

¹ **Properties of the Higgs-like state**
² **around 125 GeV in its decay to two**
³ **photons at the CMS experiment at**
⁴ **the LHC**

⁵ Matthew Kenzie
⁶ Imperial College London

⁷ A dissertation submitted to Imperial College London
⁸ for the degree of Doctor of Philosophy

DRAFT

10 **Abstract**

11 Some stuff here

DRAFT

DRAFT

Declaration

This dissertation is the result of my own work, except where explicit reference is made to the work of others, and has not been submitted for another qualification to this or any other university. Figures from other sources are labelled as CMS and the source is referenced in the figure captions. This dissertation does not exceed the word limit for the respective Degree Committee.

Matthew Kenzie

DRAFT

20 **Acknowledgements**

21 Acknows

DRAFT

DRAFT

Preface

22

23

Blah blah blah

DRAFT

Contents

24	1. Introduction and Theory	3
26	1.1. Description	3
27	1.2. The Standard Model	3
28	2. The CMS experiment	5
29	2.1. The LHC	5
30	2.2. The CMS detector	6
31	2.2.1. Tracking system	8
32	2.2.2. Electromagnetic calorimeter	10
33	2.2.3. Hadronic calorimeter	18
34	2.2.4. Muons	19
35	2.3. Particle flow and jets	19
36	2.4. Isolation	21
37	2.5. Pileup	22
38	3. Common analysis components	23
39	3.1. Boosted Decision Trees	24
40	3.2. Data samples and triggering	26
41	3.2.1. Monte Carlo	27
42	3.2.2. Pileup and beamspot reweighting	28
43	3.3. Energy measurement of photons	29
44	3.3.1. Correcting for residual discrepancies between data and Monte Carlo	33
45	3.4. Vertex reconstruction	36
46	3.4.1. Estimating the per-event probability that the correct vertex is chosen	38
47	3.5. Event preselection	43
48	3.6. Using Z decays for validation and efficiency measurements	44

49	4. Selection and Categorisation	47
50	4.1. Event selection	47
51	4.1.1. Selection using cuts in categories	48
52	4.1.2. Photon ID MVA	49
53	4.1.3. Diphoton event level MVA	52
54	4.2. Event categorisation	57
55	4.2.1. Exclusive mode tagging	58
56	4.2.2. Inclusive mode categorisation in the cut based analysis	64
57	4.2.3. Inclusive mode categorisation and VBF dijet categorisation in the mass factorised Multivariate Analysis (MVA) analysis	65
58	4.2.4. Inclusive mode categorisation in the sideband MVA analysis	66
59	4.2.5. Event categorisation summary	67
61	5. Analysis and Results	71
62	5.1. Description	71
63	5.2. Signal modelling	72
64	5.2.1. Mass factorised analysis	72
65	5.2.2. Sideband analysis	72
66	5.3. Background modelling	72
67	5.3.1. Mass factorised analysis	72
68	5.3.2. Sideband analysis	72
69	5.4. Systematic Uncertainties	72
70	5.5. Statistics	72
71	5.5.1. Use of the Likelihood function as a test statistic	72
72	5.6. Results of the mass factorised analysis	72
73	5.7. Results of the sideband analysis	72
74	A. Photon ID BDT input variables validation in $Z \rightarrow e^+e^-$	73
75	B. Diphoton BDT input variables validation in $Z \rightarrow e^+e^-$	75
76	Bibliography	79
77	List of Figures	81
78	List of Tables	87
79	List of Acronyms	88

“For my father.”

DRAFT

Chapter 1.

⁸² Introduction and Theory

⁸³ “All science is either physics or stamp collecting.”

⁸⁴ — Ernest Rutherford

⁸⁵ 1.1. Description

⁸⁶ Motivation for Higgs searches and the theory of the Standard Model. In particular, the
⁸⁷ electroweak sector and symmetry breaking.

⁸⁸ **15 pages**

⁸⁹ Some things. Ref example [1–3].

⁹⁰ 1.2. The Standard Model

⁹¹ The **sm!** (**sm!**) is very successful.

$$m_{\gamma\gamma} = 2\sqrt{E_{\gamma 1}E_{\gamma 2}(1 - \cos(\alpha))} \quad (1.1)$$

DRAFT

Chapter 2.

⁹² The CMS experiment

⁹³ “*You need something to open up a new door,
⁹⁴ To show you something you seen before,
⁹⁵ But overlooked a hundred times or more.”* — Bob Dylan

⁹⁷ 2.1. The LHC

⁹⁸ The Large Hadron Collider ([LHC](#)) is an octagonal 27km ring (large) proton-proton
⁹⁹ (hadron) particle collider. Using a multistage acceleration process two beams of protons
¹⁰⁰ are circulated in opposite directions at a centre-of-mass energy, $\sqrt{s}=7$ TeV (8 TeV) for
¹⁰¹ data collection in 2011 (2012). The beams of protons are accelerated and circulated by
¹⁰² electric and magnetic fields respectively. Further precision magnetic fields can control
¹⁰³ the position and intensity of the beams. There are four points around the ring where the
¹⁰⁴ beams can be forced to intersect producing high energy proton-proton collisions. Particle
¹⁰⁵ detectors are constructed around these points such that the collision can be reconstructed
¹⁰⁶ with the purpose of measuring physical properties and processes, calibrating the detectors
¹⁰⁷ with already known processes and searching for new physics. The remainder of this
¹⁰⁸ chapter concentrates on a description of one of these detectors, the Compact Muon
¹⁰⁹ Solenoid ([CMS](#)), which the author has worked on.

110 2.2. The CMS detector

111 The CMS detector, pictured in Fig. 2.1, is a multipurpose experiment designed for the
112 measurement of and search for a multitude of different processes. We will primarily
113 discuss its function as a Higgs finding machine, a more detailed description can be found
114 in Ref. [4]. It has a cylindrical shape consisting of a barrel segment, 21.6m long, and
115 two endcaps, 14.6m in diameter, aligned along the beam direction with its centre at the
116 beam interaction point. The endcaps are those nearer the beam line and so the materials
117 in these components typically have to be able to withstand higher amounts of radiation
118 and therefore tend to have worse performance. Many of its features exploit what one
119 would expect for measuring Higgs decays: it has almost full coverage of the area around
120 the collision point so that nearly every particle emanating from the collision can be
121 reconstructed, it has many complimentary subsystems (or layers) designed to measure
122 different specific particles so that Higgs bosons can be detected through a multitude of
123 decay modes.

124 For a Higgs with an intermediate mass (100 - 200 GeV) the high resolution (narrow
125 peak) channels are $H \rightarrow ZZ^* \rightarrow l^+l^-l^+l^-$ ¹ and $H \rightarrow \gamma\gamma$ so good energy resolution and
126 identification of electrons and muons is desirable down to very low p_T ($\sim O(10\text{GeV})$) as
127 well as good resolution and identification of high energy photons.

128 The central design feature of CMS is the very powerful superconducting magnet which
129 produces an axial magnetic field of 4T. The size of this field, as well as the density of
130 the calorimeter materials, allows for a compact and economical design (much more so
131 than its sister detector, ATLAS). Outside of the magnet lie the muon stations which
132 also serve as a return yoke for the magnetic field. The muon chambers in the barrel
133 consist of alternating layers of drift tubes and resistive plate chambers which provide both
134 accurate timing and hit location, in order to reconstruct muons down to low energies.
135 In the endcap the drift tubes are replaced with cathode strip chambers. Combining
136 information from the muon subsystem with information from the inner tracking system
137 (described below) allows muons at CMS to be reconstructed down to $p_T < 10$ GeV with
138 a resolution of $\sim 1\%$. The other three main subsystems at CMS, the tracking system
139 and the calorimeters, are located inside the magnetic field.

140 The first layer is the tracking system which is used to reconstruct the momentum
141 of any outgoing charged particles and to locate the primary and secondary vertices.

1 A * denotes that one Z can be off mass shell

This is surrounded by the calorimeters, the electromagnetic calorimeter ([ECAL](#)) and the hadronic calorimeter ([HCAL](#)). The first is a single layer of dense, transparent crystals which collects deposits of energy left by electrons and photons which shower inside the material. The second compliments this by providing a measurement of the energy deposited by hadrons (reconstructed as objects known as jets) through nuclear interactions. The [HCAL](#) is a sampling calorimeter in which the active material (plastic scintillator) is sandwiched between a dense absorbent material (brass or steel). This extends the radiation length of the calorimeter (clearly accommodating the compact design) and provides pointing information but degrades the resolution of reconstructing jets.

[CMS](#) uses a right-handed Cartesian coordinate system with the origin at the interaction point and the z -axis pointing along the beam axis. The x -axis points towards the centre of the [LHC](#) ring and the y -axis points vertically upwards. The azimuthal angle, $\phi \in [-\pi, \pi]$, is defined with respect to the x -axis in the transverse ($x - y$) plane. The polar angle θ is measured from the z -axis. Commonly, the direction of an outgoing particle is defined by ϕ and its pseudo-rapidity η ,

$$\eta = -\ln \tan\left(\frac{\theta}{2}\right). \quad (2.1)$$

The [LHC](#) is capable of producing 40M bunch collisions per second although many of these are not hard interactions, the result being that the outgoing particle debris follows the beam line. A hard (and therefore interesting) collision is characterised by the amount of energy produced in the transverse ($x - y$) plane. Therefore particles are commonly characterised by the projection of their momentum onto this plane, their transverse momentum,

$$p_T = \sqrt{p_x^2 + p_y^2}, \quad (2.2)$$

and the corresponding transverse energy, $E_T = E \sin(\theta)$.

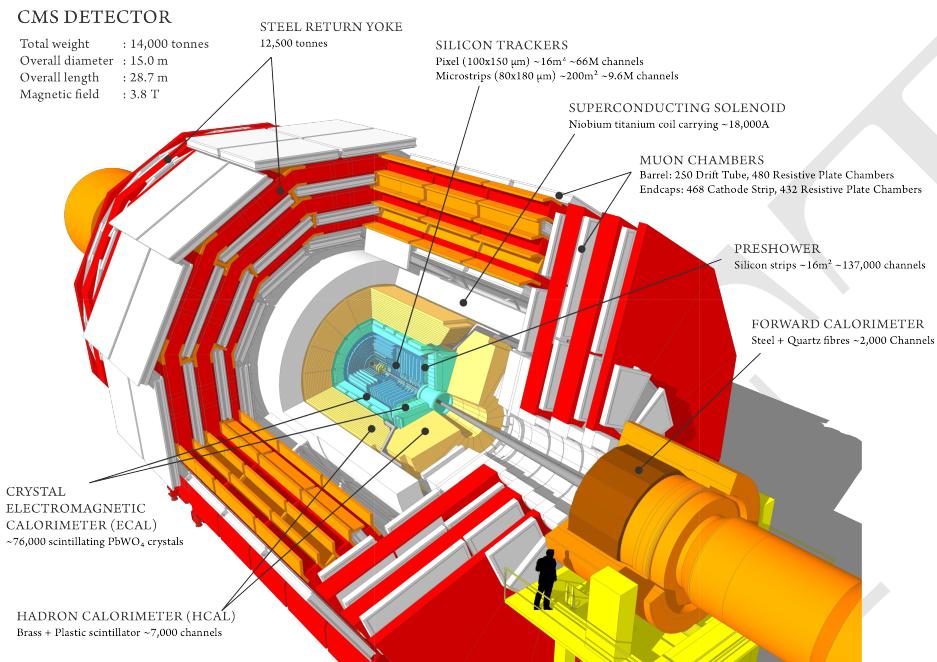


Figure 2.1.: This is a schematic representation of CMS showing the layered structure of subdetectors; tracking system, calorimeters, magnet and muon system.

165 2.2.1. Tracking system

166 During nominal LHC running conditions in 2012 there on average over 1000 particles
 167 from up to 50 overlapping $p\text{-}p$ collisions (pileup) per bunch crossing (every 50 ns). The
 168 tracker is designed to efficiently and precisely reconstruct all charged particle trajectories,
 169 thus their position and momentum, which are known as tracks. Due to the vast number
 170 of tracks emanating from multiple vertices in typical LHC collisions the tracking material
 171 and electronics are required to have high granularity, a fast response and be radiation
 172 hard. This conflicts with another important design feature of the tracker which is the
 173 aspiration to use the minimal amount of material in order to reduce multiple scattering,
 174 bremsstrahlung, photon conversion and nuclear interactions before particles reach the
 175 calorimeters. These criteria motivate the choice of silicon throughout the CMS tracking
 176 system. The structure of the CMS tracker is shown in Figure 2.2 and consists of a central
 177 pixel detector surrounded by layers of silicon strips aligned parallel to the beam line in
 178 the barrel (TIB and TOB) and perpendicular to the beam line in the endcap (TID and
 179 TEC).

180 By making multiple precise measurements of tracks as they pass through the pixel
 181 and silicon layers (hits) the track trajectory can be reconstructed and their momentum

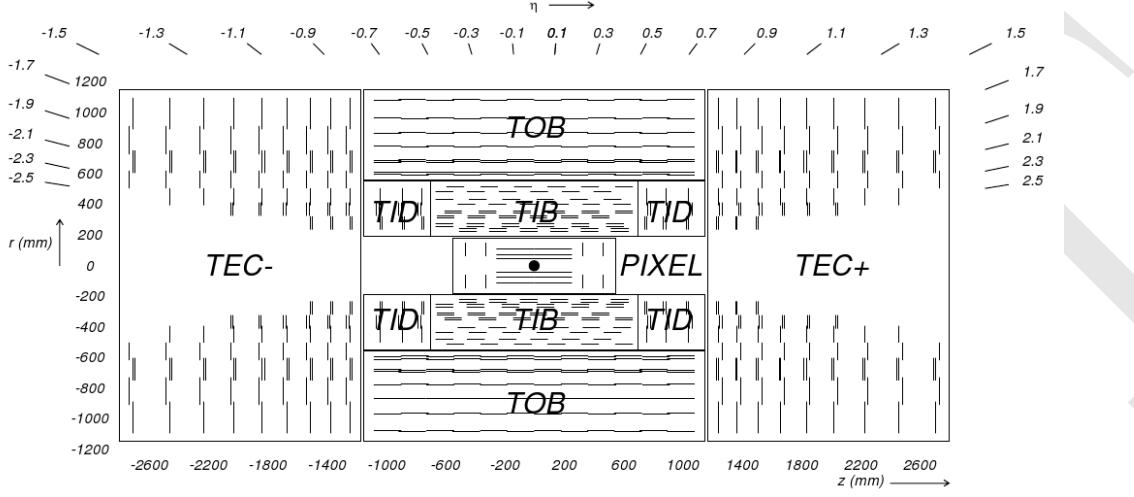


Figure 2.2.: A diagram of the CMS tracking system showing the PIXEL detector and outer silicon strip layers [4].

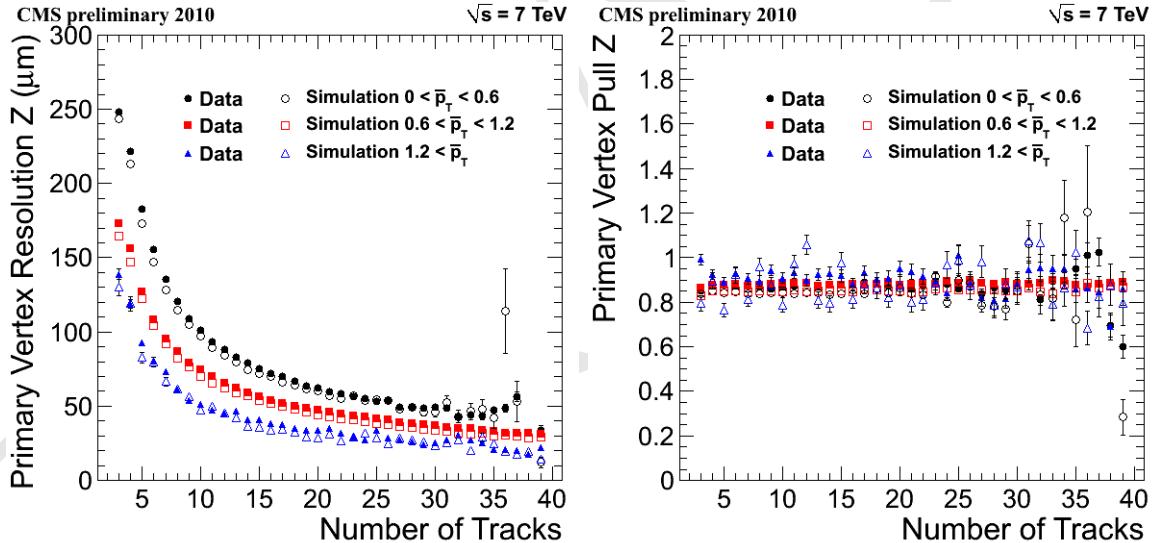


Figure 2.3.: A demonstration of the vertex position resolution in Z as a function of the number of tracks originating from that vertex. The different colors represent three different bins in the average track p_T where data is shown as solid points and simulation as open points [5]

calculated using their curvature in the ϕ plane due to the axial magnetic field. Tracks are grouped together (requiring that their separation is less than 1 cm in the z coordinate at the point of closest approach to the beamline) and assigned to a common point or origin (the primary vertex). The vertex resolution is driven both by the number of tracks originating from a particular vertex and how large their average p_T is. This is shown in Fig. 2.3 for preliminary data taken in 2010 at $\sqrt{s} = 7 \text{ TeV}$.

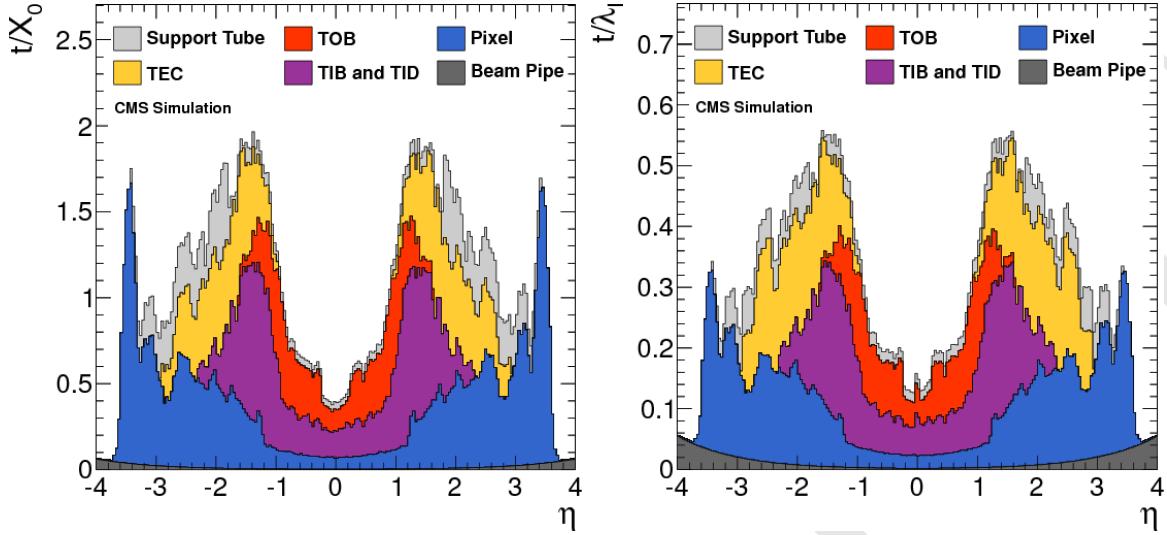


Figure 2.4.: The amount of CMS tracker material in radiation lengths (X_0) on the left and in nuclear interaction lengths (λ_I) on the right as a function of η for the different tracking subsystems.

The material budget for the tracker is shown in Fig. 2.4 which demonstrates as a function of η which subsystems of the tracker, beam pipe and servicing contribute to material inbetween the interaction point and the calorimeters in radiation lengths (X_0) and nuclear interaction lengths (λ_I). As shown in Fig. 2.2 the tracker has full coverage in ϕ and up to $|\eta| \leq 2.5$. As we will see later the tracking system is very important for the $H \rightarrow \gamma\gamma$ search at CMS as without it locating the primary vertex would be practically impossible.

2.2.2. Electromagnetic calorimeter

The ECAL is used to reconstruct the energy of electrons and photons which deposit their energy via electromagnetic showers inside the calorimeter material. The shower inside the crystal produces photons whose total energy is proportional to the energy of the incoming particle, hence the light output from the shower can be measured by photodiodes at the back of each crystal which in turn provides a measurement for the energy of the original particle. The ECAL has full hermetic coverage of the interaction point and consists of a single layer of lead tungstate (PbWO_4) crystals. The crystals are laid out in a quasi-projective geometry such that they point towards the interaction point with an offset of 3° making it much less likely that a photon, or electron, will pass straight through a gap between crystals. The ECAL consists of a barrel section and two

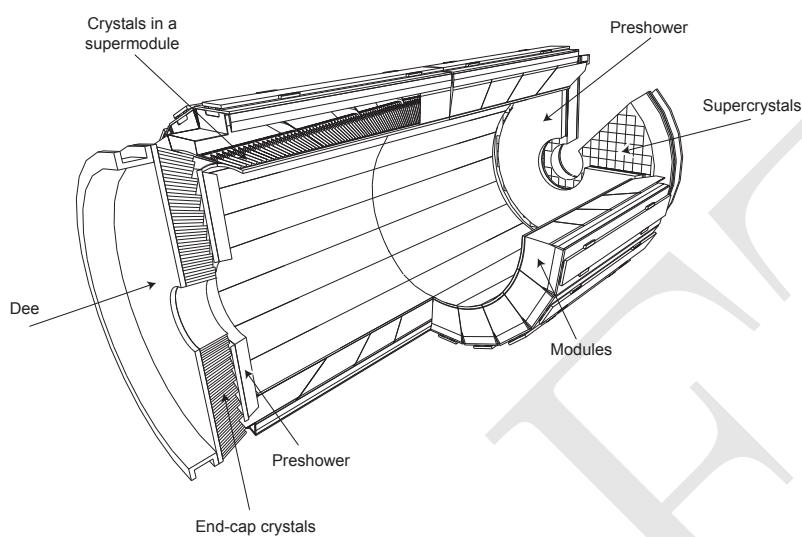


Figure 2.5.: A schematic drawing of the CMS ECAL layout showing the “modules” of crystals in the ECAL barrel and the ”Dees” of crystals in the ECAL endcap [4].

206 endcap disks which are preceded by a preshower (to aid with π^0 rejection), a schematic is
 207 shown in Fig. 2.5. There is a fiducial region between the barrel and endcap (to prevent
 208 reconstruction of showers which overlap both subsystems) yielding an ECAL coverage of
 209 $|\eta| < 2.5$ but not in the range $1.444 < |\eta| < 1.556$ and full coverage in ϕ . The crystals in
 210 both sections have a depth of 23cm (which amounts to $25.8 X_0$) implying that practically
 211 the whole depth of the shower is contained within this single layer.

212 PbWO_4 has some attractive properties for an electromagnetic calorimeter especially
 213 when considering $H \rightarrow \gamma\gamma$ decays. In order to achieve good energy resolution it is
 214 desirable to have a design in which most of the electromagnetic shower from an incoming
 215 photon or electron is contained within a single crystal, as less of the shower (and therefore
 216 energy) is lost in cracks and gaps inbetween crystals. PbWO_4 has a short radiation
 217 length ($X_0 = 0.89\text{cm}$) which complements a compact design as the entire depth of the
 218 shower can be contained in a crystal which is not very long. It has a small Molière
 219 radius (1.96cm) which means that the lateral size of the shower is small. This generally
 220 means that the cross-sectional size of a crystal can be small (yielding high granularity
 221 of the detector) whilst still containing a large percentage of the shower. It has a very
 222 short scintillation time decay constant (85% of the light is collected in 25ns), in other
 223 words it is very “fast”, allowing the energy in the shower to be collected and measured
 224 very quickly. This is clearly desirable in an environment like the LHC when collisions
 225 are happening up to every 25ns. The one draw back of PbWO_4 (apart from expense)
 226 which has crippled its use in detectors previously is its incredibly low light yield at room

temperature ($\sim 50\text{--}80$ photons/MeV). This is overcome at CMS with the use of silicon avalanche photodiodes (APDs) in the ECAL barrel and vacuum phototriodes (VPTs) in the endcap, which amplify the signal enough to make accurate measurements of the original photons energy.

231 Energy resolution

232 As we have seen previously (assume this formula will appear in the Introduction and can
233 be referenced) the diphoton invariant mass is given by,

$$m_{\gamma\gamma} = \sqrt{2E_1E_2(1 - \cos\alpha)}, \quad (2.3)$$

234 where E_1 and E_2 are the energy of the two photons and α is the angle between them.
235 Therefore the mass resolution has terms that depend on the photon energy resolution
236 and angular resolution,

$$\frac{\sigma_M}{M} = \frac{1}{2} \left[\frac{\sigma_{E_1}}{E_1} \oplus \frac{\sigma_{E_2}}{E_2} \oplus \frac{\sigma_\alpha}{\tan(\alpha/2)} \right], \quad (2.4)$$

237 where σ denotes the resolution and \oplus the quadratic sum. It is therefore desirable
238 to have both good energy resolution and good position resolution for photons (accurate
239 position measurements at the ECAL face alongside knowledge of the primary vertex can
240 be used to calculate the individual photons' direction and ergo the angle between them).
241 The energy resolution is usually then further parametrised as,

$$\frac{\sigma_E}{E} = \frac{a}{\sqrt{E}} \oplus \frac{b}{E} \oplus c, \quad (2.5)$$

242 where a is the stochastic term, b the noise term and c a constant term. In order to
243 achieve the best possible resolution all three of these terms need to be of a similar order
244 and as small as possible. The size of these terms has been determined from test beam data
245 in Ref. [4]. The stochastic term is driven by the material choice and detector type so cannot

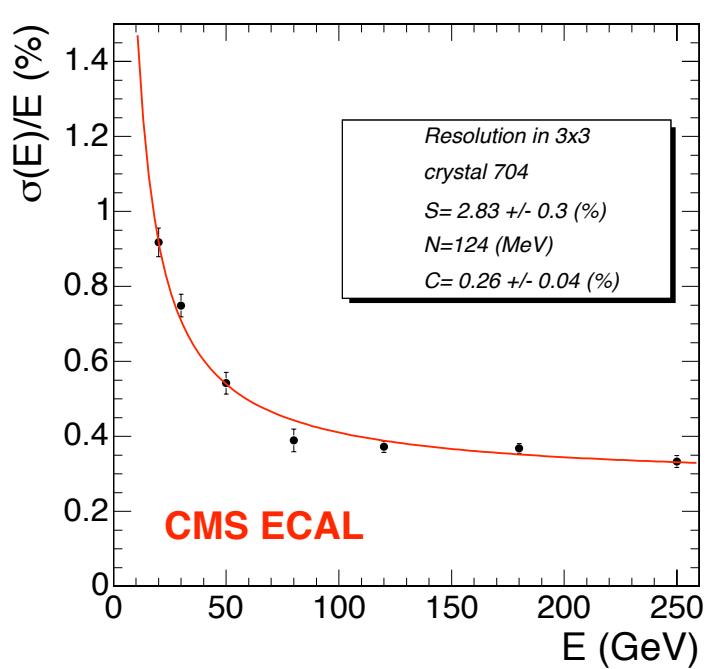


Figure 2.6.: The [ECAL](#) energy resolution, σ_E/E , as a function of electron energy measured from test beam data. The energy is measured in a 3×3 array of crystals centered on crystal of electron impact [4].

be improved once the machine is built. The main contributions to this term are lateral shower containment fluctuations, photostatistics and fluctuations in the energy deposited in the preshower absorber. For a fully active calorimeter (the [CMS ECAL](#) is not a sampling calorimeter) made of PbWO_4 the size of this term is good ($a = 2.8 \pm 0.3\% \text{ GeV}^{\frac{1}{2}}$). The constant term, which depends on non-uniformity of longitudinal light, intercalibration errors and energy leakage from the back of the crystal, can be minimized by use of *in situ* calibration of individual crystals and amounts to $c = 0.26 \pm 0.05\%$. However the $H \rightarrow \gamma\gamma$ analysis at [CMS](#) dispenses with the intercalibration constants by using a regression technique to estimate the photon energy which is discussed in Section 3.3. The noise term which has contributions from electronics noise (including signal digitisation) and event pile-up (additional particles causing overlapping signals) is measured as $b = 126 \text{ MeV}$. The ECAL energy resolution, σ_E/E , as a function of electron energy is shown in Fig. 2.6 as measured from a beam test.

²⁵⁹ **Transparency corrections**

²⁶⁰ Due to the high particle flux present at CMS the ECAL crystals and electronics have to
²⁶¹ be radiation hard, especially in the endcap. This is another motivating factor for using
²⁶² PbWO₄ as the crystal material. The crystals are additionally doped with Nb to improve
²⁶³ the induced absorption coefficient. Over time and long exposure to radiation the crystals
²⁶⁴ loose their transparency, although there is considerable natural recovery during down
²⁶⁵ periods. An important part of the ECAL monitoring and calibration comes in the form
²⁶⁶ of transparency corrections to compensate for these losses. At regular intervals during
²⁶⁷ LHC running laser pulses are injected into the crystals to measure the crystal response.
²⁶⁸ Two different wavelengths of laser are used, one blue ($\lambda = 440\text{nm}$) which is very similar
²⁶⁹ to the scintillation emission peak and therefore expected to be affected by transparency
²⁷⁰ changes in a similar way to typical scintillation light, and one red ($\lambda = 796\text{nm}$) which
²⁷¹ is far from the scintillation emission peak and affected very little by the changes in
²⁷² transparency. Hence, by comparing the red and blue laser light response, time and η
²⁷³ dependent corrections for crystal transparency loss can be calculated. A closure test for
²⁷⁴ these corrections, in 2012 data, is shown in Fig. 2.7 which shows the ratio of electron
²⁷⁵ energy (calculated from the ECAL) to electron momentum (calculated from the tracker)
²⁷⁶ before and after transparency (or laser) corrections.

²⁷⁷ **Preshower**

²⁷⁸ The dominant source of background to high energy photon signals are neutral mesons,
²⁷⁹ mainly pions (π^0), which decay into two approximately collinear photons and can therefore
²⁸⁰ look very much like a single high energy photon. The ECAL endcap is preceded by a
²⁸¹ preshower to specifically target this and provide the endcap with a higher granularity.
²⁸² The preshower is a sampling calorimeter which consists of two layers: a lead plate, to
²⁸³ initiate the shower, in front of a fine grained silicon detector which has two layers of
²⁸⁴ orthogonal strips. There are many other characteristics of π^0 decays which can help in
²⁸⁵ differentiating them from real (“prompt”) photons, these include isolation (discussed
²⁸⁶ later in this Chapter in Section 2.4) and the shower shape (discussed in Chapter 4 in
²⁸⁷ Sections 4.1.1 and ??).

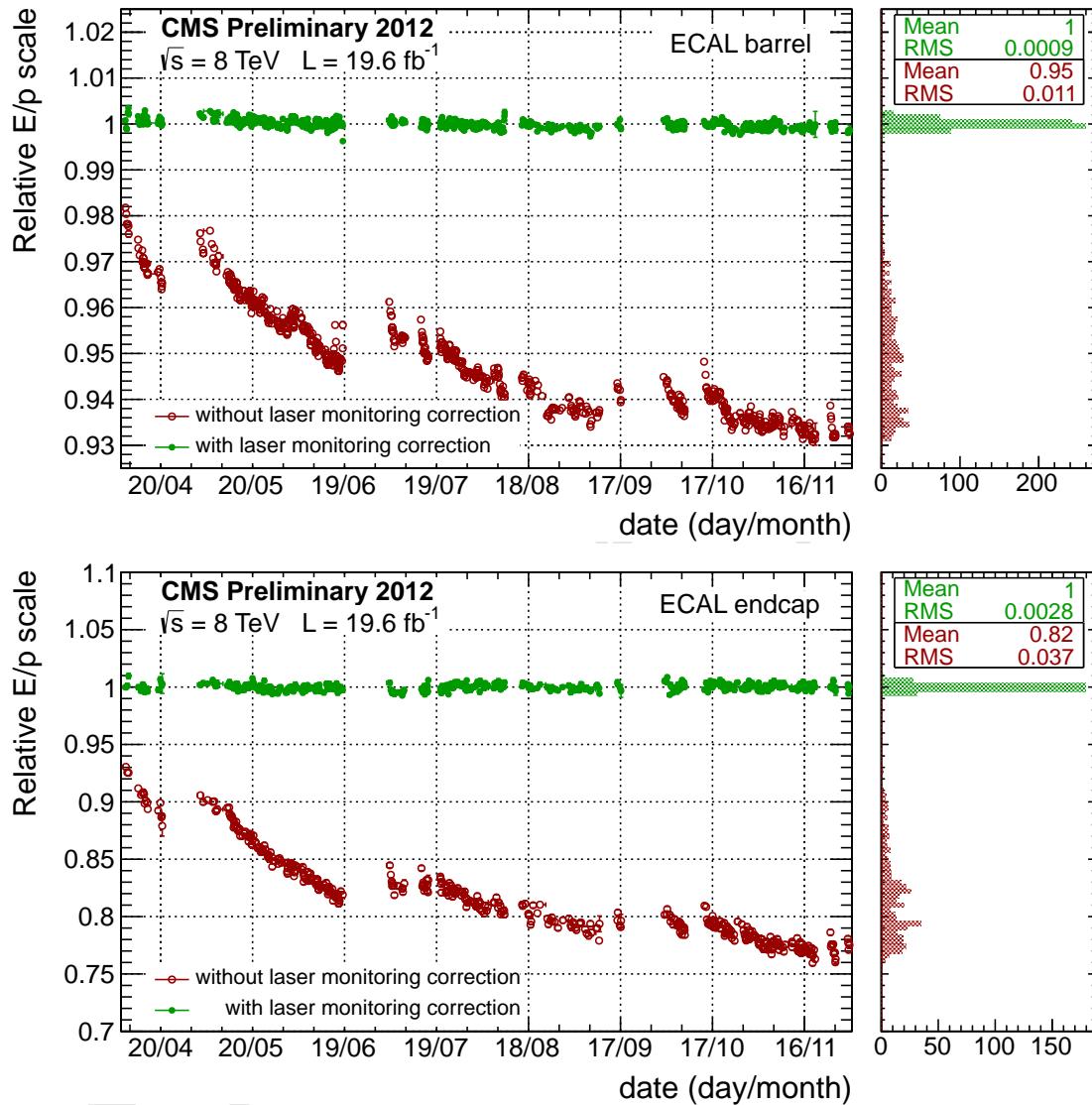


Figure 2.7.: Ratio of the electron energy, E , to the electron momentum, p , measured in the CMS barrel (top) and endcap (bottom) for 2012 data. The red (green) points show the performance before (after) the laser monitoring derived corrections [6].

²⁸⁸ **Photon reconstruction**

²⁸⁹ Calculating an incoming photons energy amounts to summing the energy deposited by
²⁹⁰ the electromagentic shower which is initiated by the photon impact at the crystal face.
²⁹¹ Due to the presence of material in the beam pipe and tracking system (see Fig 2.4)
²⁹² about 60% of photons will convert into an electron-positron pair before they reach the
²⁹³ **ECAL**. If this is the case the shower will spread out in ϕ due to the presence of the
²⁹⁴ magnetic field; the photon converts to electrons which bend and Bremsstrahlung radiate
²⁹⁵ additional photons, and can be distributed among multiple (up to hundreds of) crystals.
²⁹⁶ Consequently clustering (pattern matching) algorithms are deployed to calculate the
²⁹⁷ “raw” photon energy. Corrections to this energy are subsequently applied to account for
²⁹⁸ any energy loss as explained in Section 3.3. The shower will appear as a local maxima
²⁹⁹ amongst a spatial neighbourhood of crystal energy deposits and so the algorithms used
³⁰⁰ search first for the most energetic crystals (known as the “seed” crystals) and then extend
³⁰¹ to amass as a large a fraction as possible of the original shower energy. There are three
³⁰² cases to include which are i) photons which reach the **ECAL** without interacting with
³⁰³ any of the intermediate material in the beam pipe and tracking system (referred to as
³⁰⁴ unconverted photons), ii) photons which convert into an electron-positron pair inside the
³⁰⁵ tracker and shower in the barrel, iii) photons which convert and shower in the endcap.

³⁰⁶ We will first consider the case of converted photons in the barrel. This is so similar
³⁰⁷ to the case of a real electron that the identical algorithm is used for electrons as well.
³⁰⁸ The method used is known as the “Hybrid” algorithm, depicted in Fig. 2.8, which makes
³⁰⁹ clusters of clusters known as a “supercluster”: a cluster being a set of crystals which
³¹⁰ pick up an electron or a bremsstrahlung photon and a cluster of clusters being a set
³¹¹ of these which make up all the electrons and photons radiated from the original object.
³¹² The algorithm can be described as a five step process as follows,

- ³¹³ 1. Locate the seed crystal which is the maximum energy crystal in the search region, not
³¹⁴ already in a cluster, and which must satisfy the threshold condition, $E_T > 1$ GeV.
- ³¹⁵ 2. Extend in η to construct a “domino” which is 1×3 crystals in $\phi \times \eta$. If the energy
³¹⁶ of the central crystal in the 1×3 domino is greater than 1 GeV then extend this to
³¹⁷ a 1×5 domino in $\phi \times \eta$.
- ³¹⁸ 3. Traverse along ϕ , up to a maximum of 17 crystals in both directions, adding
³¹⁹ dominoes in this way. If a domino has less energy than 0.1 GeV then it is excluded.

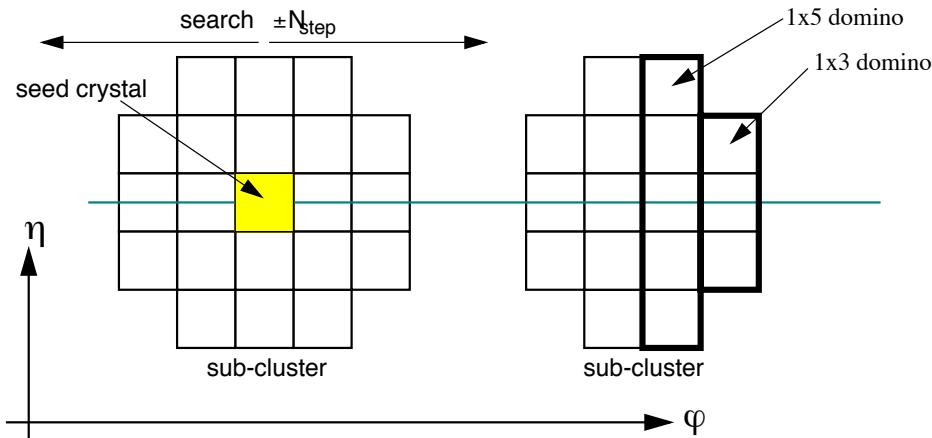


Fig. 6: Domino construction step of Hybrid algorithm

Figure 2.8.: The domino construction setup of the hybrid clustering algorithm [7].

- 320 4. Of the clusters of adjacent dominoes the seed domino (most energetic) must have
321 energy greater than 0.35 GeV.
- 322 5. Repeat starting from Step 1 to build a new supercluster.

323 In this way a collection of dominoes are clustered in ϕ creating a “supercluster” of
324 smaller clusters (which has a maximum of area of 5×35 in $\eta \times \phi$). A similar but slightly
325 modified algorithm is used for photons and electrons in the endcap. This is known as
326 the “Multi 5×5 ” algorithm and proceeds,

- 327 1. Locate the seed crystal which is the maximum energy crystal in the search region, not
328 already in a cluster, and which must satisfy the threshold condition, $E_T > 180$ MeV.
- 329 2. The 5×5 array of crystals surrounding the seed are summed to make a cluster.
- 330 3. Crystals at the edge of the cluster can seed new overlapping 5×5 clusters if they
331 are local maxima compared to their neighbouring crystals.
- 332 4. Proceed in this way in any direction building an overlapping collection of 5×5
333 clusters to create a supercluster.
- 334 5. Repeat starting from Step 1 to build a new supercluster.

335 The final case concerns photons which reach the **ECAL** without converting which have
336 a shower that is much more localised in η and ϕ , about 94% of its energy is deposited in
337 an area of 3×3 crystals and more than 97% in an area of 5×5 crystals. This provides a
338 definition for the conversion variable, R_9 , such that,

$$R_9 = \frac{E_{3 \times 3}}{E_{SC}} \begin{cases} \text{unconverted if } R_9 \geq 0.94 \\ \text{converted otherwise.} \end{cases} \quad (2.6)$$

If a photon fulfills the unconverted requirement given above (Eq. 2.6) the photon energy is reconstructed as the sum of all energy deposited in the 5×5 array of crystals which surround the most energetic crystal. It has been shown that using a fixed window for non-converting photons yields a better energy resolution than any clustering procedure [7].

The location of a supercluster is determined as the energy weighted mean position of the crystals in the supercluster. This gives a position resolution which is much smaller than the size of an individual crystal ($\sim 20 \times 20\text{mm}^2$). Values for an electron of $p_T = 35\text{ GeV}$ in the ECAL barrel, in the absence of pileup, are $\sigma_\eta = 1.0 \times 10^{-3}$, $\sigma_\phi = 1.6\text{ mrad}$.

348 Electron and photon differences

It is clear that when only considering the ECAL there is no difference between electrons and photons. Consequently at the level of the calorimetry there is no distinction between them, simply the idea of a supercluster which can apply to both. In the case of an electron, information from the tracker can be included using a Gaussian sum filter algorithm [8], where a series of compatible track hits are associated to the supercluster. This is used to provide a supplementary measurement of the electrons momentum and consequently improve the energy resolution of electrons. When considering photons for an analysis an electron veto must be applied requiring that no track hits should be found close to the interaction point near the photon direction (see Chapter 3). An important feature of supercluster reconstruction at CMS is that when all track information is ignored electrons and photons are identical. This is a principal ingredient in the $H \rightarrow \gamma\gamma$ analysis which allows data driven calibration, validation and efficiency measurements of photons using electrons (see Section 3.6).

362 2.2.3. Hadronic calorimeter

Surrounding the ECAL, but still inside the magnet, is a sampling HCAL which has geometric coverage up to $|\eta| < 5.0$ when including the specialised forward components.

365 It consists of alternating layers of brass plates and plastic scintillators (where in the very
 366 forward region the brass is replaced with steel). The **HCAL** thickness constitutes around
 367 10-15 nuclear interaction lengths (λ_I) depending on η . Any outgoing hadrons from the
 368 interaction (of which there are many for a typical event with high E_T) get reconstructed
 369 as objects known as “jets” by amalgamating information from the tracking the system,
 370 the **ECAL** and the **HCAL**. This process is described in more detail in the particle flow
 371 section below (section 2.3).

372 2.2.4. Muons

373 Outside of the magnet lie the **CMS** muon chambers. These consist of alternating layers
 374 of drift tube chambers (cathode strip chambers) in the barrel (endcap) and resistive
 375 plate chambers which also act as a return for the magentic flux. The muon detector
 376 has coverage up to $|\eta| < 2.4$ and given the particularly conspicuous signature of muons
 377 (several hits in the tracker and hits in each muon station layer) the reconstruction
 378 efficiency and momentum resolution of muons is very good even down to low p_T . The
 379 muon resolution as a function of p_T is shown in the left hand plot of Fig. 2.9 for simulated
 380 data at $\sqrt{s} = 7$ TeV.

381 2.3. Particle flow and jets

382 A traditional approach to detector based particle physics is to consider the objects
 383 we measure as opposed to the underlying physics objects. These are often known as
 384 calorimeter objects, for example, a track, an electromagnetic shower or a calorimeter
 385 jet. A more modern approach is to couple information in all of the subdetector systems
 386 together to reconstruct more physical objects. For example, a charged hadron will leave a
 387 track, deposit some energy in the **ECAL** and deposit the rest of its energy in the **HCAL**.
 388 This technique of reconstruction is known as particle flow (PF)² and is particularly
 389 useful when considering jets. Whilst an electron, photon or muon are fairly unique
 390 looking signatures hadrons are often not. The abundant number of gluons and quarks
 391 produced in **LHC** collisions hadronise via the strong interaction as they travel away from
 392 the interaction. As they have typically high momentum the hadronistation occurs in a
 393 collimated fashion leaving a signature of several tracks and a hadronic cluster. The PF

²The official name inside **CMS** is global event description (**GED**) although this is rarely used

394 reconstruction algorithm can be simplistically viewed as the following procedure (more
395 details given in Ref. [9]):

- 396 1. Make small clusters from each subdetector component; tracks, **ECAL** and **HCAL**
397 clusters to create a list of unassociated objects.
- 398 2. Match tracks and clusters together and associate them to a newly reconstructed
399 particle known as a **PF** candidate:
 - 400 • Tracks and clusters associated with hits in the muon chambers are tagged as
401 *muons* and removed from the list.
 - 402 • Tracks and clusters associated with electrons, including Bremsstrahlung pho-
403 tons, are tagged as *electrons* and removed from the list.
 - 404 • Tracks associated to an **HCAL** cluster are tagged as *charged hadrons*, assigned
405 an energy ascertained from a weighted average of the cluster energy and track
406 momentum and subsequently removed from the list.
 - 407 • Any excess cluster energy in the **HCAL** is assigned as a *neutral hadron* and
408 removed from the list.
 - 409 • If an **ECAL** cluster is associated to an **HCAL** cluster and a track, it is assigned
410 as a *charged hadron* with the appropriate weighted energy and removed from
411 the list.
 - 412 • If an **ECAL** cluster is associated to an **HCAL** cluster with no track, it is assigned
413 as either a *photon* or a *neutral hadron* depending on the **HCAL** to **ECAL** energy
414 ratio and removed from the list.
 - 415 • Any remaining unlinked candidates are assigned as *photons* or *neutral hadrons*
416 depending on whether they are **ECAL** or **HCAL** clusters.
- 417 3. In this way all information in the detector is used to create a list of candidates
418 which can be any of a muon, electron, photon, charged hadron or neutral hadron.
- 419 4. These are then used to construct composite detector objects such as jets if necessary.

420 Particle flow jets are constructed using the anti- k_T algorithm [10]. This algorithm
421 is both infrared and collinear (IRC) safe and preferentially clusters soft (low p_T) jets
422 with hard (high p_T) jets to be robust in the **LHC** pileup conditions. These jets can
423 then be additionally tagged as *b* or *c* (i.e. those containing a bottom or charm quark

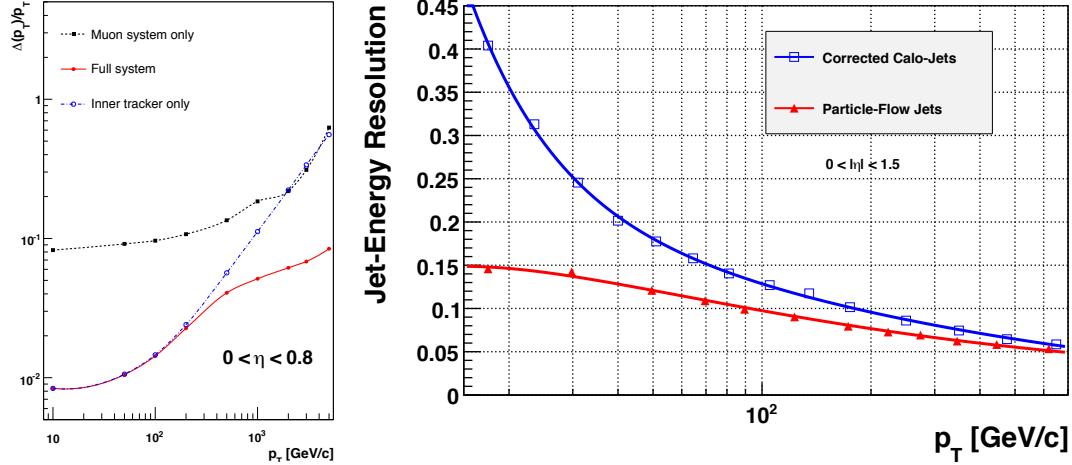


Figure 2.9.: Left: The muon p_T resolution as a function of p_T for muons in the range $0 < |\eta| < 0.8$ when using only the muon system (black), only the tracking system (blue) and using both (red) [4].

Right: The jet energy resolution as a function of jet p_T when using PF jets as compared to calo jets [13].

424 respectively) using the techniques of the type described in [11]. There are also energy
 425 corrections applied to jets to account for pileup (ρ subtraction technique as described
 426 in section 2.5), non-uniform detector response (p_T and η dependent corrections derived
 427 from Monte Carlo (MC)) and data MC differences (residual p_T and η corrections derived
 428 from γ +jet and Z +jet samples in data), see Ref. [12] for details. The clustering, and
 429 subsequent energy correction, of jets in this way also provides a measure of the missing
 430 transverse energy (\cancel{E}_T), the amount of energy in an event taken away by undetectable
 431 particles such as neutrinos. A comparison of the jet energy resolution, as a function of
 432 jet p_T , for calorimeter jets and particle flow jets is shown in Fig. 2.9 for simulated data
 433 at $\sqrt{s} = 7$ TeV.

434 Should possibly include a note here stating that a PF photon is very different to the
 435 photons used in the analysis. PF is useful for tagging physics objects like tau's or b's
 436 and for isolation sums but is not necessarily the best for reconstructing well measured
 437 objects like photons, electrons and muons.

438 2.4. Isolation

439 One way of differentiating between real (prompt) photons and fakes is the use of isolation.
 440 One would expect that in the absence of pileup a real photon would be isolated, which is

to say there are no other particles (detector activity) in its vicinity. For a jet faking a photon (which is nearly always a π^0) this is not the case and one would expect the π^0 to be surrounded by additional hadronised particles (detector activity in the tracker, ECAL and HCAL). In the CMS $H \rightarrow \gamma\gamma$ analysis isolation sums are used to distinguish prompt photons from fakes. Three variables are used which consider each photons isolation relative to activity in the surrounding environment. The procedure is to create a hollow cone around the photon candidate (of outer radius ΔR_O and inner radius ΔR_I , where $\Delta R = \sqrt{\Delta\eta^2 + \Delta\phi^2}$) and sum the energy contained in that cone of PF candidates; charged hadrons, neutral hadrons and electrons/photons. These three variables are defined in the following way:

- **Charged hadron isolation:** Sum of charged hadron PF candidates E_T in cone of $\Delta R_O = 0.3$ and $\Delta R_I = 0.02$.
- **Neutral hadron isolation:** Sum of neutral hadron PF candidates E_T in cone of $\Delta R_O = 0.3$ and $\Delta R_I = 0.0$.
- **e/γ isolation:** Sum of e/γ PF candidates E_T in cone of $\Delta R_O = 0.3$ with a central η strip of 0.070 (0.015) removed for barrel (endcap) photons.

2.5. Pileup

There can be up to 1.1×10^{11} protons in each bunch at the LHC which results in multiple interactions (can be as many as 50 primary vertices) per bunch crossing. This effect is known as pileup and can present a challenge in finding the primary vertex and also produces additional energy in each event which originates from somewhere other than the primary vertex. To combat the latter of these two effects a technique called ρ subtraction is used to corrected the energy of jets and isolation sums for pileup. ρ is defined as a per event quantity and is computed by summing the energy in all the calorimeters and dividing by the calorimeter area and thus represents the average energy density in the detector per event.

Chapter 3.

⁴⁶⁷ Common analysis components

⁴⁶⁸ “I don’t have a quote”

⁴⁶⁹ — Matthew Kenzie, 1785–1854

⁴⁷⁰ This thesis describes three complementary analysis regimes in the Higgs to two
⁴⁷¹ photons search at CMS. These differ in their photon selection, event selection, event
⁴⁷² classification (or categorisation) and statistical methods for extracting results. They are
⁴⁷³ described in the following chapter (Chapter 4). However, there are many components
⁴⁷⁴ which they share. These are detailed below.

⁴⁷⁵ As we have seen in Eq. 1.1, repeated below in Eq. 3.1 for convenience, the diphoton
⁴⁷⁶ invariant mass is constructed from the two photon energies and the angle between them.
⁴⁷⁷ Consequently, important considerations for this analysis are photon energy resolution and
⁴⁷⁸ good opening angle resolution which is completely dominated by the vertex resolution,
⁴⁷⁹ as the position resolution of the photons (location they hit the ECAL) is negligible in
⁴⁸⁰ comparison. Details of how this is exploited in the analyses are given at the end of
⁴⁸¹ this chapter in Sections 3.3 and 3.4. The selection of events is described in the chapter
⁴⁸² after this (Chapter 4) alongside the categorisation, or binning, scheme into different
⁴⁸³ classes of event which take advantage of areas of phase space which share similar signal
⁴⁸⁴ to background ratios.

$$m_{\gamma\gamma} = \sqrt{E_1 E_2 (1 - \cos(\alpha))} \quad (3.1)$$

485 After a preliminary discussion of multivariate analysis techniques, the datasets, the
486 triggering and the MC are discussed.

487 3.1. Boosted Decision Trees

488 MVAs are commonly used in High Energy Physics analyses to extract the maximum
489 possible signal sensitivity in cases where the background rates are high. The advantage
490 of MVAs is that given a set of input variables a network of sequential cuts can be built,
491 to classify or correct events, in a multidimensional phase space to exploit differences
492 between the signal and background in these variables and importantly in the correlations
493 between these variables. A particular type of MVA which is used widely in this analysis
494 is the Boosted Decision Tree (BDT). BDTs are preferred because they are more robust
495 to the inclusion of variables which have little or no discriminating power. There are
496 two broad types of BDT used, one is known as a regression BDT and the other as a
497 classification BDT.

498 A classification BDT will, given a set of input variables, assign a value (typically
499 between -1 and 1) to each event based on how signal like that event is. This serves to
500 collapse all the event information into one discriminating variable which can be used to
501 classify differences between the signal and background. The input is provided as the
502 probability distributions (which can be supplied as binned or unbinned data samples
503 or as a functional form) of the background and signal for a set of “input variables”.
504 The process involves construction of a series of Decision Tree (DT)s complemented by a
505 “boosting” step which serves to mitigate against “overtraining” on fluctuations within
506 the training samples. This analysis chooses a particular type of decision tree boosting
507 known as “gradient” boosting because it is more robust against outliers or mislabeled
508 data points [14].

509 The DT is built by applying sequential cuts to the input variables and assessing the
510 relative signal purity, p , in the sub-sample remaining after each cut.

$$p = \frac{N_s}{N_s + N_b}, \quad (3.2)$$

511 where N_s and N_b are the sum of weights of the signal and background remaining
 512 in each sub-sample. A threshold criterion, known as the Gini index [14] $p(1 - p)$, is
 513 applied to decide whether to split the sample further. The process continues and the
 514 splitting is curtailed when either the threshold or the user defined maximum tree depth
 515 (number of subsamples allowed) is reached. The value of each cut is varied such that the
 516 signal purity, p , in each sub-sample is maximised. An event is assigned a value of -1 or
 517 +1 depending on whether it falls into a sub-sample with $p > 0.5$ or not. Clearly some
 518 fraction of events will be misclassified where the actual number which get misclassified
 519 will depend on the discriminatory power available from the chosen input variables. In
 520 order to reduce this effect a series of DTs are trained and each assigned a weight derived
 521 by the “boosting” process. If we assign each DT as a member of a family of M functions,
 522 $f(\vec{x}; \vec{a}_m)$, which depend on the input variables, \vec{x} and the set of cuts in that tree, \vec{a}_m .
 523 The object is to construct an overall decision tree which consists of the weighted average
 524 of each DT,

$$F(\vec{x}; \vec{\beta}, \vec{a}) = \sum_{m=0}^M \beta_m f(\vec{x}, \vec{a}_m) \quad \text{where } \vec{\beta} = (\beta_0, \beta_1 \dots \beta_m) \quad (3.3)$$

525 The boosting procedure is implemented by adjusting the weights $\vec{\beta}$ in order to minimise
 526 the deviation in the loss function (Eq. 3.4) between the weighted tree response $F(\vec{x}; \vec{\beta}, \vec{a})$
 527 and the true output y obtained from the training sample.

$$L(F, y) = \ln(1 + e^{-2F(x)y}) \quad (3.4)$$

528 A common procedure when constructing a BDT to check for overtraining is to split
 529 both the background and signal into two independent samples. One is used to *train*
 530 the BDT and one is used to *test* the response of the output. Clearly one requires that
 531 both the training and independent test sample look the same in the output variable.
 532 This is usually quantified by use of a Kolmogorov-Smirnov test, which broadly speaking
 533 ascertains the probability that the training and test samples originate from the same
 534 underlying distribution [15].

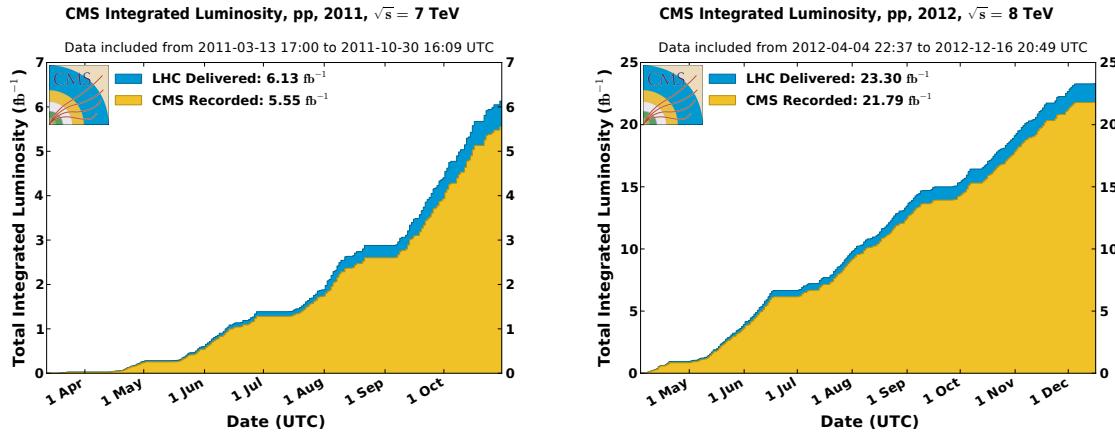


Figure 3.1.: The total integrated luminosity delived and recorded by CMS during the 2011 (left) and 2012 (right) run periods.

535 In this way the output of $F(\vec{x}; \vec{\beta}, \vec{a})$ for a classification **BDT** will be a “semi-continuous”
 536 output from -1 to 1 with signal events in general given a higher score than background
 537 events.

538 A regression **BDT** is used to estimate the true value of some variable given the values
 539 and correlations of several other variables. They are commonly used for correcting the
 540 energy of a particular object, for example a photon. Given a **MC** source of photons the
 541 “true” energy is regressed from the position, shape and raw energy of the supercluster.
 542 For regression **BDTs** the output $F(\vec{x}; \vec{\beta}, \vec{a})$ represents the estimated corrected energy
 543 and the boosting procedure targets minimising the deviation between this and the true
 544 energy in **MC**.

545 3.2. Data samples and triggering

546 The data consists of two independent samples of proton-proton collisions collected by the
 547 CMS experiment at the **LHC** in 2011 and 2012 with a centre-of-mass energy (\sqrt{s}) of 7
 548 and 8 TeV respectively. The total integrated luminosity of the two samples is 5.1 fb^{-1}
 549 and 19.7 fb^{-1} in 2011 and 2012 respectively and collectively referred to as **LHC** Run 1.
 550 The response of the detector has changed considerably over this period and much of the
 551 variation is modelled by the **MC** simulation. Figure 3.1 shows the integrated luminosity
 552 delivered and recorded by **CMS** during **LHC** Run 1.

Events are selected for the analysis by requiring they pass an asymmetric diphoton trigger with E_T thresholds of 26 (18) and 36 (22) GeV for the leading (trailing) photon in the 2011 and 2012 runs respectively. The candidates are also required to have either a high value of R_9 or to pass a loose calorimetric identification and isolation requirement. High trigger efficiency is achieved by selecting photon candidates which pass either requirement. The efficiency of the trigger for the analysis preselection is 99.5%.

3.2.1. Monte Carlo

Accurate simulation of detector effects and efficiency is highly important. Knowledge of the expected Higgs signal shape is clearly essential and although the size and shape of the $m_{\gamma\gamma}$ background is entirely data driven when extracting results, simulating the kinematics, shower shape and resolution properties of the background is essential when training the selection and optimising the binning.

As explained in Chapter 1 the two main production mechanisms for a Standard Model Higgs boson at the LHC are gluon fusion ([ggH](#)) and vector boson fusion ([VBF](#)). Typically the latter is produced at much higher Higgs p_T and this feature is exploited in the analysis. Consequently, it is important to model the p_T distribution of these two production modes accurately. The signal samples for these two processes are generated using POWHEG [16, 17] at NLO interfaced with PYTHIA [18] including a reweighting factor which matches their p_T spectrum to that when including the NNLO and NNLL terms. For the associated production modes, with a W^\pm, Z or t -quarks, (vector boson associated production ([VH](#)) and top quark associated production ([t̄H](#))) only PYTHIA is used. The Standard Model Higgs boson cross sections and branching fractions are taken from Ref. [19]

The spin-2 graviton with minimal couplings, 2_m^+ , has two production mechanisms, one via gluon-fusion (ggX) and one via quark-antiquark annihilation ($q\bar{q}X$). The graviton samples are generated using the JHU generator [20] in which the p_T spectrum of these samples is matched to the SM so that any bias obtained from mismodelling of the p_T is avoided.

The simulated background samples are used solely for cut and category optimisations and training of multivariate discriminants. The background which contains the quantum chromodynamics ([QCD](#)) continuum of prompt photons (refering back to Chapter 1 these are produced by Born and Box diagrams) is simulated using SHERPA [21] at 8TeV and

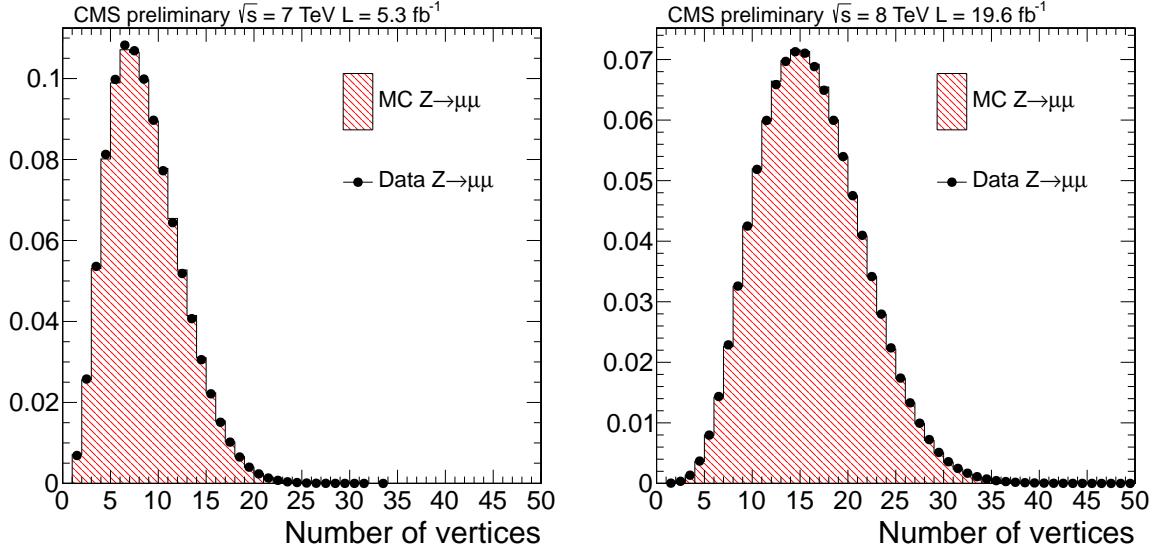


Figure 3.2.: Distribution of the number of reconstructed vertices in the 2011 (left) and 2012 (right) run periods. Calculated using the Deterministic Annealing algorithm in [REF!] for $Z \rightarrow \mu^+\mu^-$ events in data (black dots) and MC (red histogram) after reweighting.

585 MADGRAPH [22] at 7TeV. The prompt-fake and fake-fake backgrounds, in which one
 586 or more photons are faked by a neutral meson (usually a π^0) reconstructed as a jet are
 587 generated using PYTHIA. Samples of $Z \rightarrow e^+e^-$, $Z \rightarrow \mu^+\mu^-$ and $Z \rightarrow \mu^+\mu^-\gamma$ used for
 588 data/MC comparisons are generated with POWHEG.

589 All of these *generator level* samples are then run through the full CMS detector
 590 simulation using GEANT4 [23]. This includes the effect of overlapping vertices (pileup)
 591 and detector effects (such as noise and crystal degradation) in four bins of time (Run2011,
 592 Run2012AB, Run2012C, Run2012D).

593 3.2.2. Pileup and beamspot reweighting

594 An important difference between the simulated samples and the data which can have a
 595 large impact on the analysis is the distribution of the number of primary vertices. The
 596 *pileup* in the event effects many important analysis variables, for example photon shower
 597 shape and photon isolation as well as the diphoton invariant mass if the chosen vertex is
 598 wrong. Consequently the MC is reweighted such that the pileup distribution matches
 599 that in data. The reweighting technique is validated using $Z \rightarrow \mu^+\mu^-$ events as shown
 600 in Figure 3.2 for the 7 and 8 TeV samples respectively.

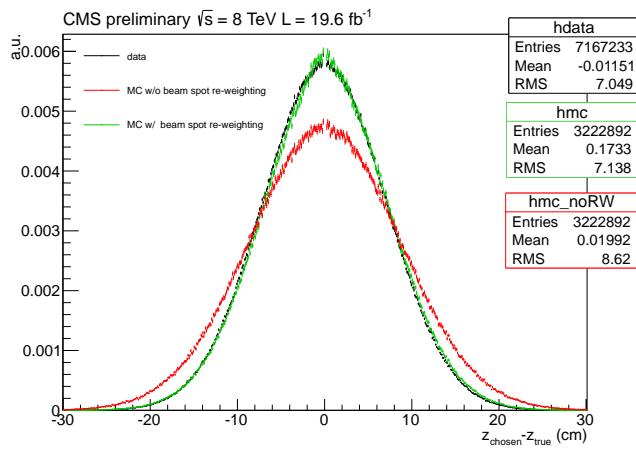


Figure 3.3.: Distribution of Δz (the distance between the chosen vertex and the true vertex in the z direction) for data (black), the MC (red) and the MC after beam spot reweighting (green) for $Z \rightarrow \mu^+ \mu^-$ events.

When the chosen vertex is incorrect the mass resolution is dominated by the spread in position of the pileup vertices (known as the beamspot width, $\sigma_z^{beamspot}$). Accurate modelling of this spread is important so that the resolution of wrong vertex events in simulation matches that in data. The MC overestimates the beamspot spread by some 20% so a simple reweighting is implemented for MC events in which the wrong vertex chosen (as the effect is negligible for events in which the chosen vertex is correct) such that the distribution of the distance between the chosen vertex position and the true vertex position, $\delta z = z_{chosen} - z_{true}$, match between data and MC. The effect with and without reweighting compared to data is shown in Figure 3.3.

3.3. Energy measurement of photons

The photon energy obtained from the supercluster sum described in Section 2.2.2 even when including the intercalibration and transparency corrections shown in Figure 2.7 does not give the most optimal resolution for the energy measurement of photons at CMS. On top of this energy (known as the *raw* supercluster energy, E_{raw}) it is also valuable to correct for additional energy losses. These arise from *bremstrahlung*; where the photon converts in the material upstream of the ECAL and the two electron legs radiate additional photons and thus the some of the photon shower is missed, and for local containment of the shower; where some energy is lost through small gaps between ECAL crystals and larger gaps between “modules” or sections of crystals. These corrections are

obtained using a specialised regression **BDT** (see Sec. 3.1) trained on a **MC** source of prompt photons from a sample containing photons and jets. The **BDT** targets accurate measurements of individual photons' energy by correcting the raw supercluster energy and provides an estimate for the energy resolution of each photon given the position and shower shape of the supercluster. The training is done separately for barrel and endcap photons (as the cluster shapes look very different for these two distinct regions) and is also performed separately for the 7 and 8 TeV datasets. The following input variables are used:

- The global position of the supercluster in η and ϕ
- A collection of shower shape variables which aim at providing information on the likelihood and location of a photon conversion and the degree of showering in the material:
 - The R_9 of the supercluster (as previously described in Section 2.2.2).
 - The ratio of the 5×5 crystal energy to the raw supercluster energy (equivalent of R_{25}).
 - The energy weighted η -width and ϕ -width of the supercluster (in other words the spread of the shower).
 - The number of basic clusters.
 - The ratio of energy in the **HCAL** behind the supercluster to the **ECAL** energy of the supercluster, H/E .
 - The ratio of the preshower energy to the raw supercluster energy (endcap only).
- A collection of the seed crystal and the seed cluster variables which aim at providing information about energy lost through gaps and cracks between crystal and crystal modules:
 - The relative energy and position of the seed cluster.
 - The local energy covariance matrix.
 - Energy ratios between the seed and the 3×3 and 5×5 areas around the seed.
 - The η and ϕ index of the seed crystal and the relative position of the seed cluster to the crystal centre.

- 649 • Additionally the number of primary vertices and the median energy density, ρ , (see
 650 Sec. 2.5) are included to account for residual energy scale effects from pileup.

651 The regression is trained using an additional piece to that described in Section 3.1
 652 whereby the target is to predict the full probability distribution of the ratio of the true
 653 energy to the raw energy, E_{true}/E_{raw} . The target is a double crystal ball distribution
 654 which consists of a Gaussian core and power law tails on either side. This can be fully
 655 parametrised by six variables, the Gaussian mean and width (μ, σ), the power parameters
 656 (n_L, n_R) and the power law tail cutoff parameters (α_L, α_R). Each of these parameters
 657 has a non-parametric dependence on the input variables, \vec{x} , and this is *learned* by the
 658 regression training whilst simultaneously minimising the likelihood,

$$-\ln \mathcal{L} = - \sum_{MCphotons} \ln p(E_{true}/E_{raw} | \mu(\vec{x}), \sigma(\vec{x}), \alpha_L(\vec{x}), \alpha_R(\vec{x}), n_L(\vec{x}), n_R(\vec{x})), \quad (3.5)$$

659 for the double crystal ball distribution, p . The most probable value for the true energy
 660 estimate of each photon is then given by,

$$E(\vec{x}, E_{raw}) = \mu(\vec{x}) E_{raw} \quad (3.6)$$

661 and the per-photon energy resolution is given by,

$$\frac{\sigma_E(\vec{x}, E_{raw})}{E(\vec{x}, E_{raw})} = \frac{\sigma(\vec{x})}{\mu(\vec{x})}. \quad (3.7)$$

662 In this way the regression predicts the full probability density of E_{true}/E_{raw} and
 663 provides an estimate of the optimal energy correction and the energy resolution per
 664 photon. A comparison of this distribution with a statistically independent MC sample is
 665 shown for the 8 TeV training in Figure 3.4. The performance of the regression compared
 666 to the default photon reconstruction is shown in Figure 3.5.

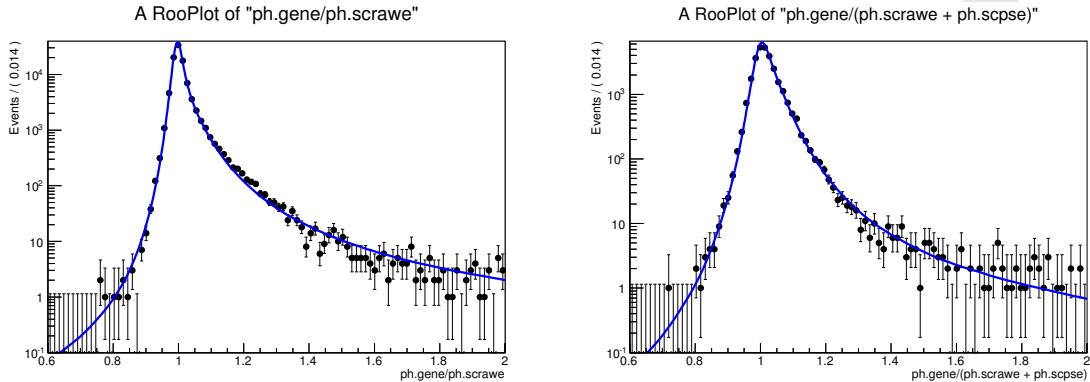


Figure 3.4.: A comparison of the predicted probability density of E_{true}/E_{raw} from the regression training (blue line) to the distribution in a statistically independent MC sample (black points) for barrel photons (left) and endcap photons (right).
PLOTS NEED TIDYING

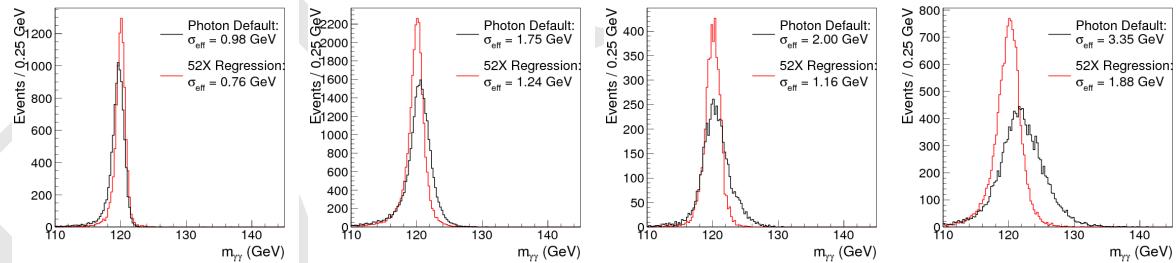


Figure 3.5.: A comparison of the regression performance compared to the photon default for the diphoton invariant mass on a MC source of Higgs decays to two photons. This is shown when both photons are in the barrel (left two plots) when both photons have $R_9 > 0.94$ (far left) and when at least one photon has $R_9 < 0.94$ (middle left). When at least one photon is in the endcap (right two plots) when both photons have $R_9 > 0.94$ (middle right) and when at least one photon has $R_9 < 0.94$ (far right). The value shown in the legend of each is the effective width, σ_{eff} , which is defined as half of the narrowest interval which contains 68.3% of the distribution. PLOTS NEED TIDYING.

667 **3.3.1. Correcting for residual discrepancies between data and
668 Monte Carlo**

669 After application of the energy regression correction there are some remaining discrep-
670 ancies between data and MC. These residual effects are accounted for using $Z \rightarrow e^+e^-$
671 events in data and simulation to correct the energy scale in the data and to apply an
672 additional smearing term to the MC with systematic uncertainties propagated through
673 the analysis to account for the uncertainties on these corrections.

674 **Energy scale corrections to the data**

675 The supercluster energy is identical for electrons and photons so by correcting the
676 supercluster energy scale to a known source, namely the mass of the Z-boson, in
677 dielectron decays the smaller residual energy scale effects are accounted for. This can
678 be done several times to account for various different effects. In the first stage scale
679 corrections are derived in bins of time (run range) and η , after applying these corrections
680 further, much smaller, residual effects are accounted for in bins of R_9 (the size of the
681 effect is different for converted and non-converted photons). After applying both of these
682 a further step is taken for the 8 TeV data in the barrel to derive residual corrections in
683 bins of E_T . Consequently the total scale correction is a product of three corrections in 59
684 bins of time \times 4 bins in η \times 2 bins in R_9 \times 6 bins in E_T (for the 8TeV barrel photons
685 alone).

686 The strategy for deriving these corrections is to take $Z \rightarrow e^+e^-$ events in data and
687 MC and extract the invariant dielectron mass in the relevant bin of interest. This
688 mass distribution is fitted with a convolution of a Breit-Wigner (designed to handle
689 the underlying Z line shape [24]) and a Crystal-Ball which models the calorimeter
690 resolution effects and bremsstrahlung losses in the material upstream of the ECAL.
691 The Breit-Wigner parameters are fixed to the PDG values of $M_Z = 91.188$ GeV and
692 $\Gamma_Z = 2.495$ GeV [24] whilst the Crystal-Ball parameters which model the detector effects
693 are allowed to float. The scale correction, ΔE , is then defined as the relative difference
694 between the Crystal-Ball peak in data and simulation,

$$\Delta E = \frac{m_{data} - m_{MC}}{m_Z}. \quad (3.8)$$

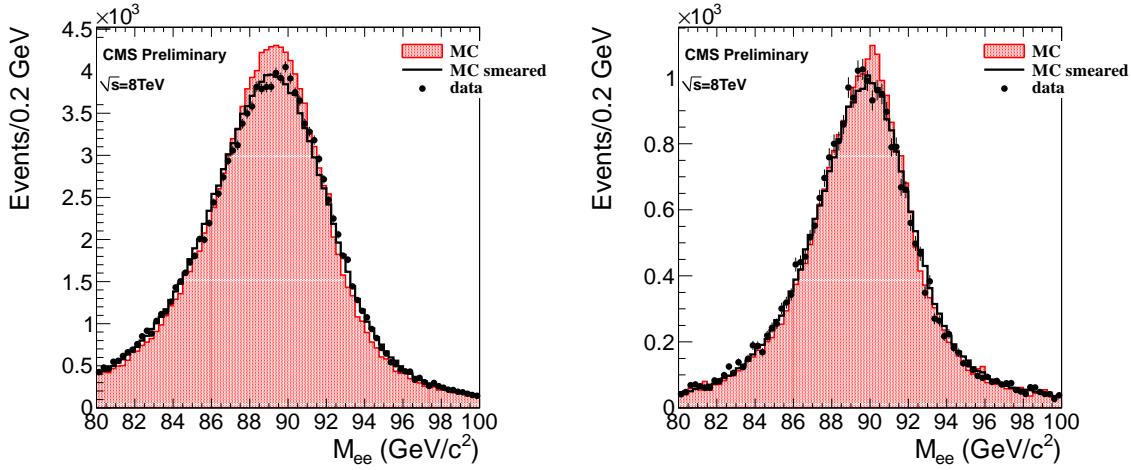


Figure 3.6.: The $Z \rightarrow e^+e^-$ invariant mass shape comparison before and after the scale and smearing corrections are applied. Shown for 8 TeV for photons with $1. < |\eta| < 1.444$ and $R_9 < 0.94$ on the left and for photons with $|\eta| < 1.$ and $R_9 > 0.4$ on the right.

695 Energy resolution smearing the Monte-Carlo

696 A similar method is used to extract a smearing factor that can be applied to the MC
 697 such that the width of the invariant mass distribution in $Z \rightarrow e^+e^-$ decays matches
 698 between data and MC. This is done in 4 bins of $\eta \times 2$ bins of R_9 and is parametrised as
 699 the quadratic sum of two resolution components: a constant term, ΔC , and a stochastic
 700 term, ΔS , which aims to model the expected resolution effects explained in Eq. 2.5 in
 701 Sec. 2.2.2. The smearing term, $\Delta \sigma$, is parametrised as,

$$\Delta \sigma = \frac{\Delta S}{\sqrt{E_T}} \oplus \Delta C. \quad (3.9)$$

702 The effect of the scale and smearing corrections is shown for the $Z \rightarrow e^+e^-$ data
 703 and MC samples in Figure 3.6 and for the analysis dataset after preselection with the
 704 electron veto inverted in Figures 3.7 and 3.8. Each of these corrections has an associated
 705 uncertainty and these uncertainties are propagated per photon through the analysis.
 706 There are also additional uncertainties included which account for differences between
 707 electrons and photons and the difference between the Z mass scale (around 90 GeV) and
 708 the Higgs mass scale (around 125 GeV). Systematic uncertainties are described in more
 709 detail in Section 5.4.

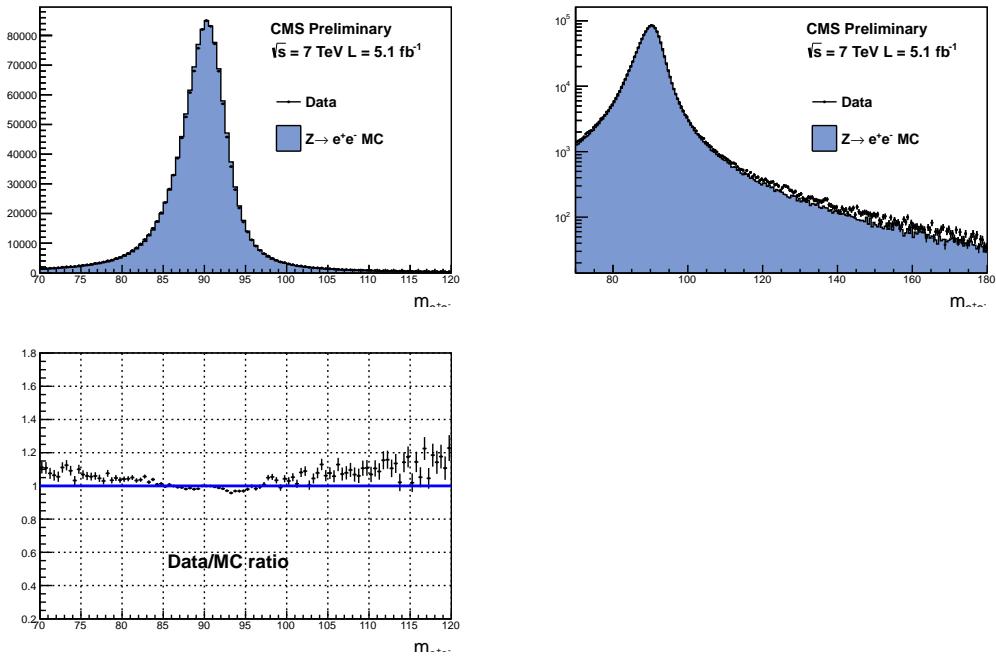


Figure 3.7.: The $Z \rightarrow e^+e^-$ invariant mass distribution at 7TeV in data (black points) and MC (blue histogram) for events which pass the analysis preselection in which the electron veto is inverted.

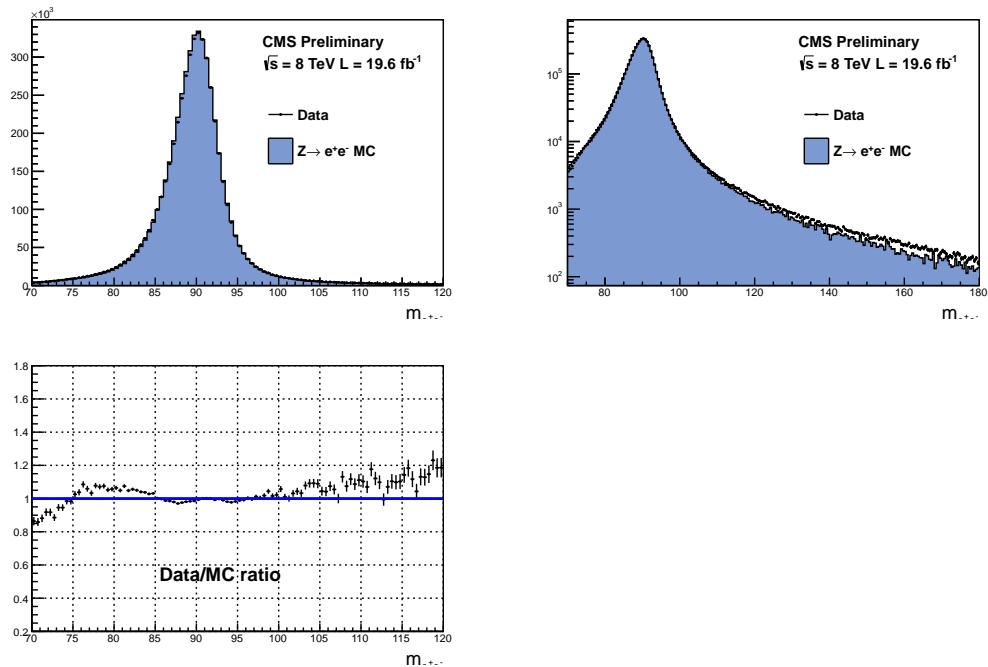


Figure 3.8.: The $Z \rightarrow e^+e^-$ invariant mass distribution at 8TeV in data (black points) and MC (blue histogram) for events which pass the analysis preselection in which the electron veto is inverted.

710 3.4. Vertex reconstruction

711 The resolution on the opening angle has a negligible effect if the correct vertex can be found
 712 within 10mm of the true interaction point. As seen in Sec. 3.2.2 the beamspot has an
 713 RMS spread of about 5 cm in the z direction and there is an average of ~ 20 vertices per
 714 bunch crossing. Because the beam direction is along the z -axis the spread of the vertex
 715 in the x and y directions is tiny (number?) and consequently mismeasurement of the
 716 primary vertex in the x - y plane is small and has no impact on the mass resolution. By
 717 assigning the correct vertex to the diphoton pair, using other information in the tracking
 718 system, most of the mass resolution can be preserved. The method used to extract the
 719 primary vertex is a BDT which exploits the correlation between the diphoton pair and
 720 the recoiling tracks from the underlying interaction as well as additional information
 721 in the tracking system from a photon conversion pair. The output of this per vertex
 722 BDT is evaluated for each vertex in the event and the primary vertex is assigned as the
 723 one with the highest value of BDT output (i.e. the value nearest 1.). In addition, it is
 724 possible to construct another BDT whose output is proportional to the probability that
 725 the chosen vertex is the correct one (described in Sec. 3.4.1). This probability becomes a
 726 useful discriminating variable for the analyses later on.

727 The vertex BDT uses the following input variables:

- 728 • $\sum_i |\vec{p}_T^i|^2$ - the sum of the transverse momentum squared of all of the tracks which
 729 originate from this vertex, representing how hard the interaction is at this vertex.
- 730 • $\sum_i (\vec{p}_T^i \cdot \frac{\vec{p}_T^{\gamma\gamma}}{|\vec{p}_T^{\gamma\gamma}|})$ - the sum of the magnitude of the transverse momentum of each track
 731 originating from this vertex relative to the transverse momentum of the diphoton
 732 system, representing the recoil of the tracks to the diphoton system.
- 733 • $(|\sum_i \vec{p}_T^i| - \vec{p}_T^{\gamma\gamma}) / (|\sum_i \vec{p}_T^i| + \vec{p}_T^{\gamma\gamma})$ - the asymmetry between the diphoton system and
 734 the other tracks originating from this vertex.
- 735 • $|z_v - z_c|/\sigma_c$ - this is added for events which contain at least one photon conversion
 736 where z_v is the z position of the vertex in question and z_c and σ_c are the estimated
 737 z position of the vertex from conversion information and its approximate error as
 738 defined below.

739 For events which contain at least one photon conversion, the conversion tracks and/or
 740 the conversion momentum can be used to point back to the beam line and estimate the

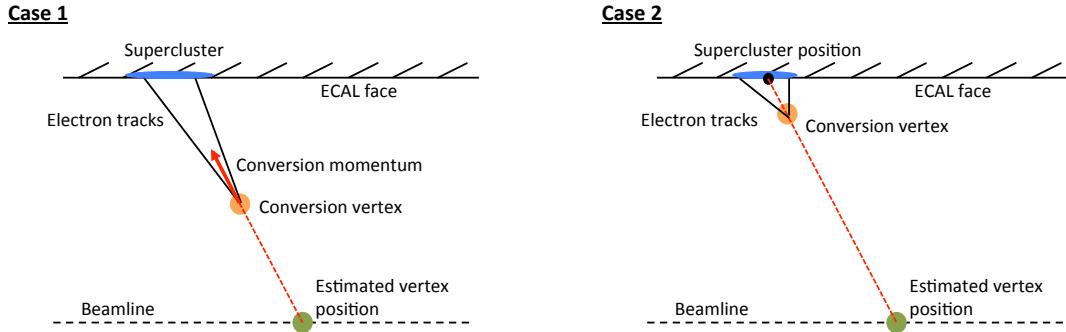


Figure 3.9.: A representation of the two methods for locating the primary vertex using photon conversion information. The left plot is for cases where the conversion occurs early enough in the tracker that the two electron tracks can be used to construct the converted pair momentum which is combined with the conversion vertex position to point back to the beam line. The right plot is for cases where the conversion occurs late in the tracker and the energy weighted supercluster position and the conversion vertex position are used to point back to the beam line.

vertex position. This can be achieved in one of two ways. In cases where the conversion occurs early, i.e. in one of the first layers of the tracking system, then the two electron legs of the conversion will leave two clean and distinct tracks. This means that the momentum of the conversion pair can be accurately reconstructed and used to point from the conversion vertex position back to the beam line and thus the nearest primary vertex. In cases where the conversion occurs late in the tracking system there are not enough track hits to accurately reconstruct the momentum of the conversion pair, however the incident position of the photon at the **ECAL** face is well known in this case, so the line which connects the **ECAL** position with the conversion vertex can be used to point back to the beam line. This is diagrammatically represented in Figure 3.9 for both cases.

For Case 1 conversions the primary vertex z position is calculated as,

$$z_c = z_{conv} - r_{conv} \cot(\alpha), \quad (3.10)$$

where z_{conv} is the z position of the conversion vertex, r_{conv} is the distance of the conversion vertex from the beam line and α is the angle between the beam line and the conversion momentum.

For Case 2 conversions the primary vertex z position is calculated as,

$$z_c = \frac{z_{conv} - r_{conv}}{(r_{SC} - r_{conv})(z_{SC} - z_{conv})}, \quad (3.11)$$

756 where z_{conv} and z_{SC} are the z positions of the conversion vertex and supercluster respec-
 757 tively, and r_{conv} and r_{SC} are the distance of the conversion vertex and the supercluster
 758 from the beamline.

759 There are 6 regions of the tracking system (refer back to Figure 2.2). When the
 760 conversion vertex is located in one of the inner regions; Pixel Barrel, Pixel Forward, TID,
 761 the Case 1 conversion information is included in the **BDT**, otherwise the Case 2 conversion
 762 information is used. The resolution on the primary vertex position in conversions is
 763 estimated per tracking region by calculating the effective width ¹ of the distribution of the
 764 difference between the z position of the primary vertex without conversion information
 765 and the z position of the primary vertex using conversion information alone, $\Delta z = z_v - z_c$.
 766 Consequently the fourth input variable to the **BDT**, shown in the list above as $|z_v - z_c|/\sigma_c$,
 767 is effectively a pull distribution for the conversion vertex. The **BDT** will favour vertices
 768 whose value of this variable is near zero.

769 The **BDT** is trained on a sample of $H \rightarrow \gamma\gamma$ **MC** events. It is tested with a statistically
 770 independent sample and further validated using $Z \rightarrow \mu^+\mu^-$ decays in data and **MC**.
 771 The efficiency is measured in data using the $Z \rightarrow \mu^+\mu^-$ channel where the muon tracks
 772 are removed from the **BDT** variables to simulate a diphoton like situation in data.
 773 The distributions of the input variables are shown for the $H \rightarrow \gamma\gamma$ training sample in
 774 Figure 3.10 (probably unnecessary). The **BDT** response is shown for $Z \rightarrow \mu^+\mu^-$ data
 775 and **MC** for both the signal (right vertex) and background (wrong vertex) in Figure 3.11.
 776 The chosen primary vertex is the one which gives the highest score **BDT** output. The
 777 efficiency of the vertex selection as a function of the $Z p_T$ and the number of reconstructed
 778 vertices as measured in $Z \rightarrow \mu^+\mu^-$ data and **MC** samples is shown in Figure 3.12.

779 3.4.1. Estimating the per-event probability that the correct 780 vertex is chosen

781 The total efficiency of assigning the correct vertex using the method described in the
 782 preceeding section is at the level of 75% during 2012 running conditions, where the

¹Half the narrowest interval which contains 68.3% of the distribution

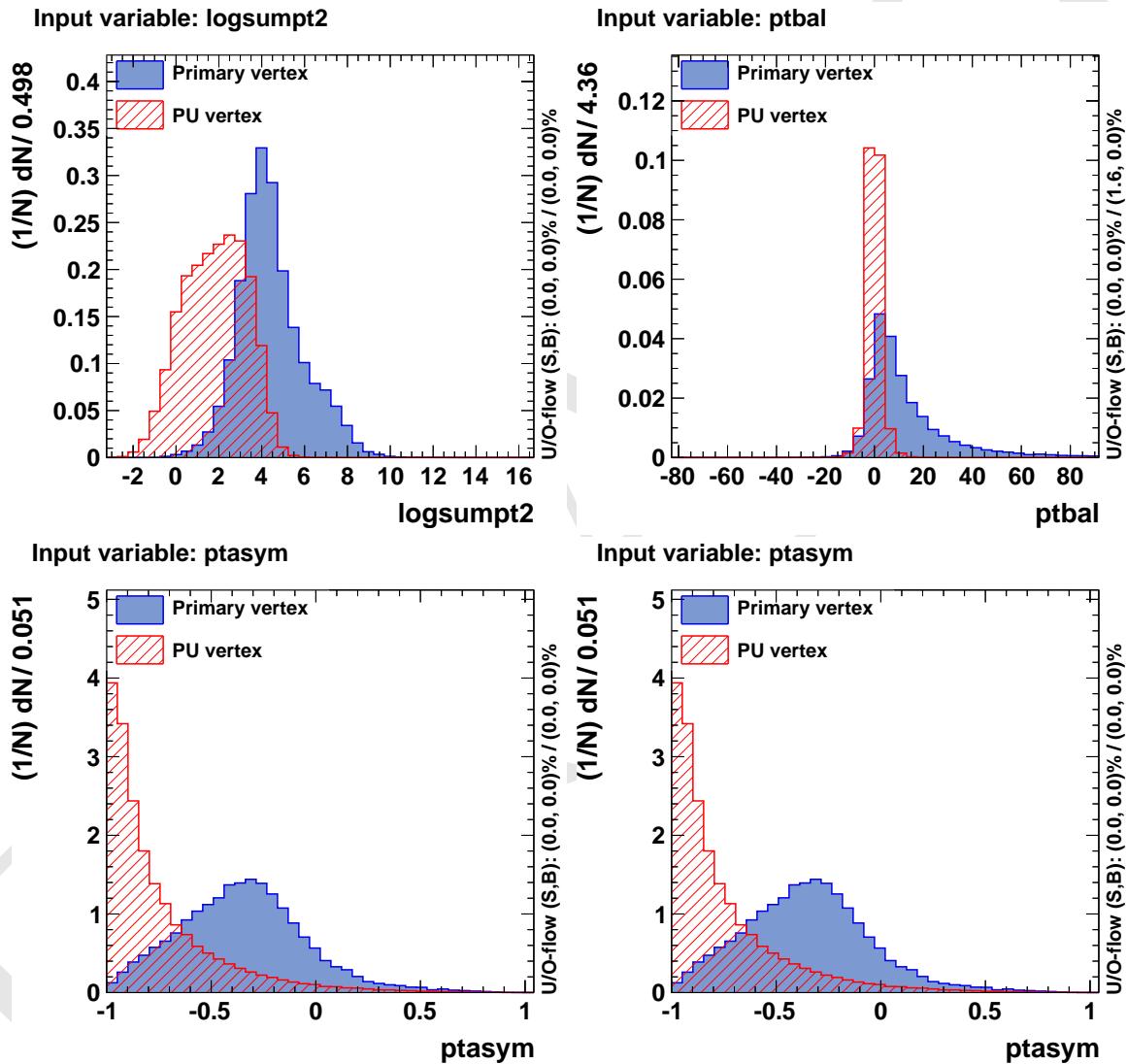


Figure 3.10.: Distributions of the input variables for the vertex BDT in the MC $H \rightarrow \gamma\gamma$ training (points) and test (filled) samples at 8TeV. Shown for the target primary vertex (blue histograms) and the background pileup vertices (red histograms). Plots need updating, tidying and labelling correctly.

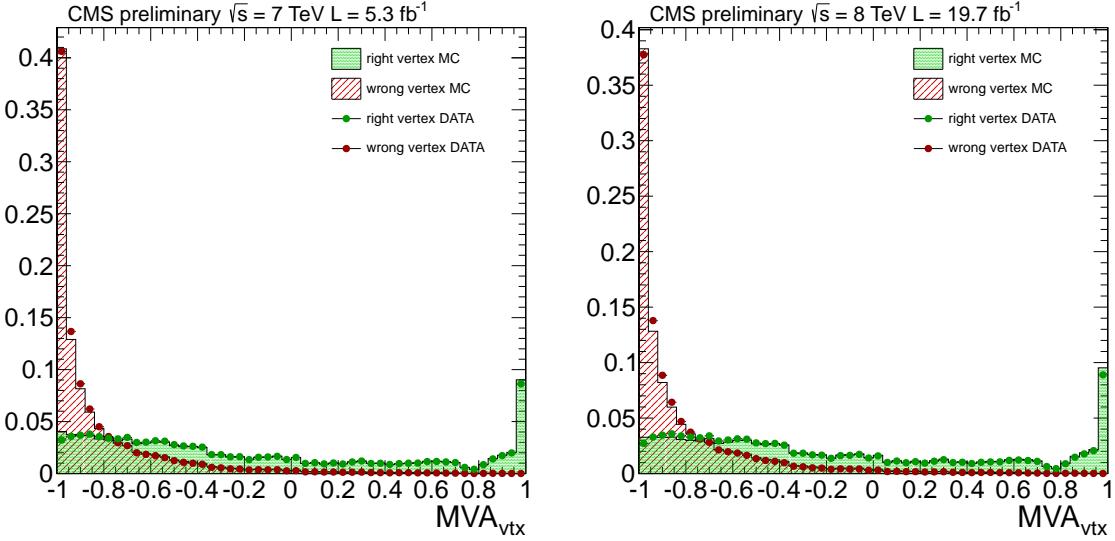


Figure 3.11.: The vertex **BDT** response for $Z \rightarrow \mu^+\mu^-$ events in data (points) and MC (filled) for the primary vertex (green) and the background pileup vertices (red).

783 correct vertex is defined as being within 10 mm of the true vertex. This means that
 784 for around 25% of preselected events the mass resolution is dominated by the vertex
 785 resolution. Consequently, it is important to ascertain the probability that the chosen
 786 vertex is the correct one. An additional specific **BDT** is constructed to address exactly
 787 this topic. The input variables used for this **BDT** are:

- 788 • The p_T of the diphoton system
- 789 • The number of vertices in each event
- 790 • The value of the per-vertex **BDT** described above
- 791 • The z distance, Δz , between the chosen vertex and the second and third choice
792 vertices.
- 793 • The number of photon conversions used, either 0, 1 or 2.

794 There is a linear relation between the response of this **BDT** and the correct vertex
 795 efficiency (or probability). This relationship is demonstrated in Figure 3.13 and is fitted
 796 with a single parameter straight line in order to analytically obtain the per-event correct
 797 vertex probability for a given event. Figure 3.14 shows that this estimation reproduces
 798 the required vertex efficiency as a function of Higgs p_T and number of reconstructed
 799 vertices.

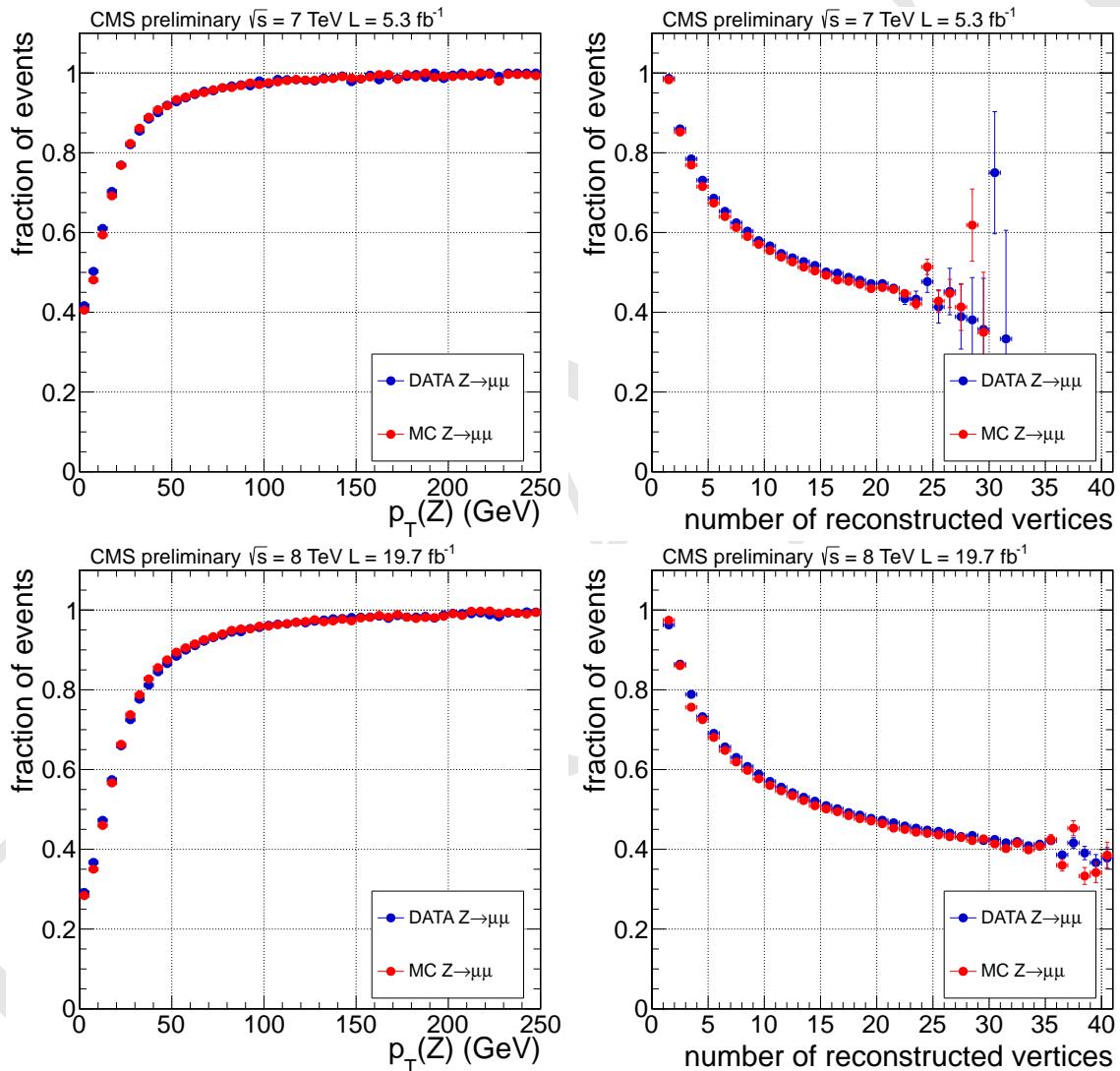


Figure 3.12.: The chosen vertex efficiency as measured in $Z \rightarrow \mu^+\mu^-$ data and MC as a function of $Z p_T$ (left) and number of reconstructed vertices (right) for the 7 TeV (top row) and 8 TeV (bottom row) data samples.

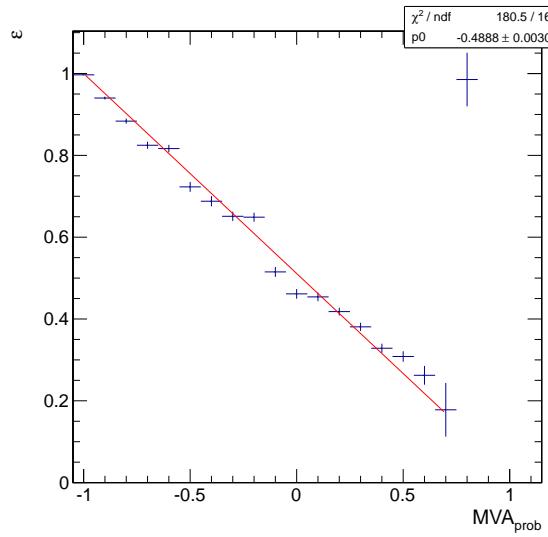


Figure 3.13.: A demonstration of the linearity relation between the per-event vertex probability BDT output distribution and the correct vertex probability

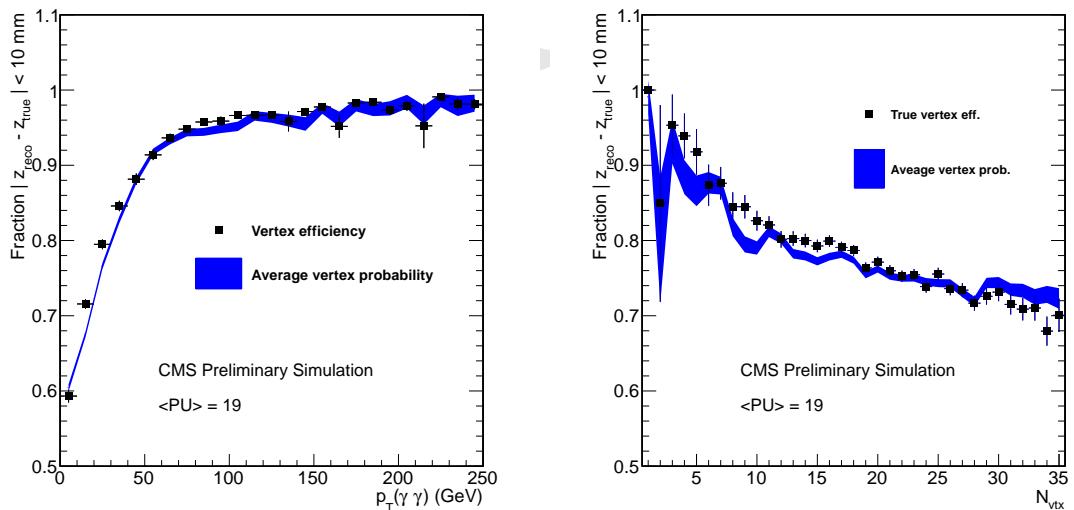


Figure 3.14.: A comparison of the true vertex efficiency (black points) and the average vertex probability (blue band) a statistically independent MC Higgs sample simulated with 2012 running conditions.

800 3.5. Event preselection

801 A simple and loose preselection is applied to all photons before entering the analysis.
 802 The preselection requirements are identical for all analysis approaches and are designed
 803 to remove some *fake* photons whilst maintaining near 100% trigger efficiency. The variables
 804 used for preselection are defined below and the preselection cuts are described in Table 3.1.

- 805 • **H/E - The ratio of hadronic energy in the **HCAL** tower behind the supercluster to the electromagnetic energy in the supercluster.** Neutral jets
 806 which fake photons typically leave a fraction of their energy in the **HCAL** so there
 807 is a requirement that the value of this variable is small.
- 808 • **$\sigma_{in\eta}^2$ - The RMS spread of the shower in the eta direction.** Multiple showers
 809 of a π^0 , or more than one π^0 , result in a wider shower in η (as the π^0 decay product
 810 photons are separated in space). This cannot be exploited in the ϕ direction because
 811 conversion electrons get separated by the magnetic field, however in η single photons,
 812 even when converted, occupy a narrow region. The separation of the two photons
 813 from a π^0 is minimal when they share the energy equally and given that typically
 814 $p_T \gg m_\pi$ the separation is close to minimal for most. Taking the transverse plane
 815 in the barrel, the separation $d = 2Rm_\pi/p_T$ (where R is the radius of the barrel),
 816 which for $p_T = 40\text{GeV}$ gives a value of $d = 8\text{mm} \sim 0.006\eta$. By referring to Table ??
 817 it is clear that the preselection requirement is quite loose.
- 818 • **ISO_{ECAL} - The total ρ corrected electromagnetic energy in a cone of
 819 radius 0.4 in (η, ϕ) around the photon candidate** - see Sec. 2.4
- 820 • **ISO_{HCAL} - The total ρ corrected hadronic energy in a cone of radius 0.4
 821 in (η, ϕ) around the photon candidate** - see Sec. 2.4
- 822 • **ISO_{Tracks} - The total ρ corrected track energy in a cone of radius 0.4 in
 823 (η, ϕ) around the photon candidate** - see Sec. 2.4
- 824 • **ISO_{PFC} - The total ρ corrected particle flow charged hadron energy in
 825 a cone of radius 0.4 in (η, ϕ) around the photon candidate** - see Sec. 2.4

826 In addition to the above an electron veto is applied to prevent contamination of the
 827 photon sample with electrons which originate from Drell Yan decays. This is achieved by
 828 removing photon candidates whose supercluster is matched to an electron track which
 829 has no missing hits in the innermost tracking region.

Table 3.1.: Preselection cut values.

		Barrel		Endcap	
R9	H/E	$\sigma_{in\eta}^2$	HoE	$\sigma_{in\eta}^2$	
≤ 0.9	< 0.075	< 0.014	< 0.075	< 0.034	
> 0.9	< 0.082	< 0.014	< 0.075	< 0.034	
Both Barrel and Endcap					
R9	ISO_{ECAL}	ISO_{HCAL}	ISO_{Tracks}	ISO_{PFCh}	
≤ 0.9	$< 4 \text{ GeV}$	$< 4 \text{ GeV}$	$< 4 \text{ GeV}$	$< 4 \text{ GeV}$	
> 0.9	$< 50 \text{ GeV}$	$< 50 \text{ GeV}$	$< 50 \text{ GeV}$	$< 4 \text{ GeV}$	

3.6. Using Z decays for validation and efficiency measurements

Whilst no known “standard candles” with high statistics exist for high p_T photons in the LHC environment a powerful control source for the $H \rightarrow \gamma\gamma$ decay in both data and MC is the $Z \rightarrow e^+e^-$ decay. From a detector view point electrons are very similar to photons and the Z is relatively near the relevant Higgs search range in mass. The differences between the Z and the Higgs, in both their mass and p_T distribution, and also the differences between electrons originating from a Z and photons originating from a Higgs are important systematic uncertainties on the Higgs mass scale and resolution. By inverting the electron veto usually applied in the preselection the Higgs to two photon analysis can be identically replicated but the very pure diphoton sample replaced with a pure dielectron sample. One additional process which can be used as a direct control for photons is the $Z \rightarrow \mu^+\mu^-\gamma$ decay although the statistics, even with the LHC luminosity, are very low. Many of the input variables used in training the BDTs and cuts are validated with both $Z \rightarrow e^+e^-$ and $Z \rightarrow \mu^+\mu^-\gamma$ data/MC comparison plots. An example is shown in Figure 3.15 for the reconstructed di-electron mass for events passing the preselection described in Table 3.1.

As previously shown (in Sec. 3.3.1) data/MC comparison of the $Z \rightarrow e^+e^-$ decay are used to derive scale corrections for the data and smearing of the MC. Discrepancies between data and MC in $Z \rightarrow e^+e^-$ decays of important analysis variables are accounted for by introducing systematic uncertainties to cover them. In addition the “tag and probe” method [25] is used on $Z \rightarrow e^+e^-$ decays to evaluate the signal efficiency for the preselection and analysis cuts. Several stages of the analysis are validated in this

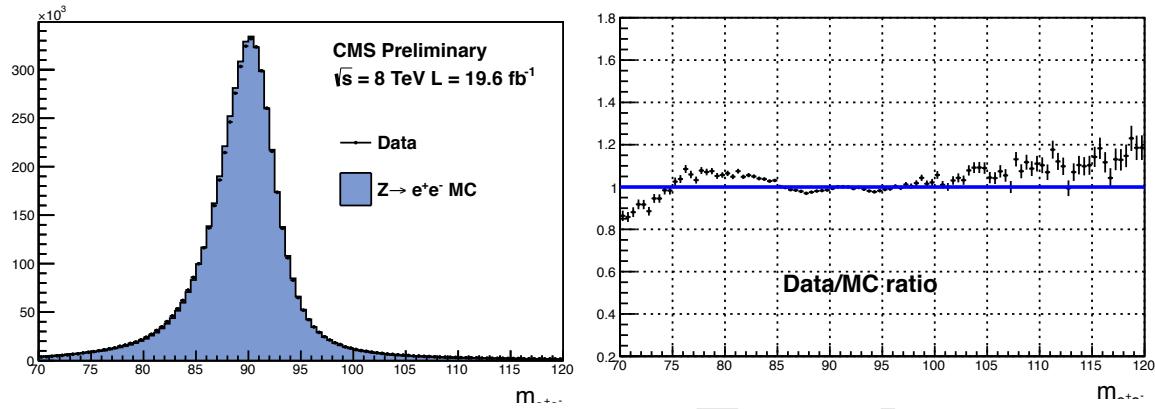


Figure 3.15.: The di-electron mass distribution after applying the preselection (described in Sec. 3.5) and the scale and smearing corrections (described in Sec. 3.3.1) whilst inverting the electron veto. The left plot shows the data at 8TeV as the black points with the Drell Yan MC sample as the blue histogram. The Data/MC ratio is shown in the right hand plot. **Plot needs updating**

way and where appropriate systematic uncertainties are included to account for any data/MC differences. Although the numbers and uncertainties themselves are derived from $Z \rightarrow e^+e^-$ samples (because of the much higher statistics), they are always cross checked with the $Z \rightarrow \mu^+\mu^-\gamma$ sample.

DRAFT

Chapter 4.

858 Selection and Categorisation

859 “Science never solves a problem without creating ten more.”

860 — George Bernard Shaw, 1856 – 1950

861 This thesis presents results of three different analyses used in the Higgs to two photon
862 search at CMS. The nominal results and properties are obtained from the so called Mass
863 Factorized MVA ([MFM](#)) analysis, which uses multivariate methods optimised specifically
864 to search for a Standard Model Higgs boson to select and categorise events. Two further
865 analyses are presented as cross checks to the baseline [MFM](#) analysis. The first, the Cuts
866 in Categories ([CiC](#)) analysis, is designed for simplicity and robustness as a cut based
867 approach and, owing to its low level of model dependence, is used for the spin analysis.
868 The second, the Sideband MVA ([SMVA](#)) analysis, serves to cross check the background,
869 the most significant unknown in a search like this, and categorisation by extracting
870 the background under the signal region from sidebands in the diphoton invariant mass
871 spectrum.

872 Both the [CiC](#) and [MFM](#) analyses have a fully parametric definition of the diphoton
873 invariant mass spectrum where the signal shape is derived from [MC](#) simulation. The
874 [SMVA](#) uses a cut and count method.

875 4.1. Event selection

876 There are two complimentary event selections used in the three analyses. The [CiC](#)
877 analysis uses a cut based photon selection whilst the [MFM](#) and [SMVA](#) analyses use

878 a series of **BDTs** to first select photons and then select events. The ultimate aim is
 879 a selection which accepts two prompt photons (i.e. does not contain any fakes) and
 880 exploits regions of phase space which have a narrow mass resolution and high signal to
 881 background ratio.

882 All events must contain at least two photons which pass $p_T^{\gamma 1} > m_{\gamma\gamma}/3$ (for the leading
 883 photon) and $p_T^{\gamma 2} > m_{\gamma\gamma}/4$ (for the subleading photon) and the invariant mass of the
 884 diphoton pair must be in the range $100 < m_{\gamma\gamma} < 180$. In the case where there are more
 885 than two photon candidates in an event, the two photons which give the highest scalar
 886 sum p_T are chosen.

887 4.1.1. Selection using cuts in categories

888 Cuts are optimised for photons in four non-overlapping categories to take advantage of
 889 the different photon energy resolution between the barrel and the endcap and between
 890 converted and non-converted photons. The categories are defined as $|\eta| < 1.444$ or
 891 $|\eta| > 1.556$ (i.e. either barrel or endcap) and $R_9 < 0.94$ or $R_9 \geq 0.94$ (i.e. either converted
 892 or non-converted) and the cuts are chosen by targeting a specific signal to background
 893 ratio (S/B). The procedure to define the cuts is as follows,

- 894 • A set of loose cuts are defined as the starting values.
- 895 • The “N-1” distribution of each cut variable is produced. This is the distribution of
 896 each variable after the cuts on all other variables have been applied.
- 897 • A smooth curve is fitted to the distribution of the B/S ratio versus the cut variable.
- 898 • The cut value is chosen by evaluating the curve for the required value of S/B .
- 899 • Do the same for all the other variables.

900 Consequently a stable set of cut values is obtained by iterating this procedure a few
 901 times. Each cut will then select events with the same purity (S/B) and thus the efficiency
 902 for selected photons is maximised for the chosen purity level. As a result the cuts in
 903 the endcap are much tighter than in the barrel, similarly the cuts for low R_9 photons
 904 are much tighter than the cuts for high R_9 photons. The cut variables are described
 905 below and the cut values shown in Table 4.1. The cut setting procedure is optimised on
 906 a signal sample of $H \rightarrow \gamma\gamma$ MC (with $m_H = 120\text{GeV}$) and a background sample of $\gamma+\text{jet}$

	Barrel		Endcap	
	$R_9 > 0.94$	$R_9 < 0.94$	$R_9 > 0.94$	$R_9 < 0.94$
PF isolation sum, chosen vertex	<6	<4.7	<5.6	<3.6
PF isolation sum, worst vertex	<10	<6.5	<5.6	<4.4
Charged PF isolation sum	<3.8	<2.5	<3.1	<2.2
$\sigma_{i\eta i\eta}$	<0.0108	<0.0102	<0.028	<0.028
H/E	<0.124	<0.092	<0.142	<0.063
R_9	>0.94	>0.298	>0.94	>0.24

Table 4.1.: Photon ID selection cut values. The cuts are applied to both the leading and subleading photons.

907 events. The photon efficiency of the cuts for the four different classes of photon is shown
 908 in Figure 4.1 as a function of the supercluster position, η , and the photon p_T .

- 909 • **PF isolation sum, chosen vertex** - Sum of particle flow photon, charged and
 910 neutral hadron isolation sums defined in Sec. 2.4 where all PF candidates originate
 911 from the primary vertex selected in Sec. 3.4
- 912 • **PF isolation sum, worst vertex** - As above but where all PF candidates originate
 913 instead from the vertex which gives the largest charged hadron PF isolation sum.
 914 This add protection for cases when the primary vertex is incorrectly assigned.
- 915 • **Charged PF isolation sum** - Described above in Sec. 2.4
- 916 • **$\sigma_{i\eta i\eta}$** - Described above in Sec. 3.5
- 917 • **H/E** - Described above in Sec. 3.5
- 918 • **R_9** - Described above in Sec. 2.2.2

919 4.1.2. Photon ID MVA

920 For the multivariate approach a BDT is trained to discriminate between prompt photons
 921 and jets. The desire is to factorise the photon selection, which is required to distinguish
 922 prompt photons from neutral mesons faking photons, and the event selection, which
 923 is required to consider kinematics, resolution etc. to distinguish $H \rightarrow \gamma\gamma$ from the
 924 $pp \rightarrow \gamma\gamma, \gamma + jet, jet + jet$ background. The input variables used are designed specifically

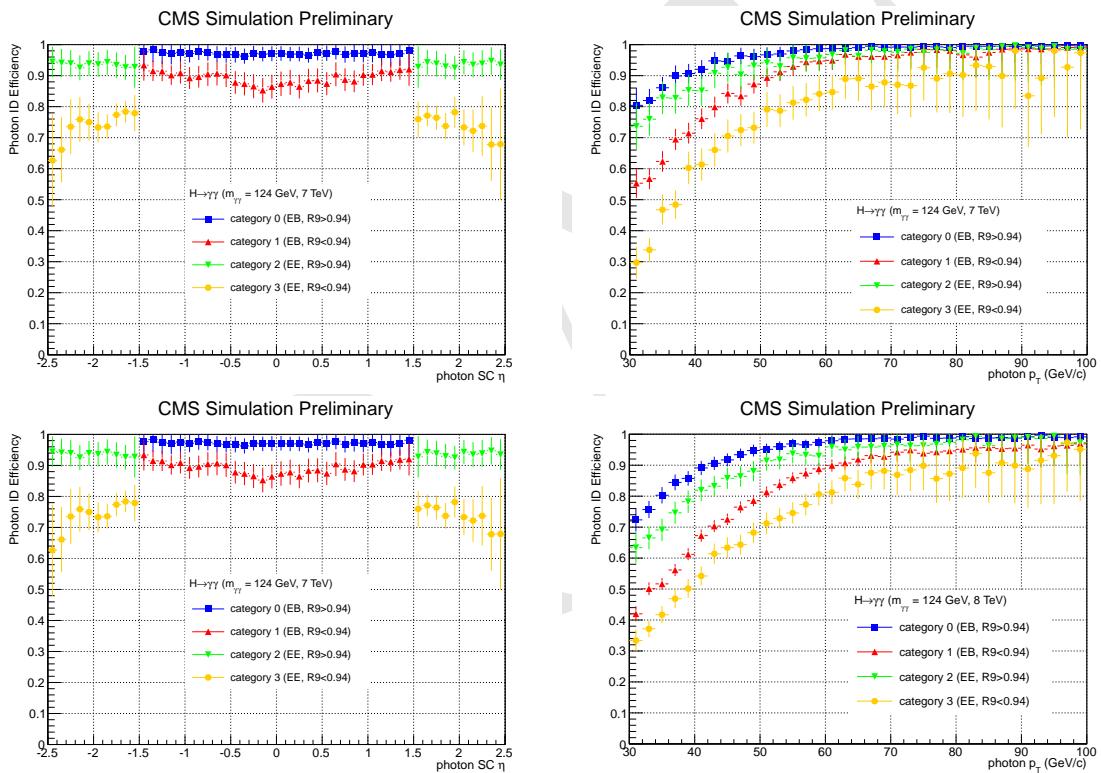


Figure 4.1.: Cut based photon ID efficiency as measured in $Z \rightarrow e^+e^-$ tag and probe. MORE DESCRIPTION

to distinguish between photons and fakes and to decouple from any event or photon kinematics which make the identification Higgs specific. Consequently, the training samples used are γ +jet samples where the identification **BDT** signal is the prompt γ and the background is the fake jet. The p_T and supercluster η distribution of the prompt photon are reweighted to match the background non-prompt photons in order to further negate any photon kinematics which the **BDT** could exploit. The photon ID **BDT** is trained separately for the barrel and endcap in both 7 and 8 TeV as these regions of phase space are so different. The result is four separate trainings in total.

The input variables aim to exploit differences in the shower shape and isolation between prompt and non-prompt photons and the correlation between these variables and the supercluster position and energy. They are,

Shower shape variables

- $\sigma_{i\eta i\eta}$ - Explained above in Sec. 3.5
- $\sigma_{i\eta i\phi}^2$ - The equivalent diagonal spread (in η, ϕ) of the shower.
- $E_{2 \times 2}/E_{5 \times 5}$ - Ratio of energy in the most energetic 2×2 cluster which contains the seed to the energy in the 5×5 cluster.
- R_9 - Explained above in Sec. 2.2.2.
- σ_η - The energy weighted standard deviation of single crystal eta within the supercluster.
- σ_ϕ - The energy weighted standard deviation of single crystal phi within the supercluster.
- σ_{xy} (for endcap only) - The standard deviation of the shower spread in the x, y planes of the preshower.

Isolation variables

- PF Photon ISO - Particle flow photon isolation sum
- PF Charged Hadron ISO (selected vertex) - Particle flow charged hadron isolation sum for candidates originating from selected vertex.
- PF Charged Hadron ISO (worst vertex) - Particle flow charged hadron isolation sum for candidates originating from the vertex with the largest isolation sum.

954 **Correlation variables**

- 955 • ρ - The median energy density in the event.
- 956 • η - The η position of the photon supercluster.
- 957 • E_{raw} - The raw energy of the photon supercluster.

958 The testing sample used to verify the output of the photon identification **BDT** is a
 959 **MC** $H \rightarrow \gamma\gamma$ sample ($m_H = 124\text{GeV}$). The photon identification **BDT** output (which is
 960 shown for the four different trainings in Figure 4.2) provides a measure of an individual
 961 photons’ “quality” and is used as an input to the event level **BDT** (described in the next
 962 section). Even so, a considerable amount of background is cut out by defining a loose
 963 cut (the **BDT** output must be > -0.2) on the photon ID **BDT** output value which is
 964 more than 99% signal efficient.

965 The photon ID **BDT** response for each photon is used as a direct input to the event
 966 level **MVA**. Given that imperfect modeling of the detector response can result in a small
 967 change in the photon ID response which has a direct impact on the event level **MVA**
 968 response, which is used to classify events, a systematic on the photon quality is applied
 969 and propagated through to the event level **MVA**. Validation of the photon ID **BDT**
 970 response in the $Z \rightarrow e^+e^-$ decay is shown in Figure 4.3 with the size of the systematic
 971 shown. Validation is done with the $Z \rightarrow e^+e^-$ decay on all of the input variables, these
 972 distributions are shown in Appendix A (this may not be entirely necessary)

973 **4.1.3. Diphoton event level MVA**

974 Whilst the CiC analysis selects events based on photon identification, the **MVA** analysis
 975 approach is to first select photons using the photon identification **BDT** described in the
 976 section above and then pass all the relevant event information through an event level
 977 **BDT**. The event level classifier, referred to as the diphoton **BDT**, is constructed to give
 978 a high score to events which fulfill the following criteria:

- 979 1. The event kinematics should be compatible with a Higgs decay.
- 980 2. The event has good mass resolution.
- 981 3. The event contains two “high quality” photons (i.e. they have a high score from the
 982 photon ID **BDT**).

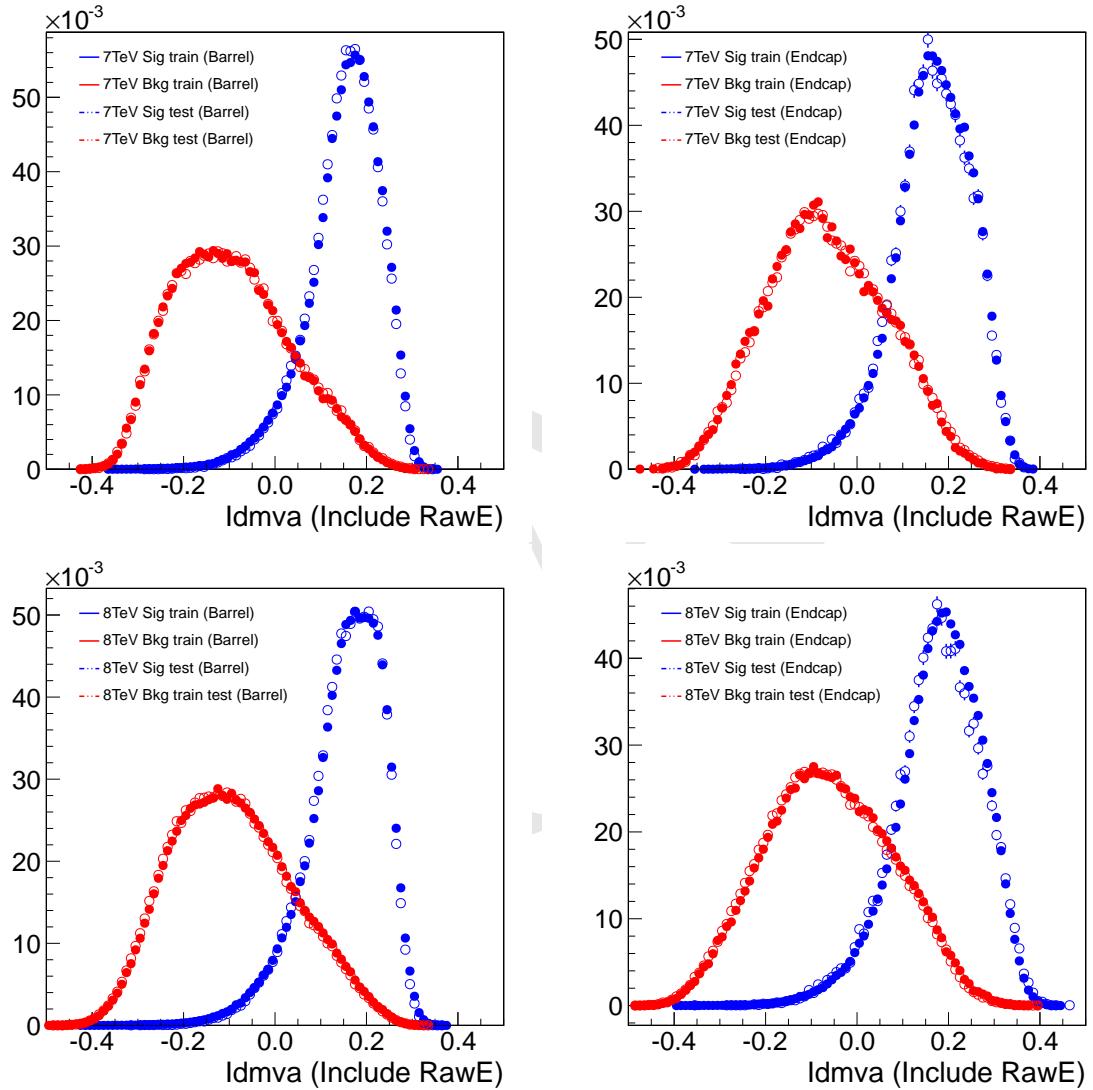


Figure 4.2.: The output distribution of the photon identification BDT for 7TeV barrel (top left), 7TeV endcap (top right), 8TeV barrel (bottom left) and 8TeV endcap (bottom right). The solid points show the $\gamma + \text{jet}$ training sample distributions and the hollow points show the $H \rightarrow \gamma\gamma$ test sample distributions for prompt signal photons in blue and fake background photons in red. A cut of > -0.2 is made on all photons.

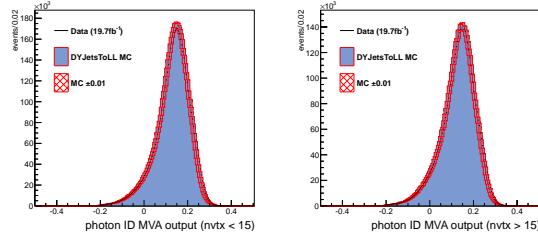


Figure 4.3.: The output distribution of the photon identification **BDT** for the 8 TeV training as validated by the $Z \rightarrow e^+e^-$ decay. The data is shown as the black points with the **MC** as the blue histogram. The systematic uncertainty on the output as applied to the **MC** is shown as the red band. **This is literally just a place holder for the plot that I want to show.**

It is highly important that the **BDT** is completely independent of Higgs mass and that the input variables have no, or at the least very little, dependence on the Higgs mass. This is essential to have a fair training. For example if the **BDT** included the Higgs mass, or a variable highly correlated with it, it would preferentially select events with this mass therefore biasing the selection towards events which have a mass near the mass of the signal used to train with. The input variables used are,

Event kinematics

- $p_T^{1(2)}/m_{\gamma\gamma}$ - The mass relative transverse momenta of each photon.
- $\eta^{1(2)}$ - The pseudorapidity of each photon.
- $\cos(\phi_1 - \phi_2)$ - The cosine of the angle between the two photons in the transverse plane. This variable reflects the p_T of the diphoton system (in other words the reconstructed Higgs candidate) without introducing a mass dependence.

Mass resolution

- $\sigma_m^{right}/m_{\gamma\gamma}$ - The mass resolution of the event assuming the correct primary vertex has been selected. In this case the angular resolution is negligible so the variable is calculated using just the two photon energy resolution values as,

$$\frac{\sigma_m^{right}}{m_{\gamma\gamma}} = \frac{1}{2} \left(\frac{\sigma_{E_1}}{E_1} \oplus \frac{\sigma_{E_2}}{E_2} \right). \quad (4.1)$$

- $\sigma_m^{wrong}/m_{\gamma\gamma}$ - The mass resolution of the event assuming the wrong vertex is selected. The vertex position in z is distributed as a Gaussian with a width equivalent to $\sqrt{2}\sigma_z^{beamspot}$ and so the angular resolution σ_m^{vtx} can be analytically calculated given

the ECAL impact positions of the two photons. Consequently the wrong vertex variable is calculated as,

$$\frac{\sigma_m^{wrong}}{m_{\gamma\gamma}} = \frac{\sigma_m^{right}}{m_{\gamma\gamma}} \oplus \frac{\sigma_m^{vtx}}{m_{\gamma\gamma}}. \quad (4.2)$$

- p_{vtx} - The probability that the selected primary vertex is correct. In order to tie together the mass resolution information given the right vertex hypothesis and the wrong vertex hypothesis, the probability that the vertex is correct is used in addition.
- It is also important to specify in the training that the signal to background ratio is inversely proportional to the mass resolution. Accordingly the signal events in the training are weighted by a factor,

$$w = \frac{p_{vtx}}{\sigma_m^{right}/m_{\gamma\gamma}} + \frac{1 - p_{vtx}}{\sigma_m^{wrong}/m_{\gamma\gamma}}. \quad (4.3)$$

1000 Photon quality

- $phoID^{1(2)}$ - The photon ID BDT output value of each photon.

The training is performed separately for 7 and 8 TeV and the samples used for the signal are all of the Standard Model $H \rightarrow \gamma\gamma$ MC samples (**ggH**, **VBF**, **VH**, **tH**) appropriately weighted by cross section and the samples used for the background are the cross section weighted mixture of Standard Model backgrounds which include contributions from $pp \rightarrow \gamma\gamma$ (prompt-prompt), $pp \rightarrow \gamma + \text{jet}$ (prompt-fake) and $pp \rightarrow \text{jet} + \text{jet}$ (fake-fake) as described in Sec. 3.2.1. The training is performed on only half of the event samples (selected by even event number) so that the BDT response can be tested on the other half (selected by odd event number).

A cut is placed on the BDT response in order to remove almost all the events which contain two fake photons. The remaining events which pass this cut are categorised in coarse bins based on the BDT response. The strategy for optimising this cut value and the category boundaries is explained in Section 4.2. The BDT response in data, background and signal is shown in Figure 4.4. The cut values are > 0.19 for the 7 TeV training and > -0.78 for the 8 TeV training. It can be seen from this figure that the cut on the event level BDT removes nearly all of the fake-fake contribution to the background whilst the remainder consists of about 70% prompt-prompt and 30% prompt-fake.

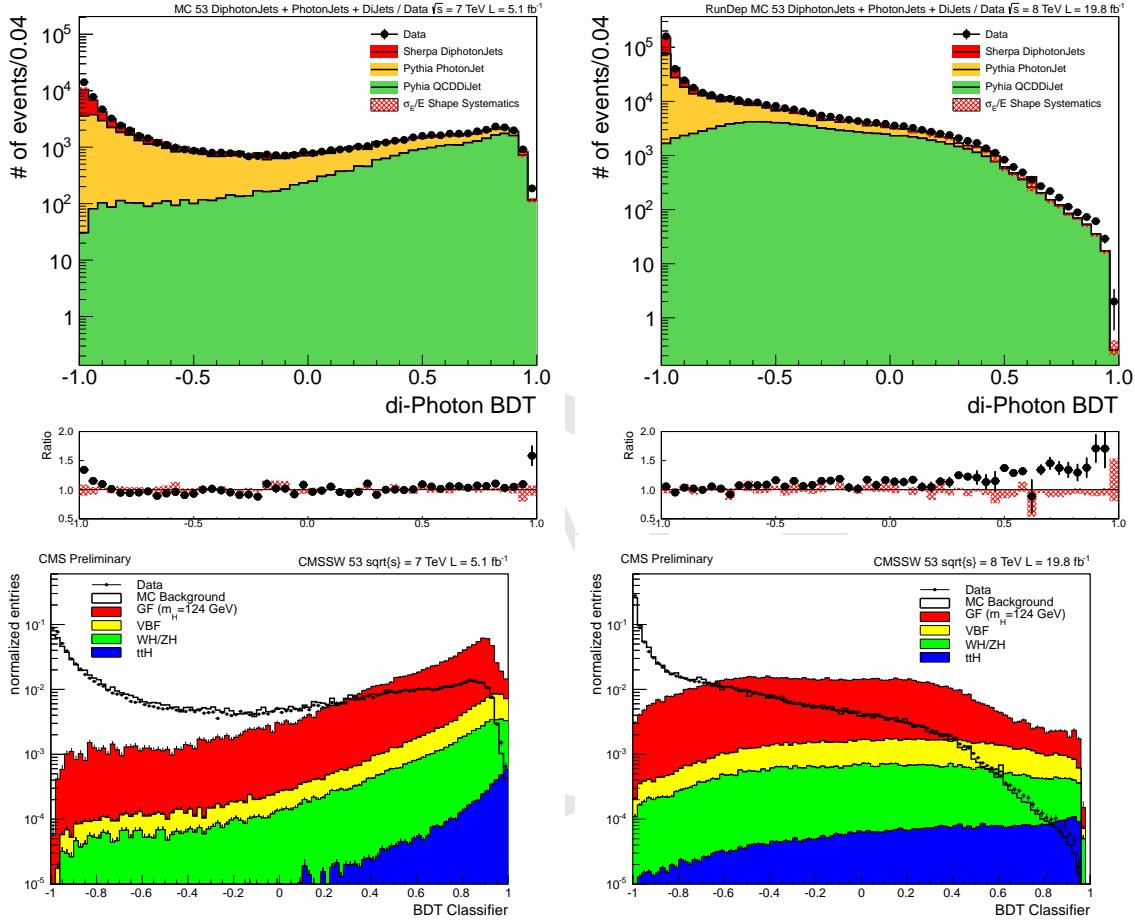


Figure 4.4.: The diphoton BDT response for the 7 TeV training (left column) and 8 TeV training (right column). The data and background distributions are shown in the plots on the top row for data in the range $100 < m_{\gamma\gamma} < 120 \text{ GeV}$ and $130 < m_{\gamma\gamma} < 180 \text{ GeV}$ (black points) and for the prompt-prompt background (green), prompt-fake background (yellow) and fake-fake background (red). The signal distributions are shown in the plots on the bottom row for gluon fusion (red), vector boson fusion (yellow), associated W, Z production (green) and associated $t\bar{t}$ production (blue) alongside the data in the range $100 < m_{\gamma\gamma} < 120 \text{ GeV}$ and $130 < m_{\gamma\gamma} < 180 \text{ GeV}$ (black points) and the total background (hollow histogram). **PLOTS NEED UPDATING**

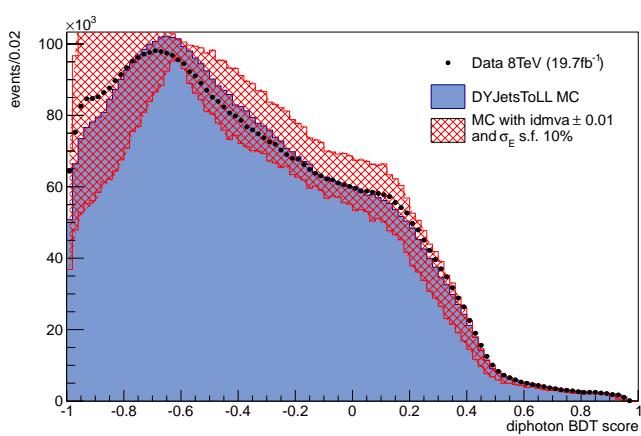


Figure 4.5.: The diphoton BDT response for the 8 TeV training in the $Z \rightarrow e^+e^-$ decay. The data is shown as the black points and the MC as the blue histogram. The systematic applied to account for variation in the BDT response from mismodelling in the photon quality response and the photon energy resolution estimate are shown as the red band. **This is just a placeholder for the plot I want to show**

Whilst the background in this analysis is obtained in a completely data driven way the signal model is obtained from MC. Consequently any uncertainties which effect the shape of the output distribution of this BDT in the signal result in event migrations between the final analysis categories. The input variables whose uncertainties have the largest effect on the BDT response in signal are the photon ID quality and the photon energy resolution estimate. This is because these variables have both i) relatively large uncertainties because of imperfect detector response modelling in the simulation and ii) are highly discriminative and hence can have a relatively large impact on the BDT response. As these variables both vary monotonically with the BDT response the systematic uncertainty on them is propagated through the analysis as an event migration. Systematic uncertainties are described in more detail in Sec. 5.4. The size of this effect is shown in the BDT response validation plot with the $Z \rightarrow e^+e^-$ decay in Figure 4.5. Validation plots using the $Z \rightarrow e^+e^-$ channel for each of the event level BDT input variables is shown in Appendix B (may not be necessary).

4.2. Event categorisation

In order to exploit different regions of phase space with similar signal to background ratios the events are split into categories. Furthermore, additional categories can be

designed to enrich the selection with events characteristic of particular Higgs production modes. The **VBF** production mode is typically accompanied by a pair of jets with a large pseduorapidity separation, the **VH** production may be accompanied by a charged lepton, missing transverse energy or jets originating from the decay of the associated W or Z boson. Similarly, **t $\bar{t}H$** production may be accompanied by b-quarks and or charged leptons. The predominant production mode, which accounts for about 88% of the signal, is **ggH** which is inclusive. The amount of signal produced from exclusive modes is approximately 8% for **VBF**, 4% for **VH** and <1% for **t $\bar{t}H$** . By including a series of event tags, and separating events accordingly, all four production modes of the Higgs at the **LHC** are harnessed in this analysis. This not only helps to increase the overall sensitivty to a Standard Model Higgs boson (given the very low background rates expected for the exclusive modes) but also significantly helps to reduce the error on measurements of the observed bosons respective couplings with fermions and bosons as the signal from the relative production modes gets split into distinct categories. The inclusive mode categorisation is done differently for each of the three analyses, a summary of all the event classifications is given in Table 4.4 at the end of the chapter.

4.2.1. Exclusive mode tagging

All events which make it to this stage will have passed the photon preselection (Sec. 3.5), as well as the basic requirement of two high p_T photons with $100 < m_{\gamma\gamma} < 180$ (Sec. 4.1), and either the CiC selection (Sec. 4.1.1) or the MVA selection (Sec. 4.1.3) and these make up the final event sample. They now pass through the “tagging” procedure described below in which they are organised into a set of non overlapping event classes. The tagging is done in a specific order to ensure there is no overlap between classes and the order is chosen such that preference is given to categories with a higher expected signal to background ratio. This order is shown alongside a summary of the relevant cut values in Table 4.4 at the end of the chapter. The dijet selection for the **MVA** analyses uses a specific combined dijet-diphoton **BDT** to exploit the very particular kinematic properties of **VBF** whereas for the **CiC** analysis this is left as a cut based criteria. If an event does not meet the requirements of a particular tag it is passed onto the next tag and if it fails all tag requirements it is placed in one of the inclusive categories, whose structure is described below (Secs. 4.2.2- 4.2.4), meaning that no event can now not be included in the analysis.

Variable	tight category	loose category
$p_T^{\gamma_1}/m_{\gamma\gamma}$	> 0.5	> 0.5
$p_T^{\gamma_2}$	$> 25 \text{ GeV}$	$> 25 \text{ GeV}$
$p_T^{j_1}$	$> 30 \text{ GeV}$	$> 30 \text{ GeV}$
$p_T^{j_2}$	$> 30 \text{ GeV}$	$> 20 \text{ GeV}$
$ \Delta\eta_{j_1 j_2} $	> 3.0	> 3.0
$ Z $	< 2.5	< 2.5
$M_{j_1 j_2}$	$> 500 \text{ GeV}$	$> 250 \text{ GeV}$
$ \Delta\phi(jj, \gamma\gamma) $	> 2.6	> 2.6

Table 4.2.: Final selection cuts for the VBF selection. Events from the first category are removed from the second one.

1067 Dijet tagged categories for VBF

1068 The following variables are used to exploit the specific topology of the jet pairs associated
 1069 to VBF Higgs production,

- 1070 • $p_T^{\gamma}/m_{\gamma\gamma}$ - The mass relative transverse momenta of the leading and subleading
 1071 photons,
- 1072 • $p_T^j/m_{\gamma\gamma}$ - The mass relative transverse momenta of the leading and subleading jets,
- 1073 • m_{jj} - The dijet invariant mass,
- 1074 • $\Delta\eta_{j_1 j_2}$ - The pseudorapidity difference between the two jets.
- 1075 • $Z = \eta(\gamma_1 + \gamma_2) - (\eta(j_1) + \eta(j_2))/2$ - The so called *Zeppenfeld* variable [26],
- 1076 • $\Delta\phi_{j_1 j_2}$ - The angular difference between the two jets in the transverse plane.

1077 Additionally the leading photon p_T requirement is raised to $p_T^{\gamma_1} > m_{\gamma\gamma}/2$. The energy
 1078 measurement of jets in the event are calibrated to correct for detector effects [27] and
 1079 additional energy in the jets from pileup is removed using the FASTJET jet areas technique
 1080 described in [28–30]. Jets are required to be within the pseudorapidity range, $\eta < 4.7$.

1081 For the CiC analysis the dijet tagging is done in a cut based way and tagged events
 1082 are placed into two categories as per the cuts defined in Table. 4.2.

For the **MVA** analyses the dijet tagging is done with use of two additional **MVA**s. The first is designed to exploit the **VBF** kinematic properties and the second is used to combine this information with the diphoton **BDT**. The input variables for the kinematic dijet **BDT** are identical to those used in the cut-based **VBF** selection where candidates are required to pass a **VBF** preselection of two jets with $p_T^{j1} > 30$ GeV and $p_T^{j2} > 20$ GeV and invariant mass, $m_{jj} > 75$ GeV. The signal sample used for training is the Standard Model **MC** with just **VBF** production, whilst the Standard Model gluon fusion **MC** is included as background along with the usual prompt-prompt, prompt-fake and fake-fake contributions. This helps produce an output in which a high score gives a very pure **VBF** sample.

The combined dijet-diphoton **BDT** has inputs of the kinematic dijet **BDT** output, the diphoton **BDT** output and $p_T^{\gamma\gamma}/m_{\gamma\gamma}$ in order to discriminate **VBF** from both the other signal types and the background, utilising all the information available in the event. The mass relative transverse momenta of the diphoton system, $p_T^{\gamma\gamma}/m_{\gamma\gamma}$, is included because of its considerable correlation with both the dijet **BDT** output and the diphoton **BDT** output.

One finds that the background rejection for **VBF** is significantly improved by the use of the combined dijet-diphoton **BDT** whilst the **VBF** purity (i.e. separation from **ggH**) is not as good when collapsing the kinematic dijet **BDT** and the combined **BDT** into one step. Consequently the trainings are performed separately and the **VBF** categories are defined by picking out regions which have a high score in the combined dijet-diphoton **BDT** response. The optimisation procedure for deciding the category boundaries is analogous to the one used for the inclusive categories in the **MFM** analysis where the target is to minimise the expected uncertainty of the signal strength from the **VBF** process alone when moving the boundaries around. It is explained further in Section 4.2.3. At 8 TeV there are three **VBF** categories and at 7 TeV there are two. The output distributions for the signal and background with the cut boundaries are shown for the two dijet **BDT**s are shown in Figures 4.6 and 4.7. Each successive **BDT** training uses statistically independent **MC** samples to avoid selection bias from fluctuations in the simulation.

1113 Lepton, jet and \cancel{E}_T tagged categories for **VH**

1114 The selection for the four categories designed to tag **VH** production are optimised by
1115 minimising the expected uncertainty on the signal strength of this process alone. Two

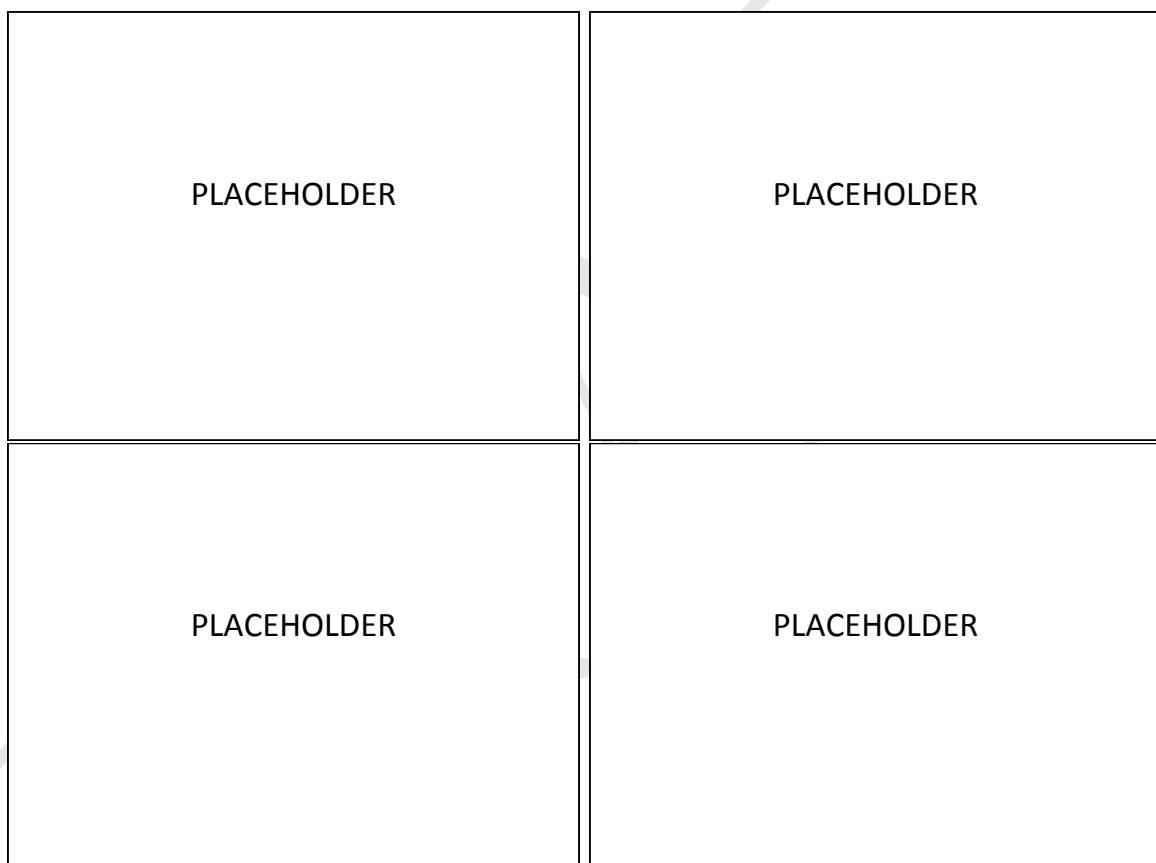


Figure 4.6.: The distributions of the kinematic dijet **BDT** response for the **VBF** signal (left) and the background (right) for the 7TeV (top row) and 8TeV (bottom row) training.

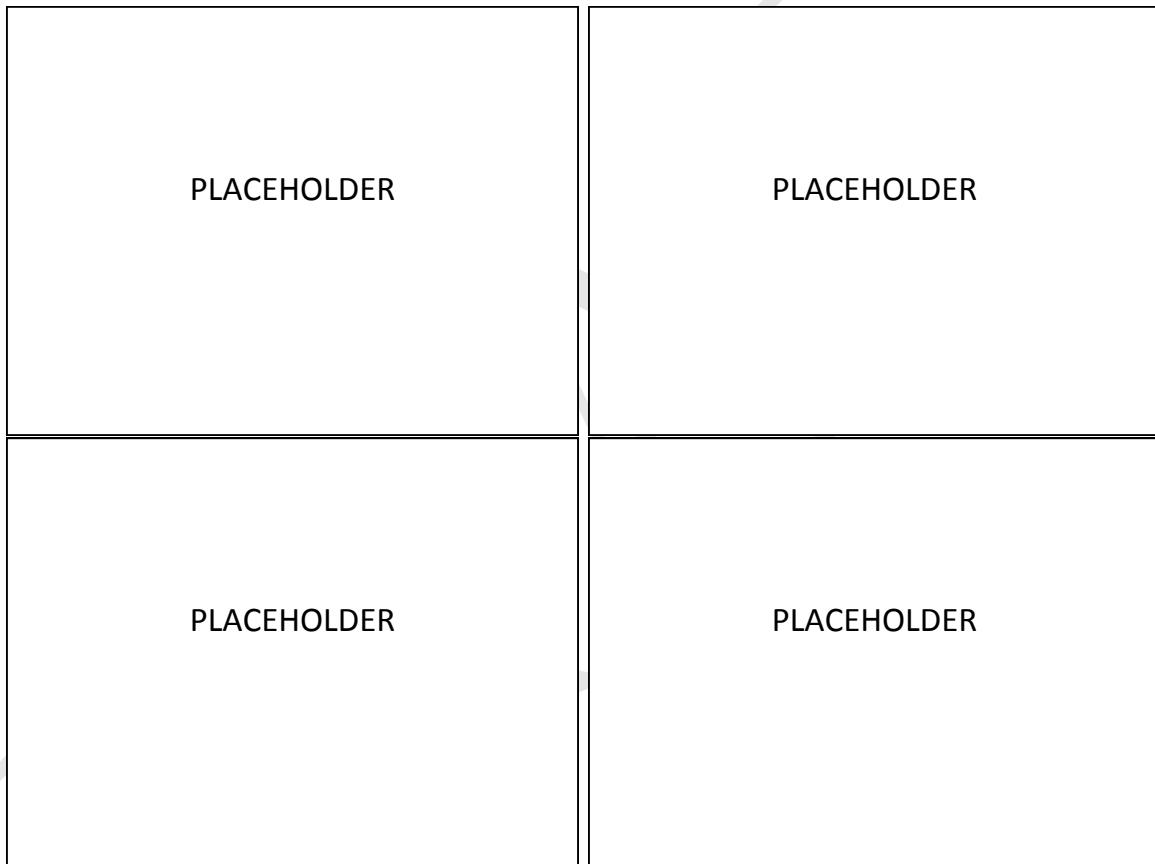


Figure 4.7.: The distributions of the combined dijet-diphoton **BDT** response for the **VBF** signal (left) and the background (right) for the 7TeV (top row) and 8TeV (bottom row) training. The red dashed lines corresponded to the location of the category boundaries.

of the classes require at least one charged muon or electron and are split into a tight selection category and a looser selection category, the third is for events consistent with large \cancel{E}_T and the fourth for events consistent with two or more jets. The leading photon cut is raised to $p_T^{\gamma_1} > 3m_{\gamma\gamma}/8$ for the lepton and \cancel{E}_T tagged categories and $p_T^{\gamma_1} > m_{\gamma\gamma}/2$ for the jet tagged category. The category requirements are as follows,

- **VH Tight l :** The tightly selected lepton category is characterised by the signature of a leptonically decaying W or Z boson and as such requires the presence of $\cancel{E}_T > 35$ GeV or another lepton of the same flavour and opposite charge as the first. In the first case (of a lepton + \cancel{E}_T) the lepton is required to have $p_T > 20$ GeV, in the latter case (of two leptons) the requirement is $p_T > 10$ GeV for both leptons whilst the invariant mass of the dilepton pair must be $70 < m_{ll} < 110$ GeV. The diphoton **BDT** output is required to be $> 0.1 (> -0.6)$ for the 7 (8) TeV datasets.
- **VH Loose l :** For the loosely selected lepton category the lepton p_T must satisfy $p_T > 20$ GeV. The selection requirements are designed to reduce the background from leptonic Z bosons (not associated with a Higgs) that contain initial or final state radiation faking the diphoton signal. Consequently, leptons are required to be separated by at least $\Delta R > 1.0$ from the closest photon and the invariant mass of any electron-photon pairs must be more than 10 GeV away from the Z boson mass. In addition a conversion veto is applied to the electrons to reduce the rate of misidentified photons. The diphoton **BDT** output is required to be $> 0.1 (> -0.6)$ for the 7 (8) TeV datasets.
- **VH \cancel{E}_T tag:** Accurate measurement and simulation of the \cancel{E}_T vector has been studied in detail at CMS and a set of standard corrections (for both data and simulation) are applied [31]. The corrected \cancel{E}_T is required to pass $\cancel{E}_T > 70$ GeV whilst the angular separation in the transverse plane between the diphoton system and the \cancel{E}_T direction must pass $\Delta\phi(\gamma\gamma, \cancel{E}_T) > 2.1$, and similarly the angle between the diphoton system and the leading jet must pass $\Delta\phi(\gamma\gamma, \text{jet}) < 2.7$. The diphoton **BDT** output is required to be $> 0.8 (> 0.0)$ for the 7 (8) TeV datasets.
- **VH jet tag:** The event must contain at least one jet pair in which both jets have $p_T > 40$ GeV and $|\eta| < 2.4$ and have an invariant mass within the range $60 < m_{jj} < 120$ GeV. The diphoton transverse momenta must satisfy $p_T^{\gamma\gamma} > 13m_{\gamma\gamma}/12$. Additionally the angular correlation between the diphoton system and the dijet system from VH associated Higgs decays can be exploited. The angle, θ^* , between the diphoton direction in the diphoton-dijet rest frame and the lab frame is flat for

events from VH decays whereas for the background and gluon fusion produced Higgs decays the distribution peaks at $|\cos(\theta^*)| = 1$. Consequently, there is a requirement that $|\cos(\theta^*)| < 0.5$. The diphoton BDT output is required to be > 0.6 (> 0.2) for the 7 (8) TeV datasets.

1154 Lepton and jet tagged categories for $t\bar{t}H$

1155 There are two categories for tagging production from $t\bar{t}H$ decays, one which is lepton
 1156 based and one which is jet based. The total fraction of signal expected from $t\bar{t}H$ is
 1157 $< 1\%$ so only a handful of events are expected. Consequently for the 7 TeV dataset the
 1158 two categories are merged into one class. As for the VH tagged categories the cuts are
 1159 optimised to minimise the expected uncertainty of the signal strength measurement of
 1160 the $t\bar{t}H$ process alone.

1161 For both classes the leading photon p_T cut is raised to $p_T^{\gamma_1} > m_{\gamma\gamma}/2$, all jets are
 1162 required to have $p_T > 25$ GeV and there must be at least one b-tagged jet present. The
 1163 specific requirements of each category are as follows,

- 1164 • **$t\bar{t}H$ multijet tag:** The requirement is at least four additional jets in the event
 1165 and no lepton. The diphoton BDT output is required to be > 0.6 (> -0.2) for the
 1166 7 (8) TeV datasets.
- 1167 • **$t\bar{t}H$ lepton tag:** At least one more jet in the event and one muon or electron
 1168 which has $p_T > 20$ GeV. The diphoton BDT output is required to be > 0.6 (> -0.6)
 1169 for the 7 (8) TeV datasets.

1170 4.2.2. Inclusive mode categorisation in the cut based analysis

1171 Any event which passes the cut based photon selection described in Section 4.1.1 and
 1172 does not fall into one of the exclusive categories described above is split into one of eight
 1173 inclusive categories depending on the supercluster position, η , of the two photons, the
 1174 conversion variable, R_9 , of the two photons and the mass relative diphoton transverse
 1175 momenta, $p_T/m_{\gamma\gamma}$. The categories are defined in Table 4.3.

$p_T/m_{\gamma\gamma}$	Maximum η	Minimum R_9
Untagged 0	<1.444	> 0.94
Untagged 1 $> 40 \text{ GeV}$		< 0.94
Untagged 2	<2.5	> 0.94
Untagged 3		< 0.94
Untagged 4	<1.444	> 0.94
Untagged 5 $< 40 \text{ GeV}$		< 0.94
Untagged 6	<2.5	> 0.94
Untagged 7		< 0.94

Table 4.3.: The definition of the inclusive categories for the [CiC](#) analysis

[1176](#) [4.2.3. Inclusive mode categorisation and VBF dijet](#) [1177](#) [categorisation in the mass factorised MVA analysis](#)

[1178](#) Firstly the combined dijet-diphoton [BDT](#) score is used to define a set of [VBF](#) categories.
[1179](#) The goal is to find the configuration of category boundaries which minimise the expected
[1180](#) uncertainty on the signal strength for [VBF](#) production alone, allowing the number of
[1181](#) categories and where the category boundaries lie to be completely free floating, with
[1182](#) the additional requirement that the efficiency \times acceptance of the categories matches
[1183](#) between 7 and 8 TeV. This results in a tight [VBF](#) category for events with a very
[1184](#) high combined dijet-diphoton [BDT](#) score, a somewhat looser category and then one or
[1185](#) more very loose categories. One finds that dropping the loosest category has a negligible
[1186](#) impact ($< 1\%$) on the expected uncertainty and as such the upper boundary for the
[1187](#) loosest category is turned into a lower cut. All events which pass the [VBF](#) preselection
[1188](#) (described in Sec. 4.2.1) and fail the lower cut are then classified somewhere in the
[1189](#) inclusive categories defined below.

[1190](#) Once the [VBF](#) category boundaries have been found the same procedure is deployed
[1191](#) using the diphoton [BDT](#) score. This time the target is to minimise the expected
[1192](#) uncertainty on the total signal strength, allowing the number of categories and the
[1193](#) category boundary values to be completely free floating. One finds a rather similar
[1194](#) structure and that dropping events in the very loosest category has a negligible impact
[1195](#) on the performance and consequently this dictates the lower cut value for the diphoton
[1196](#) [BDT](#).

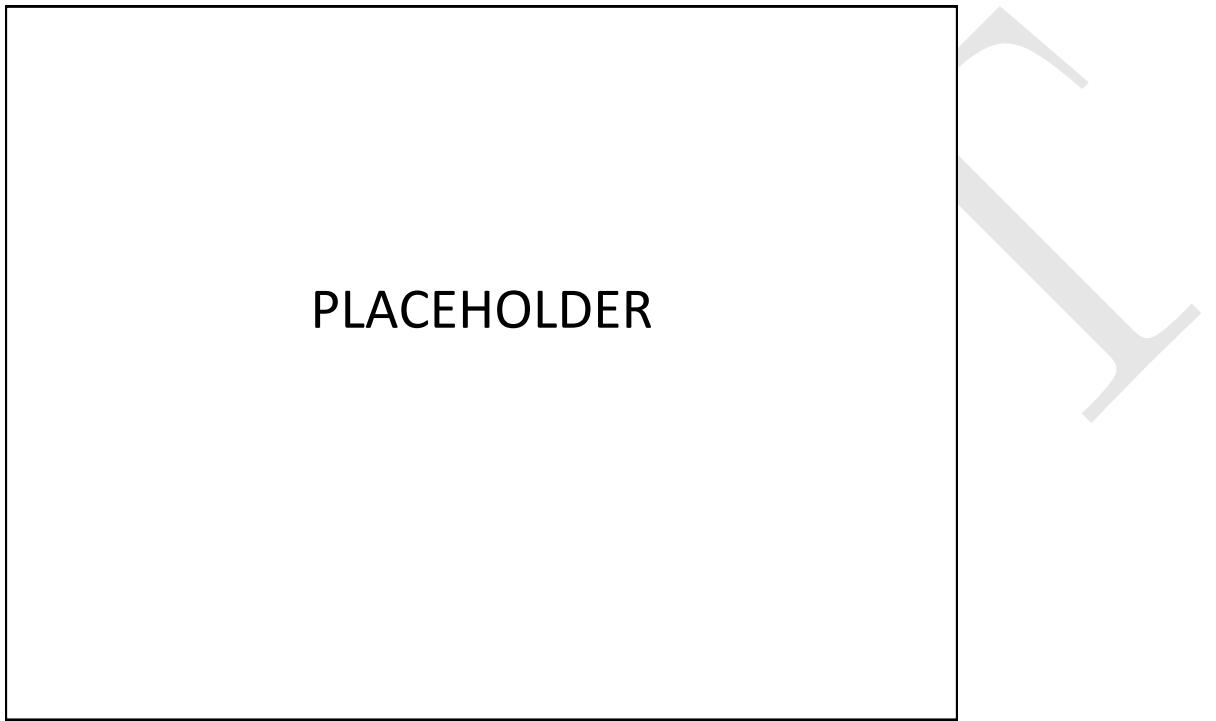


Figure 4.8.: A pictoral representation of how the VBF and inclusive category boundaries are optimised

1197 Due to the amount of simulation MC available to do these studies the precision to
 1198 which the boundary values are important is quite broad. Furthermore, the expected
 1199 uncertainty minima are very shallow so the exact position of the boundaries has a very
 1200 small impact on the performance of the analysis. For the 7 (8) TeV datasets there are 4
 1201 (5) inclusive categories with a lower diphoton BDT cut of 0.19 (-0.78) and 2 (3) VBF
 1202 categories. The boundary optimisation procedure is pictorially represented in Figure 4.8.

1203 4.2.4. Inclusive mode categorisation in the sideband MVA 1204 analysis

1205 In the SMVA analysis all the exclusive categories are identical to the MFM analysis
 1206 (including the VBF categories). However the inclusive categorisation is done slightly
 1207 differently. In the sideband analysis the signal region is defined in a $\pm 2\%$ window around
 1208 the hypothesis Higgs mass. The analysis is performed as a cut and count in the signal
 1209 window over several bins. There is one bin for each exclusive category and then several
 1210 more for the inclusive. The binning scheme for the inclusive events is defined as follows,

- Make two dimensional distributions of the diphoton **BDT** score and the distance of the invariant mass from the hypothesised Higgs mass, $\Delta m/m_H$, in the $\pm 2\%$ window for signal and background, as shown in Figure 4.9, where,

$$\Delta m/m_H = \frac{m_{\gamma\gamma} - m_H}{m_H}. \quad (4.4)$$

- ₁₂₁₁ • Select bins by isolating regions of this 2D phase space which have similar S/B ratios
₁₂₁₂ and optimise the boundaries to give the maximum expected signal significance.

₁₂₁₃ Clearly the most sensitive bins will be the ones which have a high diphoton **BDT**
₁₂₁₄ score and have a low value of $|\Delta m/m_H|$ (i.e. are near the signal peak). The category
₁₂₁₅ boundaries in this 2D plane are shown as different shades in Figure 4.10. In total there
₁₂₁₆ are 8 (10) inclusive bins for the 7 (8) TeV samples in the **SMVA**.

₁₂₁₇ 4.2.5. Event categorisation summary

₁₂₁₈ All the categories and the tagging order are summarised in the table below (Table 4.4).

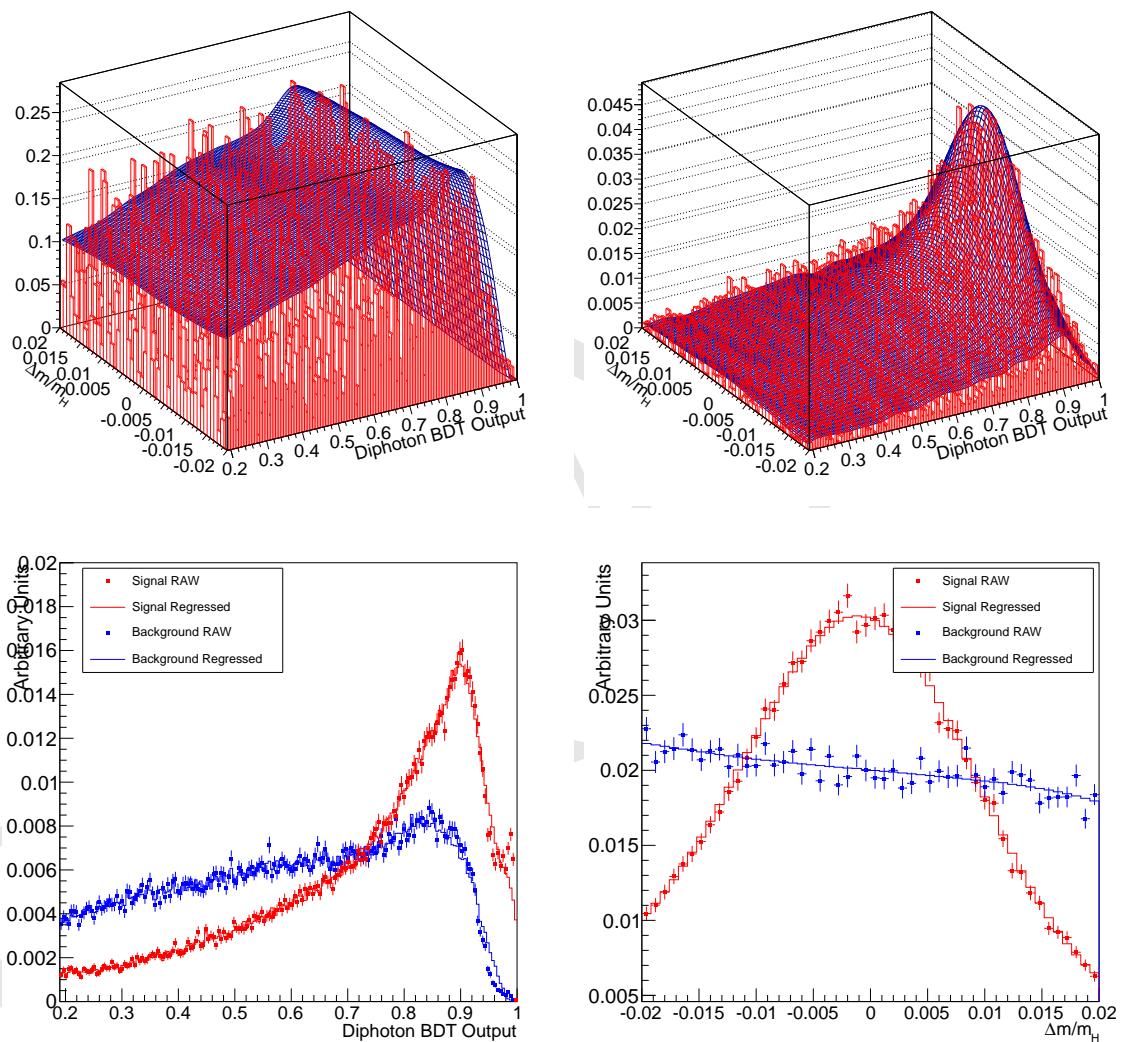


Figure 4.9.: Two dimensional distributions of the diphoton BDT output and $\Delta m/m_H$ are shown on the top row for the background (left) and signal (right) for the 7 TeV sample. The bottom rows show the projection for signal (red) and background (blue) in the two variables.

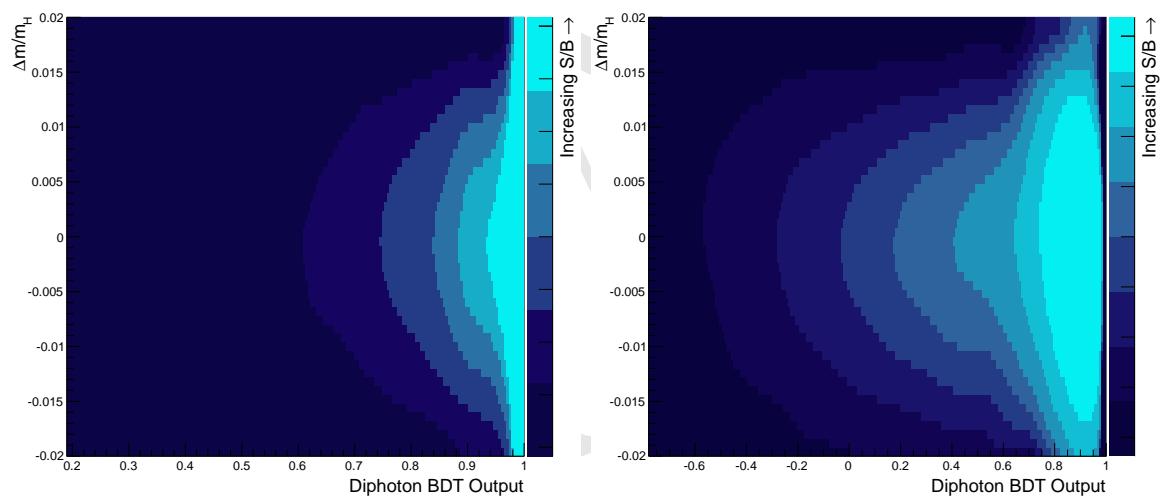


Figure 4.10.: The inclusive category bin definitions for the SMVA analysis. Shown for the 7 TeV dataset on the left and the 8 TeV dataset on the right.

Table 4.4.: The event classes at 7 and 8TeV and some of their main selection requirements. Events are tested against the selection requirements of the classes in the order they are listed here.

Label	No. of classes		Main requirements
	7GeV	8GeV	
$t\bar{t}H$ lepton tag	*	1	$p_T^\gamma(1) > m_{\gamma\gamma}/2$ 1 b-tagged jet + 1 electron or muon For MVAs: diphoton BDT $> 0.6(-0.6)$ at 7(8) TeV
VH tight ℓ tag	1	1	$p_T^\gamma(1) > 45 \cdot m_{\gamma\gamma}/120$ e or μ , $p_T > 20GeV$, and $\cancel{E}_T > 45GeV$ OR $2e$ or 2μ , $p_T > 10GeV$; $70 < m_{ll} < 110GeV$ For MVAs: diphoton BDT $> 0.1(-0.6)$ at 7(8) TeV
VH loose ℓ tag	1	1	$p_T^\gamma(1) > 45 \cdot m_{\gamma\gamma}/120$ e or μ , $p_T > 20GeV$ For MVAs: diphoton BDT $> 0.1(-0.6)$ at 7(8) TeV
VBF dijet tag	2	2/3 [†]	$p_T^\gamma(1) > m_{\gamma\gamma}/2$ For MVAs: 2 jets; dijet and combined diphoton-dijet BDTs used For CiC: 2 jets, cut based dijet selection used
VH \cancel{E}_T tag	1	1	$p_T^\gamma(1) > 45 \cdot m_{\gamma\gamma}/120$ $\cancel{E}_T > 70GeV$ For MVAs: diphoton BDT $> 0.8(0.0)$ at 7(8) TeV
$t\bar{t}H$ multijet tag	*	1	$p_T^\gamma(1) > m_{\gamma\gamma}/2$ 1 b-tagged jet + 4 more jets For MVAs: diphoton BDT $> 0.6(-0.2)$ at 7(8) TeV
VH dijet tag	1	1	$p_T^\gamma(1) > 45 \cdot m_{\gamma\gamma}/120$ jet pair, $p_T > 40GeV$ and $60 < m_{jj} < 120GeV$ For MVAs: diphoton BDT $> 0.6(0.2)$ at 7(8) TeV
Untagged	4/8 [‡]	5/8/10 [‡]	The remaining events, For MVAs: classified using diphoton BDT For CiC: classified using η and R_9 of photons and $p_T^{\gamma\gamma}/m_{\gamma\gamma}$

* For the 7TeV dataset, events in the $t\bar{t}H$ lepton tag and multijet tag classes are combined, after selection, to form a single event class.

† For the CiC (MVA) analysis there are 2 (3) dijet categories at 8 TeV.

‡ For the CiC analysis there are 8 inclusive categories at 7 and 8 TeV, for the MFM there are 4 (5) categories at 7 (8) TeV and for the SMVA there are 8 (10).

Chapter 5.

¹²¹⁹ Analysis and Results

¹²²⁰ “I’m a bus”

¹²²¹ — Darren Burton 1987–2013

¹²²² 5.1. Description

¹²²³ This would include a detailed description of the signal and background models used and
¹²²⁴ the systematics included. I would describe the statistical tools used such as hypothesis
¹²²⁵ testing with a test statistic and show the results which would include exclusion limits,
¹²²⁶ pvalues and SM best fit parameters (mass, signal strength, couplings etc.). One of the
¹²²⁷ important considerations of the Hgg analysis is the estimation of the background. This
¹²²⁸ second part of this section will concentrate on presenting an alternate analysis which
¹²²⁹ takes a completely different approach to modeling the background and cross-checks the
¹²³⁰ result above.

₁₂₃₁ **5.2. Signal modelling**

₁₂₃₂ **5.2.1. Mass factorised analysis**

₁₂₃₃ **5.2.2. Sideband analysis**

₁₂₃₄ **5.3. Background modelling**

₁₂₃₅ **5.3.1. Mass factorised analysis**

₁₂₃₆ **5.3.2. Sideband analysis**

₁₂₃₇ **5.4. Systematic Uncertainties**

₁₂₃₈ **5.5. Statistics**

₁₂₃₉ **5.5.1. Use of the Likelihood function as a test statistic**

₁₂₄₀ **5.6. Results of the mass factorised analysis**

₁₂₄₁ **5.7. Results of the sideband analysis**

₁₂₄₂

Appendix A.

¹²⁴³ Photon ID BDT input variables

¹²⁴⁴ validation in $Z \rightarrow e^+e^-$

DRAFT

Appendix B.

¹²⁴⁵ Diphoton BDT input variables
¹²⁴⁶ validation in $Z \rightarrow e^+e^-$

DRAFT

DRAFT

¹²⁴⁸ Bibliography

- ¹²⁴⁹ [1] S. Weinberg, Phys. Rev. Lett. **19**, 1264 (1967).
- ¹²⁵⁰ [2] S. L. Glashow, J. Iliopoulos, and L. Maiani, Phys. Rev. **D2**, 1285 (1970).
- ¹²⁵¹ [3] S. Willenbrock, (2004), hep-ph/0410370.
- ¹²⁵² [4] The CMS Collaboration, JINST **3**, S08004 (2008).
- ¹²⁵³ [5] The CMS Collaboration, Eur. Phys. J. C **70**, 1165 (2010), hep-ex/1007.1988.
- ¹²⁵⁴ [6] The CMS Collaboration, (2012), CMS-DP-2013-007.
- ¹²⁵⁵ [7] E. Meschi, T. Monteiro, C. Seez, and P. Vikas, (2001), CMS-AN-2001-034.
- ¹²⁵⁶ [8] W. Adam, R. Früwirth, A. Strandlie, and T. Todorov, J. Phys. G: Nucl. Part. Phys. **31**, N9 (2005).
- ¹²⁵⁸ [9] J. Weng, (2008), hep-ex/0810.3686v1.
- ¹²⁵⁹ [10] M. Cacciari, G. P. Salam, and G. Soyez, J. High Energy Phys. **04**, 063 (2008).
- ¹²⁶⁰ [11] C. Ferro, (2012), hep-ex/1201.5292.
- ¹²⁶¹ [12] The CMS Collaboration, J. Phys. Conf. Ser. **404**, 012014 (2012).
- ¹²⁶² [13] The CMS Collaboration, (2009), CMS-PFT-09-001.
- ¹²⁶³ [14] A. Hoecker *et al.*, (2007), physics/0703039.
- ¹²⁶⁴ [15] F. J. Massey, JSTOR **46**, 68 (1951).
- ¹²⁶⁵ [16] S. Alioli, P. Nason, C. Oleari, and E. Re, JHEP **04**, 002 (2009).
- ¹²⁶⁶ [17] P. Nason and C. Oleari, JHEP **02**, 037 (2010).
- ¹²⁶⁷ [18] T. Sjöstrand, S. Mrenna, and P. Z. Skands, JHEP **0605**, 026 (2006).
- ¹²⁶⁸ [19] LHC Higgs Cross Section Working Group, S. Heinemeyer *et al.*, (2013), 1307.1347.

- ₁₂₆₉ [20] Y. Gao *et al.*, Phys. Rev. D **81**, 242 (2010).
- ₁₂₇₀ [21] T. Gleisberg *et al.*, JHEP **0902**, 007 (2009), 0811.4622.
- ₁₂₇₁ [22] J. Alwall, M. Herquet, F. Maltoni, O. Mattelaer, and T. Stelzer, JHEP **1106**, 128 (2011), 1106.0522.
- ₁₂₇₂
- ₁₂₇₃ [23] GEANT4, S. Agostinelli *et al.*, Nucl. Instrum. Meth. A **506**, 250 (2003).
- ₁₂₇₄ [24] J. Beringer *et al.*, Phys. Rev. D **86**, 010001 (2012).
- ₁₂₇₅ [25] The CMS Collaboration, JHEP **2011**, 1 (2011).
- ₁₂₇₆ [26] R. Szalapski and D. Zeppenfeld, Phys. Rev. D **54**, 6680 (1996), arXiv:hep-ph/9605444.
- ₁₂₇₇
- ₁₂₇₈ [27] CMS, S. Chatrchyan *et al.*, JINST **06**, P11002 (2011).
- ₁₂₇₉ [28] M. Cacciari and G. P. Salam, Phys. Lett. **B659**, 119 (2008).
- ₁₂₈₀ [29] M. Cacciari, G. P. Salam, and G. Soyez, JHEP **04**, 005 (2008).
- ₁₂₈₁ [30] M. Cacciari, G. P. Salam, and G. Soyez, (2011), 1111.6097.
- ₁₂₈₂ [31] T. C. Collaboration, CMS PAS **CMS-PAS-JME-12-002** (2012).

List of Figures

1284	2.1. CMS diagram	8
1285	2.2. CMS tracker	9
1286	2.3. Vertex resolution	9
1287	2.4. CMS tracker material budget	10
1288	2.5. CMS ECAL schematic	11
1289	2.6. ECAL energy resolution	13
1290	2.7. ECAL laser corrections	15
1291	2.8. The domino construction setup of the hybrid clustering algorithm[7].	17
1292	2.9. Particle flow jet resolution	21
1293	3.1. The total integrated luminosity delivered and recorded by CMS during the 2011 (left) and 2012 (right) run periods.	26
1294		
1295	3.2. Distribution of the number of reconstructed vertices in the 2011 (left) and 2012 (right) run periods. Calculated using the Deterministic Annealing algorithm in [REF!] for $Z \rightarrow \mu^+ \mu^-$ events in data (black dots) and MC (red histogram) after reweighting.	28
1296		
1297		
1298		
1299	3.3. Distribution of Δz (the distance between the chosen vertex and the true vertex in the z direction) for data (black), the MC (red) and the MC after beam spot reweighting (green) for $Z \rightarrow \mu^+ \mu^-$ events.	29
1300		
1301		

1302	3.4. A comparison of the predicted probability density of E_{true}/E_{raw} from the regression training (blue line) to the distribution in a statistically independent MC sample (black points) for barrel photons (left) and endcap photons (right). PLOTS NEED TIDYING	32
1306	3.5. A comparison of the regression performance compared to the photon default for the diphoton invariant mass on a MC source of Higgs decays to two photons. This is shown when both photons are in the barrel (left two plots) when both photons have $R_9 > 0.94$ (far left) and when at least one photon has $R_9 < 0.94$ (middle left). When at least one photon is in the endcap (right two plots) when both photons have $R_9 > 0.94$ (middle right) and when at least one photon has $R_9 < 0.94$ (far right). The value shown in the legend of each is the effective width, σ_{eff} , which is defined as half of the narrowest interval which contains 68.3% of the distribution. PLOTS NEED TIDYING	32
1316	3.6. The $Z \rightarrow e^+e^-$ invariant mass shape comparison before and after the scale and smearing corrections are applied. Shown for 8 TeV for photons with $1. < \eta < 1.444$ and $R_9 < 0.94$ on the left and for photons with $ \eta < 1.$ and $R_9 > 0.4$ on the right.	34
1320	3.7. The $Z \rightarrow e^+e^-$ invariant mass distribution at 7TeV in data (black points) and MC (blue histogram) for events which pass the analysis preselection in which the electron veto is inverted.	35
1323	3.8. The $Z \rightarrow e^+e^-$ invariant mass distribution at 8TeV in data (black points) and MC (blue histogram) for events which pass the analysis preselection in which the electron veto is inverted.	35
1326	3.9. A representation of the two methods for locating the primary vertex using photon conversion information. The left plot is for cases where the conversion occurs early enough in the tracker that the two electron tracks can be used to construct the converted pair momentum which is combined with the conversion vertex position to point back to the beam line. The right plot is for cases where the conversion occurs late in the tracker and the energy weighted supercluster position and the conversion vertex position are used to point back to the beam line.	37

1334	3.10. Distributions of the input variables for the vertex BDT in the MC $H \rightarrow \gamma\gamma$ training (points) and test (filled) samples at 8TeV. Shown for the target primary vertex (blue histograms) and the background pileup vertices (red histograms). Plots need updating, tidying and labelling correctly.	39
1335		
1336		
1337		
1338	3.11. The vertex BDT response for $Z \rightarrow \mu^+\mu^-$ events in data (points) and MC (filled) for the primary vertex (green) and the background pileup vertices (red).	40
1339		
1340		
1341	3.12. The chosen vertex efficiency as measured in $Z \rightarrow \mu^+\mu^-$ data and MC as a function of $Z p_T$ (left) and number of reconstructed vertices (right) for the 7 TeV (top row) and 8 TeV (bottom row) data samples.	41
1342		
1343		
1344	3.13. A demonstration of the linearity relation between the per-event vertex probability BDT output distribution and the correct vertex probability	42
1345		
1346	3.14. A comparison of the true vertex efficiency (black points) and the average vertex probability (blue band) a statistically independent MC Higgs sample simulated with 2012 running conditions.	42
1347		
1348		
1349	3.15. The di-electron mass distribution after applying the preselection (described in Sec. 3.5) and the scale and smearing corrections (described in Sec. 3.3.1) whilst inverting the electron veto. The left plot shows the data at 8TeV as the black points with the Drell Yan MC sample as the blue histogram. The Data/MC ratio is shown in the right hand plot. Plot needs updating	45
1350		
1351		
1352		
1353		
1354	4.1. Cut based photon ID efficiency as measured in $Z \rightarrow e^+e^-$ tag and probe. MORE DESCRIPTION	50
1355		
1356	4.2. The output distribution of the photon identification BDT for 7TeV barrel (top left), 7TeV endcap (top right), 8TeV barrel (bottom left) and 8TeV endcap (bottom right). The solid points show the $\gamma + \text{jet}$ training sample distributions and the hollow points show the $H \rightarrow \gamma\gamma$ test sample distributions for prompt signal photons in blue and fake background photons in red. A cut of > -0.2 is made on all photons.	53
1357		
1358		
1359		
1360		
1361		

1362	4.3. The output distribution of the photon identification BDT for the 8 TeV training as validated by the $Z \rightarrow e^+e^-$ decay. The data is shown as the black points with the MC as the blue histogram. The systematic uncertainty on the output as applied to the MC is shown as the red band. This is literally just a place holder for the plot that I want to show.	54
1367	4.4. The diphoton BDT response for the 7 TeV training (left column) and 8 TeV training (right column). The data and background distributions are shown in the plots on the top row for data in the range $100 < m_{\gamma\gamma} < 120$ GeV and $130 < m_{\gamma\gamma} < 180$ GeV (black points) and for the prompt-prompt background (green), prompt-fake background (yellow) and fake-fake background (red). The signal distributions are shown in the plots on the bottom row for gluon fusion (red), vector boson fusion (yellow), associated W, Z production (green) and associated $t\bar{t}$ production (blue) alongside the data in the range $100 < m_{\gamma\gamma} < 120$ GeV and $130 < m_{\gamma\gamma} < 180$ GeV (black points) and the total background (hollow histogram). PLOTS NEED UPDATING	56
1378	4.5. The diphoton BDT response for the 8 TeV training in the $Z \rightarrow e^+e^-$ decay. The data is shown as the black points and the MC as the blue histogram. The systematic applied to account for variation in the BDT response from mismodelling in the photon quality response and the photon energy resolution estimate are shown as the red band. This is just a placeholder for the plot I want to show	57
1384	4.6. The distributions of the kinematic dijet BDT response for the VBF signal (left) and the background (right) for the 7TeV (top row) and 8TeV (bottom row) training.	61
1387	4.7. The distributions of the combined dijet-diphoton BDT response for the VBF signal (left) and the background (right) for the 7TeV (top row) and 8TeV (bottom row) training. The red dashed lines corresponded to the location of the category boundaries.	62
1391	4.8. A pictoral representation of how the VBF and inclusive category boundaries are optimised	66

1393	4.9. Two dimensional distributions of the diphoton BDT output and $\Delta m/m_H$		
1394	are shown on the top row for the background (left) and signal (right) for		
1395	the 7 TeV sample. The bottom rows show the projection for signal (red)		
1396	and background (blue) in the two variables.	68	
1397	4.10. The inclusive category bin definitions for the SMVA analysis. Shown for		
1398	the 7 TeV dataset on the left and the 8 TeV dataset on the right.	69	

DRAFT

List of Tables

1399	3.1. Preselection cut values.	44
1400		
1401	4.1. Photon ID selection cut values. The cuts are applied to both the leading	
1402	and subleading photons.	49
1403	4.2. Final selection cuts for the VBF selection. Events from the first category	
1404	are removed from the second one.	59
1405	4.3. The definition of the inclusive categories for the CiC analysis	65
1406		
1407	4.4. The event classes at 7 and 8TeV and some of their main selection require-	
1408	ments. Events are tested against the selection requirements of the classes	
	in the order they are listed here.	70

List of Acronyms

1409		
1410		
1411	BDT	- Boosted Decision Tree
1412	DT	- Decision Tree
1413	CMS	- Compact Muon Solenoid
1414	ECAL	- electromagnetic calorimeter
1415	HCAL	- hadronic calorimeter
1416	LHC	- Large Hadron Collider
1417	MVA	- Multivariate Analysis
1418	CiC	- Cuts in Categories
1419	MFM	- Mass Factorized MVA
1420	SMVA	- Sideband MVA
1421	PbWO ₄	- lead tungstate
1422	APDs	- avalanche photodiodes
1423	VPTs	- vacuum phototriodes
1424	GED	- global event description
1425	PF	- particle flow
1426	MC	- Monte Carlo
1427	QCD	- quantum chromodynamics
1428	ggH	- gluon fusion
1429	VBF	- vector boson fusion

¹⁴³⁰	VH	- vector boson associated production
¹⁴³¹	$t\bar{t}H$	- top quark associated production