

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	15
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.	16
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.	18
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.	18
68	5.5	Results of study of the Z+jets to $\gamma$ +jets cross-section ratio for both data and MadGraph simulation.	23
70	5.6	Jet multiplicity and the associated $S_\gamma$ scale factor in the loose photon con- trol region before any corrections are applied.	24
72	5.7	$N_{jet}$ distribution in the tight $\mu\mu$ control region after $S_\gamma$ corrections.	24
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.	25
76	5.9	$N_b$ distribution before (left) and after (right) applying the $S_\gamma(N_j)$ scale factor.	25
78	5.10	$H_T$ and $m_{T2}$ distributions applying the $S_\gamma(N_j)$ scale factor.	26

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 ( $R_{norm}$ ) . . . . .	27
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 ( $R_{norm}$ ) . . . . .	27
83		
84		
85	5.13 Systematic uncertainty in the final prediction, as a function of the search bin, associated to the MC statistics. . . . .	29
86		
87	5.14 $Z \rightarrow \nu\bar{\nu}$ background prediction for all search bins, including the breakdown of the various uncertainties. . . . .	30
88		

# **Contents**

89	<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 . . . . .	16
103	5.0.2 Photon Purity and Fake Rate . . . . .	17
104	5.1 The $Z \rightarrow \mu^+ \mu^-$ Control Region . . . . .	20
105	5.1.1 Muon ID and Isolation . . . . .	20
106	5.1.2 Muon Selection in the Tight Control Region . . . . .	21
107	5.2 Analysis . . . . .	22
108	5.2.1 Shape Correction Using the $\gamma + \text{jets}$ Control Sample . . . . .	22
109	5.2.2 Normalization Correction Using the tight $Z \rightarrow \mu^+ \mu^-$ Control Sam- ple . . . . .	25
111	5.3 Results . . . . .	28
112	5.3.1 Systematics . . . . .	28
113	5.3.2 $Z \rightarrow \nu \bar{\nu}$ Estimation for the Search Bins . . . . .	29
114	<b>6 References</b>	31

<sup>115</sup> **Chapter 1**

<sup>116</sup> **Introduction**

<sup>117</sup> The work for this thesis was performed at CERN on CMS Experiment. CERN stands  
<sup>118</sup> for European Organization for Nuclear Research. It was founded in 1954 and is located  
<sup>119</sup> at the Franco-Swiss border near Geneva. At CERN, physicists and engineers are probing  
<sup>120</sup> the fundamental structure of the universe. They use the world's largest and most complex  
<sup>121</sup> scientific instruments to study the basic constituents of matter – the fundamental parti-  
<sup>122</sup> cles. The instruments used at CERN are purpose-built particle accelerators and detectors.  
<sup>123</sup> Accelerators boost beams of particles to high energies before the beams are made to col-  
<sup>124</sup> lide with each other or with stationary targets. Detectors observe and record the results  
<sup>125</sup> of these collisions. The accelerator at CERN is called the Large Hadron Collider (LHC),  
<sup>126</sup> the largest machine ever built by humans and it collides particles (protons) at close to the  
<sup>127</sup> speed of light. The process gives the physicists clues about how the particles interact, and  
<sup>128</sup> provides insights into the fundamental laws of nature. Seven experiments at the LHC use  
<sup>129</sup> detectors to analyze particles produced by proton-proton collisions. The biggest of these  
<sup>130</sup> experiments, ATLAS and CMS, use general-purpose detectors designed to study the fun-  
<sup>131</sup> damental nature of matter and fundamental forces and to look for new physics or evidence  
<sup>132</sup> of particles that are beyond the Standard Model. Having two independently designed de-  
<sup>133</sup> tectors is vital for cross-confirmation of any new discoveries made. The other two major  
<sup>134</sup> detectors ALICE and LHCb, respectively, study a state of matter that was present just  
<sup>135</sup> moments after the Big Bang and preponderance of matter than antimatter. Each experi-

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

138

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

140

141     ?? gives a brief motivation

142

143     ?? is dedicated to a study optimizing

144

145     ?? ptimated.

146

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

148

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

151

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

154     1. **Andrés Abreu** gave a talk “*Estimation of the Z Invisible Background for Searches*

155         *for Supersymmetry in the All-Hadronic Channel*” at “APS April 2018: American  
156         Physical Society April Meeting 2018, 14-17 Apr 2018”, Columbus, OH

157     2. **Andrés Abreu** gave a talk “*Phase-2 Pixel upgrade simulations*” at the “USLUA

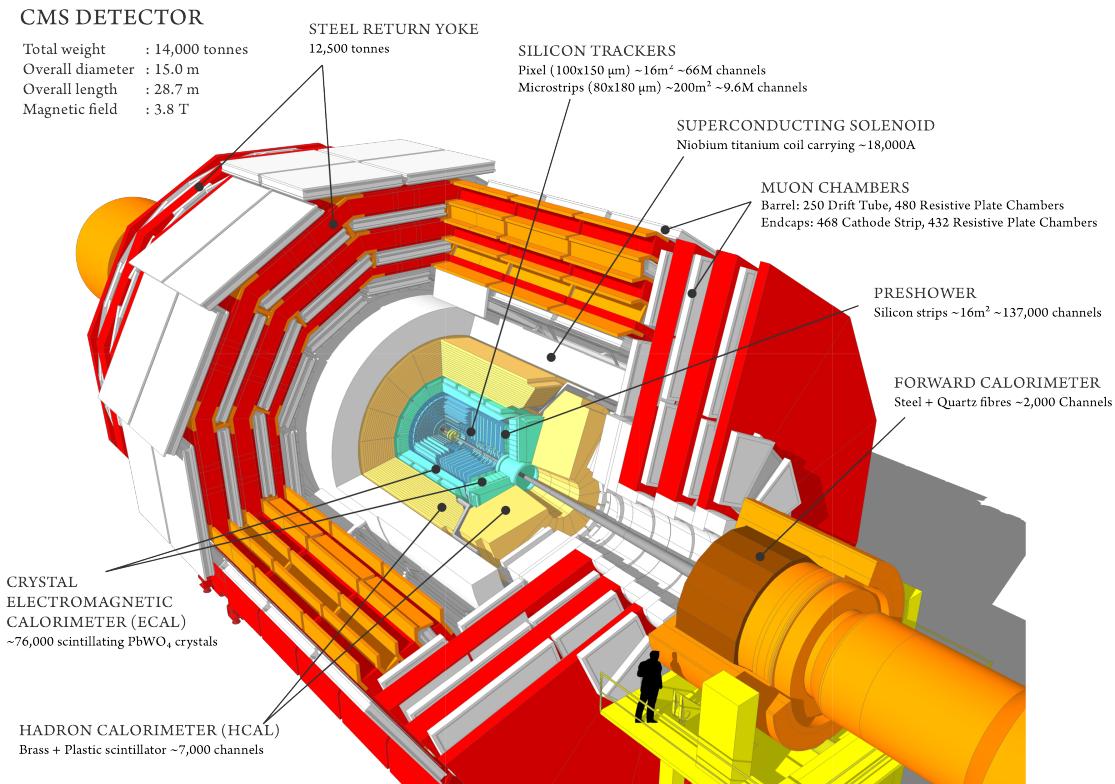
158         Annual meeting: 2017 US LHC Users Association Meeting, 1-3 Nov 2017”, Fer-  
159         milab, Batavia, IL

<sup>160</sup> **Chapter 2**

<sup>161</sup> **The CMS Experiment**

<sup>162</sup> The Compact Muon Solenoid (CMS) detector is a general purpose particle detector  
<sup>163</sup> designed to investigate various physical phenomena concerning the SM and beyond it,  
<sup>164</sup> such as Supersymmetry, Extra Dimensions and Dark Matter. As its name implies, the  
<sup>165</sup> detector is a solenoid which is constructed around a superconducting magnet capable of  
<sup>166</sup> producing a magnetic field of 3.8 T. The magnetic coil is 13m long with an inner diameter  
<sup>167</sup> of 6m, making it the largest superconducting magnet ever constructed. The CMS detector  
<sup>168</sup> itself is 21m long with a diameter of 15m and it has a weight of approximately 14,000  
<sup>169</sup> tons. The CMS experiment is one of the largest scientific collaborations in the history  
<sup>170</sup> of mankind with over 4,000 participants from 42 countries and 182 institutions. CMS is  
<sup>171</sup> located at one of these points and it essentially acts as a giant super highspeed camera  
<sup>172</sup> that makes 3D images of the collisions that are produced at a rate of 40 MHz (40 million  
<sup>173</sup> times per second). The detector has an onion-like structure to capture all the particles that  
<sup>174</sup> are produced in these high energy collisions most of them being unstable and decaying  
<sup>175</sup> further to stable particles that are detected. CMS detector was designed with the following  
<sup>176</sup> features (as shown in [Figure 2.1](#)) :

- <sup>177</sup> 1. A **magnet** with large bending power and high performance muon detector for good  
<sup>178</sup> muon identification and momentum resolution over a wide range of momenta and  
<sup>179</sup> angles.

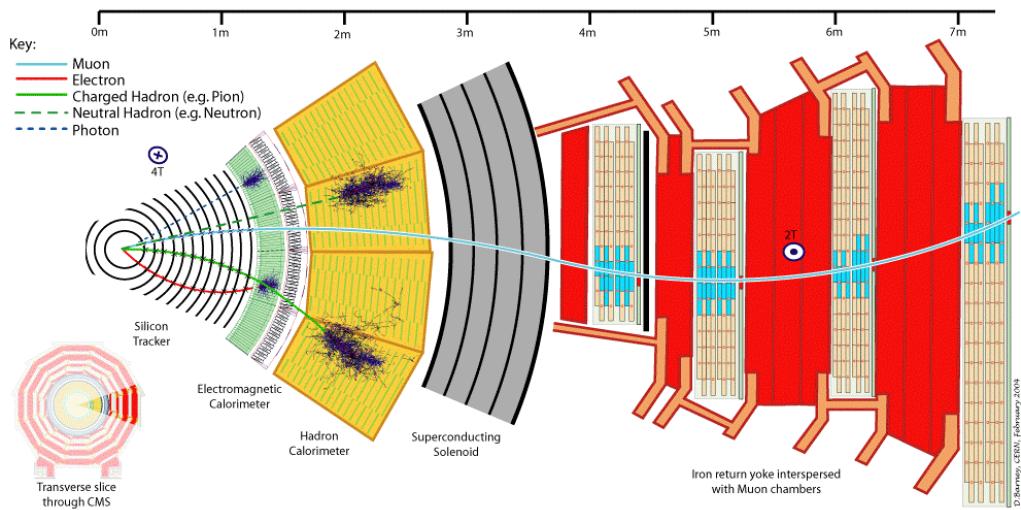


**Figure 2.1:** CMS Detector

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

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

the particle as they move. The Calorimeters measure the energy of particles that pass through them by absorbing their energy with the intent of stopping them. The particle identification detectors work by detecting radiation emitted by charged particles and using this information they can measure the speed, momentum, and mass of a particle. After the information is put together to make the “snapshot” of the collision one looks for results 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

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

# <sup>213</sup> Chapter 3

## <sup>214</sup> Data Collection and Data Quality

### <sup>215</sup> Monitoring

#### <sup>216</sup> 3.1 What is Data Collection for CMS?

<sup>217</sup> During data taking there are millions of collisions occurring in the center of the detec-  
<sup>218</sup> tor every second. The data per event is around one million bytes (1 MB), that is produced  
<sup>219</sup> at a rate of about 600 million events per second [1], that's about 600 MB/s. Keeping  
<sup>220</sup> in mind that only certain events are considered "interesting" for analysis, the task of de-  
<sup>221</sup> ciding what events to consider out of all the data collected is a two-stage process. First,  
<sup>222</sup> the events are filtered down to 100 thousand events per second for digital reconstruction  
<sup>223</sup> and then more specialized algorithms filter the data even more to around 100 200 events  
<sup>224</sup> per second that are found interesting. For CMS there is a Data Acquisition System that  
<sup>225</sup> records the raw data to what's called a High-Level Trigger farm which is a room full  
<sup>226</sup> of servers that are dedicated to processing and classify this raw data quickly. The data  
<sup>227</sup> then gets sent to what's known as the Tier-0 farm where the full processing and the first  
<sup>228</sup> reconstruction of the data are done. [2]

## <sup>229</sup> 3.2 What is Data Quality Monitoring?

<sup>230</sup> To operate a sophisticated and complex apparatus as CMS, a quick online feedback on  
<sup>231</sup> the quality of the data recorded is needed to avoid taking low quality data and to guarantee  
<sup>232</sup> a good baseline for the offline analysis. Collecting a good data sets from the collisions  
<sup>233</sup> is an important step towards search for new physics as deluge of new data poses an extra  
<sup>234</sup> challenge of processing and storage. This all makes it all the more important to design  
<sup>235</sup> algorithms and special software to control the quality of the data. This is where the Data  
<sup>236</sup> Quality Monitoring (DQM) plays a critical in the maintainability of the experiment, the  
<sup>237</sup> operation efficiency and performs a reliable data certification. The high-level goal of  
<sup>238</sup> the system is to discover and pinpoint errors, problems occurring in detector hardware  
<sup>239</sup> or reconstruction software, early, with sufficient accuracy and clarity to maintain good  
<sup>240</sup> detector and operation efficiency. The DQM workflow consists of 2 types: **Online** and  
<sup>241</sup> **Offline**.

<sup>242</sup> The **Online** DQM consists of receiving data taken from the event and trigger his-  
<sup>243</sup> tograms to produce results in the form of monitoring elements like histogram references  
<sup>244</sup> and quality reports. This live monitoring of each detector's status during data taking gives  
<sup>245</sup> the online crew the possibility to identify problems with extremely low latency, mini-  
<sup>246</sup> mizing the amount of data that would otherwise be unsuitable for physics analysis. The  
<sup>247</sup> scrutiny of the Online DQM is a 24/7 job that consists of people or shifters that work at the  
<sup>248</sup> CMS control center constantly monitoring the hundreds of different plots and histograms  
<sup>249</sup> produced by the DQM software. This consumes a lot of manpower and is strenuous work.

<sup>250</sup> The **Offline** DQM is more focused on the full statistics over the entire run of the  
<sup>251</sup> experiment and works more on the data certification. In the offline environment, the  
<sup>252</sup> system is used to review the results of the final data reconstruction on a run-by-run basis,  
<sup>253</sup> serving as the basis for certified data used across the CMS collaboration in all physics  
<sup>254</sup> analyses. In addition, the DQM framework is an integral part of the prompt calibration  
<sup>255</sup> loop. This is a specialized workflow run before the data are reconstructed to compute and  
<sup>256</sup> validate the most up-to-date set of conditions and calibrations subsequently used during

257 the prompt reconstruction.

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

# <sup>261</sup> Chapter 4

## <sup>262</sup> What is Machine Learning?

<sup>263</sup> Machine Learning (ML) can be defined as an application of Artificial Intelligence that  
<sup>264</sup> permits the computer system to learn without being told explicitly. In ML a computer  
<sup>265</sup> program is said to learn from experience E with respect to some class of tasks T and  
<sup>266</sup> performance measure P, if its performance at tasks in T, as measured by P, improves  
<sup>267</sup> with experience E [3]. ML has made tremendous strides in the past decades and has  
<sup>268</sup> become very popular recently due to its multifaceted applications. It is being used on  
<sup>269</sup> social media, marketing, and in the scientific community as well. Some examples of  
<sup>270</sup> ML applications are: the algorithms used on application in smartphones to detect human  
<sup>271</sup> faces, self-driving cars, computer games, stock prediction, and voice recognition. An  
<sup>272</sup> interesting characteristic of ML algorithms is that the more data one inputs the better is  
<sup>273</sup> the performance. The ML application has a very wide spectrum covering almost every  
<sup>274</sup> aspect of human endeavor that involves a lot of data. Scientific analysis today generates  
<sup>275</sup> enormous data and is hence a perfect use case to apply ML techniques. In this work  
<sup>276</sup> we use enhanced ML techniques based on progress in the recent past.

<sup>277</sup> In general, there are two main categories to classify machine learning problems: **Su-**  
<sup>278</sup> **pervised Learning** (SL) and **Unsupervised Learning** (UL). SL is the most used ML  
<sup>279</sup> approach and has proven to be very effective for a wide variety of problems. Examples  
<sup>280</sup> of common SL problems are: spam filters, predicting housing prices, identifying a ma-  
<sup>281</sup> lignant or benign tumor, etc. These types of problems are characterized by providing a

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

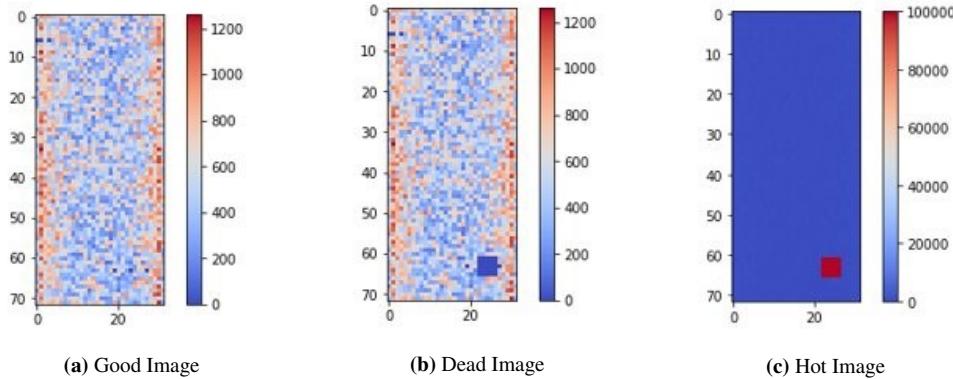
UL algorithms are different in the sense that they do not have the “right answers” given to the machine. Instead, UL algorithms are used for finding patterns and make clusters from the given data. That is what also forms the basis of a search engine (e.g. Google news). Clicking on a link to a news article,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

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

## 4.1 Developing the Algorithm

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

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



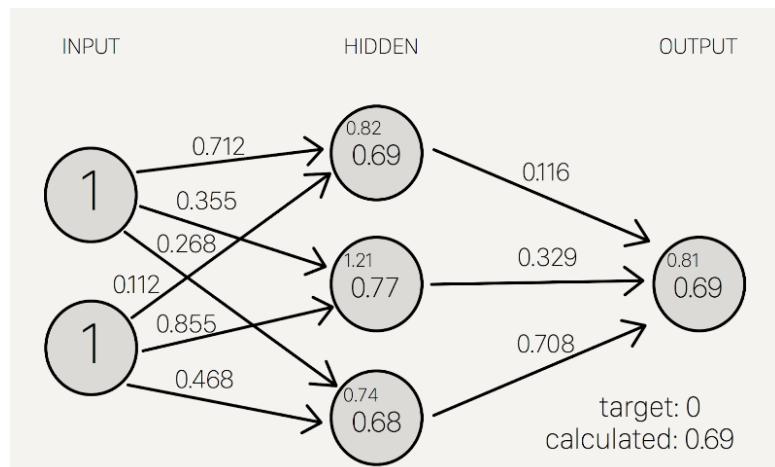
**Figure 4.1:** Occupancy maps with 5x5 affected regions

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

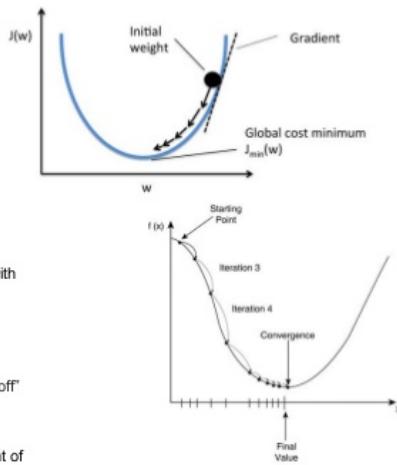
$\eta$ : Learning Rate (typically 0.01 to 0.001)

$\eta$  :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

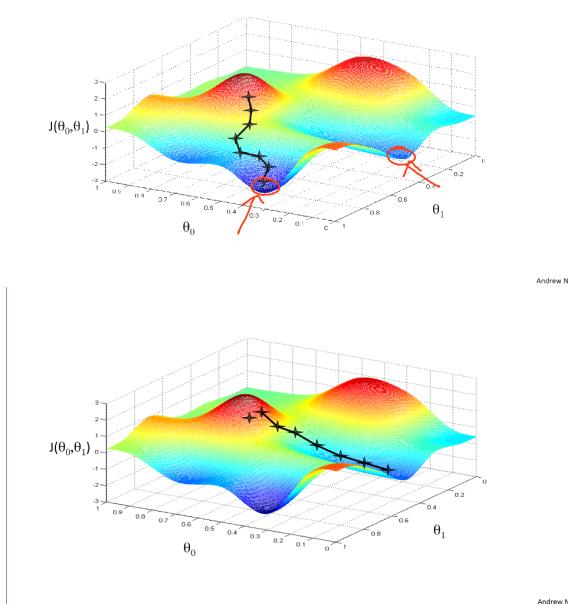
We seek to update the weights and biases by a value indicating how “off” they were from their target.

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



**Figure 4.3:** Gradient Descent algorithm

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

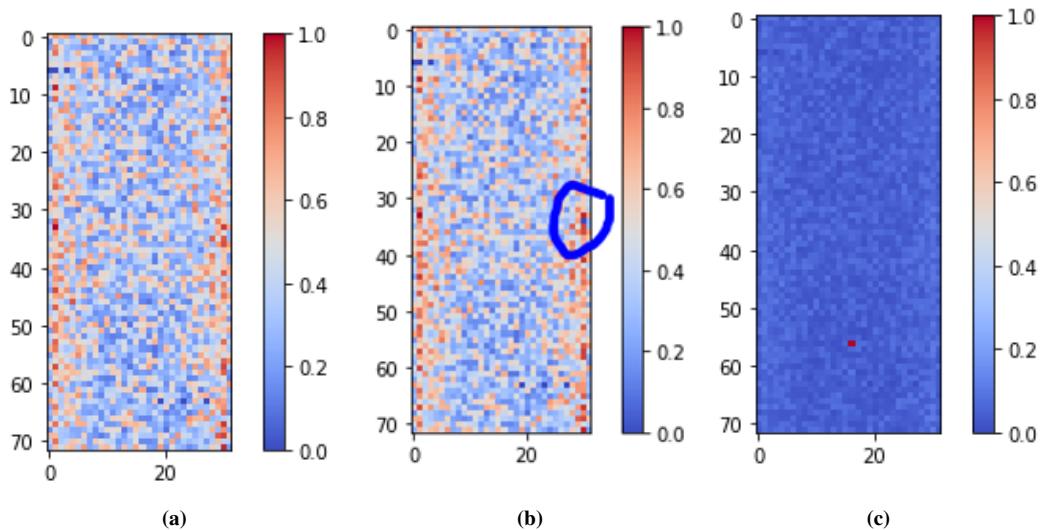


**Figure 4.4:** Loss Function surface

355 **Chapter 5**

356 **Results**

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



364 **Figure 5.1:** Occupancy Maps with 1x1 bad regions. A) Good image B) Dead image  
d vice-versa. ?? shows the cut values that are applied to photons that are found within  
365 both the ECAL barrel and endcap range. The associated values to the efficiency and the

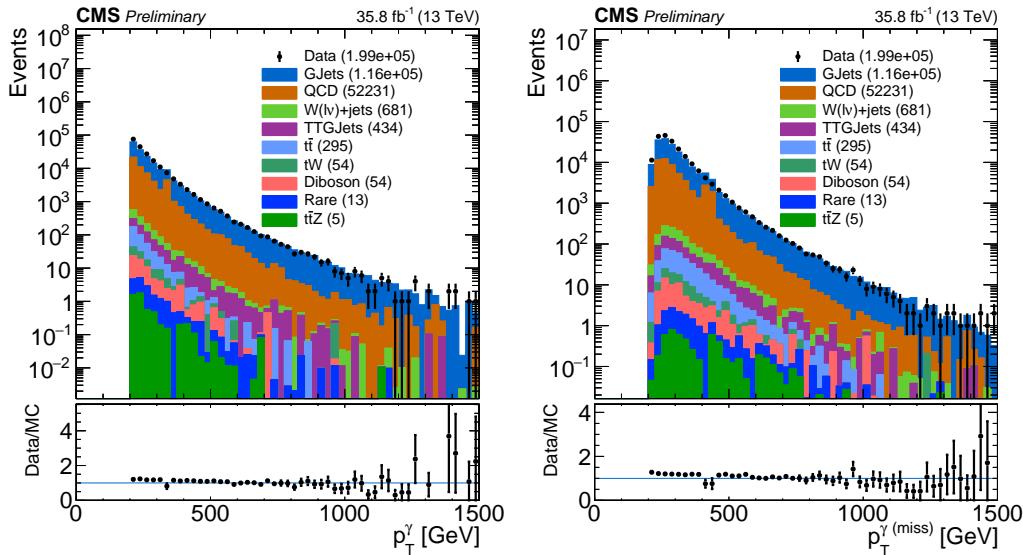
366 background rejection rate are shown for each of the three different photon ID selections.

367

368 In order to obtain the high efficiency and background rejection rae contributions.

369 **5.0.1 Photon Selection**

370 The event selection process for the  $\gamma$ +jets control region starts with photon candidates  
371 that have a  $p_T > 200$  GeV and are within the acceptance range of the CMS ECAL (given  
372 by  $|\eta| < 1.4442$  for the barrel and  $1.566 < |\eta| < 2.5$  for the endcaps). The photons are  
373 subjected to pass the loose ID/isolation cuts described in ?? in order to remove  $\sim 85\%$   
374 of the background processes and obtain a prompt phototly improves the prompt photon  
375 selection by removing many of the events in the simulated samples where a lepton gets  
376 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.

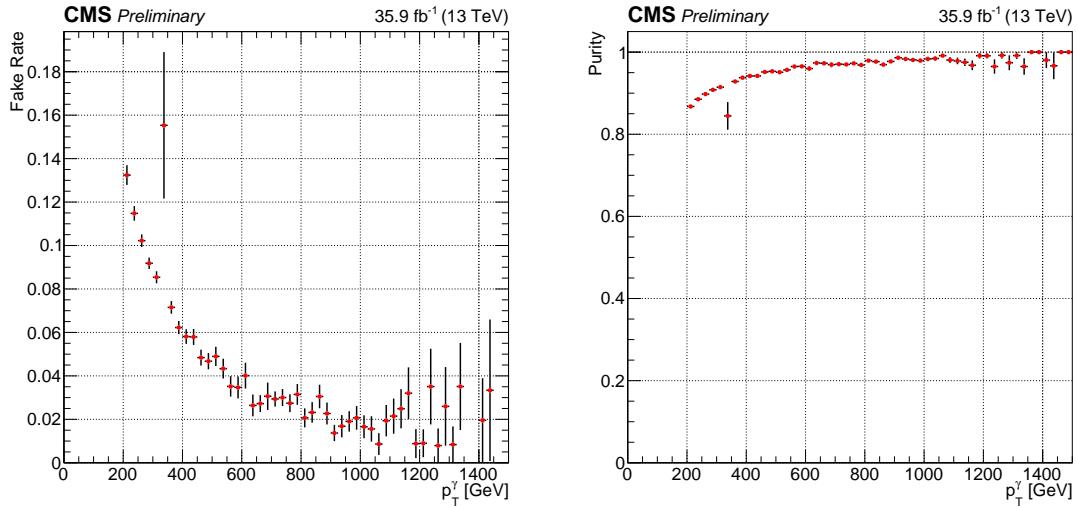
377 To further emulate the  $Z \rightarrow \nu\bar{\nu} + \text{jets}$  background, a variable in which the photons are  
378 treated as  $p_T^{\text{miss}}$  is defined. We call this variable  $p_T^{\gamma(\text{miss})}$  and we obtain it by adding the  
379  $p_T^\gamma$  for every event to the total  $p_T^{\text{miss}}$  in the event. Both the  $p_T^\gamma$  and the resulting  $p_T^{\gamma(\text{miss})}$   
380 distributions are shown in Figure 5.2 as data/MC comparison plots, where the simulated  
381 backgrounds are stacked in order of ascending contribution.

382

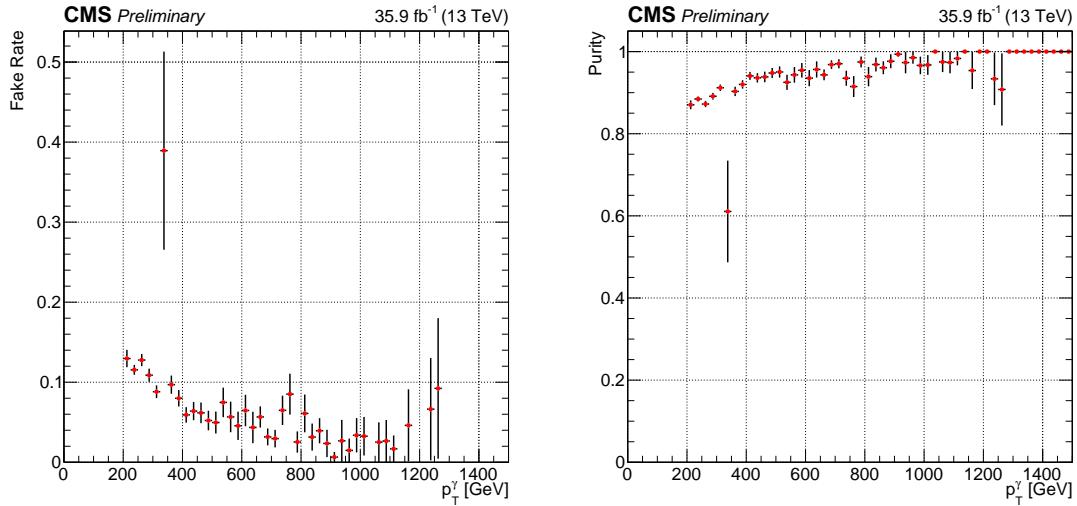
383 The main contributions from simulation arise from the  $\gamma$ +jets, QCD and to a lesser  
384 extent,  $t\bar{t}\gamma$ . Other non-dominant backgrounds in the control region include contributions  
385 from  $W(l\nu)$ +jets,  $t\bar{t}$ , Diboson,  $tW$ ,  $t\bar{t}Z$  and rare processes. Most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

### 394 5.0.2 Photon Purity and Fake Rate

395 Three different types of photons make up the  $\gamma$ +jets CS: prompt photons, produced  
396 either directly or through fragmentation, and fake photons. Prompt photons are defined  
397 as photons which are formed shortly after the proton-proton collision (i.e. before the pro-  
398 duced quarks and gluons have had enough time to form hadrons). Two types of photons  
399 fit in this category. The first type, which we designate as direct photons, are photons that  
400 are produced directly from the proton-proton interaction [5]. A secondary type of prompt  
401 photon, that is virtually indistinguishable from the direct photons at the detector level,  
402 originates from the decay of  $\pi^0$  mesons and are called fragmentation photons. The final  
403 type of photon found in the CS corresponds to fake (or non-prompt) photons. The fake  
404 photon contribution typically arises from leptons (mostly electrons) whose tracks are not  
405 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.

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

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

417

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

428

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

431

$$432 p_\gamma = \frac{\text{prompt}}{\text{prompt} + \text{fake}},$$

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

439

$$440 \quad f = \frac{fake}{prompt+fake} ,$$

441     Figure 5.3 shows the purity and fakerate for photons that pass the loose ID/selection,  
 442    have a  $p_T > 200$  GeV and are within the ECAL acceptance range. A sample is obtained  
 443    in which 77% of the photons are direct, 12% are fragmentation and 11% are fakes. This  
 444    implies an average purity of  $\sim 89\%$  for this sample, well within the value that is expected.  
 445    Figure 5.4 shows the same ratios for the loose  $\gamma +$  jets control region described in subsection  
 446    5.0.1. Although the amount of statistics has decreased due to the additional cuts, a  
 447    similar trend can be observed.

## 448    5.1 The $Z \rightarrow \mu^+ \mu^-$ Control Region

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

### 456    5.1.1 Muon ID and Isolation

457    The muons are selected using the “medium muon” selection [6], per the recommen-  
 458    dation of the muon POG. The muon candidates in this selection satisfy  $p_T > 10$  GeV and  
 459     $|\eta| < 2.4$ . Other additional cuts are applied to aid in the muon candidate selection, such as  
 460    an impact parameter cut. Muons are also subjected to a PF relative-isolation (also referred  
 461    to as mini-isolation) in which the cone size is inversely proportional to the muon  $p_T$ . This  
 462    requirement enforces the  $p_T$  within the isolation cone to be at most 20% of the muon  $p_T$

463 in order to eliminate events with an isolated muon. Details of the medium photon selec-  
 464 tion are included in [Table 5.1](#) and [Table 5.2](#), while details of the impact parameter cut are  
 465 summarized in [Table 5.3](#).

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 5.1:** 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 5.2:** Muon Medium ID HIP Safe Good Global Muon

Muon Impact Parameter	
d0 <	0.2
dz <	0.5

**Table 5.3:** Additional Impact Parameter cut on Muons

### 466 5.1.2 Muon Selection in the Tight Control Region

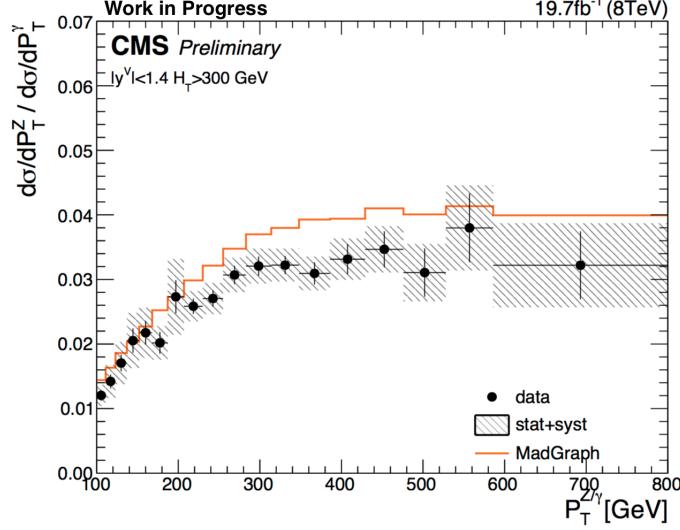
467 Events are selected from data samples that contain exactly two oppositely charged  
 468 muons ( $\mu^+ \mu^-$ ), which fall within the invariant mass  $81 < m_{ll} < 101$  GeV window for  
 469 the Z boson. Additional cuts for the tight muon selection include baseline requirements  
 470 such as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

## 477 5.2 Analysis

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

### 485 5.2.1 Shape Correction Using the $\gamma$ + jets Control Sample

486 In this section the validation of the  $\gamma$ +jets simulation is discussed in terms of the  
 487 shape of the loose photon control region. As it was shown in [subsection 5.0.2](#), this con-  
 488 trol region has high purity for  $\gamma$ +jets events, particularly in regions of high  $p_T$  ( $\gtrsim 300$ ).  
 489 In order to apply this correction factor it is assumed that the shape differences between  
 490 data and simulation are similar between  $Z \rightarrow \nu\bar{\nu}$  and  $\gamma$ +jets events. This assumption is  
 491 validated in studies which compare the cross-section ratio of  $Z$ +jets to  $\gamma$ +jets events [7].  
 492 [Figure 5.5](#) shows the results of this study, conducted in 2014, for both data and Mad-  
 493 Graph simulation with an integrated luminosity of  $19.7\text{fb}^{-1}$  and a center-of-mass energy  
 494 of 8 TeV. It can be seen that for values of  $p_T^{Z/\gamma} \gtrsim 300$  GeV, the ratio of the cross-section of  
 495 both processes becomes nearly constant. It is then a matter of applying a factor to account  
 496 for the difference in the amount of events between the  $Z$ +jets and  $\gamma$ +jets events in order  
 497 to obtain the total amount of  $Z \rightarrow \nu\bar{\nu}$ +jets events. This factor is obtained from the tight  
 498  $\mu\mu$  control region, as shown in [subsection 5.2.2](#).

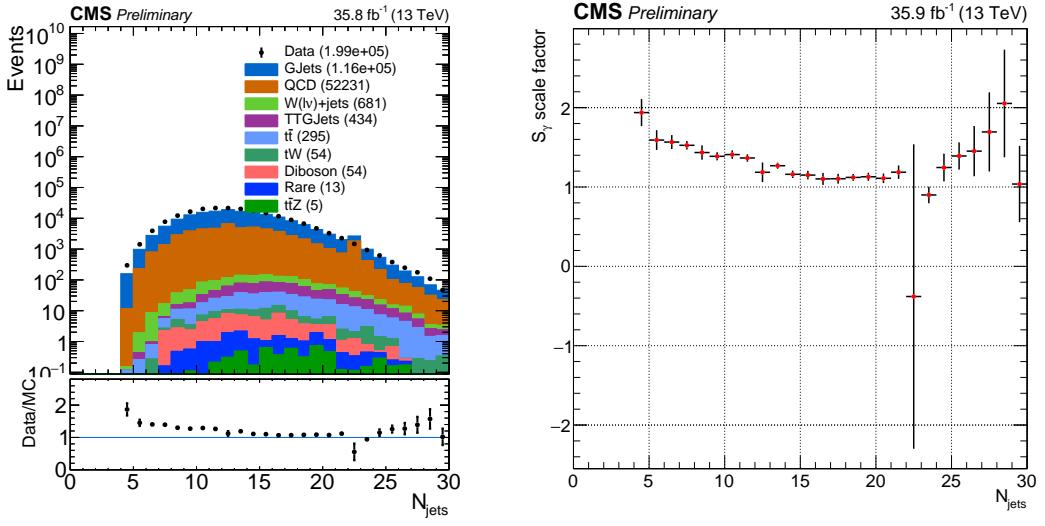


**Figure 5.5:** Results of study of the  $Z$ +jets to  $\gamma$ +jets cross-section ratio for both data and MadGraph simulation. For high values of the vector boson transverse momentum, the ratio between these processes is observed to be nearly constant.

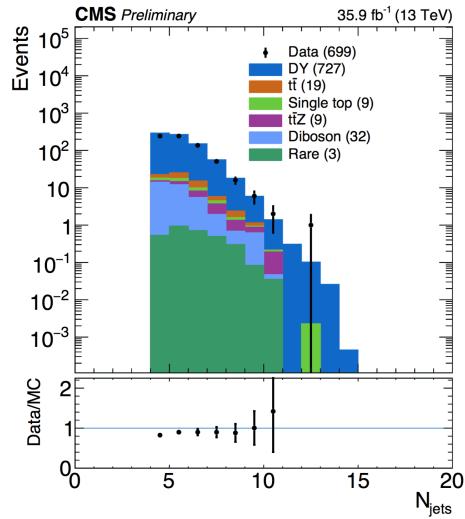
499 In order to obtain the shape corrections, the ratio between data and simulation of the  
 500 jet multiplicity distribution is used (Figure 5.6). This is due to it exhibiting the highest  
 501 difference between the observed data and MC. The re-weight for the  $\gamma$ +jets simulation  
 502 sample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

$$505 \quad S_\gamma^i = \frac{\text{Data}^i - \text{MC}_{\text{other}}^i}{\text{MC}_{\gamma+\text{jets}}^i},$$

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

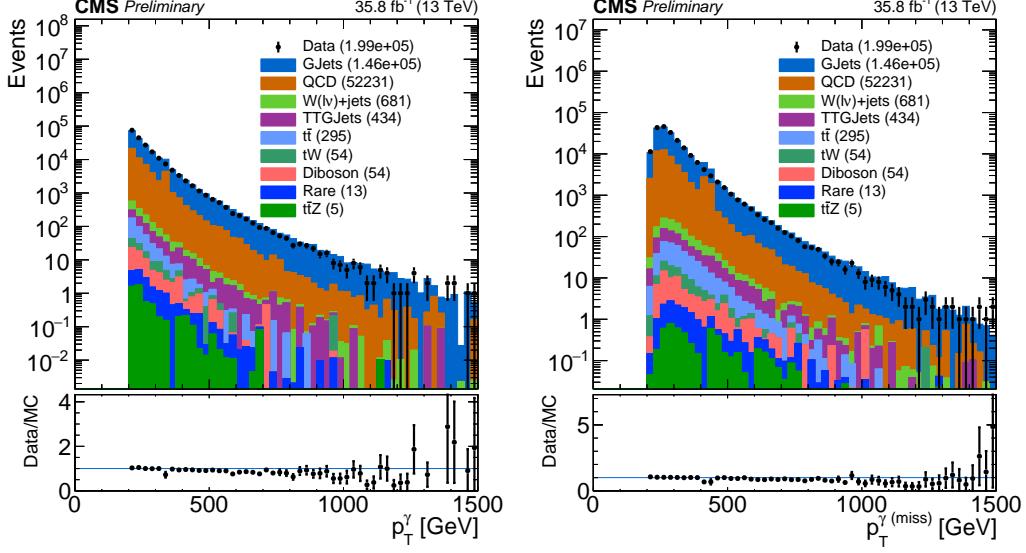


**Figure 5.6:** Jet multiplicity and the associated  $S_\gamma$  scale factor in the loose photon control region before any corrections are applied.

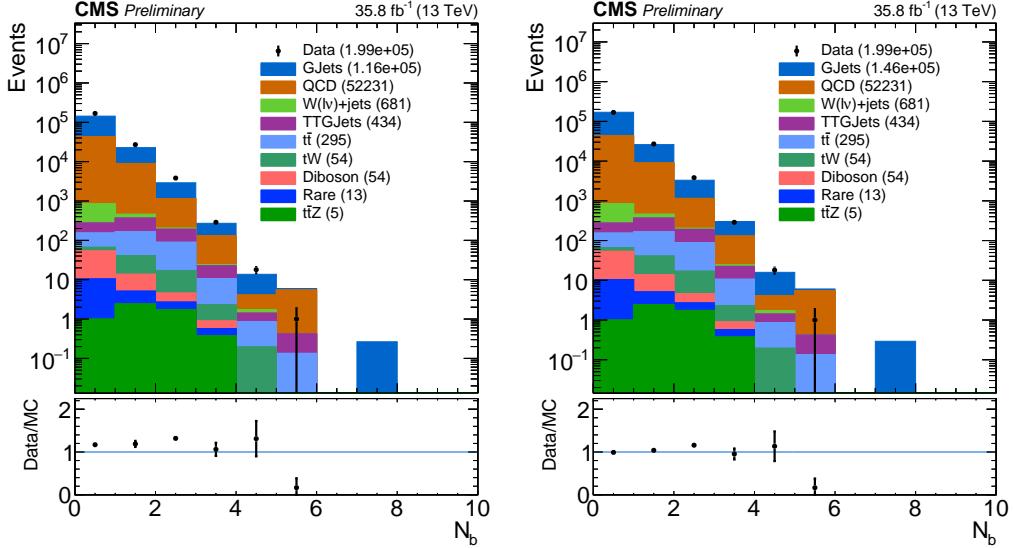


**Figure 5.7:**  $N_{jet}$  distribution in the tight  $\mu\mu$  control region after  $S_\gamma$  corrections.

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



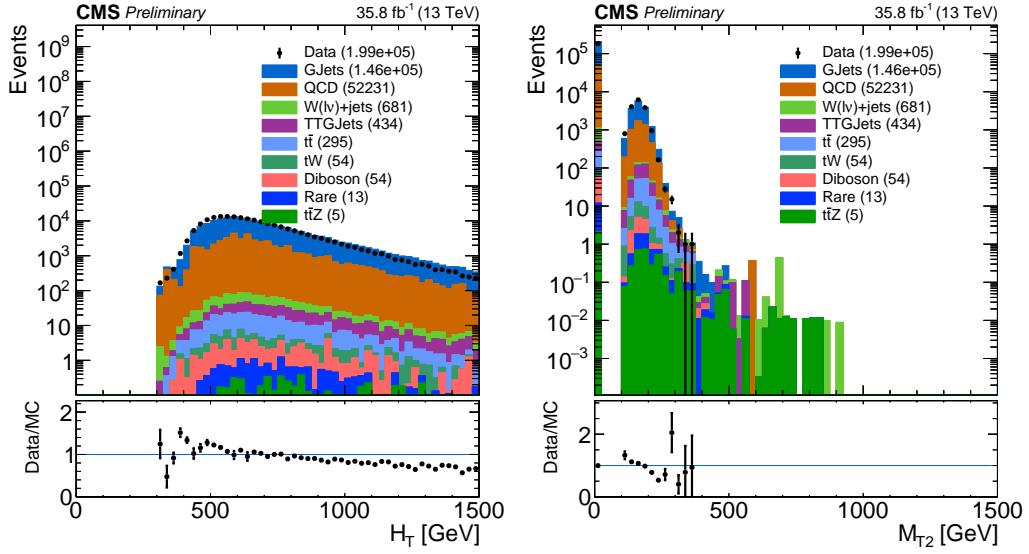
**Figure 5.8:**  $p_T^\gamma$  (left) and  $p_T^{\gamma(\text{miss})}$  (right) distributions after applying the  $S_\gamma(N_j)$  scale factor. Comparing to Figure 5.2, an improvement in the agreement between data/MC can be observed.



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

518    **5.2.2 Normalization Correction Using the tight  $Z \rightarrow \mu^+ \mu^-$  Control**  
 519    **Sample**

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

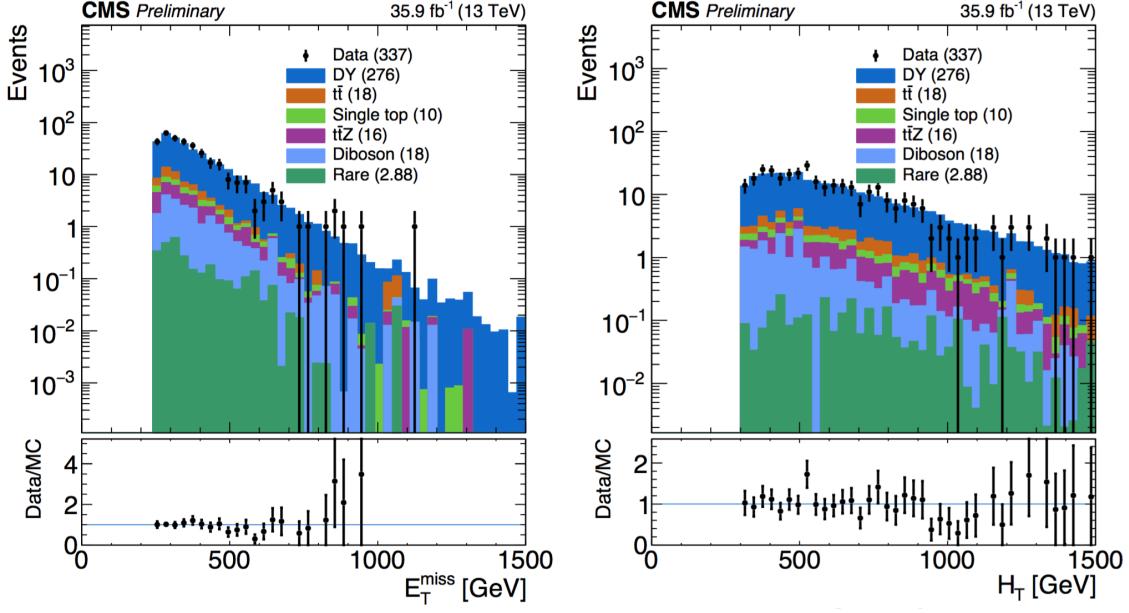


**Figure 5.10:**  $H_T$  and  $m_{T2}$  distributions applying the  $S_\gamma(N_j)$  scale factor.

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

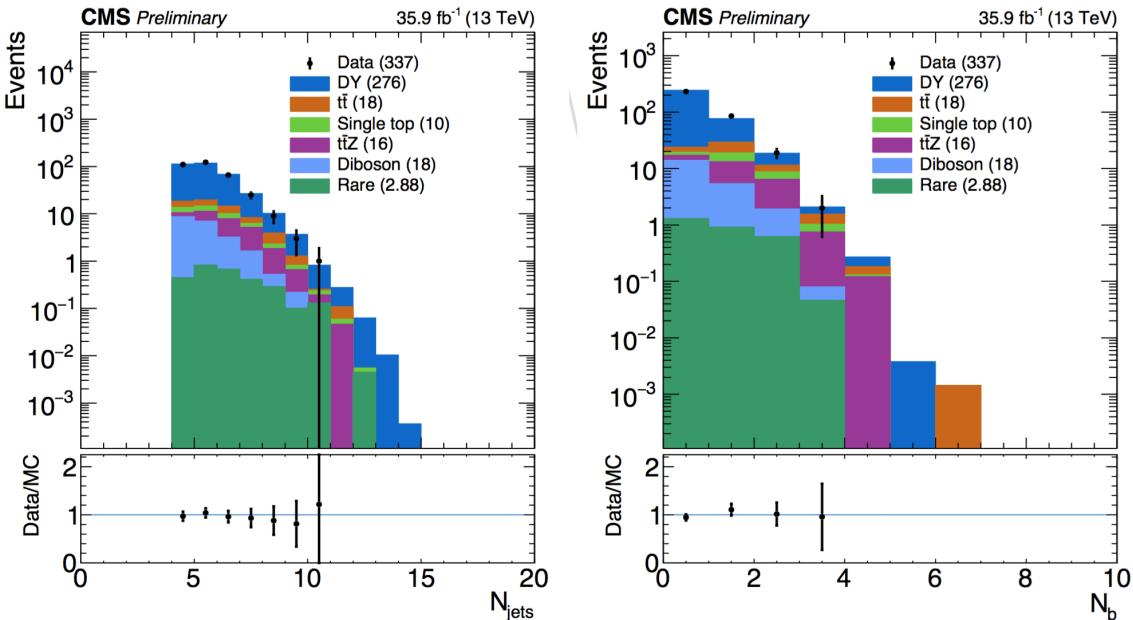
529 
$$R_{norm} = 1.070 \pm 0.085,$$

530 where the uncertainty includes only the associated statistical uncertainties on data and  
 531 simulation. This uncertainty is found to be propagated to the final background prediction,  
 532 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 ( $R_{\text{norm}}$ ).

534 Data/MC comparisons are shown in [Figure 5.11](#) and [Figure 5.12](#) after applying  $R_{\text{norm}}$   
 535 for several distributions in the study. With this final global scale factor all the required  
 536 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 ( $R_{\text{norm}}$ ).

## 537 5.3 Results

538 In this section the results for the final estimation of the of the  $Z \rightarrow \nu\bar{\nu}$  are presented.

539 The current study includes preliminary results using only data obtained at the CMS detec-

540 tor during 2016. The results for this study are intended to confirm the assumption that the

541 additional  $\gamma + \text{jets}$  control region introduced in this analysis reduce the overall uncertain-

542 ties obtained in the 2016 analyses (described in ??). Furthermore, this study is intended

543 as a benchmark for future analyses of the SUSY stop group based in Fermilab and will be

544 the method used for the 2017 CMS data.

### 545 5.3.1 Systematics

546 Two categories of uncertainties for the  $Z \rightarrow \nu\bar{\nu}$  prediction are considered: uncertain-

547 ties that are associated to the use of MC simulation and the uncertainties specifically

548 associated to the background prediction method. Several sources are acknowledged in the

549 first category mentioned such as PDF and renormalization/factorization scale choices, jet

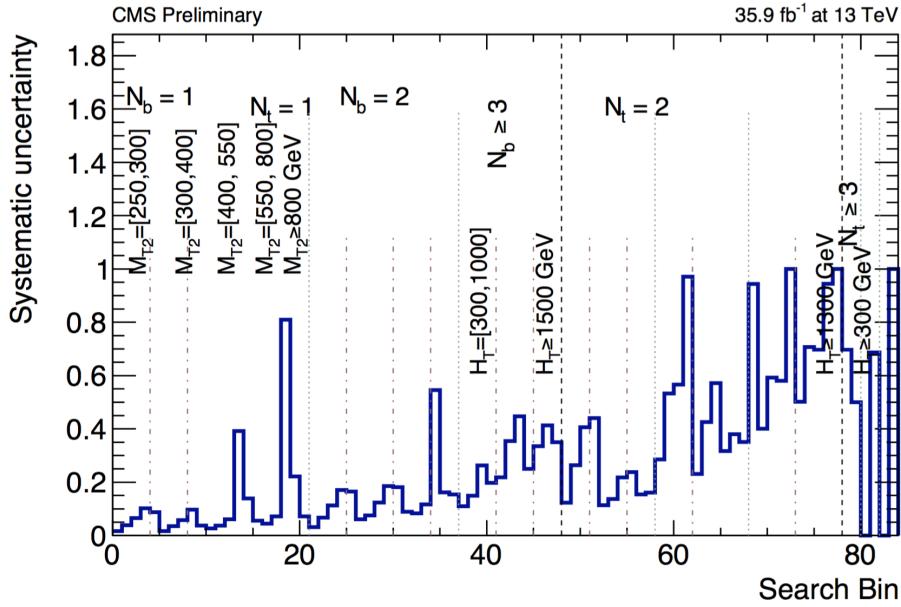
550 and  $p_T^{miss}$  energy scale uncertainties b-tag scale factor uncertainties, and trigger efficiency

551 uncertainties. Given that the simulation sample is normalized to data in the tight control

552 region, uncertainties associated with the luminosity and cross-section are excluded. In

553 addition, the overall  $Z \rightarrow \nu\bar{\nu}$  statistical uncertainty from MC simulation is also taken into

554 account.



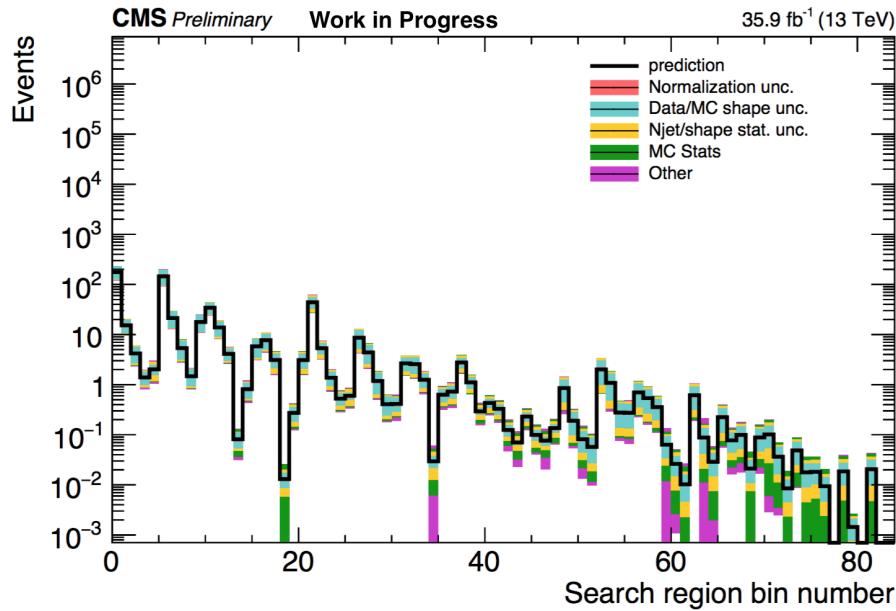
**Figure 5.13:** Systematic uncertainty in the final prediction, as a function of the search bin, associated to the MC statistics.

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

### 566     5.3.2 $Z \rightarrow \nu\bar{\nu}$ Estimation for the Search Bins

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

<sup>572</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>573</sup> **Chapter 6**

<sup>574</sup> **References**

<sup>575</sup> [1] CERN, “Processing what to record?,” 2018.

<sup>576</sup> [2] CMS, “The CMS Computing Project,” tech. rep., CERN, 2005.

<sup>577</sup> [3] Coursera, “Machine learning,” 2018.

<sup>578</sup> [4] A. S. Walia, “Types of optimization algorithms used in neural networks and ways to  
<sup>579</sup> optimize gradient descent,” 2018.

<sup>580</sup> [5] M. Klasen, C. Klein-Bosing, and H. Poppenborg, “Prompt photon produc-  
<sup>581</sup> tion and photon-jet correlations at the LHC,” *JHEP*, vol. 03, p. 081, 2018.  
<sup>582</sup> doi:10.1007/JHEP03(2018)081.

<sup>583</sup> [6] S. Folgueras, “Baseline muon selections for Run-II.” <https://twiki.cern.ch/twiki/bin/view/CMS/SWGuideMuonIdRun2>. Accessed: 2018-05-13.

<sup>585</sup> [7] CMS Collaboration, “Measurement of the Z/gamma\*+jets/photon+jets cross section  
<sup>586</sup> ratio in pp collisions at sqrt(s)=8 TeV,” 2014. [CMS-PAS-SMP-14-005](#).