

1      **Using Machine Learning Techniques for Data Quality**  
2      **Monitoring at CMS Experiment**

3      by

4      Guillermo A. Fidalgo Rodríguez

5      A thesis presented for the degree of

6      BACHELLOR'S OF SCIENCE??

7      in

8      Physics

9      UNIVERSITY OF PUERTO RICO  
10     MAYAGÜEZ CAMPUS

11     2018

12     Approved by:

13     \_\_\_\_\_  
14     Sudhir Malik, Ph.D.

15     President, Graduate Committee  
Date

16     \_\_\_\_\_  
17     Héctor Méndez, Ph.D.

18     Member, Graduate Committee  
Date

19     \_\_\_\_\_  
20     Samuel Santana Colón, Ph.D.  
21     Member, Graduate Committee

Date

22     \_\_\_\_\_  
23     Rafael A. Ramos, Ph.D.  
24     Chairperson of the Department

Date

<sup>25</sup> **Abstract**

<sup>26</sup> The Data Quality Monitoring (DQM) of CMS is a key asset to deliver high-quality  
<sup>27</sup> data for physics analysis and it is used both in the online and offline environment. The cur-  
<sup>28</sup> rent paradigm of the quality assessment is labor intensive and it is based on the scrutiny of  
<sup>29</sup> a large number of histograms by detector experts comparing them with a reference. This  
<sup>30</sup> project aims at applying recent progress in Machine Learning techniques to the automa-  
<sup>31</sup> tion of the DQM scrutiny. In particular the use of convolutional neural networks to spot  
<sup>32</sup> problems in the acquired data is presented with particular attention to semi-supervised  
<sup>33</sup> models (e.g. autoencoders) to define a classification strategy that doesn't assume previ-  
<sup>34</sup>ous knowledge of failure modes. Real data from the hadron calorimeter of CMS are used  
<sup>35</sup> to demonstrate the effectiveness of the proposed approach.

<sup>36</sup> *Keywords:* [DQM, online, offline, Machine Learning ]

<sup>37</sup> **Acknowledgments**

<sup>38</sup> I wish to thank United States State Department and University of Michigan for pro-  
<sup>39</sup> viding the opportunity to work abroad at CERN during the 2018 Winter Semester. I also  
<sup>40</sup> wish to thank CERN staff, CMS Experiment , Texas Tech University, and the University  
<sup>41</sup> of Puerto Rico at Mayagüez, with special thanks to Dr. Federico de Guio for his local  
<sup>42</sup> mentorship and Dr. Nural Akchurin, and Dr. Sudhir Malik for their guidance. A very  
<sup>43</sup> special thanks to Dr. Jean Krisch for accepting me for this great research opportunity and  
<sup>44</sup> Dr. Steven Goldfarb for being a wonderful host and overall local guidance at CERN.

# 45 List of Figures

46	2.1	CMS Detector	4
47	2.2	The trajectory of a particle traveling through the layers of the detector leaving behind it's signature footprint	5
49	4.1	Occupancy maps with 5x5 affected regions	11
50	4.2	Weights and Biases	12
51	4.3	Gradient Descent algorithm	13
52	4.4	Loss Function surface	14
53	5.1	Occupancy Maps with 1x1 bad regions. A) Good image B) Dead image C) Hot image	16
55	5.2	Shown are both the $p_T^\gamma$ (left) and $p_T^{\gamma(miss)}$ (right) distributions before ap- plying any corrections. $p_T^{\gamma(miss)}$ is obtained by adding the $p_T^\gamma$ to the total $p_T^{miss}$ in every event.	18
58	5.3	Plots for Fake Rate (left) and Purity (right) as a function of the photon $p_T$ are shown. The events are selected are required to have a $p_T > 200$ GeV, be within the ECAL acceptance range, and pass the loose ID selection cuts. This selection was produced in order to verify the values given by the E/ $\gamma$ POG. As can be seen, the efficiency (purity) is seen to agree with the values of the loose photon ID/isolation selection.	20
64	5.4	Plots for Fake Rate (left) and Purity (right) as a function of the photon $p_T$ are shown. These plots include photons with the full control region selection. Aside from exhibiting lower statistics, the plots seem to agree with the fake rate and purity before all the control region cuts are applied.	20
68	5.5	Results of study of the Z+jets to $\gamma$ +jets cross-section ratio for both data and MadGraph simulation.	25
70	5.6	Jet multiplicity and the associated $S_\gamma$ scale factor in the loose photon con- trol region before any corrections are applied.	26
72	5.7	$N_{jet}$ distribution in the tight $\mu\mu$ control region after $S_\gamma$ corrections.	26
73	5.8	$p_T^\gamma$ (left) and $p_T^{\gamma(miss)}$ (right) distributions after applying the $S_\gamma(N_j)$ scale factor. Comparing to Figure 5.2, an improvement in the agreement be- tween data/MC can be observed.	27
76	5.9	$N_b$ distribution before (left) and after (right) applying the $S_\gamma(N_j)$ scale factor.	27
78	5.10	$H_T$ and $m_{T2}$ distributions applying the $S_\gamma(N_j)$ scale factor.	28

79	5.11 Shown are data/MC comparisons for the $p_T^{miss}$ (left) and $H_T$ (right) distributions after applying both the $N_j$ -dependent shape corrections ( $S_\gamma$ ) and the global normalization scale factor ( $R_{norm}$ ) . . . . .	29
80		
81		
82	5.12 Shown are data/MC comparisons for the $N_j$ (left) and $N_b$ (right) distributions after applying both the $N_j$ -dependent shape corrections ( $S_\gamma$ ) and the global normalization scale factor ( $R_{norm}$ ) . . . . .	29
83		
84		
85	5.13 Systematic uncertainty in the final prediction, as a function of the search bin, associated to the MC statistics. . . . .	31
86		
87	5.14 $Z \rightarrow \nu\bar{\nu}$ background prediction for all search bins, including the breakdown of the various uncertainties. . . . .	32
88		

# 89 **Contents**

90	<b>Abstract</b>	i
91	<b>Acknowledgments</b>	ii
92	<b>List of Figures</b>	iii
93	<b>1 Introduction</b>	1
94	<b>2 The CMS Experiment</b>	3
95	<b>3 Data Collection and Data Quality Monitoring</b>	6
96	3.1 What is Data Collection for CMS? . . . . .	6
97	3.2 What is Data Quality Monitoring? . . . . .	7
98	<b>4 What is Machine Learning?</b>	9
99	4.1 Developing the Algorithm . . . . .	10
100	4.2 Teaching the Algorithm . . . . .	12
101	<b>5 Results</b>	15
102	5.0.1 Photon Selection . . . . .	18
103	5.0.2 Photon Purity and Fake Rate . . . . .	19
104	5.1 The $Z \rightarrow \mu^+ \mu^-$ Control Region . . . . .	22
105	5.1.1 Muon ID and Isolation . . . . .	22
106	5.1.2 Muon Selection in the Tight Control Region . . . . .	23
107	5.2 Analysis . . . . .	24
108	5.2.1 Shape Correction Using the $\gamma + \text{jets}$ Control Sample . . . . .	24
109	5.2.2 Normalization Correction Using the tight $Z \rightarrow \mu^+ \mu^-$ Control Sam- 110 ple . . . . .	27
111	5.3 Results . . . . .	30
112	5.3.1 Systematics . . . . .	30
113	5.3.2 $Z \rightarrow \nu \bar{\nu}$ Estimation for the Search Bins . . . . .	31
114	<b>A Appendix Title</b>	33
115	<b>B References</b>	34

<sup>116</sup> **Chapter 1**

<sup>117</sup> **Introduction**

<sup>118</sup> The work for this thesis was performed at CERN on CMS Experiment. CERN stands  
<sup>119</sup> for European Organization for Nuclear Research. It was founded in 1954 and is located  
<sup>120</sup> at the Franco-Swiss border near Geneva. At CERN, physicists and engineers are probing  
<sup>121</sup> the fundamental structure of the universe. They use the world's largest and most complex  
<sup>122</sup> scientific instruments to study the basic constituents of matter – the fundamental parti-  
<sup>123</sup> cles. The instruments used at CERN are purpose-built particle accelerators and detectors.  
<sup>124</sup> Accelerators boost beams of particles to high energies before the beams are made to col-  
<sup>125</sup> lide with each other or with stationary targets. Detectors observe and record the results  
<sup>126</sup> of these collisions. The accelerator at CERN is called the Large Hadron Collider (LHC),  
<sup>127</sup> the largest machine ever built by humans and it collides particles (protons) at close to the  
<sup>128</sup> speed of light. The process gives the physicists clues about how the particles interact, and  
<sup>129</sup> provides insights into the fundamental laws of nature. Seven experiments at the LHC use  
<sup>130</sup> detectors to analyze particles produced by proton-proton collisions. The biggest of these  
<sup>131</sup> experiments, ATLAS and CMS, use general-purpose detectors designed to study the fun-  
<sup>132</sup> damental nature of matter and fundamental forces and to look for new physics or evidence  
<sup>133</sup> of particles that are beyond the Standard Model. Having two independently designed de-  
<sup>134</sup> tectors is vital for cross-confirmation of any new discoveries made. The other two major  
<sup>135</sup> detectors ALICE and LHCb, respectively, study a state of matter that was present just  
<sup>136</sup> moments after the Big Bang and preponderance of matter than antimatter. Each experi-

137 ment does important research that is key to understanding the universe that surrounds and  
138 makes us.

139

140      [Chapter 2](#) presents a basic description of the Large Hadron Collider and CMS Detector

141

142      ?? gives a brief motivation

143

144      ?? is dedicated to a study optimizing

145

146      ?? ptimated.

147

148      ?? details an improvarger production cross-section than Z+jets process used before.

149

150      The conclusions and results of each chapter are presented in the corresponding chap-  
151 ter.

152

153      This thesis work has been presented at several internal meetings of the CMS Experi-  
154 ment and at the following international meetings and conferences:

155      1. **Andrés Abreu** gave a talk “*Estimation of the Z Invisible Background for Searches*  
156            *for Supersymmetry in the All-Hadronic Channel*” at “APS April 2018: American  
157            Physical Society April Meeting 2018, 14-17 Apr 2018”, Columbus, OH

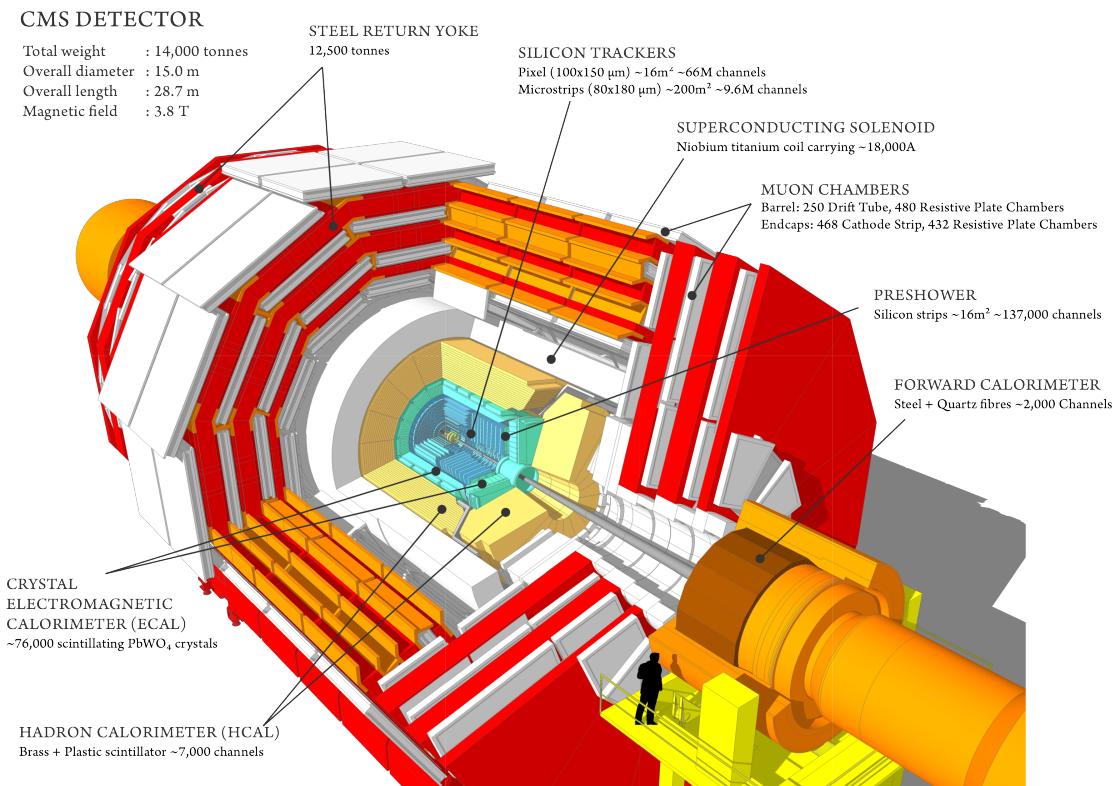
158      2. **Andrés Abreu** gave a talk “*Phase-2 Pixel upgrade simulations*” at the “USLUA  
159            Annual meeting: 2017 US LHC Users Association Meeting, 1-3 Nov 2017”, Fer-  
160            milab, Batavia, IL

<sub>161</sub> **Chapter 2**

<sub>162</sub> **The CMS Experiment**

<sub>163</sub> The Compact Muon Solenoid (CMS) detector is a general purpose particle detector  
<sub>164</sub> designed to investigate various physical phenomena concerning the SM and beyond it,  
<sub>165</sub> such as Supersymmetry, Extra Dimensions and Dark Matter. As its name implies, the  
<sub>166</sub> detector is a solenoid which is constructed around a superconducting magnet capable of  
<sub>167</sub> producing a magnetic field of 3.8 T. The magnetic coil is 13m long with an inner diameter  
<sub>168</sub> of 6m, making it the largest superconducting magnet ever constructed. The CMS detector  
<sub>169</sub> itself is 21m long with a diameter of 15m and it has a weight of approximately 14,000  
<sub>170</sub> tons. The CMS experiment is one of the largest scientific collaborations in the history  
<sub>171</sub> of mankind with over 4,000 participants from 42 countries and 182 institutions. CMS is  
<sub>172</sub> located at one of these points and it essentially acts as a giant super highspeed camera  
<sub>173</sub> that makes 3D images of the collisions that are produced at a rate of 40 MHz (40 million  
<sub>174</sub> times per second). The detector has an onion-like structure to capture all the particles that  
<sub>175</sub> are produced in these high energy collisions most of them being unstable and decaying  
<sub>176</sub> further to stable particles that are detected. CMS detector was designed with the following  
<sub>177</sub> features (as shown in [Figure 2.1](#)) :

- <sub>178</sub> 1. A **magnet** with large bending power and high performance muon detector for good  
<sub>179</sub> muon identification and momentum resolution over a wide range of momenta and  
<sub>180</sub> angles.

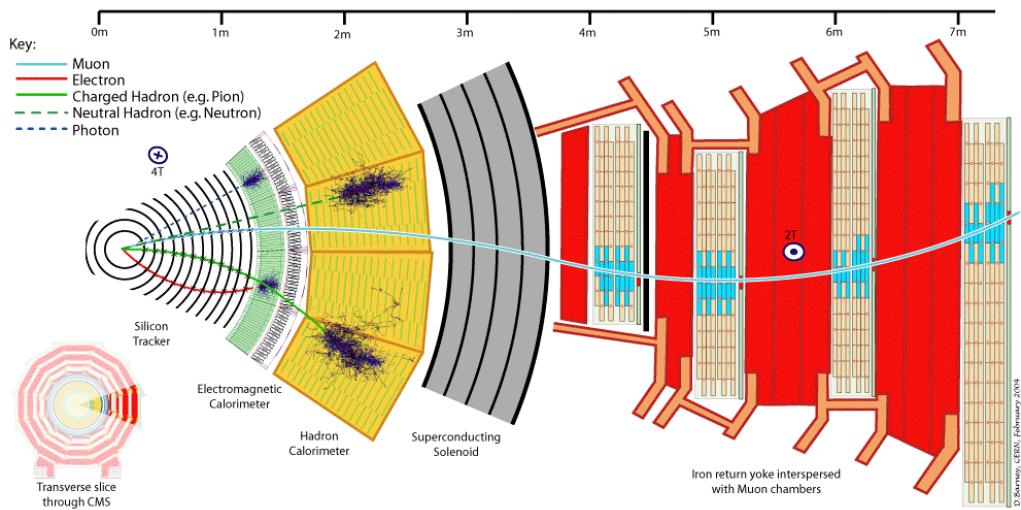


**Figure 2.1:** CMS Detector

- 181     2. An **inner tracking system** capable of high reconstruction efficiency and momen-  
182       tum resolution requiring **pixel detectors** close to the interaction region.
- 183     3. An **electromagnetic calorimeter** able to provide good electromagnetic energy res-  
184       olution and a high isolation efficiency for photons and leptons.
- 185     4. A **hadron calorimeter** capable of providing precise missing-transverse-energy and  
186       dijet-mass resolution.

187       A property from these particles that is exploited is their charge. Normally, particles  
188       produced in collisions travel in a straight line, but in the presence of a magnetic field,  
189       their paths are skewed and curved. Except the muon system, the rest of the subdetectors  
190       lie inside a 3.8 Tesla magnetic field . Due to the magnetic field the trajectory of charged  
191       particle produced in the collisions gets curved (as shown in [Figure 2.2](#) ) and one can  
192       calculate the particle's momentum and know the type of charge on the particle. The  
193       Tracking devices are responsible for drawing the trajectory of the particles by using a  
194       computer program that reconstructs the path by using electrical signals that are left by

195 the particle as they move. The Calorimeters measure the energy of particles that pass  
 196 through them by absorbing their energy with the intent of stopping them. The particle  
 197 identification detectors work by detecting radiation emitted by charged particles and using  
 198 this information they can measure the speed, momentum, and mass of a particle. After the  
 199 information is put together to make the “snapshot” of the collision one looks for results  
 200 that do not fit the current theories or models in order to look for new physics.



**Figure 2.2:** The trajectory of a particle traveling through the layers of the detector leaving behind its signature footprint

201 The project focusses specifically on data collected from one of the Calorimeters, - the  
 202 Hadron Calorimeter (HCAL). The HCAL, as its name indicates, is designed to detect and  
 203 measure the energy of hadrons or, particles that are composed of quarks and gluons, like  
 204 protons and neutrons. Additionally, it provides an indirect measurement of the presence  
 205 of non-interacting, uncharged particles such as neutrinos (missing energy) . Measuring  
 206 these particles is important as they can tell us if new particles such as the Higgs boson or  
 207 supersymmetric particles (much heavier versions of the standard particles we know) have  
 208 been formed. The layers of the HCAL are structured in a staggered fashion to prevent any  
 209 gaps that a particle might pass through undetected. There are two main parts: the barrel  
 210 and the end caps. There are 36 barrel wedges that form the last layer of the detector inside  
 211 the magnet coil, there is another layer outside this, and on the endcaps, there are another  
 212 36 wedges to detect particles that come out at shallow angles with respect to the beam  
 213 line.

# <sup>214</sup> Chapter 3

## <sup>215</sup> Data Collection and Data Quality

### <sup>216</sup> Monitoring

#### <sup>217</sup> 3.1 What is Data Collection for CMS?

<sup>218</sup> During data taking there are millions of collisions occurring in the center of the de-  
<sup>219</sup> tector every second. The data per event is around one million bytes (1 MB), that is produced  
<sup>220</sup> at a rate of about 600 million events per second [1], that's about 600 MB/s. Keeping  
<sup>221</sup> in mind that only certain events are considered "interesting" for analysis, the task of de-  
<sup>222</sup> ciding what events to consider out of all the data collected is a two-stage process. First,  
<sup>223</sup> the events are filtered down to 100 thousand events per second for digital reconstruction  
<sup>224</sup> and then more specialized algorithms filter the data even more to around 100 200 events  
<sup>225</sup> per second that are found interesting. For CMS there is a Data Acquisition System that  
<sup>226</sup> records the raw data to what's called a High-Level Trigger farm which is a room full  
<sup>227</sup> of servers that are dedicated to processing and classify this raw data quickly. The data  
<sup>228</sup> then gets sent to what's known as the Tier-0 farm where the full processing and the first  
<sup>229</sup> reconstruction of the data are done. [2]

## 230 3.2 What is Data Quality Monitoring?

231 To operate a sophisticated and complex apparatus as CMS, a quick online feedback on  
232 the quality of the data recorded is needed to avoid taking low quality data and to guarantee  
233 a good baseline for the offline analysis. Collecting a good data sets from the collisions  
234 is an important step towards search for new physics as deluge of new data poses an extra  
235 challenge of processing and storage. This all makes it all the more important to design  
236 algorithms and special software to control the quality of the data. This is where the Data  
237 Quality Monitoring (DQM) plays a critical in the maintainability of the experiment, the  
238 operation efficiency and performs a reliable data certification. The high-level goal of  
239 the system is to discover and pinpoint errors, problems occurring in detector hardware  
240 or reconstruction software, early, with sufficient accuracy and clarity to maintain good  
241 detector and operation efficiency. The DQM workflow consists of 2 types: **Online** and  
242 **Offline**.

243 The **Online** DQM consists of receiving data taken from the event and trigger his-  
244 tograms to produce results in the form of monitoring elements like histogram references  
245 and quality reports. This live monitoring of each detector's status during data taking gives  
246 the online crew the possibility to identify problems with extremely low latency, mini-  
247 mizing the amount of data that would otherwise be unsuitable for physics analysis. The  
248 scrutiny of the Online DQM is a 24/7 job that consists of people or shifters that work at the  
249 CMS control center constantly monitoring the hundreds of different plots and histograms  
250 produced by the DQM software. This consumes a lot of manpower and is strenuous work.

251 The **Offline** DQM is more focused on the full statistics over the entire run of the  
252 experiment and works more on the data certification. In the offline environment, the  
253 system is used to review the results of the final data reconstruction on a run-by-run basis,  
254 serving as the basis for certified data used across the CMS collaboration in all physics  
255 analyses. In addition, the DQM framework is an integral part of the prompt calibration  
256 loop. This is a specialized workflow run before the data are reconstructed to compute and  
257 validate the most up-to-date set of conditions and calibrations subsequently used during

<sup>258</sup> the prompt reconstruction.

<sup>259</sup> This project aims to minimize the DQM scrutiny by eye and automate the process so  
<sup>260</sup> that there is a more efficient process to monitor the detector and the quality of the data by  
<sup>261</sup> implementing Machine Learning techniques.

# <sup>262</sup> Chapter 4

## <sup>263</sup> What is Machine Learning?

<sup>264</sup> Machine Learning (ML) can be defined as an application of Artificial Intelligence that  
<sup>265</sup> permits the computer system to learn without being told explicitly. In ML a computer  
<sup>266</sup> program is said to learn from experience E with respect to some class of tasks T and  
<sup>267</sup> performance measure P, if its performance at tasks in T, as measured by P, improves  
<sup>268</sup> with experience E [3]. ML has made tremendous strides in the past decades and has  
<sup>269</sup> become very popular recently due to its multifaceted applications. It is being used on  
<sup>270</sup> social media, marketing, and in the scientific community as well. Some examples of  
<sup>271</sup> ML applications are: the algorithms used on application in smartphones to detect human  
<sup>272</sup> faces, self-driving cars, computer games, stock prediction, and voice recognition. An  
<sup>273</sup> interesting characteristic of ML algorithms is that the more data one inputs the better is  
<sup>274</sup> the performance. The ML application has a very wide spectrum covering almost every  
<sup>275</sup> aspect of human endeavor that involves a lot of data. Scientific analysis today generates  
<sup>276</sup> enormous data and is hence a perfect use case to apply ML techniques. In this work  
<sup>277</sup> we use enhanced ML techniques based on progress in the recent past.

<sup>278</sup> In general, there are two main categories to classify machine learning problems: **Su-**  
<sup>279</sup> **pervised Learning** (SL) and **Unsupervised Learning** (UL). SL is the most used ML  
<sup>280</sup> approach and has proven to be very effective for a wide variety of problems. Examples  
<sup>281</sup> of common SL problems are: spam filters, predicting housing prices, identifying a ma-  
<sup>282</sup> lignant or benign tumor, etc. These types of problems are characterized by providing a

283 “right answer” as a reference. For example, spam filter algorithms identify emails that  
284 are spams by training on a dataset that has examples of such emails. In case of predicting  
285 house prices, the algorithm is trained on a dataset of houses involving features like the  
286 area of the house, number of rooms, and the selling price of the house.

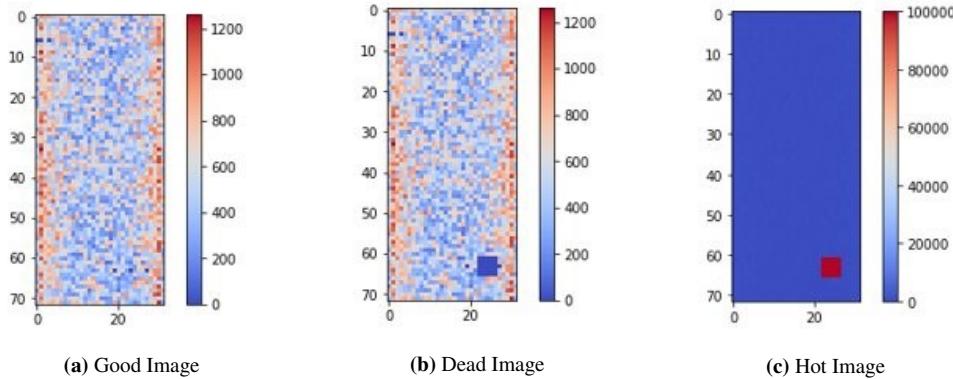
287 UL algorithms are different in the sense that they do not have the “right answers”  
288 given to the machine. Instead, UL algorithms are used for finding patterns and make  
289 clusters from the given data. That is what also forms the basis of a search engine (e.g.  
290 Google news). Clicking on a link to a news article, one gets many different stories of  
291 different journals that have some correlation with the article searched. This happens be-  
292 cause the ML algorithm is capable of learning features and shared patterns from a bunch  
293 of data without being given any specifics. Another interesting UL problem is the so-called  
294 “cocktail party” that involves distinguishing the voice of two people recording on two mi-  
295 crophones located at different places. The ML algorithm is able to separate the sources of  
296 the voices in the recordings by learning the voice features that correspond to each person,  
297 showing the power of the UL algorithm.

298 In this study, I have focused on an SL approach and a variant of the UL approach,  
299 called the **Semi-Supervised Learning** approach (SSL). The SSL is named so because  
300 the data involves looking at images that are already known to be “Good” but one doesn’t  
301 necessarily know every possible situation that produces a “Bad” image. The purpose is to  
302 define a metric for a “good” image and subsequently decide if an image is “bad” in case  
303 it deviates too much from an acceptable value.

## 304 4.1 Developing the Algorithm

305 To develop an ML algorithm the following are taken into consideration, what is the  
306 task? and what is the method to approach the task? In our case, we are looking into images  
307 that have information about the activity that the channels in the HCAL are detecting.  
308 These images are called ”occupancy maps” and they are a visual way of monitoring the  
309 health of the detector itself (see [Figure 4.1](#)). There are two common problems that can be

310 identified by viewing occupancy maps which are called "dead channels" and "hot towers".  
 311 They are referred to as "**dead**" and "**hot**" respectively in the rest of this document. Dead  
 312 channels mean that on a certain place in the occupancy map there is not any readout from  
 313 the channels on the HCAL and hot channels mean that there are channels that are being  
 314 triggered by noise or are damaged in a way that makes them readout too much activity.



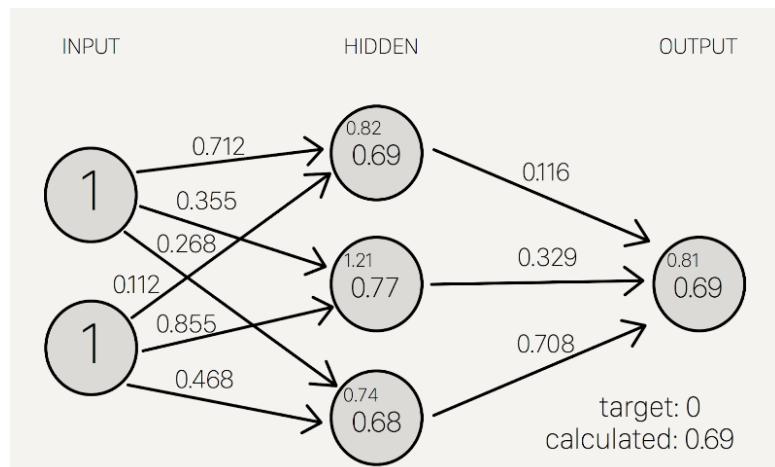
**Figure 4.1:** Occupancy maps with 5x5 affected regions

315 The problem is the following, to create a model that can detect and classify what  
 316 type of scenario is occurring on each occupancy map. For this, we want to go with a  
 317 SL approach which means that we will give the model the images as the input and it  
 318 will train on these images by learning to identify patterns or features in the image and  
 319 try to do a "fit" from the images to their corresponding labels. After the training, the  
 320 algorithm will be given a testing set for us to evaluate the model's ability to correctly  
 321 detect if there is a problem with the image and what type of problem is being detected.  
 322 The output of the model will be the predicted class of the test image. The predictions are  
 323 based on the labels and their corresponding images that were given to the model during  
 324 training. This means that if the model was trained with 3 different types of images with  
 325 their corresponding label the model will only work well for images that present similar  
 326 patterns or characteristics to those presented in the training. For example, if we only  
 327 train the model to distinguish between "good" and "hot" then when the model encounters  
 328 images that aren't either of these two, like an image labeled "dead", then the model will  
 329 not know what to do with this image and will give it an incorrect label. After the SL  
 330 model has been tested the next step is trying an SSL model. The term semi-supervised

simply means that there isn't a ground truth label that is being given to the model during training because either there isn't necessarily a ground truth, or we don't know what the ground truth is. What we do know, is what is considered as a "good" image and what this approach hopes to accomplish is to use the error in the reconstruction of the input image and use that information to discriminate between the "good" vs the "bad" images.

## 4.2 Teaching the Algorithm

The way an ML algorithm learns is by an iterative process called an optimization algorithm in which the predicted output value of the model is compared to the desired output (See [Figure 4.2](#)) and the weights and biases of the model are adjusted such that the predicted output is closer to the desired output.



**Figure 4.2:** Weights and Biases

"Optimization algorithms helps us to **minimize** (or **maximize**) an **Objective function** (*another name for Error function*)  $E(x)$  which is simply a mathematical function dependent on the Model's internal **learnable parameters** which are used in computing the target values(Y) from the set of *predictors*(X) used in the model. For example - we call the **Weights(W)** and the **Bias(b)** values of the neural network as its internal learnable *parameters* which are used in computing the output values and are learned and updated in the direction of optimal solution i.e. minimizing the **Loss** by the network's training process and also play a major role in the **training** process of the Neural Network Model." [4].

## Gradient Descent

The “Learning” in Machine Learning.

Update the values of X (punish) it when it is wrong.

$$X = X - \eta \nabla(X)$$

X: weights or biases

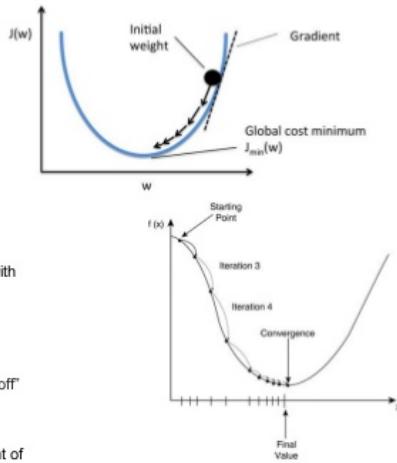
$\eta$ : Learning Rate (typically 0.01 to 0.001)

$\eta$  :The rate at which our network learns. This can change over time with methods such as Adam, Adagrad etc.  (hyperparameter)

$\nabla(X)$ : Gradient of X

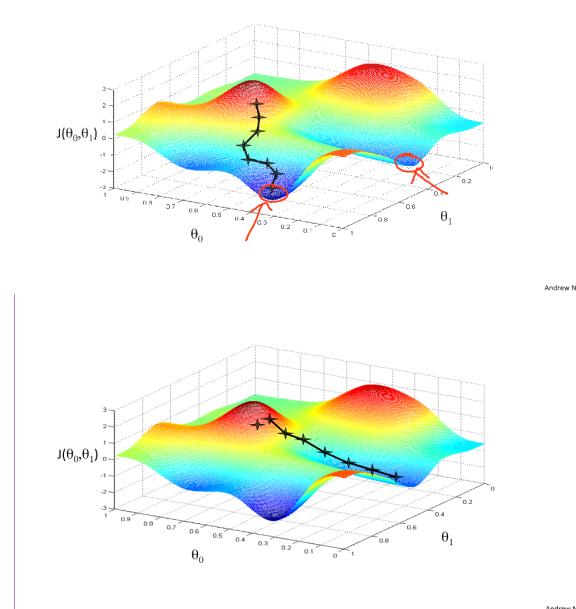
We seek to update the weights and biases by a value indicating how “off” they were from their target.

Gradients naturally have increasing slope, so we put a negative in front of it to go downwards



**Figure 4.3:** Gradient Descent algorithm

- 349 The most basic and probably the most used optimizer is called Gradient Descent (GD).
- 350 GD is based on the concept of using the gradient of a loss or cost function and moving
- 351 the weights and biases of the ML model so that the predicted value is taking a step in the
- 352 decreasing direction of this error function (See [Figure 4.3](#)). In general, the “terrain” of the
- 353 loss function is not a smooth bowl-shaped surface like the one present in the image. The
- 354 most general form of the surface is more similar to a rocky mountain (See [Figure 4.4](#)),
- 355 which presents a problem when using simple optimizers like GD.



**Figure 4.4:** Loss Function surface

356 **Chapter 5**

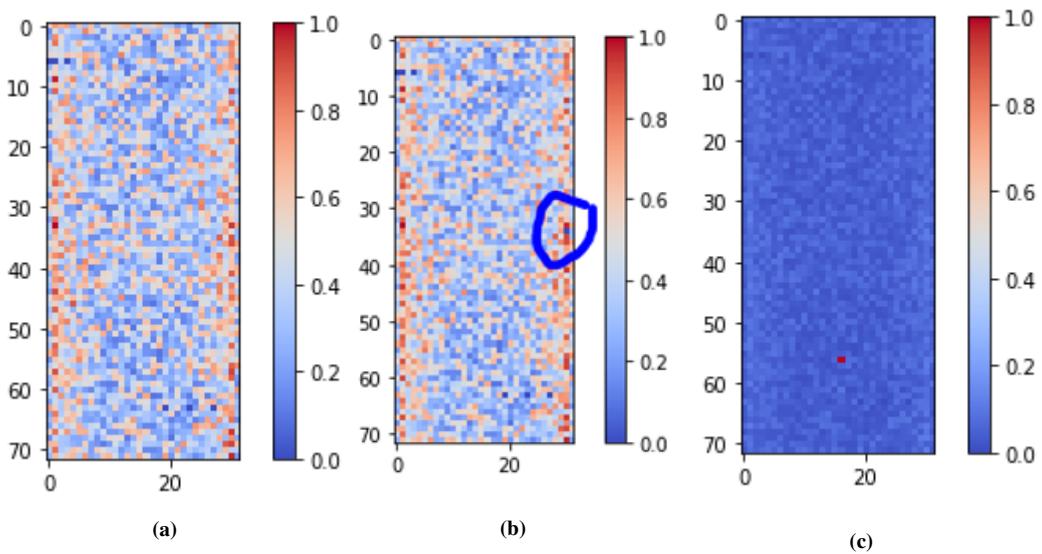
357 **Results**

358 Here first the limitations of Scikit-learn predefined ML models - Logistic Regres-  
359 sion(LR) and Multi-Layer-Perceptron(MLP), are described. The Logistic Regression  
360 Model seems to work almost perfectly with all 3 classes when the bad region size is  
361 5x5 (as in [Figure 4.1](#)) with either the same or randomized location. When the bad region  
362 size is 1x1 like in [Figure 5.1](#) the LR Model performs poorly with an accuracy of approxi-  
363 mately 20%. The MLP does not seem to work in any of the used cases that are studied as  
364 it always performs poorly with an accuracy of 40%.

365 d vice-versa. [Table 5.1](#) shows the cut values that are applied to photons that are found  
366 within both the ECAL barrel and endcap range. The associated values to the efficiency  
367 and the background rejection rate are shown for each of the three different photon ID se-  
368 lections.

369

370 In order to obtain the high efficiency and background rejection rates shown, a robust  
371 set of identification and isolation criteria are selected. A total of five parameters are  
372 used for this simple cut based method. For photon identification, the H/E and the  $\sigma I\eta I\eta$   
373 variables are found to provide the best results. The H/E parameter is defined as the ratio  
374 of the HCAL tower energy over the ECAL seed cluster energy. A threshold value is  
375 selected on H/E to remove background from electrons that are detected in both the ECAL  
376 and HCAL but have no reconstructed track [\[5\]](#). The  $\sigma I\eta I\eta$  variable is known as the



**Figure 5.1:** Occupancy Maps with 1x1 bad regions. A) Good image B) Dead image C) Hot image

377 photon shower shape variable and is defined in the ECAL as the energy weighted standard  
 378 deviation of a single crystal within the  $5 \times 5$  crystal  $\eta$  range, centered around the crystal  
 379 with maximum energy [6]. This variable is a key component in the identification of both  
 380 electrons and photons since it provides a measure of the shower width where most of the  
 381 energy has been deposited in a given ECAL crystal.

**Table 5.1:** Identification and isolation cut values for photons provided by the CMS EGM POG. Values are provided for the three working points described (loose, medium and tight) for both the ECAL barrel and endcaps. The photon selection efficiency of the three working points, as well as their associated background rejection rate, are provided.

Barrel	Loose (90.06%)	Medium (80.19%)	Tight (70.01%)
<b>Background Rejection</b>	<b>Loose (85.73%)</b>	<b>Medium (88.87%)</b>	<b>Tight (90.66%)</b>
H/E	0.105	0.035	0.020
$\sigma I\eta I\eta$	0.0103	0.0103	0.0103
$\rho$ -corrected PF charged hadron isolation	2.839	1.416	1.158
$\rho$ -corrected PF neutral hadron isolation	$9.188 + 0.0126 \cdot p_T^\gamma + 0.000026 \cdot p_T^{\gamma^2}$	$2.491 + 0.0126 \cdot p_T^\gamma + 0.000026 \cdot p_T^{\gamma^2}$	$1.267 + 0.0126 \cdot p_T^\gamma + 0.000026 \cdot p_T^{\gamma^2}$
$\rho$ -corrected PF photon isolation	$2.956 + 0.0035 \cdot p_T^\gamma$	$2.952 + 0.0040 \cdot p_T^\gamma$	$2.065 + 0.0035 \cdot p_T^\gamma$

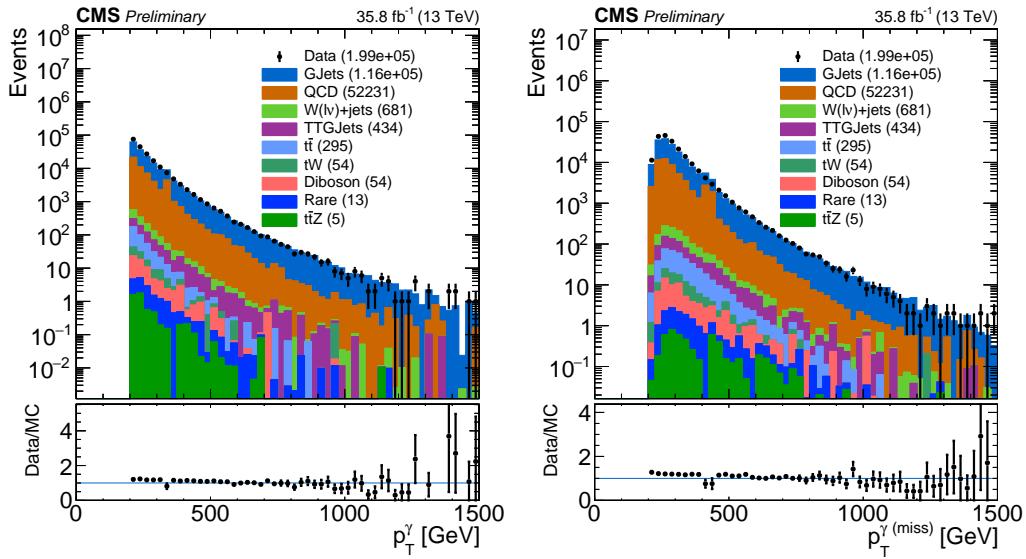
End Cap	Loose (90.81%)	Medium (80.06%)	Tight (70.11%)
<b>Background Rejection</b>	<b>Loose (76.90%)</b>	<b>Medium (81.50%)</b>	<b>Tight (84.34%)</b>
H/E	0.029	0.027	0.025
$\sigma I\eta I\eta$	0.0276	0.0271	0.0271
$\rho$ -corrected PF charged hadron isolation	2.150	1.012	0.575
$\rho$ -corrected PF neutral hadron isolation	$10.471 + 0.0119 \cdot p_T^\gamma + 0.000025 \cdot p_T^{\gamma^2}$	$9.131 + 0.0119 \cdot p_T^\gamma + 0.000025 \cdot p_T^{\gamma^2}$	$8.916 + 0.0119 \cdot p_T^\gamma + 0.000025 \cdot p_T^{\gamma^2}$
$\rho$ -corrected PF photon isolation	$4.895 + 0.0040 \cdot p_T^\gamma$	$4.095 + 0.0040 \cdot p_T^\gamma$	$3.272 + 0.0040 \cdot p_T^\gamma$

382 The other three parameters considered, comprise the isolation portion of the photon  
 383 selection cuts. These are the  $\rho$ -corrected particle flow (PF) charged hadron, neutral hadron  
 384 and photon isolation parameters. As can be seen from [Table 5.1](#), two of these parameters  
 385 (the neutral hadron and photon isolation) have a dependence on  $p_T^\gamma$ . These cuts are used  
 386 to ensure that the identified photon is well isolated within its own cone and by rejecting  
 387 photons that are identified within close proximity to either a charged or a neutral hadron  
 388 [[7](#)]. The value  $\rho$  included in the name of each of these parameters refers to the total

389 pileup density [8]. Therefore, the term “ $\rho$ -corrected” implies that these values, which are  
 390 sensitive to the residual contamination that arises from pile-up, have been corrected to  
 391 include these contributions.

### 392 5.0.1 Photon Selection

393 The event selection process for the  $\gamma$ +jets control region starts with photon candidates  
 394 that have a  $p_T > 200$  GeV and are within the acceptance range of the CMS ECAL (given  
 395 by  $|\eta| < 1.4442$  for the barrel and  $1.566 < |\eta| < 2.5$  for the endcaps). The photons are  
 396 subjected to pass the loose ID/isolation cuts described in ?? in order to remove  $\sim 85\%$   
 397 of the background processes and obtain a prompt photon sample that is  $\sim 90\%$  pure, on  
 398 average. Additional restrictions include some of the same requirements imposed on the  
 399 signal baseline selection discussed in ???. These include  $N_j \geq 4$ , an  $H_T > 300$  GeV,  
 400 the  $\Delta\phi$  requirements for leading jets and the lepton vetoes described in ???. The lepton  
 401 veto, in particular, greatly improves the prompt photon selection by removing many of  
 402 the events in the simulated samples where a lepton gets misidentified as a photon.



**Figure 5.2:** Shown are both the  $p_T^\gamma$  (left) and  $p_T^{\gamma(\text{miss})}$  (right) distributions before applying any corrections.  $p_T^{\gamma(\text{miss})}$  is obtained by adding the  $p_T^\gamma$  to the total  $p_T^{\text{miss}}$  in every event.

403 To further emulate the  $Z \rightarrow \nu\bar{\nu} + \text{jets}$  background, a variable in which the photons are  
 404 treated as  $p_T^{\text{miss}}$  is defined. We call this variable  $p_T^{\gamma(\text{miss})}$  and we obtain it by adding the

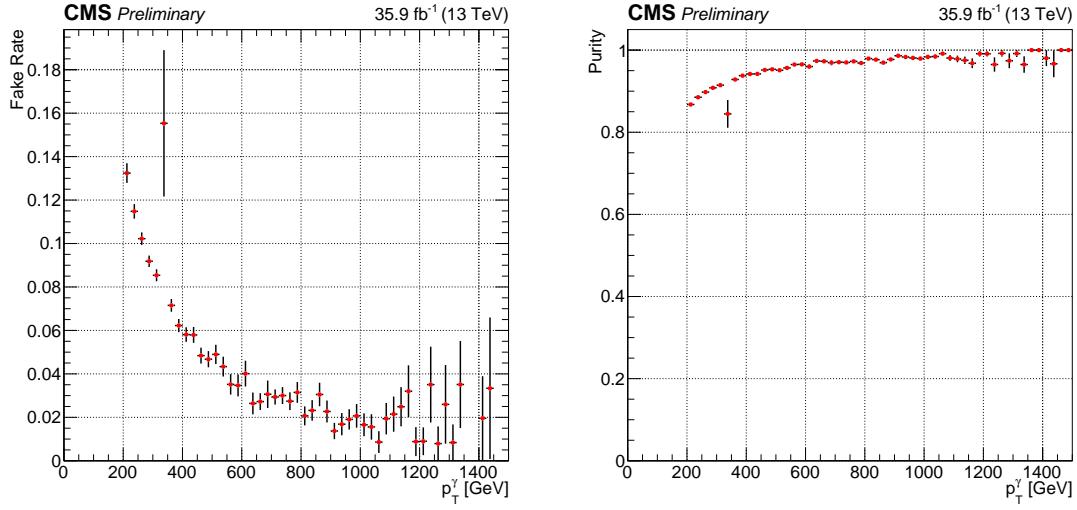
405  $p_T^\gamma$  for every event to the total  $p_T^{miss}$  in the event. Both the  $p_T^\gamma$  and the resulting  $p_T^{\gamma(miss)}$   
 406 distributions are shown in [Figure 5.2](#) as data/MC comparison plots, where the simulated  
 407 backgrounds are stacked in order of ascending contribution.

408

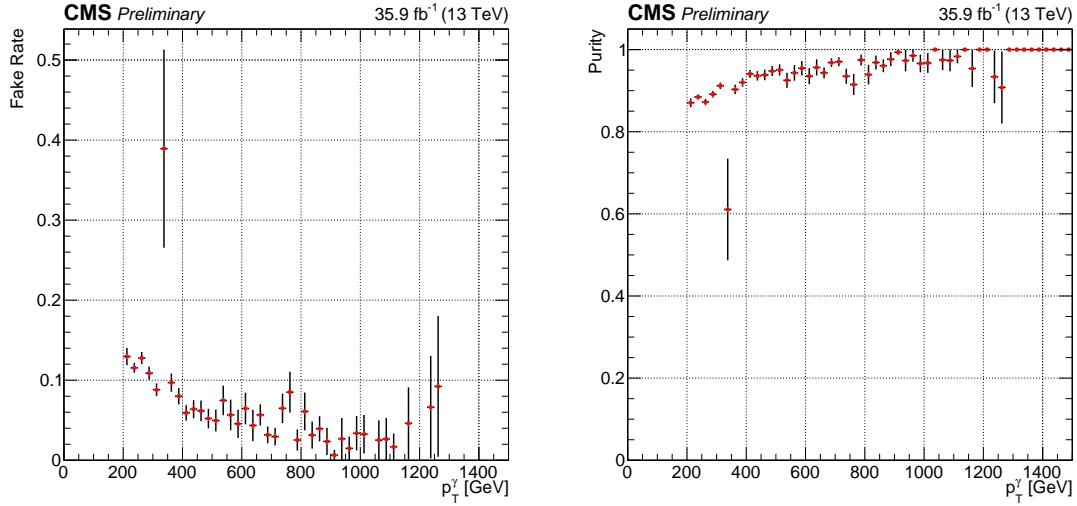
409 The main contributions from simulation arise from the  $\gamma+jets$ , QCD and to a lesser  
 410 extent,  $t\bar{t}\gamma$ . Other non-dominant backgrounds in the control region include contributions  
 411 from  $W(l\nu)+jets$ ,  $t\bar{t}$ , Diboson,  $tW$ ,  $t\bar{t}Z$  and rare processes. Most of these lesser back-  
 412 grounds are nearly negligible (several orders of magnitude lower than the dominant back-  
 413 grounds) and are considered to be mostly composed of fake photons. In addition to the  
 414 cuts described, all of the simulation samples are subjected to weights that apply correc-  
 415 tions to pileup as well as the b-tagging efficiency. Data, on the other hand, is obtained  
 416 from a sample that contains events with at least one identified photon. Photons in this  
 417 sample are also subjected to the high-level trigger HLT\_Photon175, which restricts the  
 418 selection to photons that have a  $p_T > 175$  GeV. Both simulation and data are subjected to  
 419 the same selection criteria established in this section.

## 420 **5.0.2 Photon Purity and Fake Rate**

421 Three different types of photons make up the  $\gamma+jets$  CS: prompt photons, produced  
 422 either directly or through fragmentation, and fake photons. Prompt photons are defined  
 423 as photons which are formed shortly after the proton-proton collision (i.e. before the pro-  
 424 duced quarks and gluons have had enough time to form hadrons). Two types of photons  
 425 fit in this category. The first type, which we designate as direct photons, are photons that  
 426 are produced directly from the proton-proton interaction [\[9\]](#). A secondary type of prompt  
 427 photon, that is virtually indistinguishable from the direct photons at the detector level,  
 428 originates from the decay of  $\pi^0$  mesons and are called fragmentation photons. The final  
 429 type of photon found in the CS corresponds to fake (or non-prompt) photons. The fake  
 430 photon contribution typically arises from leptons (mostly electrons) whose tracks are not  
 431 properly reconstructed, yet leave energy measurements in the ECAL.



**Figure 5.3:** Plots for Fake Rate (left) and Purity (right) as a function of the photon  $p_T$  are shown. The events are selected are required to have a  $p_T > 200$  GeV, be within the ECAL acceptance range, and pass the loose ID selection cuts. This selection was produced in order to verify the values given by the E/ $\gamma$  POG. As can be seen, the efficiency (purity) is seen to agree with the values of the loose photon ID/isolation selection.



**Figure 5.4:** Plots for Fake Rate (left) and Purity (right) as a function of the photon  $p_T$  are shown. These plots include photons with the full control region selection. Aside from exhibiting lower statistics, the plots seem to agree with the fake rate and purity before all the control region cuts are applied.

432     In order to identify prompt photons, reconstructed photons from the  $\gamma + \text{jets}$  and QCD  
 433    samples are matched to generator-level photons in space and momentum by requiring  
 434     $\Delta R(\gamma_{\text{gen}}, \gamma_{\text{reco}}) < 0.4$  and  $0.5 < p_T^{\text{gen}}/p_T^{\text{reco}} < 2.0$ , respectively. Any reconstructed photon  
 435    which fails to get matched to a generator level photon is labeled as a fake/non-prompt  
 436    photon. Direct photons are identified by further requiring that the reconstructed photons  
 437    be matched to a parton (a gluon or quark) in space as  $\Delta R(\gamma, \text{parton}) > 0.4$ . This require-

438 ment is intended to distinguish the reconstructed photons from highly boosted  $\pi^0$ 's, which  
 439 compose a large portion of the experimentally indistinguishable fragmentation photons.  
 440 Finally, fragmentation photons are obtained exclusively from QCD simulation and are re-  
 441 quired to have  $\Delta R(\gamma, \text{parton}) < 0.4$  in order to avoid double counting photons from the  
 442  $\gamma+\text{jets}$  sample.

443

444 With all three types of photons defined, a study can be carried out from simulation  
 445 to estimate their respective contributions to the defined control region. The study takes  
 446 into account that any reconstructed photon in the  $\gamma+\text{jets}$  or QCD samples can only be  
 447 categorized as prompt (through direct production or fragmentation) or non-prompt (fake).  
 448 The purity and fake rate can then be defined in terms of the relative proportions of prompt  
 449 or non-prompt photons with respect to the sum of the contributions from all three types of  
 450 photons. Identified direct photons are taken from the  $\gamma+\text{jets}$  sample exclusively. Mean-  
 451 while, the fragmentation and fake photon contributions are taken from the QCD sample.  
 452 The three quantities are then added together and their respective contributions are deter-  
 453 mined in terms of the photon  $p_T$  ([Figure 5.3](#)).

454

455 The photon purity ([Figure 5.3](#) and [Figure 5.4](#), right) is defined in terms of the prompt  
 456 and non-prompt photons as:

457

$$458 p_\gamma = \frac{\text{prompt}}{\text{prompt} + \text{fake}},$$

459 where the prompt photon portion comes from the sum of the direct photons (extracted  
 460 from the  $\gamma+\text{jets}$  sample) and the fragmentation photons (extracted from the QCD sample).  
 461 The remaining non-prompt (or fake) photons all come from photons in the QCD sample  
 462 that were not matched to truth-level photons in space and momentum with the specified  
 463 required conditions. Meanwhile, the photon fake rate ([Figure 5.3](#) and [Figure 5.4](#), left) is  
 464 defined from this same combination of samples as:

465

$$466 \quad f = \frac{fake}{prompt+fake} ,$$

467     Figure 5.3 shows the purity and fakerate for photons that pass the loose ID/selection,  
 468    have a  $p_T > 200$  GeV and are within the ECAL acceptance range. A sample is obtained  
 469    in which 77% of the photons are direct, 12% are fragmentation and 11% are fakes. This  
 470    implies an average purity of  $\sim 89\%$  for this sample, well within the value that is expected.  
 471    Figure 5.4 shows the same ratios for the loose  $\gamma +$  jets control region described in subsection 5.0.1.  
 472    Although the amount of statistics has decreased due to the additional cuts, a  
 473    similar trend can be observed.

## 474    5.1 The $Z \rightarrow \mu^+ \mu^-$ Control Region

475    The  $Z \rightarrow \mu^+ \mu^-$  control region defined in this section is in every respect identical to  
 476    the one applied in the 2016 analysis (??). The only difference between the 2016 method  
 477    and the one discussed in this chapter is that the Drell-Yan (DY) sample is only used for  
 478    the normalization correction of the  $Z \rightarrow \nu \bar{\nu}$  background. Therefore, the loose  $\mu\mu$  control  
 479    region is not used or applied in the calculation of the scale factors. In the following  
 480    subsections only the tight  $\mu\mu$  control region, and its usage to obtain the normalization  
 481    scale factor  $R_{norm}$ , is discussed.

### 482    5.1.1 Muon ID and Isolation

483    The muons are selected using the “medium muon” selection [10], per the recommen-  
 484    dation of the muon POG. The muon candidates in this selection satisfy  $p_T > 10$  GeV and  
 485     $|\eta| < 2.4$ . Other additional cuts are applied to aid in the muon candidate selection, such as  
 486    an impact parameter cut. Muons are also subjected to a PF relative-isolation (also referred  
 487    to as mini-isolation) in which the cone size is inversely proportional to the muon  $p_T$ . This  
 488    requirement enforces the  $p_T$  within the isolation cone to be at most 20% of the muon  $p_T$

489 in order to eliminate events with an isolated muon. Details of the medium photon selec-  
 490 tion are included in [Table 5.2](#) and [Table 5.3](#), while details of the impact parameter cut are  
 491 summarized in [Table 5.4](#).

Muon Medium ID	
Loose muon ID	Yes
Fraction of valid tracker hits >	0.80
Good Global muon OR Tight segment compatibility >	Yes OR 0.451

**Table 5.2:** Muon Medium ID 2016 HIP Safe

Good Global muon	
Global muon	Yes
Normalized global-track $\chi^2 <$	3
Tracker-Standalone position match <	12
Kick finder <	20
Segment compatibility >	0.303

**Table 5.3:** Muon Medium ID HIP Safe Good Global Muon

Muon Impact Parameter	
d0 <	0.2
dz <	0.5

**Table 5.4:** Additional Impact Parameter cut on Muons

### 492 5.1.2 Muon Selection in the Tight Control Region

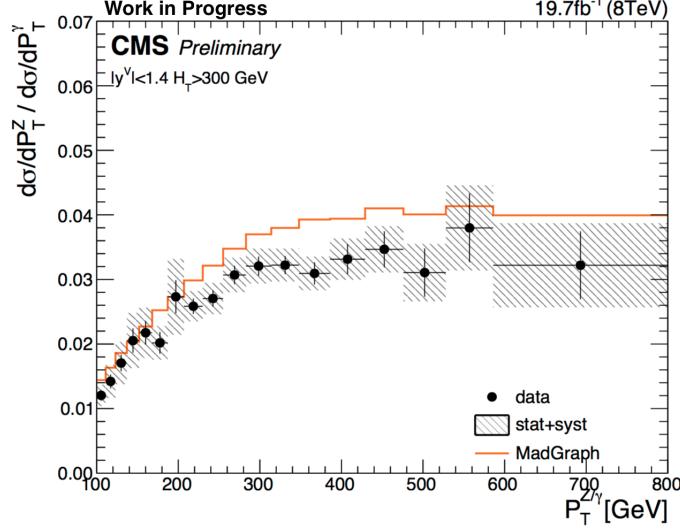
493 Events are selected from data samples that contain exactly two oppositely charged  
 494 muons ( $\mu^+ \mu^-$ ), which fall within the invariant mass  $81 < m_{ll} < 101$  GeV window for  
 495 the Z boson. Additional cuts for the tight muon selection include baseline requirements  
 496 such as an  $H_T > 300$  GeV,  $N_j \geq 4$ , the  $\Delta\phi$  baseline cut on leading jets, a  $p_T^{miss} > 250$   
 497 GeV, an  $m_{T2} > 200$  GeV and at least 1 top-tagged jet  $N_t \geq 1$ . In addition, the  $p_T$  of  
 498 the two muons are required to be  $p_T > 50$  GeV for the leading muon and  $p_T > 20$  GeV  
 499 for the sub-leading one. The only difference, when compared to the signal region is the  
 500 missing lepton veto, in addition to the dimuon events being treated as  $p_T^{miss}$ . This makes  
 501 for a region that exhibits very similar kinematics to the  $Z \rightarrow \nu\bar{\nu}$  signal region, yet suffers  
 502 from a lack of statistics.

## 503 5.2 Analysis

504 In this section a detailed explanation of the calculation of the scale factors for both  
 505 shape and normalization corrections is provided. The following methods make use of  
 506 the loose  $\gamma$ +jets and the tight  $\mu\mu$  control regions defined in the previous sections. The  
 507 procedure involves extracting the shape corrections  $S_\gamma$  from the  $\gamma$ +jets control region and  
 508 afterwards obtain a single normalization correction factor  $R_{norm}$  from the tight  $\mu\mu$  control  
 509 region. Both factors will then be applied to the final prediction of the  $Z \rightarrow \nu\bar{\nu}$  background  
 510 in each of the required search bins.

### 511 5.2.1 Shape Correction Using the $\gamma$ + jets Control Sample

512 In this section the validation of the  $\gamma$ +jets simulation is discussed in terms of the  
 513 shape of the loose photon control region. As it was shown in [subsection 5.0.2](#), this con-  
 514 trol region has high purity for  $\gamma$ +jets events, particularly in regions of high  $p_T$  ( $\gtrsim 300$ ).  
 515 In order to apply this correction factor it is assumed that the shape differences between  
 516 data and simulation are similar between  $Z \rightarrow \nu\bar{\nu}$  and  $\gamma$ +jets events. This assumption is  
 517 validated in studies which compare the cross-section ratio of  $Z$ +jets to  $\gamma$ +jets events [11].  
 518 [Figure 5.5](#) shows the results of this study, conducted in 2014, for both data and Mad-  
 519 Graph simulation with an integrated luminosity of  $19.7\text{fb}^{-1}$  and a center-of-mass energy  
 520 of  $8\text{ TeV}$ . It can be seen that for values of  $p_T^{Z/\gamma} \gtrsim 300\text{ GeV}$ , the ratio of the cross-section of  
 521 both processes becomes nearly constant. It is then a matter of applying a factor to account  
 522 for the difference in the amount of events between the  $Z$ +jets and  $\gamma$ +jets events in order  
 523 to obtain the total amount of  $Z \rightarrow \nu\bar{\nu}$ +jets events. This factor is obtained from the tight  
 524  $\mu\mu$  control region, as shown in [subsection 5.2.2](#).

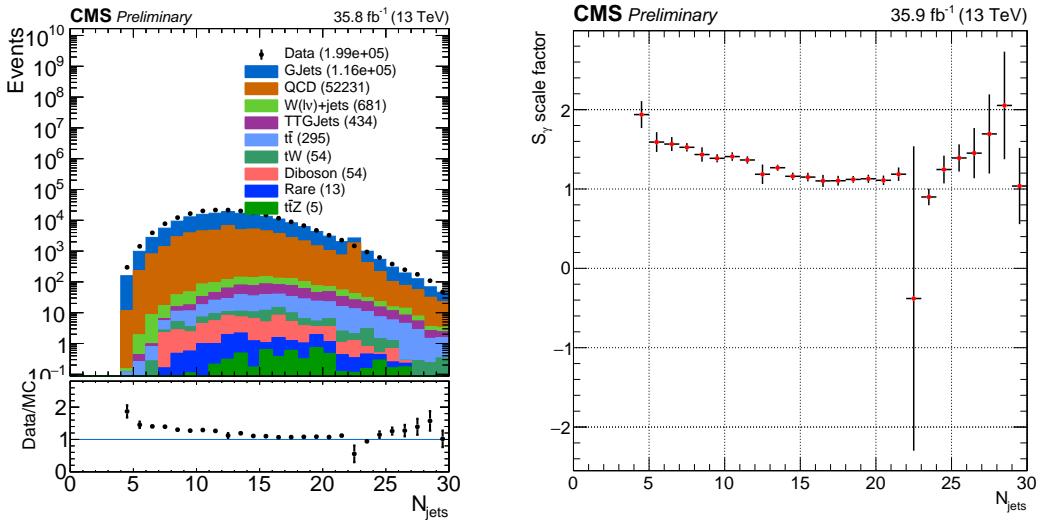


**Figure 5.5:** Results of study of the  $Z$ +jets to  $\gamma$ +jets cross-section ratio for both data and MadGraph simulation. For high values of the vector boson transverse momentum, the ratio between these processes is observed to be nearly constant.

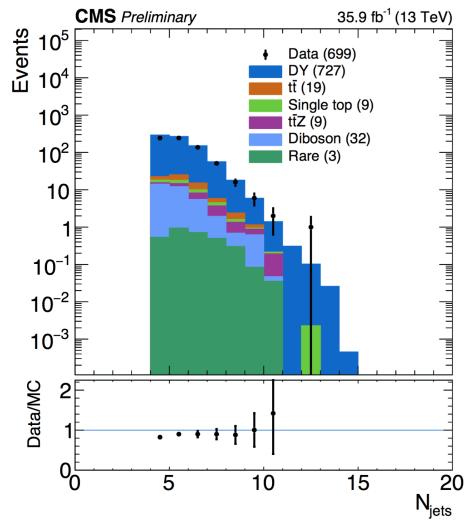
525 In order to obtain the shape corrections, the ratio between data and simulation of the  
 526 jet multiplicity distribution is used (Figure 5.6). This is due to it exhibiting the highest  
 527 difference between the observed data and MC. The re-weight for the  $\gamma$ +jets simulation  
 528 sample is then accomplished by applying the  $N_{jet}$  dependent factor  $S_\gamma(N_j)$ . This scale  
 529 factor is determined by taking the ratio of the data and simulation, after subtracting all  
 530 other MC samples from data events:

$$531 \quad S_\gamma^i = \frac{\text{Data}^i - \text{MC}_{\text{other}}^i}{\text{MC}_{\gamma+\text{jets}}^i},$$

532 where  $i$  denotes any given bin in the  $N_j$  distribution. The shape correction factors  $S_\gamma^i$   
 533 are displayed graphically in Figure 5.6 (right) for each  $N_j$  bin. These factors correct for  
 534 differences in the jet multiplicity shape, while the overall normalization is estimated from  
 535 the tight  $\mu\mu$  control region. Figure 5.7 shows the  $N_j$  distribution in the tight  $\mu\mu$  control  
 536 region after the calculated scale factors have been applied. The  $S_\gamma$  correction will be  
 537 applied to the  $Z \rightarrow \nu\bar{\nu}$  simulation final prediction for each of the analysis search bins. The  
 538 uncertainty associated with the scale factor is estimated from the event yields in the loose  
 539 photon control region. This uncertainty will form part of the total systematic uncertainty  
 540 in the final prediction.

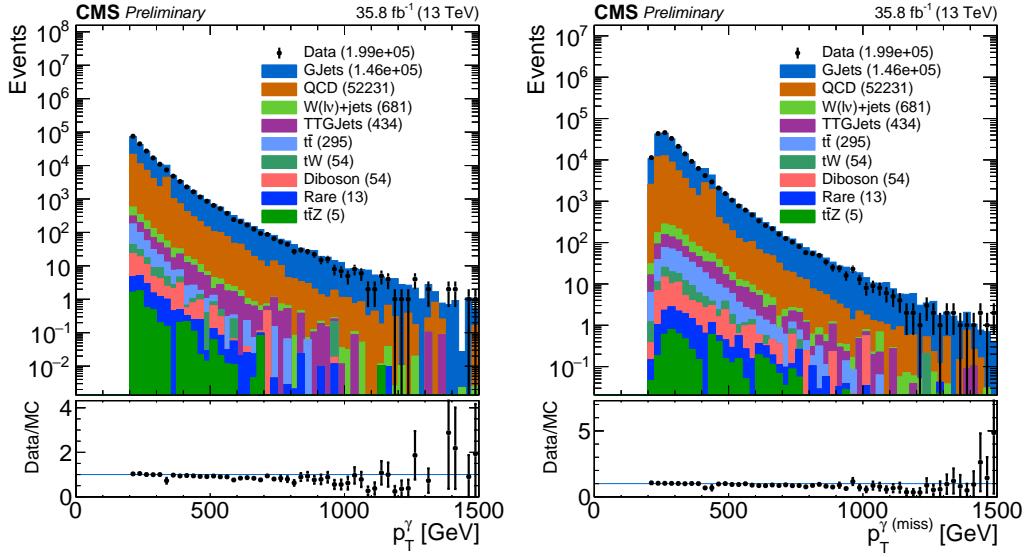


**Figure 5.6:** Jet multiplicity and the associated  $S_\gamma$  scale factor in the loose photon control region before any corrections are applied.

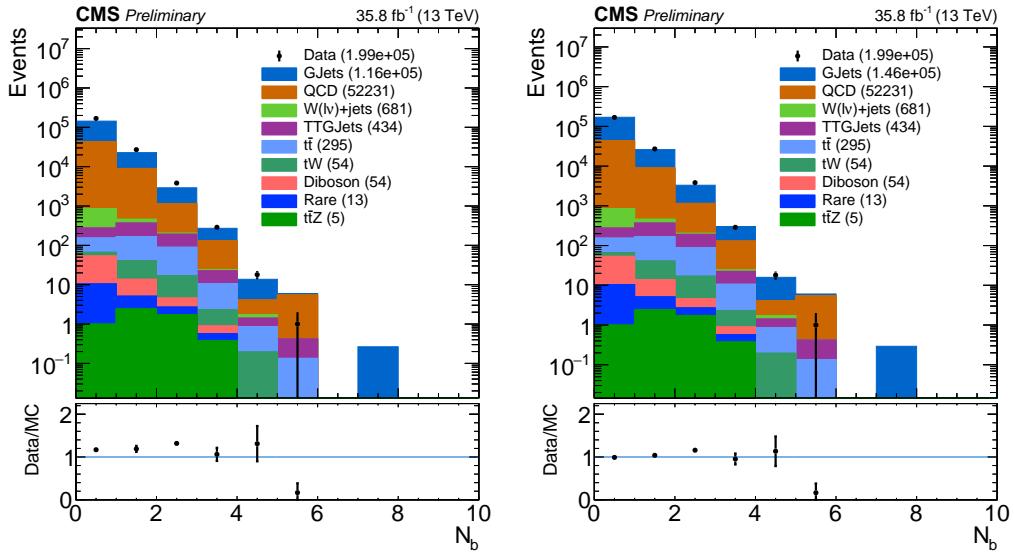


**Figure 5.7:**  $N_{jet}$  distribution in the tight  $\mu\mu$  control region after  $S_\gamma$  corrections.

541      The effect of the  $S_\gamma(N_j)$  scale factor is shown for various distributions. These results  
 542      show that the overall agreement between data and simulation improves after applying the  
 543      corresponding shape corrections.



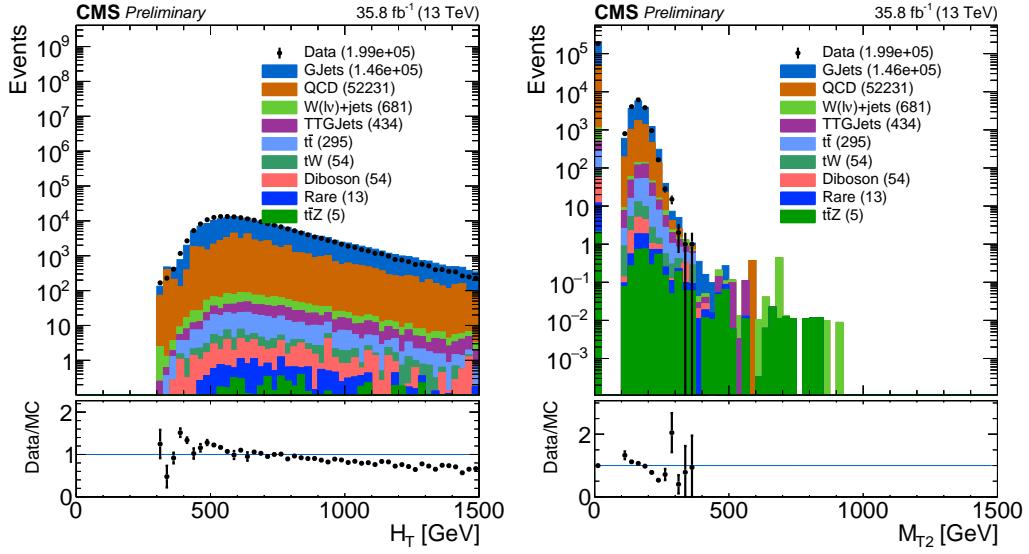
**Figure 5.8:**  $p_T^\gamma$  (left) and  $p_T^{\gamma(\text{miss})}$  (right) distributions after applying the  $S_\gamma(N_j)$  scale factor. Comparing to Figure 5.2, an improvement in the agreement between data/MC can be observed.



**Figure 5.9:**  $N_b$  distribution before (left) and after (right) applying the  $S_\gamma(N_j)$  scale factor.

544    **5.2.2 Normalization Correction Using the tight  $Z \rightarrow \mu^+ \mu^-$  Control  
545    Sample**

546    In order to constrain the normalization of the  $Z \rightarrow \nu\bar{\nu}$  simulation sample, a normal-  
547    ization correction factor  $R_{norm}$  is calculated from the tight  $\mu\mu$  control region defined in  
548    subsection 5.2.2. Two categories are considered: the zero b-tagged jet category ( $N_b = 0$ ),  
549    and the  $\geq 1$  b-tagged jet category ( $N_b \geq 1$ ). Both of these categories are statistically

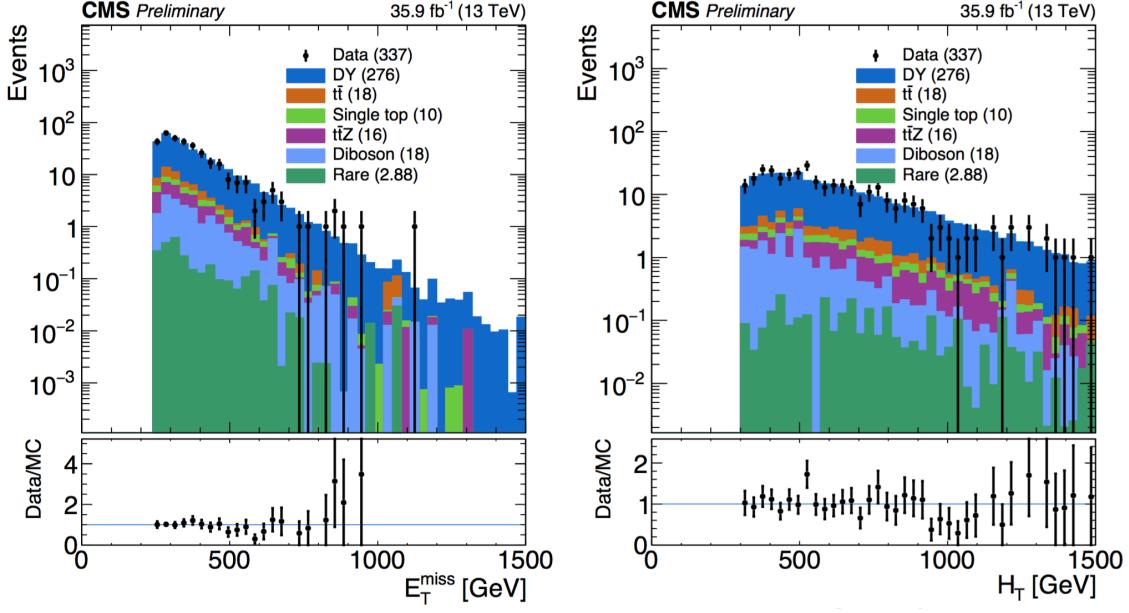


**Figure 5.10:**  $H_T$  and  $m_{T2}$  distributions applying the  $S_\gamma(N_j)$  scale factor.

550 consistent with each other but the inclusive region ( $N_b \geq 0$ ) has a lower overall un-  
 551 certainty. The method used to calculate the normalization scale factor requires that the  
 552  $N_j$ -dependent shape correction factors already be applied. Then, the  $R_{norm}$  factor can  
 553 be extracted from the ratio of the total event yield in data to that in the simulation. This  
 554 factor is found to be:

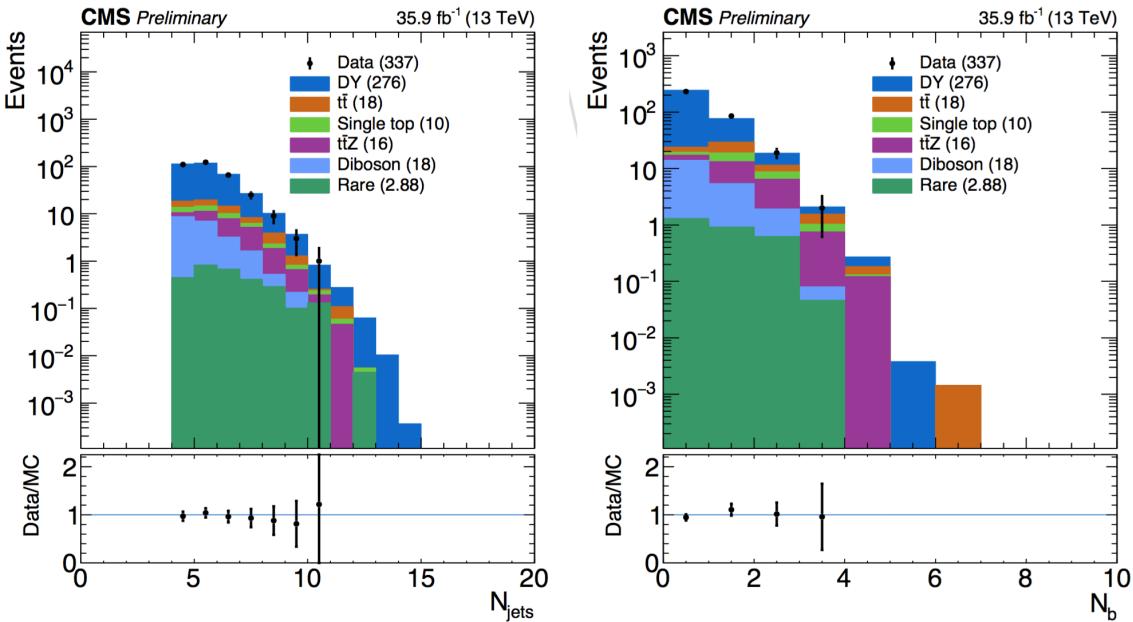
$$555 \quad R_{norm} = 1.070 \pm 0.085,$$

556 where the uncertainty includes only the associated statistical uncertainties on data and  
 557 simulation. This uncertainty is found to be propagated to the final background prediction,  
 558 see subsection 5.3.1.



**Figure 5.11:** Shown are data/MC comparisons for the  $p_T^{\text{miss}}$  (left) and  $H_T$  (right) distributions after applying both the  $N_j$ -dependent shape corrections ( $S_\gamma$ ) and the global normalization scale factor ( $R_{\text{norm}}$ ).

560 Data/MC comparisons are shown in [Figure 5.11](#) and [Figure 5.12](#) after applying  $R_{\text{norm}}$   
 561 for several distributions in the study. With this final global scale factor all the required  
 562 ingredients for the central value of the  $Z \rightarrow \nu\bar{\nu}$  background prediction are obtained.



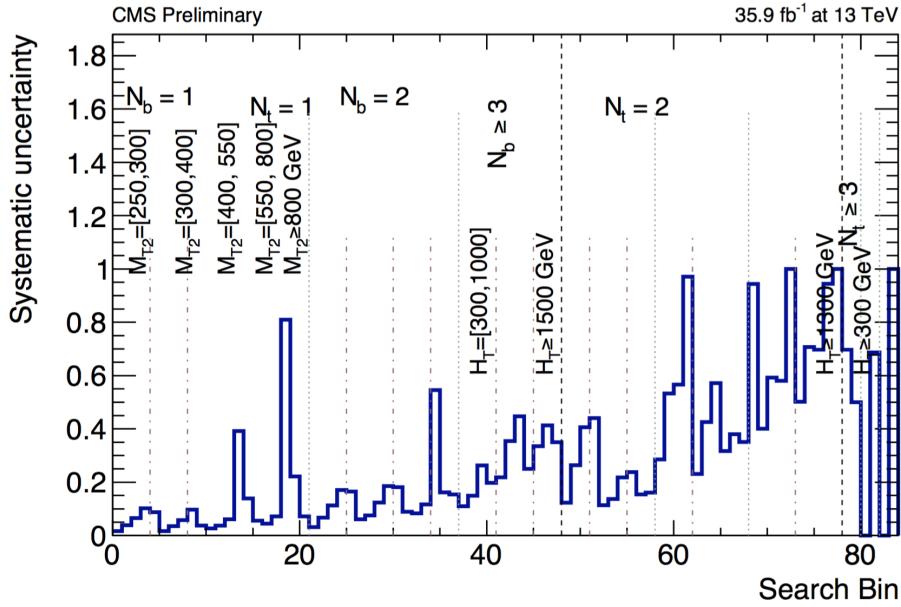
**Figure 5.12:** Shown are data/MC comparisons for the  $N_j$  (left) and  $N_b$  (right) distributions after applying both the  $N_j$ -dependent shape corrections ( $S_\gamma$ ) and the global normalization scale factor ( $R_{\text{norm}}$ ).

## 563 5.3 Results

564 In this section the results for the final estimation of the of the  $Z \rightarrow \nu\bar{\nu}$  are presented.  
 565 The current study includes preliminary results using only data obtained at the CMS detec-  
 566 tor during 2016. The results for this study are intended to confirm the assumption that the  
 567 additional  $\gamma + \text{jets}$  control region introduced in this analysis reduce the overall uncertain-  
 568 ties obtained in the 2016 analyses (described in ??). Furthermore, this study is intended  
 569 as a benchmark for future analyses of the SUSY stop group based in Fermilab and will be  
 570 the method used for the 2017 CMS data.

### 571 5.3.1 Systematics

572 Two categories of uncertainties for the  $Z \rightarrow \nu\bar{\nu}$  prediction are considered: uncertain-  
 573 ties that are associated to the use of MC simulation and the uncertainties specifically  
 574 associated to the background prediction method. Several sources are acknowledged in the  
 575 first category mentioned such as PDF and renormalization/factorization scale choices, jet  
 576 and  $p_T^{miss}$  energy scale uncertainties b-tag scale factor uncertainties, and trigger efficiency  
 577 uncertainties. Given that the simulation sample is normalized to data in the tight control  
 578 region, uncertainties associated with the luminosity and cross-section are excluded. In  
 579 addition, the overall  $Z \rightarrow \nu\bar{\nu}$  statistical uncertainty from MC simulation is also taken into  
 580 account.



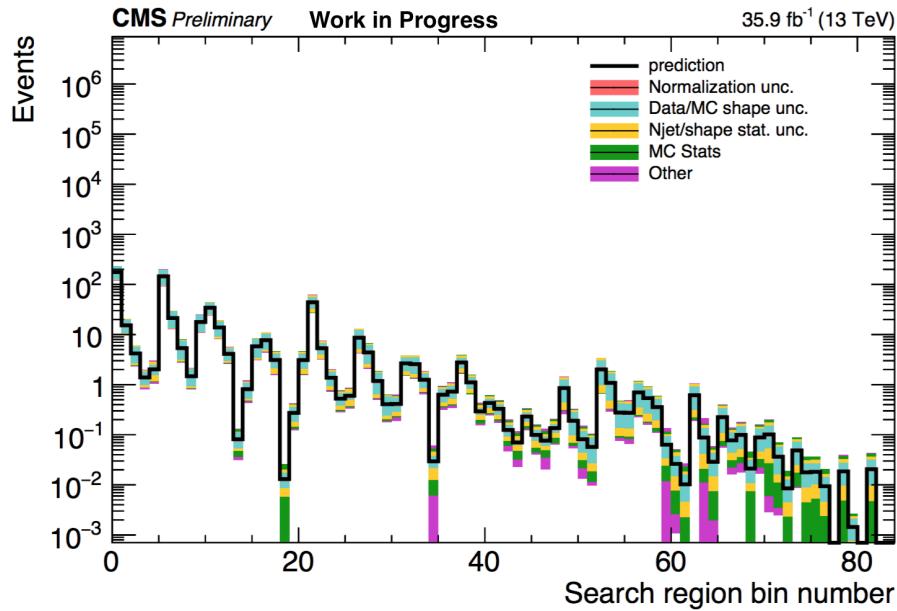
**Figure 5.13:** Systematic uncertainty in the final prediction, as a function of the search bin, associated to the MC statistics.

582 The statistical uncertainty associated with each bin in the MC is propagated as a sys-  
 583 tematic uncertainty. The relative uncertainty per bin can be see in [Figure 5.13](#). It shows  
 584 that the uncertainties for the MC vary from as low as 1% up to 81% and even 100% in  
 585 some regions. Since the final estimation is scaled using the global normalization factor  
 586 from the tight  $\mu\mu$  control region ( $R_{norm}$ ), the total uncertainty, due to limited amounts of  
 587 events in data, is propagated in the final prediction. This is also true for the  $S_\gamma(N_j)$  scale  
 588 factor, in which the residual differences in search variables other than  $N_j$  are evaluated in  
 589 the loose photon control region. Both the uncertainty arising from the  $N_j$  re-weighting  
 590 as well as the residual differences are evaluated together. The uncertainty from  $R_{norm}$  is  
 591 propagated as a flat value of 7.9% uncertainty per each search bin.

### 592 5.3.2 $Z \rightarrow \nu\bar{\nu}$ Estimation for the Search Bins

593 The final estimation for the  $Z \rightarrow \nu\bar{\nu}$  background calculated for all 84 search bins is  
 594 shown in [Figure 5.14](#). The statistical uncertainty in bins that have zero events is treated  
 595 as the average weight (the sum of the weights squared over the weight) times the poisson  
 596 error on 0 which is 1.8. This average weight is calculated on the basis of a relaxed cut in  
 597 which  $N_b \geq 2$  is required. For comparison, a cut in which  $N_t > 2$  where two tops are

<sup>598</sup> fake for the  $Z \rightarrow \nu\bar{\nu}$  is used.



**Figure 5.14:**  $Z \rightarrow \nu\bar{\nu}$  background prediction for all search bins, including the breakdown of the various uncertainties.

<sup>599</sup> **Appendix A**

<sup>600</sup> **Appendix Title**

601 **Appendix B**

602 **References**

- 603 [1] CERN, “Processing what to record?,” 2018.
- 604 [2] CMS, “The CMS Computing Project,” tech. rep., CERN, 2005.
- 605 [3] Coursera, “Machine learning,” 2018.
- 606 [4] A. S. Walia, “Types of optimization algorithms used in neural networks and ways to  
607 optimize gradient descent,” 2018.
- 608 [5] J. Berryhill, “Electron Detection at CMS.” The Egamma Physics Object Group,  
609 2009. [ppt](#).
- 610 [6] CMS Collaboration, “Performance of Photon Reconstruction and Identification with  
611 the CMS Detector in Proton-Proton Collisions at  $\sqrt{s} = 8$  TeV,” *JINST*, vol. 10,  
612 no. 08, p. P08010, 2015. [doi:10.1088/1748-0221/10/08/P08010](https://doi.org/10.1088/1748-0221/10/08/P08010).
- 613 [7] A. M. Sirunyan *et al.*, “Particle-flow reconstruction and global event description  
614 with the CMS detector,” *JINST*, vol. 12, no. 10, p. P10003, 2017. [doi:10.1088/1748-0221/12/10/P10003](https://doi.org/10.1088/1748-0221/12/10/P10003).
- 616 [8] J. de Favereau, C. Delaere, P. Demin, A. Giammanco, V. Lemaitre, A. Mertens, and  
617 M. Selvaggi, “DELPHES 3, A modular framework for fast simulation of a generic  
618 collider experiment,” *JHEP*, vol. 02, p. 057, 2014. [doi:10.1007/JHEP02\(2014\)057](https://doi.org/10.1007/JHEP02(2014)057).

- 619 [9] M. Klasen, C. Klein-Bosing, and H. Poppenborg, “Prompt photon produc-  
620 tion and photon-jet correlations at the LHC,” *JHEP*, vol. 03, p. 081, 2018.  
621 doi:10.1007/JHEP03(2018)081.
- 622 [10] S. Folgueras, “Baseline muon selections for Run-II.” <https://twiki.cern.ch/twiki/bin/view/CMS/SWGuideMuonIdRun2>. Accessed: 2018-05-13.
- 624 [11] CMS Collaboration, “Measurement of the Z/gamma\*+jets/photon+jets cross section  
625 ratio in pp collisions at  $\sqrt{s}=8$  TeV,” 2014. [CMS-PAS-SMP-14-005](#).