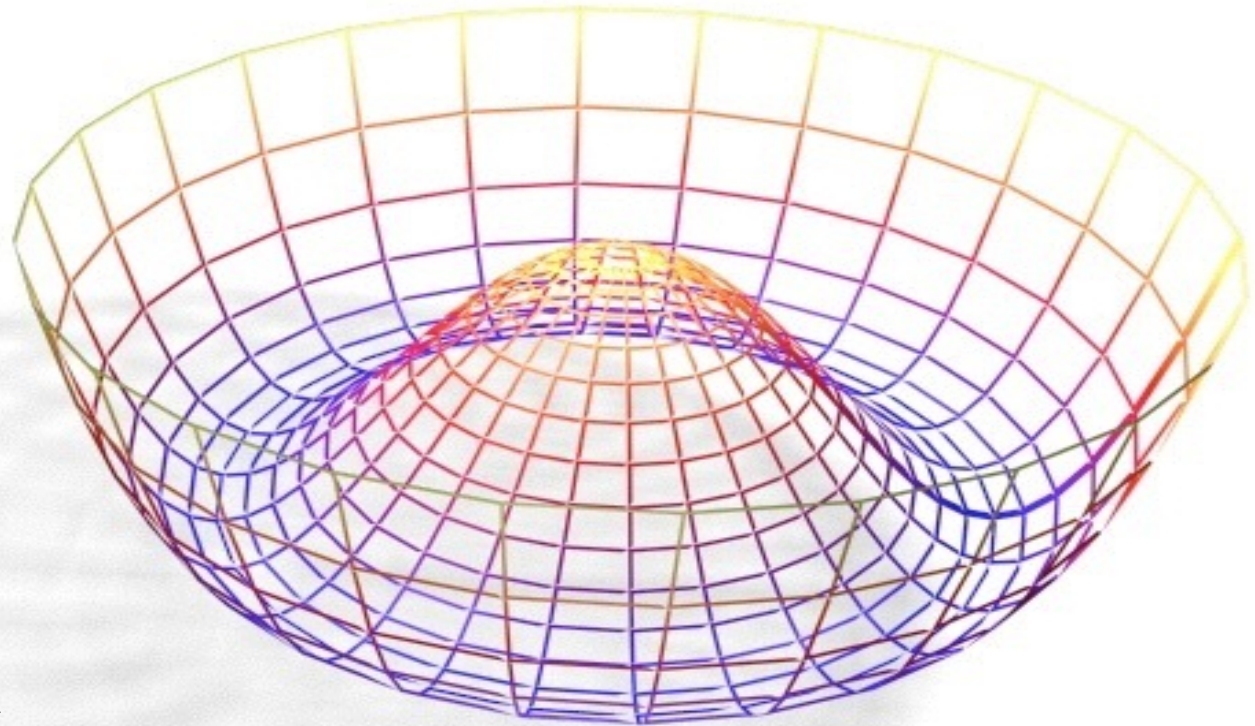




# *The Data Science Challenges of Particle Physics*



***Kyle Cranmer,***  
New York University



# BIO

- Experimental Particle [Physicist](#)
- Statistics Convener of [ATLAS](#) experiment at LHC
- Founder of [RooStats](#) framework (used for Higgs discovery)
- Co-lead [Open Science Working Group](#) for Moore-Sloan Data Science Environment at NYU

# A harbinger for things to come



Large, Distributed Collaborations  
**Big Science**



Complicated Sensor Environment  
**Big Data**  
**Big Simulation**

$$\begin{aligned}\mathcal{L}_{SM} = & \underbrace{\frac{1}{4}\mathbf{W}_{\mu\nu} \cdot \mathbf{W}^{\mu\nu} - \frac{1}{4}B_{\mu\nu}B^{\mu\nu} - \frac{1}{4}G_{\mu\nu}^a G^{\mu\nu a}}_{\text{kinetic energies and self-interactions of the gauge bosons}} \\ & + \underbrace{\bar{L}\gamma^\mu(i\partial_\mu - \frac{1}{2}g\boldsymbol{\tau} \cdot \mathbf{W}_\mu - \frac{1}{2}g'Y B_\mu)L + \bar{R}\gamma^\mu(i\partial_\mu - \frac{1}{2}g'Y B_\mu)R}_{\text{kinetic energies and electroweak interactions of fermions}} \\ & + \underbrace{\frac{1}{2}[(i\partial_\mu - \frac{1}{2}g\boldsymbol{\tau} \cdot \mathbf{W}_\mu - \frac{1}{2}g'Y B_\mu)\phi]^2 - V(\phi)}_{W^\pm, Z, \gamma, \text{ and Higgs masses and couplings}} \\ & + \underbrace{g''(\bar{q}\gamma^\mu T_a q)G_\mu^a}_{\text{gravitational interactions}} + \underbrace{(G_1\bar{L}\phi R + G_2\bar{L}\phi_c R + h.c.)}_{\text{Yukawa interactions}}\end{aligned}$$

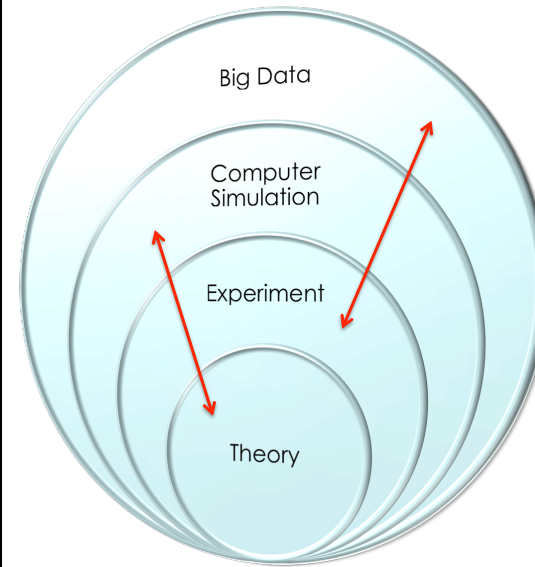
Scientifically Motivated Data Modeling  
**Big Simulation**  
**Big Model**

# Complex Models for Big Data



Max Welling  
UvA

# The Four Paradigms



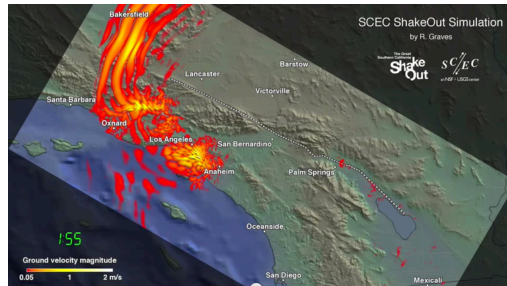
We have added big data to computer simulation, experiment and theory.



Not replaced it...

# Big Simulation

Computer simulations have become increasingly complex (e.g. weather, earthquake models)

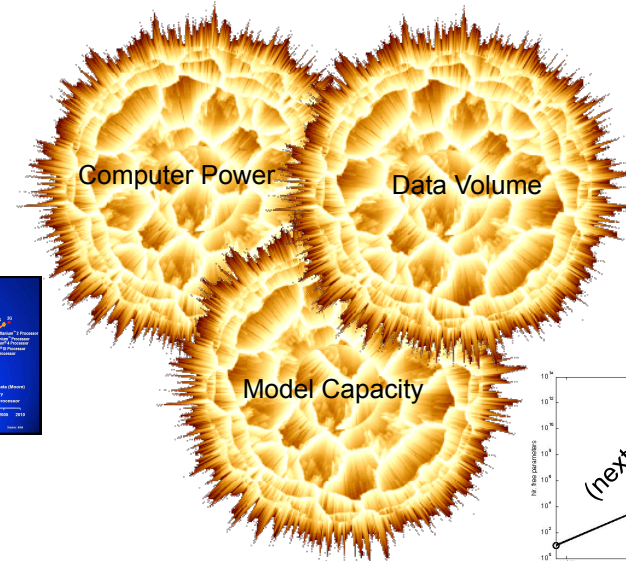
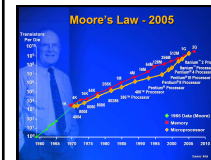


This production run producing 360 sec of wave propagation sustained 220 Tflop/s for 24 hours on NCCS Jaguar using 223,074 cores.

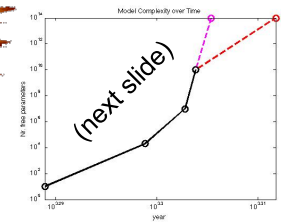
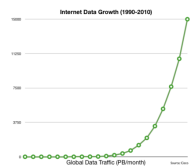
The Computational Wall: If a model has hundreds of parameters, how can we:

- 1) Find the parameter values that match the observations best?
- 2) Determine if we underfit (model too simple) or overfit (model too complex)?
- 3) Compare two models?

# 3x Exponential Growth in Machine Learning



Data is Growing Exponentially



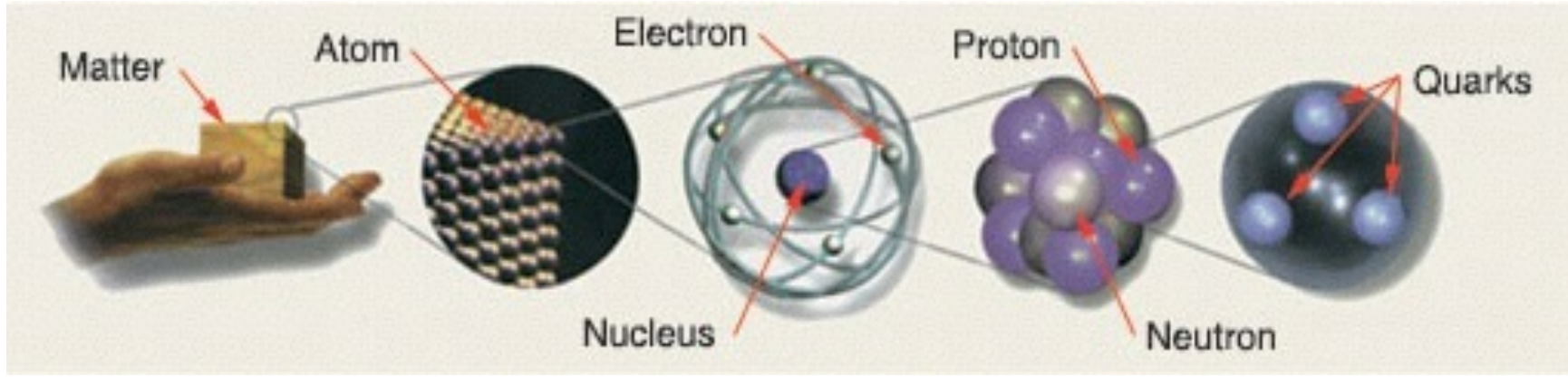
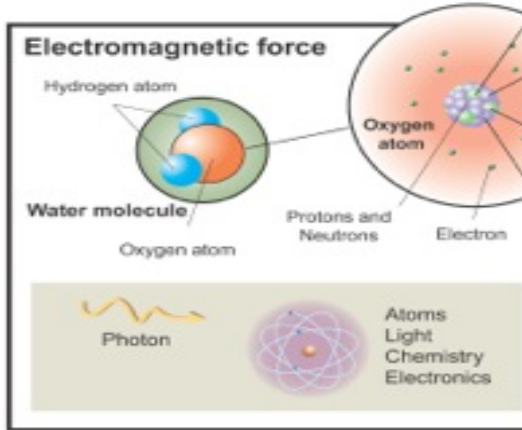
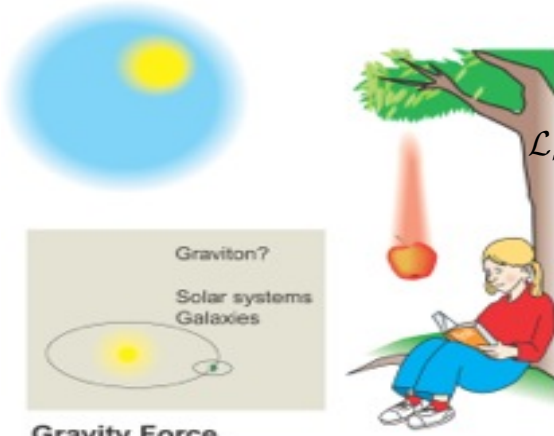


Illustration: Typoform



$$\mathcal{L}_{SM} =$$

Glueons (8)

 $\frac{1}{4} \mathbf{W}_{\mu\nu} \cdot \mathbf{W}^{\mu\nu} - \frac{1}{4} B_{\mu\nu} B^{\mu\nu}$ 

Kinetic energies and self-interactions of the gauge bosons

Mesons

 $\bar{L} \gamma^\mu (i \partial_\mu - \frac{1}{2} g_T \mathbf{W}_\mu - \frac{1}{2} g' Y B_\mu) L + \bar{R} \gamma^\mu (i \partial_\mu - \frac{1}{2} g' Y B_\mu) R$ 

Kinetic energies and electroweak interactions of fermions

$$+$$

Stron Force

 $\frac{1}{2} \bar{q} (i \partial_\mu - \frac{1}{2} g_T \mathbf{W}_\mu - \frac{1}{2} g' Y B_\mu) q + \bar{q} \gamma^\mu (i \partial_\mu - g_s \mathbf{T} \cdot \mathbf{G}_\mu) q$ 

kinetic energies and electroweak interactions of fermions

$$+$$

Bosons (W,Z)

 $g'' (\bar{q} \gamma^\mu T_a q) G_\mu^a + (G_1 \bar{L} \phi R + G_2 \bar{L} \phi R + h.c.)$ 

interactions between quarks and gluons

$$+$$

Neutron decay  
Beta decay  
Neutrino interactions  
Burning of the sun

W force carrier particle

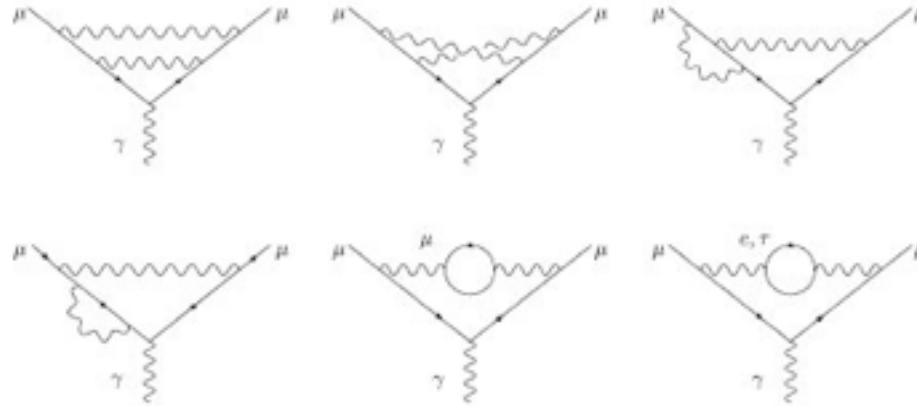
Leptons Quarks

<i>u</i> up	<i>c</i> charm	<i>t</i> top
<i>d</i> down	<i>s</i> strange	<i>b</i> bottom
<i>ν<sub>e</sub></i> e- neutrino	<i>ν<sub>μ</sub></i> μ- neutrino	<i>ν<sub>τ</sub></i> τ- neutrino
<i>e</i> electron	<i>μ</i> muon	<i>τ</i> tau

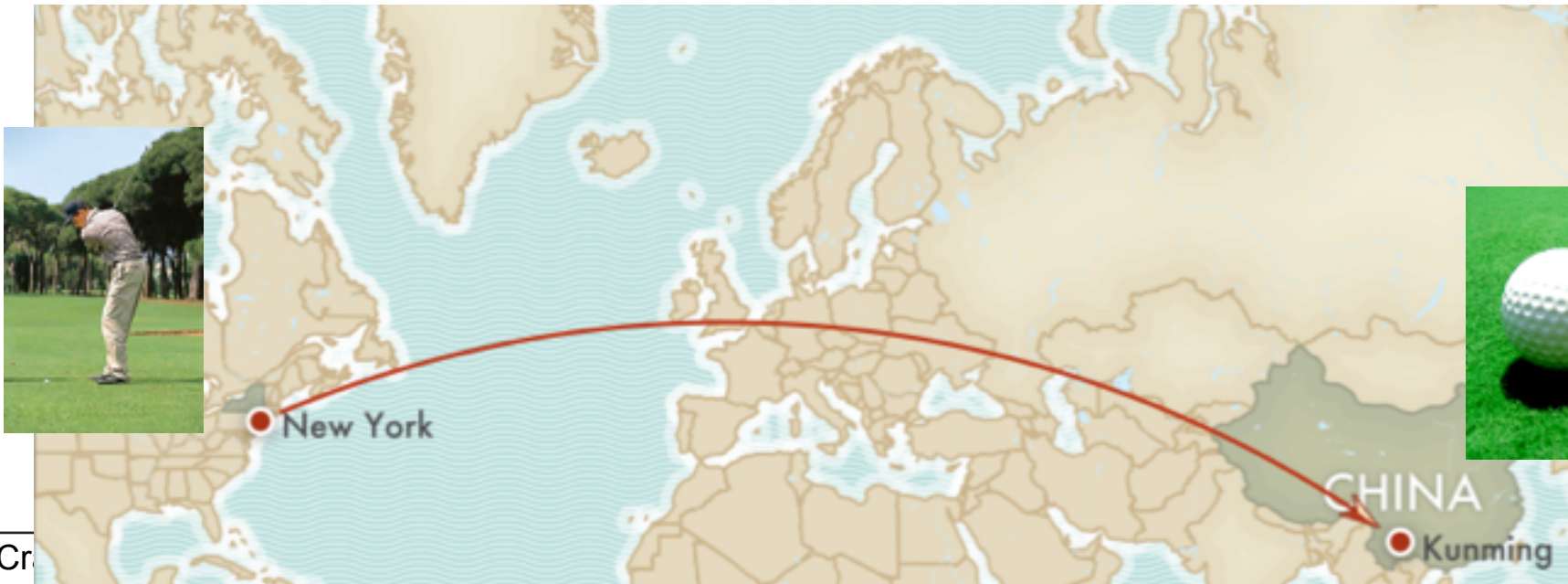
Three Generations of Matter

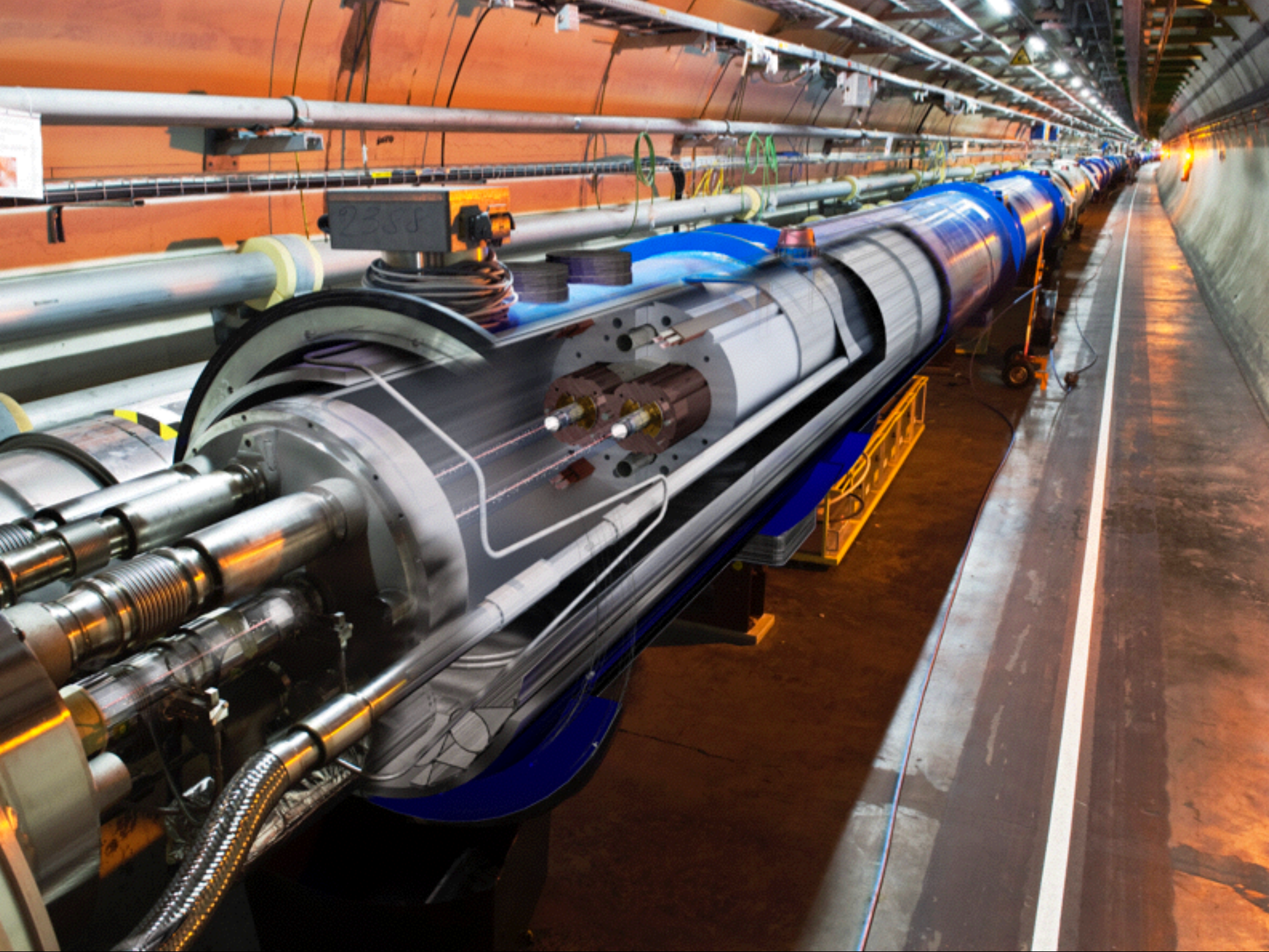
Non-trivial aspects of the theory have been tested to  $< 1$  ppm

A unique realm for reasonable statistical exploration of a scientific theory



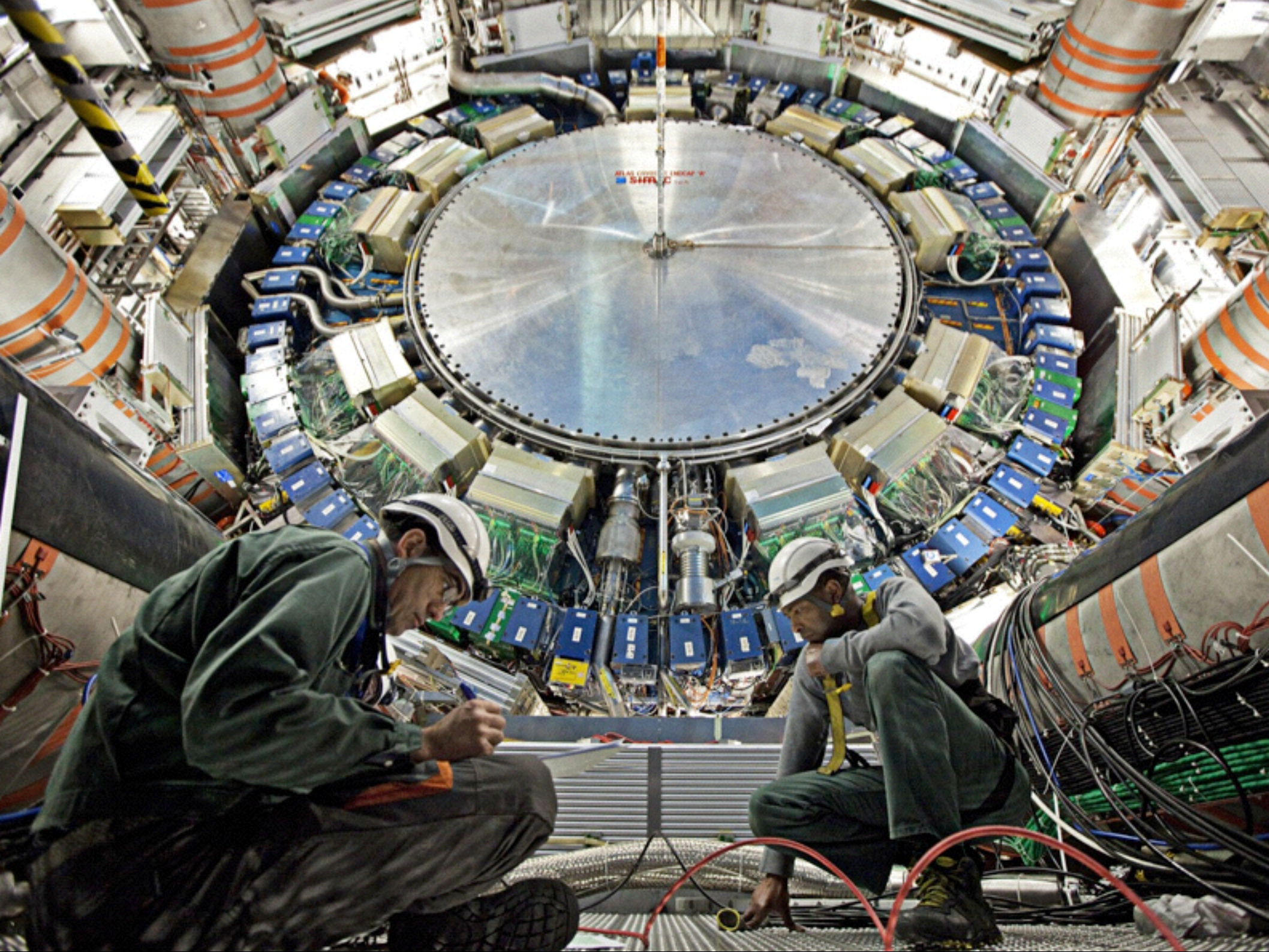
$$a_{\mu}(\text{exp}) = 11\,659\,208(6) \times 10^{-10} \text{ (0.5 ppm)}$$









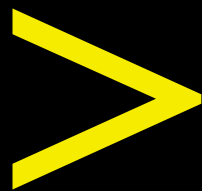


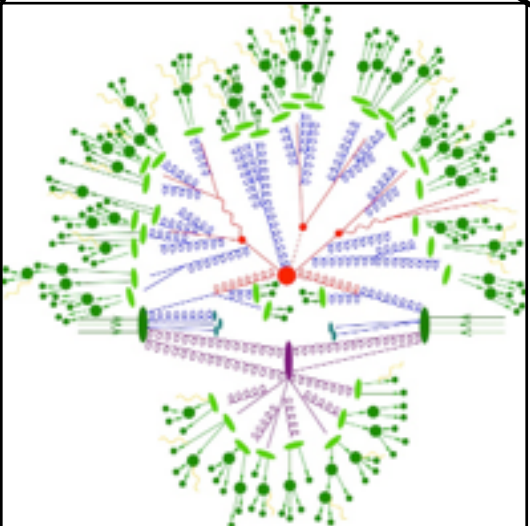
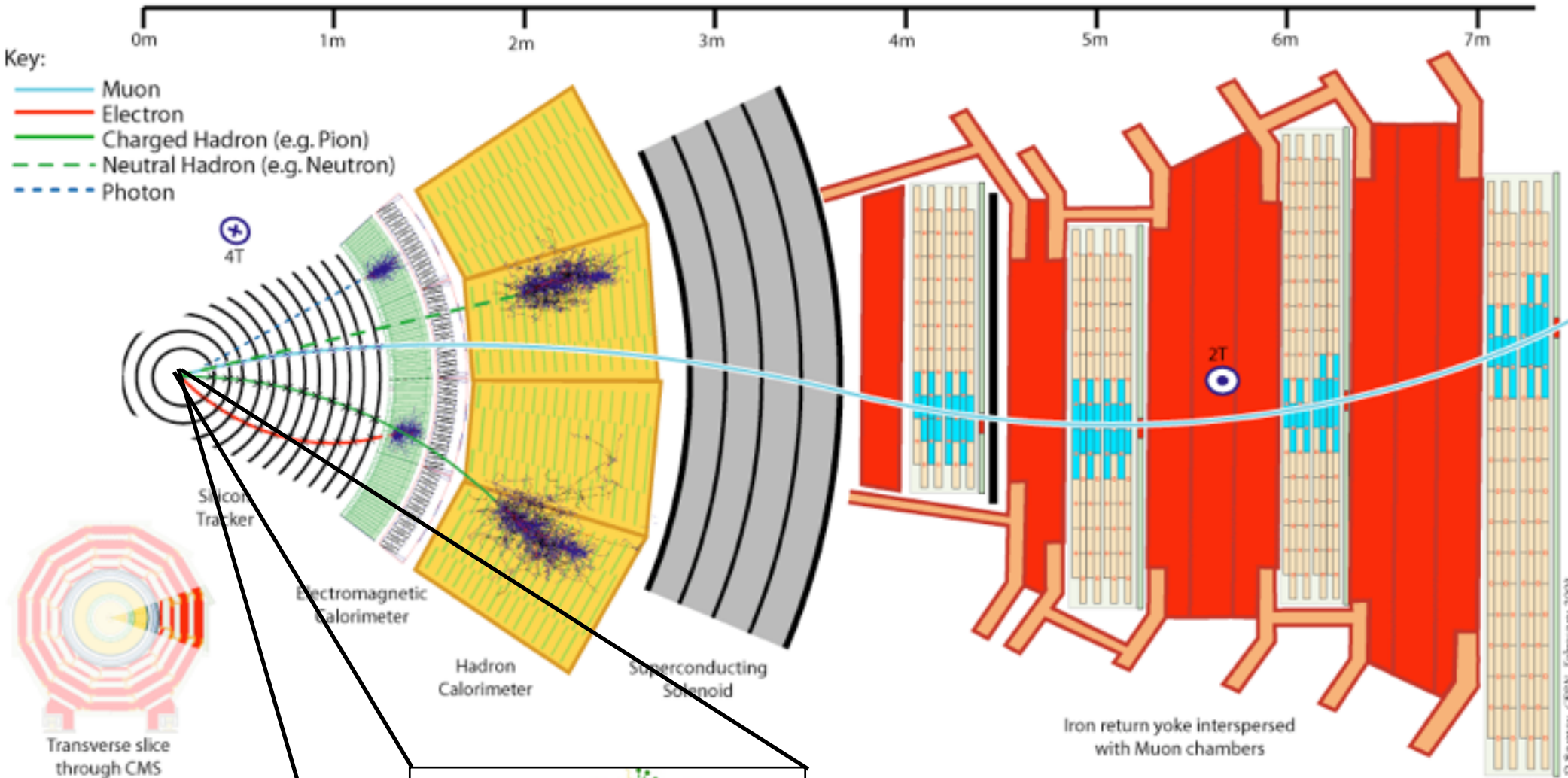
REAL TIME LINKAP 9  
SFR 12



**SLICE of  
the ATLAS  
BARREL**







$$\mathcal{L}_{SM} = \underbrace{\frac{1}{4} \mathbf{W}_{\mu\nu} \cdot \mathbf{W}^{\mu\nu} - \frac{1}{4} B_{\mu\nu} B^{\mu\nu} - \frac{1}{4} G_{\mu\nu}^a G_a^{\mu\nu}}_{\text{kinetic energies and self-interactions of the gauge bosons}}$$

$$+ \underbrace{\bar{L} \gamma^\mu (i \partial_\mu - \frac{1}{2} g_T \tau \cdot \mathbf{W}_\mu - \frac{1}{2} g' Y B_\mu) L + \bar{R} \gamma^\mu (i \partial_\mu - \frac{1}{2} g' Y B_\mu) R}_{\text{kinetic energies and electroweak interactions of fermions}}$$

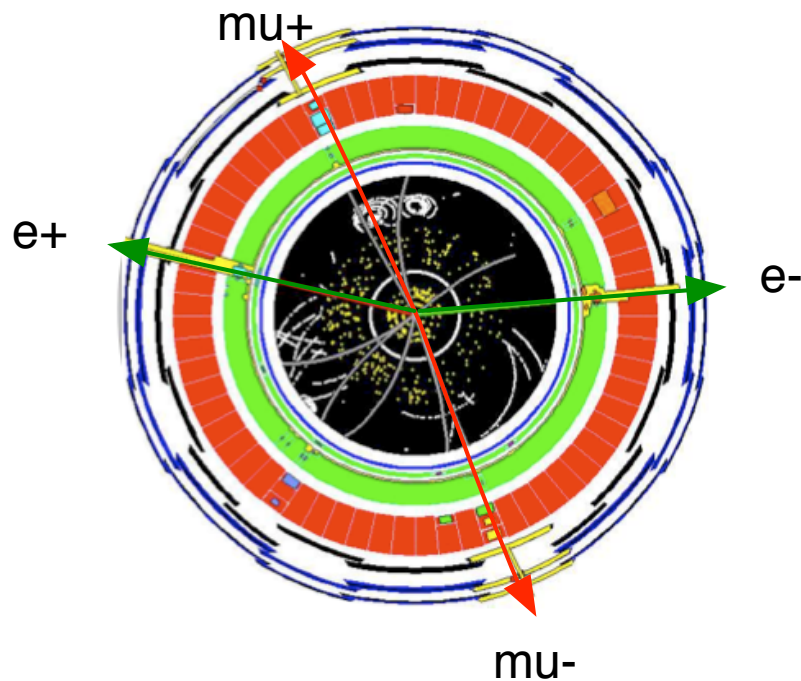
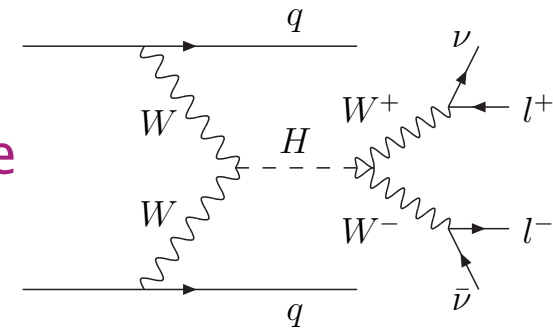
$$+ \frac{1}{2} \left[ (i \partial_\mu - \frac{1}{2} g_T \tau \cdot \mathbf{W}_\mu - \frac{1}{2} g' Y B_\mu) \phi \right]^2 - V(\phi)$$

$W^\pm, Z, \gamma$  and Higgs masses and couplings

$$+ \underbrace{g'' (\bar{q} \gamma^\mu T_a q) G_\mu^a}_{\text{interactions between quarks and gluons}} + \underbrace{(G_1 \bar{L} \phi R + G_2 \bar{L} \phi_c R + h.c.)}_{\text{fermion masses and couplings to Higgs}}$$

1) The language of the Standard Model is Quantum Field Theory

2) Perturbation Theory, Feynman Diagrams, and Factorization are used to construct Monte Carlo simulations of the interactions



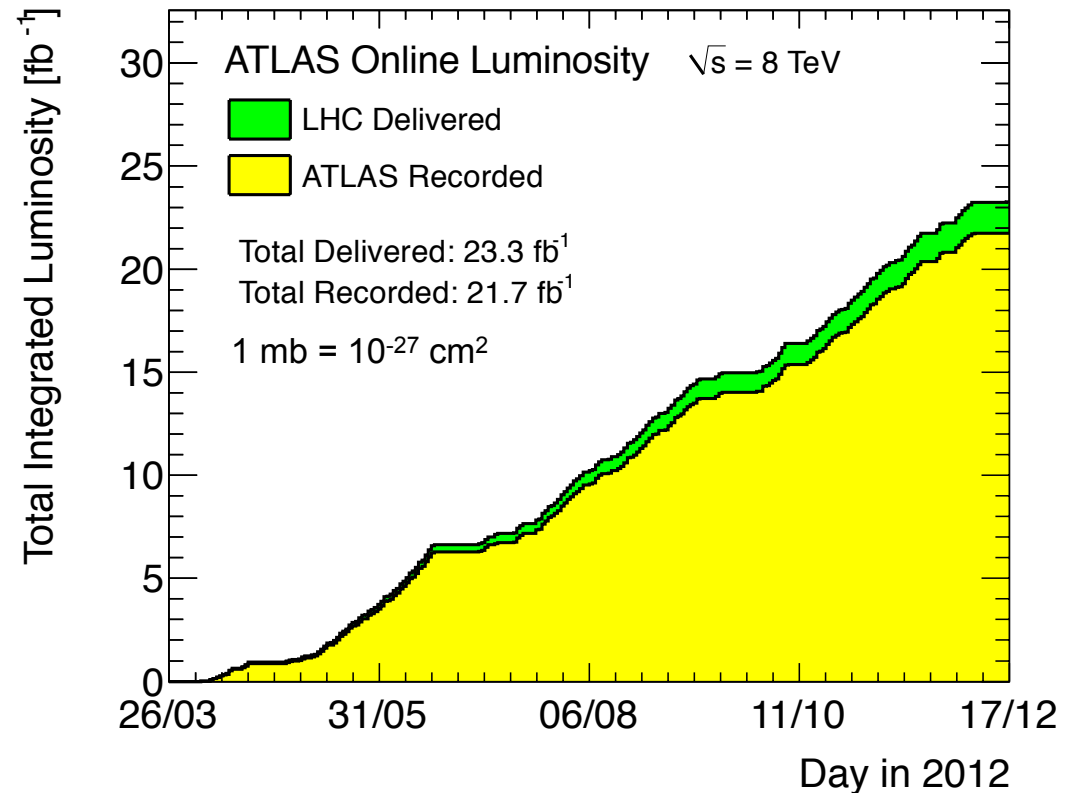
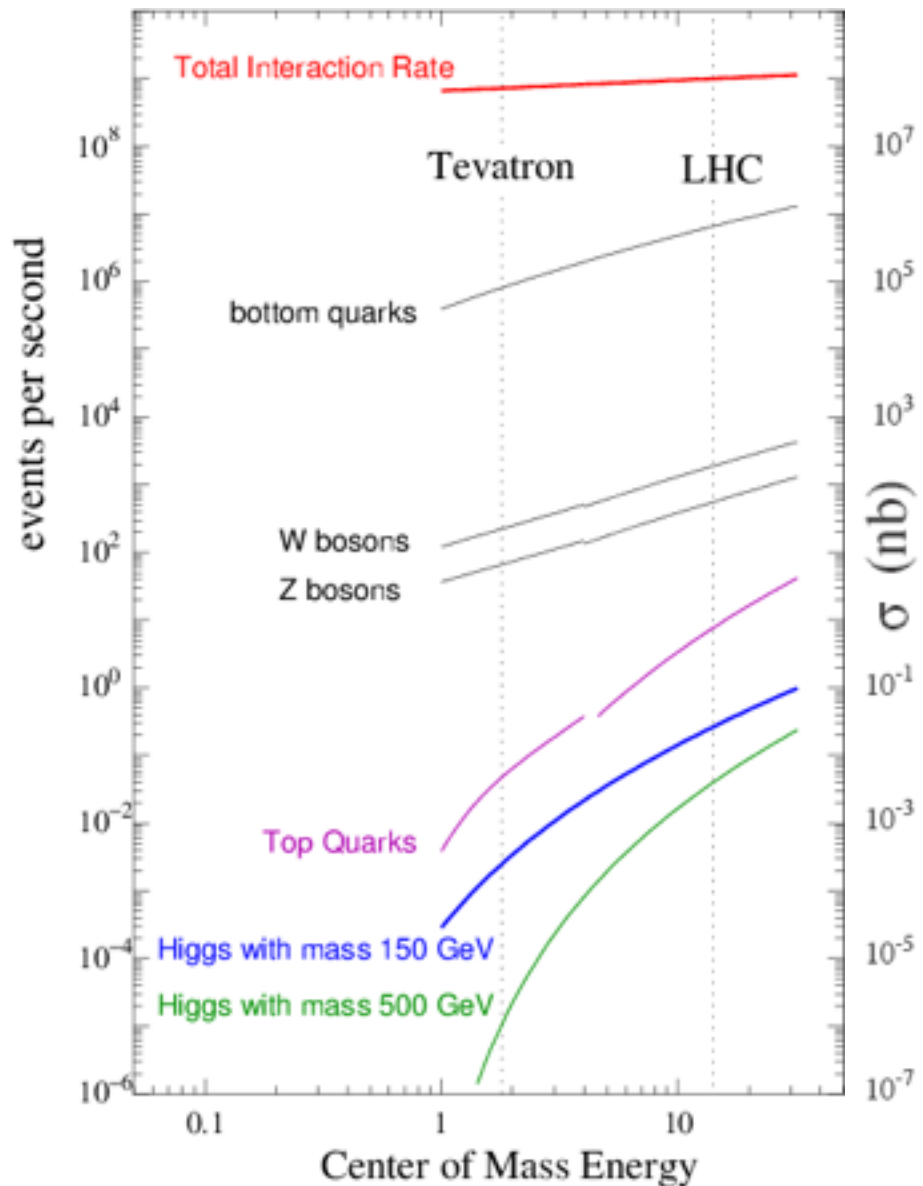
3) The interaction of outgoing particles with the detector is simulated.

4) Finally, we run algorithms on the simulated data as if they were from real collisions.

# Number of collisions

expected number of scatterings = cross section [cm<sup>2</sup>] x Luminosity [1/cm<sup>2</sup>]

$$80 \text{ mb} \cdot 25 \text{ fb}^{-1} = 2 \cdot 10^{15} \text{ collisions}$$



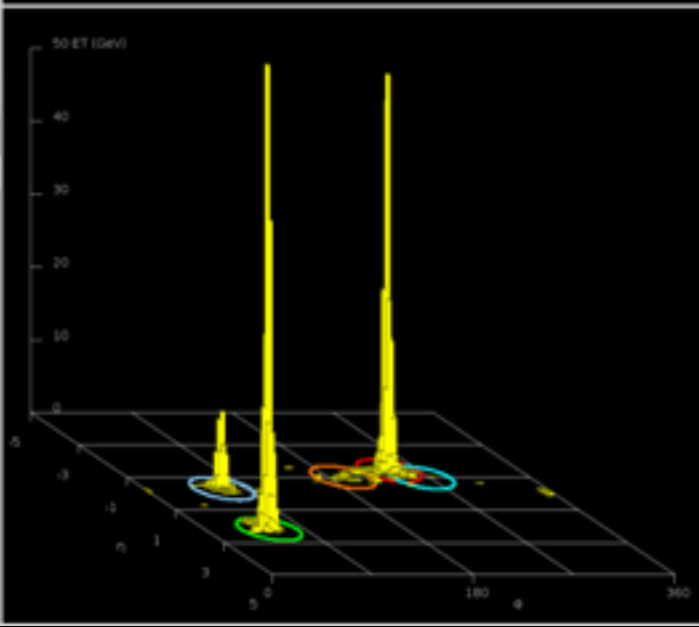
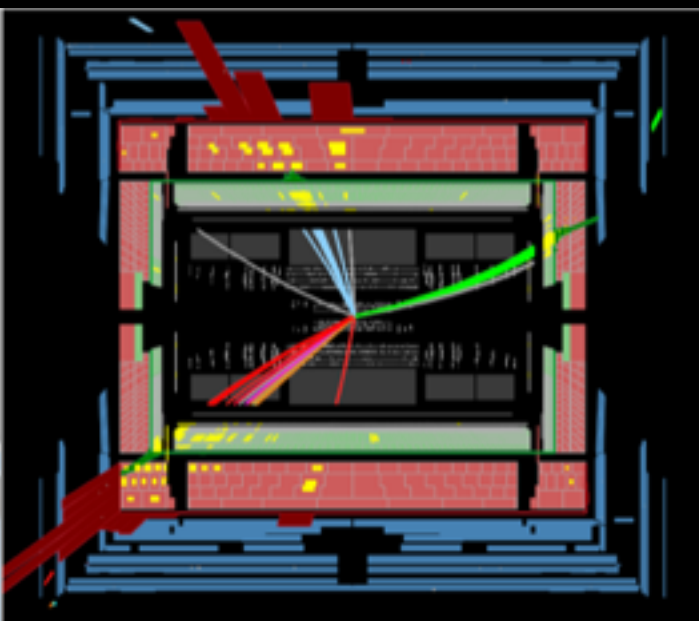
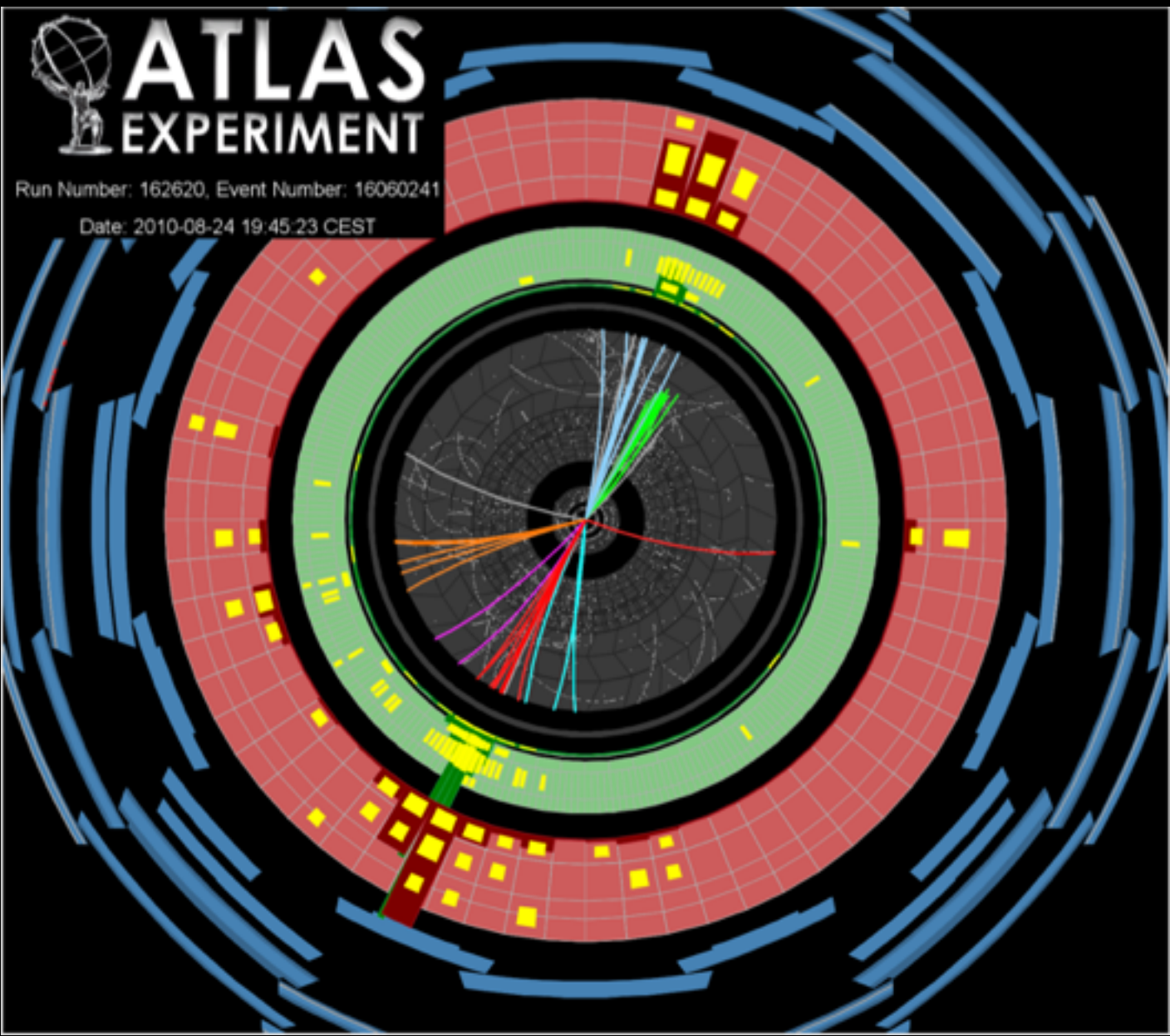
signal : background  $\sim 1 : 10^9$



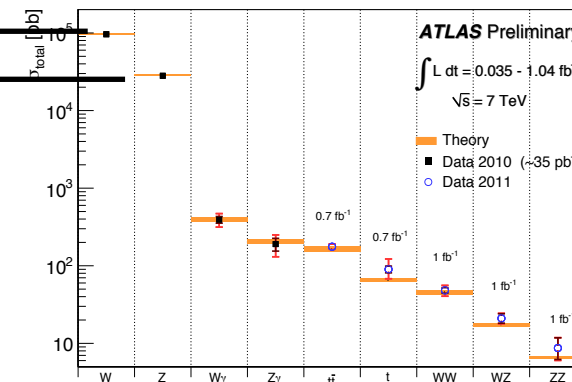
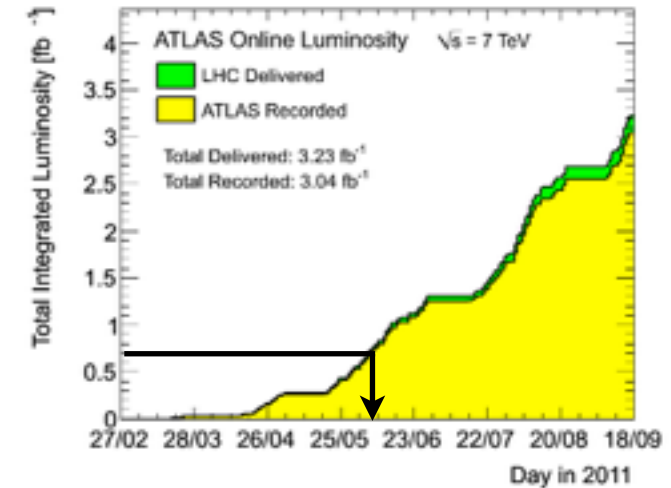
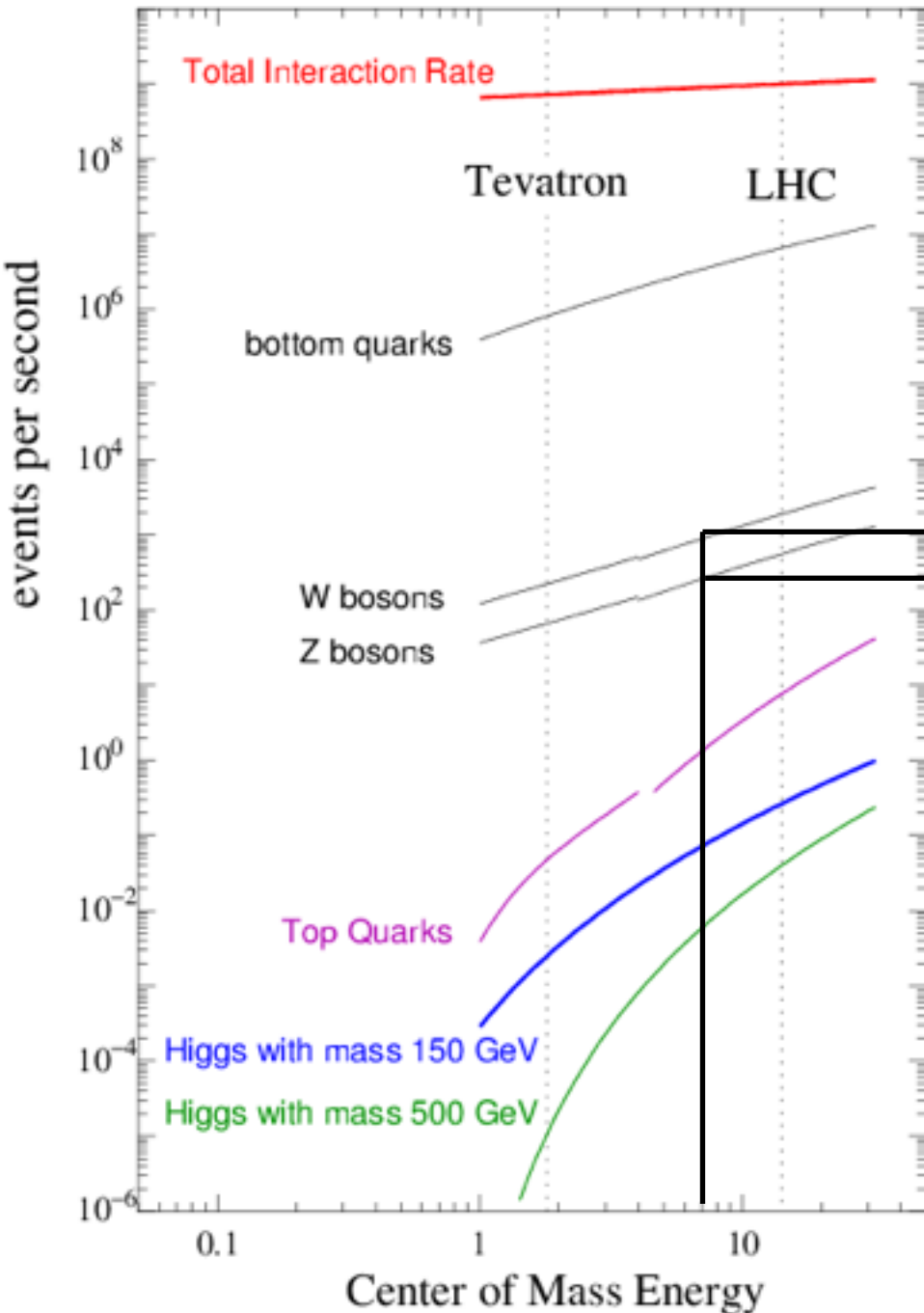
# ATLAS EXPERIMENT

Run Number: 162620, Event Number: 16060241

Date: 2010-08-24 19:45:23 CEST



# The steady march of progress

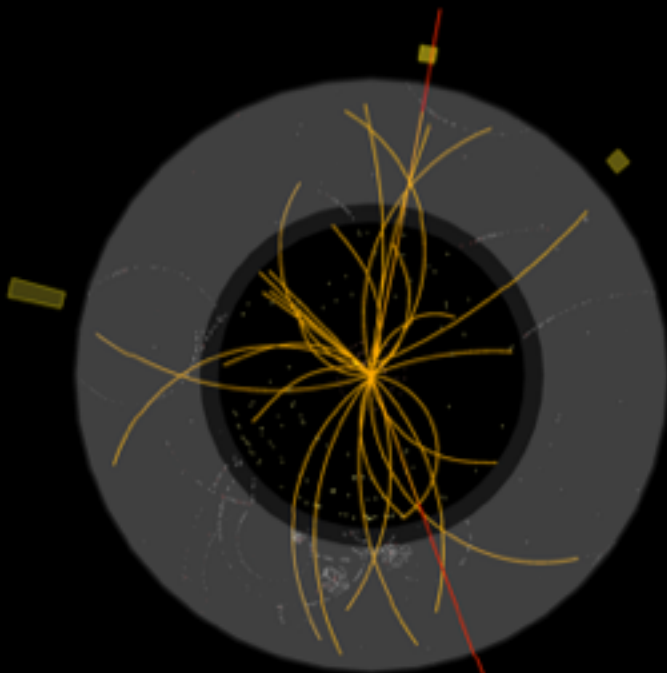






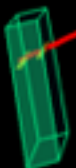
# ATLAS EXPERIMENT

Run: 154822, Event: 14321500  
Date: 2010-05-10 02:07:22 CEST

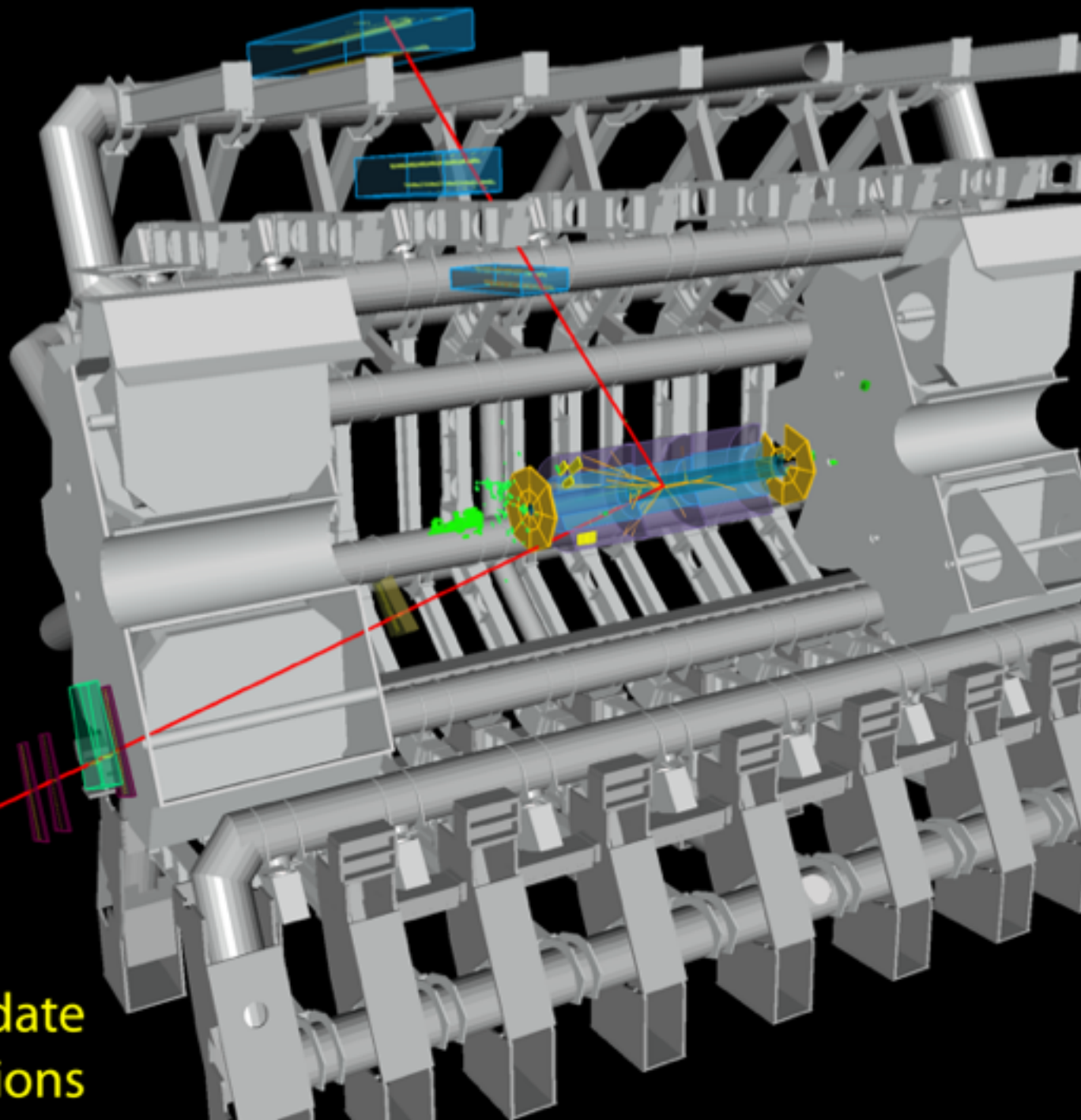


$p_T(\mu^-) = 27 \text{ GeV}$   $\eta(\mu^-) = 0.7$   
 $p_T(\mu^+) = 45 \text{ GeV}$   $\eta(\mu^+) = 2.2$

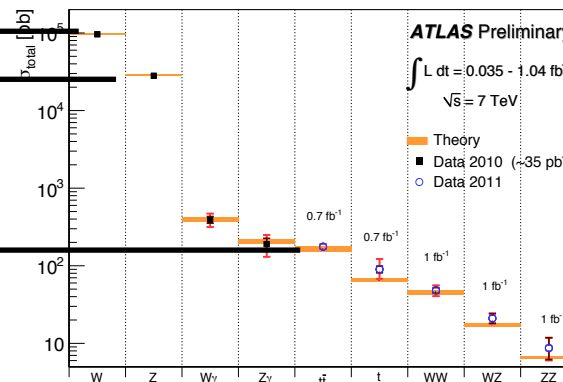
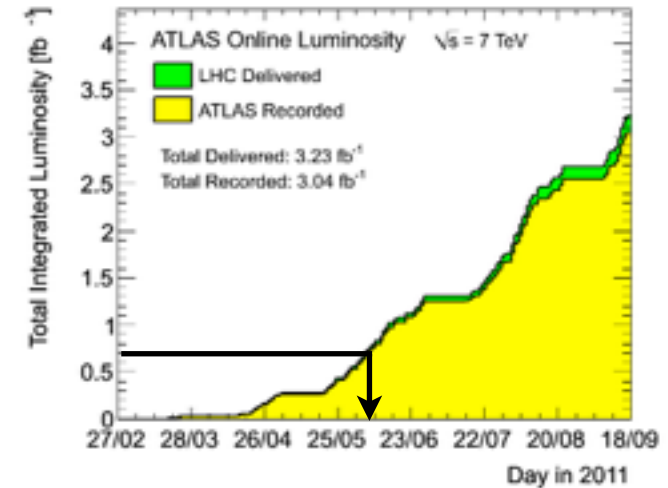
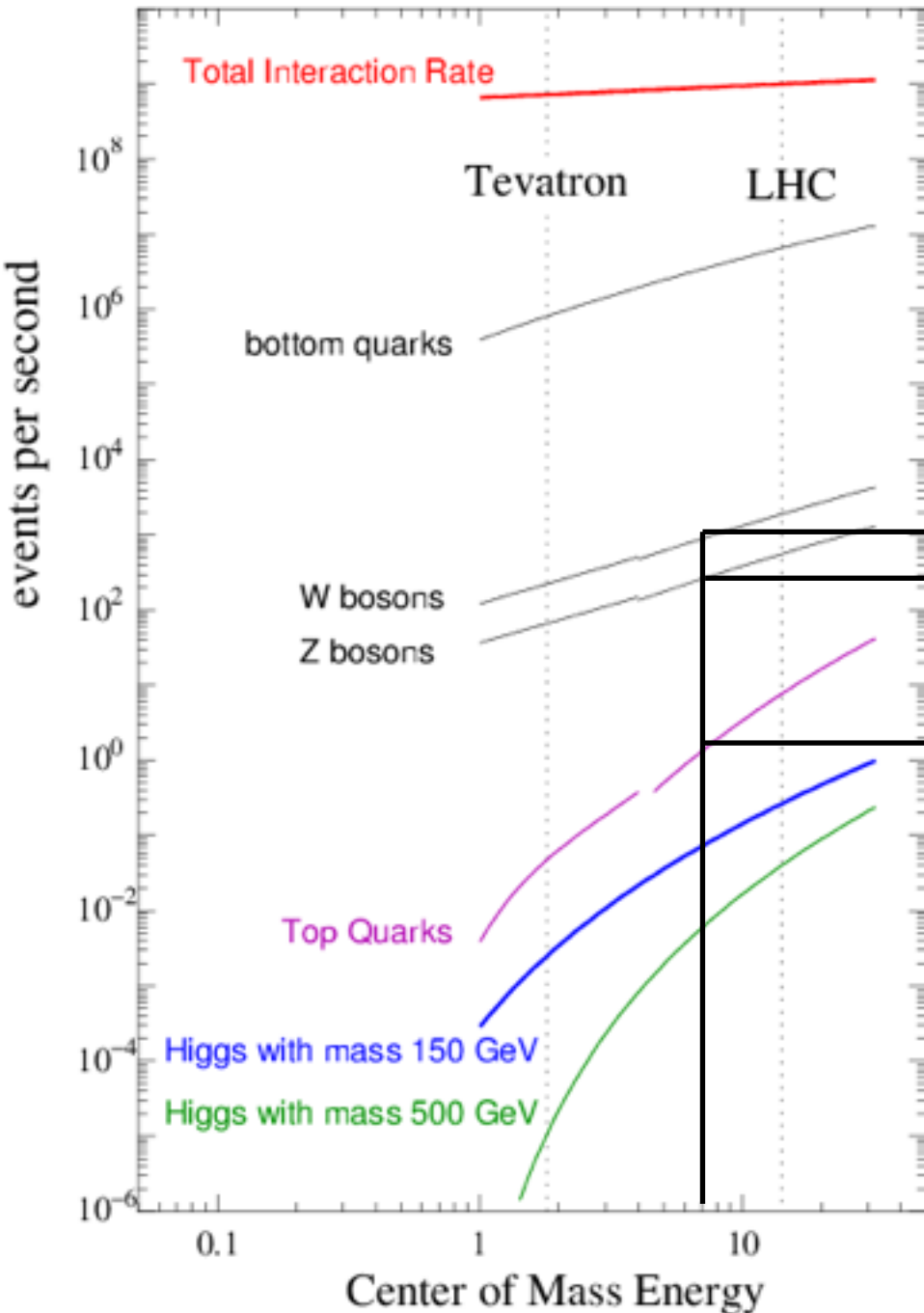
$M_{\mu\mu} = 87 \text{ GeV}$



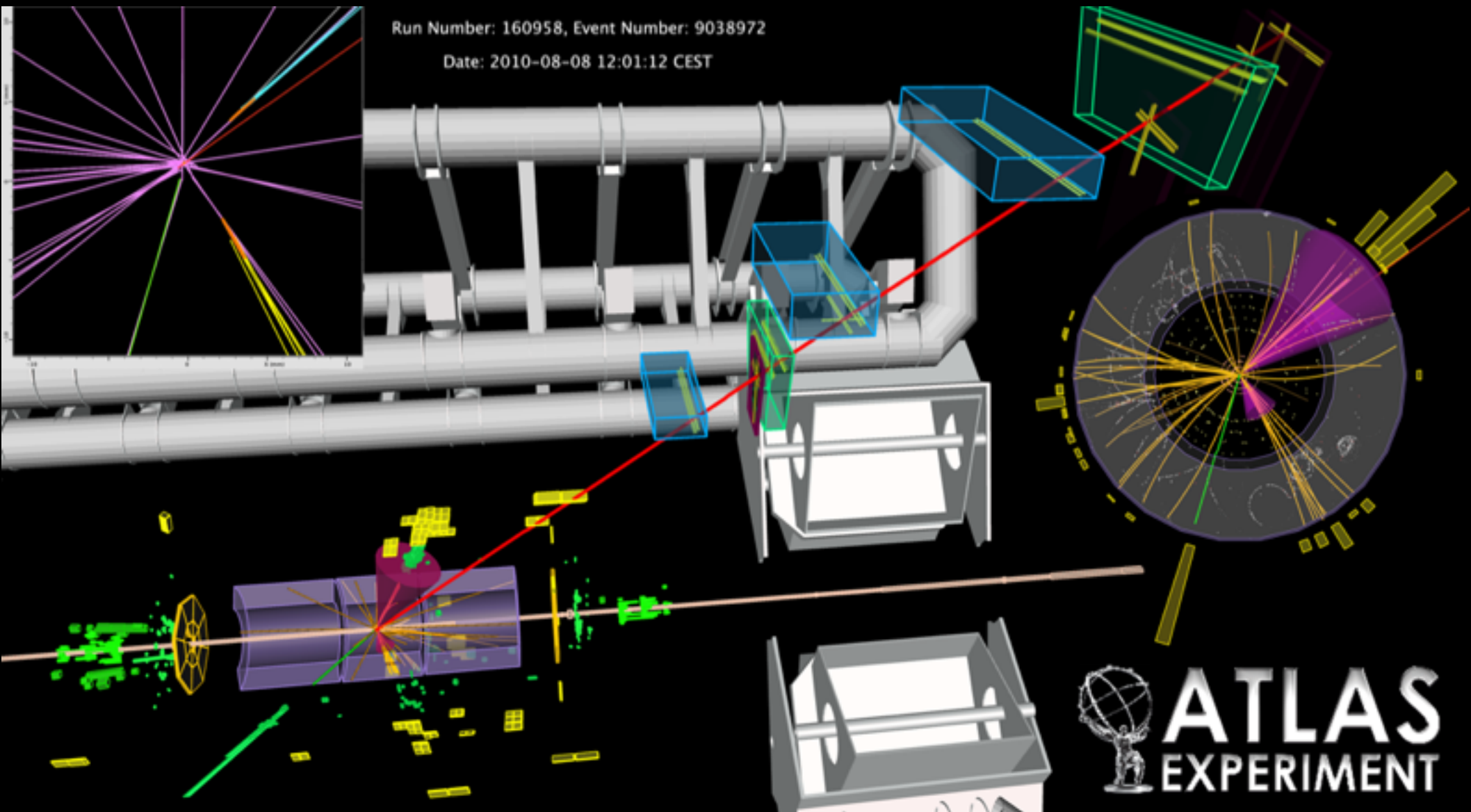
**$Z \rightarrow \mu\mu$  candidate  
in 7 TeV collisions**



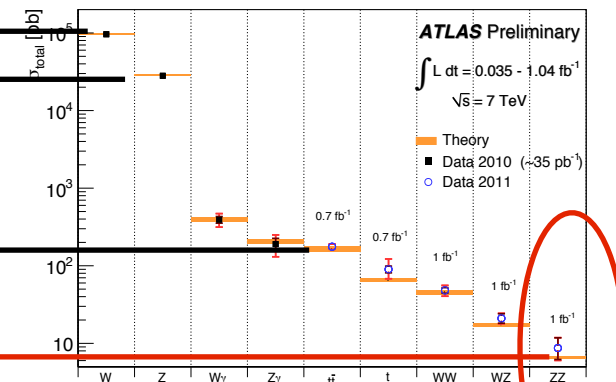
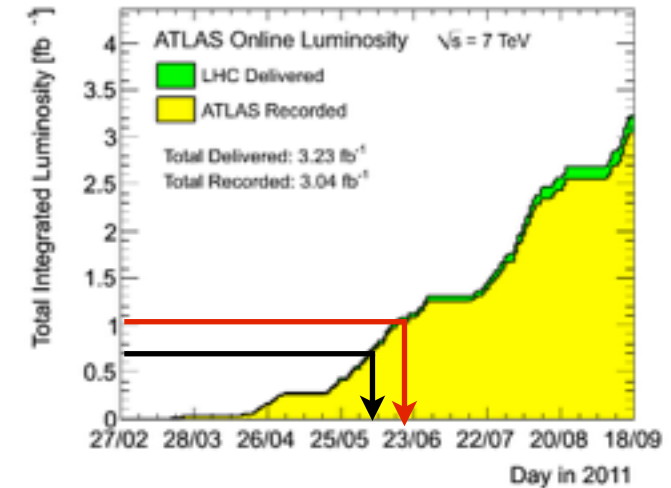
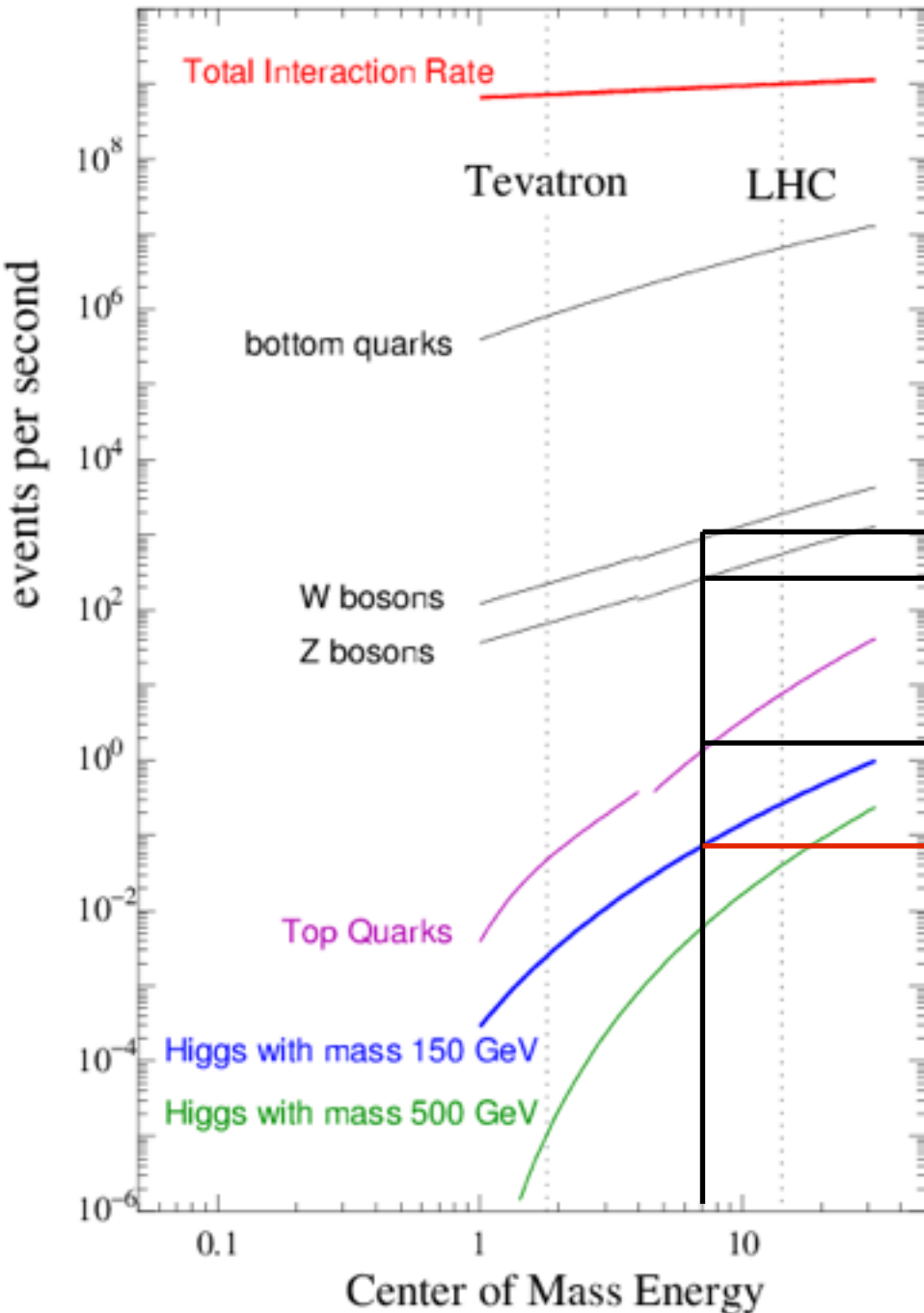
# The steady march of progress



# Top quark pair decaying to $b\bar{b} e\mu E_{T,miss}$



# The steady march of progress

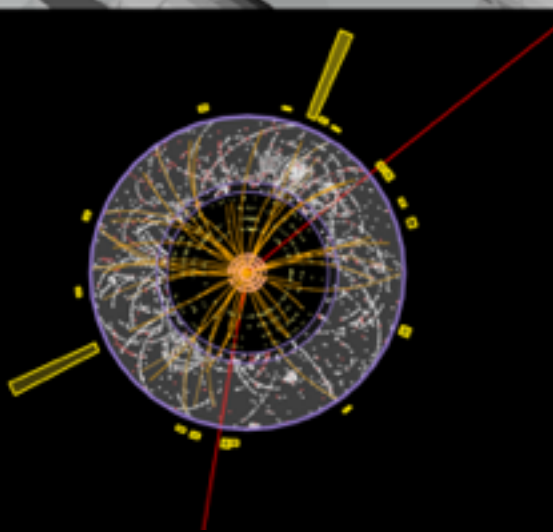
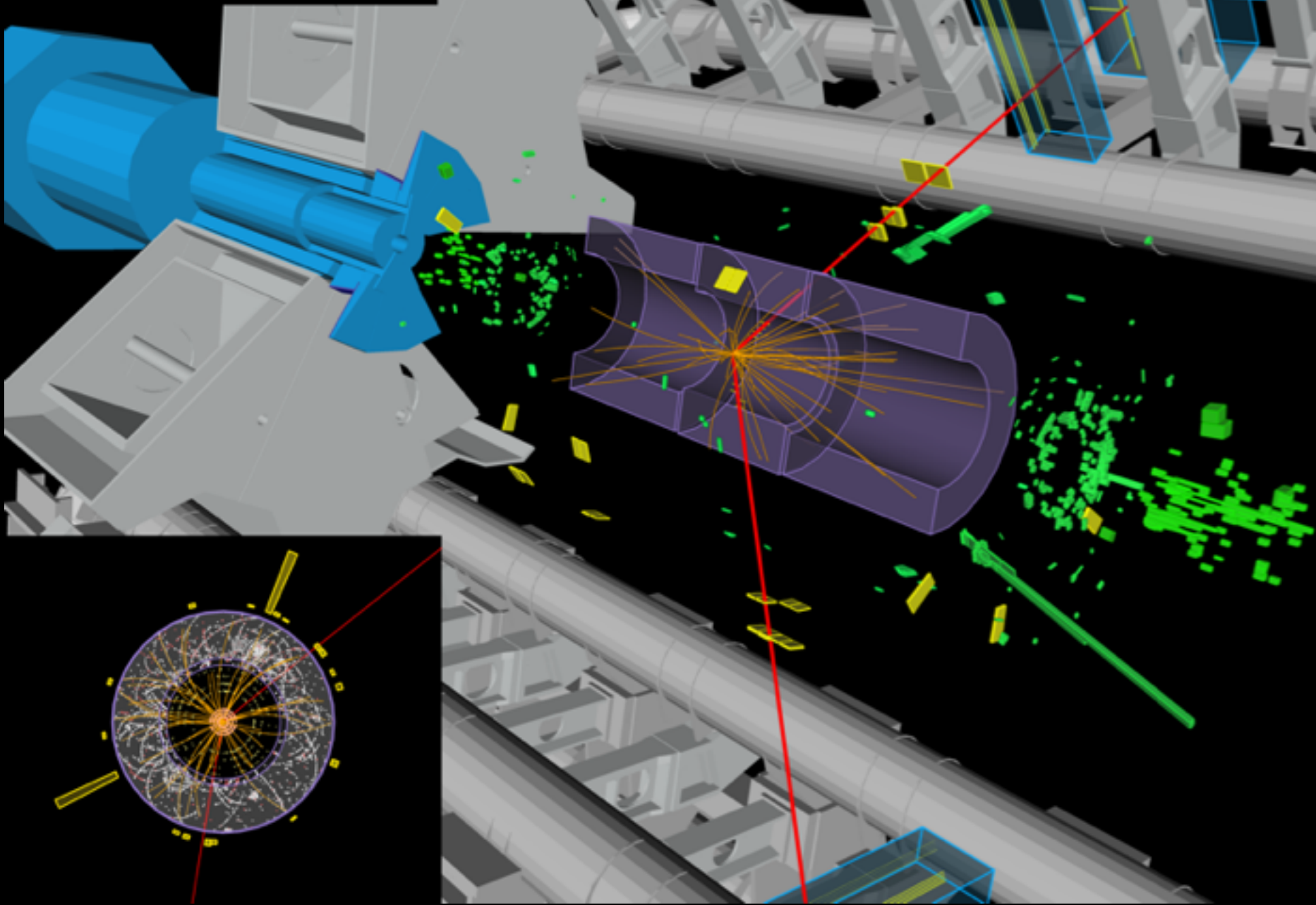




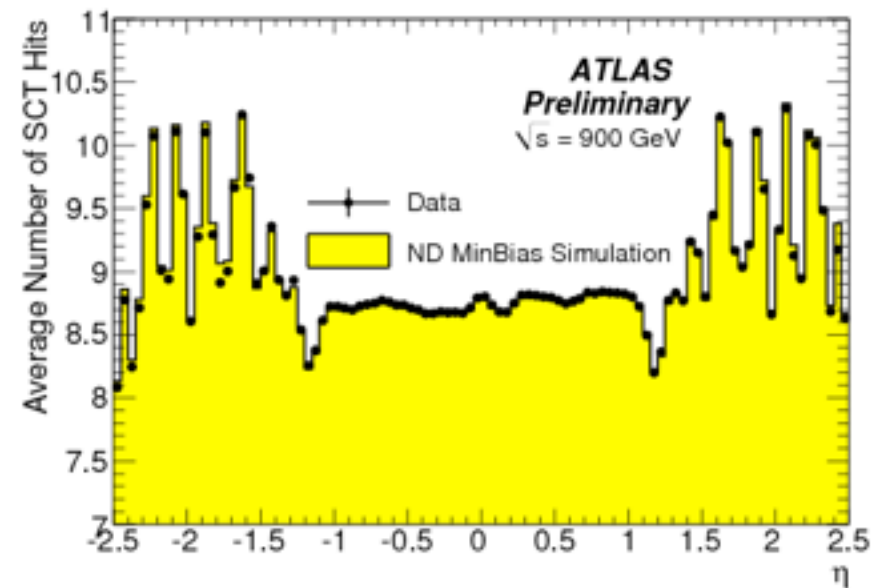
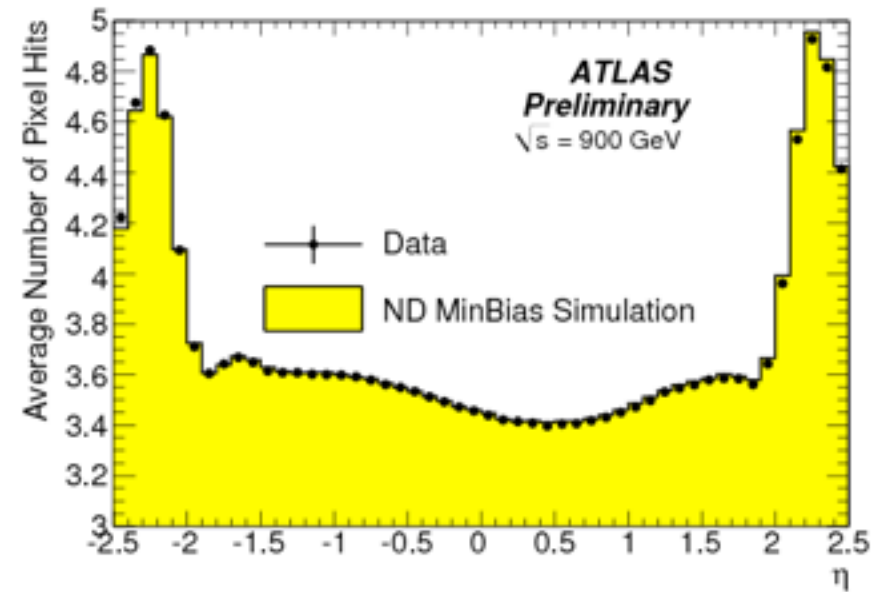
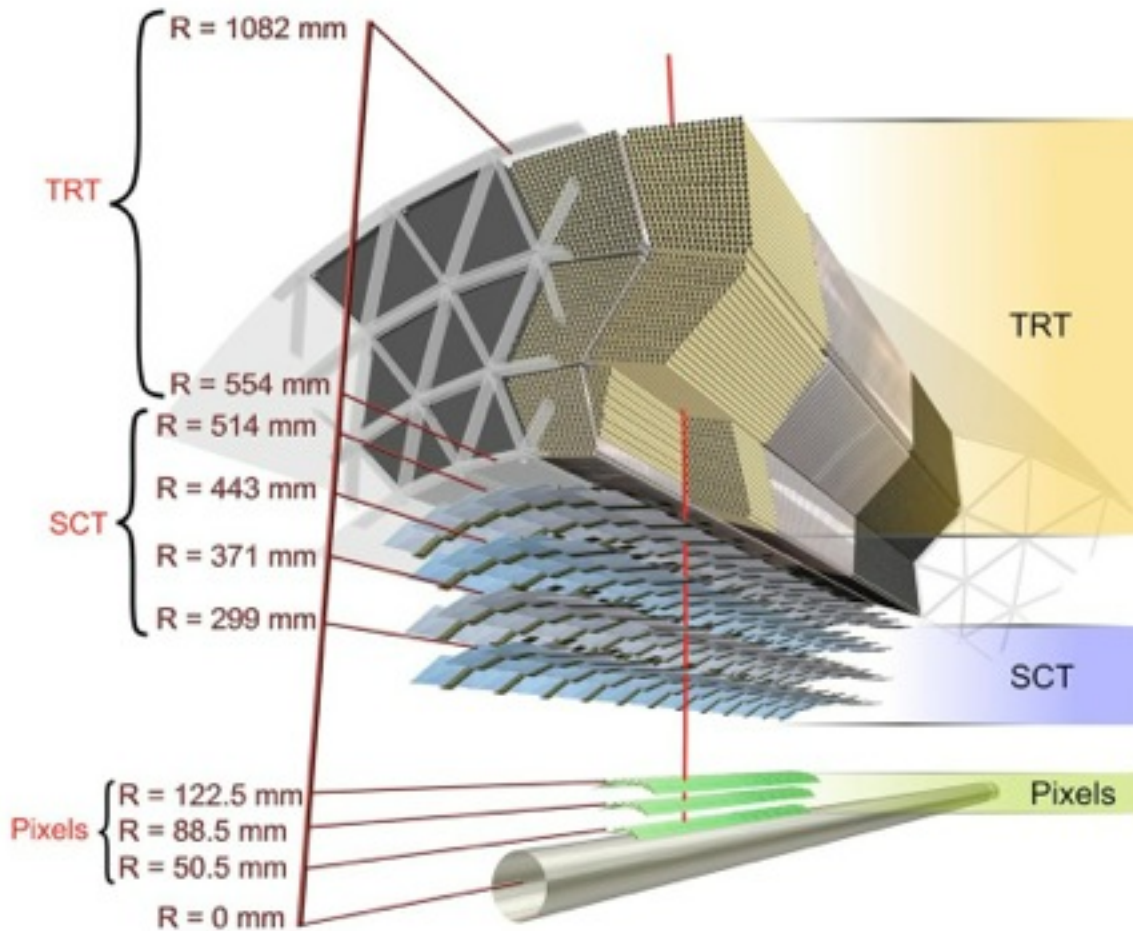
# ATLAS EXPERIMENT

Run Number: 182747, Event Number: 63217197

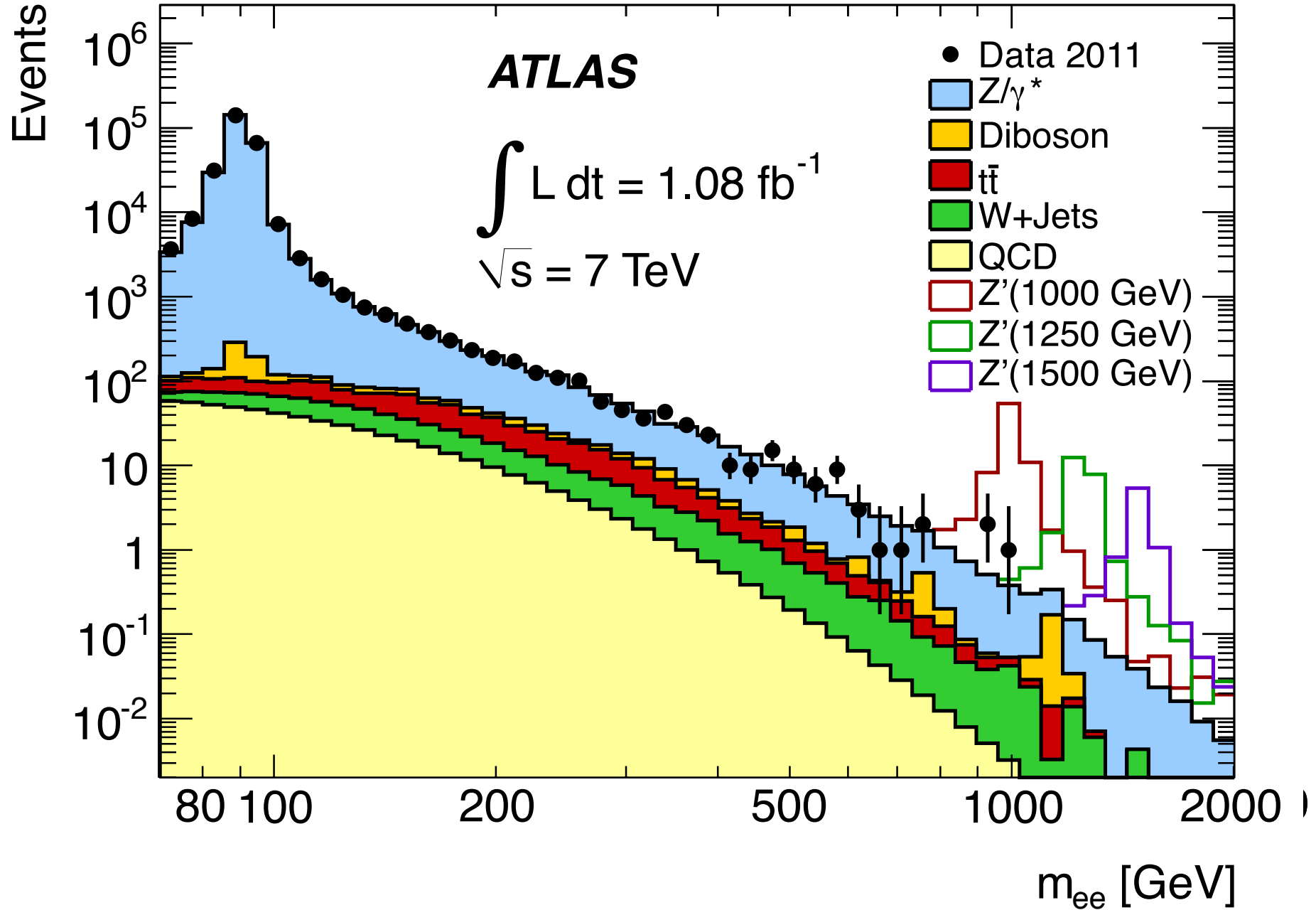
Date: 2011-05-28 13:06:57 CEST



# How good is the modeling?



# How good is the modeling?

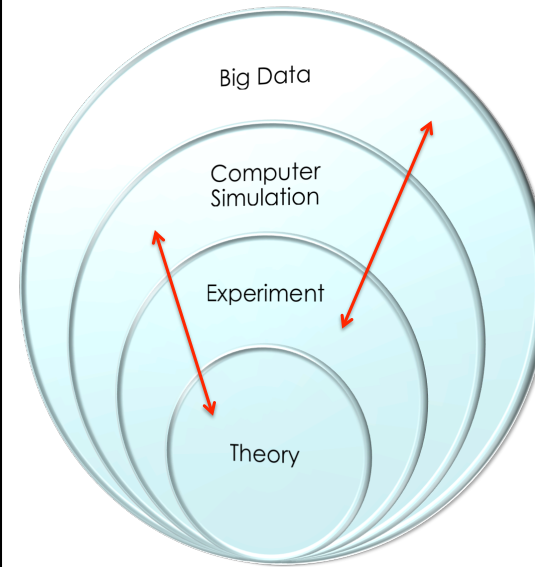


# Complex Models for Big Data



Max Welling  
UvA

# The Four Paradigms

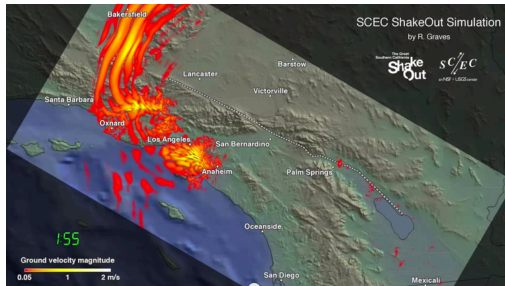


We have added big data to computer simulation, experiment and theory.



Not replaced it...

# Big Simulation



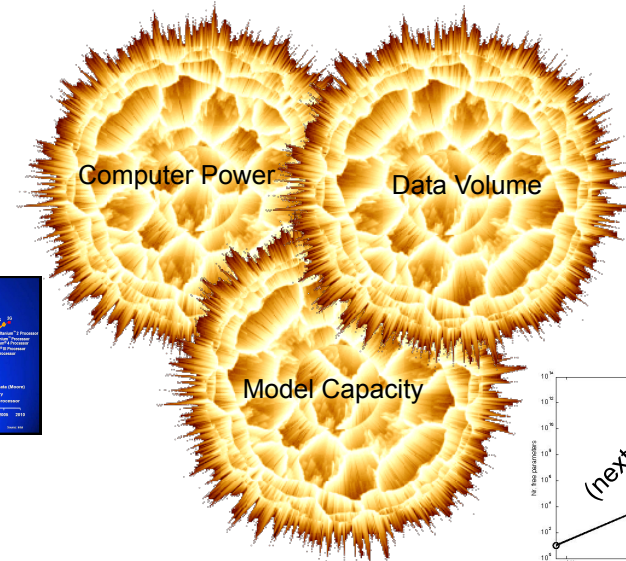
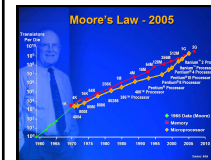
Computer simulations have become increasingly complex (e.g. weather, earthquake models)

This production run producing 360 sec of wave propagation sustained 220 Tflop/s for 24 hours on NCCS Jaguar using 223,074 cores.

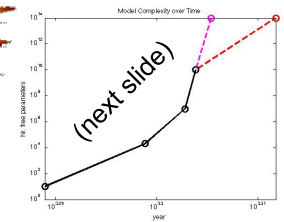
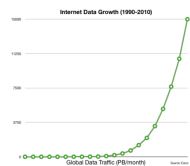
The Computational Wall: If a model has hundreds of parameters, how can we:

- 1) Find the parameter values that match the observations best?
- 2) Determine if we underfit (model too simple) or overfit (model too complex)?
- 3) Compare two models?

# 3x Exponential Growth in Machine Learning



Data is Growing Exponentially



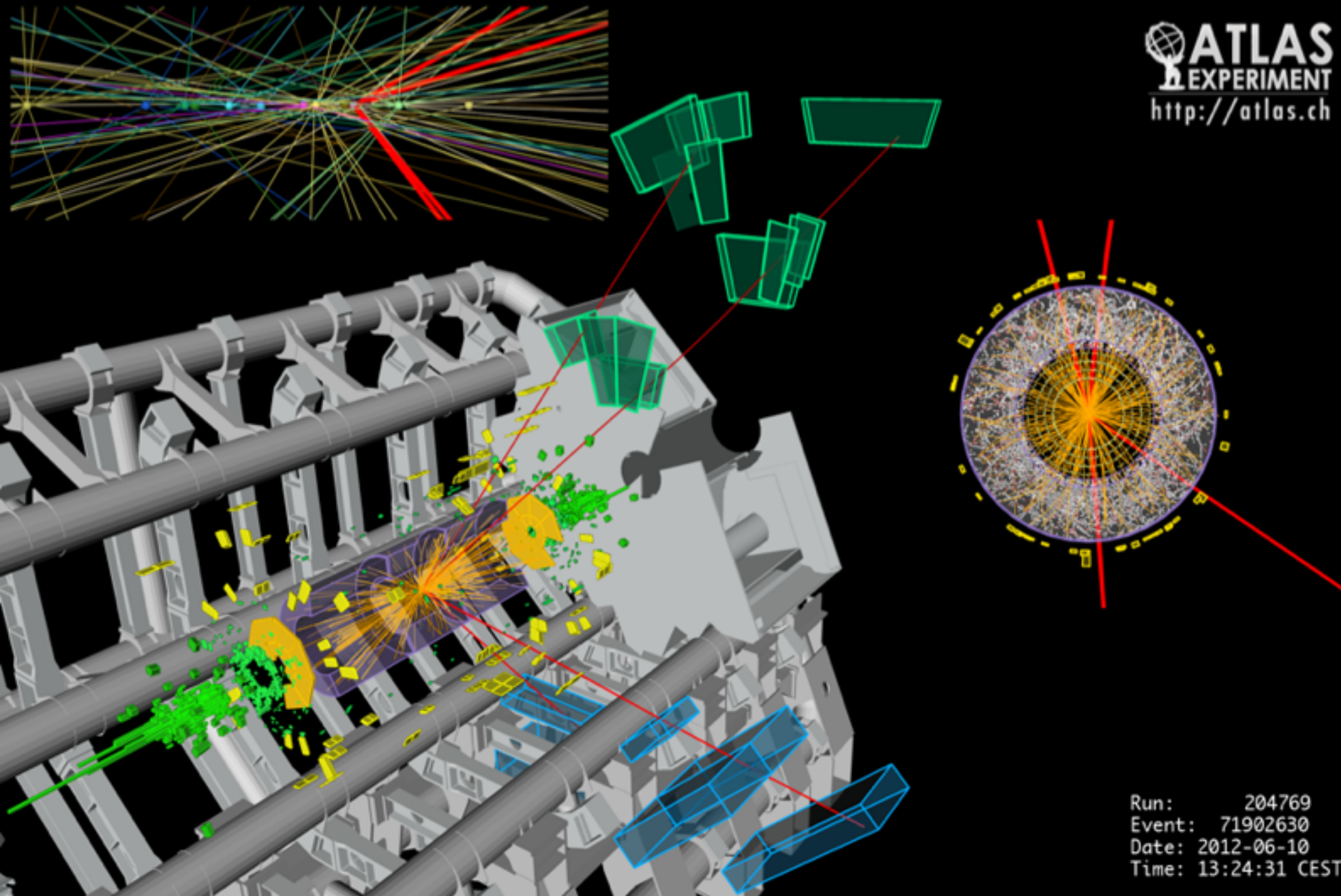


# Use of Machine Learning:

## Particle-Level and Event-Level

$$H \rightarrow ZZ \rightarrow 4l$$

 **ATLAS**  
EXPERIMENT  
<http://atlas.ch>



Run: 204769  
Event: 71902630  
Date: 2012-06-10  
Time: 13:24:31 CEST

# Putting the Higgs back together again

Don't believe the media:  $E \neq mc^2$

What Einstein really said:

$$E^2 = (mc^2)^2 + (|\vec{p}|c)^2$$

Every physics student knows energy and momentum are conserved

$$E_{\text{Higgs}} = E_{\text{before}} = E_{\text{after}} = \sum_i E_i$$
$$\vec{p}_{\text{Higgs}} = \vec{p}_{\text{before}} = \vec{p}_{\text{after}} = \sum_i \vec{p}_i$$

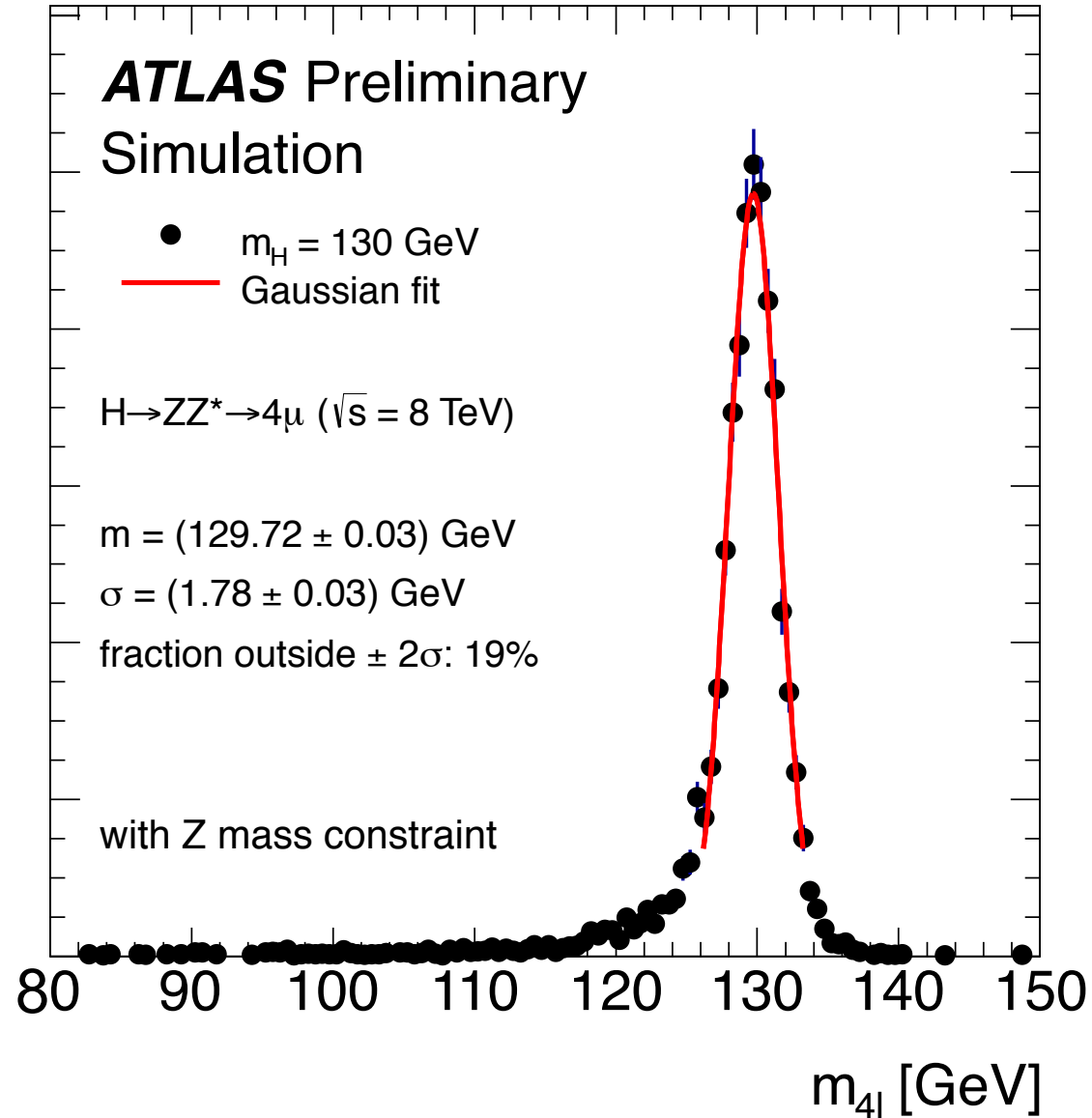
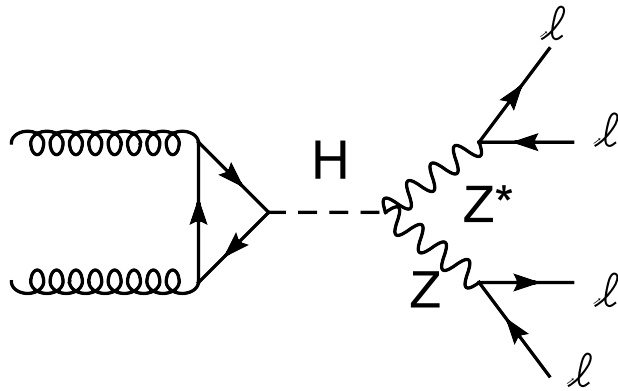
Thus, we can estimate the mass of the Higgs with

$$m_H = \sqrt{E_{\text{after}}^2/c^4 - |\vec{p}_{\text{after}}|^2/c^2}$$

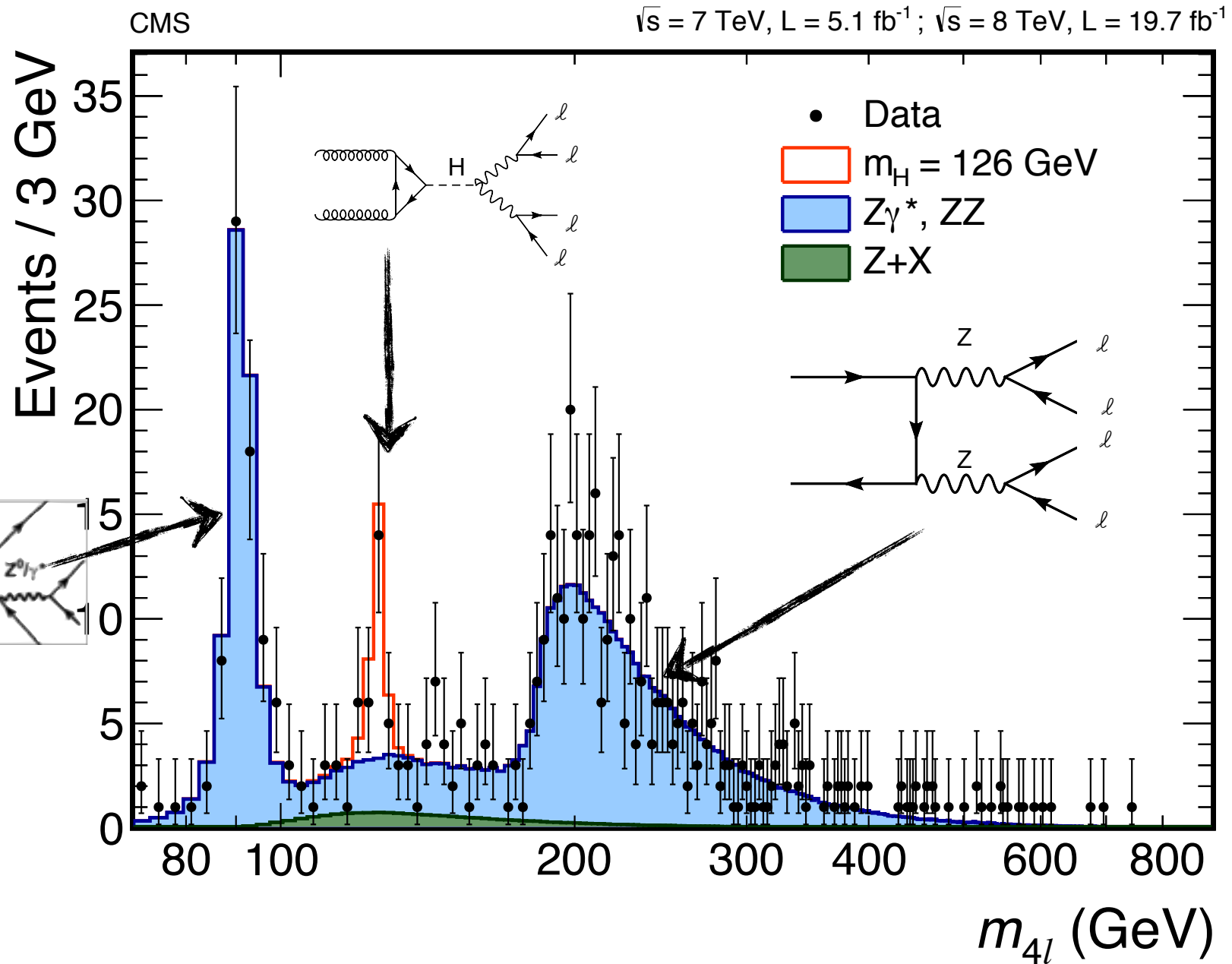
# An example high-level feature

From the 16 energies and momenta measured in this system, this particular combination gives a very sharp feature.

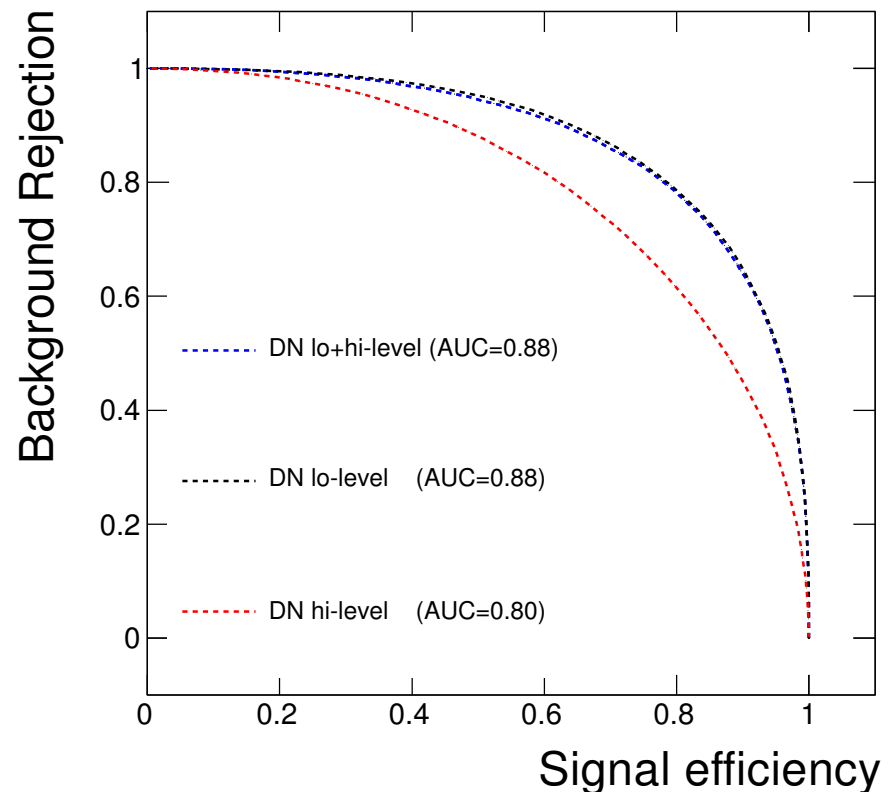
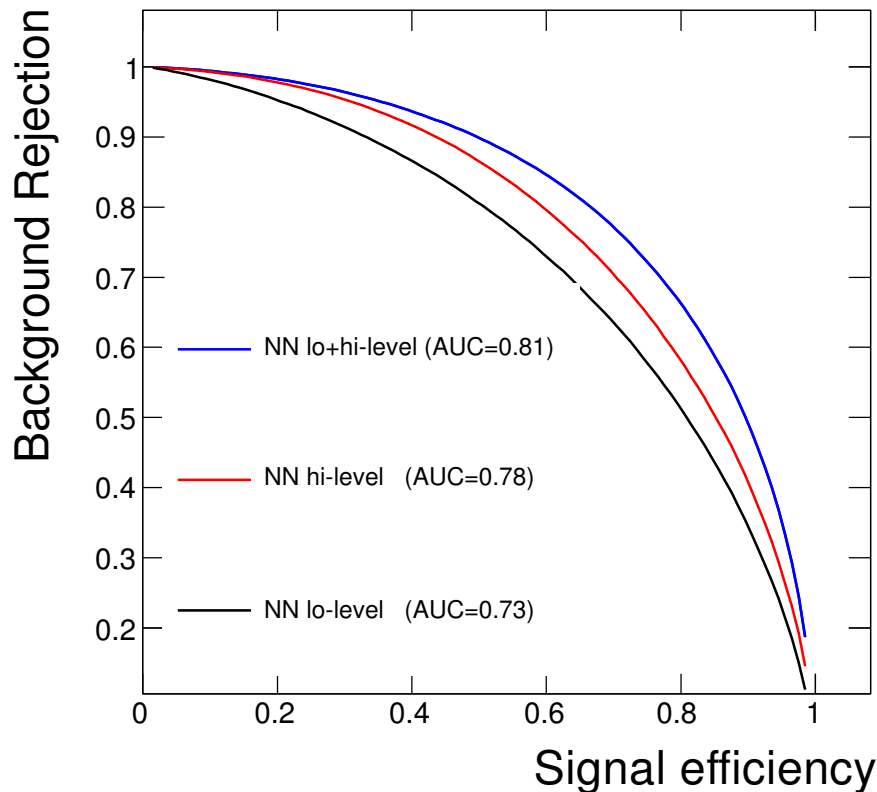
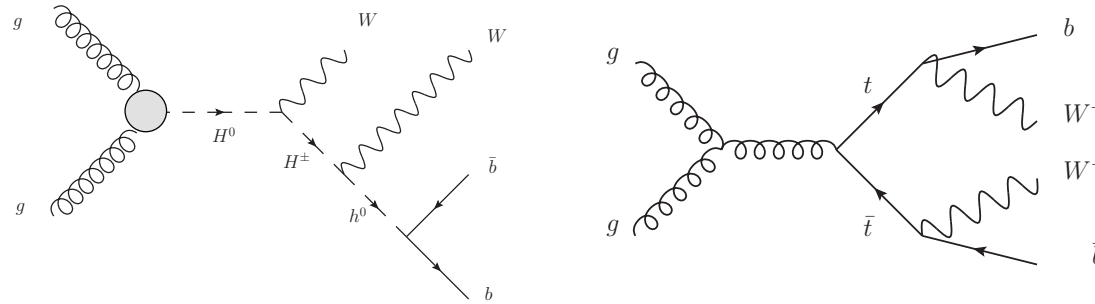
~sufficient statistic



# The observation in the 4l channel

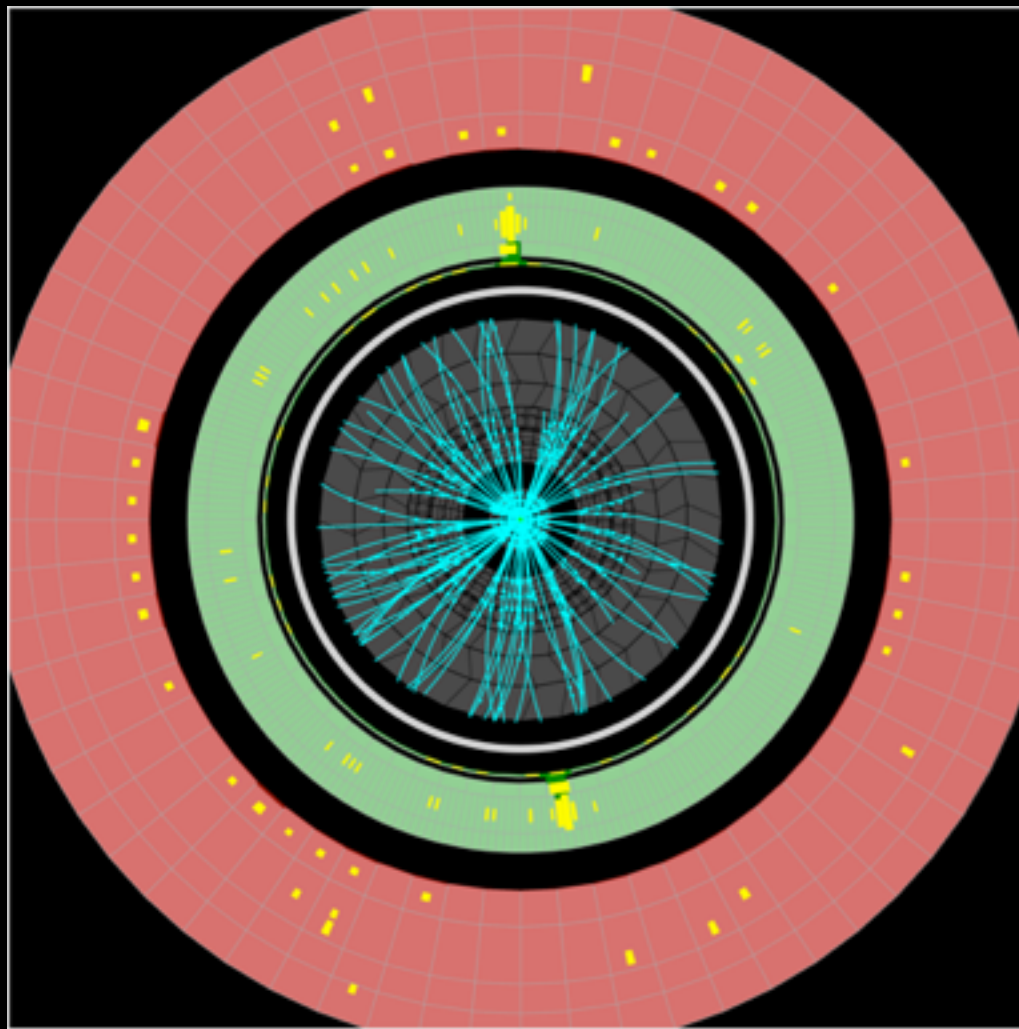


ML techniques performed poorly unless these high-level features were supplied. Deep learning techniques can discover them.



P. Baldi, P. Sadowski, and D. Whiteson [arXiv:1402.4735] GPU-accelerated Theano and Pylearn2 <https://github.com/uci-igb/higgs-susy>.

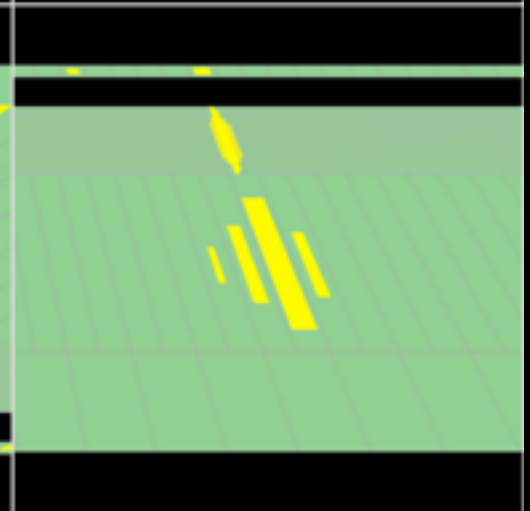
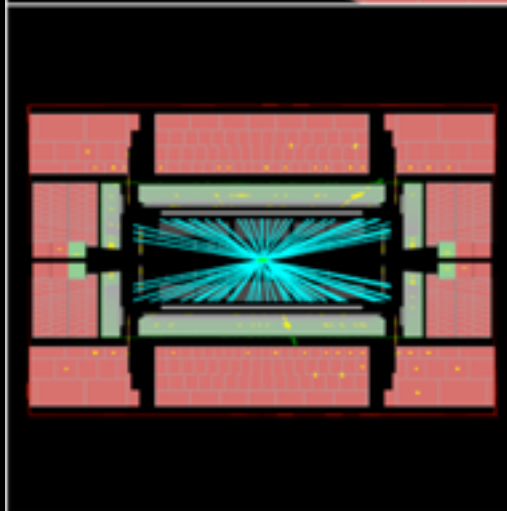
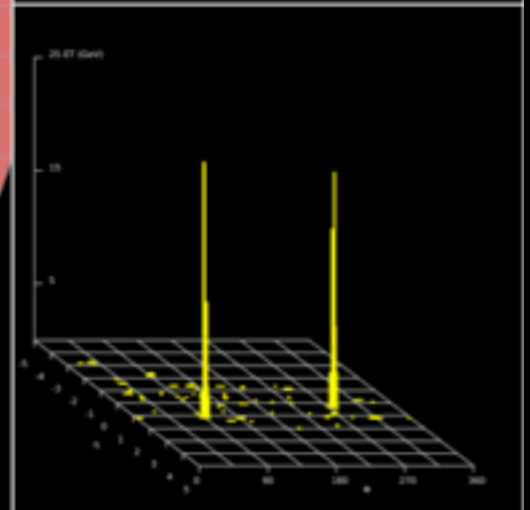
$$H \rightarrow \gamma\gamma$$



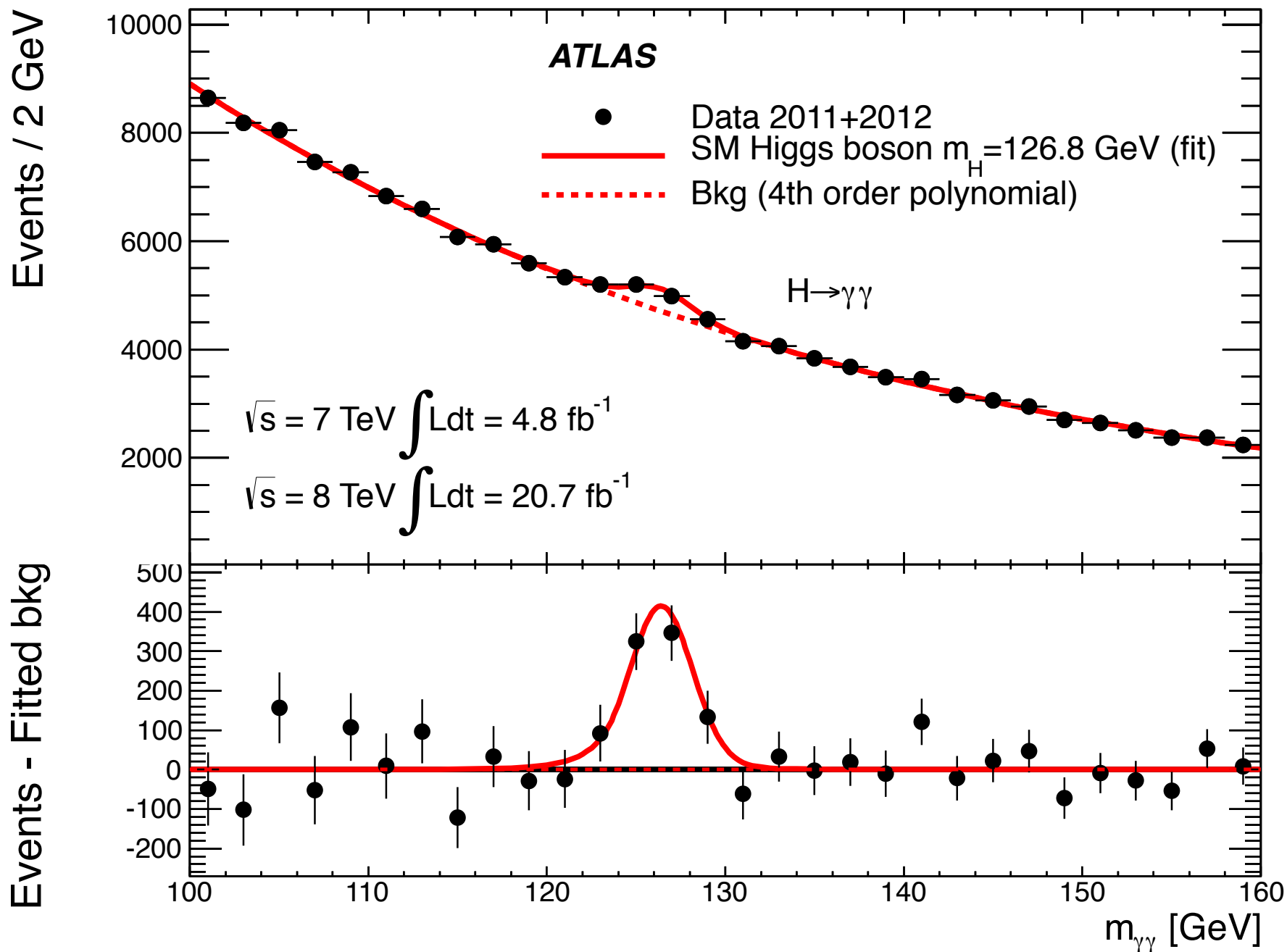
 **ATLAS**  
EXPERIMENT

Run Number: 203779, Event Number: 56662314

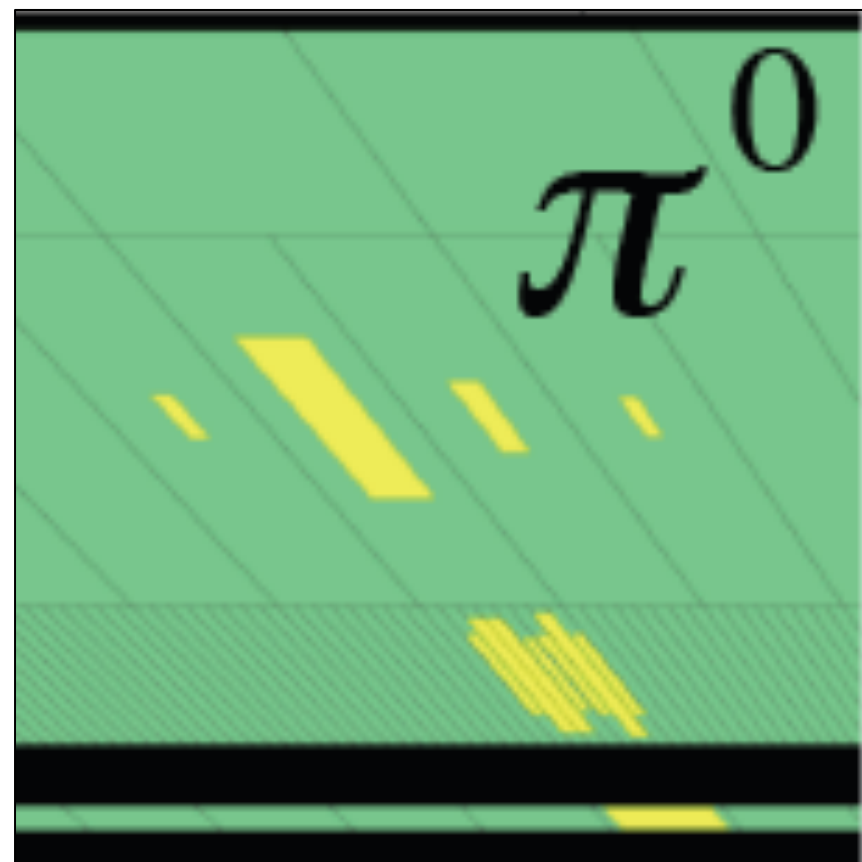
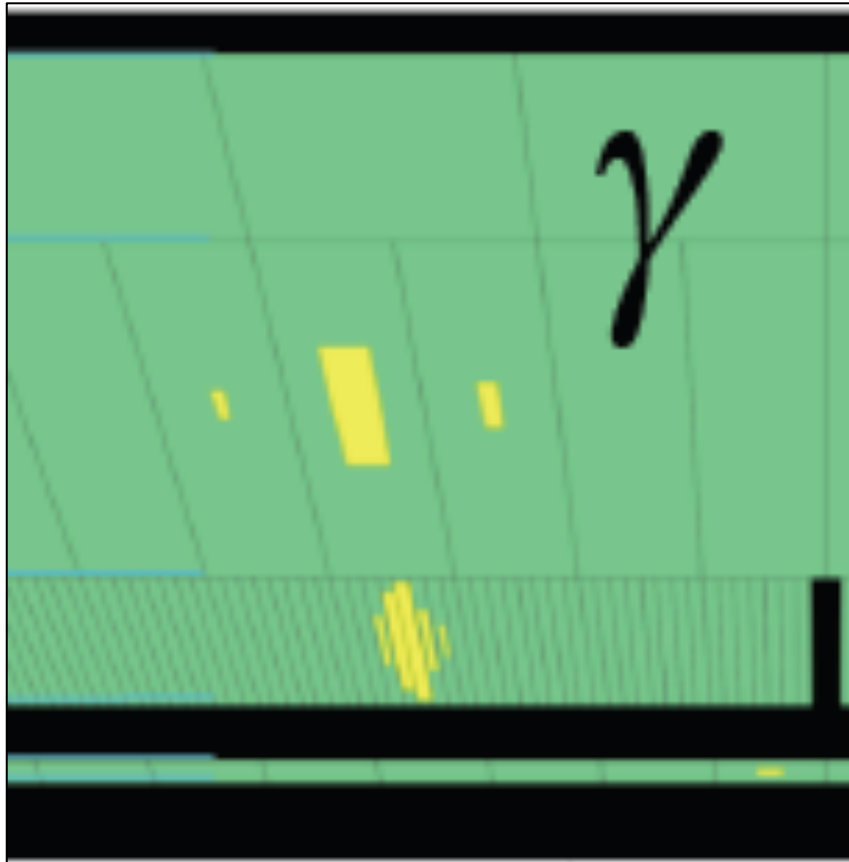
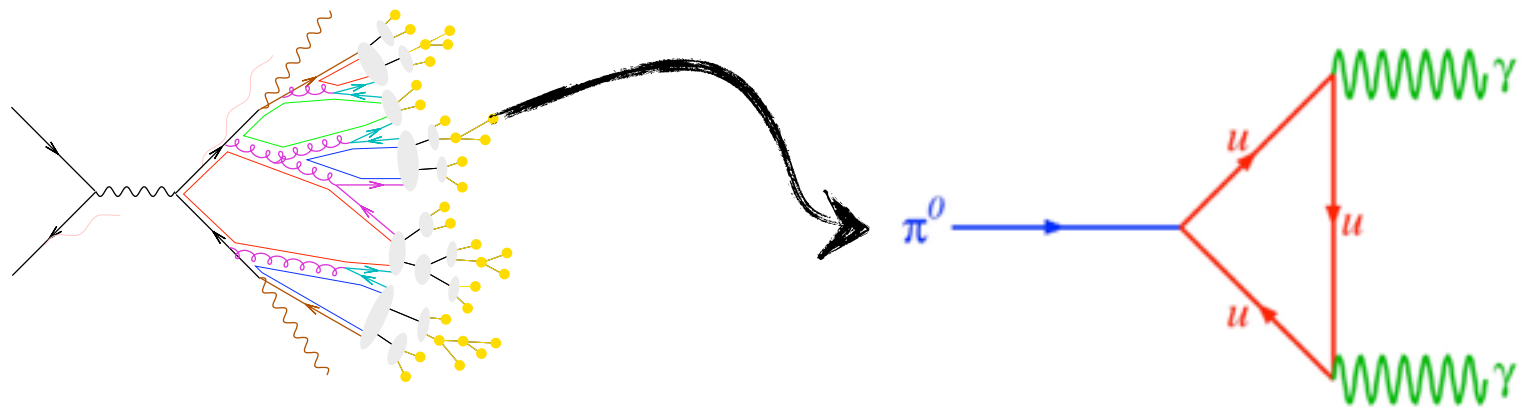
Date: 2012-05-23 22:19:29 CEST

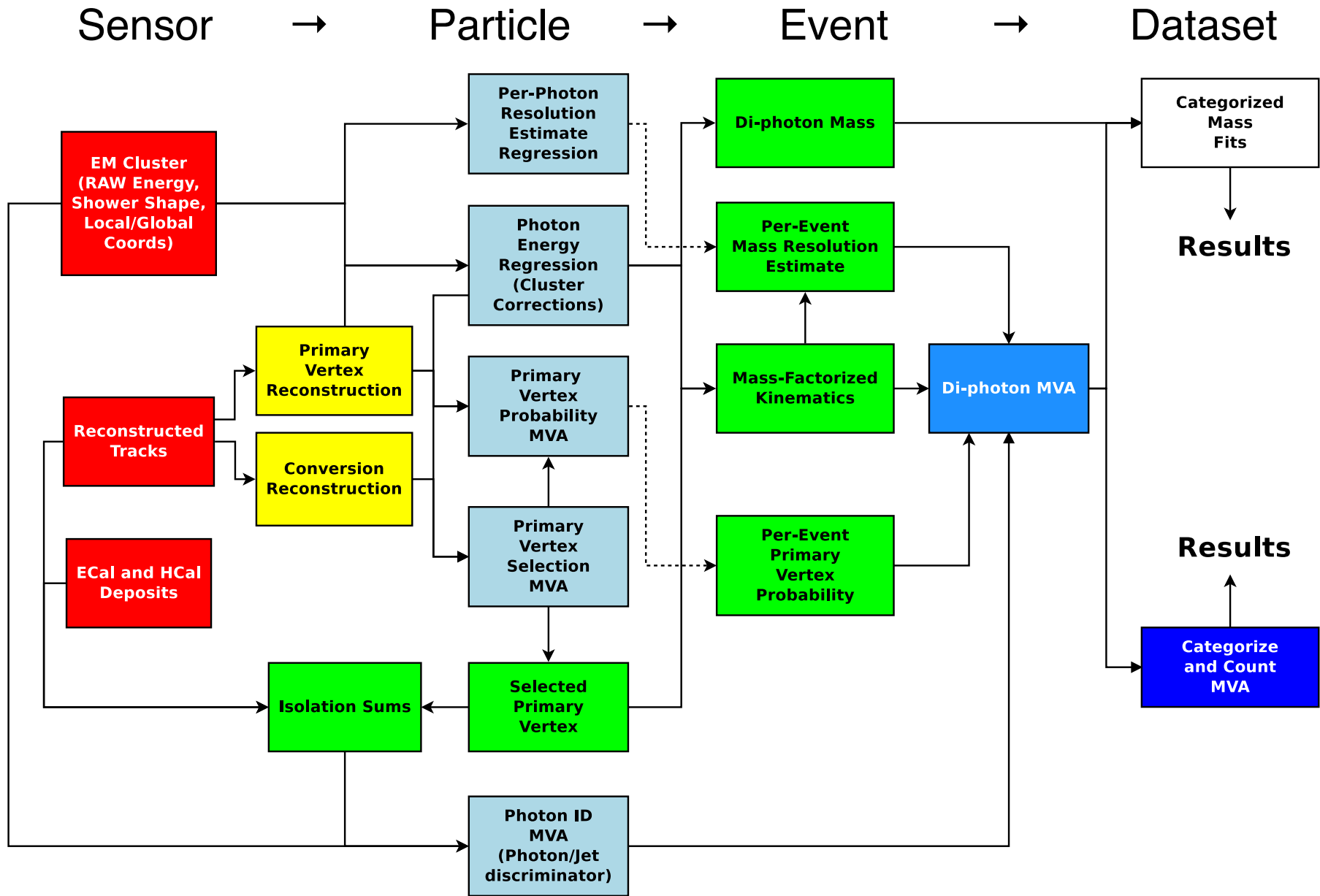


# The observation in the 2 photon channel









**\*MVA = BDT implemented in TMVA  
(Deep networks being used for particle identification)**



## Higgs Boson Machine Learning Challenge

Monday, May 12, 2014 2 months to go  
\$13,000 • 837 teams Monday, September 15, 2014

### Dashboard

- Home ↑
- Data ☰
- Make a submission ✍

### Information ⓘ

- Description
- Evaluation
- Rules
- Prizes
- About the Sponsors
- Timeline

Competition Details » [Get the Data](#) » [Make a submission](#)

## Evaluation

The evaluation metric is the *approximate median significance (AMS)*:

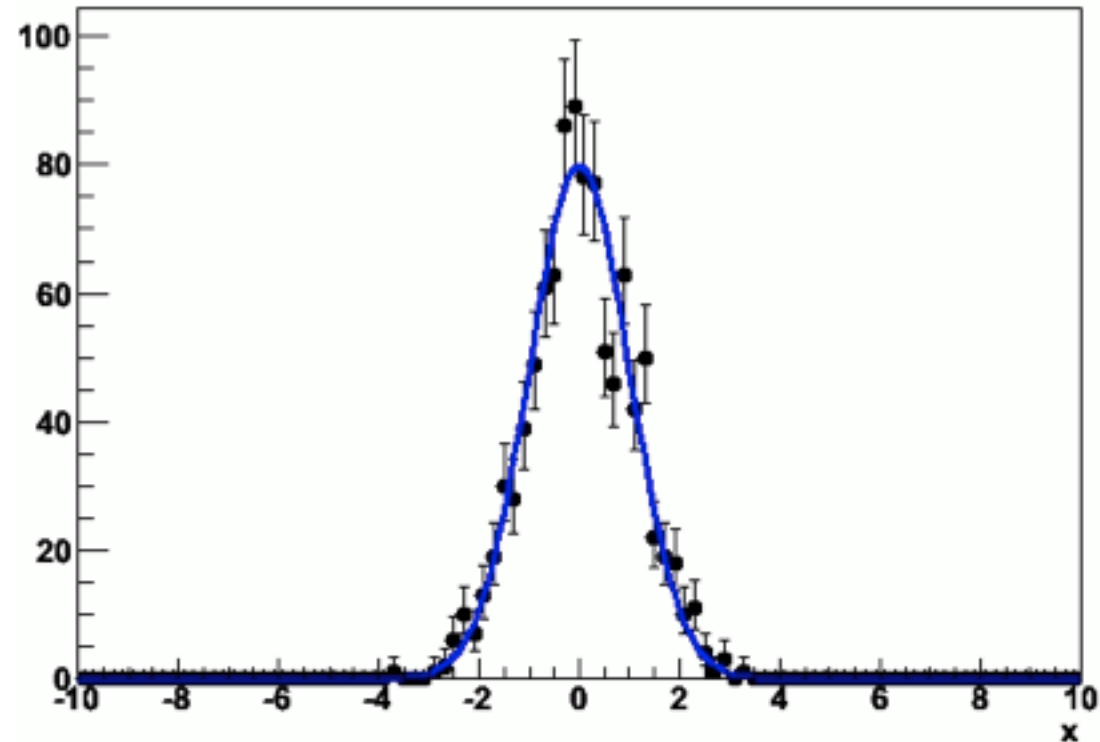
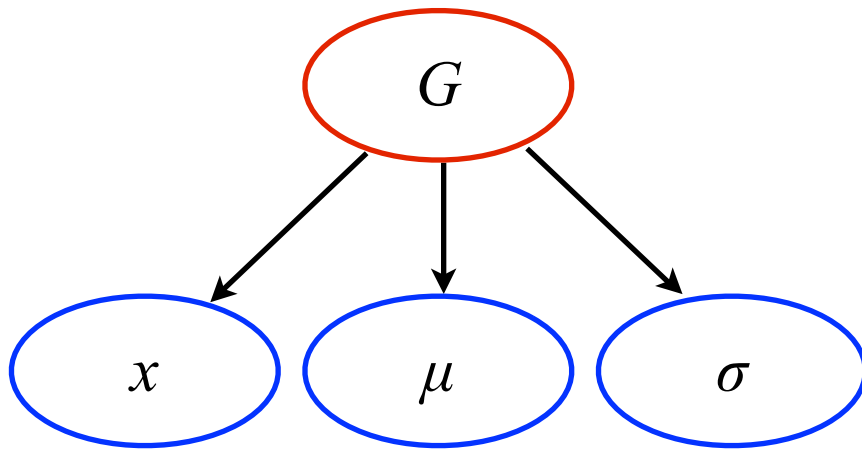
$$AMS = \sqrt{2 \left( (s + b + b_r) \log \left( 1 + \frac{s}{b + b_r} \right) - s \right)}$$

1	↑1	Gábor Melis *	3.80573	32	Thu, 26 Jun 2014 06:14:34 (-0.2h)
359	↓49	Jeje	3.25012	4	Sat, 21 Jun 2014 01:11:13
		simple TMVA boosted trees	3.24954		
360	↓49	Xiaohu SUN	3.24954	3	Tue, 03 Jun 2014 13:14:47

# Statistical Modeling for Higgs Discovery

I will represent PDFs graphically as below (directed acyclic graph)

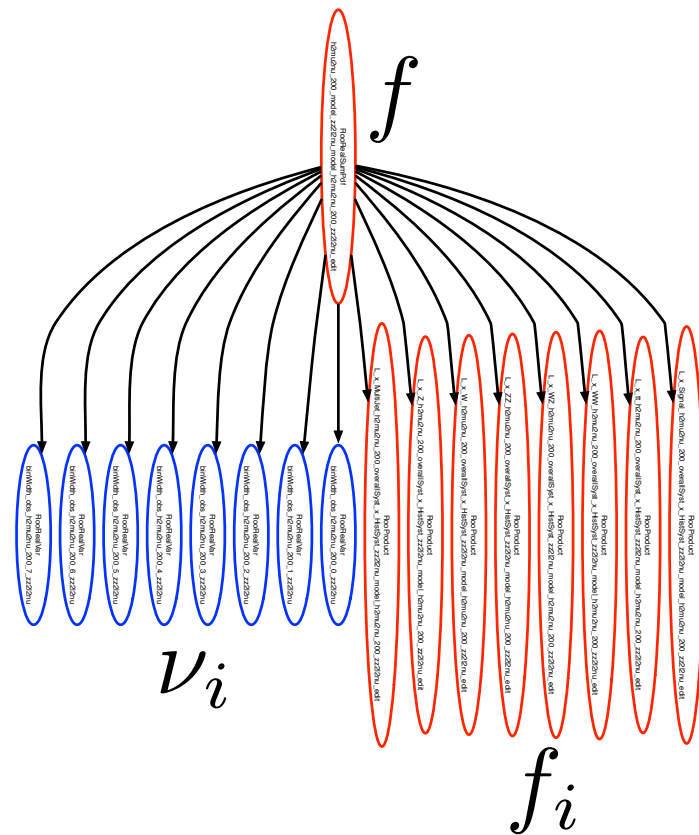
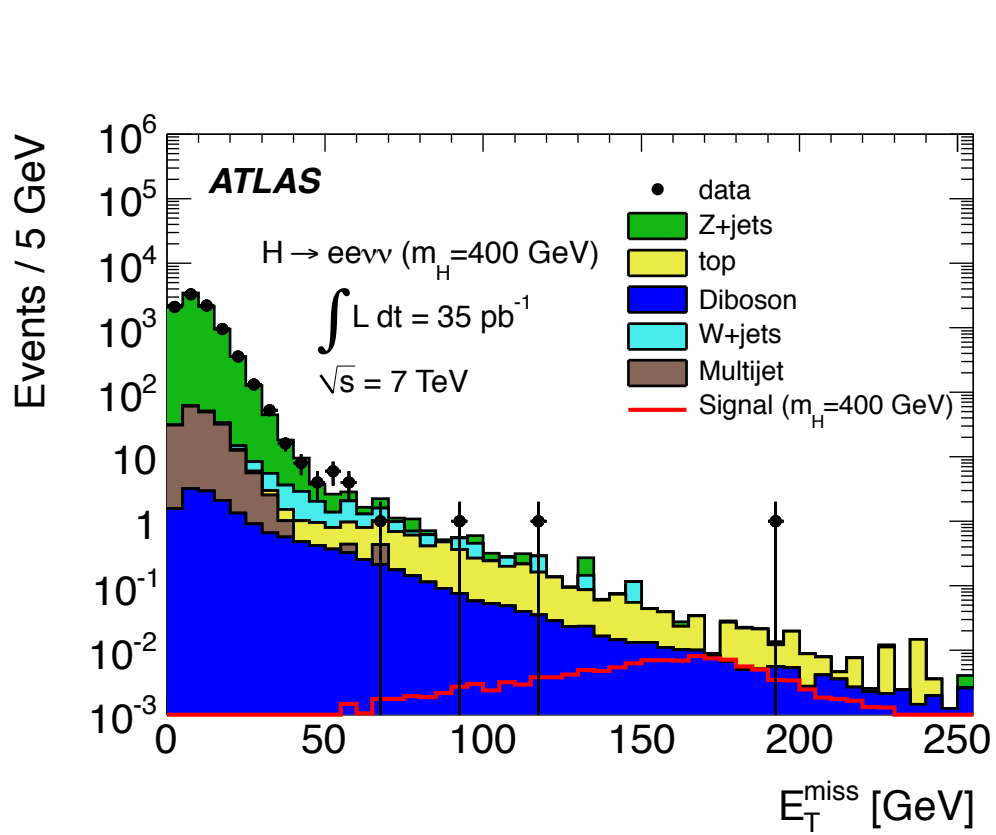
- ▶ eg. a Gaussian  $G(x|\mu, \sigma)$  is parametrized by  $(\mu, \sigma)$
- ▶ every node is a real-valued function of the nodes below



Clearly related to Graphical Models, but not the focus here.

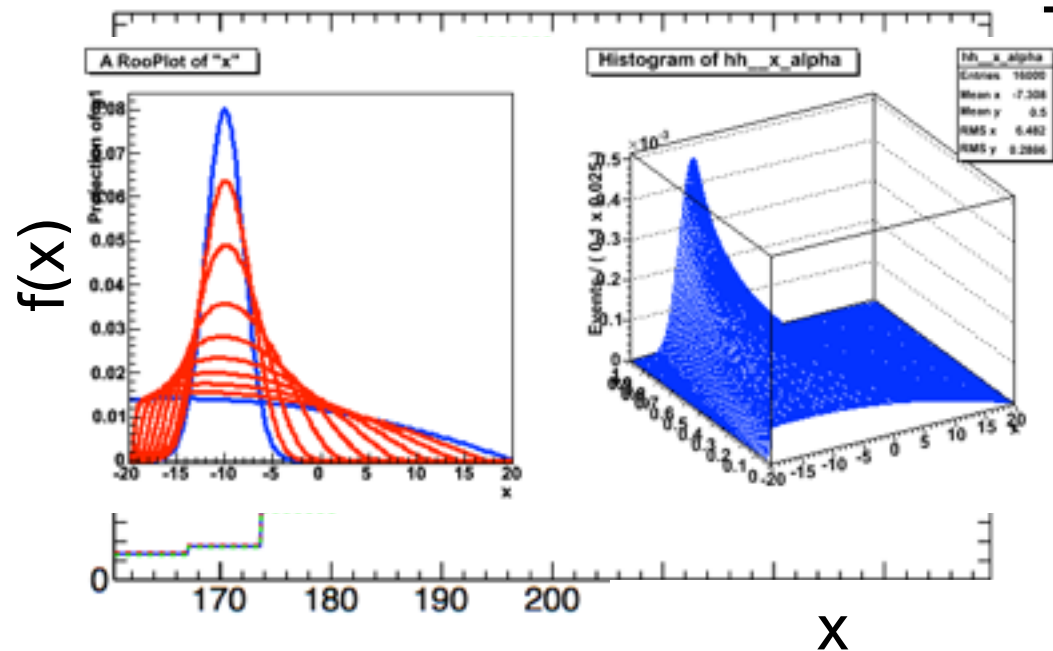
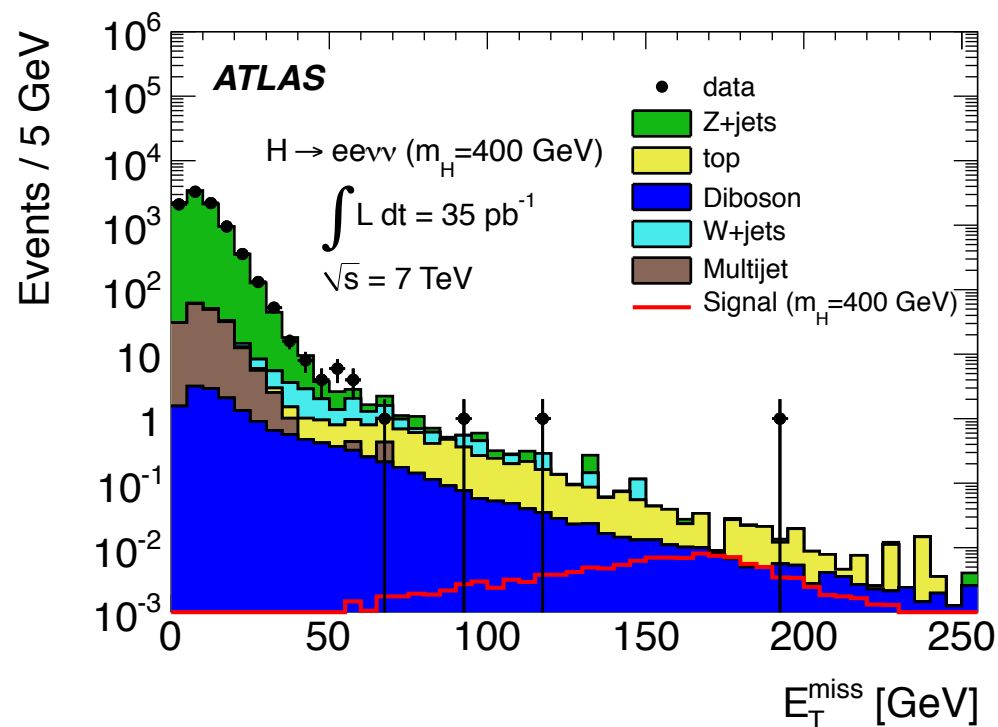
Total distribution is a mixture model with components corresponding to various signal and background interactions

$$f(x) = \frac{1}{\nu} \sum_{i \in \text{interactions}} \nu_i f_i(x), \quad \nu = \sum_{i \in \text{interactions}} \nu_i$$



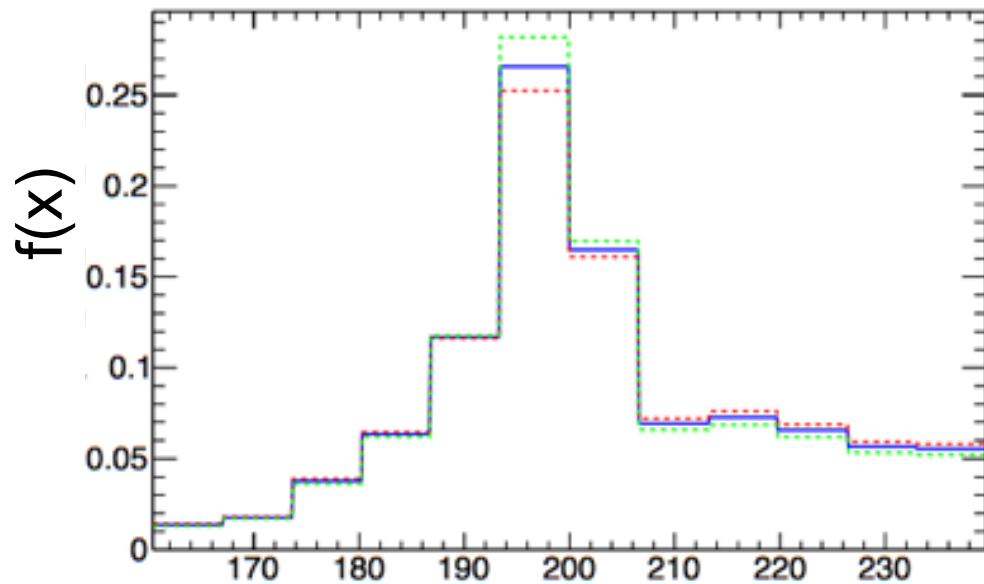
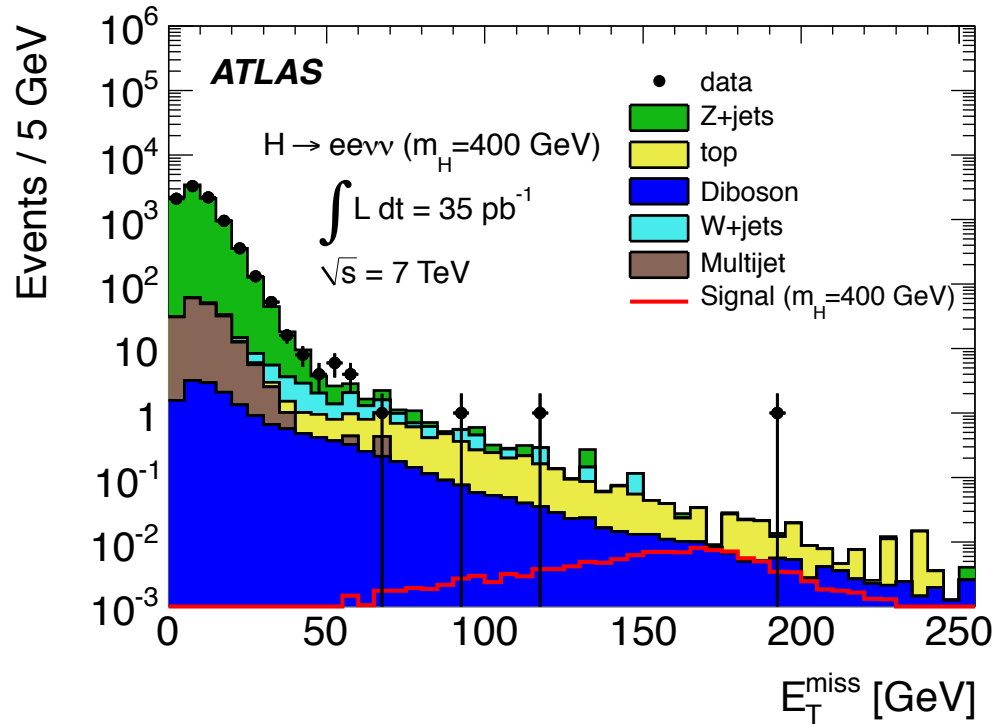
Tabulate effect of individual variations of sources of systematic uncertainty

- typically one at a time evaluated at nominal and “ $\pm 1 \sigma$ ”
- use some form of interpolation to parametrize  $p^{th}$  variation in terms of **nuisance parameter  $\alpha_p$**

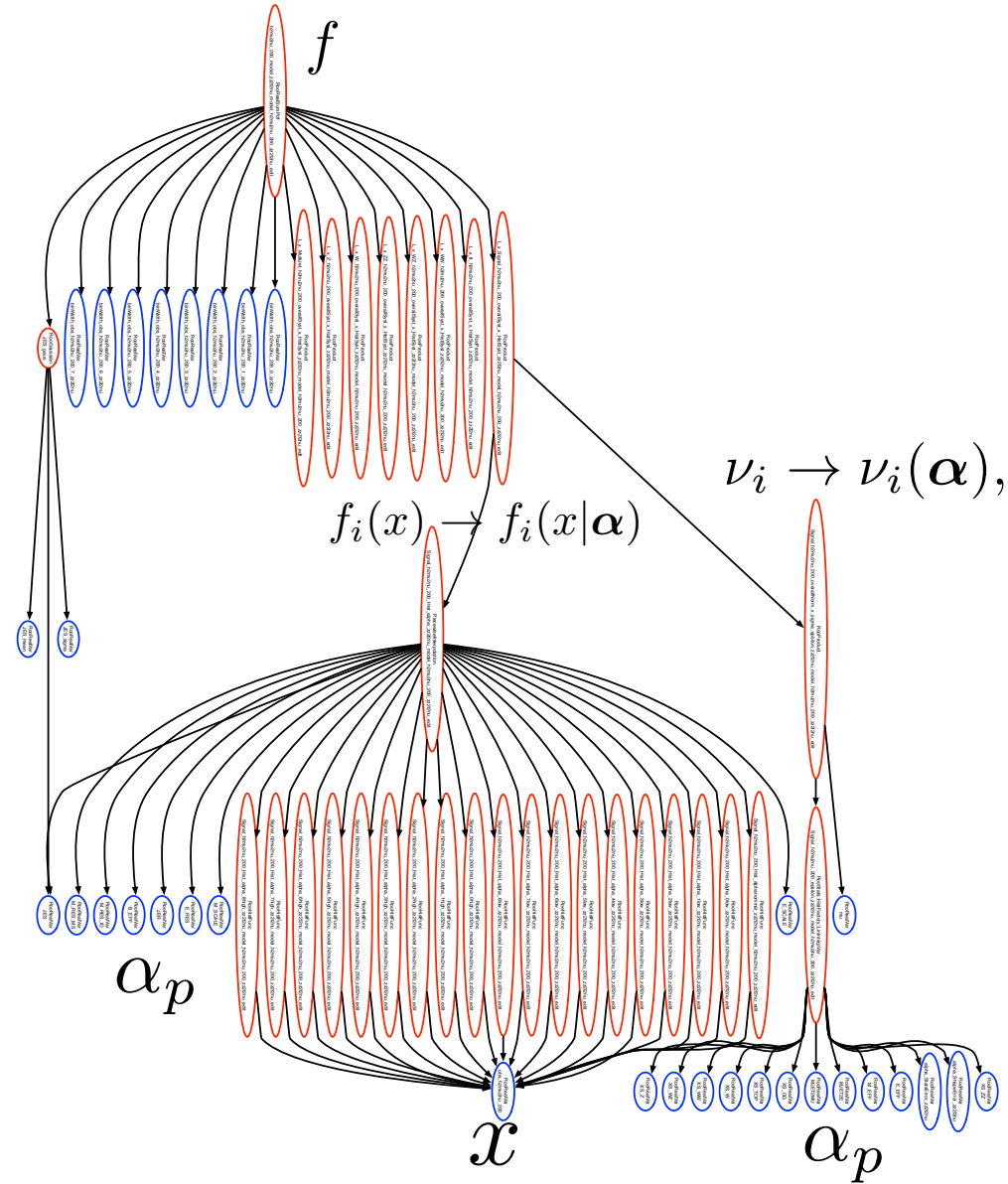


$$f(\mathcal{D}|\alpha) = \text{Pois}(n|\nu(\alpha)) \prod_{e=1}^n f(x_e|\alpha)$$

# Visualizing the model for one dataset

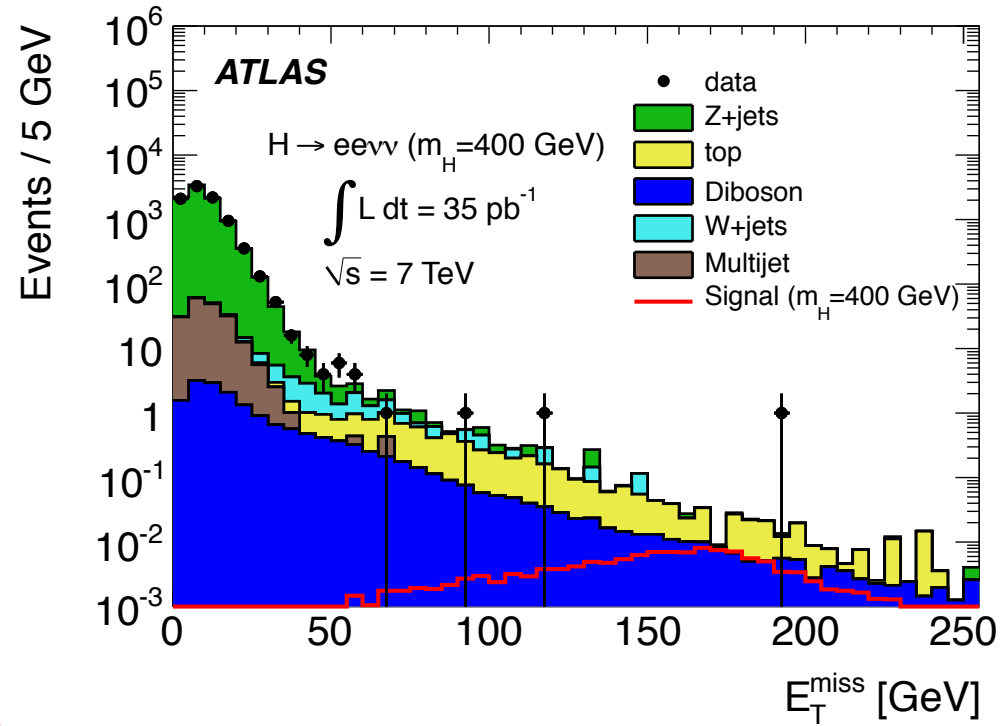
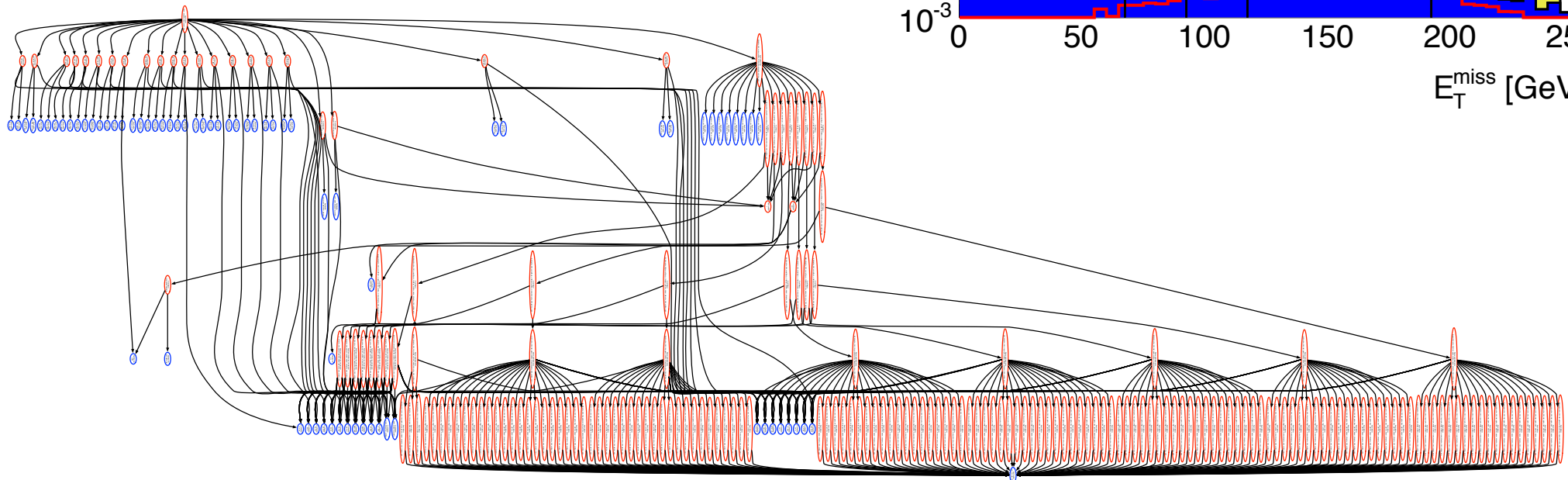


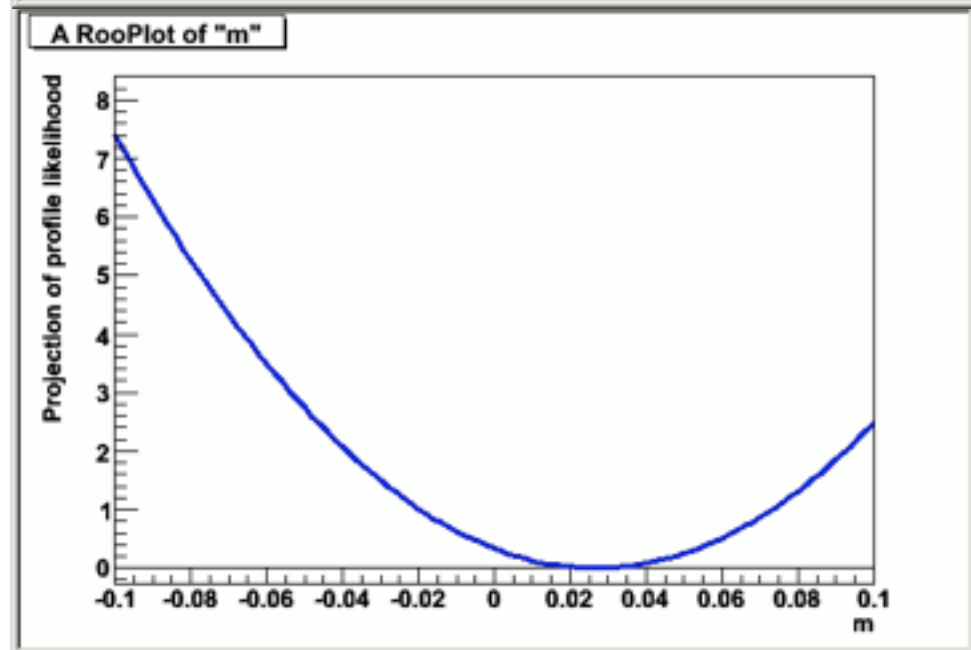
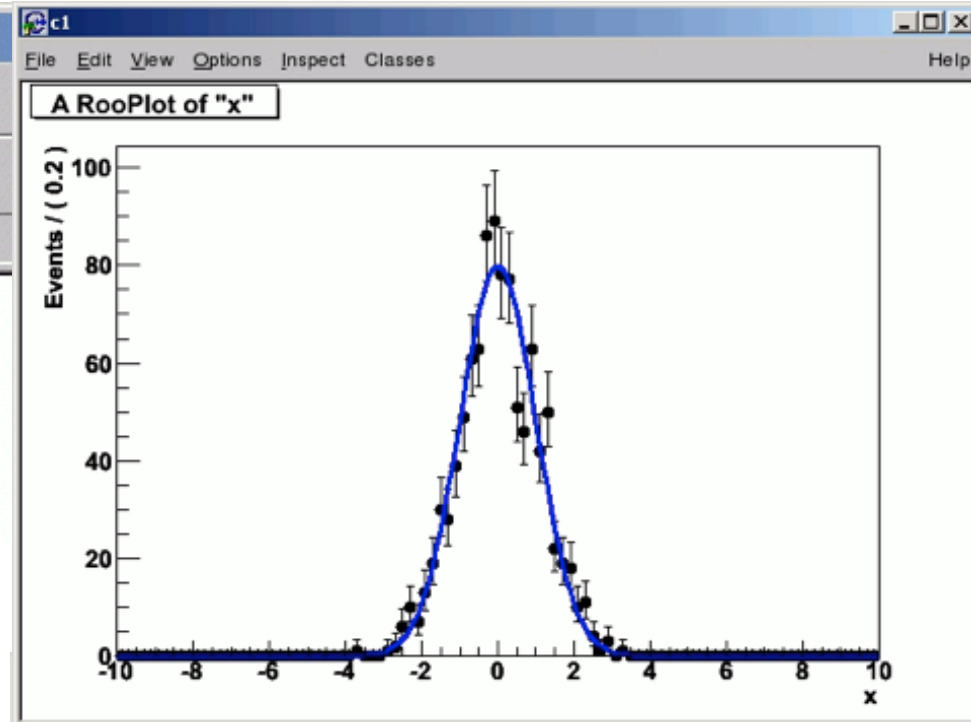
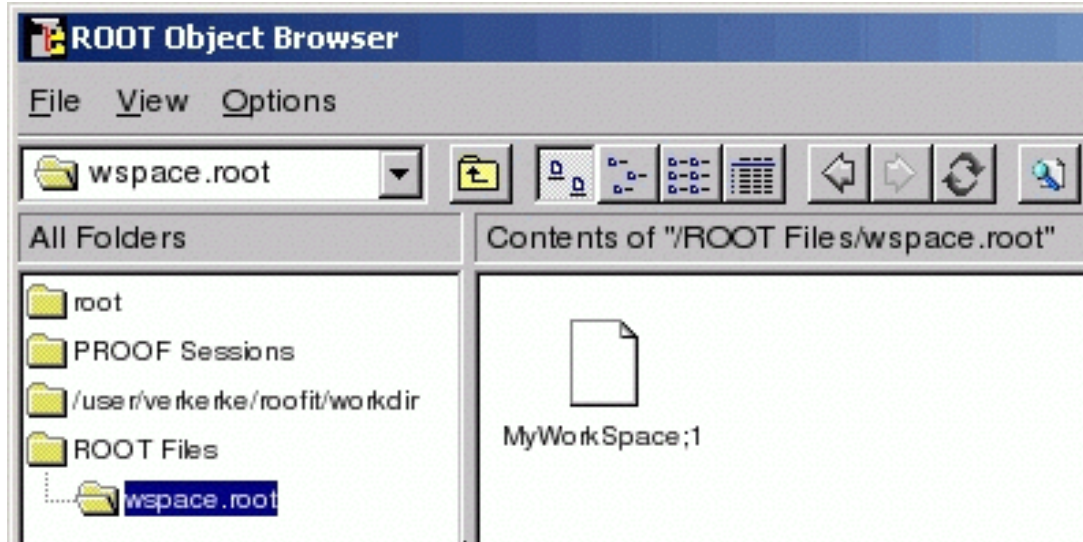
X





After parametrizing each component of the mixture model, the pdf for a single channel might look like this



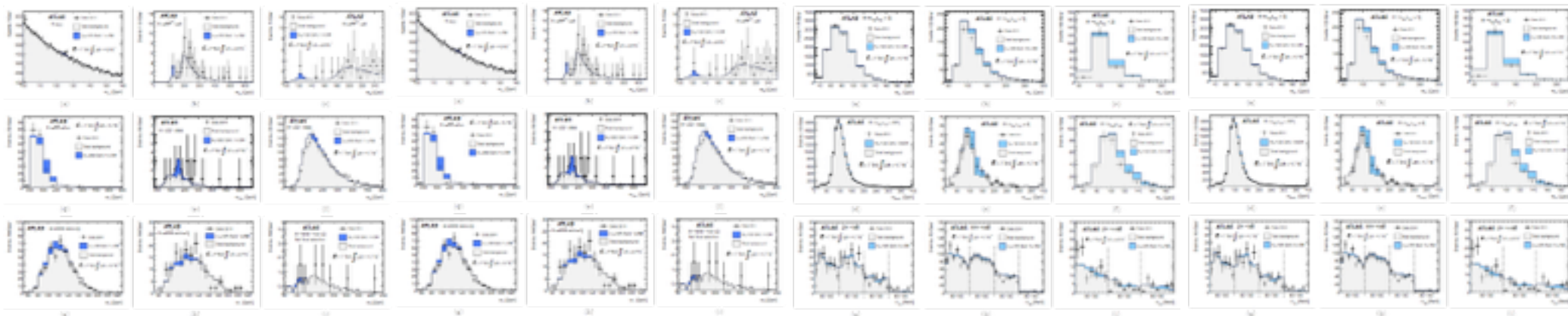


RooFit's Workspace now provides the ability to save in a file the full likelihood model, any priors you might want, and the data necessary to reproduce likelihood function.

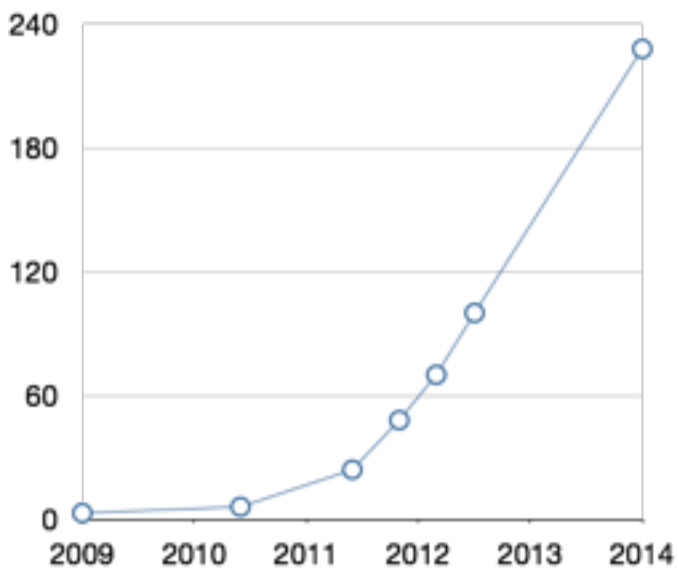
Gives flexibility in later statistical analysis (frequentist vs. bayesian) and handles for detailed meta-analysis



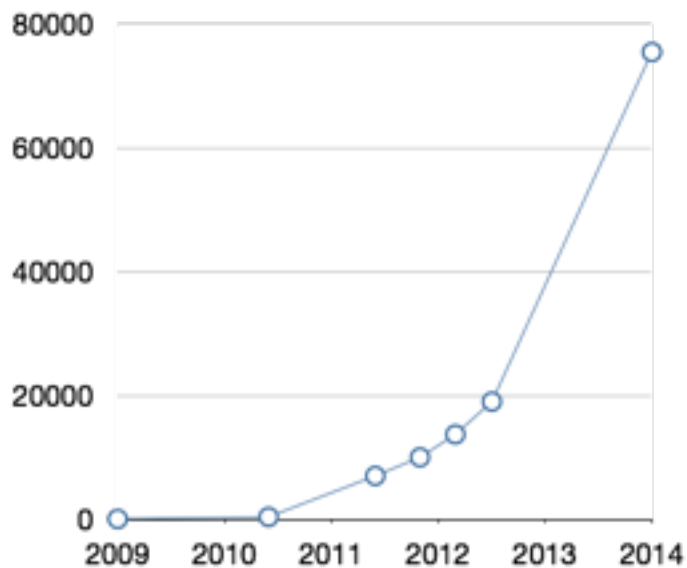
# Collaborative Statistical Modeling



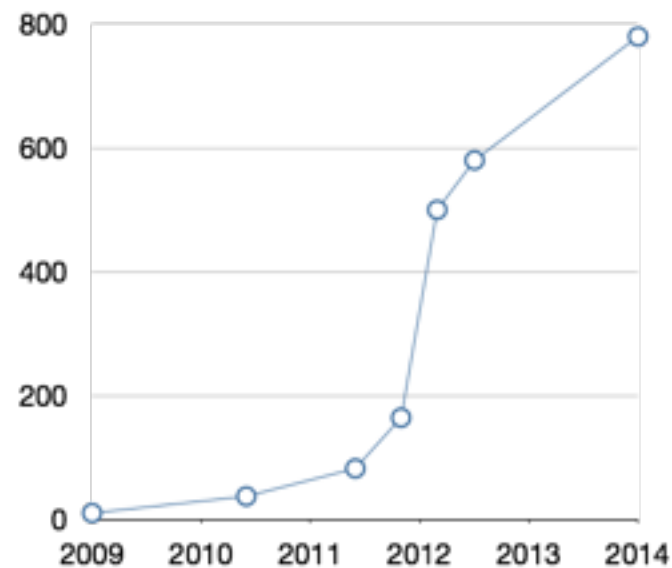
**Number of Datasets Combined**



**Number of Model Components**



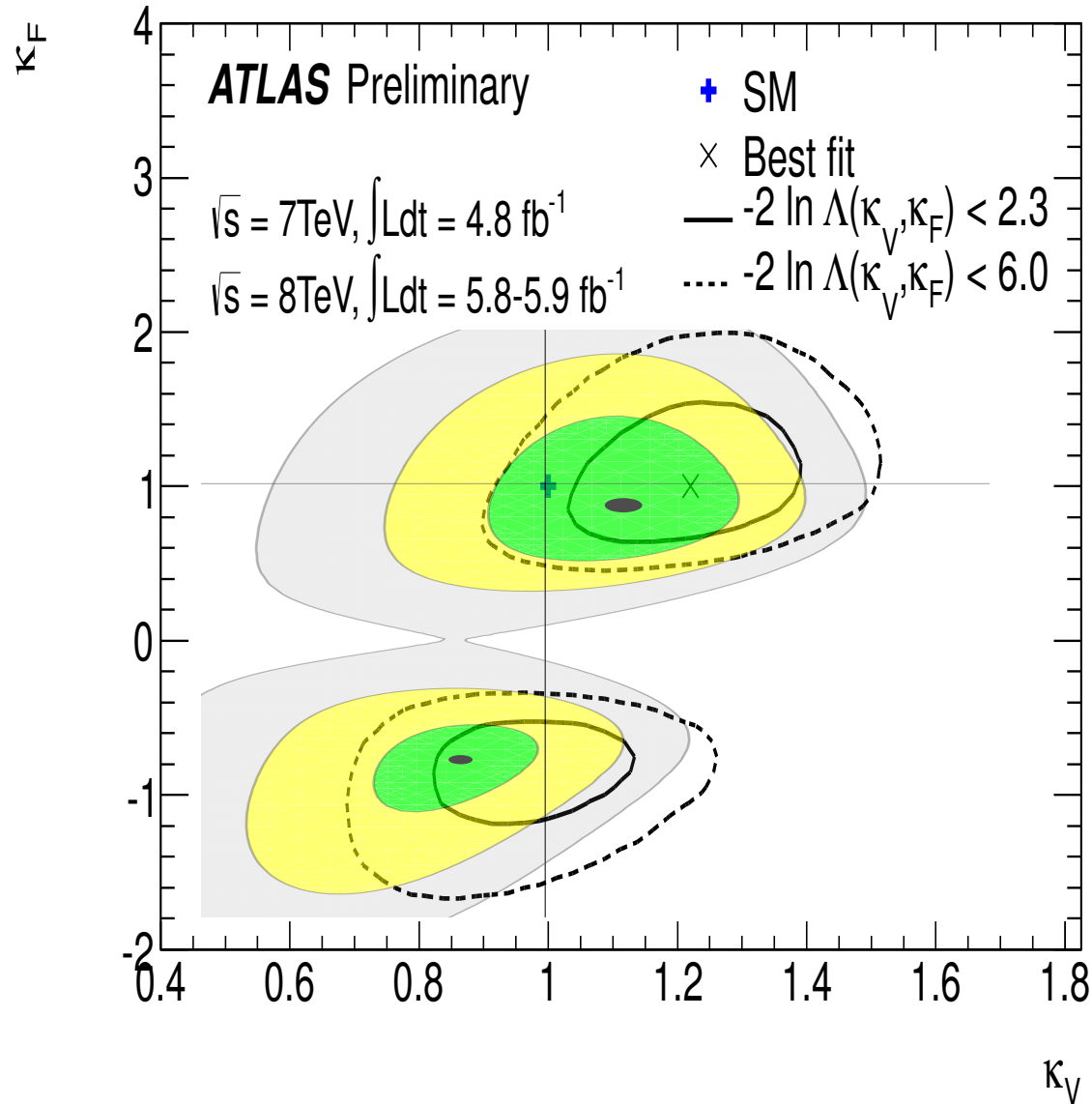
**Number of Parameters in Likelihood**





# REPRODUCIBILITY PROBLEM

*Not possible for others to reproduce results from paper.*



# PUBLISHING LIKELIHOODS

Information References (121) Citations (128) Files Plots **HepData**

## Measurements of Higgs boson production and couplings in diboson final states with the ATLAS detector at the LHC

ATLAS Collaboration (Georges Aad (Freiburg U.) et al.) [Show all 2923 authors](#)

Jul 4, 2013 - 32 pages

Phys.Lett. B726 (2013) 88-119 (2013)

DOI: [10.1016/j.physletb.2013.08.010](https://doi.org/10.1016/j.physletb.2013.08.010)

Information Citations (4) Files

## Data from Figure 7 from: Measurements of Higgs boson production and couplings in diboson final states with the ATLAS detector at the LHC

ATLAS Collaboration (Aad, Georges (Freiburg U.) [...]) [Show all 2923 authors](#)

Cite as: ATLAS Collaboration ( 2013 ) HepData, <http://doi.org/10.7484/INSPIREHEP.DATA.A78C.HK44>

Note, data record itself has 4 citation

 **figshare**   
@figshare ⚙️ Following

The Higgs Boson data is definitely the jewel in the #DataCite crown. Hopefully the first of many!

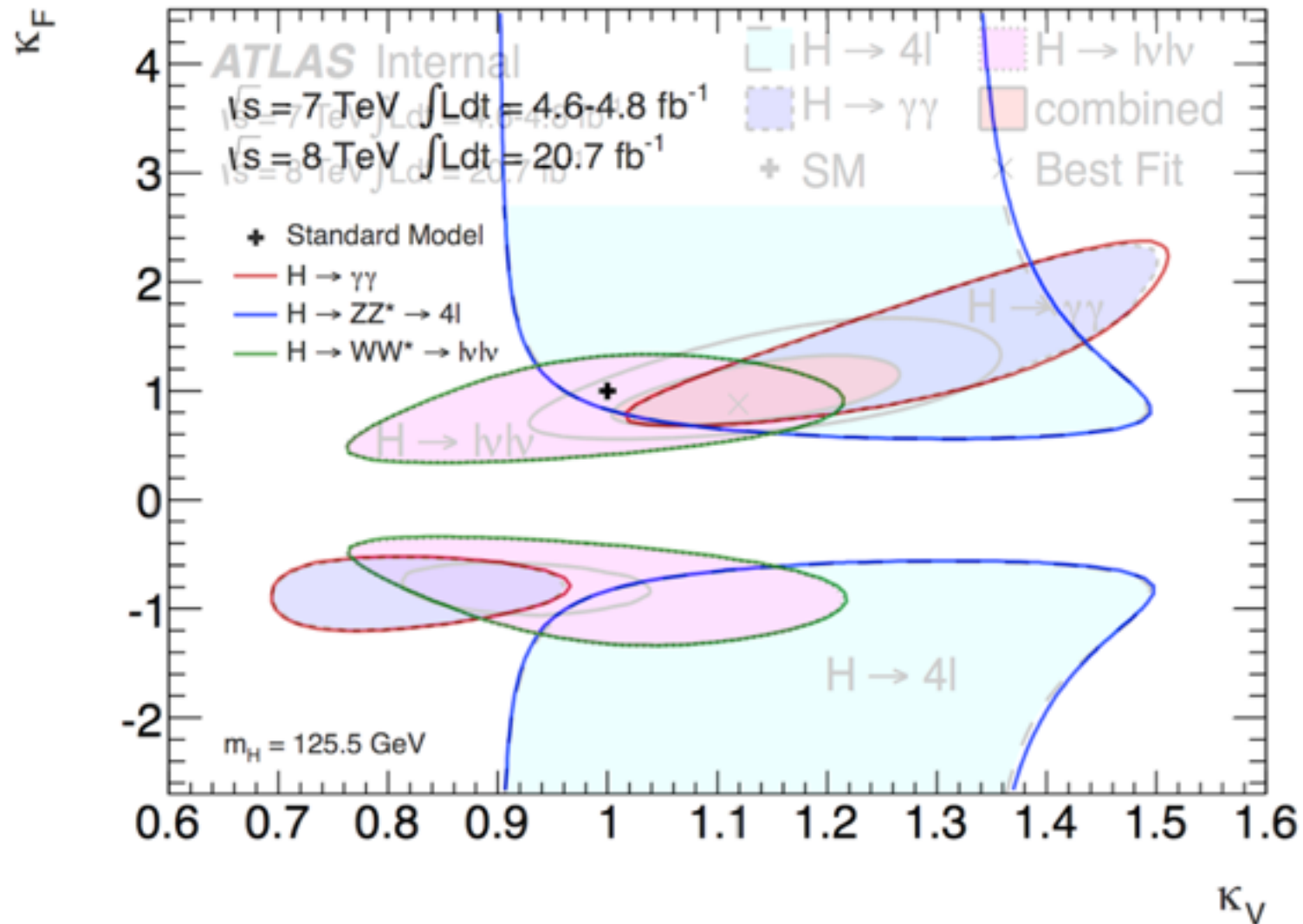
 Reply  Retweeted  Favorite  More

RETWEETS 7 FAVORITES 2

# PUBLISHING LIKELIHOODS

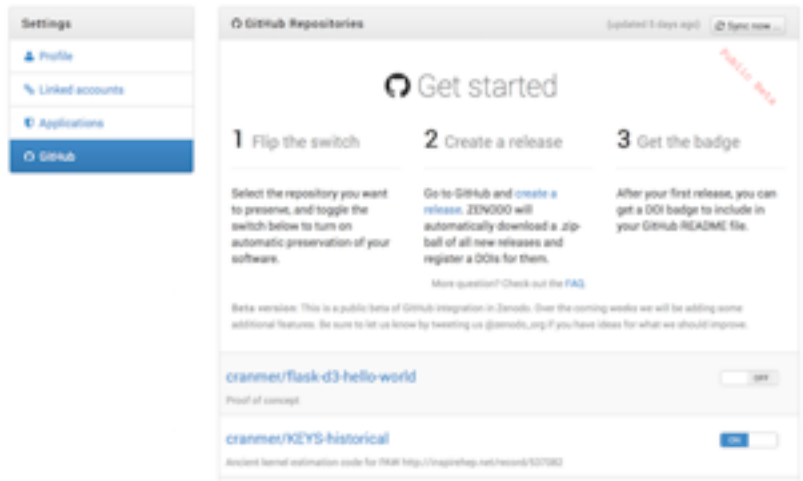
*Reproducing derived results from original paper!*





# CODE AS A RESEARCH PRODUCT

GitHub → Zenodo → INSPIRE



Mathematica → figshare → INSPIRE

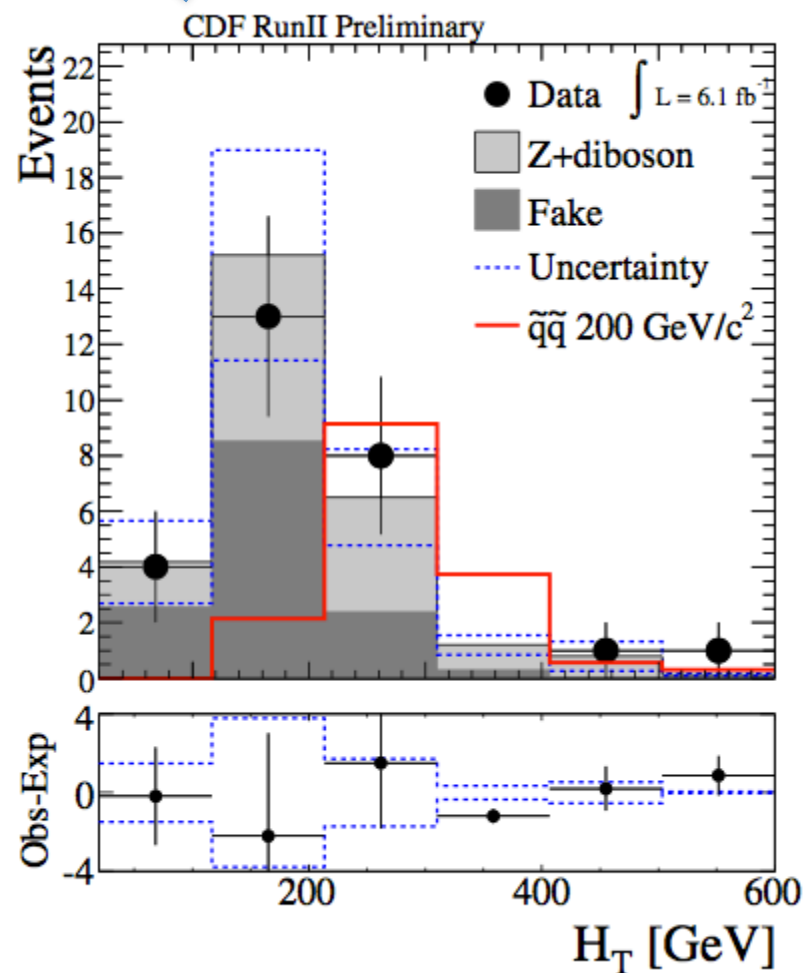
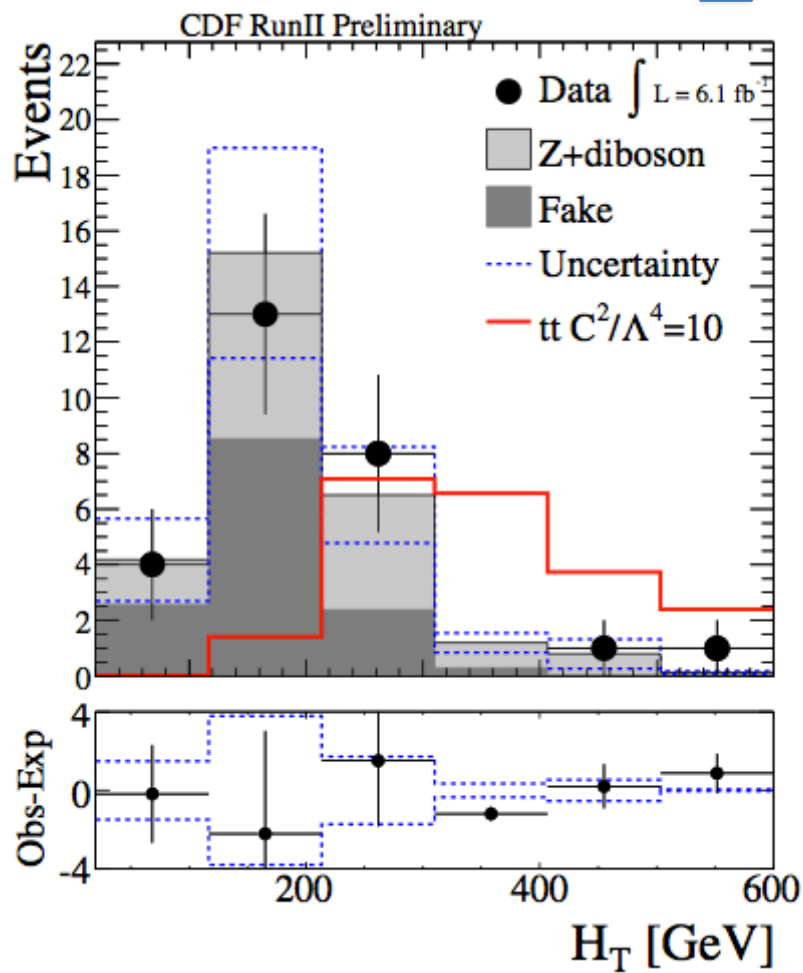
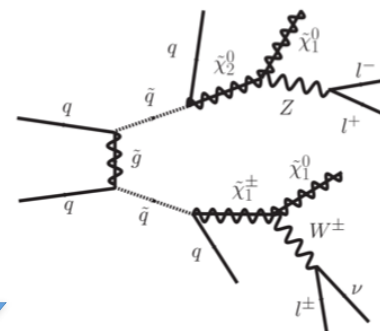
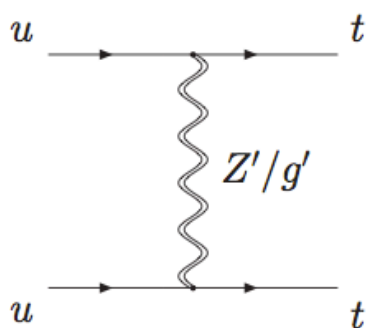
## Supplementary Material for "A Novel Approach to Higgs Coupling Measurements"

(2013) *figshare*.



- highly discussed by scholars 97 - 100 percentile  of datasets published in 2013  3 figshare shares 
- highly viewed by scholars 97 - 100 percentile  of datasets published in 2013  202 figshare views 
- highly viewed by scholars 97 - 100 percentile  of datasets published in 2013  9 figshare downloads 
- highly discussed by public 97 - 100 percentile  of datasets published in 2013  10 tweets 

# RECASTING



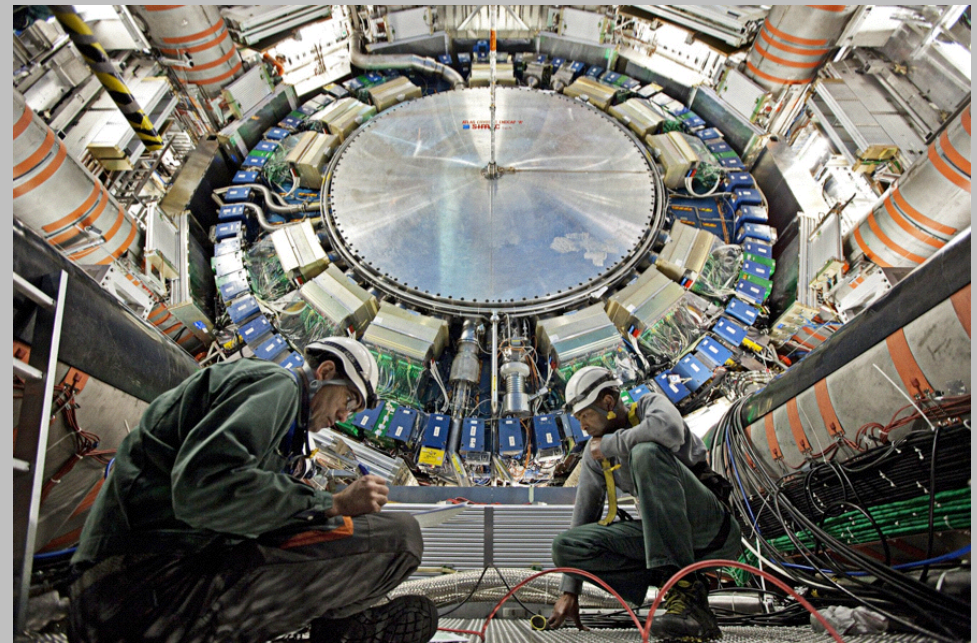
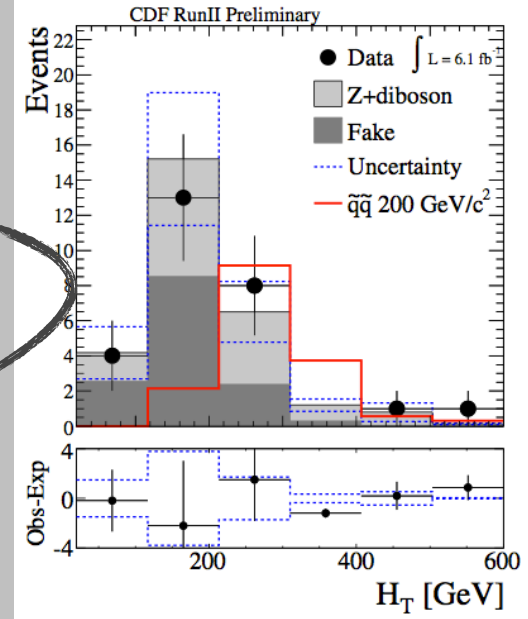
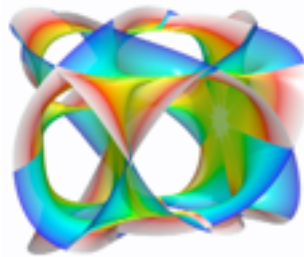
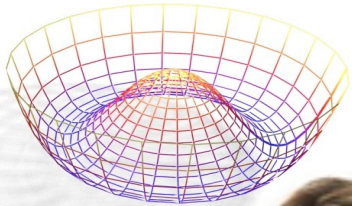
# THEORY

# SERVICE

$$\begin{aligned}
 \mathcal{L}_{SM} = & \underbrace{\frac{1}{4} \mathbf{W}_{\mu\nu} \cdot \mathbf{W}^{\mu\nu} - \frac{1}{4} B_{\mu\nu} B^{\mu\nu} - \frac{1}{4} G_{\mu\nu}^a G_a^{\mu\nu}}_{\text{kinetic energies and self-interactions of the gauge bosons}} \\
 & + \underbrace{\bar{L} \gamma^\mu (i \partial_\mu - \frac{1}{2} g \boldsymbol{\tau} \cdot \mathbf{W}_\mu - \frac{1}{2} g' Y B_\mu) L + \bar{R} \gamma^\mu (i \partial_\mu - \frac{1}{2} g' Y B_\mu) R}_{\text{kinetic energies and electroweak interactions of fermions}} \\
 & + \underbrace{\frac{1}{2} \left| (i \partial_\mu - \frac{1}{2} g \boldsymbol{\tau} \cdot \mathbf{W}_\mu - \frac{1}{2} g' Y B_\mu) \phi \right|^2 - V(\phi)}_{W^\pm, Z, \gamma, \text{ and Higgs masses and couplings}} \\
 & + \underbrace{g'' (\bar{q} \gamma^\mu T_a q) G_\mu^a}_{\text{interactions between quarks and gluons}} + \underbrace{(G_1 \bar{L} \phi R + G_2 \bar{L} \phi_c R + h.c.)}_{\text{fermion masses and couplings to Higgs}}
 \end{aligned}$$

Q

A



# Review of Challenges and Possible Research Topics

The complexity of our statistical models is growing exponentially, starting to need new algorithms to deal with them or principles for simplifying them

- graphical models, automatic differentiation, distributed processing, ...
- better optimization & sampling algorithms
- optimal statistical procedures subject to computational constraints ([link](#))

Interpolation of distributions based on simulated samples with different parameter settings a weak point

- experimental design, response surface interpolation, Gaussian processes, ...
- complication: samples often not statistically independent

Machine learning + computer simulations

- Most analyses either use computer simulations of the detector or ad hoc parametrized models.
- Little use of machine learning to learn the expensive computer simulation

Most discussion with statisticians has focused on hypothesis testing and confidence intervals for final results. Many interesting problems up-stream

- **exception:** machine learning for selecting candidate signal events
- **barriers:** collaborations do not openly share data, requires some semi-formal agreement
- **progress:** movement towards open access (link to policy)

## Importance sampling for rare events in simulation

- The simulation of our detectors is very computationally challenging and we use brute force to populate tails in cases where we can do something smarter

## Particle physics is a unique arena for data science

- well posed questions in an extreme setting
- lots of data, complicated sensor environment, strong theoretical basis

Congratulations and best wishes to

