\text{data}|\hat{\mu}, \hat{b}, \hat{v})},$$

$$L(\text{data}|\hat{\mu}, \hat{b}, \hat{v}) \qquad L(\text{data}|\mu = 0, \hat{b}, \hat{v})$$



# p values

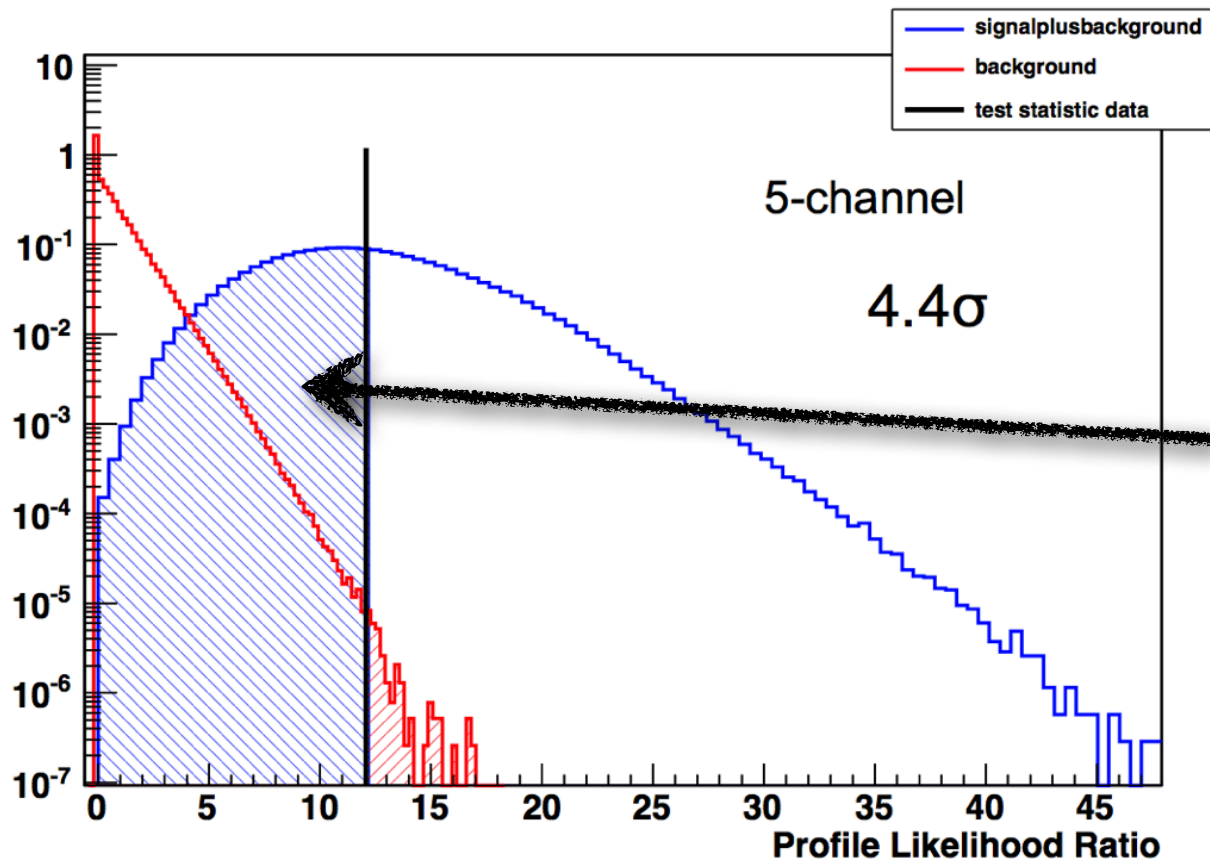


$p_0 =$   
probability  
to observe data  
or **more** signal-like  
under **background**  
hypothesis

$$p_0 = P(q_0 \geq q_o^{obs})$$

(K. Cranmer)

# p values



$p_\mu =$   
probability  
to observe data  
or **less** signal-like  
under **signal+b**  
hypothesis

# Philosophy

Bayesian  
&  
Frequentist

# Bayesian

Data: fixed

Parameter values: unknown

Probability: our lack of knowledge

PDFs over parameters: sensible

# Frequentist

Data: one example from ens.

Parameter values: **fixed (even if unknown)**

Probability: **rate of occurrence**

PDFs over parameters: **not sensible**

# Bayesian Prob.

Bayes theorem:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

rearrange:

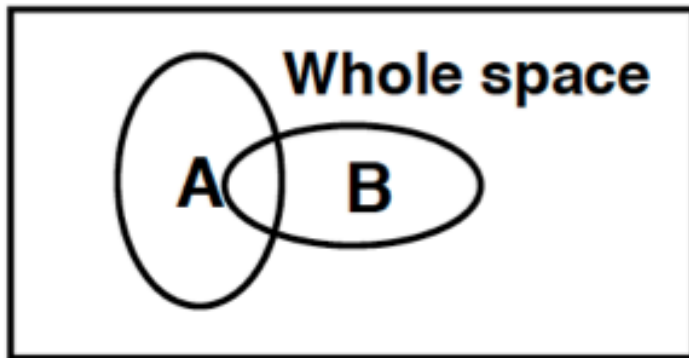
$$P(A|B)P(B) = P(A \cap B) = P(B|A)P(A)$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$



# In Pictures

## P, Conditional P, and Derivation of Bayes' Theorem in Pictures



$$P(A) = \frac{\text{Area of } A}{\text{Area of Whole space}}$$

$$P(B) = \frac{\text{Area of } B}{\text{Area of Whole space}}$$

$$P(A|B) = \frac{\text{Area of } A \cap B}{\text{Area of } B}$$

$$P(B|A) = \frac{\text{Area of } A \cap B}{\text{Area of } A}$$

$$P(A \cap B) = \frac{\text{Area of } A \cap B}{\text{Area of Whole space}}$$

$$P(A) \times P(B|A) = \frac{\text{Area of } A}{\text{Area of Whole space}} \times \frac{\text{Area of } A \cap B}{\text{Area of } A} = \frac{\text{Area of } A \cap B}{\text{Area of Whole space}} = P(A \cap B)$$

$$P(B) \times P(A|B) = \frac{\text{Area of } B}{\text{Area of Whole space}} \times \frac{\text{Area of } A \cap B}{\text{Area of } B} = \frac{\text{Area of } A \cap B}{\text{Area of Whole space}} = P(A \cap B)$$

$$\Rightarrow P(B|A) = P(A|B) \times P(B) / P(A)$$

# Example 1

$P(\text{data} \mid \text{theory}) \neq P(\text{theory} \mid \text{data})$

Theory = (male or female)

Data = (pregnant | not pregnant)

$P(\text{pregnant} \mid \text{female}) \sim 3\%$

BUT

$P(\text{female} \mid \text{pregnant}) > 99\%$

# Example 2

## Higgs search

Expected **bg** = 0.1

Expected **signal** = 10

$P(N \mid \text{no Higgs}) = 0.1$

$P(N \mid \text{Higgs}) = 10.1$

What is  $P(\text{Higgs} \mid N=8)$ ?  $P(H \mid N = 8) = \frac{P(N = 8 \mid H)P(H)}{P(N = 8)}$

Depends on  $P(H)$ !

(K Cranmer)

# Parameter estimation

## Bayesian parameter estimation:

Want to know the probability that some parameter  $\theta$  is in some range  $[\theta_0, \theta_1]$

- or -

Want to find a range  $[\theta_0, \theta_1]$  that has probability of 0.95

# Parameter estimation

## Bayesian parameter estimation

Want to know the probability that  
parameter  $\theta$  is in some range

- or -

Want to  
know the  
probability  
of

**Remember:**  
Probability reflects lack of knowledge, and  
prior information.

$[\theta_0, \theta_1]$  that has probability

# How?

The probability that the true value is inside an interval is:

$$1 - \alpha = \int_{\theta_{lo}}^{\theta_{hi}} p(\theta|x) d\theta$$

For lower or upper limits, choose zero or infinity as boundaries.  
where we integrate out the nuisance parameters:

$$p(\theta|x) = \int d\nu p(\theta, \nu|x)$$

where

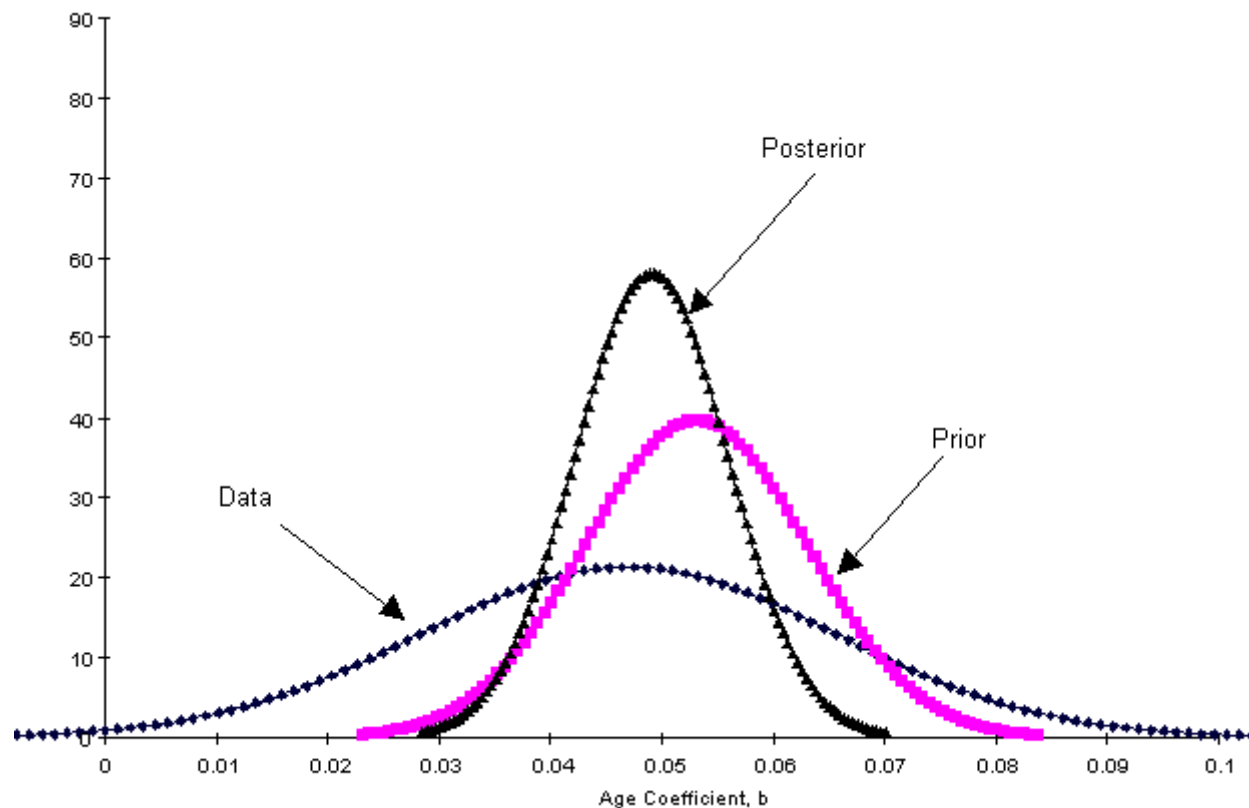
$$p(\theta, \nu|x) = \frac{p(x|\theta, \nu)p(\theta, \nu)}{p(x)}$$

These integrals can be very hard to do if the space is high dimensional.

# Priors

## Choice of prior $p(\theta)$

- important but subjective choice



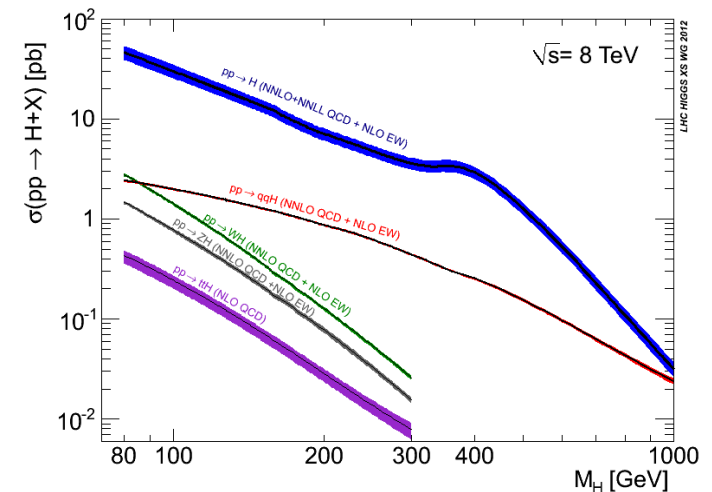
# Priors

## Choice of prior $p(\theta)$

- Example: measuring Higgs cross-section
- Want to be unbiased: choose uniform prior?

$$\sigma = [0, \Lambda] \rightarrow P = k$$

- But  $\sigma$  and mass relationship makes this prior **not flat in mass**



- no uninformative prior across all transformations