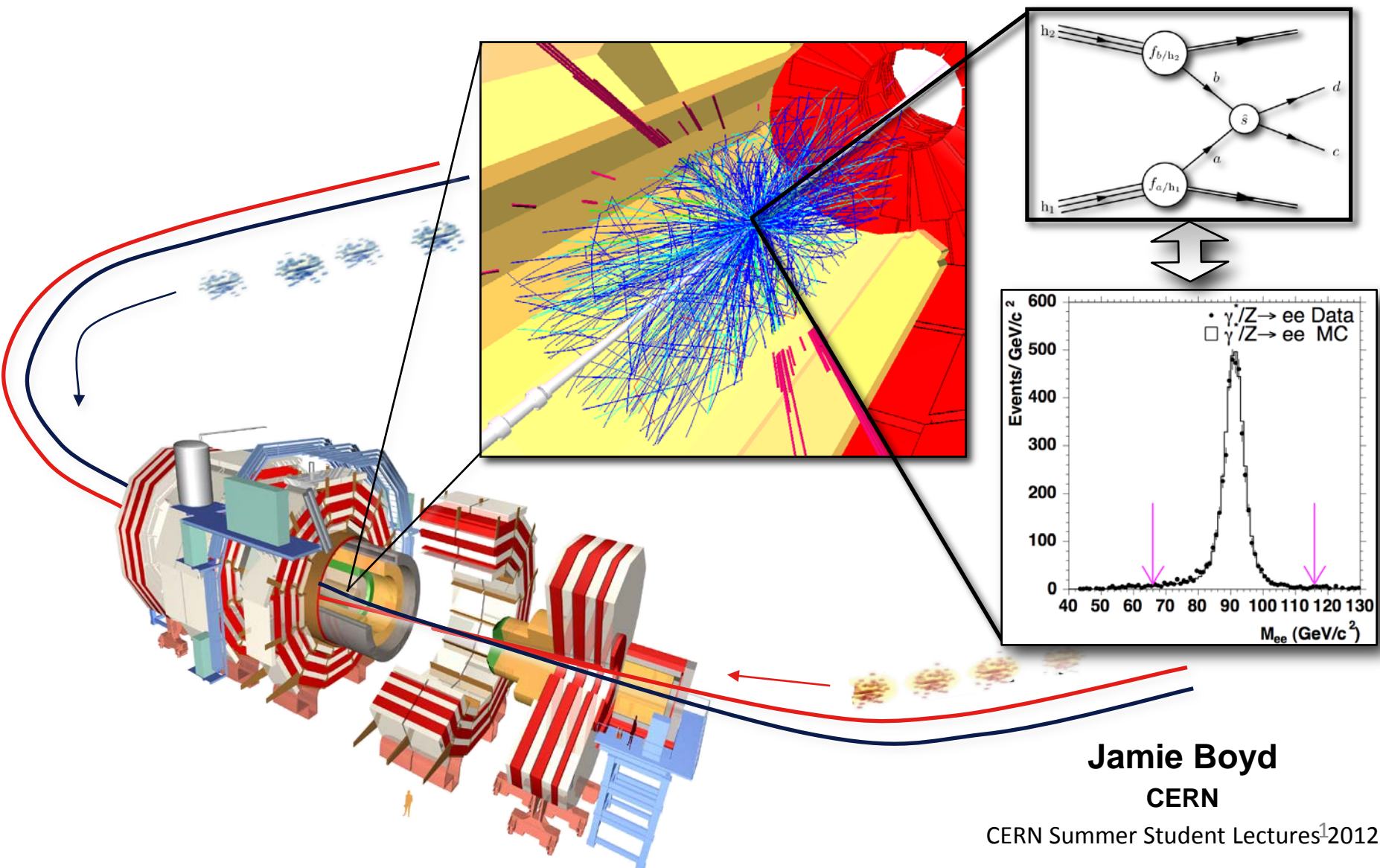


From Raw Data to Physics Results



Jamie Boyd
CERN

Outline

- Summary
 - Brief overview of the full lecture course
- A simple example
 - Measuring the Z^0 cross-section
- Reconstruction & Simulation
 - Track reconstruction
 - Calorimeter reconstruction
 - Physics object reconstruction
 - Simulation
- Physics Analysis
 - Data Quality
 - Jet cross-section
 - $Z' \rightarrow ll$
 - $H \rightarrow \gamma\gamma$
 - $H \rightarrow ZZ \rightarrow 4l$
- Computing infrastructure

Today's Lecture

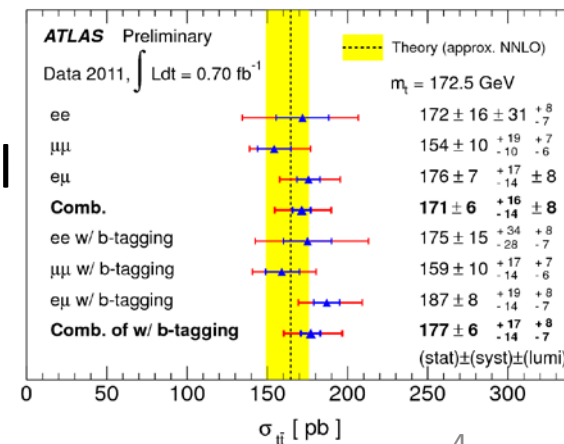
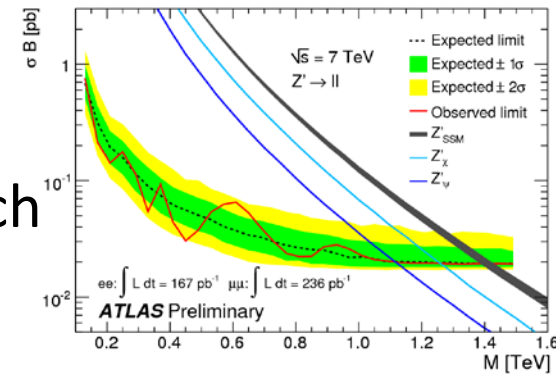
Disclaimer : Much of the content based on previous years lectures
Thanks to G. Dissertori



Section 4: Physics Analysis

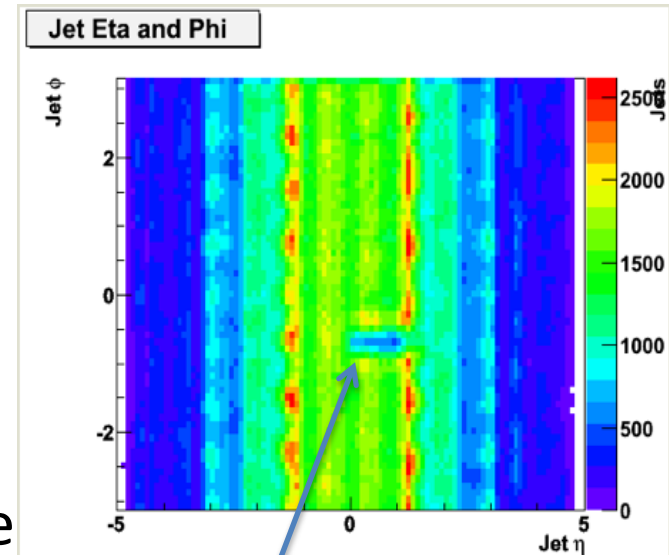
Physics Analysis

- Two main types of physics analysis at LHC
 - Searching for new particles
 - Making precision measurements
- Searches statistically limited
 - More data is the way of improving the search
 - If don't see anything new set limits on what you have excluded
- Precision measurements
 - Precision often limited by the systematic uncertainties
 - Precision measurements of Standard Model parameters allows important tests of the consistency of the theory



Data Quality

- Particle physics detectors are very complex instruments – often there can be small problems
 - Noise in the detectors
 - Regions of the detector that are not working
- These effect the physics object reconstruction and the physics analysis
- Can exclude events with problems from the analysis
- Or include them, but take into account the problem (model them in simulation if possible)
- Need a thorough system to check that the data is of good quality for physics analysis
 - Lots of people, checking lots of histograms...
- Data Quality very important to get correct physics results



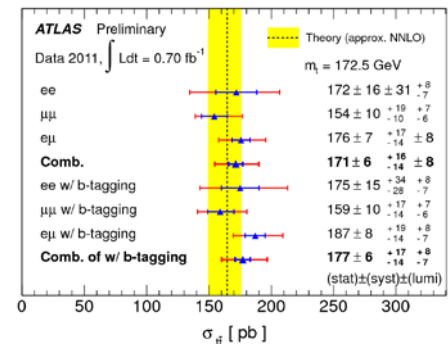
A dead region in the calorimeter in ATLAS



~TB

Physics Analysis Steps

- Start with the output of reconstruction
- Apply an event selection based on the reconstructed object quantities
 - Often calculate new information e.g masses of combinations of particles
 - Event selection designed to improve the ‘signal’ to ‘background’ in your event sample
- Estimate
 - Efficiency of selection (& uncertainty)
 - Background after selection (& uncertainty)
 - Can use simulation for these – but have to use data-driven techniques to understand the uncertainties
- Make final plot
 - Comparing data to theory
 - Correcting for efficiency and background in data
 - Include the statistical and systematic uncertainties





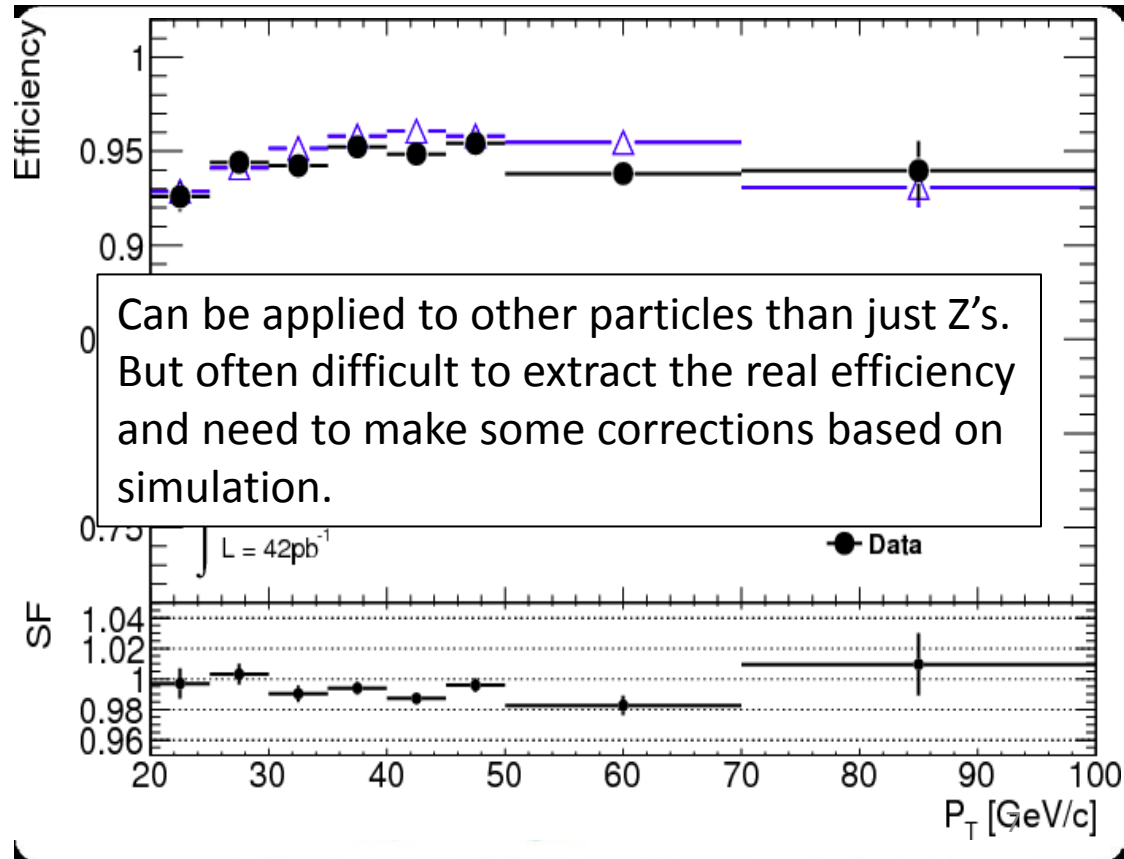
Measuring efficiencies from the data

- Can use simulation to get reconstruction efficiency
- But difficult to know how well the simulation describes the detector
- Need to try to use real data to estimate the efficiency too
- Example: Tag & Probe

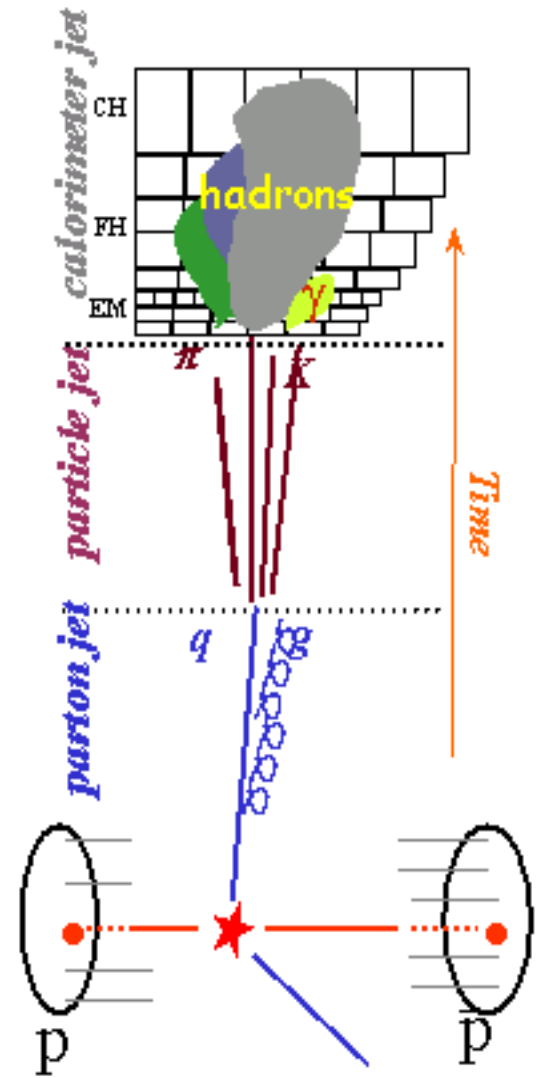
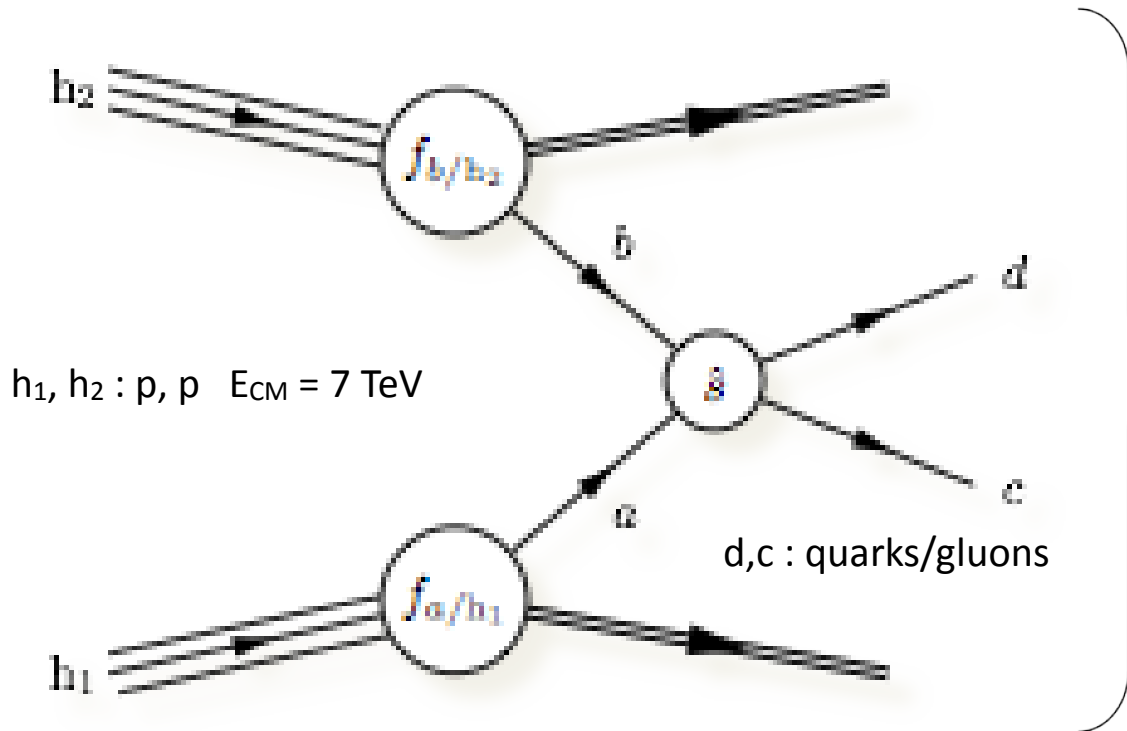
idea: use $Z \rightarrow \mu\mu$ decays to give us a pure sample of muon to use to measure the efficiency.

recipe:

1. select events with 1 reconstructed muon and 1 high momentum track
2. require that the mass of the track and the muon is consistent with Z mass
3. test if track is also reconstructed as a muon



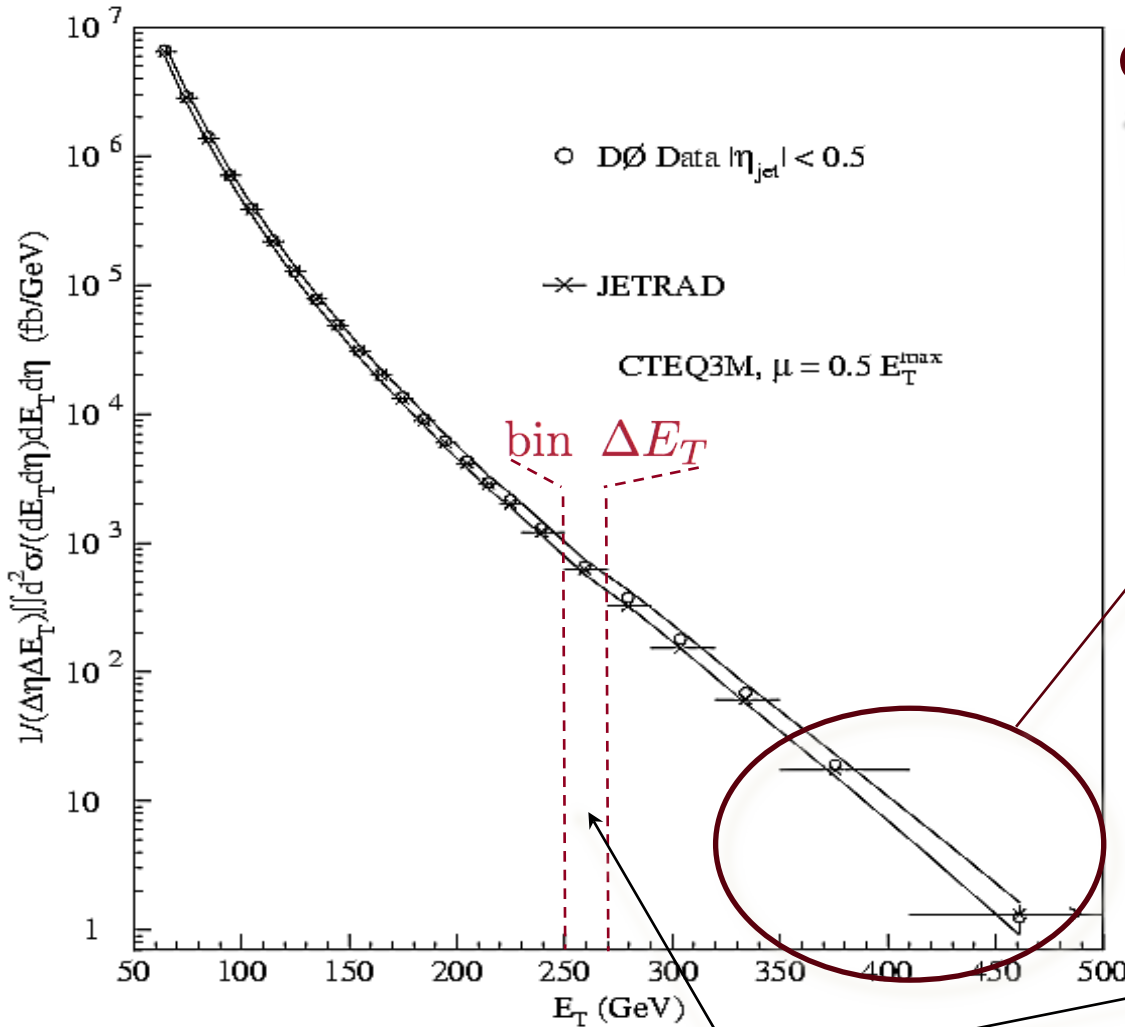
JET production at hadron colliders



Goal

- measure probability that **quarks/gluons are produced with a certain energy, at a certain angle**
- **Problem** : do not observe quarks and gluons directly, only hadrons, which appear collimated into **jets**
- Reconstruct from energy clusters in the calorimeter
 - Unfortunately don't have time to go into details of jet reconstruction

What do we have to measure?



Goal

- measure **cross section** (probability) that **jets** are produced with a certain **transverse energy** E_T , within a certain **rapidity range**
- **Test of perturbative QCD**, over many orders of magnitude!
- Look at **very high energy tail**, **new physics** could show up there in form of excess (eg. sub-structure of quarks?)

can be calculated in pert. QCD

$$\left\langle \frac{d^2\sigma}{dE_T d\eta} \right\rangle = \frac{N}{\Delta E_T \Delta \eta \epsilon \mathcal{L}}$$

• **count** number of events, N , in this bin

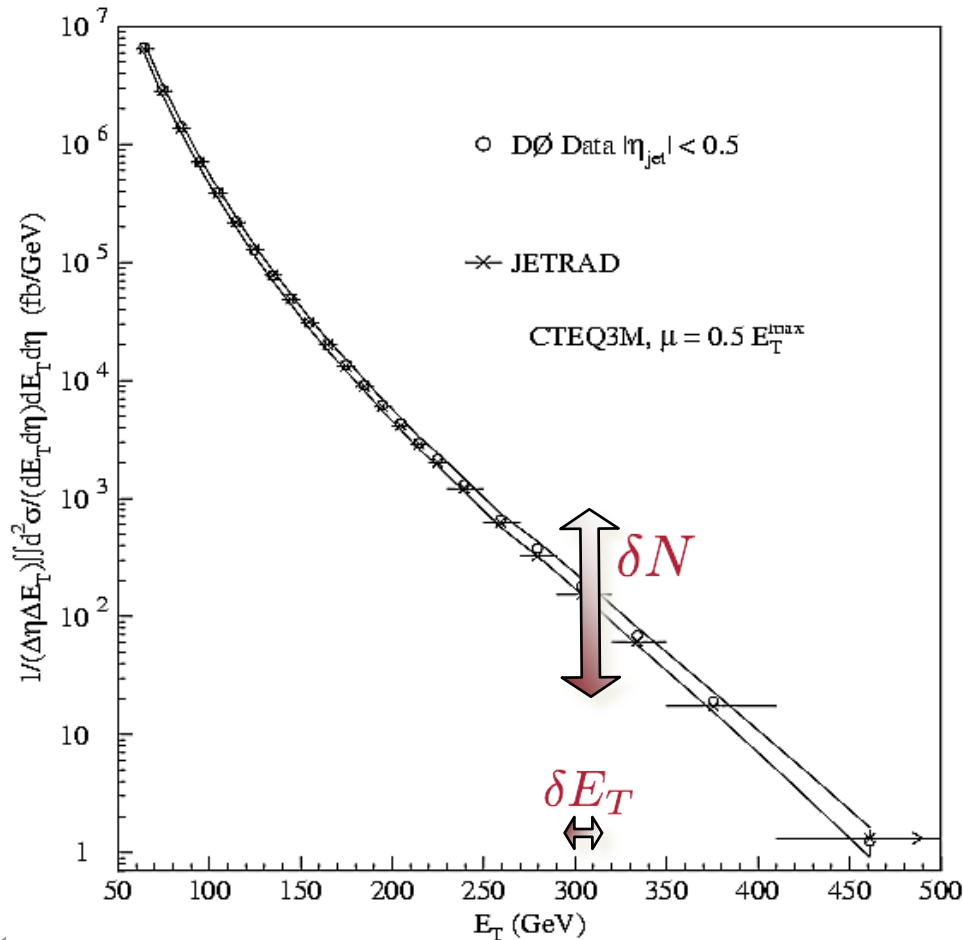
• for a certain range in **rapidity** (angle) $\Delta\eta$

efficiency to reconstruct jets

integrated accelerator **luminosity**

Problem : Energy scale

- Question : how well do we know the **energy calibration**?
- Critical because of very steeply falling spectrum!



$$\frac{d^2\sigma}{dE_T d\eta} \approx \text{const} \cdot E_T^{-6}$$



relative uncertainties

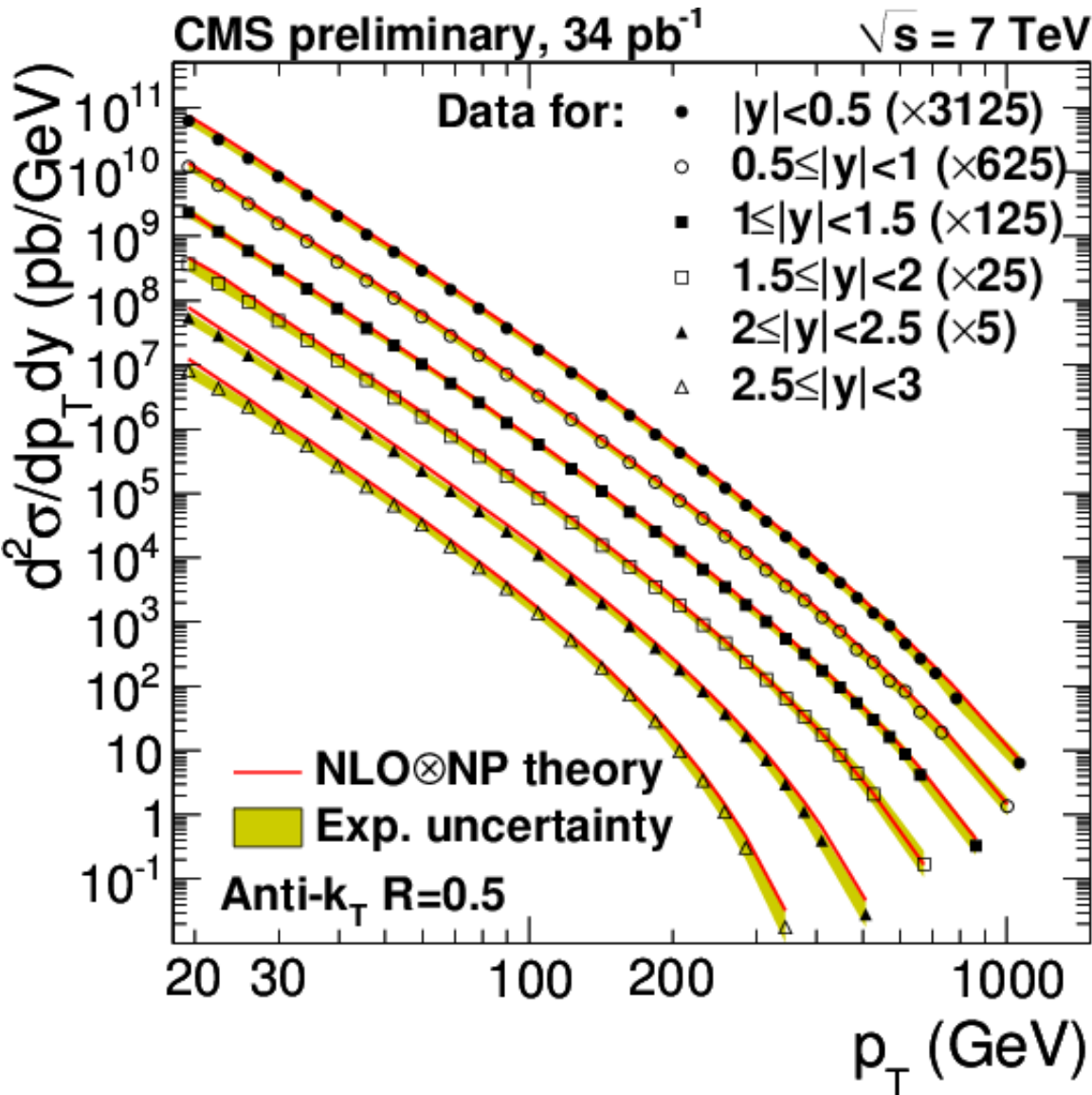
$$\frac{\delta N}{N} \approx 6 \cdot \frac{\delta E_T}{E_T}$$

so beware:

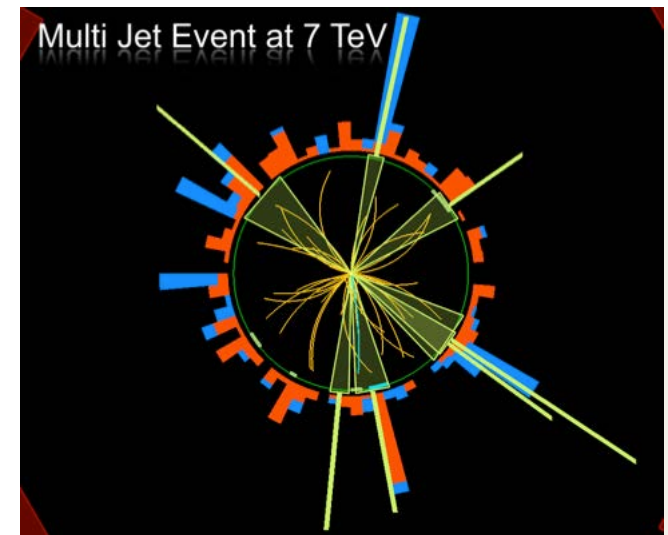
eg. an uncertainty of **5%** on absolute energy scale (calibration)

- ⦿ an uncertainty of **30%** (!) on the measured cross section

Jet cross-section at the LHC



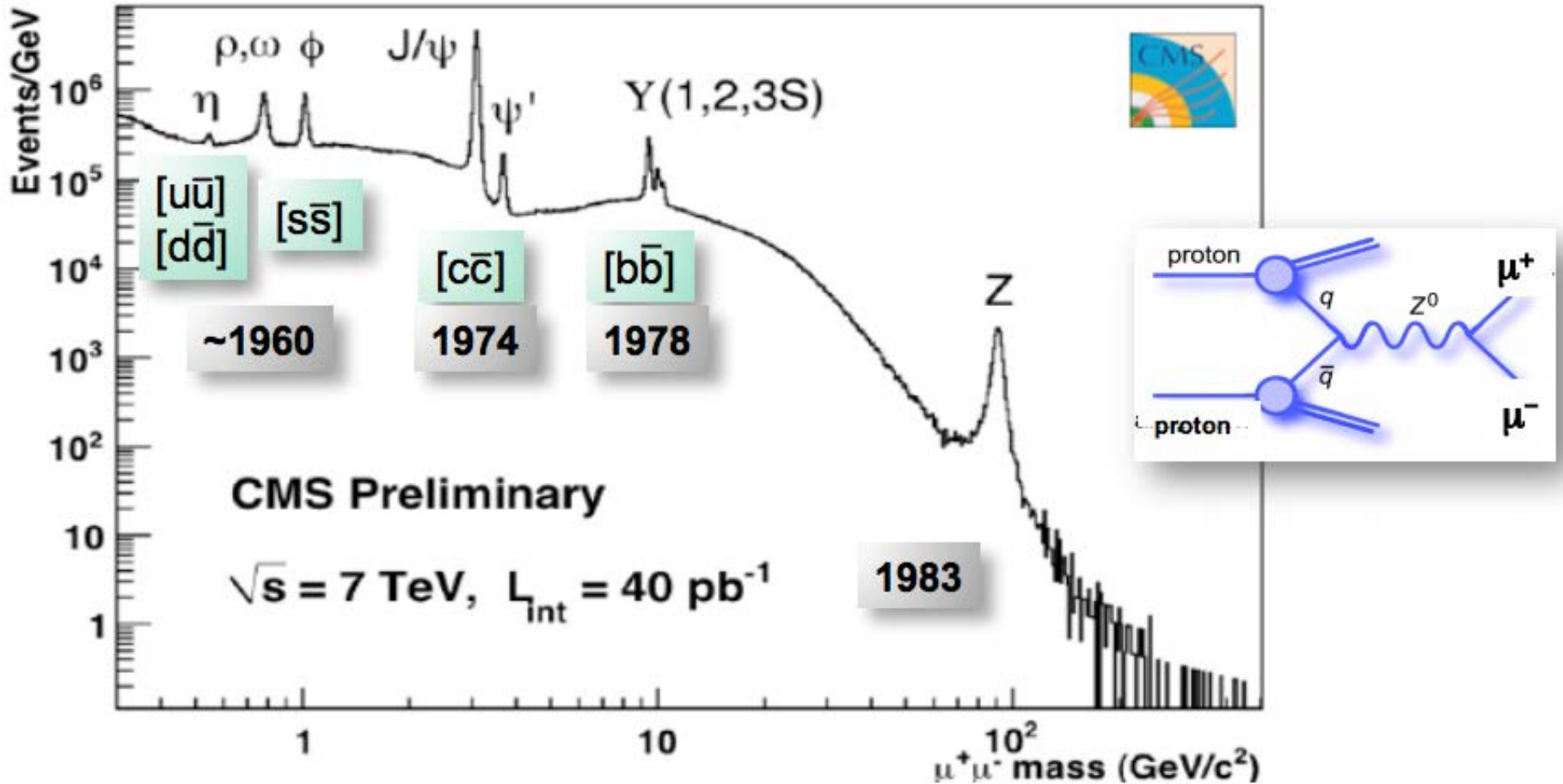
Data agrees with theoretical prediction over many orders of magnitude!





Example – search for a new heavy Z'

Many new physics models have a new heavy gauge boson which can decay to leptons. Like a Z but heavier - called Z' . Important to search for such new particles at the LHC.

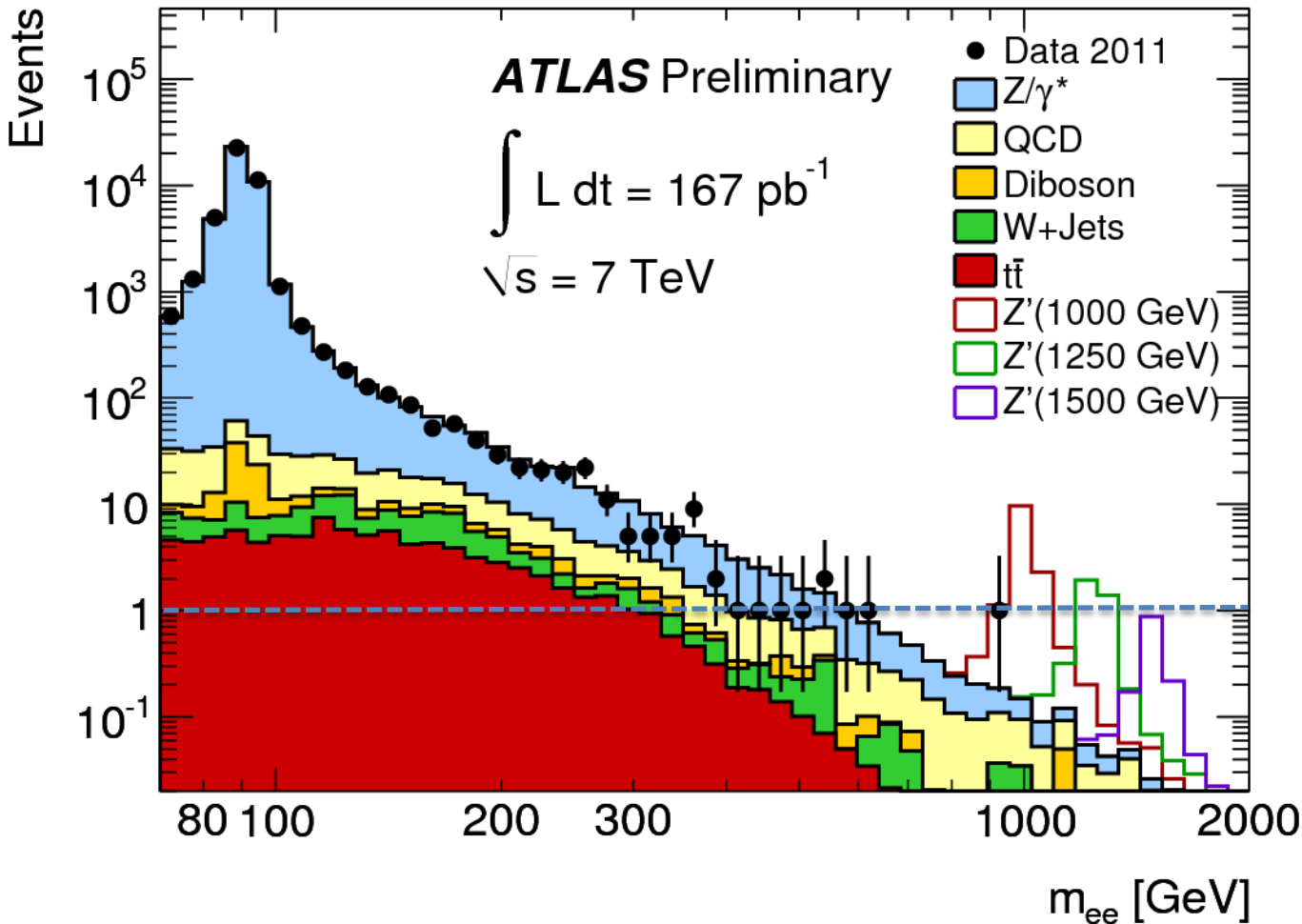


Historically many important discoveries (Nobel prizes) in di-lepton mass spectrum



Example – search for a new heavy Z'

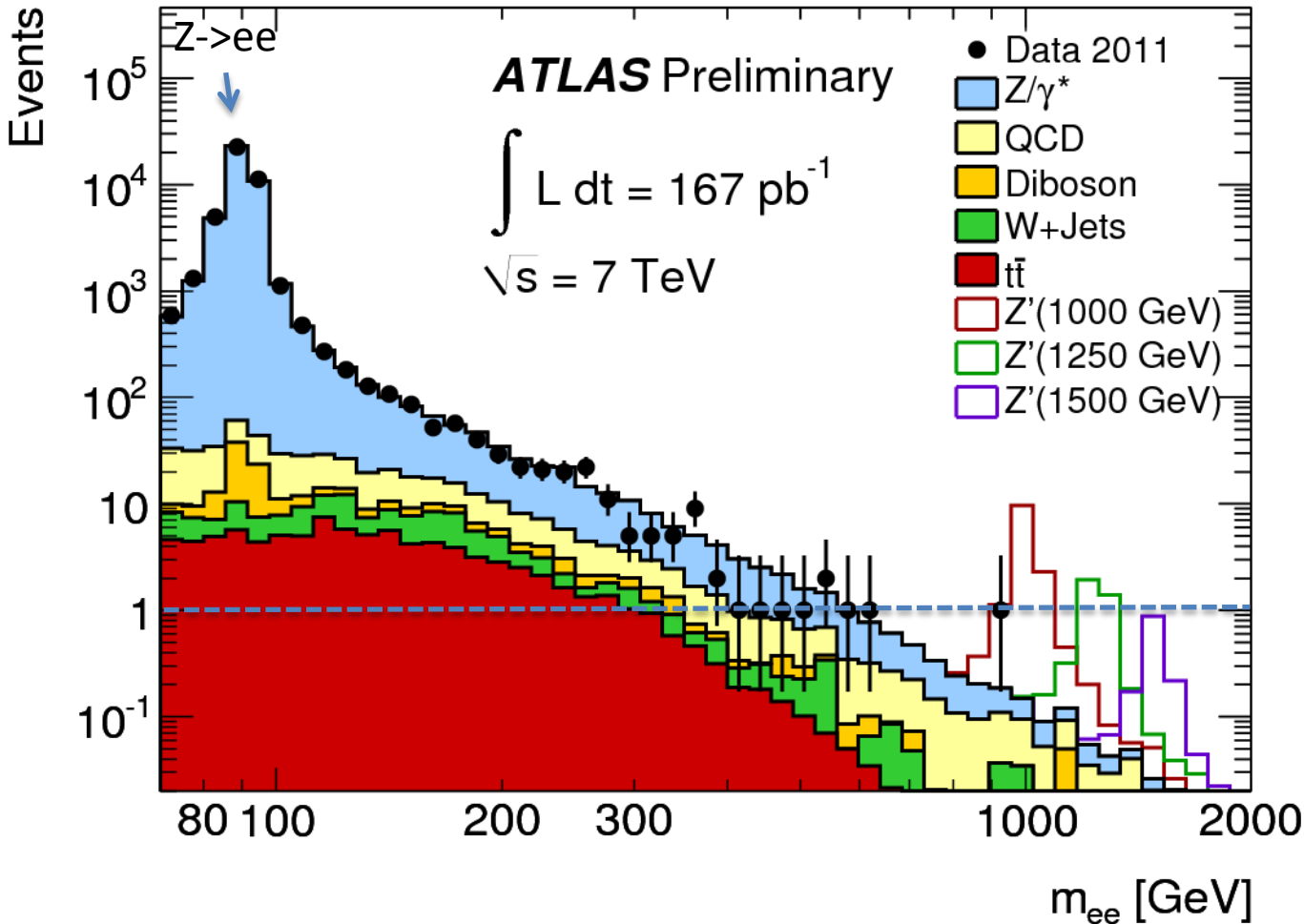
Like $Z \rightarrow ee$ but at higher mass.





Example – search for a new heavy Z'

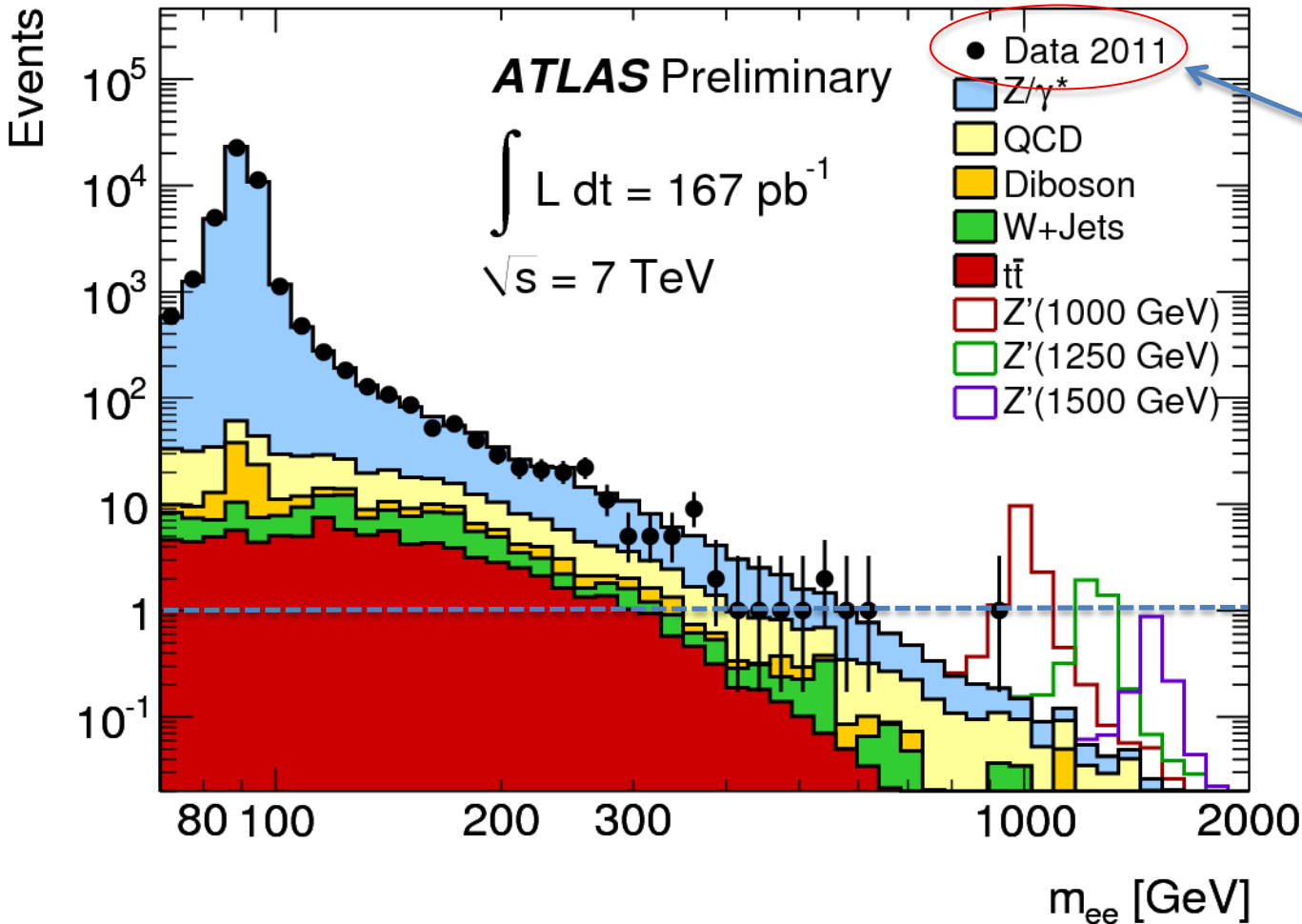
Like $Z \rightarrow ee$ but at higher mass.





Example – search for a new heavy Z'

Like $Z \rightarrow ee$ but at higher mass.

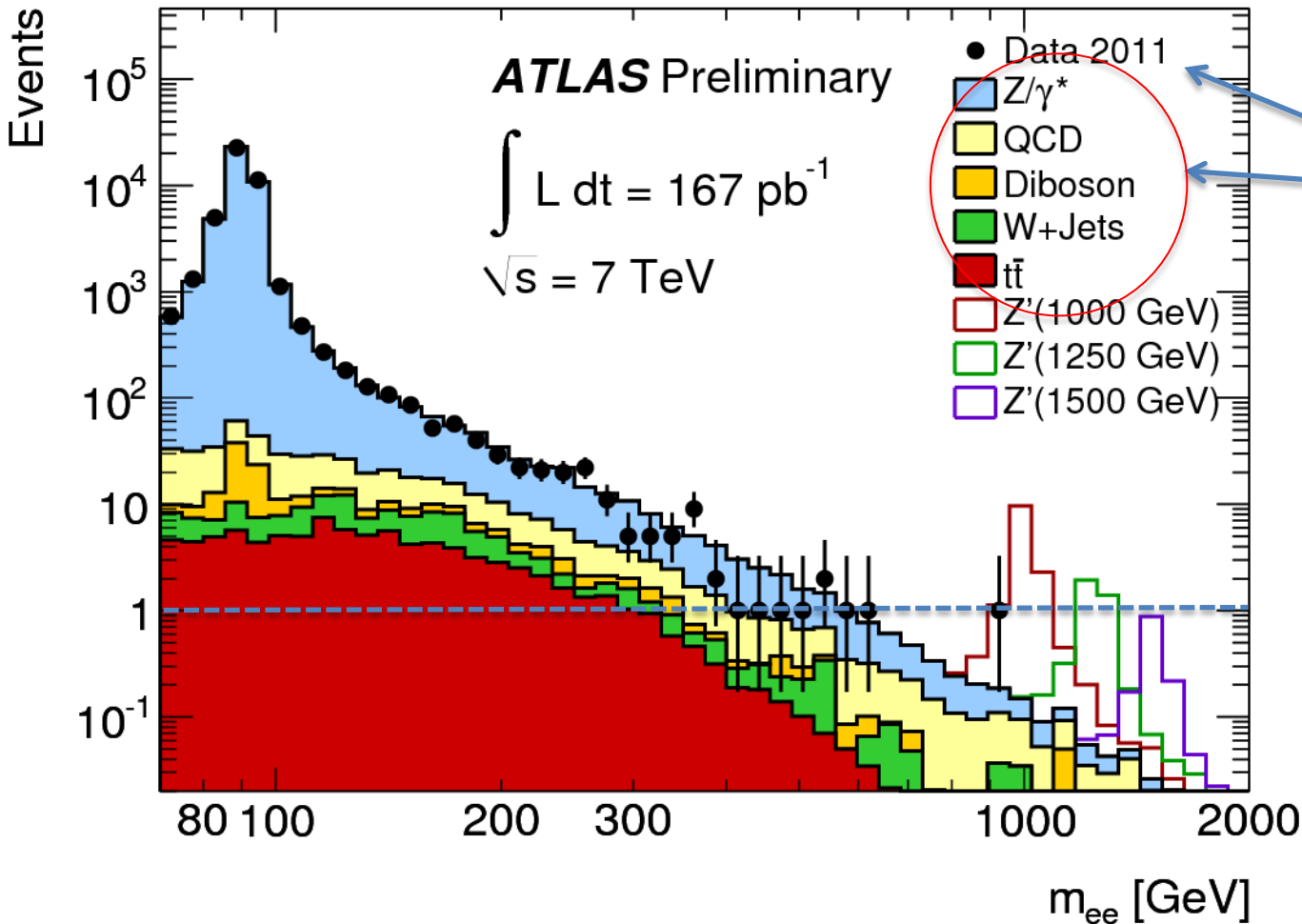


Select 2 electron candidates and plot their invariant mass for
1. Data



Example – search for a new heavy Z'

Like $Z \rightarrow ee$ but at higher mass.



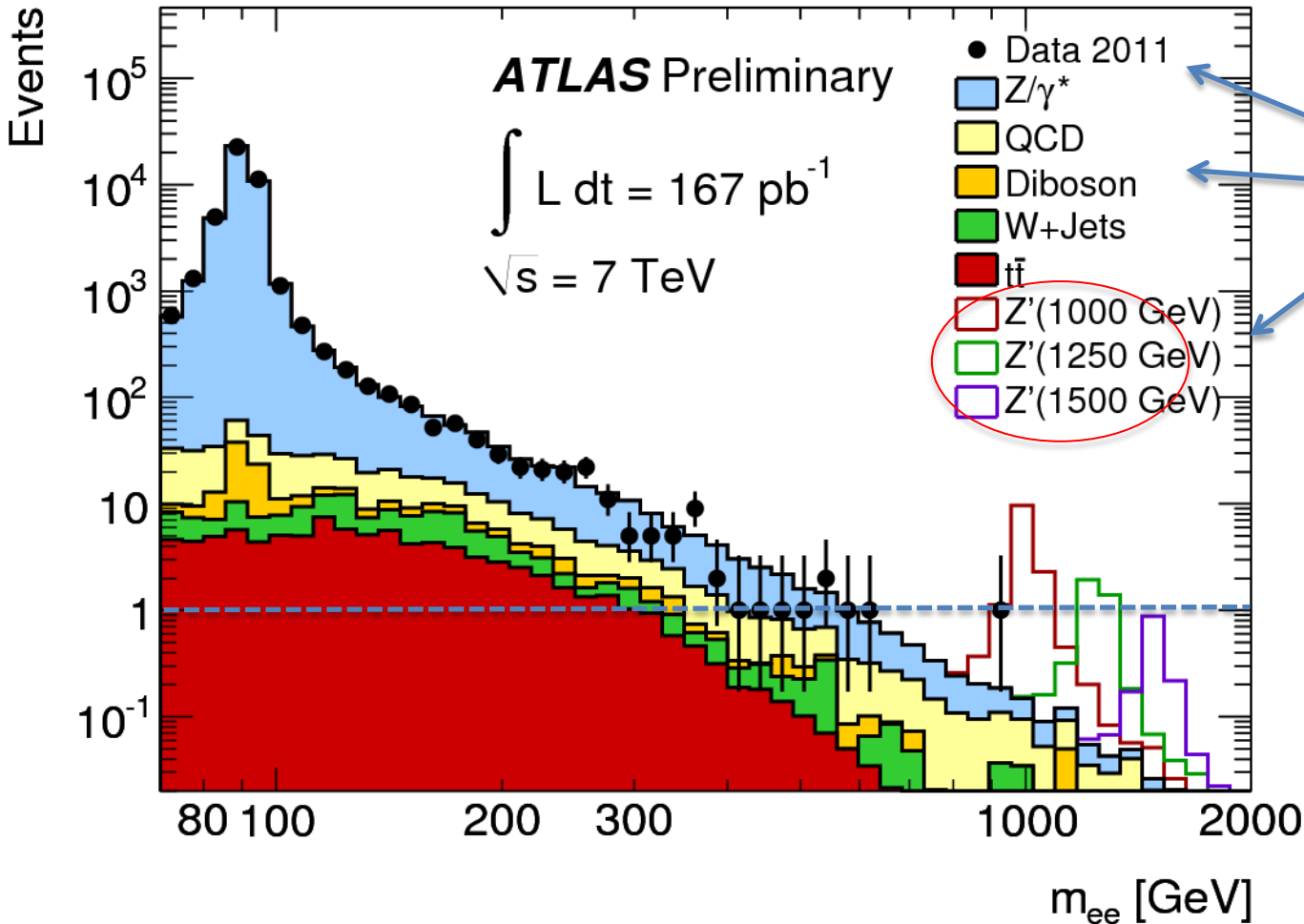
Select 2 electron candidates and plot their invariant mass for

1. Data
2. **Simulated background events**



Example – search for a new heavy Z'

Like $Z \rightarrow ee$ but at higher mass.

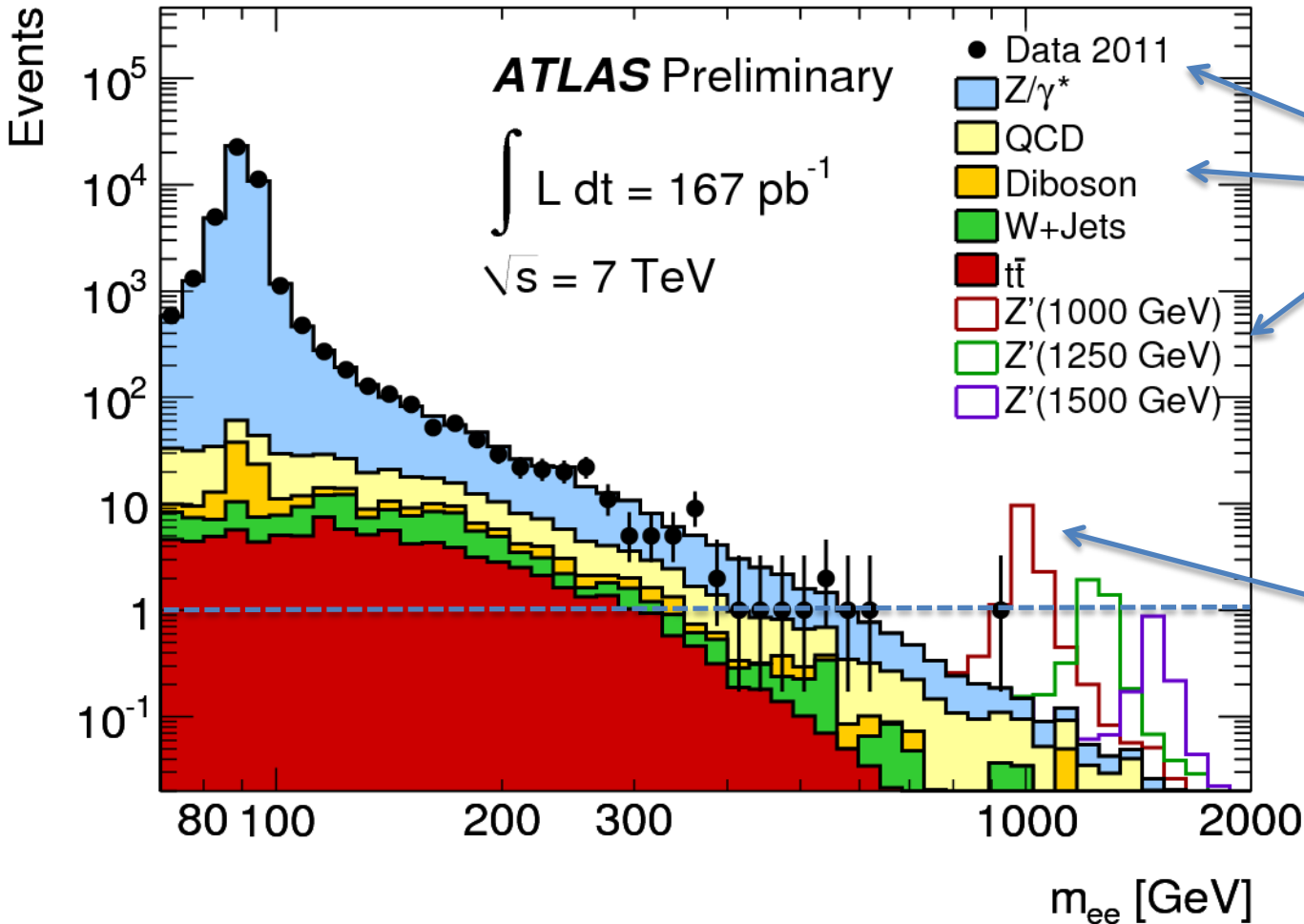


- Select 2 electron candidates and plot their invariant mass for
1. Data
 2. Simulated backgrounds events
 3. **Simulated signal (Z') with different masses**



Example – search for a new heavy Z'

Like $Z \rightarrow ee$ but at higher mass.



Select 2 electron candidates and plot their invariant mass for

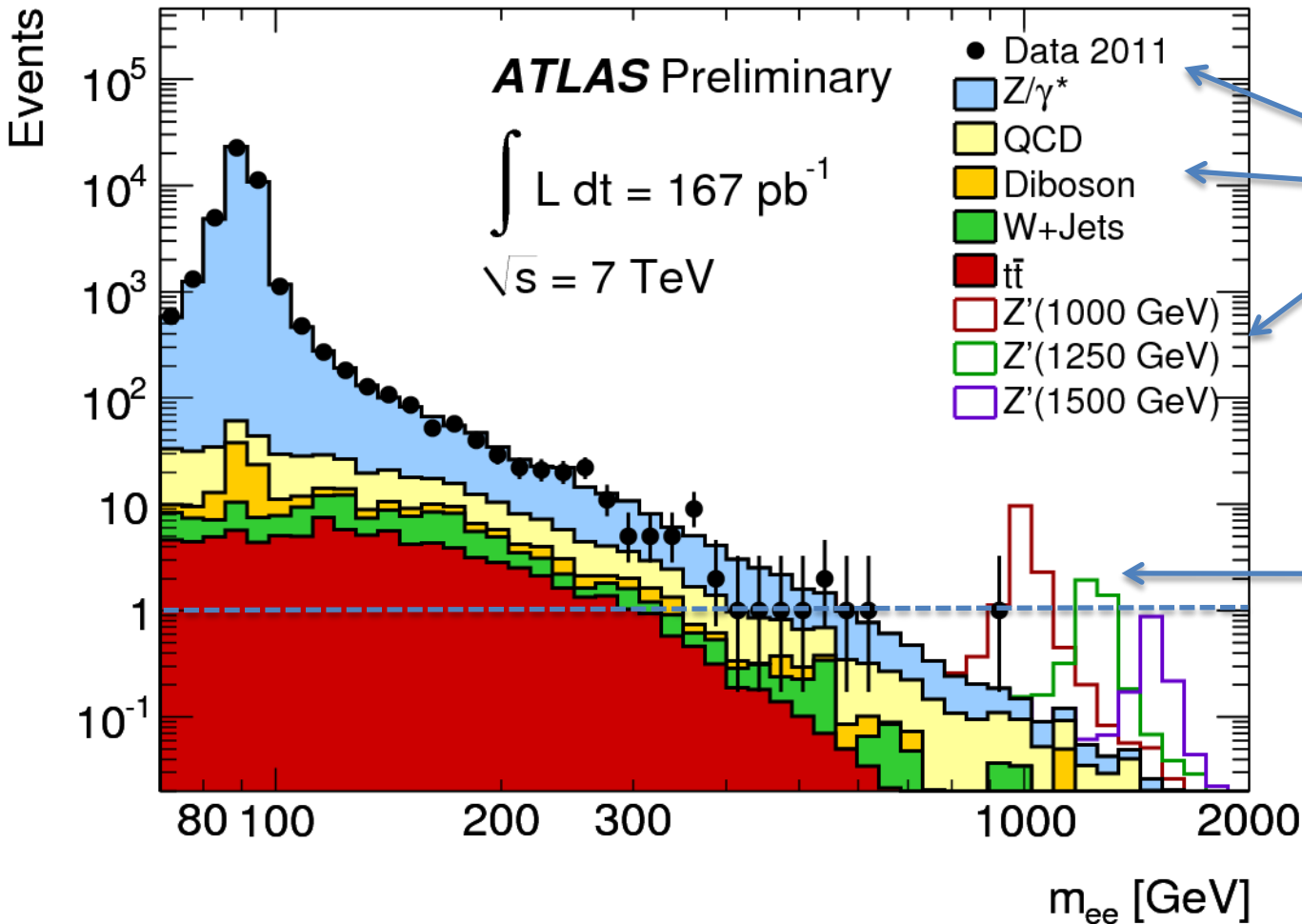
1. Data
2. Simulated backgrounds events
3. Simulated signal (Z') with different masses

Data inconsistent with a 1TeV Z'



Example – search for a new heavy Z'

Like $Z \rightarrow ee$ but at higher mass.



Select 2 electron candidates and plot their invariant mass for

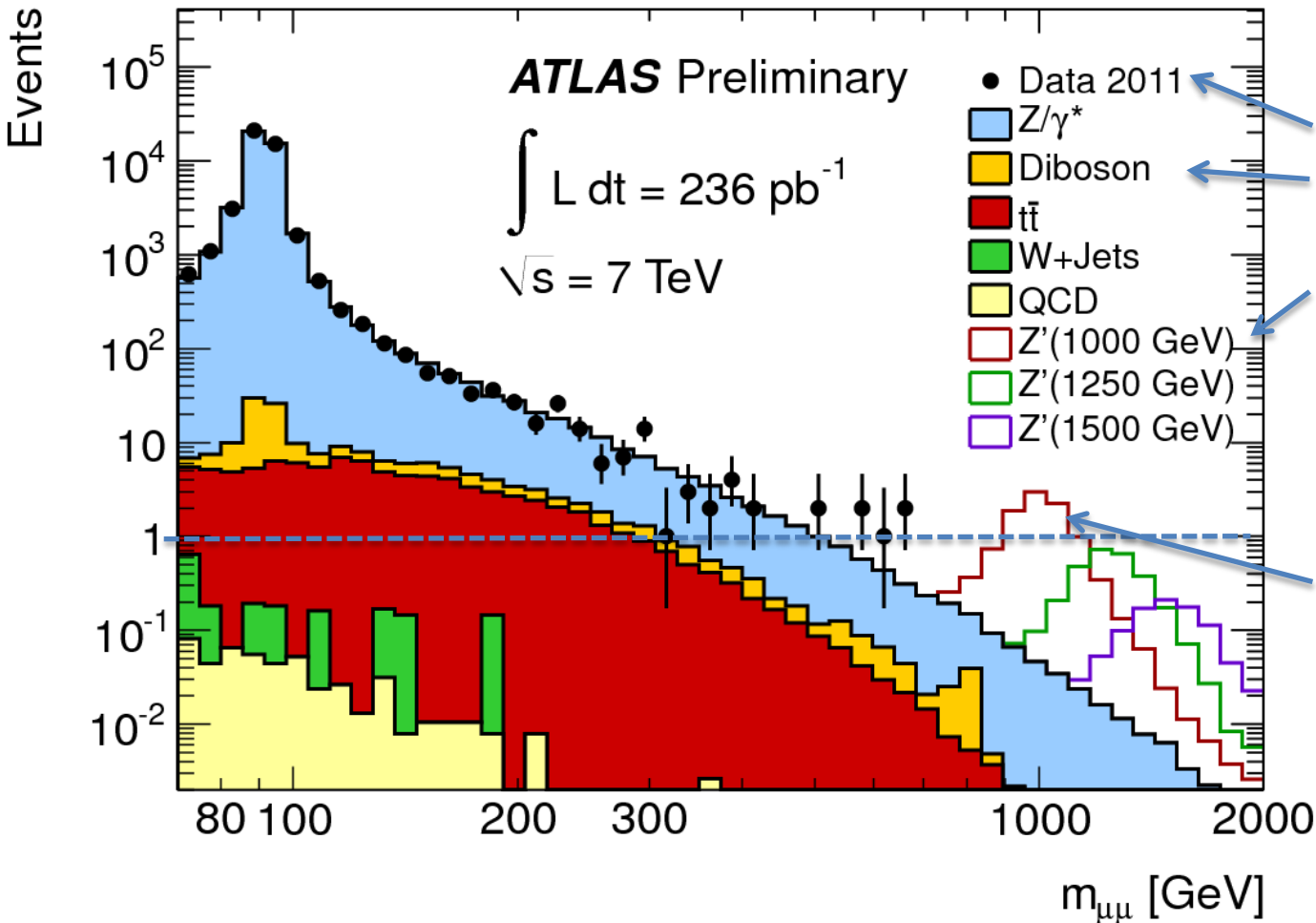
1. Data
2. Simulated backgrounds events
3. Simulated signal (Z') with different masses

Cross-section decreases with mass (higher the mass of the Z' , the more data needed to discover it)



Example – search for a new heavy Z'

Now for muons!

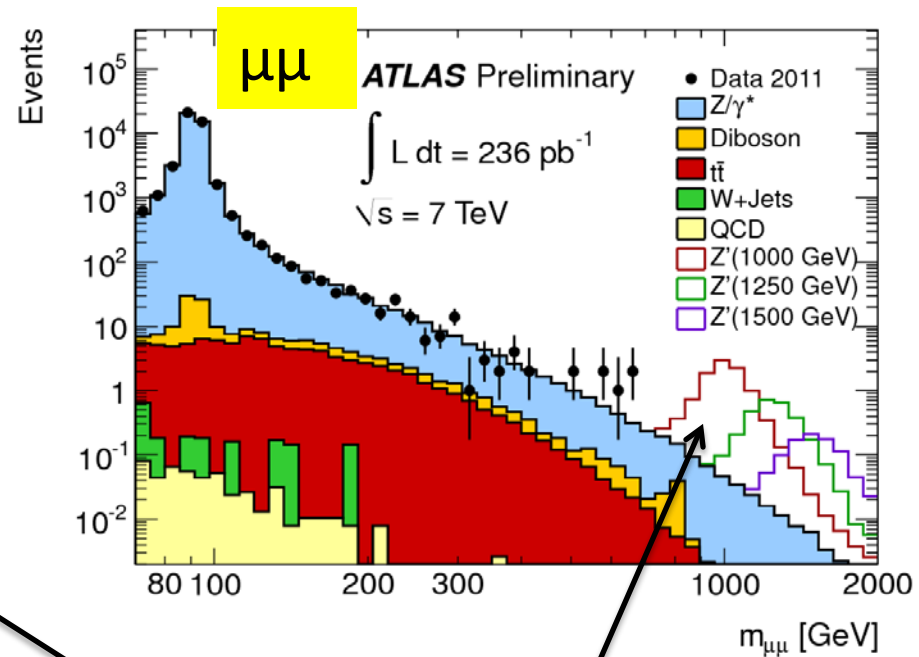
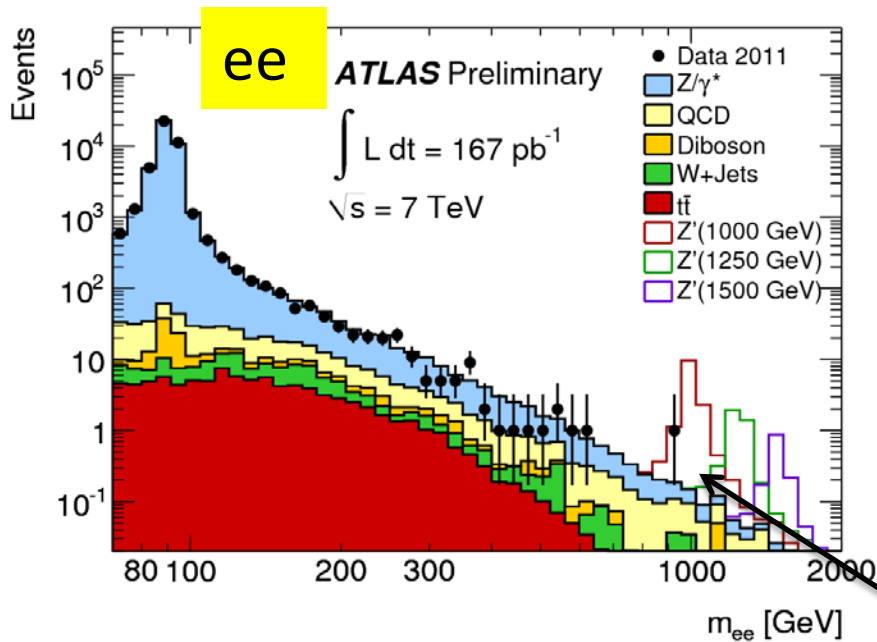


- Select 2 muon candidates and plot their invariant mass for
1. Data
 2. Simulated backgrounds events
 3. Simulated signal (Z') with different masses

Data inconsistent with a 1TeV Z'

Z' analysis

Combining the electron and muon channel the data exclude Z' upto mass of 1.4 TeV
 Need to take into account the statistical and systematic uncertainties!



Simulation tells us that the 1TeV Z' is narrower in electron decay mode than muon decay mode
 (Electron momentum resolution better at high energy)

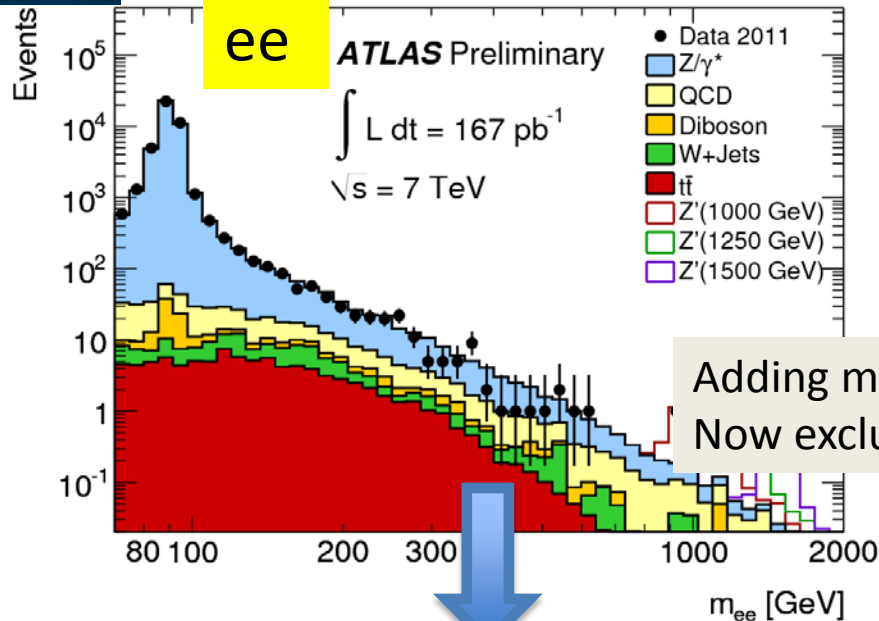
Background composition different in the electron and muon channels.

(Different amount of data in 2 plots because data-quality requirements different for e and μ)

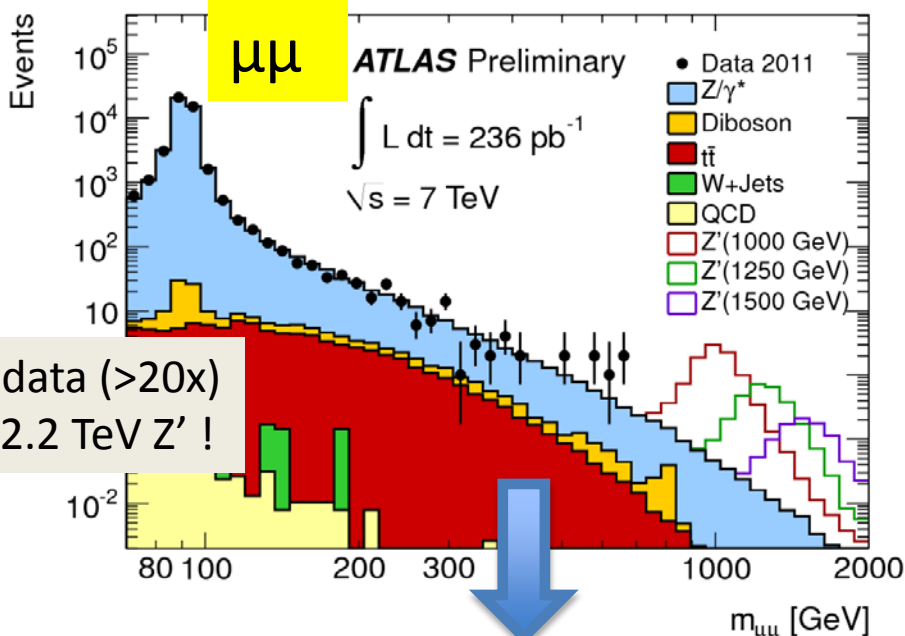


Z' analysis

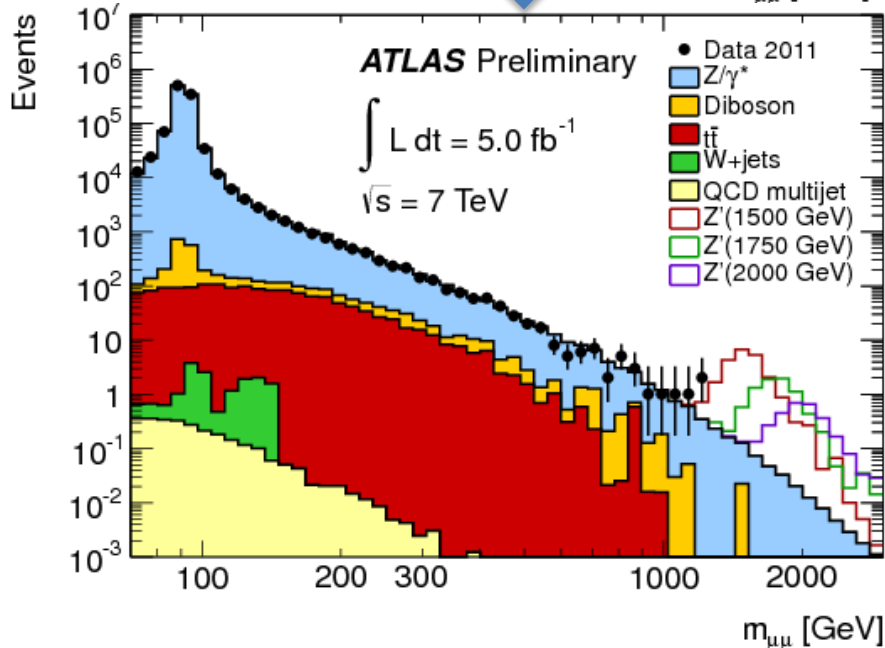
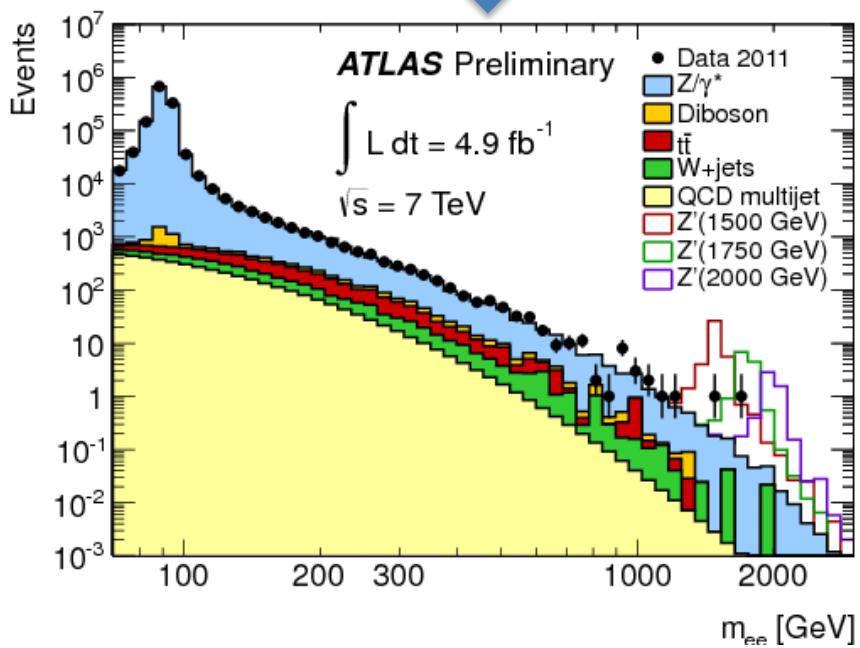
ee



$\mu\mu$

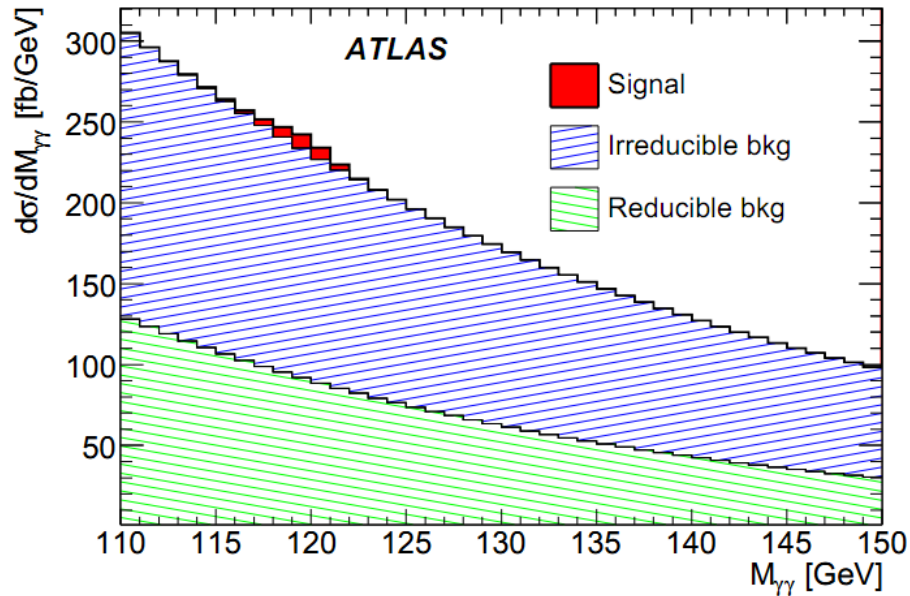


Adding more data (>20x)
Now exclude 2.2 TeV Z' !



H- \rightarrow $\gamma\gamma$

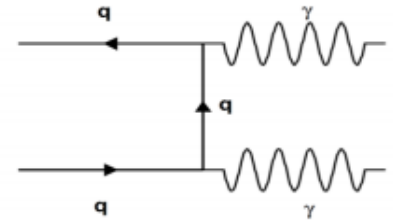
H- \rightarrow $\gamma\gamma$ is the best way of discovering a low mass Higgs at the LHC.



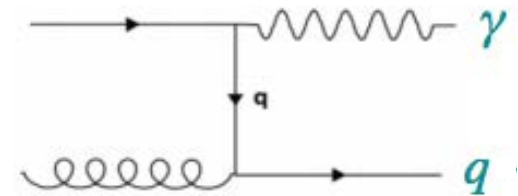
Tiny signal on top of large background.
Good resolution is essential to being able to observe this.

Signal extracting by fitting the 2 photon mass spectrum.

Irreducible background:
-2 real photons



Reducible background:
-photon + jet
(where jet fakes a photon)



Need good photon/jet separation to minimize the reducible background.



Importance of Resolution

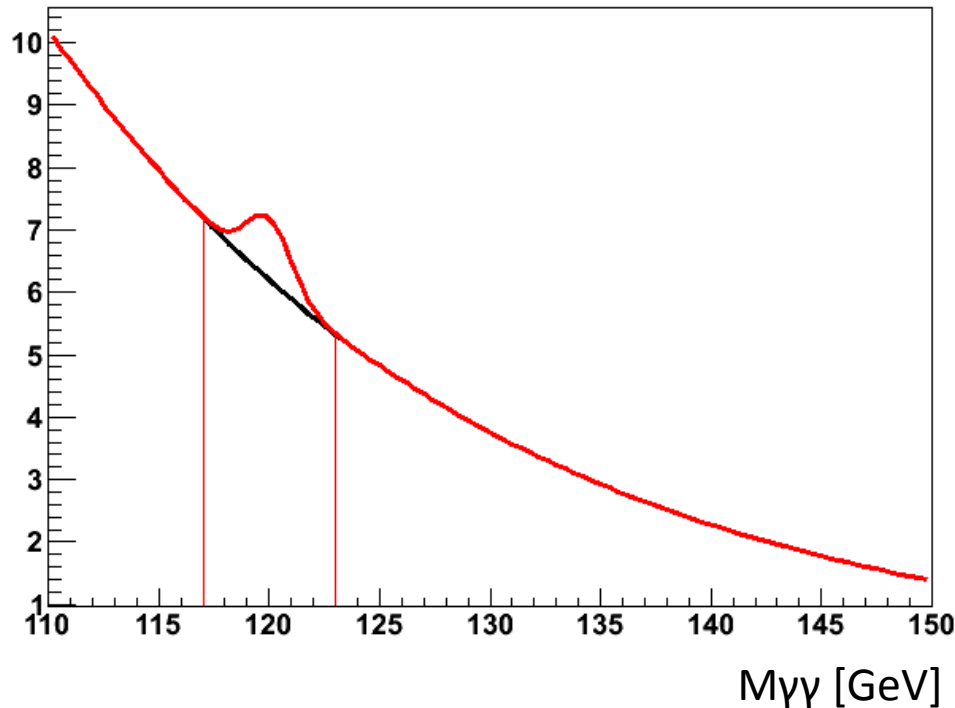
Toy example: Signal peak on exponential background.

2 different signal resolutions. Same number of signal events in each peak

Would discover the left hand signal much quicker!

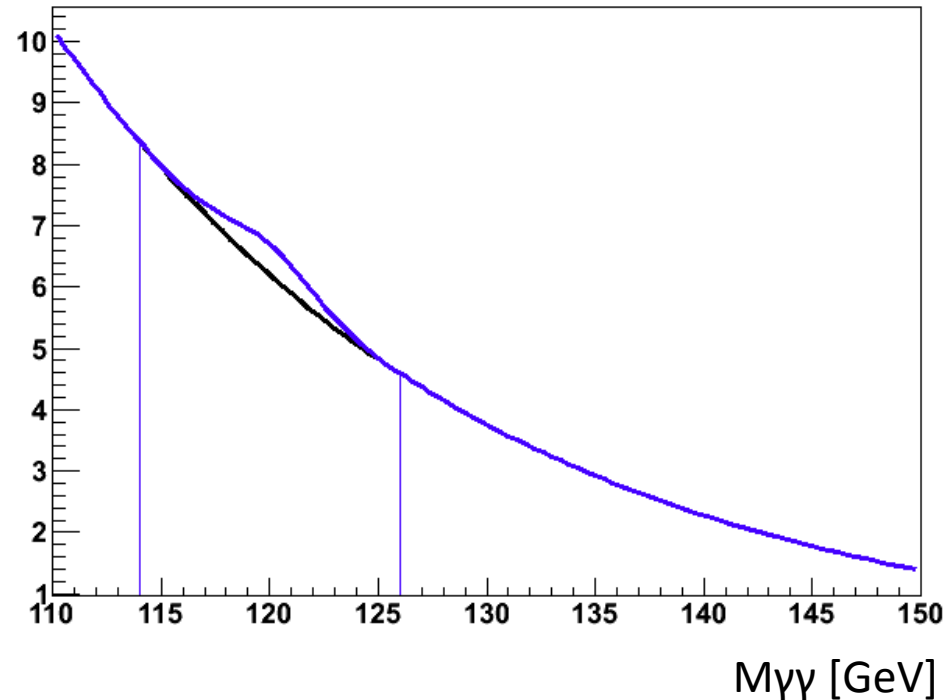
Mass resolution 1 GeV

Signal over background in cut range $\sim 10\%$



Mass resolution 2 GeV

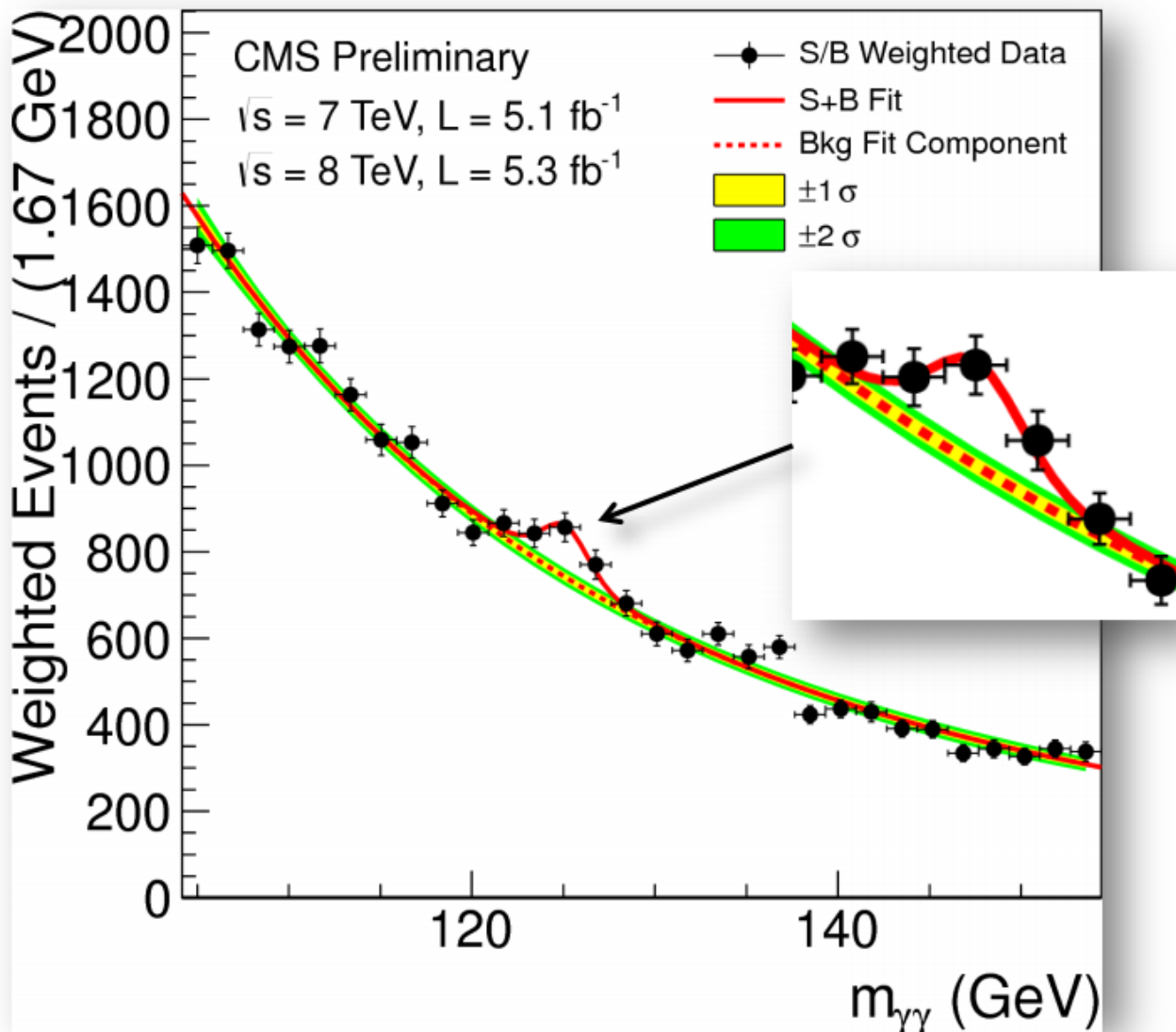
Signal over background in cut range $\sim 5\%$



Very important to build the detector to give you the best resolution.

But also to optimize the reconstruction algorithms and calibrations to give the best resolution possible for that detector.

H-> $\gamma\gamma$ Results



Analysis carried out in different categories of events with expected different resolution and S/B.

Left plot shows mass plot where events are WEIGHTED by the category (more like what the fit “see’s”).

ATLAS & CMS both see significant ($>4\sigma$) peak at $\sim 126\text{GeV}$. Most important result in the recent Higgs observation.

In both experiments a huge amount of work went into getting best resolution which allowed this result to happen so quickly!



Another Example: $H \rightarrow ZZ$

- Searching for the Higgs boson in the decay $H \rightarrow ZZ$
- The Z can decay like
 - $Z \rightarrow qq$ (quarks seen as jets in the detector)
 - $Z \rightarrow ll$ (electrons, muons or taus)
 - $Z \rightarrow \nu\nu$ (neutrino's do not interact with the detector and so are only 'seen' as missing energy)
- $H \rightarrow ZZ \rightarrow l^+l^-l^+l^-$ is by far the easiest to detect experimentally (we only look for $l = \text{electron or muon}$, as these are the easiest)
- $H \rightarrow 4l$ is called the 'golden mode' as experimentally it is by far the easiest
 - Leptons have low background
 - Leptons reconstruction has high efficiency
 - Leptons have good momentum resolution

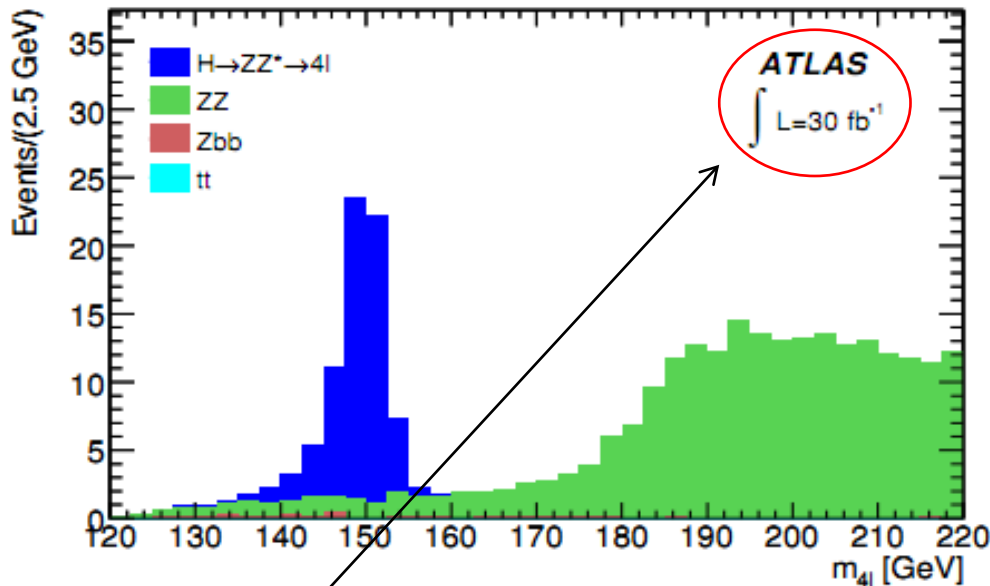


Simple H->ZZ analysis steps

1. Select events which contain reconstructed: $e^+ e^- \mu^+ \mu^-$ (or 4e or 4 μ)
2. make sure mass of the lepton pairs is consistent with the Z mass
3. Histogram the mass of the 4 leptons
4. See a peak corresponding to the Higgs (hopefully!)

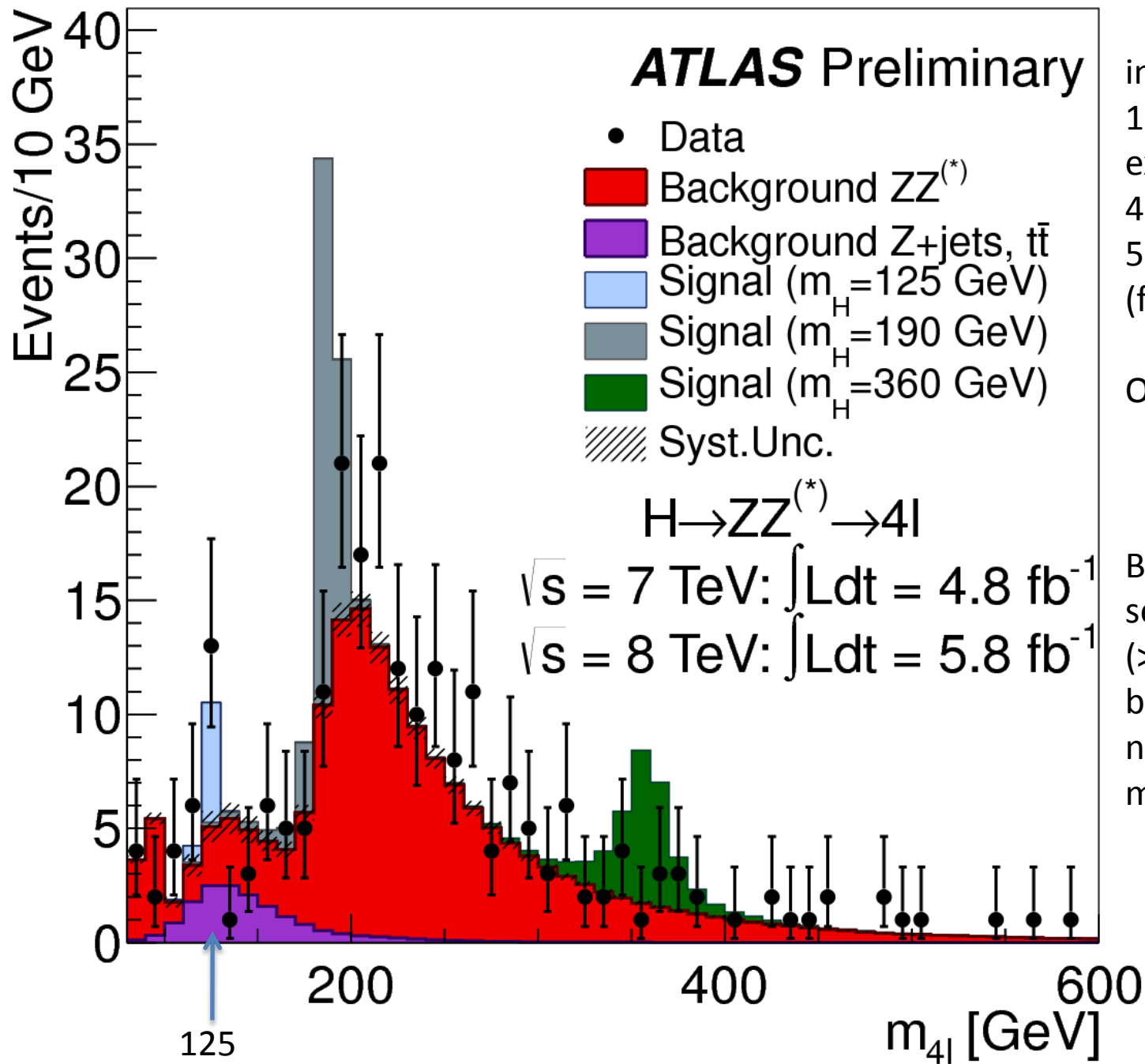
(size of peak depends on Higgs mass, for low mass Higgs ($<2M_Z$) one of the Z's is off shell – doesn't have Z mass!)

Simulated data!



Lepton efficiency very important for this mode.
If we have a 90% efficiency per lepton that gives $0.90^4=66\%$ efficiency for the Higgs!

~3 times more data than now!

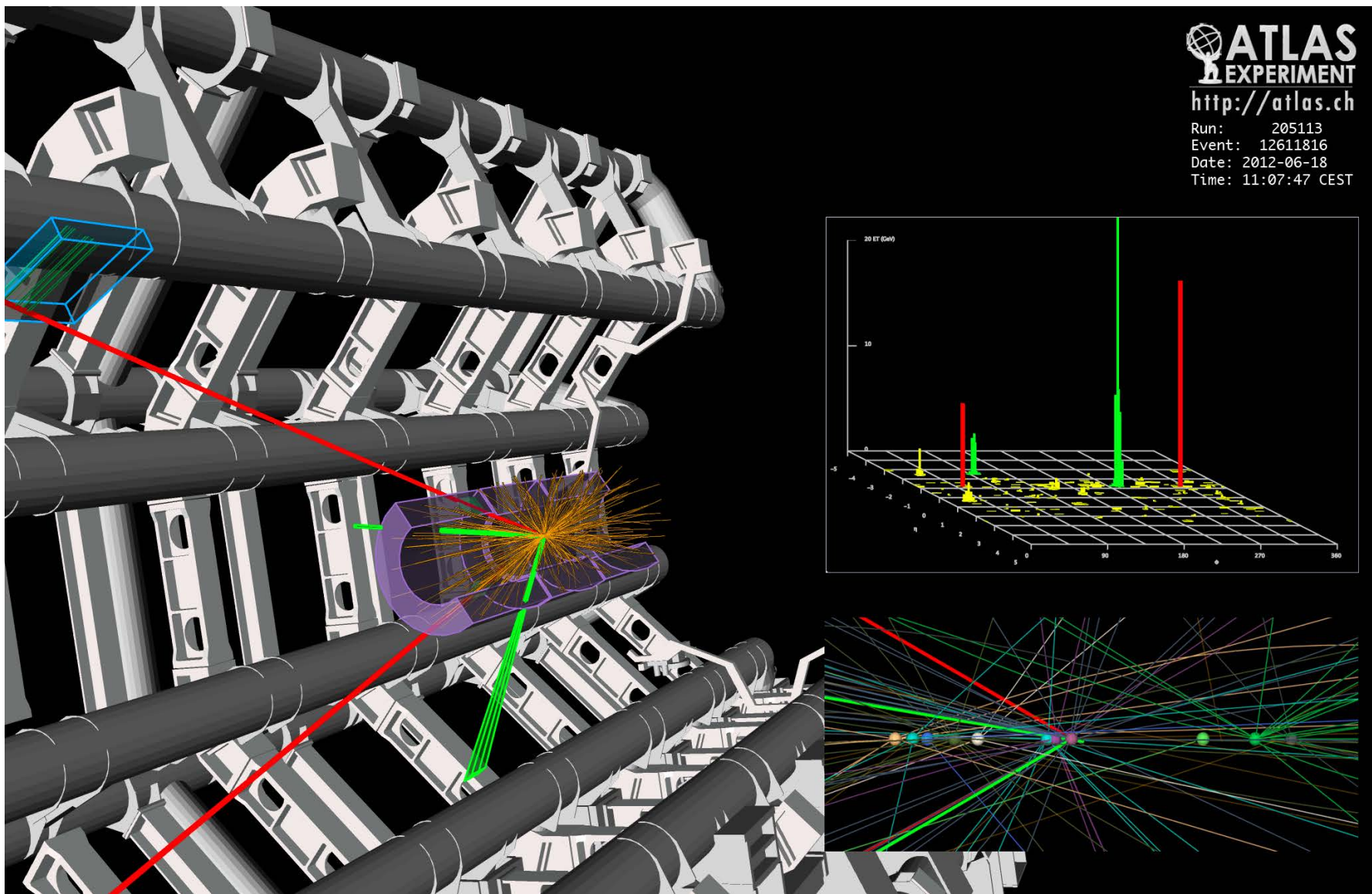


in ATLAS in range:
 $120 < m_{4l} < 130$ GeV
 expect:
 4.8 ± 0.3 Bkgd events
 5.3 ± 0.5 Signal events
 (for 125 GeV Higgs)

Observe 13 events!

Both ATLAS & CMS
 see significant excess
 ($>3\sigma$) over expected
 background in
 number of events at
 $m_{4l} \sim 125$ GeV

H \rightarrow 2 μ 2e candidate

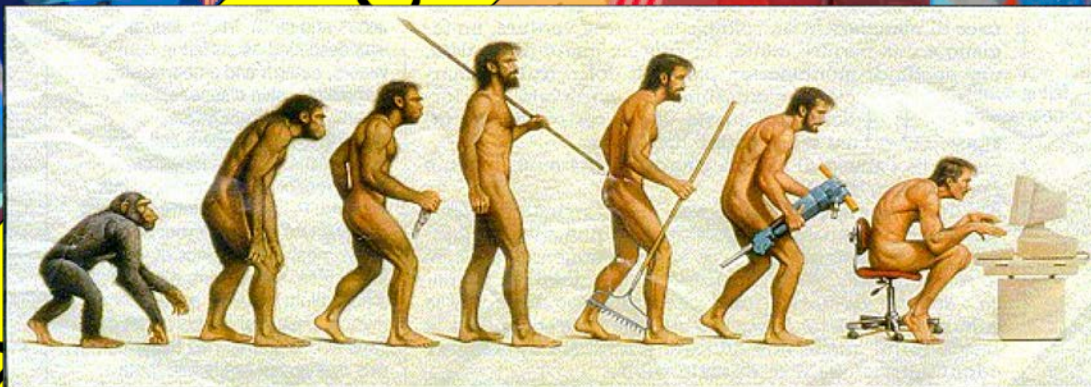




Summary from Physics examples

- Data Quality very important to not include junk data in the analysis
- Jet cross-section
 - Energy calibration uncertainty leads to a big uncertainty due to the sharply falling spectrum
- Z' :
 - use of simulation to see what a new physics signal would look like
- $H \rightarrow \gamma\gamma$:
 - importance of resolution
 - Importance of background rejection
- $H \rightarrow ZZ \rightarrow 4l$:
 - importance of high efficiency
- Many other types of physics analyses (measuring cross-sections, masses, lifetimes)
 - Require also accurate knowledge of the efficiency, and the resolutions and the background rates

Computing behind this all...



Somewhere, something went terribly wrong

S



Some numbers

Examples from ATLAS

- **Rate** of events streaming out from High-Level Trigger farm ~ 400 Hz
- each event has a size of the order of **1.5MB**
- about 10^7 events in total per day
- will have roughly 150 “physics” days per year
- thus about 10^9 evts/year, a few **Pbyte**



“prompt” processing

- Reco time per event on std. CPU: < 30 sec (on lxplus)
- increases with pileup (more combinatorics in the tracking)



simulating a few billions of events

- are mostly done at computing centers outside CERN
- Simulation very CPU intensive

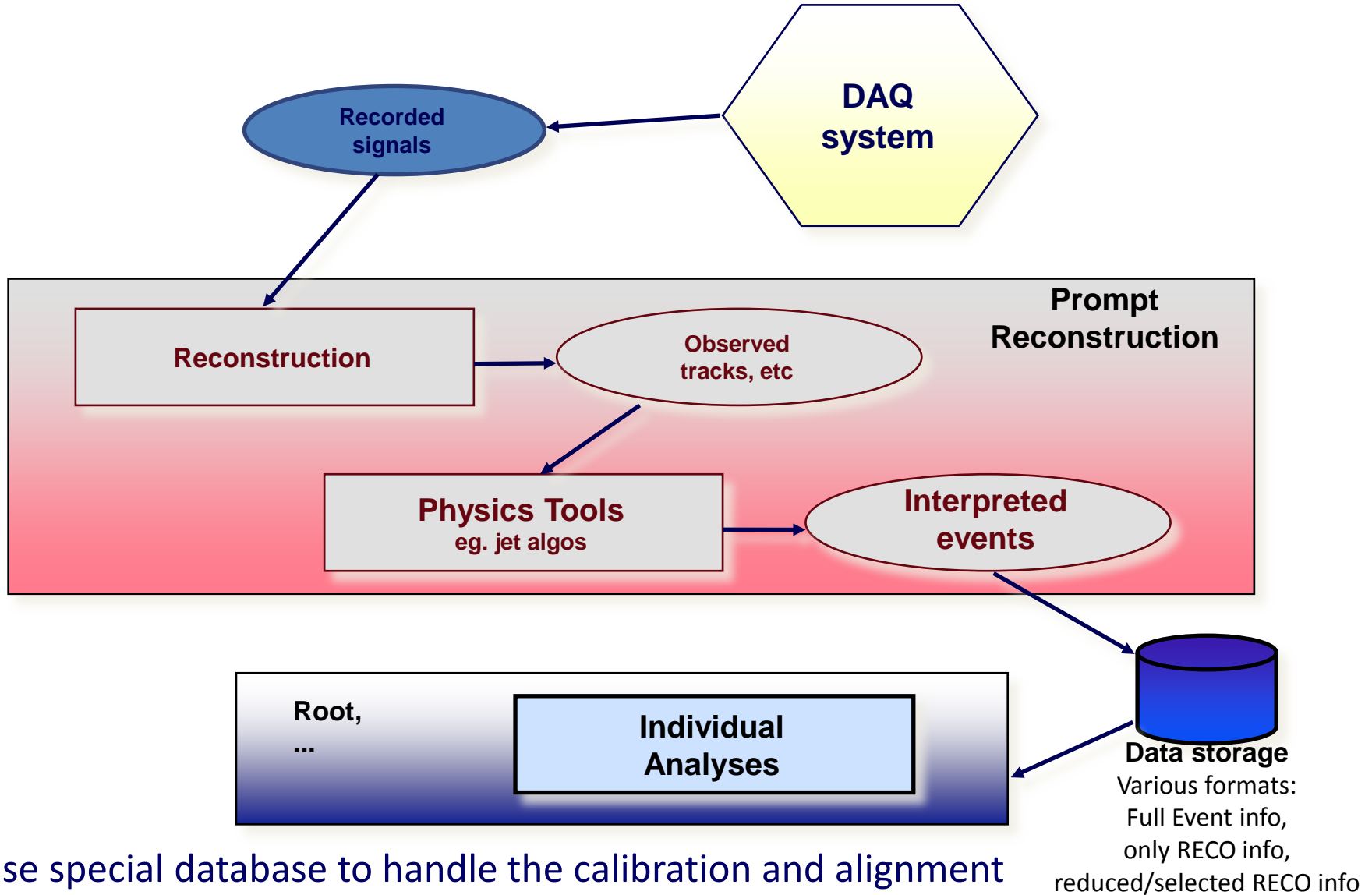
~ 4 million lines of code (reconstruction and simulation)

- ~ 1000 software developers on ATLAS





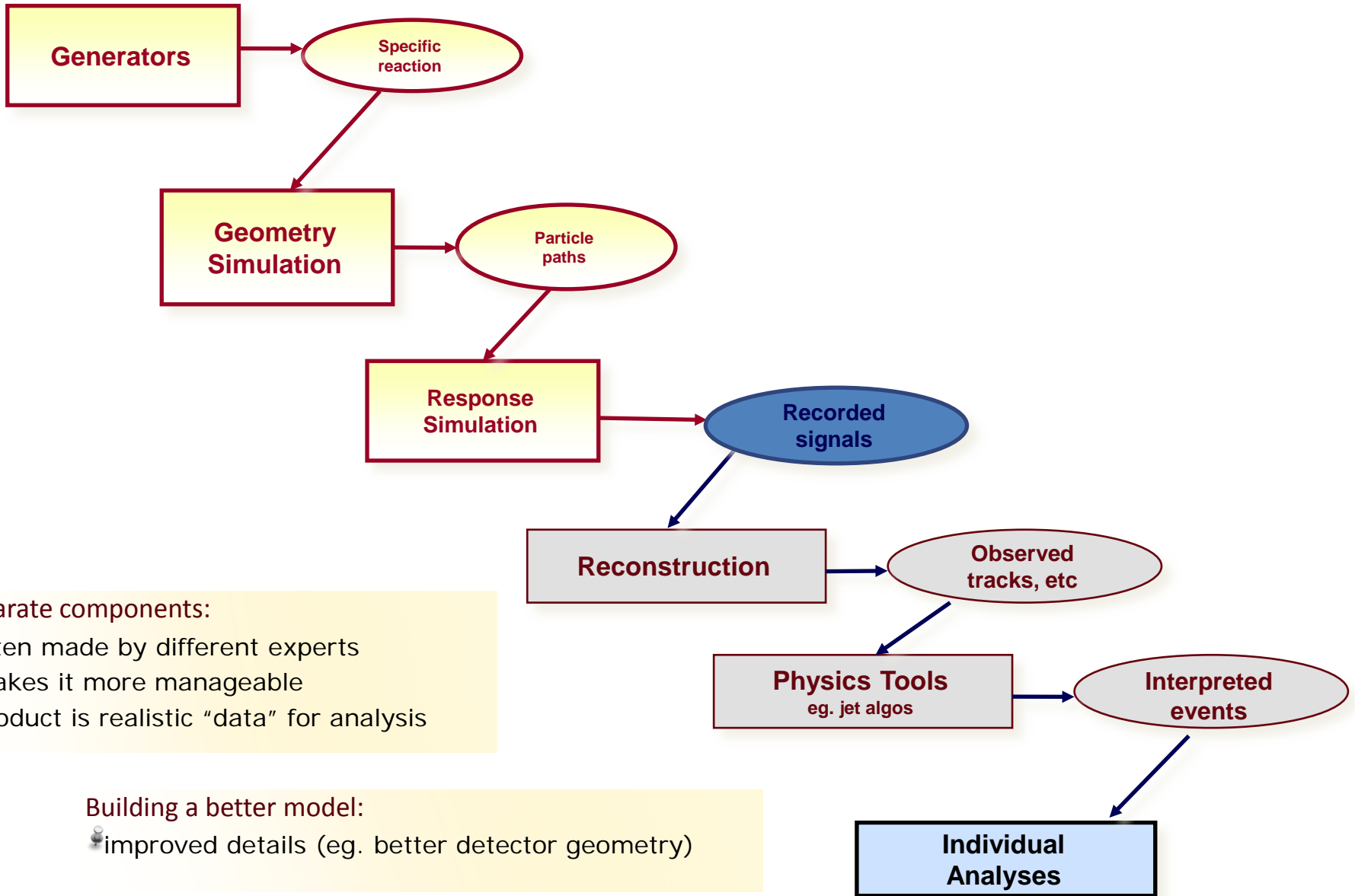
RECO flow



Use special database to handle the calibration and alignment data needed in reconstruction



Flow of simulated data



Separate components:

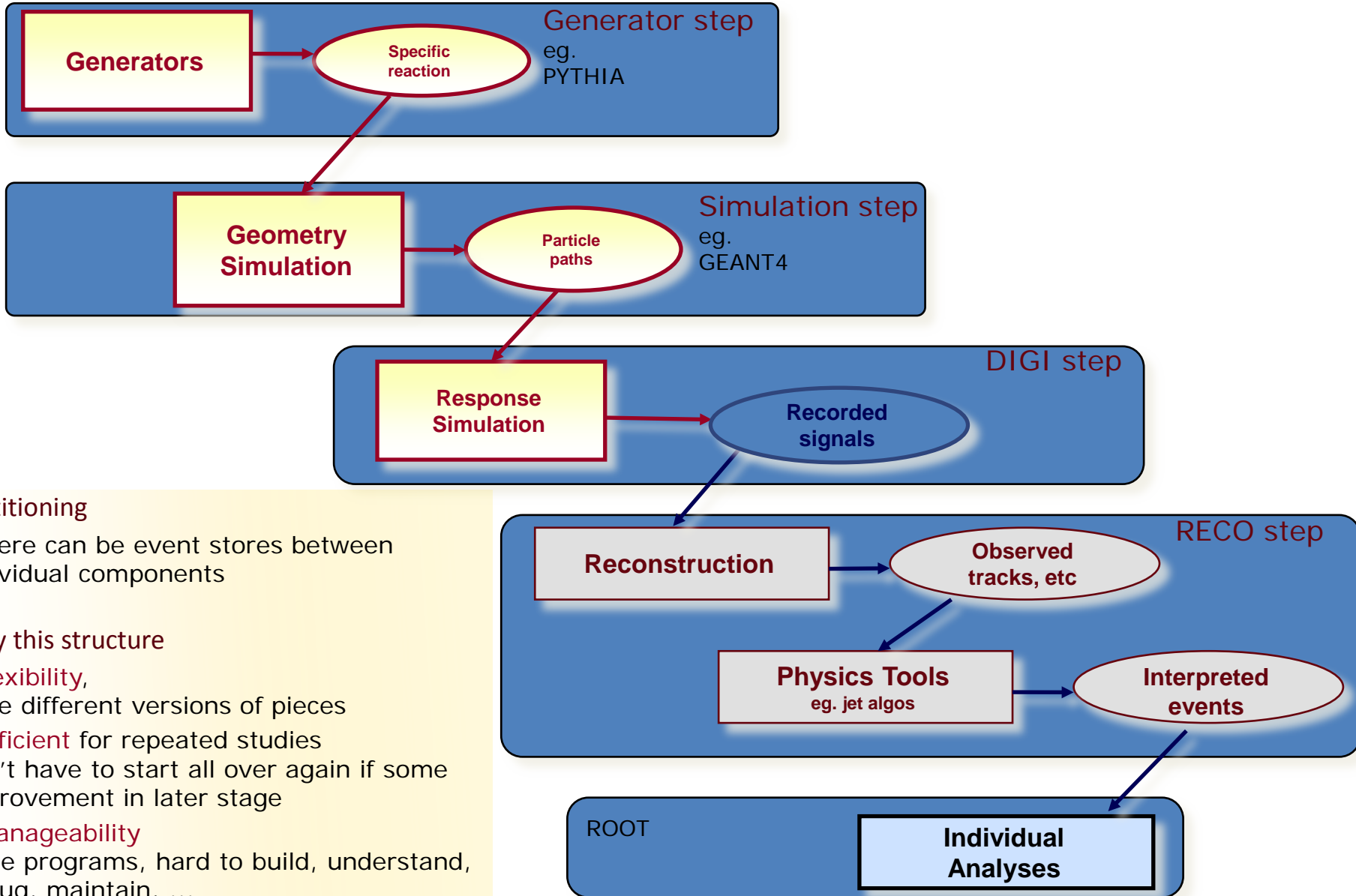
- often made by different experts
- makes it more manageable
- Product is realistic "data" for analysis

Building a better model:

- improved details (eg. better detector geometry)



Partitioning production systems



Partitioning

- there can be event stores between individual components

Why this structure

- **flexibility**, have different versions of pieces
- **efficient** for repeated studies don't have to start all over again if some improvement in later stage
- **Manageability** large programs, hard to build, understand, debug, maintain, ...

Conclusions

- **Reconstruction and Analysis**
is how we get from raw data to physics papers
- Sophisticated reconstruction algorithms + calibrations and alignments needed
 - High efficiencies, good resolutions and low fake rates
 - Important to get the best physics out of the experiment
- Detailed simulation also plays a key role
- Complex software infrastructure needed to be able to obtain the final physics results
- ~~Any~~ All discoveries at the LHC will rely on the data-quality, simulation, reconstruction and analysis chain working well
- Even to me it is often a miracle that we can generate wonderful results from these complicated instruments!

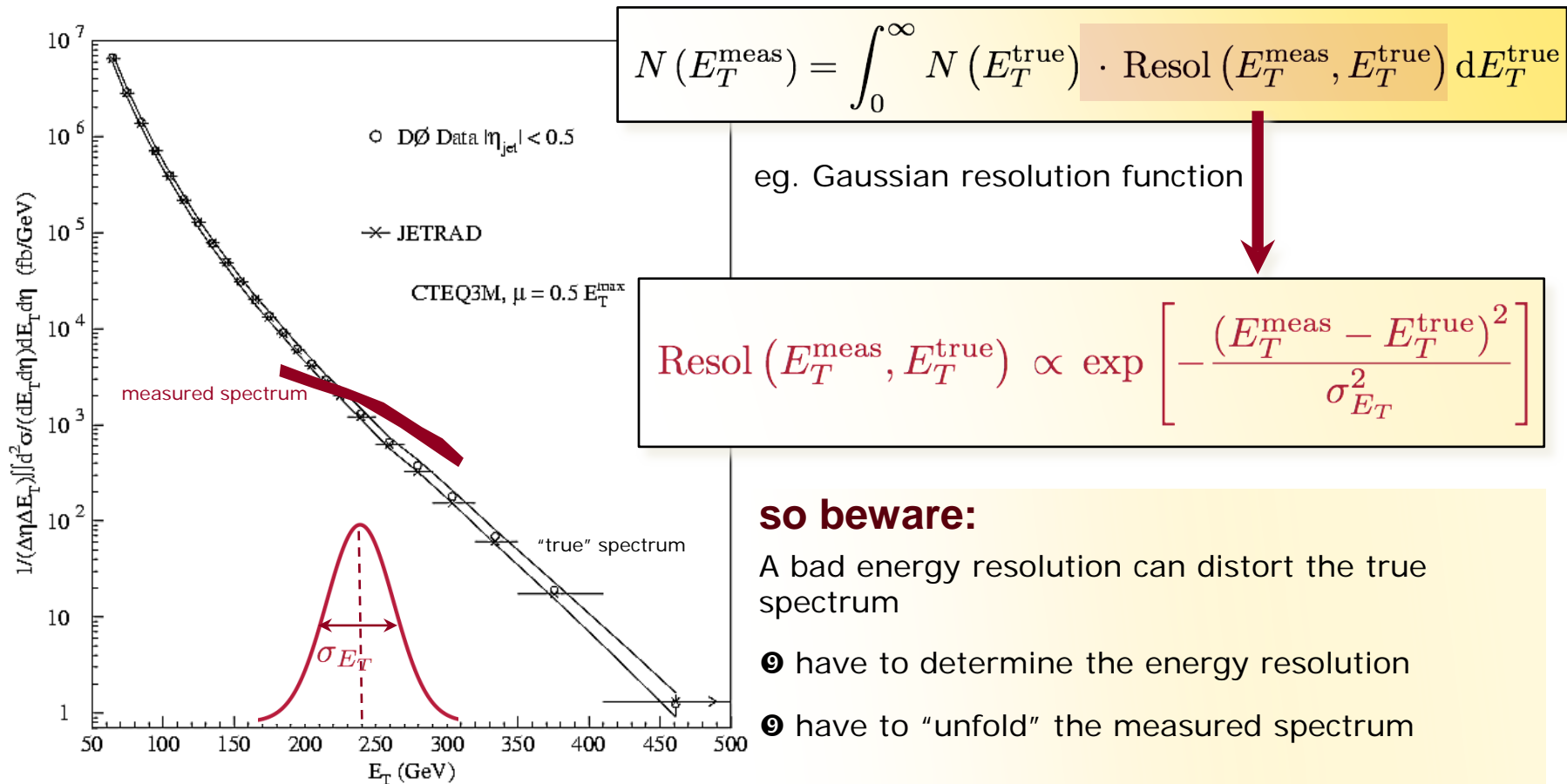


Final remarks

- This is **GREAT** time to be a CERN summer student
- The experiments have a lot of high quality data
- The LHC is working great
- The experiments all are working very well
- ~~There could be an exciting physics discovery at any time~~
(Already happened!!!)
- Work hard
 - Learn what you can!
- And most of all - Enjoy yourself!!!

Problem 2 : Energy resolution

- The **energy resolution** can distort the spectrum
- Again : Critical because of very steeply falling spectrum!



so beware:

A bad energy resolution can distort the true spectrum

- have to determine the energy resolution
- have to “unfold” the measured spectrum

● problem is minimized if **bin width** $\sim \sigma_{E_T}$

Estimated energy in the ECAL:

$$E_{e,\gamma} = F \times \sum_{\text{clusters}} G c_i A_i$$

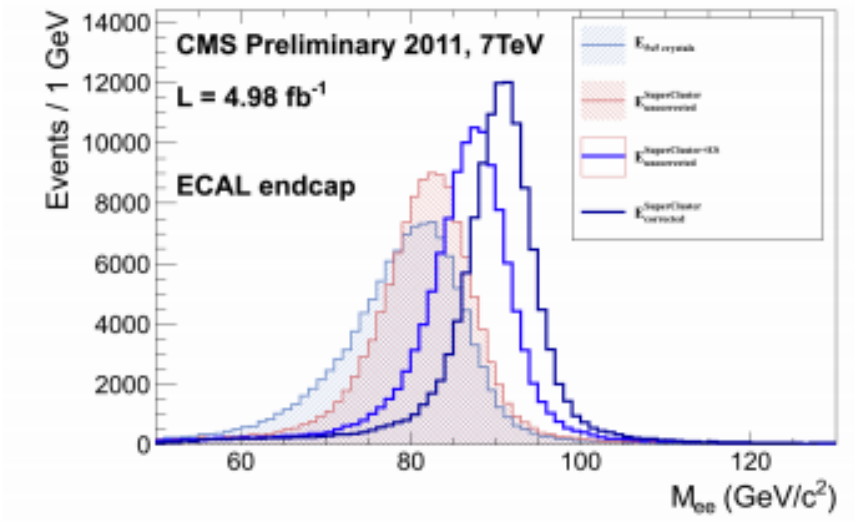
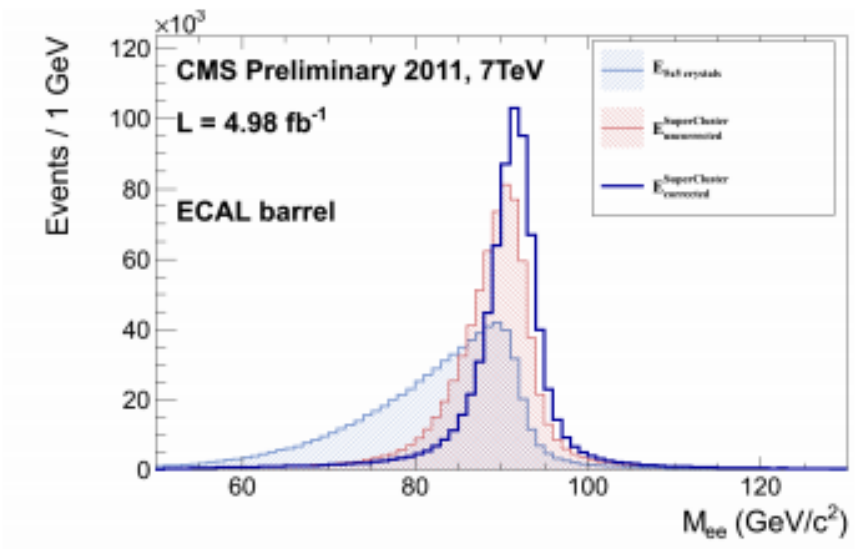
Corrections Calibration

Energy correction scheme

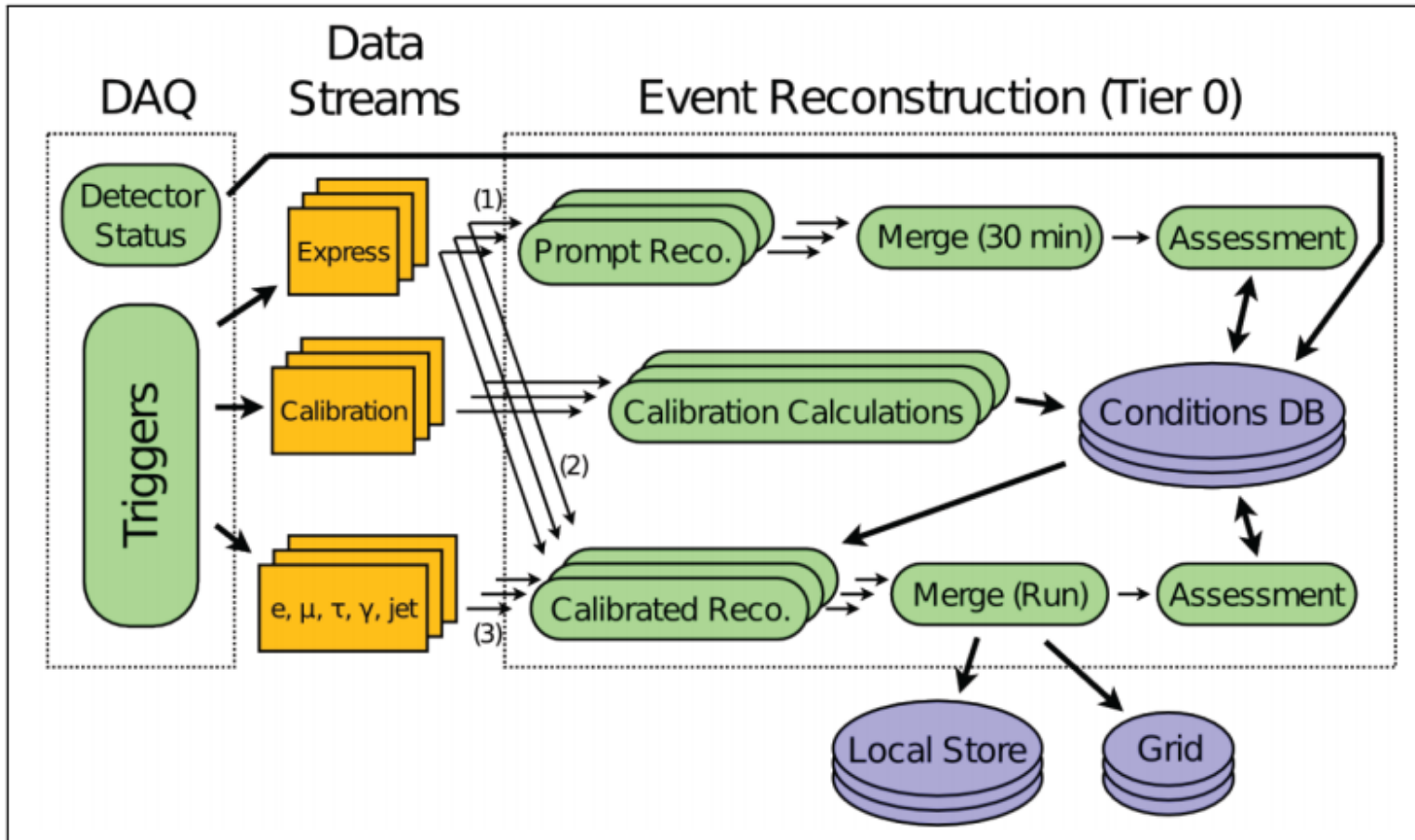
- ⇒ $F = 1$ for 5x5 crystal sum for the energy of unconverted photons;
- ⇒ c_i - intercalibration constants (π^0)
- ⇒ transparency correction with laser monitoring (LM)

ECAL cluster energies corrected using an MC trained multivariate regression

- ⇒ performed after individual crystal transparency correction and intercalibration
- ⇒ also provides per photon energy resolution estimate



Prompt reconstruction



In ATLAS we reconstruct the data ~ 36 hours after it is recorded. This time is used to derive updated calibrations from the data, that are needed in the reconstruction.

Once a year we reprocess all the data with updated software and calibrations.