

# Statistical Tools in Collider Experiments

## Multivariate analysis in high energy physics

### Lecture 4

Pauli Lectures - 09/02/2012

**Nicolas Chanon - ETH Zürich**



# Outline

1. Introduction
2. Multivariate methods
3. Optimization of MVA methods
4. Application of MVA methods in HEP
5. Understanding Tevatron and LHC results

# **Lecture 4. Application of multivariate methods in HEP**

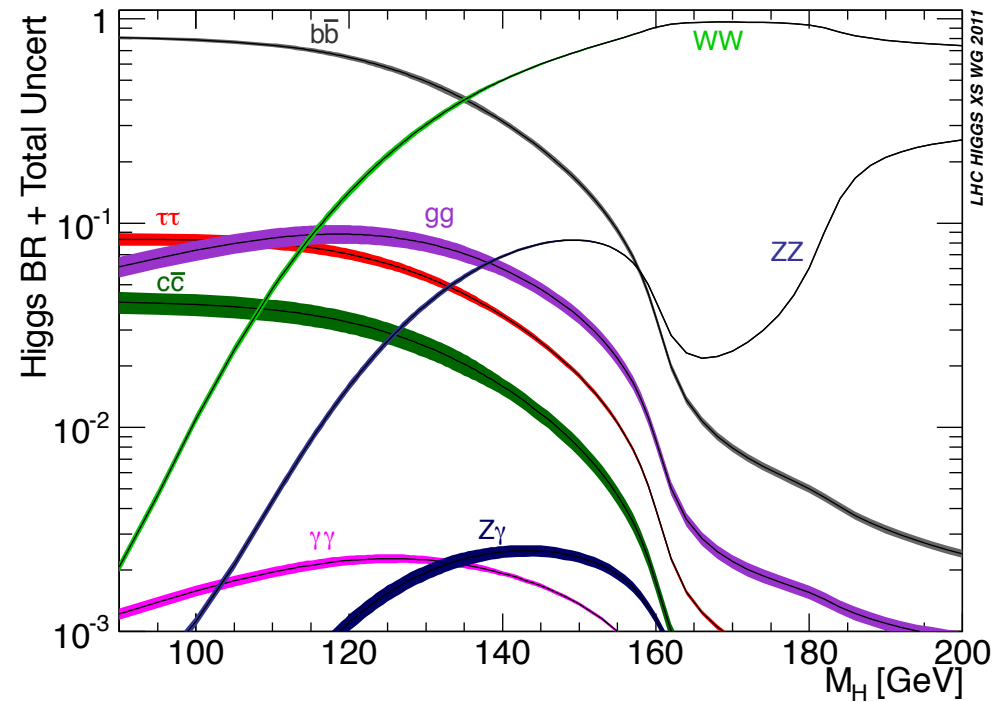
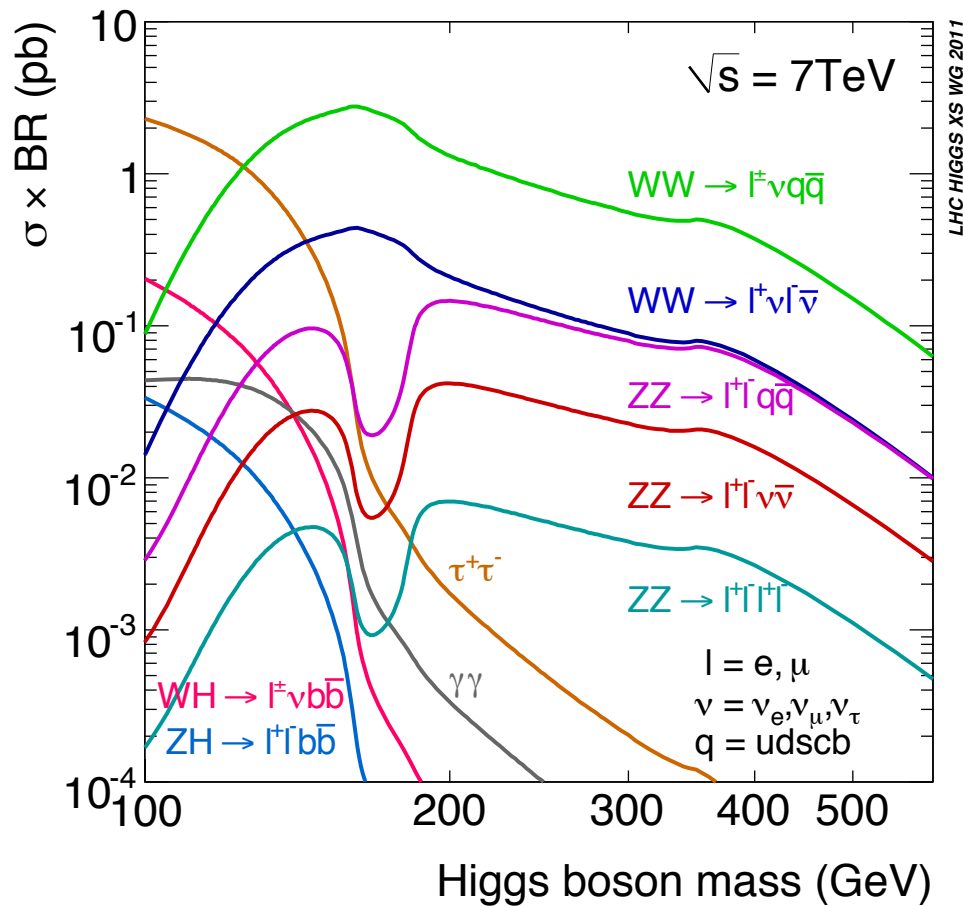
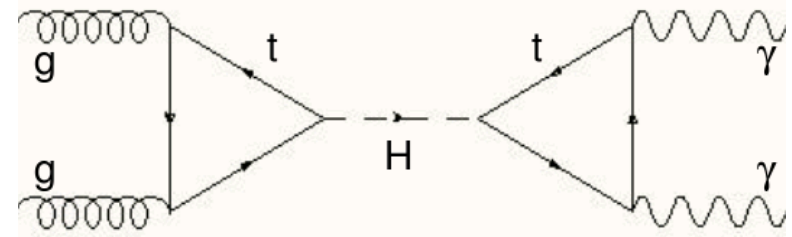
# Outline of the lecture

## How are applied multivariate methods in high energy physics ?

- We will take the example of  $H \rightarrow \gamma\gamma$  searches at LHC
- Details on the physics and experimental problems for this channel
- It will be the occasion to introduce the exercises

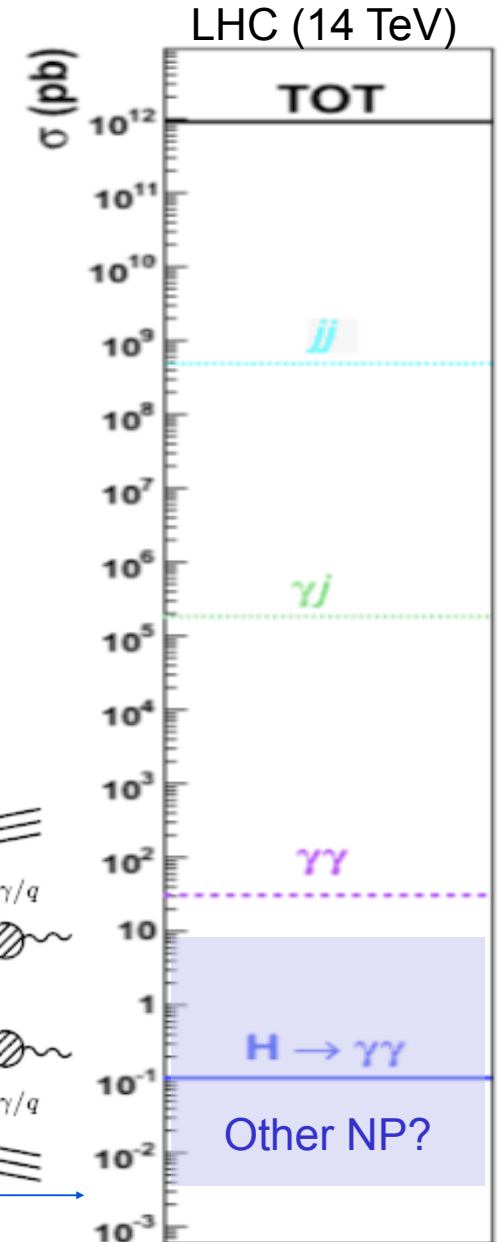
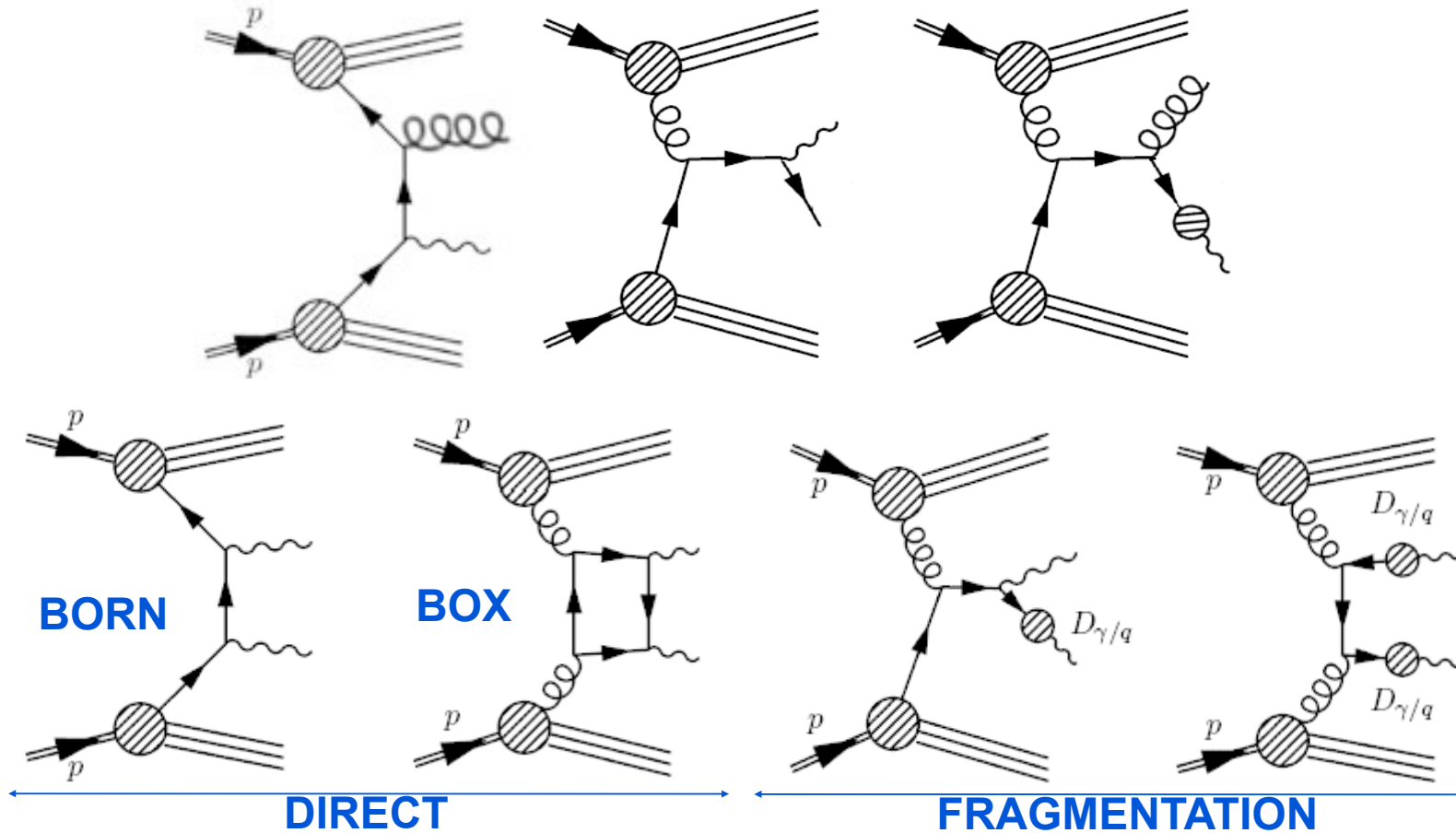
# H → γγ at LHC : signal

- H → γγ produced mainly via gluon fusion
- Branching ratio ~0.2%



# H → $\gamma\gamma$ at LHC : background

- Huge background over 9 order of magnitude in cross-section
- Dijet, gamma-jet processes with jets faking photons
- Diphoton continuum

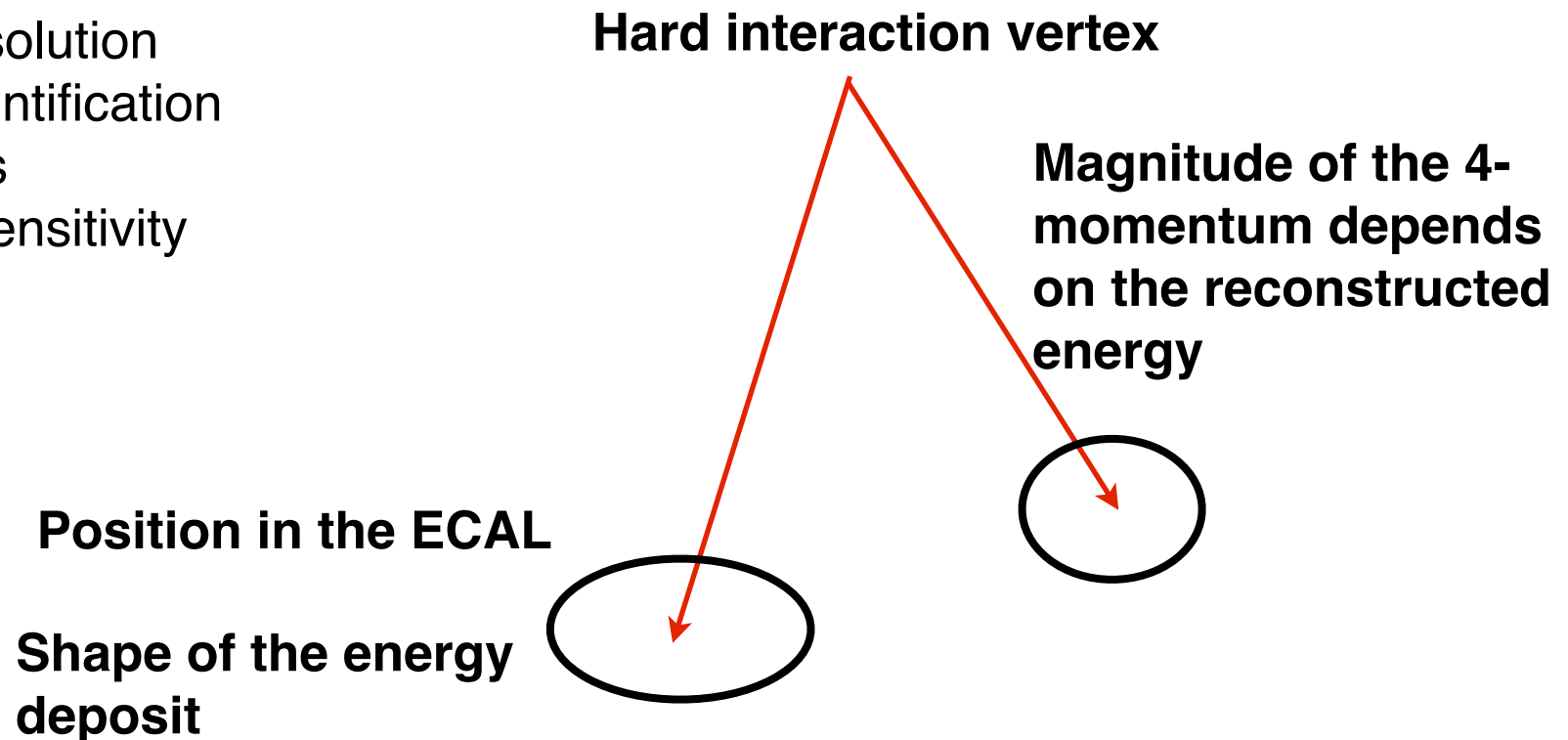


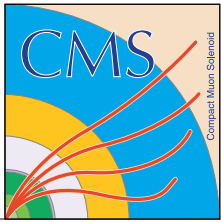
# $H \rightarrow \gamma\gamma$ at LHC : issues

- This channel suffers from small branching ratio and huge background.
- But it has the best sensitivity at low mass
- Reason : CMS and ATLAS have very good resolution on the  $\gamma\gamma$  invariant mass

## Main issues for $H \rightarrow \gamma\gamma$ :

- Vertex identification
- Energy resolution
- Photon identification
- Kinematics
- Analysis sensitivity





# CMS electromagnetic calorimeter (ECAL)

The **ECAL** is made of scintillating crystals of  $\text{PbWO}_4$  :

- **Barrel** : 36 “supermodules” with 1700 crystals each (coverage  $|\eta| < 1.48$ )

- **Endcaps** : 268 “supercrystals” with 25 crystals each (coverage  $1.48 < |\eta| < 3.0$ )

Furthermore, a **preshower** made of silicon strip sensors is located in front of the endcaps ( $1.65 < |\eta| < 2.6$ )

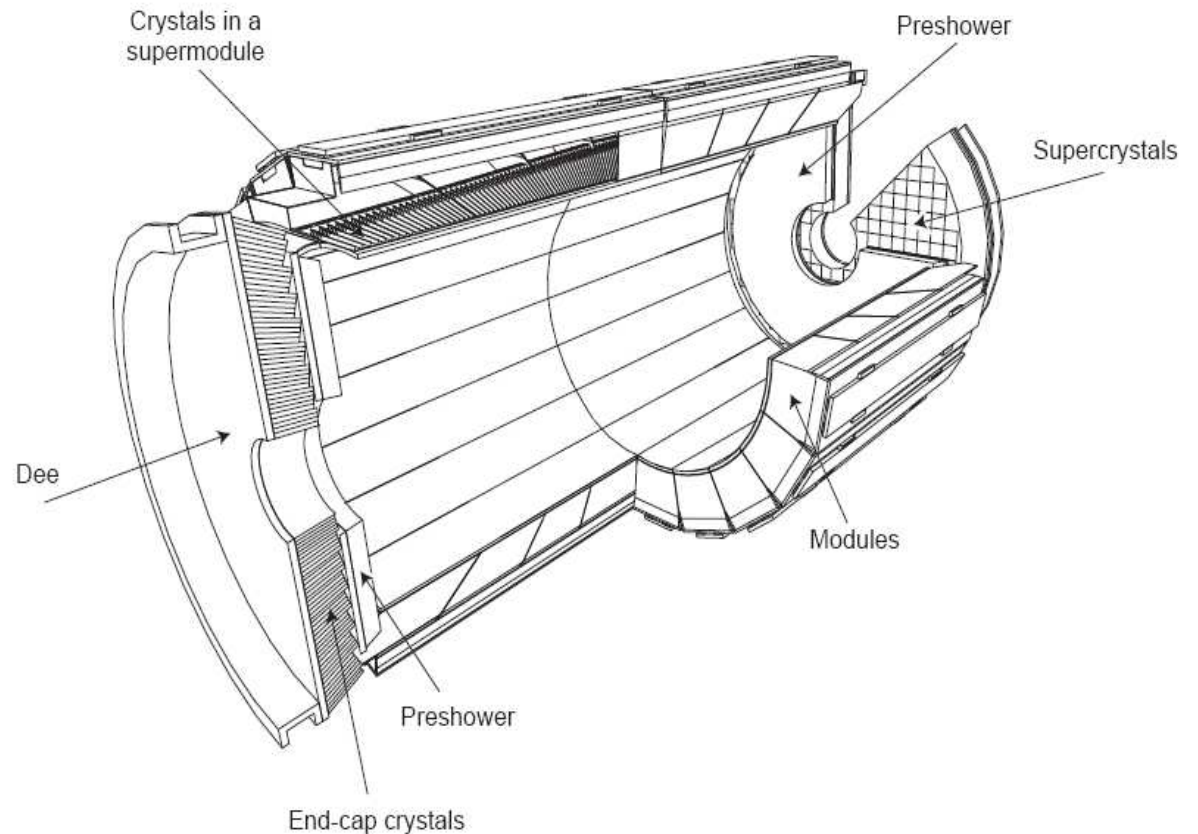
**Energy resolution** (measured in electron test beam) :

$$\frac{\sigma(E)}{E} = \frac{a}{\sqrt{E(\text{GeV})}} \oplus \frac{b}{E(\text{GeV})} \oplus c$$

$a = 2.8\%$  stochastic term

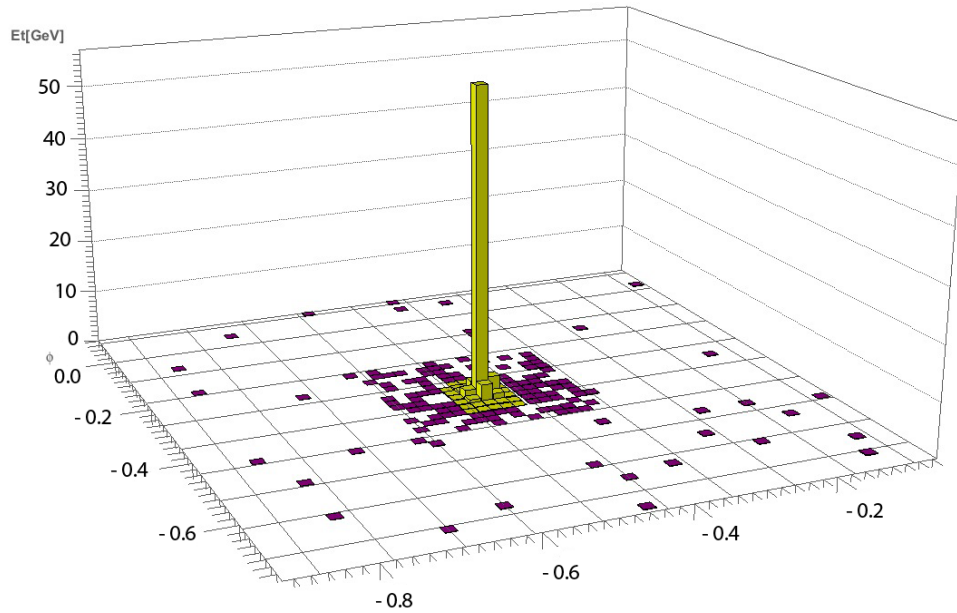
$b = 12\%$  noise term

$c = 0.3\%$  constant term





# Photon event display

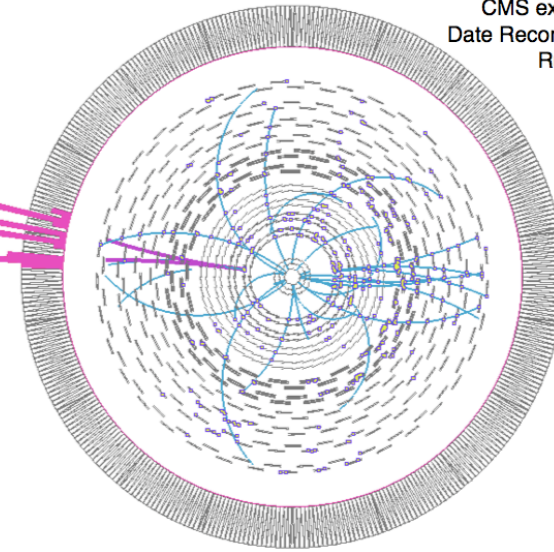


CMS Experiment at LHC, CERN  
Data recorded: Thu Jul 1 09:08:48 2010 CEST  
Run/Event: 139103 / 222480885



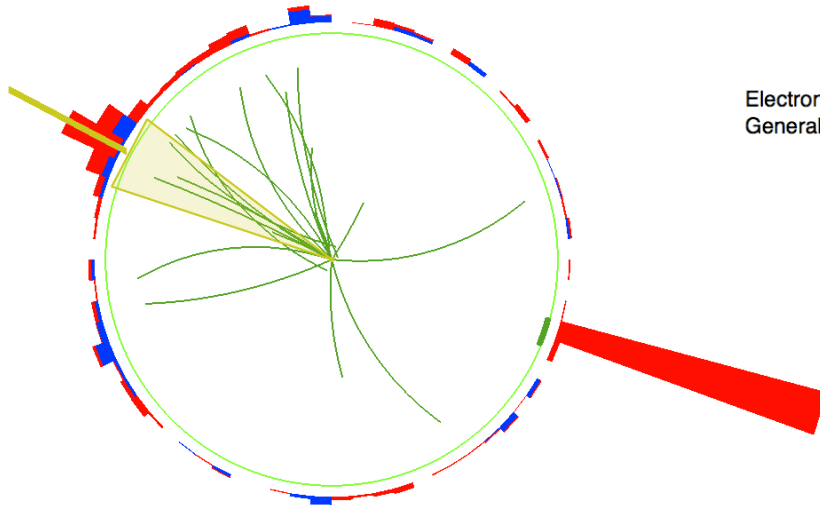
$E_{SC} = 21.45$  GeV

$E_{SC} = 11.92$  GeV



CMS experiment at the LHC, CERN  
Date Recorded: 2009-12-12 16:58 CET  
Run/Event: 124024/14608879  
Conversion candidate event  
 $\sqrt{s} = 900$  GeV

Electron tracks are shown in purple, and their superclusters in pink in the ECAL.  
General tracks are in blue and tracker clusters (silicon strips) are shown by small squares.

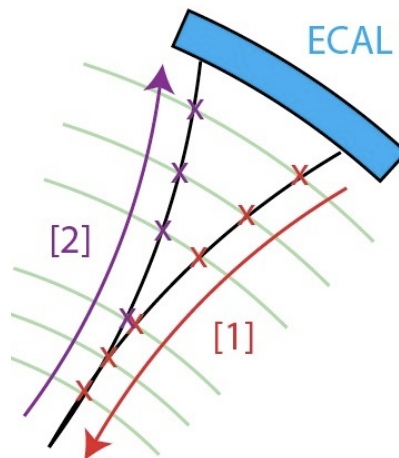




# Photon reconstruction

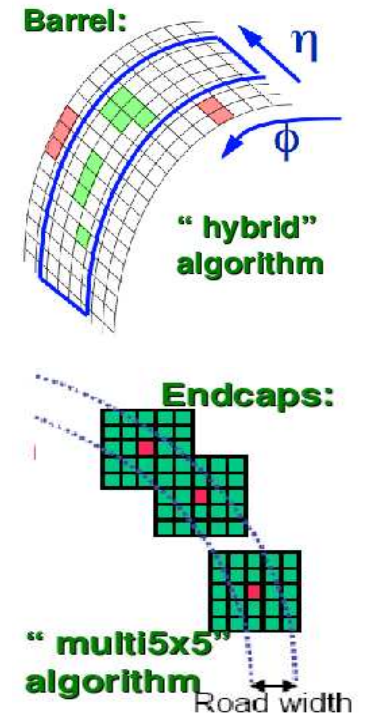
**Photons** are reconstructed with energy deposits in **ECAL** crystals

- **Barrel** : take advantage of the 3.8 T magnetic field which bends the charged particles trajectory (in case of a photon conversion)
- **Endcap** : merge contiguous  $5 \times 5$ -crystal matrices around the most energetic crystals



## Converted photons :

- Start from **energy deposits in ECAL**
- **Track finding** proceeds inward and outwards, taking into account electron energy loss by bremsstrahlung
- Select the  $e^+/e^-$  pair with the best vertex fit  $\chi^2$

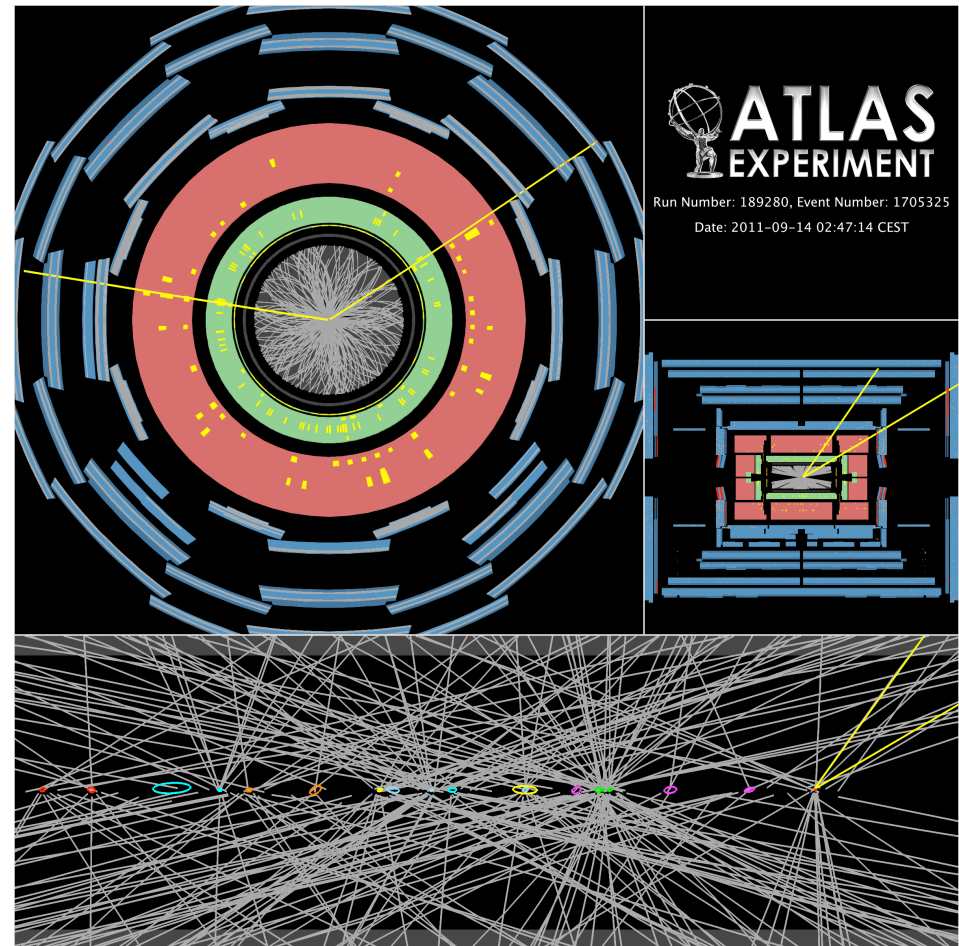


# $H \rightarrow \gamma\gamma$ at LHC : vertexing

- Up to  $\sim 20$  pile-up events per bunch crossing in 2011
- How to identify the hard interaction vertex ?
- Usual vertexing algorithm uses reconstructed tracks. Choose the vertex having the **highest sum pt squared**.

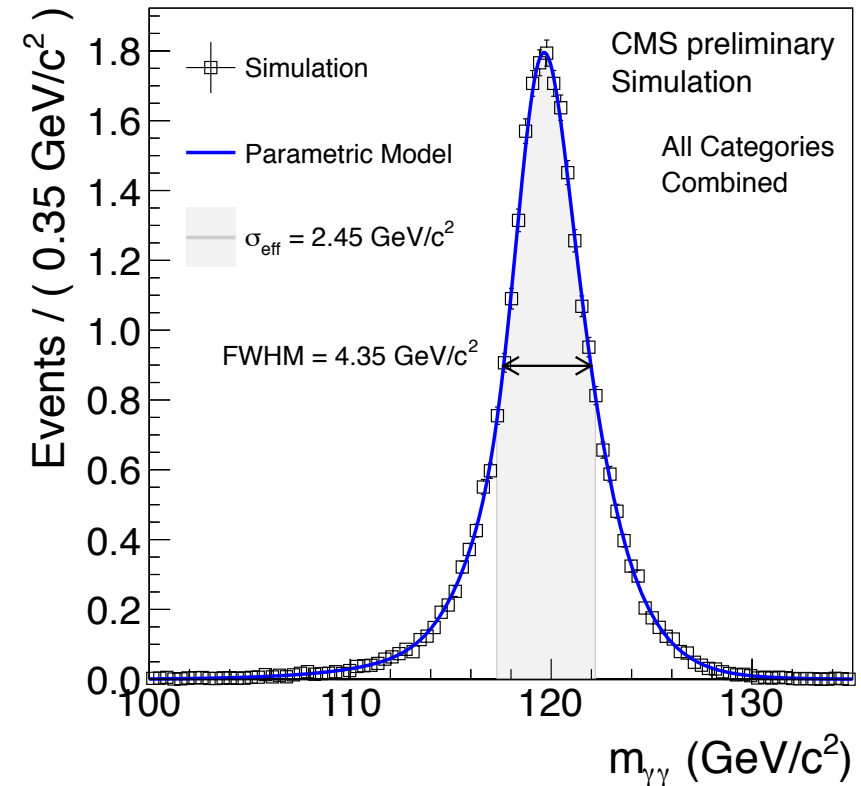
For  $H \rightarrow \gamma\gamma$  we have additional information :

- ATLAS : calorimeter pointing (photon conversion tracks pointing)
- CMS : multivariate method using tracks + diphoton kinematics, combined with conversion information



# H → $\gamma\gamma$ at LHC : energy resolution

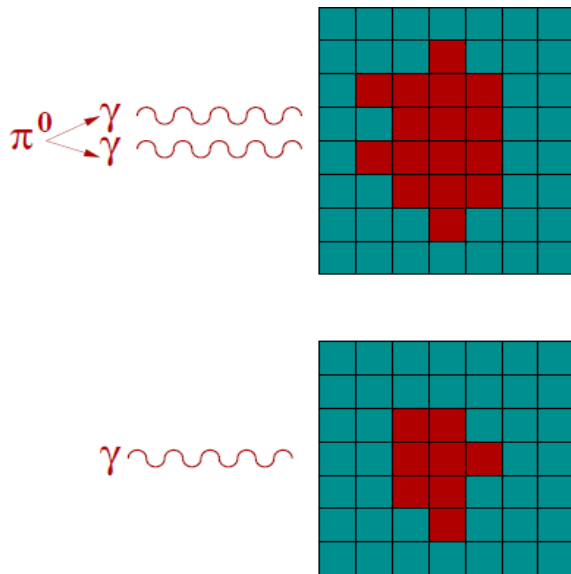
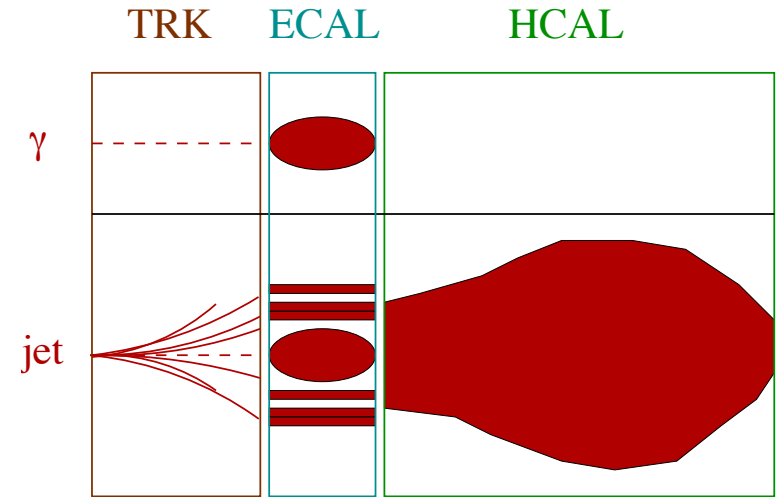
- Higgs natural width is zero from an experimental point of view in the  $\gamma\gamma$  channel
- So the **experimental width** is driven by **how well the photon energy is reconstructed** (once measured the position in the ECAL and the vertex found)
- CMS : PbWO<sub>4</sub> crystals calorimeter, subject to **loss of transparency**
- Clustering of the energy deposited is affected by the **tracker material** in front of the ECAL
- Corrections to get back the reconstructed energy to the energy at the vertex might not be optimal
- CMS : energy regression



# H → $\gamma\gamma$ at LHC : photon identification

## Why jets can fake photons ?

- Isolated boosted  $\pi^0$  decaying to 2 photons can be reconstructed in one single supercluster

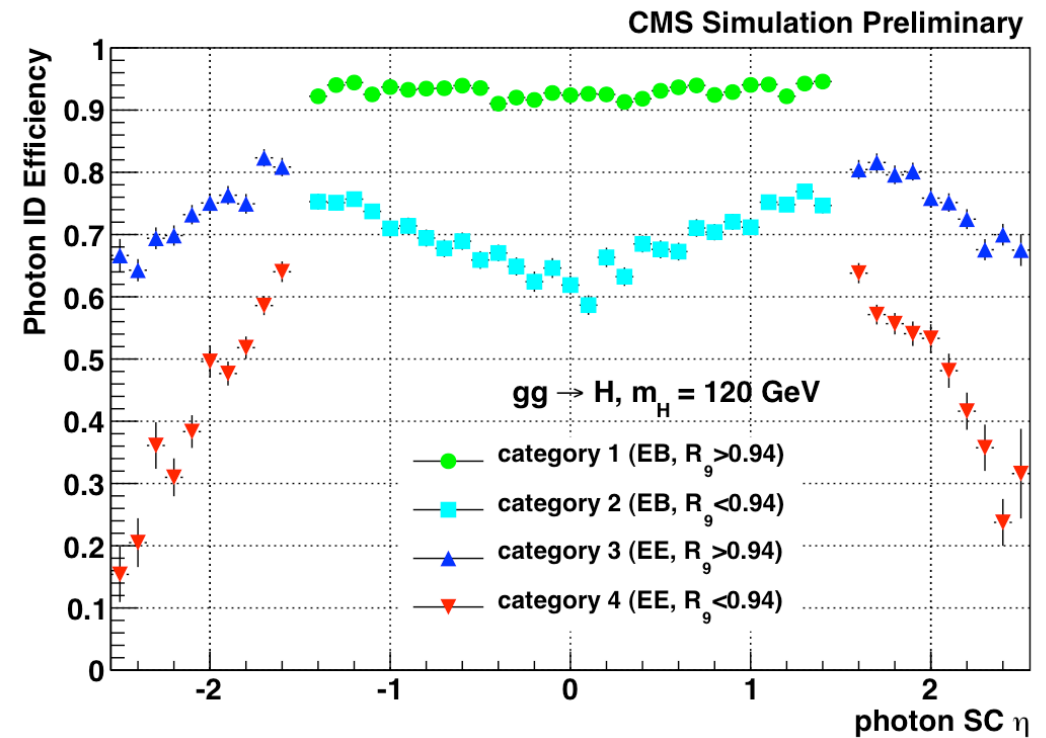
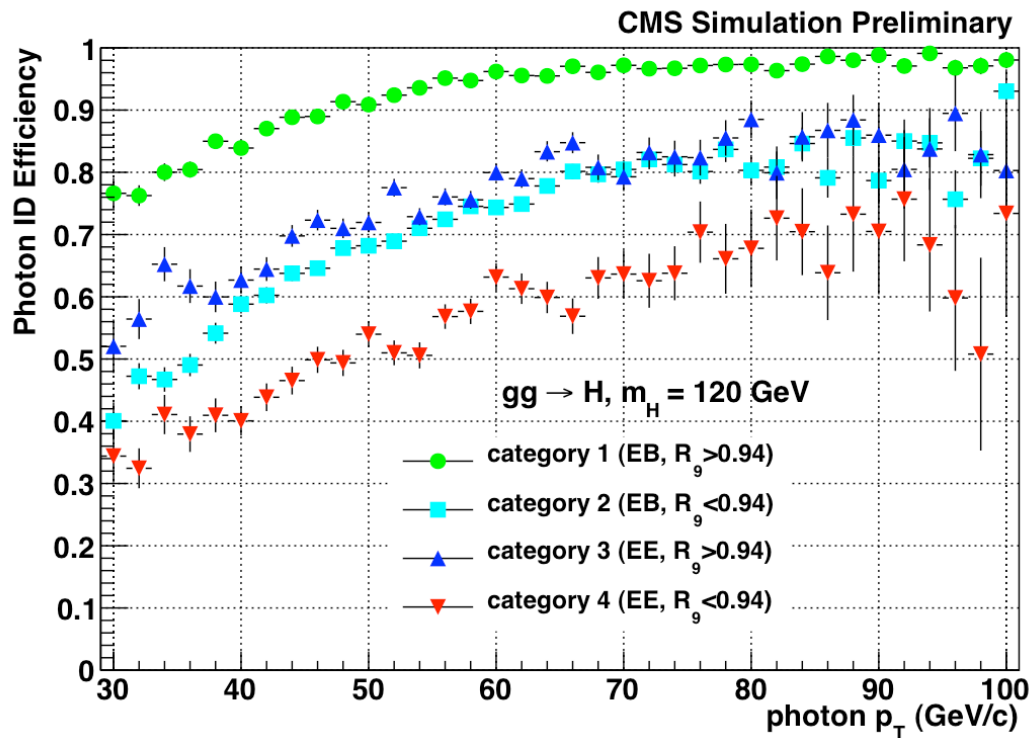


## Photon identification :

- **Electron rejection** : the energy deposit should not be matched to hits in the pixel detector
- The **transverse shape** of the energy deposits in ECAL should be compatible with a single photon shower
- **Isolation** : in a cone  $\Delta R < 0.4$  around the photon, use  $\sum E_T$  of energy deposits in **ECAL**, **HCAL** and  $\sum p_T$  of the charged particles measured in the **tracker**

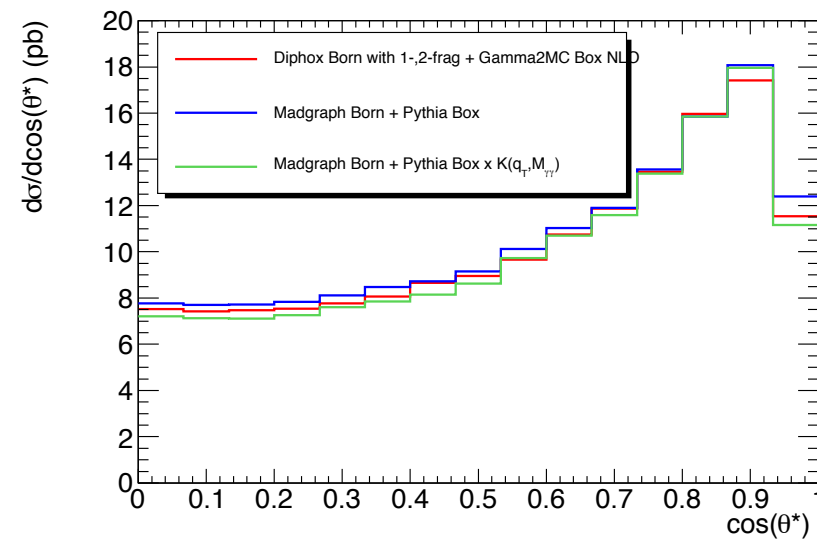
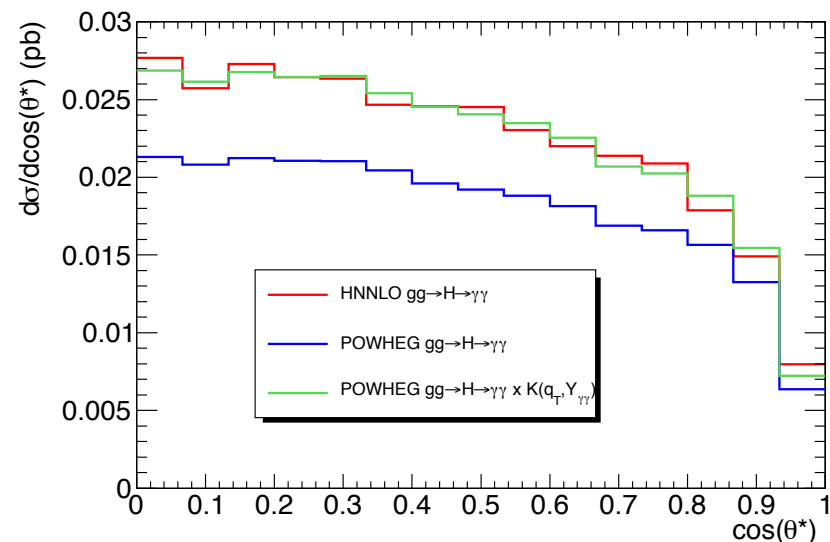
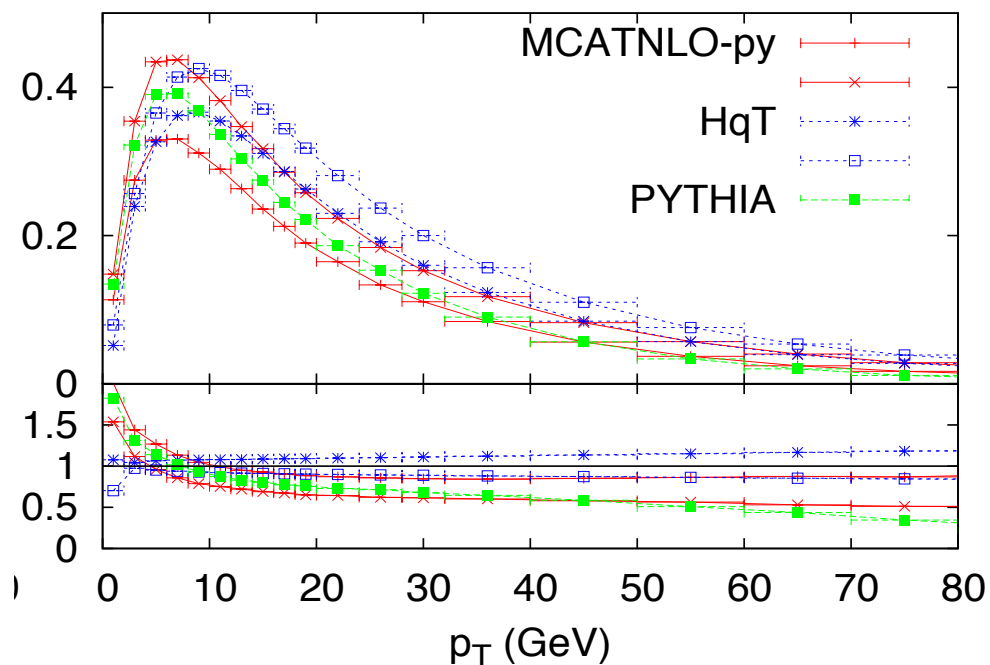
# H → $\gamma\gamma$ at LHC : photon identification

- In CMS photon identification is achieved using cuts on :
- **3 cluster shape variables** : H/E, transverse shape of the electromagnetic deposit,  $R9 = E_{3\times 3}/E_{\text{supercluster}}$
- **3 Isolation variables** : ECAL+HCAL+tracker in 0.3, 0.4 cones according to the wrong and right vertex hypothesis, Tracker isolation alone



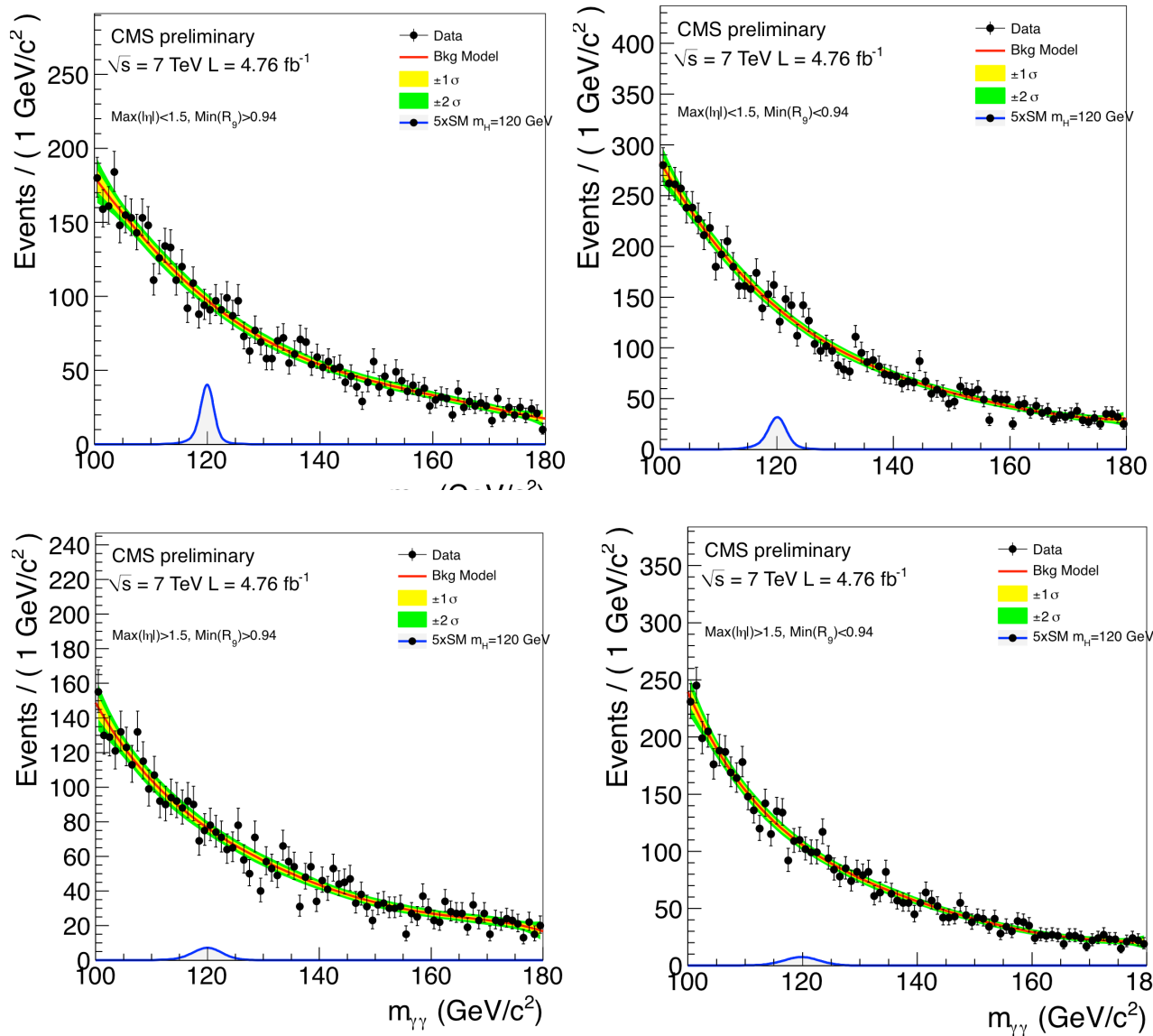
# H → γγ at LHC : kinematics

- **Photon p<sub>T</sub> threshold** : usually asymmetric, p<sub>T</sub> > 40, 30 or 25 GeV
- **cos(θ\*)** : can be discriminant in some kinematical regime
- **Diphoton p<sub>T</sub>** as discriminant variable : a myth for the gluon fusion



# H → $\gamma\gamma$ at LHC : diphoton categories

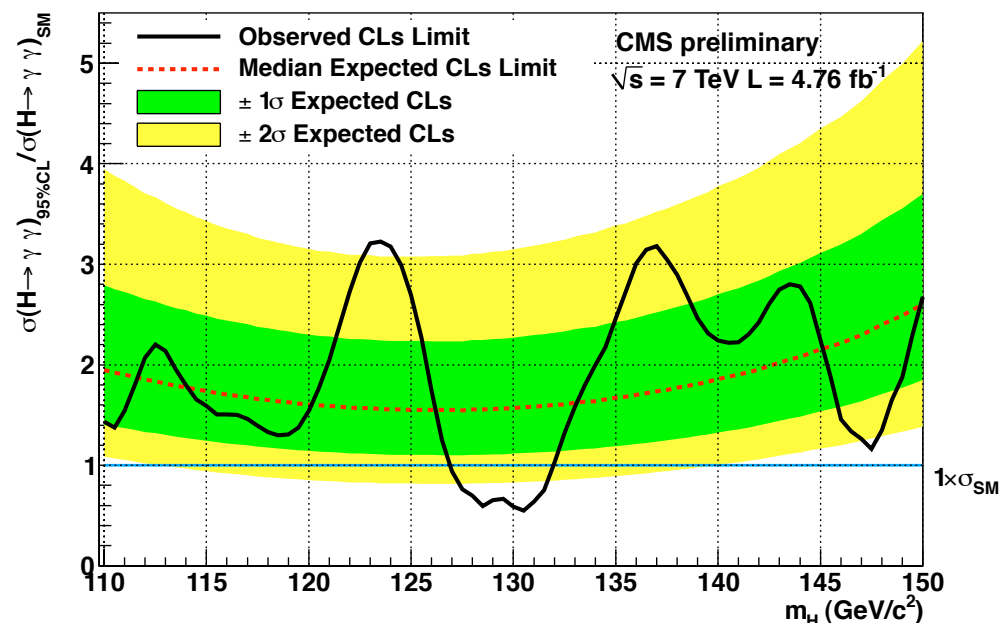
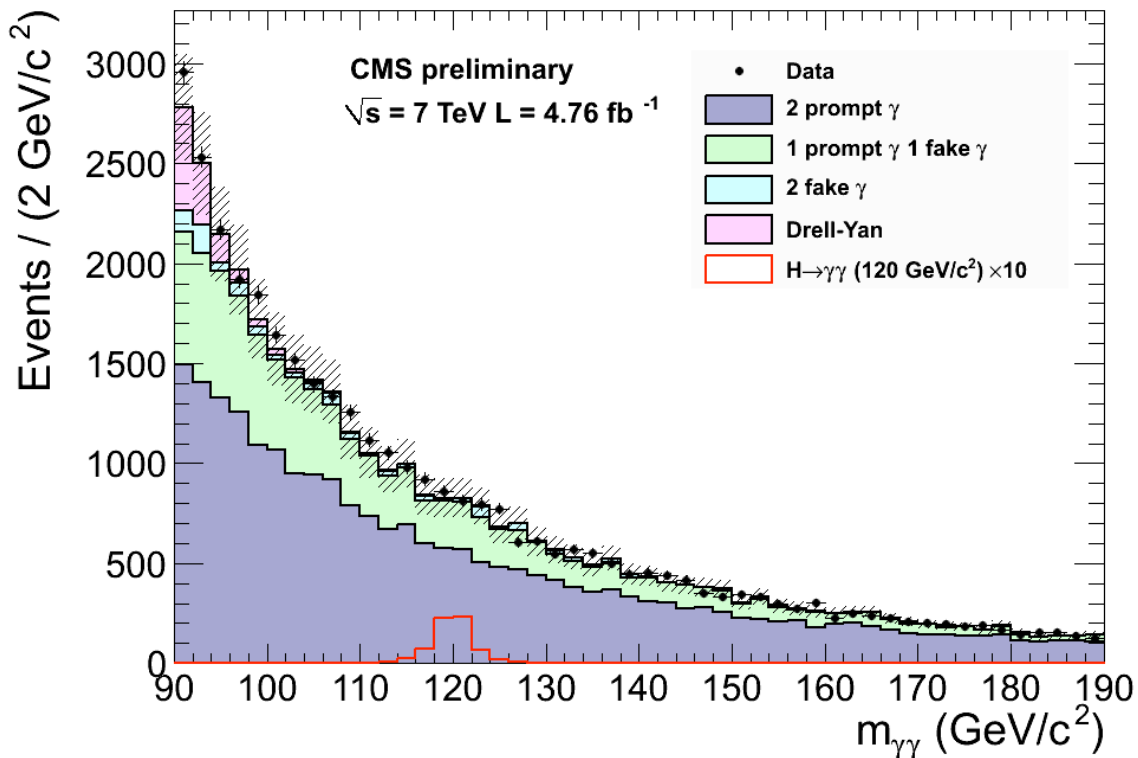
- CMS : 4 eta-r9 categories to improve mass resolution and s/b ratio
- ATLAS : 9 categories (eta / conversion / pt thrust)





# H $\rightarrow\gamma\gamma$ at LHC : analysis sensitivity

- Fit of the diphoton invariant mass distribution in data (how to choose the fit function ?)
- MC is not used to derive the sensitivity
- Unbinned CLs method



# Exercises

## - Inspired by H->2photons searches in CMS

- Can be downloaded from the lecture webpage
- Provide signal and background samples
- Variables : kinematics, photon identification, energy resolution

## 3 steps :

- **TMVA basics**
- **MVA application in the analysis**
- Sensitivity estimation (next lecture)

# Exercise I

## Installing ROOT

- Simplest option is probably to download the binaries (just unpack it)
- Do not forget to source bin/thisroot.sh

## Download the exercises on the webpage

- Pdf with instructions and questions
- The samples in the ROOT format

## Having a look to the samples :

- root -l Sample.root

## Running TMVA

- Go to the directory tmva
- Classifier training can be launched using TMVAClassification.C
- Once the classifiers are trained, one can investigate with TMVAGui.C
- One can also have a look to the training output : TMVA.root

# Samples

## Samples provided were generated using Pythia :

- **gg**→**H**→**γγ** **mH=120 GeV** (100kevt generated) - forget other production mechanisms
- **γγ** **Born** (1Mevt)
- **γγ** **Box** (1Mevt)
- **γ+Jet** (20Mevt - lack of statistics)
- Dijet background was not generated (1000x more events would have been needed due to the small jet→γ misidentification rate)

## Experiment simulation

- Events have been passed into a (home-made) program which emulates the experiment
- Energy smearing due to finite detector effects
- Energy deposits variables
- Important correlations taken into account

# Variables

## List of the variables :

### Diphoton variables :

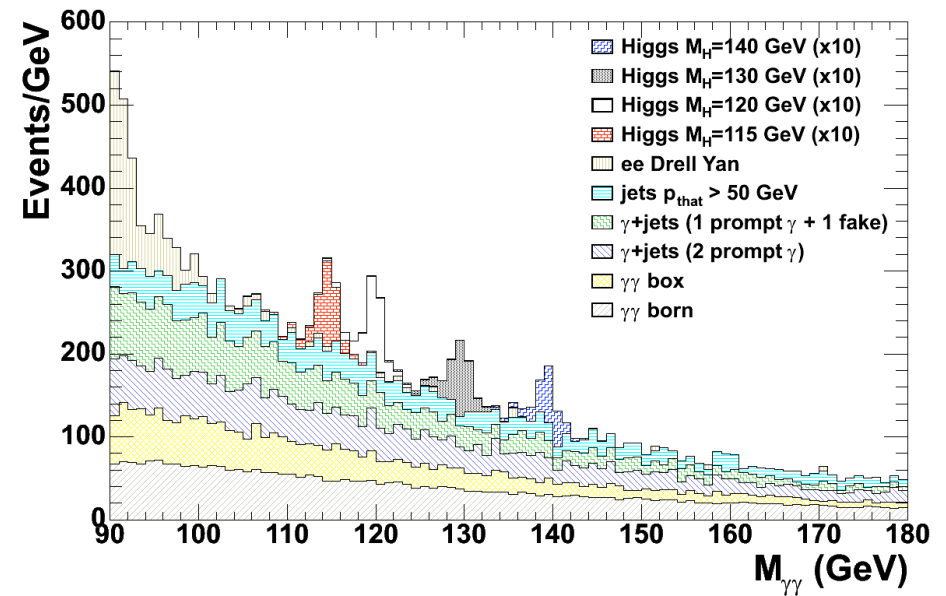
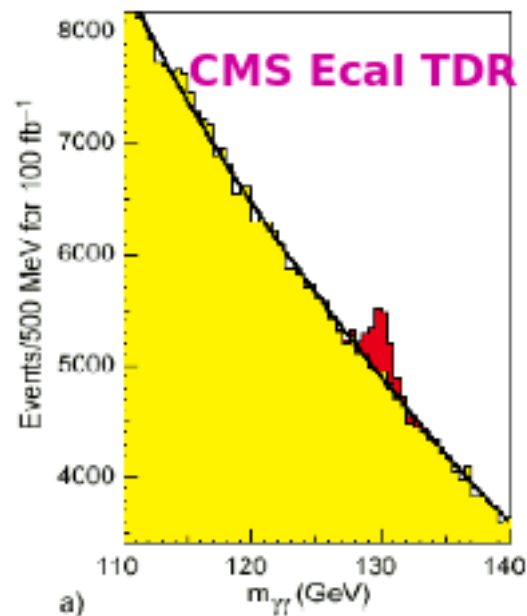
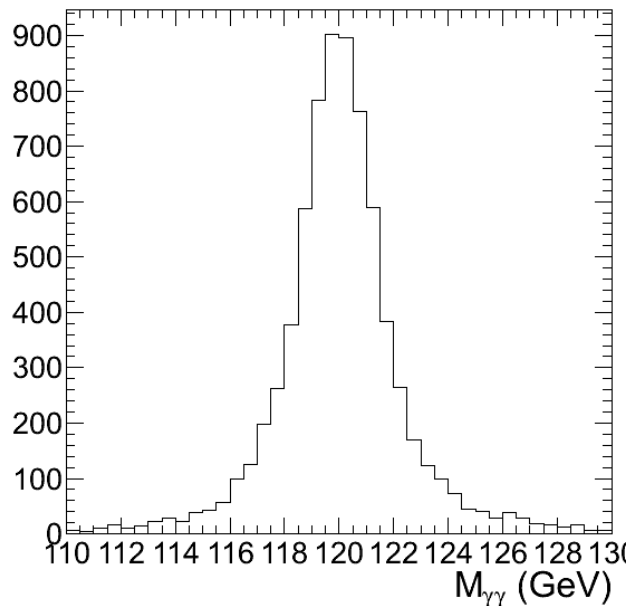
- Invariant mass
- pT of the diphoton system
- $\cos(\theta^*)$

### Variables for the highest pt and second highest pt photons :

- 4-Momentum
- Eta
- Cluster shape variables
- Isolation variables
- pdgId : photon or meson ?

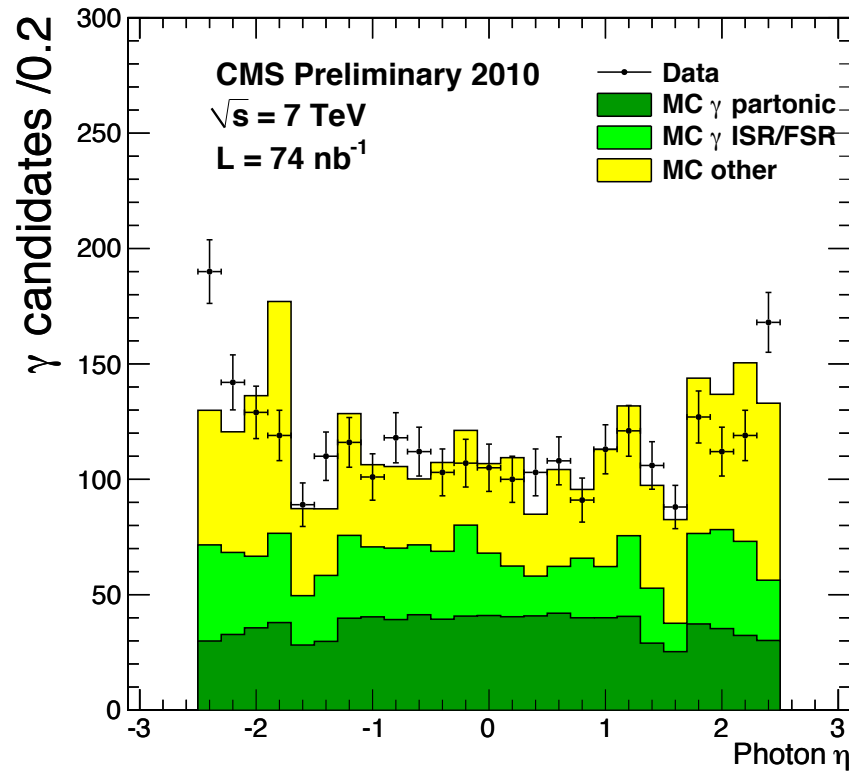
# Invariant mass

- In the exercise, the diphoton mass resolution is different from the one we get in reality, but the order of magnitude is the right one
- Look for a sharp peak in a steeply falling background
- After photon identification, the jet-jet and gamma+jet background is much reduced



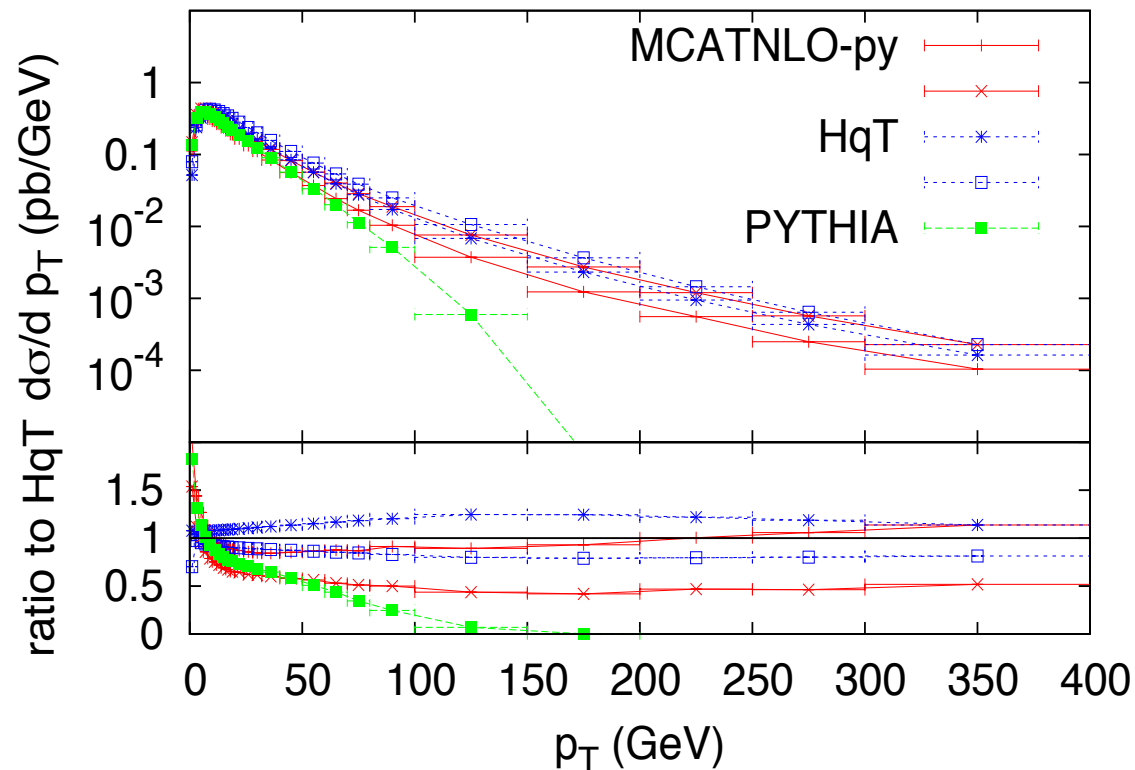
# Photon kinematics

- In the exercise, the photon  $P_t$  is smeared
- The reconstruction efficiency ( $\eta$ -dependent) is not taken into account. This gives more photons in the barrel-endcap transition region than expected experimentally



# qT, $\cos(\theta^*)$

- The diphoton transverse momentum is only LO+LL from Pythia here
- Can be used for the purpose of demonstration, but the discriminating power is much reduced in reality
- $\cos(\theta^*)$  can also be used, but it is difficult to make it very discriminant with the trigger thresholds actually used

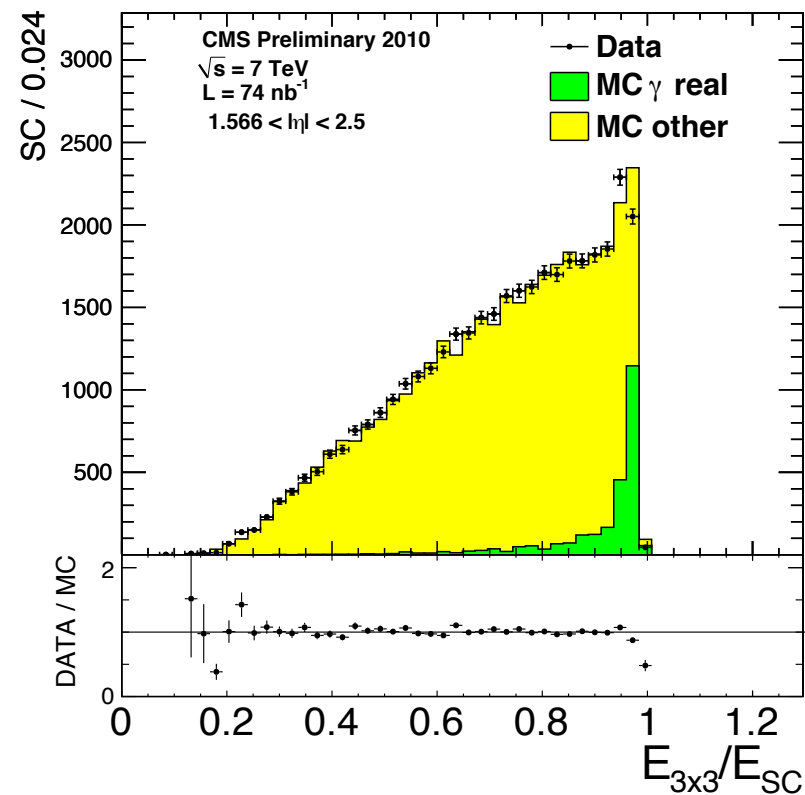
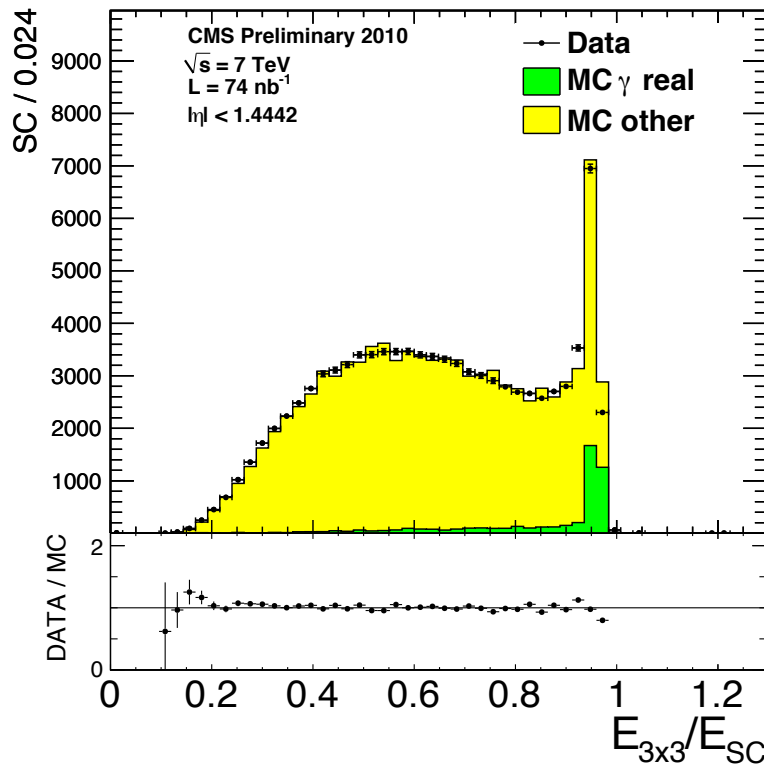




# R9 cluster shape variable

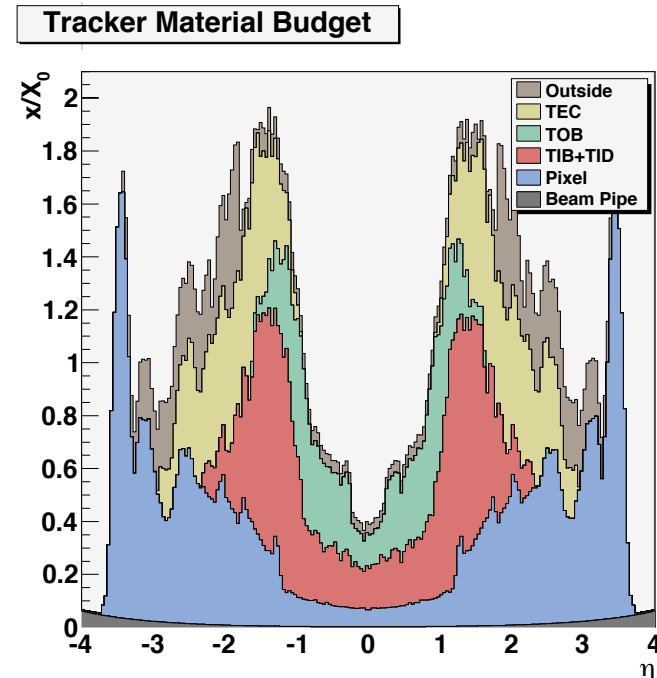
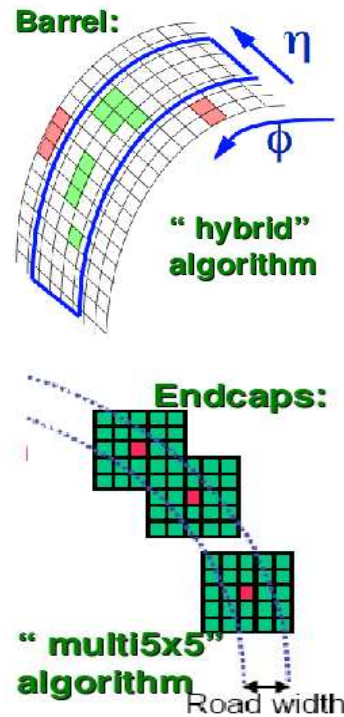
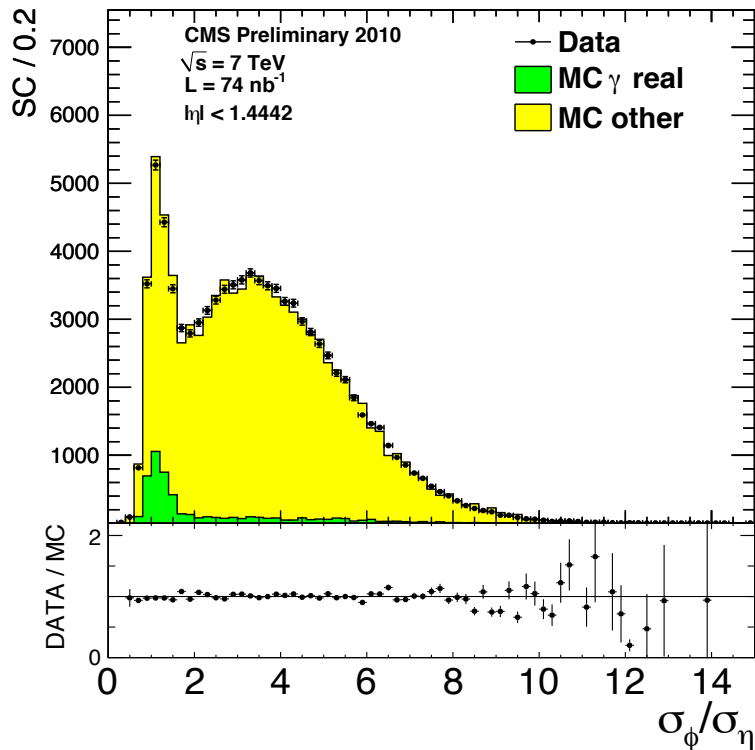
## R9 = E<sub>3x3</sub>/E<sub>supercluster</sub>

- High R9 : unconverted photon, very good energy resolution
- Low R9 : converted photon, poor energy resolution
- $\pi^0$  also located at low R9
- R9 is very  $\eta$ -dependent



# Brem cluster shape variable

- A photon energy deposit is broader in  $\phi$  than in  $\eta$ , due to the magnetic field which bent the conversion trajectory around the z axis
- The  $\eta$ -width is broader for a  $\pi^0$  than a photon
- The clustering algorithm is affected by the material in front of the ECAL : strongly  $\eta$ -dependant
- The photon energy resolution is strongly dependent on  $\sigma_\phi/\sigma_\eta$



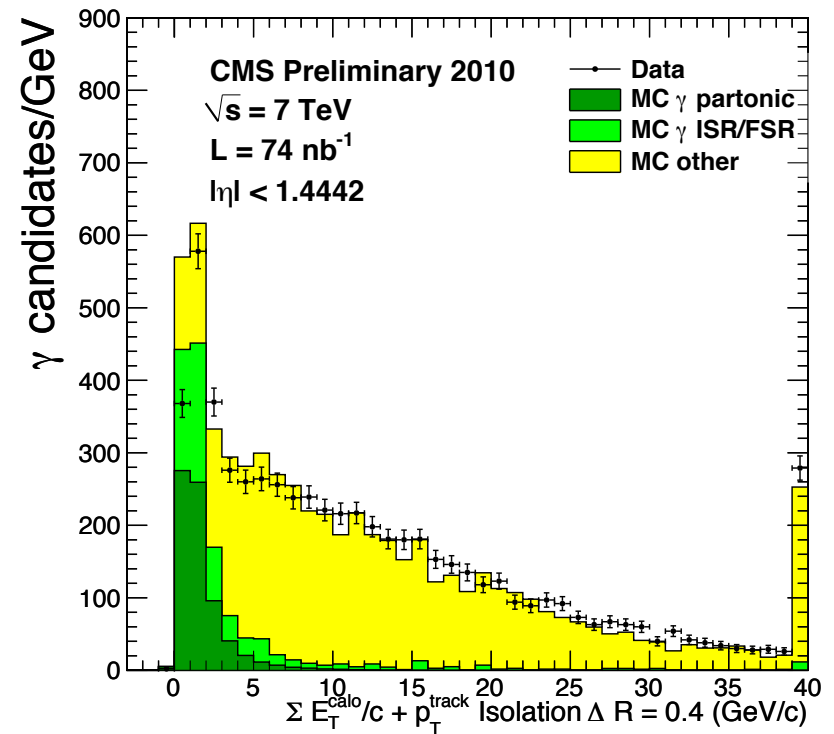
# Isolation energy

Isolation energy is defined in a  $\Delta R$  cone of 0.3 or 0.4 around the photon

- Tracker isolation : Sum  $p_T$  of the tracks reconstructed inside the cone
- ECAL, HCAL isolation : Sum  $E_t$  of the energy deposits inside the cone

Isolation energy is coming from :

- Underlying event
- Pile-up
- QCD/QED radiation
  
- Prompt photons are isolated
- Neutral mesons within jets are less isolated
  
- In the exercise, 0.3 and 0.4 cones are used



# Possible multivariate methods

To improve the  $H \rightarrow \gamma\gamma$  analysis sensitivity, one can use several multi-variate methods :

## Vertexing MVA

- Used in CMS results since Summer 2011
- In the exercises, no pile-up. The vertex is assumed to be correctly reconstructed.

## Energy regression

- Used in CMS results since Dec 13
- Can be tried with the samples provided in the exercises

## Photon identification with MVA

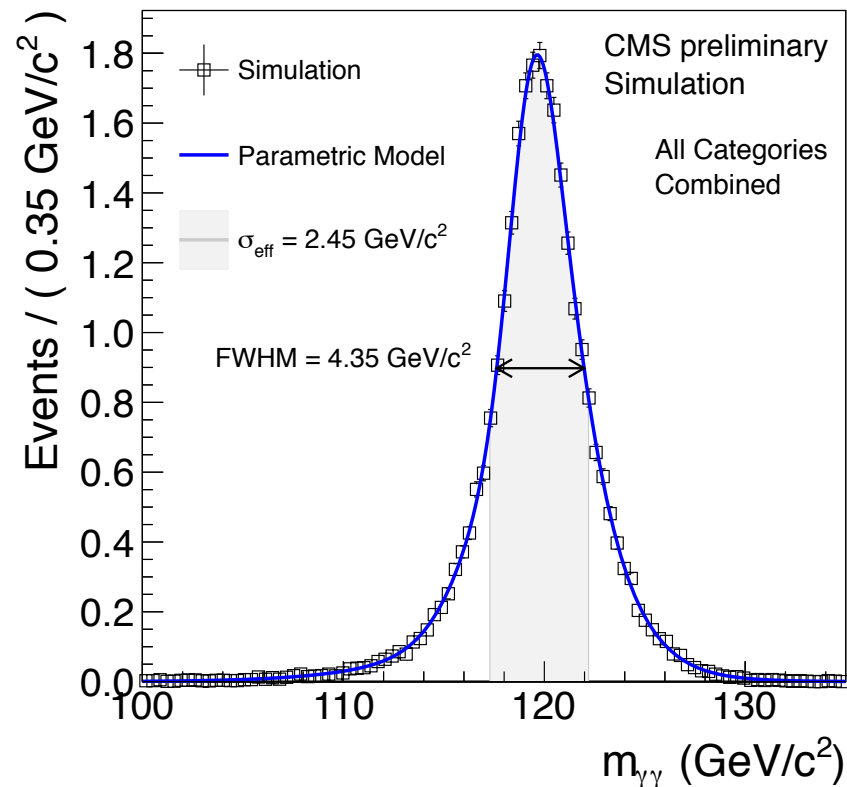
- Photon identification performed with rectangular cuts for the moment
- Can be tried in the exercises

## Kinematics MVA

- Only the invariant mass is used for the moment - no MVA
- Can be tried in the exercises

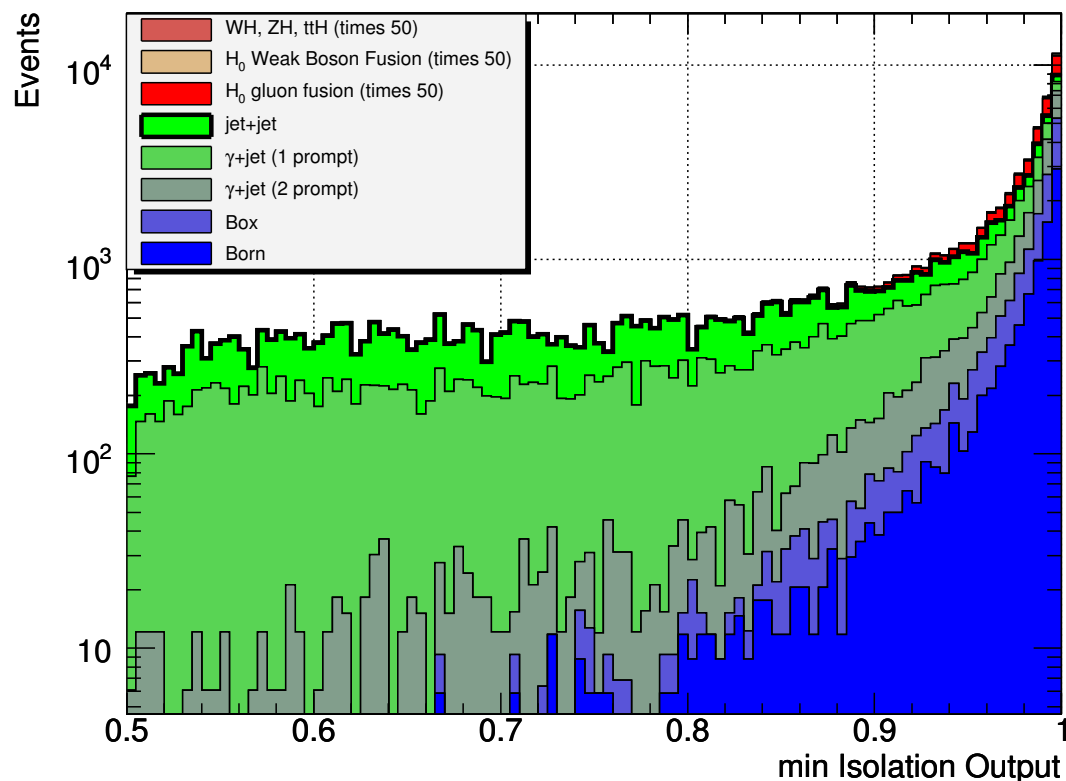
# Energy regression

- Last CMS results (Dec 13) were using a multivariate technique to improve the photon energy resolution
- Perform a regression from the reconstructed energy to the generated energy, using many geometrical variables and cluster shape variables
- This improved a lot the invariant mass resolution



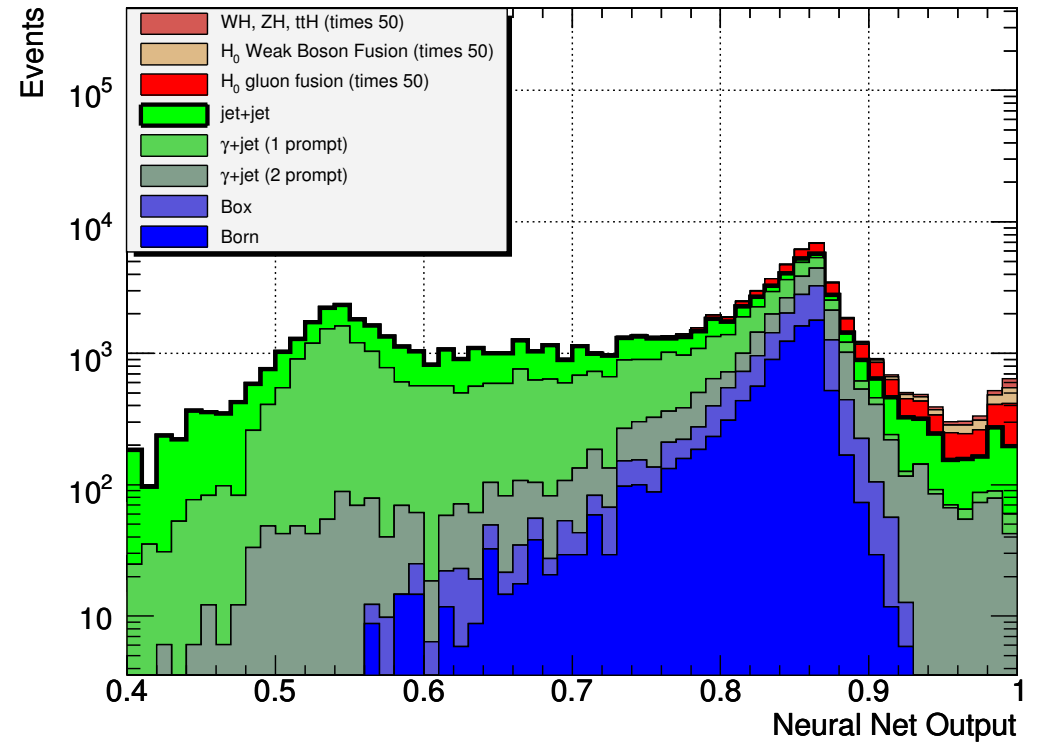
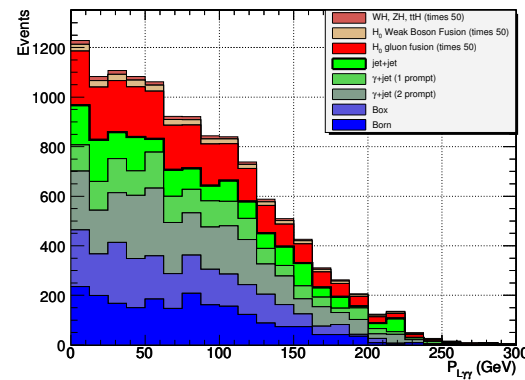
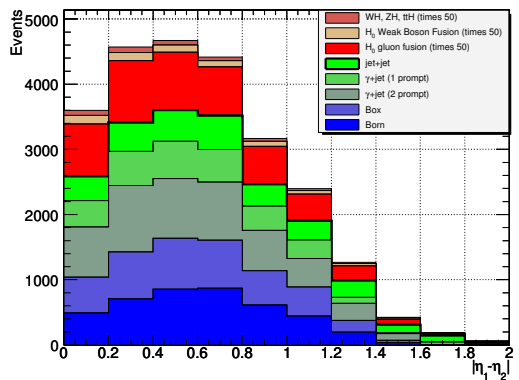
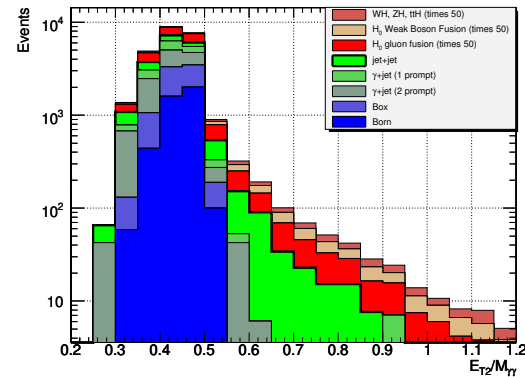
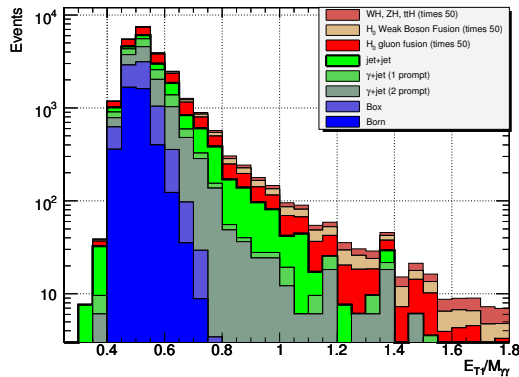
# Photon identification MVA

- In CMS Physics TDR vol. II, a photon identification NN was used :
- Uses ECAL, HCAL, Tracker isolation
- And R9 cluster shape variable



# Kinematics MVA

- In **CMS Physics TDR vol. II**, a **global NN** was used :
- ET/M of the two photons, pseudo-rapidity difference, pT of the diphoton system
- The two outputs of the NNisol



# Kinematics MVA

- **ATLAS** was also foreseeing a multivariate analysis for  $H \rightarrow \gamma\gamma$  using kinematics
- No classifier, use rather invariant mass, diphoton  $p_T$  and  $\cos(\theta^*)$  in a **3-dimensional likelihood** (as 3 different channels)
- Also considered 0, 1 and 2-jet bin cases

Fit variables	Categories	Higgs boson mass fixed		Higgs boson mass floating	
		$\langle \Delta \ln \mathcal{L} \rangle$	Significance [ $\sigma$ ]	$\langle \Delta \ln \mathcal{L} \rangle$	Significance [ $\sigma$ ]
$m_{\gamma\gamma}$	–	$2.67 \pm 0.04$	$2.31 \pm 0.02$	$3.54 \pm 0.05$	$1.44 \pm 0.02$
$m_{\gamma\gamma}$	$\eta$	$3.18 \pm 0.05$	$2.52 \pm 0.02$	–	–
$m_{\gamma\gamma}$	$\eta$ , Conversions	$3.32 \pm 0.05$	$2.58 \pm 0.02$	–	–
$m_{\gamma\gamma}$	$\eta$ , Conversions, Jets	$5.99 \pm 0.07$	$3.46 \pm 0.02$	$6.66 \pm 0.07$	$2.64 \pm 0.02$
$m_{\gamma\gamma},  \cos\theta^* $	$\eta$ , Conversions, Jets	$7.33 \pm 0.08$	$3.83 \pm 0.02$	–	–
$m_{\gamma\gamma}, P_{T,H}$	$\eta$ , Conversions, Jets	$7.03 \pm 0.08$	$3.75 \pm 0.02$	–	–
$m_{\gamma\gamma}, P_{T,H},  \cos\theta^* $	$\eta$ , Conversions, Jets	$8.49 \pm 0.08$	$4.12 \pm 0.02$	$9.25 \pm 0.09$	–

