



FACULTAD DE CIENCIAS  
UNIVERSIDAD DE CANTABRIA

---

**Search for dark matter production in  
association with top quarks in the  
dilepton final state at  $\sqrt{s} = 13$  TeV**

---

A thesis submitted in fulfillment of the requirements for the  
**Degree of Doctor of Philosophy**

---

Written by  
**Cédric Prieëls**

Under the supervision of  
**Jónatan Piedra Gómez**  
**Pablo Martínez Ruíz del Árbol**

Santander, July 2021





FACULTAD DE CIENCIAS  
UNIVERSIDAD DE CANTABRIA

---

**Búsqueda de materia oscura en  
asociación con quarks top en el estado  
final dileptónico a  $\sqrt{s} = 13$  TeV**

---

Memoria para optar al

**Grado de doctor**

---

Escrita por

**Cédric Prieëls**

Bajo la supervisión de  
**Jónatan Piedra Gómez**  
**Pablo Martínez Ruíz del Árbol**

Santander, Julio 2021



---

# Abstract

A search for dark matter particles produced in association with one or two top quarks is presented in this work, by studying in particular the dilepton decay channel of both production modes. This analysis was done by considering the full  $(137.1 \pm 2.0) \text{ fb}^{-1}$  of proton-proton collisions data collected by the Compact Muon Solenoid (CMS) detector during the Run II of operation of the Large Hadron Collider (LHC), at a center of mass energy of  $\sqrt{s} = 13 \text{ TeV}$ . This is the first time that such a search combining the  $t/\bar{t}$  and  $t\bar{t}+\text{DM}$  models is performed in the dilepton final state.

This search relies on the implementation of advanced Machine Learning techniques, allowing us to isolate the different signals studied with respect to the backgrounds, mainly estimated from Monte-Carlo and mostly made out of Standard Model  $t\bar{t}$  and single top processes, and whose discriminating power was combined into a single output variable used to perform a shape analysis thanks to the definition and previous optimization of a neural network.

Upper limits on the production cross section of different signal models, both scalar and pseudoscalar, were then extracted. At the end of the day, no evidence for the existence of dark matter has been found, but upper limits on the signal strength have been obtained by considering different production models and channels. This analysis allowed us in this sense to achieve an expected (observed) exclusion for both scalar and pseudoscalar mediators up to 200 (XXX) GeV, improving by a factor of more than 2 the previous results obtained in 2016 for the  $t\bar{t}+\text{DM}$  model alone.



---

## Acknowledgments

First of all, I would like to thank the group of high energy physics of the Instituto de Física de Cantabria (IFCA) in general for the opportunity they gave me by offering me such an interesting PhD position within the group. This was a great opportunity for me to start working as a CERN collaborator, which had been a dream of mine for many years. All these years spent at IFCA were for sure extremely interesting, brought me a lot in both the personal and professional points of view, and something I will definitely benefit from in the future.

I would also like to personally thank Jónatan Piedra for supervising this thesis, but also for all your help and patience during all these years, especially at the beginning. You taught me a lot and I cannot count the number of bugs and issues you helped me solve, so I am really thankful that you kept your door open at all times anyway.

Pablo Martínez, I am grateful for the great work you did as the co-director of this thesis, as I am sure this work would not have been possible without your help and dedication to this analysis. You also became a friend over the years and I will never be able to thank you enough for your help, especially in this last year which has undoubtedly been quite challenging for me.

I would also like to thank all the people who helped me at some point: Alicia, Rocío, Pupi, Chus and Celso, always there to help me fight the administrative related issues I faced. This work would not have been possible either without the help of all the people involved in this particular analysis, such as Deborah, Nicole, Alexander and Dominic. Thank you as well to all the people involved in the development of the Latino framework, which quickly became a great tool to work with. I am also deeply grateful to all the students and seniors from the University of Oviedo for all these interesting meetings, making the Wednesday mornings special each week.

Thank you also to my past and present (PhD) coworkers who fought alongside me: Pedro, Nicolò, Celia, Clara, Andrea, Pablo, Barbara and Juan, you made the students office a better place to work in and I wish you all the best, even though I am sure you will all achieve great things in the future.

Of course, I would especially like to thank my parents Philippe and Chantal, and my brothers Antoine and Lucas for giving me the opportunity to come and live in Spain for so long. I never could have reached this far without your presence and support that I felt daily, even more than a thousand kilometers away. Thank you to Inés as well, as I would never have started this thesis if I hadn't met you. Finally, I would like to especially thank Aline for her incredible support over the past few months, which helped me a lot and for which I am extremely grateful. Even though one chapter now ends, I cannot wait to get started on the next one, and I sure hope that this is only the beginning with you.



---

## Acronyms used

<b>ADMX</b>	Axion Dark Matter Experiment	<b>DY</b>	Drell-Yan
<b>ALICE</b>	A Large Ion Collider Experiment	<b>ECAL</b>	Electromagnetic Calorimeter
<b>AMS</b>	Alpha Magnetic Spectrometer	<b>EDM</b>	Event Data Model
<b>ANN</b>	Analysis Neural Network	<b>EFT</b>	Effective Field Theory
<b>AOD</b>	Analysis Object Data	<b>EWK</b>	Electroweak
<b>ATLAS</b>	A Toroidal LHC ApparatuS	<b>FR</b>	Fake Rate
<b>BDT</b>	Boosted Decision Tree	<b>FSR</b>	Final State Radiation
<b>BR</b>	Branching Ratio	<b>GEM</b>	Gas Electron Multiplier
<b>BSM</b>	Beyond the Standard Model	<b>GSF</b>	Gaussian Sum Filter
<b>BW</b>	Breit-Wigner	<b>HCAL</b>	Hadronic Calorimeter
<b>CAST</b>	CERN Axion Solar Telescope	<b>HLT</b>	High-Level Trigger
<b>CERN</b>	European Council for Nuclear Research	<b>HO</b>	Hadron Outer
<b>CL</b>	Confidence Level	<b>IACT</b>	Imaging Atmospheric Cherenkov Telescopes
<b>CMB</b>	Cosmic Microwave Background	<b>IAXO</b>	International AXion Observatory
<b>CMS</b>	Compact Muon Solenoid	<b>IFCA</b>	Instituto de Física de Cantabria
<b>CSC</b>	Cathode Strip Chamber	<b>ISR</b>	Initial State Radiation
<b>CR</b>	Control Region	<b>JEC</b>	Jet Energy Correction
<b>CSV</b>	Combined Secondary Vertex	<b>JES</b>	Jet Energy Scale
<b>CTA</b>	Cherenkov Telescope Array	<b>KF</b>	Kalman Filter
<b>DAQ</b>	Data AcQuisition system	<b>L1</b>	Level-1 Trigger
<b>DAS</b>	Data Aggregation System	<b>LAT</b>	Fermi Large Telescope
<b>DCS</b>	Detector Control System	<b>LEP</b>	Large Electron Positron collider
<b>DQM</b>	Data Quality Monitoring	<b>LHC</b>	Large Hadron Collider
<b>DM</b>	Dark Matter	<b>LNGS</b>	Laboratori Nazionali del Gran Sasso
<b>DMWG</b>	Dark Matter Working Group	<b>LO</b>	Leading Order
<b>DT</b>	Drift tube	<b>LS</b>	Long Shutdown

<b>LSP</b>	Lightest Supersymmetric Particle	<b>QFT</b>	Quantum Field Theory
<b>MACHO</b>	Massive Compact Halo Object	<b>RMS</b>	Root Mean Square
<b>MC</b>	Monte Carlo	<b>ROC</b>	Receiver Operating Characteristic
<b>MET</b>	Missing Transverse Momentum	<b>RPC</b>	Resistive Plate Chamber
<b>MFV</b>	Minimal Flavour Violation	<b>SC</b>	Super Cluster
<b>ML</b>	Machine Learning	<b>SD</b>	Spin Dependent
<b>MPI</b>	Multiple Parton Interaction	<b>SF</b>	Scale Factors
<b>MSSM</b>	Minimal Supersymmetric Standard Model	<b>SI</b>	Spin Independent
<b>MVA</b>	Multi-Variate Analysis	<b>SM</b>	Standard Model
<b>NFW</b>	Navarro-Frenk-White	<b>SPS</b>	Super Proton Synchrotron
<b>NLO</b>	Next to Leading Order	<b>SR</b>	Signal Region
<b>PDF</b>	Parton Density Function	<b>TEC</b>	Tracker EndCap
<b>PF</b>	Particle Flow	<b>TIB/TBD</b>	Tracker Inner Barrel and Disks
<b>POG</b>	Physics Object Group	<b>TOB</b>	Tracker Outer Barrel
<b>PR</b>	Prompt Rate	<b>TPG</b>	Trigger Primitive Generators
<b>PS</b>	Proton Synchrotron	<b>UE</b>	Underlying Event
<b>PU</b>	Pile up	<b>UED</b>	Universal Extra Dimensions
<b>PUPPI</b>	Pileup Per Particle Identification	<b>VBF</b>	Vector Boson Fusion
<b>PV</b>	Primary Vertex	<b>WIMP</b>	Weakly Interactive Massive Particle
<b>QCD</b>	Quantum ChromoDynamics	<b>WP</b>	Working Point

---

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>The dark matter case</b>	<b>5</b>
2.1	The Standard Model (SM) . . . . .	5
2.2	At the origins of dark matter . . . . .	6
2.2.1	Zwicky and the virial theorem . . . . .	7
2.2.2	Spiral galaxies rotation curves . . . . .	7
2.2.3	Cosmic Microwave Background (CMB) anisotropies . . . . .	8
2.2.4	Gravitational lensing . . . . .	10
2.3	Dark matter properties . . . . .	11
2.4	Dark matter candidates . . . . .	13
2.5	Dark matter searches . . . . .	18
2.5.1	Direct searches . . . . .	19
2.5.2	Indirect searches . . . . .	22
2.5.3	Collider production . . . . .	25
2.6	Dark matter production at the Large Hadron Collider (LHC) . . . . .	28
2.7	The focus of this thesis . . . . .	30
2.8	Previous relevant results . . . . .	32
<b>3</b>	<b>The experimental setup</b>	<b>37</b>
3.1	The Large Hadron Collider (LHC) . . . . .	37
3.1.1	The LHC in a nutshell . . . . .	38
3.1.2	Key parameters . . . . .	39
3.2	The CMS detector . . . . .	41
3.2.1	Tracker . . . . .	44
3.2.2	Electromagnetic Calorimeter (ECAL) . . . . .	46
3.2.3	Hadronic Calorimeter (HCAL) . . . . .	48
3.2.4	Solenoid . . . . .	49
3.2.5	Muon system . . . . .	50
3.2.6	Trigger system . . . . .	53

3.2.7	Data AcQuisition system (DAQ)	55
3.3	CMS main goals	56
<b>4</b>	<b>Event reconstruction</b>	<b>57</b>
4.1	Particle Flow (PF) algorithm	57
4.2	Primary vertex definition	58
4.3	Leptons reconstruction	60
4.3.1	Muons	60
4.3.2	Electrons	62
4.4	Jets reconstruction	64
4.4.1	B-tagging	65
4.5	Missing Transverse Momentum (MET)	67
4.6	Top reconstruction	70
4.6.1	Numerical and analytical top reconstruction	70
4.6.2	Top reconstruction with additional dark matter	72
4.6.3	Top reconstruction in practice	73
<b>5</b>	<b>Data, signals and backgrounds</b>	<b>77</b>
5.1	The Monte-Carlo simulation method	77
5.2	Files format	80
5.3	Analysis code	81
5.4	Data samples	82
5.5	Signal samples	82
5.6	Backgrounds prediction	82
5.6.1	Top production	84
5.6.2	Drell-Yan estimation	87
5.6.3	$t\bar{t} + W/t\bar{t} + Z$	89
5.6.4	Non prompt leptons contamination	90
5.6.5	Smaller backgrounds	95
5.6.6	Weights and corrections applied	96
<b>6</b>	<b>Event selection</b>	<b>99</b>
6.1	Objects selection	99
6.1.1	Triggers selection	99
6.1.2	Electrons selection	100
6.1.3	Muons selection	101
6.1.4	Leptons selection	102
6.1.5	Jet selection	102
6.2	Signal regions	103
6.3	Control regions	105

6.3.1	Inclusive control region . . . . .	105
6.3.2	Top control region . . . . .	105
6.3.3	DY control region . . . . .	114
6.3.4	$t\bar{t} + W/t\bar{t} + Z$ control region . . . . .	114
6.3.5	Same sign control region . . . . .	114
<b>7</b>	<b>Signal extraction</b>	<b>119</b>
7.1	Discriminating variables . . . . .	119
7.2	Multivariate analysis . . . . .	126
7.2.1	Methods used . . . . .	127
7.2.2	Training process . . . . .	135
7.2.3	Evaluation process . . . . .	135
7.2.4	Shape analysis . . . . .	136
<b>8</b>	<b>Results and interpretations</b>	<b>143</b>
8.1	Statistical interpretation . . . . .	143
8.2	Systematics and uncertainties . . . . .	146
8.3	Results . . . . .	149
<b>9</b>	<b>Conclusions</b>	<b>153</b>
9.1	Future prospects . . . . .	154
<b>Appendices</b>		<b>155</b>
<b>A</b>	<b>Resumen en español</b>	<b>157</b>
A.1	Materia oscura . . . . .	157
A.2	El dispositivo experimental . . . . .	159
A.3	Reconstrucción de objetos . . . . .	160
A.4	Análisis de datos . . . . .	161
A.5	Resultados obtenidos . . . . .	163
<b>B</b>	<b>Samples used</b>	<b>165</b>
B.1	Data samples . . . . .	165
B.2	Signal samples . . . . .	165
B.3	Backgrounds samples . . . . .	165
<b>C</b>	<b>Multi-Variate Analysis (MVA) optimization</b>	<b>175</b>
<b>D</b>	<b>Pulls and impacts plots</b>	<b>185</b>
<b>Bibliography</b>		<b>200</b>



---

---

# Chapter 1

---

## Introduction

The Standard Model (SM) of particle physics [1] is nowadays the most accepted mathematical model used to describe the elementary particles and three of the 4 fundamental forces of nature (electromagnetic, weak and strong interactions, while the gravitational interaction is out of reach of this model). Even though quite simple in concept, it has been able to describe most of the phenomena observed in nature so far with an incredible level of precision, and has made a lot of predictions that have now been proven to be true, such as the postulate of the Brout-Englert-Higgs mechanism [2, 3] followed by the discovery of the Higgs boson itself [4, 5] by the CMS [6] and A Toroidal LHC ApparatuS (ATLAS) [7] experiments analyzing the proton-proton collisions produced by the Large Hadron Collider (LHC) at a center of mass energy  $\sqrt{s} = 7$  and 8 TeV, announced at the European Council for Nuclear Research (CERN) on the 4th of July 2012.

However, as accurate as it seems to be, this theory, introduced in Section 2.1, is known to have several shortcomings which require further investigation. Eventual exotic particles which do not fit in the current model could in this sense be the sign of new physics and have therefore been extensively searched for over the course of the last decades in order to enhance our understanding of the Universe and all its constituents.

In this context, the first serious Dark Matter (DM) hypothesis was introduced in the 1930s because of gravitational anomalies observed by several astrophysicists, as a way to explain the apparent non-luminous missing mass in the Universe [8]. Indeed, the visible mass in most galaxies appears to be way too low to explain several astrophysical processes, such as the rotation curves of the galaxies [9], which seems to be incompatible with the well established laws of gravitation. Some additional measurements of the gravitational lensing (in the Bullet Cluster, for example [10]) and the anisotropies observed in the Cosmic Microwave Background (CMB) [11] are other evidences for the existence of DM, as explained in Section 2.2.

As far as we currently know from cosmological measurements, ordinary baryonic matter only constitutes around 5% of the Universe, while DM represents around 26% of the total energy density of the Universe (the rest is being considered as dark energy) [12]. Trying to understand the nature and fundamental properties of this new kind of exotic matter, as done in Section 2.3, is therefore crucial to try and understand the laws of physics in the Universe, with many theorists and large experiments around the world currently involved in such searches.

Nowadays, the existence of DM is well motivated in the physics community, even though it has

never been observed directly, since our only evidences so far for its existence come from its large-scale gravitational effects. While its mass, spin, nature and basic properties are still unknown and extensively studied, one of the best DM candidates are the so-called Weakly Interactive Massive Particles (WIMPs), introduced in 2.4 predicted to interact both gravitationally and weakly with SM particles. This would allow direct and indirect direction of such candidates, used as the driving process of many experiments over the last decades, trying to find hints for a possible interaction between standard baryonic particles and DM particles, or even between several DM particles themselves. Dark matter production through the use of a particle accelerator colliding SM particles together, such as the LHC, is also a possibility and will be considered as the main channel towards the eventual detection of this exotic matter throughout this work. The production through colliders is actually able to provide constraints on low dark matter masses as well, in a region where both the direct and indirect searches are less efficient, which makes the LHC a perfect tool to study this kind of Beyond the Standard Model (BSM) physics. All these searches will be summarized in Section 2.5.

However, observing DM is still extremely difficult, mainly because it barely interacts with ordinary baryonic matter, except through gravity (we have to assume that it does interact with SM at least weakly for the sake of this work though, as we would not be able to discover it as an individual particle if it were not the case). This means that nowadays, all the experiments searching for DM have only been able to put constraints on the DM particle mass and on the interaction cross sections between the dark and standard sectors. Actually, even if the collisions between protons produced by the LHC do have a sufficient amount of energy to produce this kind of particles, we would not expect them to interact with our detector, making their actual detection even harder. The presence of such matter has to be inferred from the study of the interaction between SM particles and CMS itself, since a typical DM-like event consists of at least one energetic SM particle produced in association with a large imbalance in the momentum due to the presence of an eventual DM candidate that was able to escape our detection. This kind of collider searches considered in this work can be grouped under the name of *mono-X* searches, where the *X* stands for any kind of SM particle (a jet, a lepton or a photon for example).

In the context of this work, DM is searched for in association with one or two top quarks which play the role of the SM particle allowing us to actually trigger the event, as described in Section 2.7. The top quark, the most massive of all the fundamental particles observed by far, is indeed an excellent object to study in this context, mainly because of its high mass and because of the expected Yukawa-like coupling structure of the new physics model [13]. However, this also means that the phenomenology of this quark is mostly driven by its large mass and that it decays before hadronization can occur, almost always into a W boson and a bottom quark. The final state of the process we are interested in is then made of some b-jets, leptons and/or quarks and is mostly categorized depending on the decay of the W itself. This work will actually be focused on the two leptons final state, also known as the dileptonic channel, mostly because leptons are by reconstruction much cleaner than jets and because this channel does not have lots of background processes raising to a similar final state, even though its branching ratio is the smallest, as will be explained in Section 2.6.

The LHC has now been running for a decade, and several similar searches have already been carried out and published in the past by the CMS and ATLAS collaborations, at different center of mass energies, as described in Section 2.8. First of all, at 8 TeV, several searches for a pair of top quarks were published by the CMS (in association with DM in the semileptonic [14] and dileptonic [15] final states) and ATLAS collaborations [16]. Then, at 13 TeV, the ATLAS collaboration published on one hand several studies, considering different final states and different luminosities (ranging from  $13.3 \text{ fb}^{-1}$  to  $137.1 \text{ fb}^{-1}$ ) [17, 18, 19, 20, 21]. On the other hand, the CMS collaboration published a few extremely important papers for this study [22, 23]. For the first time in 2019, the results obtained by the  $t/\bar{t}+\text{DM}$  and  $t\bar{t}+\text{DM}$  analyses have also been combined and published using

the data collected during the year 2016 [24]. Our main objective is now to repeat and improve this analysis while considering the full Run II dataset, globally improving the analysis strategy and including the dileptonic final state for the first time in this combination.

After a general introduction about dark matter in Chapter 2, the experimental setup will be detailed in Chapter 3. This will include a discussion about the LHC, along with a complete description of CMS, the detector used to collect the data that will be analyzed throughout this work. The data has been collected during the years 2016, 2017 and 2018 and corresponds to an integrated luminosity of  $(137.1 \pm 0.2) \text{ fb}^{-1}$ , collected during the Run II of operation of the LHC and at a center of mass energy  $\sqrt{s} = 13 \text{ TeV}$ . A particular care will be given to the explanation of the Particle Flow (PF) algorithm, used to reconstruct the different objects of the analysis and that will be defined in Chapter 4, while the estimation of the different backgrounds and the selection of interesting events will be detailed throughout Chapters 5 and 6.

Distinguishing between the signals we are studying and backgrounds having a much higher cross-section and kinematically really close, such as the SM  $t\bar{t}$  without production of DM is not a straightforward task (sometimes a production of missing transverse momentum due to the presence of physical neutrinos is even obtained). To enhance the signals and to obtain some discrimination between these kind of processes, an algebraic reconstruction of the event and top-notch Machine Learning (ML) techniques are used in this work, in order to train a network of neurons to perform this task. The main objective is to make them learn how to combine the discriminating power of a set of input variables in order to create a single output variable describing the probability of a single event to be classified as signal or background. All this process will be detailed in Chapter 7.

Finally, a statistical interpretation of our data will be performed and different sources of systematic uncertainties will be considered in Chapter 8. This will allow us to set upper limits on the signal strength of DM particles produced in our particular channel and for the simplified models considered in this analysis. The conclusions of this work and some additional prospects for the future will then finally be presented in Chapter 9.



---

---

# Chapter 2

---

## The dark matter case

In this chapter, some general explanations about the SM will first of all be given in Section 2.1 as a general introduction about today's most accepted mathematical model describing all the known particles and their interactions. The case for DM will then be presented in Section 2.2, along with a summary of the main evidences, mostly astrophysical, which lead to the introduction of this kind of Beyond the Standard Model (BSM) physics. Then, the main properties expected by such exotic matter will be introduced in Section 2.3 and nowadays most accepted DM candidates, such as the Weakly Interactive Massive Particle (WIMP) will be presented in Section 2.4. The main ways we have to search (direct, indirect and collider searches) for this new exotic physics along with the main experiments dedicated to such searches will then be shown in Section 2.5. Finally, the searches performed in colliders such as the LHC and our particular channels of interest (DM produced in association with either one or two top quarks) will be detailed in Section 2.6 and the latest similar results and exclusion limits published over the course of the previous years by the ATLAS and CMS collaborations will be shown in Section 2.8.

### 2.1 The Standard Model (SM)

The SM of particle physics is a relativistic Quantum Field Theory (QFT) able to mathematically summarize our current understanding of all the known particles and to describe three out of the four main interactions between these particles (the weak, strong and electromagnetic forces, while an explanation of the origin of the gravitational interaction is out of the reach of this theory).

If the neutrinos are considered to be normal Dirac fermions in the sense the antineutrinos are expected to be different than neutrinos (the question whether or not they could actually be Majorana particles is still well under discussion [26]), then the SM contains 26 free parameters, mainly the masses of the 12 predicted fermions, as shown in Figure 2.1, along with the couplings describing the strengths of the three interactions, two parameters describing the Higgs potential, eight mixing angles and, maybe, the phase of the eventual strong CP violation that will be mentioned in Section 2.4 (but in any case this parameter is expected to be close to 0). This is a quite high number of free parameters, but the value of most of them has now been derived experimentally.

The SM Lagrangian density in a differential volume element  $\mathcal{L} = \int L d^3x$ , accounting for the

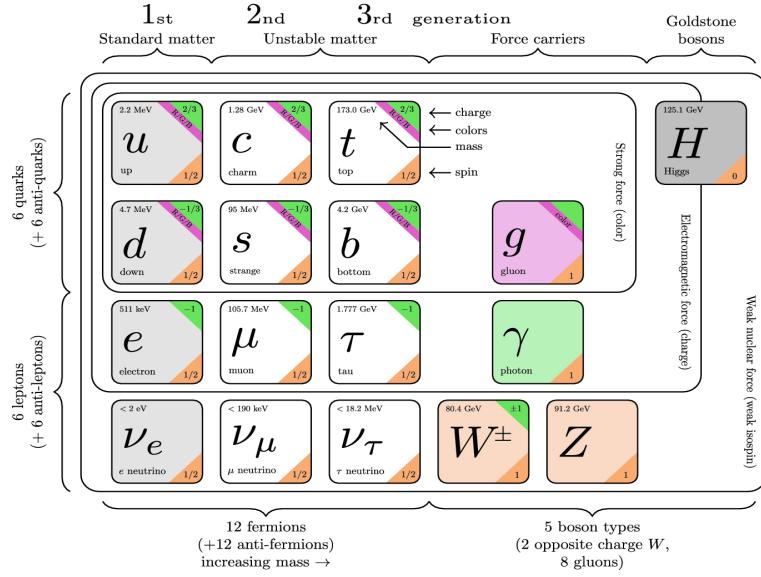


Figure 2.1: Representation of the 12 fermions of the SM [25] along with the main force carriers and the Higgs boson, discovered in 2012 and completing the SM.

kinetic and potential energy of a system, takes the (very) simplified form given in Equation 2.1, where  $F_{\mu\nu}$  is the field strength tensor accounting for the different interactions,  $\psi$  is the interacting field describing quarks and leptons and  $\phi$  is the Brout-Englert-Higgs field while  $y_{ij}$  are the Yukawa couplings to this field, which depend on the mass of the particle considered and which will be described later on in Section 2.5.3.

$$\mathcal{L} = -\frac{1}{4}F_{\mu\nu}F^{\mu\nu} + i\bar{\psi}\mathcal{D}\psi + \bar{\psi}_i y_{ij} \psi_j \phi + |D_\mu \phi|^2 - V(\phi) \quad (2.1)$$

A complete description of this model is out of the scope of this work but it is important to note that the SM, although mostly experimental, is still working extremely well today. Indeed, it managed to successfully describe the outcome of all the experimental data and to make predictions on new phenomena, so far always confirmed (we can quote as example that it successfully predicted the existence of the gluons, the top quarks, along with the  $W$ ,  $Z$  and Higgs bosons [27]).

However, we know that the SM is not a complete theory of Particle Physics as it fails to explain some known phenomena in Nature. There are in this sense some open questions and many BSM theories trying to explain such observations, such as an eventual inclusion of the gravitational interaction within this model. We can quote for example as BSM theories the possible existence of the supersymmetry [28], telling us that each particle should have a super partner whose spin differs by 1/2 or the possible existence of DM particles, the main subject of the following sections.

## 2.2 At the origins of dark matter

The origin of the concept of dark matter can be traced back to the 19th century, even though this concept has changed quite a lot over the years. Back then, DM was more considered to be ordinary matter which simply did not emit any kind of electromagnetic radiation, being therefore invisible and dark, but which does have a strong gravitation impact because of its mass. In the 20th century, astronomical and cosmological observations established the existence of such a species of matter embedded on the core of giant astronomical objects such as galaxies and clusters [29].

### 2.2.1 Zwicky and the virial theorem

In the 20th century, the first experimental evidences for the existence of dark matter were found. In 1933, Fritz Zwicky managed to determine the mass of the Coma Cluster using the virial theorem [30], which states that in a cluster in equilibrium under its own gravitation the kinetic energy must be comparable to its gravitational binding energy.

Mathematically, the virial theorem can be written in Equation 2.2, where the brackets represent the mean value of the quantity obtained over time or position, the universal constant of gravitation is  $G = 6.67 \cdot 10^{-11} \text{ m}^3 \text{ kg}^{-1} \text{ s}^{-2}$  and where the gravitation potential energy expression can be simplified assuming a spherical distribution of the masses and the same average density everywhere in the cluster considered for the calculation.

$$2\langle T \rangle + \langle V \rangle = 0, \text{ where } \begin{cases} T = \frac{1}{2} \sum_i M_i v_i^2 = \frac{1}{2} M \langle v^2 \rangle \\ V = -4\pi G \int_0^R M \rho r dr \propto \frac{GM^2}{R} \end{cases} \quad (2.2)$$

Solving these simple equations gives us an approximate value of the mass of the cluster in Equation 2.3, where  $R$  is the radius of the cluster and  $\langle\langle v^2 \rangle\rangle$  is the squared velocity of all the galaxies averaged over position and time.

$$M \propto \frac{\langle\langle v^2 \rangle\rangle R}{G} \quad (2.3)$$

Zwicky studied the velocity dispersion of the galaxies in the Coma Cluster and used this formula to estimate the total mass of the cluster, concluding that this value was around 400-500 times larger than the mass previously estimated by Edwin Hubble, who simply considered the number of visible galaxies within this cluster for his calculations. One plausible explanation for this discrepancy is to introduce the concept of DM, which contributes to the mass of the cluster without increasing the galactic luminosity.

Zwicky's results were actually quite controversial since they were based on statistical calculations relying on different hypotheses not always justified, such as the fact that the galaxies in the cluster must be gravitationally bound with each other and they were actually proven to be overestimated later on [31], but additional observations came to enforce the validity of his conclusions anyway.

### 2.2.2 Spiral galaxies rotation curves

Despite being controversial and slightly off, Zwicky's results were soon followed by a series of additional astronomical observations leading to the same conclusion, the possible existence of non-luminous matter in all the galaxies, called dark matter. The most famous of these results is the study of the observed and expected rotation curves of the stars within spiral galaxies such as the Milky Way in the 1970s [9].

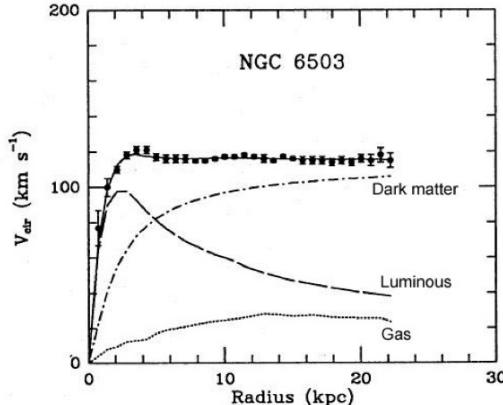
According to this study, if we assume that we can apply Newton's universal laws of gravitation at the galactic scale, then the stars within this kind of galaxies should rotate with a velocity depending on the radius to the galactic center obtained by the usual equation for centripetal acceleration in

a gravitational field and represented in Equation 2.4, where  $M(r)$  accounts for the total mass encountered in a radius  $r$ .

$$v_{\text{rotation}}(r) = \sqrt{\frac{GM(r)}{r}} \quad (2.4)$$

At first approximation, one can assume that most of the mass within this kind of galaxies comes from the inner core, which means that, at large radius, the velocity of individual stars is expected to decrease proportionally to  $r^{-1/2}$ . Any deviation to this rule suggests that either our understanding of gravity at large scales or our basic understanding of galaxies as a celestial body made of stars, dust and gas, has to be revised.

Actually, observations made by Vera Rubin and her team in the early 1970s with a new spectrograph designed to study the velocity curves of spiral galaxies with a degree of accuracy never reached before, did not confirm these expectations [32]. Indeed, according to these results, from a given value of the radius, the velocity curve appears to be flat instead of decreasing, as illustrated in Figure 2.2. This is another hint that can motivate the introduction of the concept of DM.



K.G. Begeman, A.H. Broeels, R.H. Sanders. 1991. Mon.Not.RAS 249, 523.

Figure 2.2: Expected and observed rotation curves of the galaxy NGC 6503 [9]. The black dots correspond to the data and the *luminous* line corresponds to the rotation curve decreasing as  $r^{-1/2}$  expected from Newtonian dynamics.

### 2.2.3 Cosmic Microwave Background (CMB) anisotropies

The CMB is a mostly uniform background of primary radio waves emitted when the Universe became transparent around 380 000 years after the Big Bang and was discovered accidentally in the 1940s [33]. Studying it is extremely important, as it is actually made of the oldest and cleanest electromagnetic radiation we can find in the Universe. Precise measurements of this radiation are actually critical in many different fields of physics, since any proposed model of the Universe must be able to explain this radiation, its temperature and anisotropies.

Recent measurements determined that the CMB can be considered as emitting a thermal black body spectrum at a temperature of  $(2.72548 \pm 0.00057)\text{K}$  [34]. However, we now know that this temperature is actually not constant as some small anisotropies can be observed (at the  $10^{-5}$  level), depending on the value of the solid angle of observation, as represented in Figure 2.3.

We see these fluctuations projected over a 2D sphere, and it is therefore natural to introduce at this point Laplace's spherical harmonics,  $Y_{lm}(\theta, \phi)$ , a complete set of orthogonal functions obtained by solving Laplace's equation  $\nabla\psi = 0$  on a sphere and defined by a few parameters such as  $l$ , the

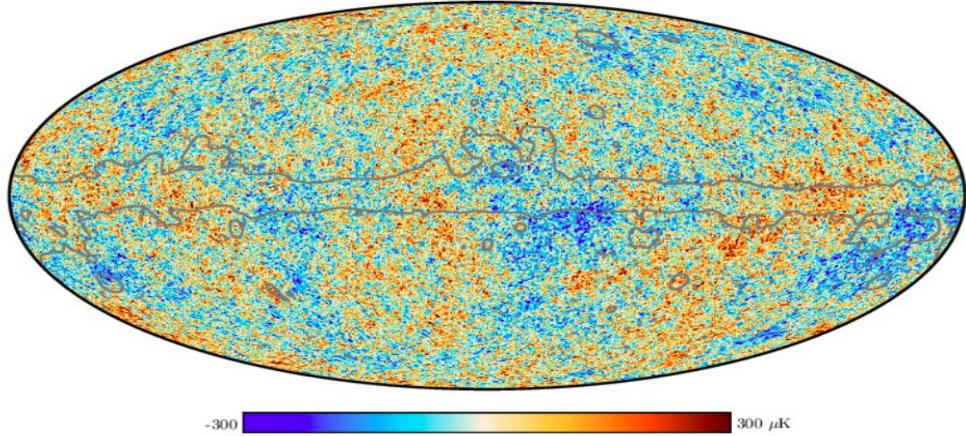


Figure 2.3: Anisotropies at the  $10^{-5}$  level in the temperature of the CMB, as observed by the Planck satellite in 2018 [35].

multipole, representing a given solid angle in the sky ( $l = 100$  corresponds to  $\sim 1^\circ$ ) and  $m$ , the number of poles, such as  $-l \leq m \leq l$  [36].

It is possible to show that these spherical harmonics form a complete orthonormal basis on this space and therefore that any function that can be defined on the sphere may be expanded into these harmonics using coefficients called  $a_{lm}$ . The temperature fluctuations, whose value depend on the two usual spherical angles  $\theta$  and  $\phi$  can therefore be expanded using these generic functions, according to Equation 2.5.

$$\frac{\Delta T}{T}(\theta, \phi) = \sum_{l=0}^{\infty} \sum_{m=-l}^{m=l} a_{lm} Y_{lm}(\theta, \phi) \quad (2.5)$$

The information about the anisotropies can actually be extracted from the variance of these harmonic coefficients  $a_{lm}$  of the expansion since the CMB is assumed to be a Gaussian random field. The power spectrum of the CMB can be extracted according to Equation 2.6, and from which most of the cosmological information of the CMB is derived.

$$D_l = \frac{l(l-1)C_l}{2\pi} = \sum_m |a_{lm}^2| \quad (2.6)$$

Interestingly enough, this spectrum is directly affected by the value of the density of the dark matter in the Universe. Indeed, if we remove the effect of the  $l = 0$  and  $l = 1$  poles, which do not have anything to do with the origin of the Universe, the  $l = 2$  poles can be related to the fluctuations happening before the recombination phase. Such fluctuations depend on the mass but also on the radiative pressure at this epoch, and should be null for dark matter but not for baryonic matter. Doing a multi-parameters fit on the observed data represented in the power spectrum (cf. Figure 2.4) is then able to give us directly the energy density of baryonic  $\Omega_b$  and dark  $\Omega_\chi$  matter, along with other important parameters of the  $\Lambda CDM$  cosmological model. Today's most precise measurements have been obtained in 2018 using the Planck satellite, and lead to the determination of these two quantities:  $\Omega_b h^2 = (0.02220 \pm 0.00020)$  and  $\Omega_\chi h^2 = (0.1185 \pm 0.0015)$  [37]. This is one of the strongest constraint we have on DM so far, since any DM candidate will need to comply with this measurement.

By dividing these results with the value of the scaling factor for the Hubble expansion rate  $h = 0.674$

[38], we can obtain from these numbers a proportion of 4.9% of ordinary baryonic matter and 26.1% of dark matter in the Universe.

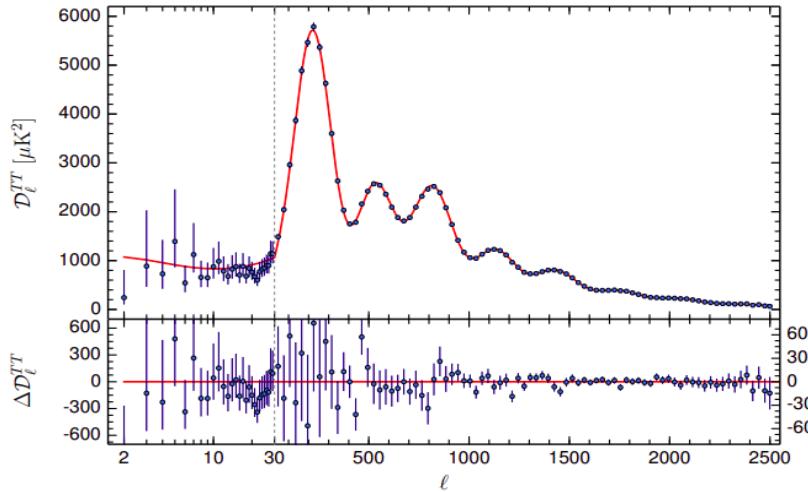


Figure 2.4: Power spectrum of the CMB obtained by Planck, representing the fluctuations of the temperature of the radiation with respect to the angular angle of observation [37].

## 2.2.4 Gravitational lensing

The last evidence supporting the existence of dark matter has been obtained by observing several clusters of galaxies in the Universe, such as the Bullet Cluster, and by studying their mass distribution through gravitational lensing.

The gravitational lensing effect is a consequence of the general relativity, a theory developed by Einstein as a way to represent gravity using the geometry of space-time, stating that massive objects lying between distant sources and an observer should act as a lens and bend the light emitted by the source. This deviation of the light is actually proportional to the mass of the object in between the source and the observer, meaning that the gravitational lensing can give us a way to measure the mass distribution of massive objects, such as galaxy clusters. This mass distribution can then be compared to the luminous distribution of the cluster, to see if we can observe a discrepancy between the two measurements, which could be another hint of the existence of DM.

The bullet Cluster is particularly interesting in this context since it actually provides an evidence for the eventual existence of DM which does not rely on any mathematical assumption (other than the general relativity principle) and cannot at principle be explained by alternate laws of gravitation. In this case, some observations clearly showed that the spatial deviations between the center of the total mass and the center of the baryonic luminous mass cannot be explained with an alteration of the gravitational force law alone, with a statistical significance of  $8\sigma$  [39].

As seen in Figure 2.5, the image taken by Chandra clearly shows an offset between the visible plasma of the cluster and the actual mass distribution measured through gravitational lensing (green contours). The center of the luminous mass of the cluster does not seem to match the one obtained considering its non-luminous counterpart as well, which is another evidence of the possible existence of DM within galaxy clusters.

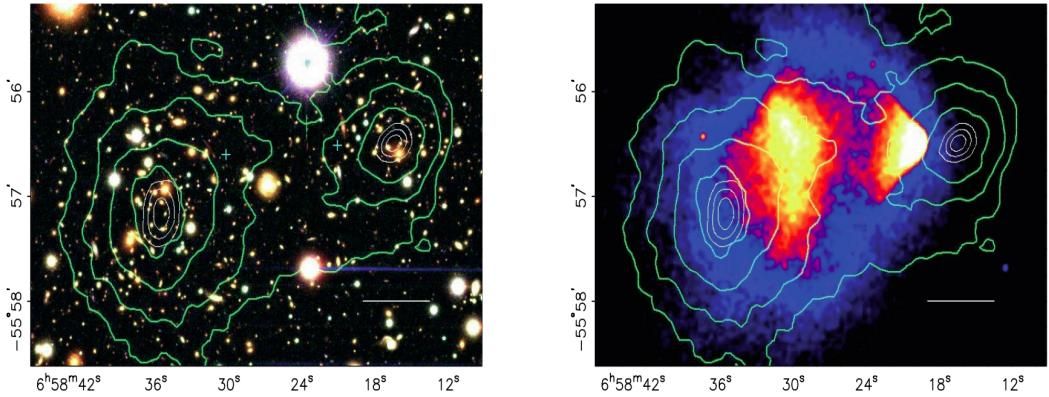


Figure 2.5: Mass distributions obtained by the Magellan telescope in the visible (on the left) and the Chandra telescope on the X-rays spectrum (on the right) telescopes of the Bullet Cluster. Being shifted compared to each other, this is yet another clear evidence for the existence of DM [39].

## 2.3 Dark matter properties

All the previous observations allow us to list some of the most important properties that any dark matter candidate should have. Even though several theories exist, each giving slightly different properties to the DM, we will consider in this work the following mostly accepted properties for such particles:

- First of all, we will assume that DM is a **particle**. As far as we know, the Universe is simply composed of particles so there is no objective reason to think that DM, being matter with a certain mass, might not be made of some kind of indivisible particles at some level.
- Then, the DM candidate should of course be **dark**. This means that it should not interact at all with electromagnetic radiation such as light, and that it should therefore be **electrically neutral**. However, it has to interact at least gravitationally because of the observations for the evidence of such a particle explained before, mostly relying on gravitational effects, and we assume in this work that it interacts weakly as well to have a chance to discover it with particle accelerators.
- It is **non-baryonic**, mainly because the energy density for the baryonic matter obtained by observing the power spectrum of the CMB is too low to account for DM as well, as explained in Section 2.2.3. Indeed, according to these results, baryonic matter accounts for around 5% while dark matter accounts for more than 25% of the energy density of the Universe.
- We will also only consider **cold** dark matter, since the widely accepted  $\Lambda_{CDM}$  cosmological model is actually based on this assumption. By cold, we do not refer to the temperature of these particles but actually to their size, and therefore to the velocity by which they can travel in the Universe. Large scale structures of the Universe such as we can observe them today cannot actually be explained if DM is made of hot/relativistic particles, as represented in Figure 2.6. However, although not really as popular, it is important to note that alternative models with warm DM have also been developed and still exist today.
- It is interesting to report as well that DM particles are expected to be found near the electroweak symmetry breaking scale, between **10 GeV and 1 TeV**. This is a consequence of the expected production mechanism of such particles, the so-called thermal freeze-out [40]. This principle tells us that at some point in history, DM was supposed to be in thermal

equilibrium with other primordial SM particles, meaning that its production and annihilation rates were equal. However, because of the expansion of the Universe, at some point DM particles were simply too far apart from each other and these reactions maintaining this equilibrium were not efficient enough any more. At this stage, the abundance of DM became fixed: this is the freeze-out, as represented in Figure 2.7.

This principle is interesting because, as a rule of thumb, we can say that if a particle interacts heavily, it will stay longer in equilibrium and its freeze-out abundance will therefore be smaller, so there is a mathematical relation between the strength of the SM/DM interaction  $g$ , the mass of the DM particle  $m_\chi$  and its relic abundance  $\Omega_\chi$  that has been precisely measured, as expressed in Equation 2.7 [41].

$$\Omega_\chi \propto \frac{m_\chi^2}{g^4} \quad (2.7)$$

By using a typical value for  $g$  of the order of the Fermi coupling constant  $G_0^F \simeq 4.54 \cdot 10^{14} \text{ J}^{-2}$  we can see that, in order to observe a freeze-out abundance comparable to the one observed recently by the Planck satellite, the DM candidate should have a mass between 10 GeV and 1 TeV as previously stated. The measurement of the CMB power spectrum is therefore able to put constraints on the DM cross-section with the baryonic sector, and all the DM candidates quoted in Section 2.4 will have to respect this criteria.

In any case, DM is not expected to have a mass lower than 300 eV since at this scale, the phase space density that would be needed to explain the relic abundance of DM would simply violate the Pauli-exclusion principle [44].

- Finally, the DM particles should be **long-lived**. Indeed, we expect that they were produced 13.8 billion years ago during the Big Bang, but it seems they are still present in the Universe since we still see their effect today. They should then be stable particles, or their lifetime should at least be larger than the age of the Universe itself.

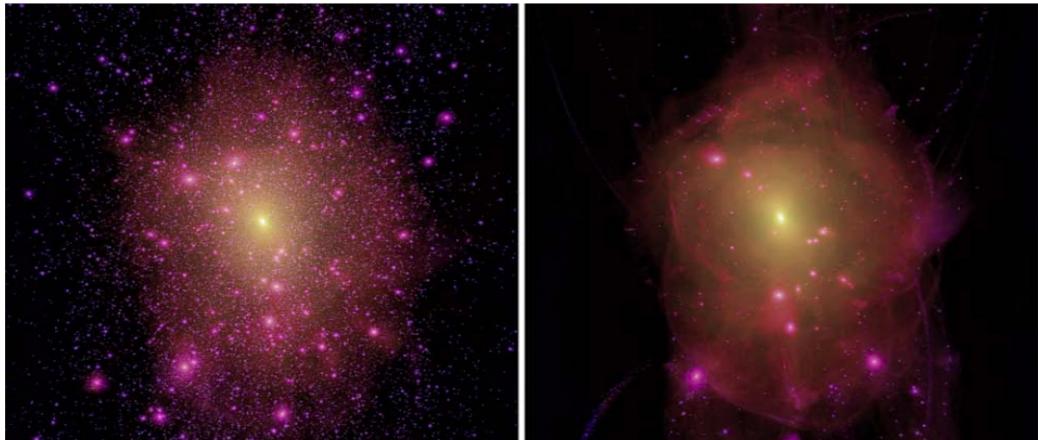


Figure 2.6: Computer simulations for cold (on the left) and warm (on the right) DM scenarios and their impact on a galactic halo at 0 redshift [42].

All these properties narrow quite a lot the number of possible DM candidates.

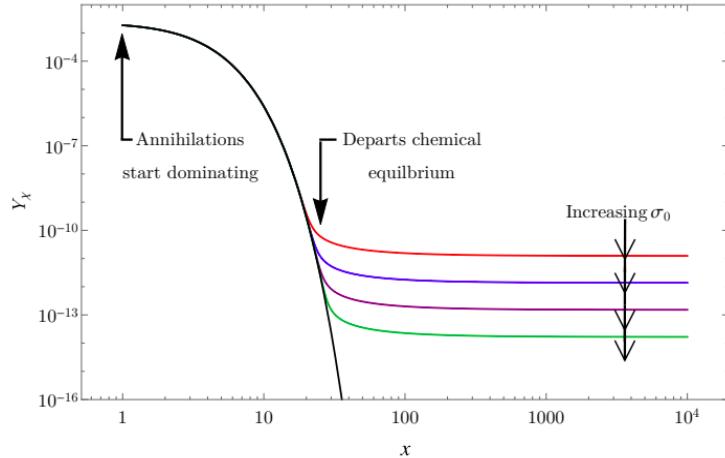


Figure 2.7: Schematic representation of the freeze-out process, representing the abundance of a 500 GeV DM as  $Y_\chi$  with respect to the time and the impact of increasing cross-section annihilation values on this freeze-out abundance [43].

## 2.4 Dark matter candidates

Several different categories of particles could pretend to be good candidates for dark matter but only the most interesting ones will be quoted here, since an extensive list of all the different possible candidates is out of the scope of this work. Two SM particles will first of all be investigated, before introducing some BSM theories providing additional DM candidates with the expected properties.

### Massive Compact Halo Objects (MACHOs)

The first obvious DM candidates are the so-called MACHOs. These objects are massive astronomical non-luminous bodies (such as black holes) made of baryonic matter and very hard to detect, that could be responsible of the gravitational lensing observed and that could explain the apparent missing mass in the Universe. However, as we saw in Section 2.3, DM is not expected to be made of such ordinary matter, mainly because observational data of the CMB and the deduced baryonic density of energy in the Universe is able to rule out this possibility.

Several different experiments did try to search for such DM anyway, and managed to constrain the properties of this kind of objects. The main way to search for such massive objects is through their gravitational microlensing effect since, according to the general relativity principle, they should bend the light of luminous objects located behind them, such as stars, and this bending actually depends on the mass of the eventual MACHO. Experiments like the MACHO project and EROS observed in this context  $\Theta(10^7)$  stars for several years, looking for microlensing events in order to constrain this particular DM model. Results published in 2000 from the MACHO project, after studying almost 12 million stars, actually observed between 13 and 17 such events, lower than expected if DM was only made of MACHOs.

The collaboration actually managed to exclude at the 95% Confidence Level (CL) the possibility of the dark halo to be entirely made of such baryonic particles [45]. On the other hand, the EROS collaboration only observed 1 microlensing after studying more than 30 millions of bright stars during 6.7 years, while  $\sim 39$  events were expected [46]. Both results have also been combined in order to obtain the exclusion plot represented in Figure 2.8. From this study, MACHOs with low

masses of the order of a fraction of the mass of the Sun  $M_\odot$  ( $10^{-7}M_\odot < m < 10^{-3}M_\odot$ ) should make up less than 25% of the dark matter halo for most models considered at the 95% CL [47].

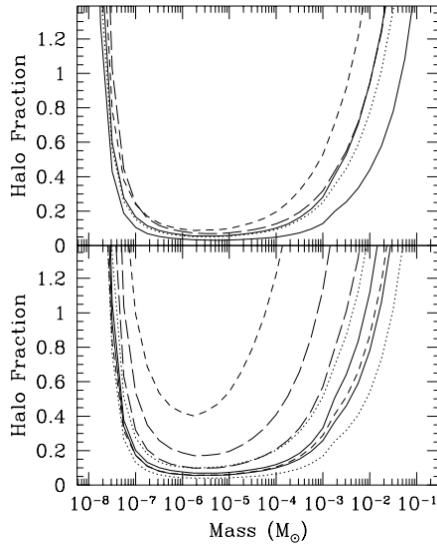


Figure 2.8: Halo fraction upper limits at the 95% CL compared to the mass of the lensing object for different MACHO models considered by the EROS (on the top) and the MACHO (on the bottom) collaborations [47].

## Active neutrinos

SM *active* neutrinos  $\nu$  (as opposed to *sterile* neutrinos, which will be the subject of the discussion in the next section) have been considered as good DM candidates for a long time as well, since they are electrically neutral and long-lived SM particles, two important properties of any DM candidate. They have a few particular properties that might be interesting in this context.

- First of all, even though it has still not been measured precisely, the sum of their masses has recently been measured to be lower than 0.17 eV at the 95% CL from cosmological studies [51]. Even though this is not fully understood, such value is incredibly low compared to other SM particles, this particularity usually being referred to as the *mass puzzle* of the neutrinos.
- Their low mass has a consequence in the sense that it means that the gravitational interaction between two neutrinos is usually considered to be negligible and that we can consider that they only interact weakly, making it hard to study their properties. Their actual cross-section of interaction, as represented in Figure 2.9 and which of course depends on their energies and on the channel of interaction (neutral or charged current considered), is therefore usually extremely low, making it hard to detect neutrinos.

This figure clearly shows that an approximate value of the cross section for such neutrinos is of the order of magnitude of  $10^{-38} \text{ cm}^2 \text{ GeV}^{-1}$ , which is typically tens of orders of magnitude lower than the interaction cross section of a photon ( $\sigma_\gamma \sim 10^{-25} \text{ cm}^2$  [49]). This means that the typical neutrinos of a few MeV produced by nuclear reactors have a mean free path of approximately 10 light years in steel.

- Neutrinos are also the only SM particle only observed in their left-handed chirality state, while anti-neutrinos can only be observed in their right-hand state. This could mean two things: either right-handed neutrinos do not exist in nature for some reason, or we have just not been able to detect them because their interaction with baryonic matter is too weak.

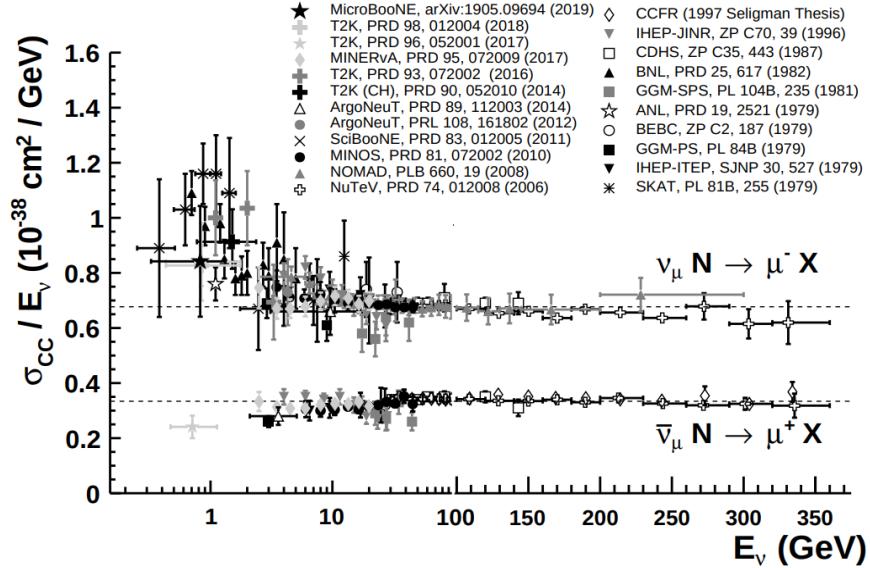


Figure 2.9: Neutrino cross section of interaction from the charged current as measured by different experiments over a large range of energies, for both neutrinos  $\nu$  and antineutrinos  $\bar{\nu}$  [48].

Right-handed neutrinos, also referred to as *sterile* neutrinos, do not fit in the current SM but could actually also be a strong BSM DM candidate, as we will discuss in the next subsection.

However, two physical reasons can explain why we do not really believe that DM could be made of neutrinos any more. First of all, their relative abundance does not match the expected one for DM from the freeze-out mechanism explained in Section 2.3. Indeed, their freeze-out abundance can be computed quite easily from Equation 2.8 [50], where the sum of the masses of the three neutrino flavors has been calculated to be lower than 0.17 eV [51] instead of the 11.5 eV expected to obtain the correct DM relic abundance as observed today from the power spectrum of the CMB.

$$\Omega_\nu h^2 = \sum_{i=1}^3 \frac{m_{\nu_i}}{93 \text{ eV}} \quad (2.8)$$

Additionally, for several reasons explained in Section 2.3, a good DM candidate is expected to be cold or in other words, non-relativistic. However, being extremely light, neutrinos are expected to be ultra-relativistic and could therefore not be responsible of the emergence of large scale structures as observed today: we can therefore most probably rule out the possibility of DM being made of SM neutrinos.

### Sterile neutrinos

The most obvious SM particles being ruled out as DM candidates, we can now introduce other BSM theories introducing additional particles that could have the properties searched for. The first one of these theories introduces the so-called sterile neutrinos, usually referring to right-handed SM neutrinos, as discussed in the previous subsection.

If they exist, sterile neutrinos are expected to interact in an even weaker way than SM active neutrinos, they could be very long-lived as well and in principle, nothing prevents us from considering that they could have a mass superior to 0.4 keV, giving therefore the correct DM relic abundance

[44]. A superior bound of 50 keV on such particles can also be obtained considering limits on the observation of the monochromatic decay  $\gamma$  line originating from the one loop radiative decay  $N \rightarrow \nu + \gamma$  of such particles.

Several experiments are already searching for this kind of particles at this level of energy. Most of these experiments focus on the analysis of  $\gamma$ -rays and are actually searching for this particular monochromatic line resulting of the decay of sterile neutrinos. Two independent groups actually announced in 2014 the observation of an unidentified emission line at 3.57 keV (Figure 2.10) which did not match any known atomic emission line and which is actually consistent with an eventual DM signal [52, 53], since most of the possible instrumental contamination effects have been ruled out over the course of the last few years.

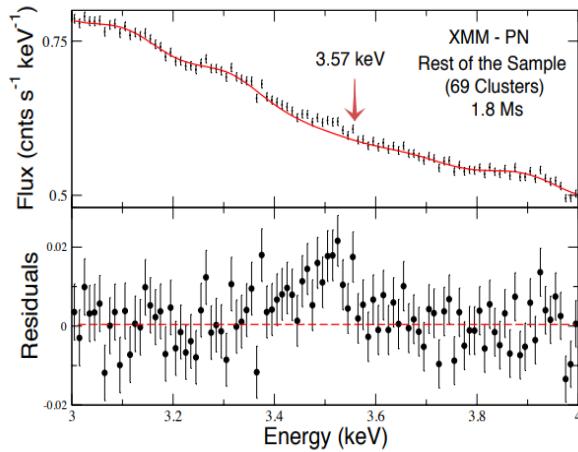


Figure 2.10: 3.57 keV emission line detected with a  $4.5\sigma$  CL by the XMM-Newton telescope in 2014, which could be a hint of the presence of DM [52].

However, some additional studies of the galactic center pointed out the fact that this observation might actually come from the observation of a potassium K XVIII transition line [54]. Recent observations actually ruled out at the 99% CL an eventual DM origin for this particular line [55], but further studies are still ongoing.

## Axions

Axions could also explain the particle nature of DM, since their existence is enough to explain 100% of the DM in the Universe, unlike most of the other candidates presented so far. Axions are hypothetical stable neutral particles, with masses of the order of the meV, introduced as a consequence of the strong-CP violation issue of Quantum ChromoDynamics (QCD). This issue is the following: the usual  $\Theta$  term of the QCD Lagrangian, the QCD vacuum angle shown in Equation 2.9 [56] should be responsible of breaking the CP symmetry, but this effect has actually never been observed so far: this is the so-called the strong CP problem.

Two ways to explain that we have never observed this phenomena exist: the first is to assume that one of the quarks of the SM is massless but this does not match the current observations and measurements. The second consists in assuming that the parameter  $\Theta$ , the QCD vacuum angle, is small enough so that this term becomes negligible. However, by definition, the  $\Theta$  angle should be between 0 and  $2\pi$ , so there is no physical reason for this parameter to be small, unless some new physics can be introduced, such as the theory developed in 1977 by Peccei and Quinn [57] by

relaxing  $\Theta$  from a parameter to a dynamic variable and absorbing it through the introduction of a new pseudoscalar particle that was called the axion.

$$\mathcal{L}_\Theta = \frac{\Theta}{32\pi^2} \epsilon_{\mu\nu\rho\sigma} G_a^{\mu\mu} G_a^{\rho\sigma} \quad (2.9)$$

By definition, it is possible to show that axions satisfy two of the previous criteria for a good DM candidate, since they are non-relativistic and their abundance might be enough to account for the dark matter energy density observed, because their actual abundance can easily be computed from Equation 2.10 [58], from which we could conclude that an axion having a mass of  $\sim 20$   $\mu$ eV could account for the DM relic density of the Universe, as observed today.

$$\Omega_a \simeq \left( \frac{6\mu\text{eV}}{m_a} \right)^{7/6} \quad (2.10)$$

Several axion search experiments have therefore been set up, such as the Axion Dark Matter Experiment (ADMX), a resonant microwave cavity installed at the University of Washington, the CERN Axion Solar Telescope (CAST), a CERN experiment observing the Sun which came online in 2002 and which managed in 2014 to turn up definitely the existence of solar axions [59], or the International AXion Observatory (IAXO), whose aim will be to search for axions with a much better signal to ratio noise than CAST. All the results obtained by these experiments along with their future projections are represented in Figure 2.11.

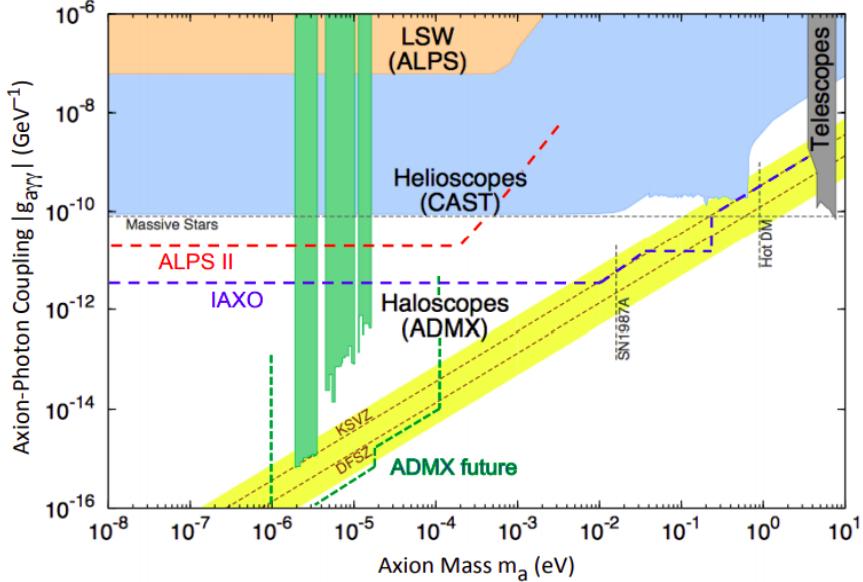


Figure 2.11: Axions exclusion summary plot and projected coverage of axion searches experiments, such as ADMX, CAST and IAXO [58].

### Weakly Interactive Massive Particles (WIMPs)

The actual DM candidates that will be mostly considered throughout this work are the so-called WIMPs, which are expected to interact, even though very weakly, with ordinary baryonic matter and which have an expected mass in the range of 100 GeV to 1 TeV for reasonable electroweak production cross-section values, right where we expect DM to be found from its relic density: this is the so-called "WIMP miracle", an important concept that can be translated mathematically

as well. Indeed, because of the freeze-out scenario explained in Section 2.3, we can find an expression relating the relic abundance of DM  $\Omega_\chi$  with its annihilation cross section  $\langle\sigma_A v\rangle$  through Equation 2.11 [67].

$$\Omega_\chi h^2 \sim \frac{3 \cdot 10^{-27} \text{ cm}^3 \text{ s}^{-1}}{\langle\sigma_A v\rangle} \quad (2.11)$$

This equation implies that, since we do know the current abundance of DM in the Universe, the total annihilation cross section of DM should be equal to  $\sim 3 \cdot 10^{-27} \text{ cm}^3 \text{ s}^{-1}$ , which corresponds to the typical value given by WIMPs for a range of dark matter masses matching the expected one.

Several strategies can be used to detect such particles, as we will see in Section 2.5. This kind of particle basically arises in various BSM theories, such as the Lightest Supersymmetric Particle (LSP) in SUSY. According to this theory, each SM particle should have a superpartner whose quantum numbers would be identical except for their spins, which would differ by one half. All of these superpartners would then be potentially new and undiscovered particles, giving us a perfect DM candidate in most of the Minimal Supersymmetric Standard Model (MSSM) theories, the neutralino  $\chi$  [60].

The WIMPs are also interesting in the sense that introducing them in the terms of this supersymmetric theory would not only give us a strong DM candidate, but would also solve the hierarchy problem, the apparent large discrepancy among the masses of the different SM particles ranging several orders of magnitude.

## 2.5 Dark matter searches

As previously stated, several cosmological evidences allow us to introduce the concept of dark matter, but its properties such as its mass, coupling and interaction cross-section are difficult to study in this context. Several different ways can then be used to try and detect DM particles in order to study them, as represented in Figure 2.12, strategies which can usually be divided into three categories: the direct and indirect searches, mostly relying on the production of baryonic matter from the interaction between two DM particles or on the observation of the interaction between the dark and baryonic sectors, and the production in colliders, usually able to probe lower DM candidates masses and which will actually be the main focus of this work.

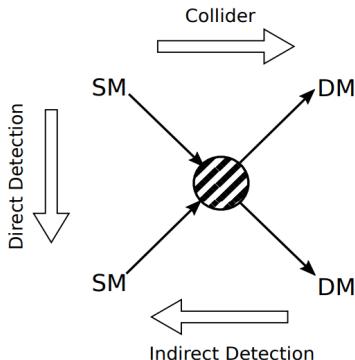


Figure 2.12: Schematic view of the three main DM detection strategies: direct, indirect and collider production searches [61].

### 2.5.1 Direct searches

From cosmological observations, we know that we live in a halo of dark matter. In this case, WIMPs should cross the Earth every day and, even if they interact only weakly, we should be able to directly detect them through their interaction with ordinary baryonic matter, for example because of their scattering with the nuclei of the atoms. Indeed, the transfer of momentum between these two particles in this case might be detectable with the correct experimental device, typically placed deep underground to have the lowest possible background, which is the main source of perturbations of such experiments.

To study this particular category of searches, let's first of all introduce the rate of expected WIMP scattering off a target nucleus of mass  $m_N$  with Equation 2.12, rate which ends up being described by a simple steeply falling exponential function as shown in Figure 2.13 [62], where  $E_{nr}$  is the nuclear recoil energy measured,  $m_\chi$  is the WIMP mass,  $\sigma$  its cross section,  $\rho_0 = 0.3 \text{ GeV cm}^{-3}$  is the local dark matter density and  $f(v)$  the normalized WIMP velocity distribution.

$$\frac{dR}{dE_{nr}} = \frac{\rho_0 M}{m_N m_\chi} \int_{v_{\min}}^{v_{\text{esc}}} v f(v) \frac{d\sigma}{dE_{nr}} dv \propto \exp\left(-\frac{E_{nr}}{E_0} \frac{4m_\chi m_N}{(m_\chi + m_N)^2}\right) \quad (2.12)$$

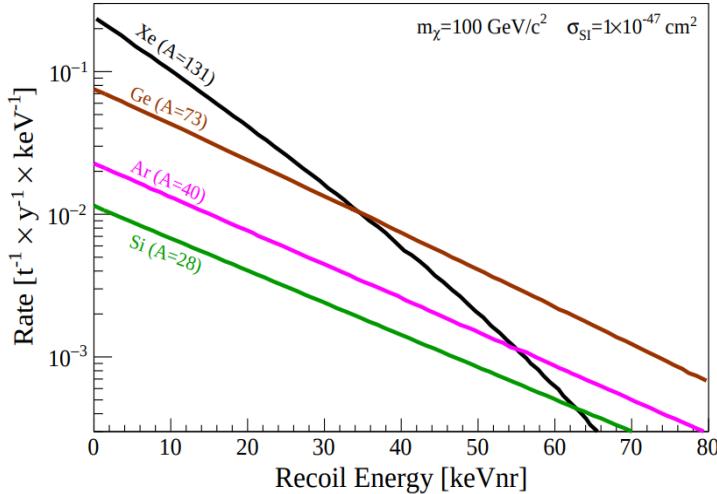


Figure 2.13: Nuclear recoil spectra induced in different materials for a given DM WIMP of 100 GeV, assuming a WIMP-nucleon Spin Independent (SI) cross section [62].

From this relation, the Equation 2.13 can be easily derived, representing this time the number of expected DM events in an experiment running during a time  $T$ , where  $\epsilon(E_{nr})$  is the efficiency of the detector for a given recoil energy.

$$N = T \int_{E_{\min}}^{E_{\max}} \epsilon(E_{nr}) \frac{dR}{dE_{nr}} dE_{nr} \quad (2.13)$$

The maximal velocity  $v_{\text{esc}}$  used as superior bound of the integral in Equation 2.12 has actually been measured to be in the range [498 – 608] km/s at the 90% CL [63], since any particle having a velocity higher than this would not be bound any more to the gravitational potential of a galaxy. This has an important consequence: all the direct and indirect detection experiments actually need to take into account the annual modulation of the observed count rate, due to the movement

of the Earth around the Sun, as shown in Figure 2.14 [64], since this velocity is not negligible compared to the escape velocity  $v_{\text{esc}}$ .

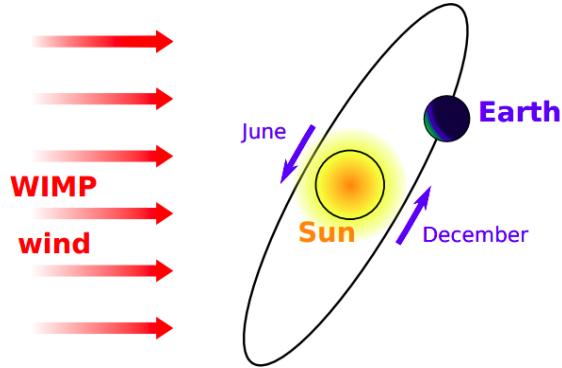


Figure 2.14: Schematic representation of the annual modulation of the WIMP wind introduced by the motion of rotation of the Earth around the Sun [64].

From our perspective, it seems indeed that the velocity of the speed of WIMP particles arriving is changing depending on the month of the year, since the Earth is sometimes moving in the direction of the WIMP source, and is sometimes moving away: the maximal velocity is reached around June. It is extremely important to take into account this effect since, as we saw on the previous equations such as Equation 2.11, the count rate of incoming particles  $N(t)$  actually depends on this velocity, and this modulation then introduces a periodical modulation that we need to take into account, as shown in Equation 2.14, where the periodical part usually introduces a  $\sim 5\%$  deviation [62].

$$N(t) = B + N_0 + N_m \cos(\omega(t - t_0)) \quad (2.14)$$

This effect is also important because an experiment performed during a long period of time can actually help us finding an eventual hint of DM particles, since our signal is expected to follow this periodical deviation while the background is expected to be constant. Moreover, this WIMP wind is expected to come from a particular region of the sky while the backgrounds are expected to be distributed uniformly, so this gives a clear way to isolate the signal.

Finally, it is also important to note that two different kinds of direct searches can be defined, depending on the category of the scattering between the DM and the nucleus: the Spin Independent (SI) (proceeding through the scalar term) and Spin Dependent (SD) (proceeding through the axial term of the Lagrangian) searches, since the interaction cross section  $\sigma$  of Equation 2.12 is expected to be different for DM particles having a spin 0 or not, as shown in Equation 2.15, where  $F$  is a factor accounting for the dependence of the scattering on the energy. This means that results obtained by either hypothesis cannot be compared with each other.

$$\frac{d\sigma}{dE_{nr}} \propto \sigma_{SI} F_{SI}^2(E_{nr}) + \sigma_{SD} F_{SD}^2(E_{nr}) \quad (2.15)$$

As previously stated, many experiments are dedicated to the direct search of DM particles, but in order to isolate an eventual DM signal, an environment with an ultra-low background is usually required, which is usually reduced either by placing the detector underground (to reduce the contamination due to cosmic rays), by increasing the statistics (by repeating the experiment and/or observing for a larger amount of time) or by choosing carefully the active material of the detector

(to reduce the internal background coming from the detector itself). The impact these kind of parameters have on the final limits can be seen in Figure 2.15.

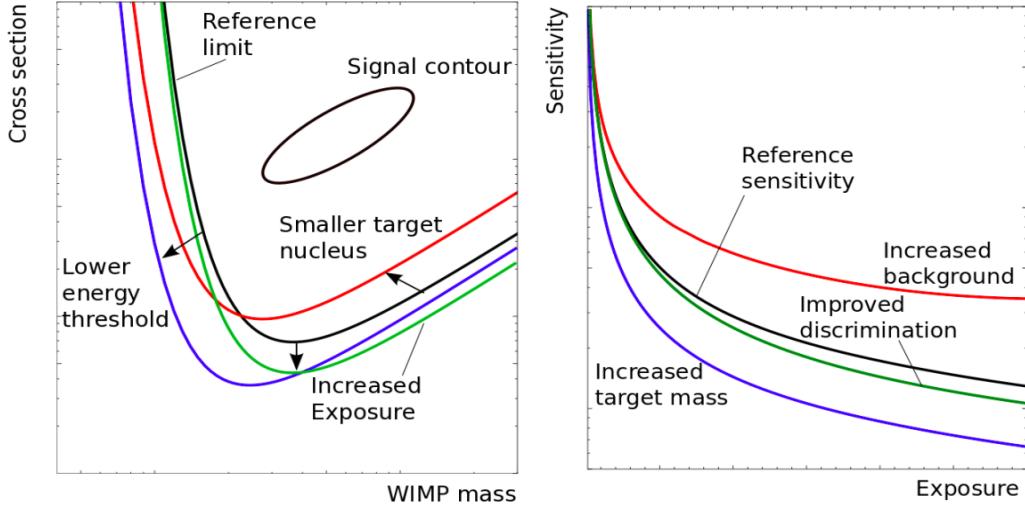


Figure 2.15: Impact of different experimental parameters on the final limits depending on the cross section and WIMP mass (on the left) or sensitivity and exposure (on the right), with respect to the expected limits (black curve) [65].

These detectors try to detect the scattering of an unknown exotic particle with an ordinary nucleus, which typically can give rise to different categories of signals. Some detectors try for example to detect the direct ionization of the target atom, while others focus on the emission of light coming from the de-excitation of the scattered nucleus, and some even search for the heat produced by these collisions under the form of phonons (collective excitation phenomena in condensed matter) in a crystal. These different search strategies have been summarized schematically in Figure 2.16.

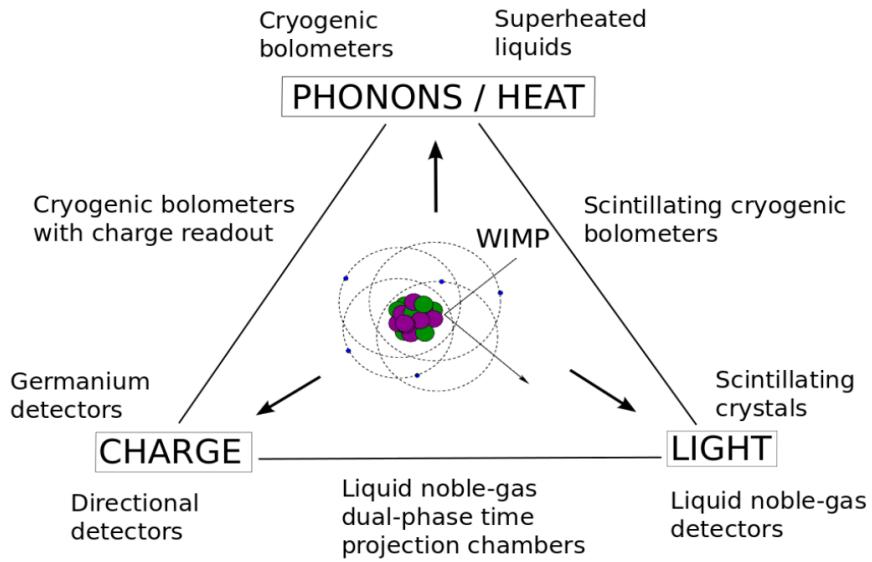


Figure 2.16: Schematic representation of the three main strategies to detect directly the interaction between DM particles and an ordinary nucleus [65].

As of today, no direct experiment has been able to detect serious hints for the existence of DM, and they have only been able to set limits on the scattering cross section depending on the models parameters, as seen in Figure 2.17 for multiple experiments at once.

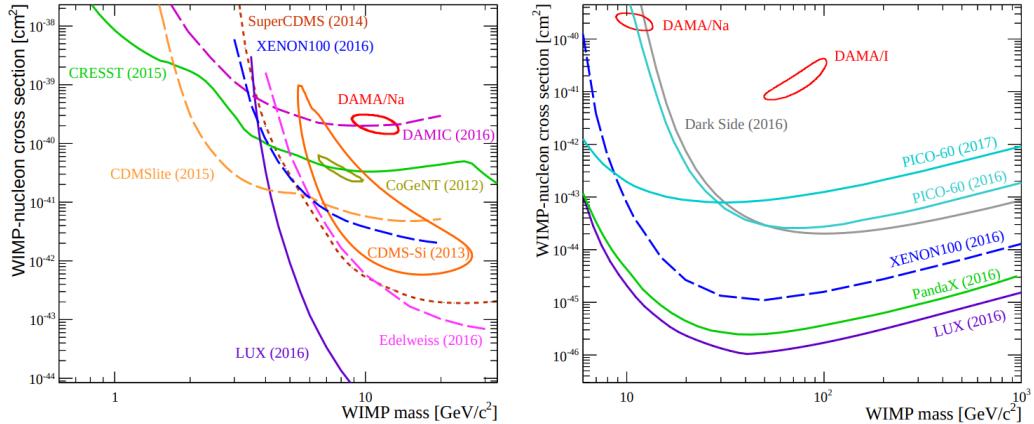


Figure 2.17: Exclusion limits obtained by various direct detection experiments considering a SI interaction cross section for low WIMP (on the left) or high WIMP masses (on the right) [65].

However, the DAMA experiment at the Laboratori Nazionali del Gran Sasso (LNGS) did find an interesting result by showing the hints of an annual modulation signal compatible to the expected one due to the WIMP wind in the 2-6 keV energy range, as seen on Figure 2.18 [66]. Further investigations about this modulation are still ongoing today, since no systematic effect able to account for the observed modulation amplitude and to simultaneously satisfy all the requirements of the signature has been found so far.

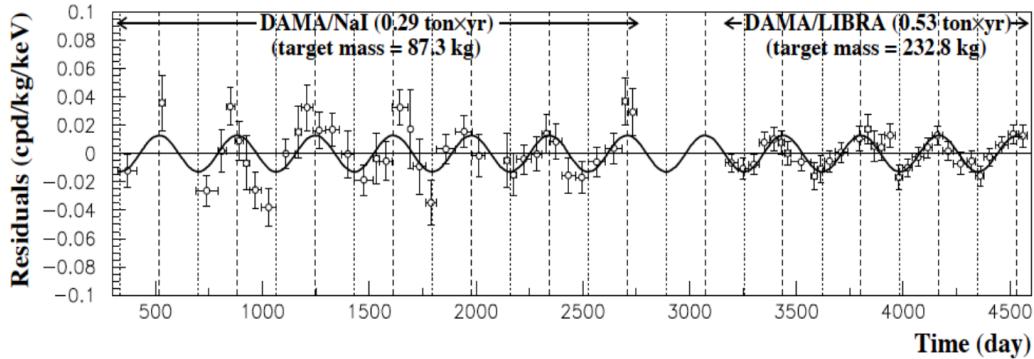


Figure 2.18: Observed and expected annual modulation in single hits events in the 2-6 keV energy range by the DAMA experiment [66].

## 2.5.2 Indirect searches

The indirect detection of DM particles consists basically in searching for SM products coming from the annihilation of two DM particles or from its eventual direct decay, usually under the form of a flux of  $\gamma$ -rays, neutrinos, cosmic-rays or anti-matter appearing as an excess over the expected background. Indeed, many extensions to the SM, such as the supersymmetry SUSY or Universal Extra Dimensions (UED) do provide solid DM candidates (the lightest supersymmetric particle and the lightest Kaluza-Klein state, respectively) expected to interact with each other, resulting in the immediate production of SM (un)stable particles that can be detected by telescopes either placed on the ground or directly in space. Another point to make is that indirect searches are also usually affected by the annual fluctuation induced by the movement of the Earth around the Sun, as explained in Section 2.5.1.

Indirect searches are actually extremely useful since they are sensitive to the DM annihilation cross section, mass and the density profile of DM halos  $\rho(\vec{r}(s, \Omega))$ , usually represented by a Navarro-Frenk-White (NFW) profile, as shown in Equation 2.16, where  $r_s = 20$  kpc is the scale radius of the DM halo for the Milky Way and  $r$  is the distance to the center of the cluster considered, assumed to be spherical in this case [68].

$$\rho(r) \propto \frac{r_s}{r \left(1 + \left(\frac{r}{r_s}\right)^2\right)} \quad (2.16)$$

The flux coming from the annihilation of two DM particles is then expected to be proportional to its annihilation cross section  $\sigma v$ , the solid angle of observation  $\Omega$  and the number of particles emitted by this annihilation  $\frac{dN}{dE}$ , according to Equation 2.17, where the integration is done over the line of sight  $l$  and the solid angle of observation  $\Omega$ .

$$\frac{d\Phi}{d\Omega dE} = \frac{\sigma v}{8\pi m_\chi^2} \cdot \frac{dN}{dE} \iint_{l, \Omega} \rho^2(\vec{r}(s, \Omega)) dl d\Omega \equiv P \cdot J(\Delta\Omega) \quad (2.17)$$

This equation is extremely important for two reasons. First of all, it shows that, if a signal is found in direct detection, we could use this detection to determine the new object mass and scattering cross section and then use this information in order to obtain the DM density profile this way: the data obtained by the different strategies of detection are actually complementary. The second reason is that as we can see, the flux of incoming particles can actually be divided into two factors:  $P$ , entirely dependent on the physics of the DM particle, and the  $J$ -factor  $J(\Delta\Omega)$ , depending only on the distribution of DM within the system considered. This  $J$ -factor is in this sense actually a measurement of the quality of an astronomical object for an indirect measurement, since the higher the flux received, the better the measurement will be in general (even though this is not the only factor which matters, since for example the galactic center has a higher  $J$ -factor than the best dwarf galaxy observable, but also has a lot of backgrounds affecting the measurement).

As different channels of observation are available for us to analyze the eventual annihilation of DM, several strategies can be used in order to detect DM indirectly, by searching different kinds of SM particles. Anyway, all these strategies have one goal in common: try to reduce the background, since the signal searched for is usually quite low while the uncertainties associated to the background in astrophysics are usually quite high.

### Through $\gamma$ -rays detection

The golden channel for such searches is through the production of  $\gamma$ -rays by DM annihilation  $\chi\chi \rightarrow \gamma + X$  or decay  $\chi \rightarrow \gamma + X$ , mainly because the energy scale of the WIMPs implies that most of the annihilation and decay radiation will be emitted in this range of energies and because  $\gamma$ -rays are usually not deflected when traveling to the observer (this means that the exact source of this kind of radiation can be quite easily and precisely pin-pointed). However, they do have one drawback as well: the Earth's atmosphere is usually opaque to this kind of radiation at this level of energy. This means that most of experiments searching for them simply cannot be performed from the ground and have to be sent to space.

One of the most famous detectors in this category is the Fermi Large Telescope (LAT), a pair production detector launched in June 2008 and mostly sensitive to  $\gamma$ -rays between 20 MeV and 300 GeV [69]. This experiment has managed to exclude a large portion of the phase space, as seen in Figure 2.19. The GAMMA-400 experiment, whose launch is scheduled in 2020, will pick

up the work of LAT, by studying a similar range but with much improved angular and energy resolutions [70].

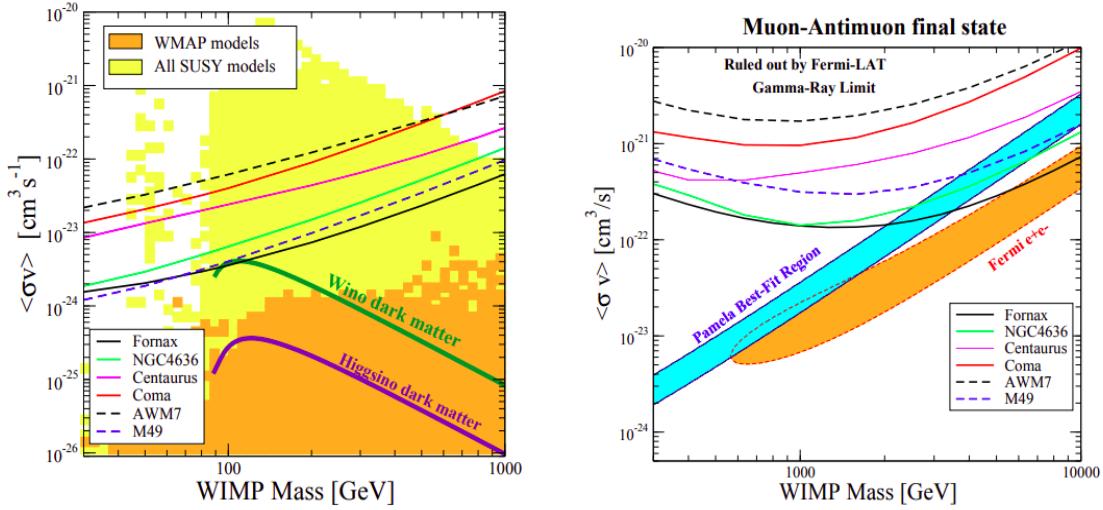


Figure 2.19: Upper limits on the DM annihilation cross section considering  $b\bar{b}$  (on the left) and  $\mu^+\mu^-$  (on the right) final states as a function of the WIMP mass, for different clusters studied [69].

It is much harder for DM to produce high energy  $\gamma$ -rays and therefore, the flux of such particles decreases quite quickly with energy, making it harder to study since telescopes then need a much larger effective area to pick up the same quantity of signal, and have therefore to be put on the ground. Such telescopes do exist, are called Imaging Atmospheric Cherenkov Telescopes (IACT) and have to take into account the atmospheric perturbations to work in an optimal way. They are usually sensitive to a range of energies going from  $\sim 10$  GeV up to  $\sim 100$  TeV, but can usually only study a small portion of the sky (up to a few degrees), forcing these experiments to choose carefully the objects to be studied. The Cherenkov Telescope Array (CTA) is a brand new telescope of this kind whose construction is supposed to start in 2020, that should improve greatly the sensitivity of high masses indirect DM searches.

### Through neutrinos detection

Interacting only weakly, neutrinos are another reliable source of data in the Universe since they are not supposed to be altered when traveling large distances as well, even though detecting neutrinos is much harder than detecting  $\gamma$ -rays and usually involves huge tanks of water in which neutrinos can be detected with the Cherenkov effect, which consists of the emission of an electromagnetic radiation when a charged particle moves through a dielectric medium with a speed greater than the speed of light in this medium.

The most famous detector of this kind, the IceCube neutrino observatory, actually uses the ice of the South Pole instead of water to detect these particles with photo-detectors, mainly because of the low interaction cross section of the neutrinos which then requires the installation of a huge volume of material to increase the probability of interaction. Super-Kamiokande (Super-K), in Japan, is another large Cherenkov experiment dedicated to the detection of cosmic neutrinos. Both detectors are also largely involved as direct searches experiments, since they are also sensitive to the eventual recoil between DM and ordinary matter nuclei.

The problem with this kind of experiment is the difficulty of actually detecting some neutrinos, along with the background levels from atmospheric neutrinos, as represented in Figure 2.20, typ-

ically several orders of magnitude larger than the signal. Multiple strategies therefore need to be put in place in order to reduce the background level in such experiments, such as the study of the directionality of the source and an appropriate choice of angle of observation, since most of the contamination is coming from tau neutrinos, themselves originating from muon neutrinos oscillation, strongly suppressed around the zenith.

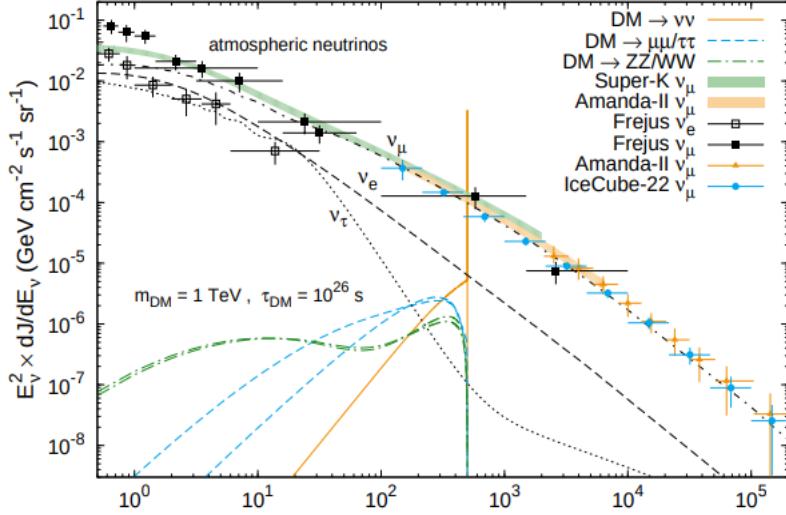


Figure 2.20: Neutrino spectra for a scalar DM candidate of 1 TeV for different indirect detection experiments and the corresponding background level expected [71].

### Through cosmic rays detection

Searching for anti-matter in the cosmic rays presents the advantage of being highly sensitive, because of the low levels of backgrounds this kind of searches implies. However, cosmic rays are affected when traveling through the Universe, and determining the exact location of their emission can be quite challenging.

Among the most famous detectors of this kind, we can quote PAMELA, a spatial telescope dedicated to the study of such cosmic rays since 2006. CERN's Alpha Magnetic Spectrometer (AMS), installed in the International Space Station has also studied such radiation from a range of a few hundred MeV up to 1 TeV. The data collected by this detector is compared to the exclusion limits obtained by the IceCube detector in Figure 2.21.

### 2.5.3 Collider production

In this particular kind of searches, we are interested in the eventual direct production of DM candidates following the collision between two highly energetic SM particles, a perfectly viable scenario if we keep assuming that DM should at least interact weakly with ordinary matter if we want to be able to produce or detect it in a laboratory.

In this case, two main models for the interaction between SM and DM particles can be considered, each with a different level of complexity and different possible applications:

- The Effective Field Theory (EFT) approach is usually considered to be the easiest, even though it can still give us plenty of information about this kind of interaction. It relies

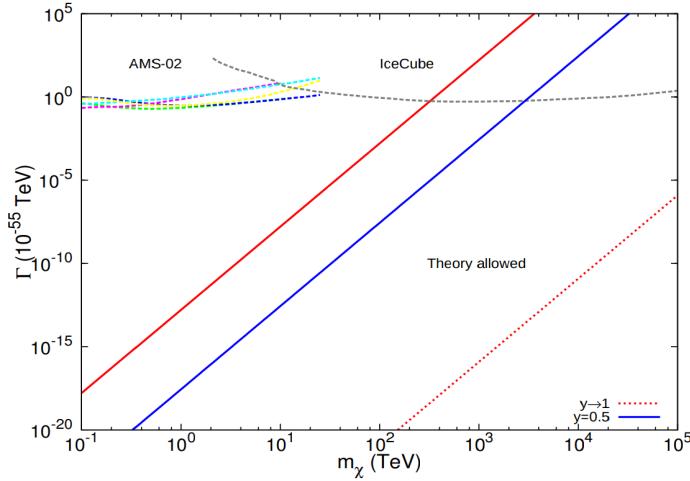


Figure 2.21: Limits of the decay width of the interaction with respect to the DM mass obtained by both IceCube and AMS [72].

the assumption that the energy scale of the new exotic physics is much larger than the actual energy accessible with our experiment, since the momentum transfer of this interaction needs to be much smaller than the mediator mass in this case. According to this approach, represented in Figure 2.22, the interaction can be described using simplified operators (the most simple ones being either scalar  $\bar{\chi}\chi\bar{q}q$ , pseudoscalar  $\bar{\chi}\gamma^5\chi\bar{q}\gamma^5q$ , vector  $\bar{\chi}\gamma^\mu\chi\bar{q}\gamma_\mu q$  or axial-vector  $\bar{\chi}\gamma^\mu\gamma^5\chi\bar{q}\gamma_\mu\gamma^5q$ ) since, according to the assumption made, its mediator can usually be integrated out [61].

Despite of this strong assumption, this approach is actually quite useful anyway in the sense that it is able to provide us bounds on the new physics scale  $\Lambda$ , which can be related to the different couplings of the interaction as in Equation 2.18, where  $g_q$  and  $g_\chi$  are the coupling between the mediator and the SM or DM particle, and  $m_{\text{med}}$  is the mass of the mediator.

$$\frac{1}{\Lambda^2} = \frac{g_q g_\chi}{m_{\text{med}}} \quad (2.18)$$

Additionally, direct searches experiments typically also introduce this kind of assumptions to extract constraints from their measurements, making it straightforward to compare the results obtained in both approaches. However, the transfer of momentum in the direct searches experiments is usually of the order of a few keV as seen in Equation 2.19, while in collider experiments such as the LHC, this is of the order of  $\Theta(\text{GeV-TeV})$ .

$$E = \frac{m_\chi v_\chi^2}{2} \simeq 100 \text{ GeV} \cdot 10^{-6} \simeq 50 \text{ keV} \quad (2.19)$$

This means that, especially due to the increasing center of mass energy given by the LHC over the last few years, the basic EFT assumption is usually not respected and gives information about an out of reach phase space region anyway, making its usefulness quite relative in most cases. This is why alternative models need to be developed as well.

- The simplified models attempt to solve the issue related to the approximation made by the EFT approach, by increasing the level of details regarding the interaction between the dark and baryonic sectors. This is usually done by explicitly taking into account the mediator of the interaction, which can be considered of two different types in the context of this work: either scalar  $\phi$  or pseudoscalar  $a$  (both having a spin 0), described by the Lagrangians in Equation 2.20, considering the DM candidate to be a Dirac fermion coupling to the SM

through the mediator considered and under the assumption of Minimal Flavour Violation (MFV) (a proposal made to characterize the effects of flavor transitions in new theories of particle physics). In this equation, the sum runs over the three SM families and the parameters  $y_i^f = \sqrt{2} \frac{m_i^f}{v}$  are the Yukawa couplings, much larger for the top quarks because of their mass, which will allow us to simplify the following equations [73].

$$\begin{cases} \mathcal{L}_{\text{fermion},\phi} \propto -g_\chi \phi \bar{\chi} \chi - \frac{\phi}{\sqrt{2}} \sum_i (g_u y_i^u \bar{u}_i u_i + g_d y_i^d \bar{d}_i d_i + g_l y_i^l \bar{l}_i l_i) \\ \mathcal{L}_{\text{fermion},a} \propto -ig_\chi a \bar{\chi} \gamma^5 \chi - \frac{ia}{\sqrt{2}} \sum_i (g_u y_i^u \bar{u}_i \gamma^5 u_i + g_d y_i^d \bar{d}_i \gamma^5 d_i + g_l y_i^l \bar{l}_i \gamma^5 l_i) \end{cases} \quad (2.20)$$

An important parameter in this case is the decay width of this mediator  $\Gamma_{\text{med}}$ , given by Equation 2.21 for either a scalar mediator  $\phi$  or Equation 2.22 for a pseudoscalar mediator  $a$ , where the first term corresponds to the mediator decay to SM particles, the second to its decay to DM particles and the last term its possible decay to gluons.

$$\begin{cases} \Gamma_\phi = \sum_f N_C \frac{y_f^2 g_\nu^2 m_\phi}{16\pi} \left(1 - \frac{4m_f^2}{m_\phi^2}\right)^{3/2} + \frac{g_\chi^2 m_\phi}{8\pi} \left(1 - \frac{4m_f^2}{m_\phi^2}\right)^{3/2} + \frac{\alpha_S^2 g_\nu^2 m_\phi^3}{32\pi^3 \nu^2} \left| f_\phi \left( \frac{4m_t^2}{m_\phi^2} \right) \right|^2 \\ f_\phi(\tau) = \tau \left( 1 + (1 - \tau) \arctan^2 \left( \frac{1}{\sqrt{\tau - 1}} \right) \right) \end{cases} \quad (2.21)$$

$$\begin{cases} \Gamma_a = \sum_f N_C \frac{y_f^2 g_\nu^2 m_a}{16\pi} \left(1 - \frac{4m_f^2}{m_a^2}\right)^{1/2} + \frac{g_\chi^2 m_a}{8\pi} \left(1 - \frac{4m_f^2}{m_a^2}\right)^{1/2} + \frac{\alpha_S^2 g_\nu^2 m_a^3}{32\pi^3 \nu^2} \left| f_a \left( \frac{4m_t^2}{m_a^2} \right) \right|^2 \\ f_a(\tau) = \tau \arctan^2 \left( \frac{1}{\sqrt{\tau - 1}} \right) \end{cases} \quad (2.22)$$

In the case of the simplified models, the minimal set of parameters describing the interaction is therefore  $\{m_\chi, m_{\text{med}}, g_\chi, g_q\}$ .

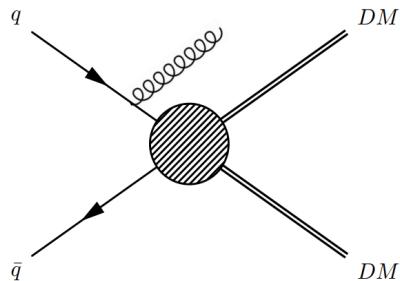


Figure 2.22: Schematic representation of a typical EFT modelization of an LHC event with an Initial State Radiation (ISR) object used to trigger the event [61].

An additional categorization of the DM production models can be done, by separating the so-called s-channel and t-channel models, as shown in Figure 2.23. In the first case, the most common one in collider searches, the mediator between the SM and DM is assumed to be a boson, and usually decays directly into a pair of DM particles. On the other hand, in the t-channel models, the mediator couples to one quark and one DM particle, with a colored exchange particle required.

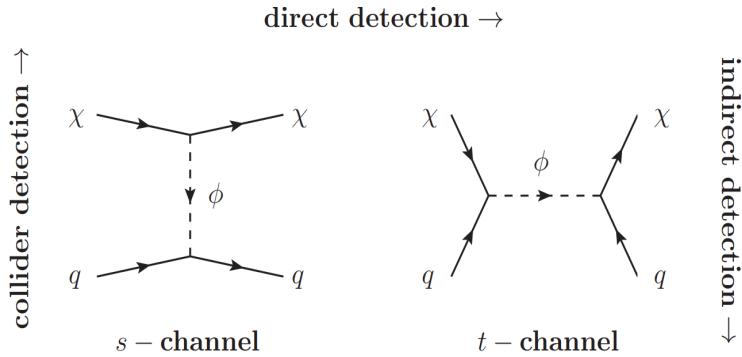


Figure 2.23: Schematic representation of a typical collider DM production through *s*-channel and *t*-channel processes [74].

In this work, simplified models following the ATLAS-CMS Dark Matter Forum recommendations [75], with either scalar or pseudoscalar mediators will always be considered because of their relative simplicity and the lack of strong assumptions behind them.

## 2.6 Dark matter production at the Large Hadron Collider (LHC)

The LHC, colliding protons at a center of mass energy of 13 TeV, is actually a perfect machine to study such processes, because of the expected range of masses (100 GeV to 1 TeV) for the best DM candidates, as discussed in Section 2.3. However, because of the weak interactions of such particles, they are not expected to interact at all with the detector, which will then basically search for missing transverse momentum along with a SM particle triggering the event.

Several major categories of DM searches exist at the LHC, performed mainly by the CMS and ATLAS collaborations, depending on the strategy applied for such searches:

- First of all are the so-called **mono-X searches**, where X stands for a detectable SM particle used to trigger the event while the DM mediator usually decays into a pair of particles escaping the detector, leaving behind some MET, a key variable to all these searches that will be described in Section 4.5. Depending on the nature of the X particle, several searches can be performed: we can for example mention the mono- $\gamma$  [76, 77] or mono-jet [78, 79] 13 TeV searches performed by the ATLAS and CMS collaborations.

Additionally, if the DM mass is high enough, a coupling to the Higgs boson is also possible, or an eventual decay of the Higgs itself to a couple of DM particles  $H \rightarrow \chi\bar{\chi}$  is kinematically not impossible. This channel, known as the mono-Higgs, is sensitive to all the decays of the Higgs, even though the searches excluding most of the phase spaces are performed using the  $H \rightarrow b\bar{b}$  and  $H \rightarrow \gamma\gamma$  decays [80, 81].

- The searches for DM production in **association with heavy quark(s)** belong to the second category, such as the one performed in this work. Other analyses such as the one probing the  $b\bar{b} + \text{DM}$  to its different final states depending on the decay of the bottom quark (at 8 or 13 TeV) also belong to this group [16, 20].

These searches have to combine the discriminant power of several variables to separate the signal from the backgrounds, since the MET distribution on its own is usually not enough,

but they present the advantage of being favored by the higher Yukawa coupling to more massive particles implied by the MFV assumption.

- **Dijet searches** provide the best exclusion limits for most of the spin-1 mediated models considered (up to a few TeV for typical coupling choices) [82, 83]. In this case, the sensitivity is obtained by searching for either narrow or large resonances on the exponentially falling QCD background, while the other searches were mostly dedicated in searching for bumps in the MET spectrum.
- **Multi-cascade searches**, involving the production of large cascade chains in which a stable dark matter particle is produced at the end. This kind of signature is very popular in Supersymmetry models, which not only solve the problem of DM but also the hierarchy problem while giving us perfect DM candidates such as the neutralino  $\chi$ , defined as the lightest stable supersymmetric particle, obtained in many of the MSSM theories, having an hypothetical mass below the TeV scale [84].
- **Searches with Higgs driven states**: several models such as the Higgs portal to the dark sector is another interesting strategy. In some specific cases, when considering spin-0 mediated interactions between the dark and baryonic sectors, the Higgs could be considered as the mediator of the interaction as well, which only requires a minimal modification of the SM Lagrangian. It is then necessary to study the different Higgs production modes, such as the gluon fusion and Vector Boson Fusion (VBF) mechanisms, to search for an eventual invisible decay of the Higgs into a couple of DM particles, assuming that its mass is lower than  $\sim 62.5$  GeV,  $m_H/2$  [85].
- Finally, and this is quite new, **long-lived searches** can also be performed. These are interesting because they can extend the current searches performed to also consider the eventual creation of long-lived particles which would decay a few centimeters further than the primary vertex of the  $pp$  collision [86]. Typical SM signatures do not usually include such events, making this channel relatively background-free, even though the reconstruction of the different objects is much harder in this case.

All the results from the different searches performed by the CMS collaboration using the full  $(35.9 \pm 0.9) \text{ fb}^{-1}$  dataset collected at 13 TeV can be summarized in Figure 2.24 for spin-0 mediators and in Figure 2.25 when considering spin-1 mediators.

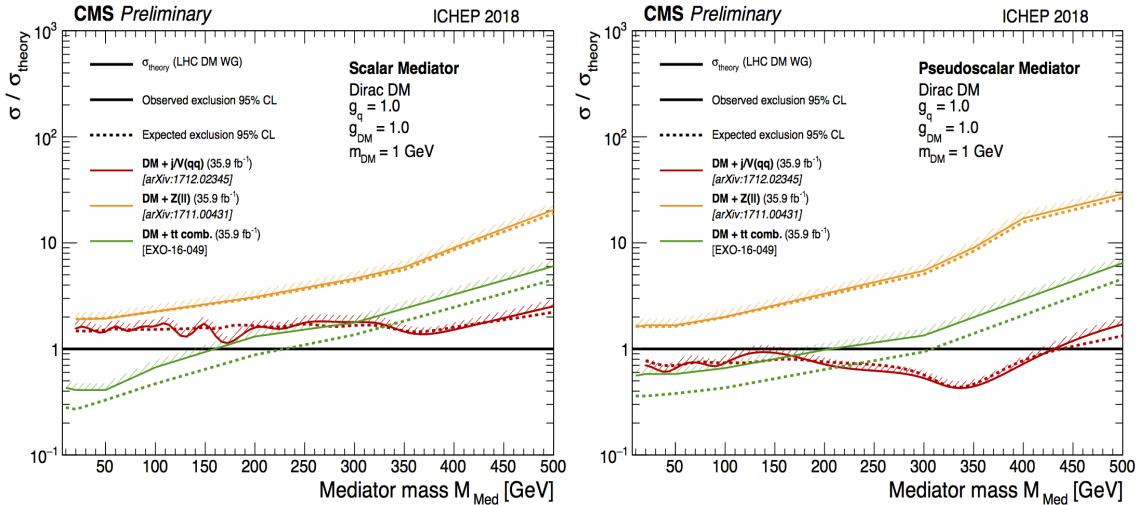


Figure 2.24: Observed and expected 95% exclusion limits obtained by different searches of the CMS collaboration as a function the spin-0 scalar (on the left) or pseudoscalar (on the right) mediator.

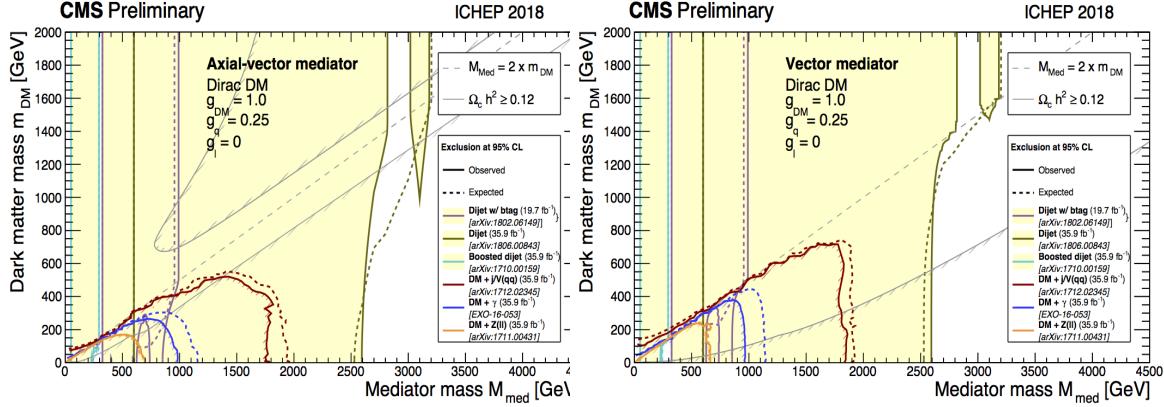


Figure 2.25: Observed and expected 95% exclusion limits obtained by different searches of the CMS collaboration as a function of the spin-1 mediator, considering axial-vector (on the left) and axial (on the right) interactions.

These results have also been compared to the results obtained by several direct detection experiments in Figure 2.26, considering both the Spin Dependent (SD) and Spin Independent (SI) cases, as explained in Section 2.5.1.

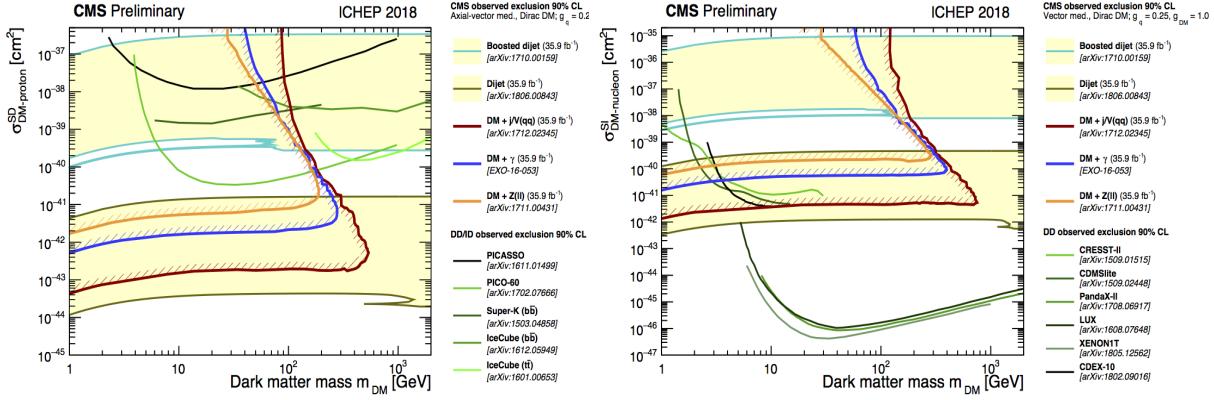


Figure 2.26: CMS 90% exclusion limits compared to the most famous direct detection experiments for the SD (on the left) and SI (on the right) scenarios, obtained using similar couplings.

This particular analysis and our signal of interest, the search for DM production together with a single or a pair of top quarks falls into the second category. It is already important to note that this analysis has been performed considering the DM candidate to be a Dirac fermion with all the couplings equal and set to 1, as recommended by the Dark Matter Working Group (DMWG) [87].

## 2.7 The focus of this thesis

This thesis focuses on a particular search for DM, produced in association with either one or two top quarks, considering only the dilepton final state of such processes. In particular, two different categories of signal, later referred to as the  $t/\bar{t}+\text{DM}$  and  $t\bar{t}+\text{DM}$  processes, will be considered and combined in order to compute upper limits on the DM production signal strength.

### The single top production channel

The first kind of signal considered in this analysis is the production of DM in association with a single top quark, known as the  $t/\bar{t}$ +DM analysis. This process is expected to be mediated by a spin-0 mediator, either scalar or pseudoscalar, and is associated with a light quark and a W boson (the mono-top analysis is dealing with the case where a single top is created without any additional particle).

Three different Feynman diagrams can be associated to this particular analysis depending on the model considered, as shown in Figure 2.27.

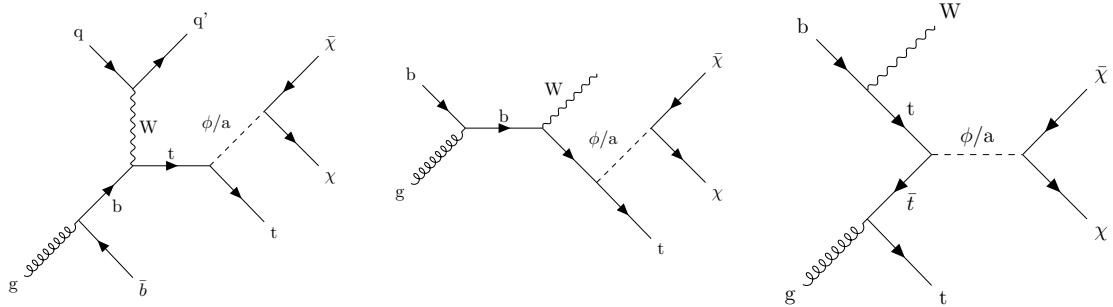


Figure 2.27: Feynman diagrams involving the production of DM with a single top quark according to its t-channel W boson (on the left), or tW (on the center and on the right) production modes.

As discussed previously, the top quark will be dynamically favored due to its high mass and therefore high Yukawa coupling, but this has another consequence as well, since the lifetime of this quark is extremely low, of the order of  $10^{-15}$  s [88]. This means that this particle usually decays before being able to form hadrons, so we can only detect the products of its decay, not the top quark itself. In almost 100% of the cases, the top actually decays into a bottom quark and a W boson, which decays itself before being detected into quarks and/or leptons. Even though this will be detailed in Chapter 6, we can already conclude that the typical final state of such signature is therefore made out of MET coming from the DM, one b-tagged jet along with one or two W bosons, seen as a combination of jets and leptons, depending on the channel considered.

### The $t\bar{t}$ production channel

The  $t\bar{t}$ +DM analysis is really similar, except that in this case, we have two top quarks in the final state, as represented in Figure 2.28.

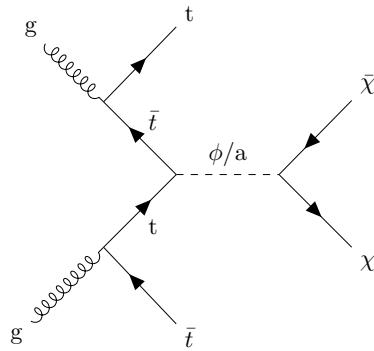


Figure 2.28: Schematic representation of a typical  $t\bar{t}$ +DM event.

The final state expected by such an event is therefore made of two b-tagged jets, along with two

W bosons and some MET coming from the pair of DM particles. In this case as well, both scalar  $\phi$  and pseudoscalar  $a$  mediators will be considered, with a mass range from 10 to 500 GeV.

### The dilepton final state

As previously explained, most of the models considered will produce exactly two W bosons, which are not stable long enough and therefore decay before reaching the detector, meaning that we cannot directly detect such bosons. However, we can detect the results of the decays of these W, since they can decay either to hadrons ( $67.6 \pm 0.27\%$  branching ratio) or to a lepton and a neutrino, giving us an additional contribution of MET ( $10.8 \pm 0.09\%$  for each lepton) [88].

This means that this kind of analyses featuring two W bosons in the final state can focus on different channels: either completely hadronic (if both the W decay into quarks), semileptonic or dileptonic (when both W bosons decay into two neutrinos and two leptons of opposite charge).

Given the Branching Ratio (BR) of the W decay, it is easy to see that the dileptonic channel, which will be the focus of this work, is not favored statistically. However, this channel features less backgrounds than other channels, resulting in a better signal isolation, and because leptons can usually be reconstructed in a better way than jets (cf. Chapter 4), leading to lower uncertainties and improved limits.

## 2.8 Previous relevant results

This analysis being performed at a center of mass energy  $\sqrt{s} = 13$  TeV, only the most relevant results to this energy obtained by both the CMS and ATLAS collaborations will be quoted.

First of all, the **ATLAS collaboration** published interesting results at this center of mass energy, considering an integrated luminosity of  $13.3 \text{ fb}^{-1}$ , and obtained the corresponding exclusion limits at the 95% CL, considering the  $t\bar{t}+\text{DM}$  model and both scalar and pseudoscalar spin-0 mediators for the interaction, as shown in Figure 2.29. In this case, and for the couplings considered, an exclusion up to  $\sim 375$  GeV has been achieved [17].

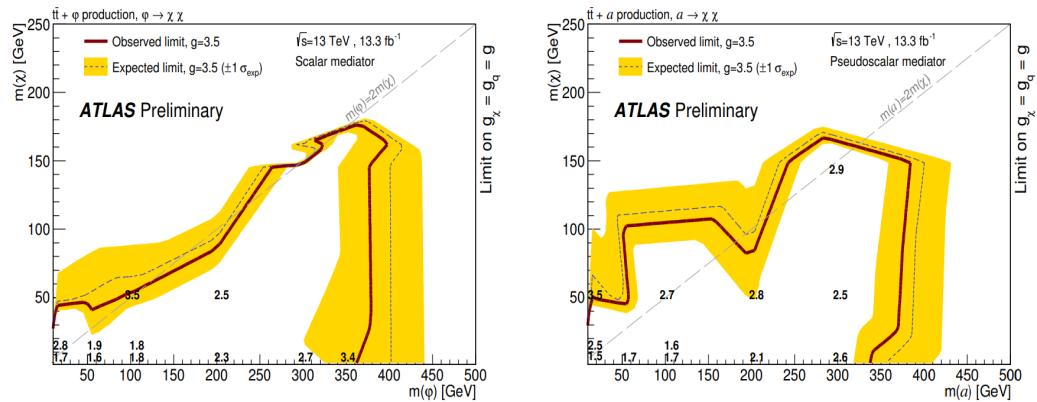


Figure 2.29: Limits on the DM and mediator masses obtained by ATLAS using  $13.3 \text{ fb}^{-1}$  of 13 TeV data, considering scalar (on the left) and pseudoscalar (on the right) mediators [17].

Considering the full 2016 dataset of  $36.1 \text{ fb}^{-1}$ , similar results have been obtained, as shown in

Figure 2.30. In this case, and for lower coupling values, an exclusion up to around 100 GeV has been obtained considering scalar and pseudoscalar mediators [20].

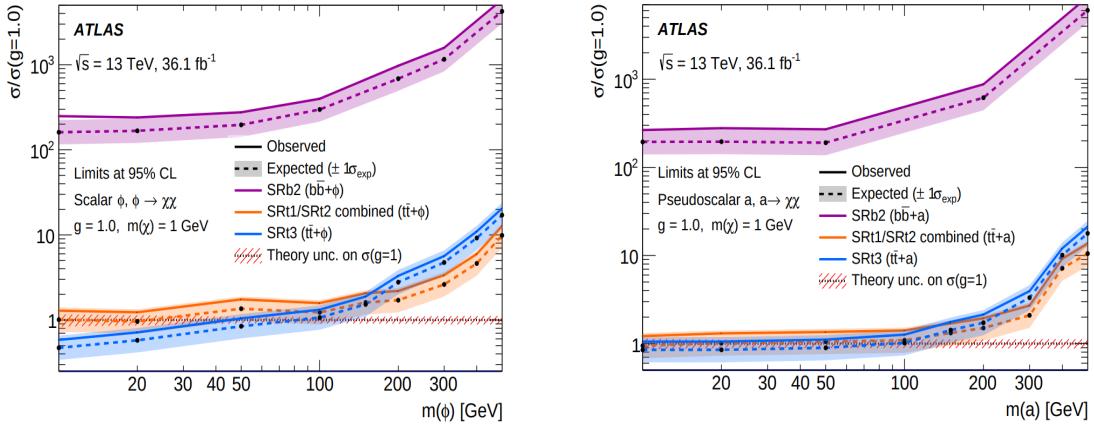


Figure 2.30: Exclusion limits at the 95% CL obtained by ATLAS considering scalar (on the left) and pseudoscalar (on the right) mediators, for a DM mass of 1 GeV [20].

Finally, ATLAS recently presented in ICHEP 2020 the first results for a combined search of stop and for a similar final state in the dilepton final state, and using the full Run II dataset, as shown in Figure 2.31 [21]. According to these latest results, the collaboration has now been able to achieve an exclusion of scalar (pseudoscalar) mediators with masses up to 250 (300) GeV.

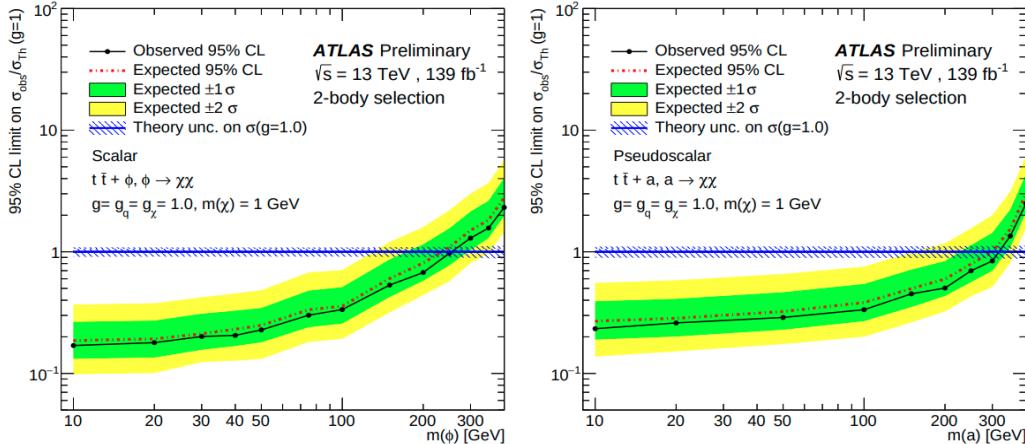


Figure 2.31: Exclusion limits at the 95% CL obtained by ATLAS considering scalar (on the left) and pseudoscalar (on the right) mediators, for a DM mass of 1 GeV, considering the dilepton final state of such model and the full Run II dataset [20].

Finally, recent summary plots for ATLAS dark matter searches, showing the exclusion limits obtained for the three possible final states have been published, as shown in Figure 2.32. The collaboration did not perform any combination between the different channels though.

The **CMS collaboration** published in 2018 a similar analysis, using the  $35.9 \text{ fb}^{-1}$  of data collected during the year 2016. This analysis combined the three different final states possible (hadronic, semileptonic and dileptonic) and computed the limits on the signal strength for different mediator and dark matter masses, considering both scalar and pseudoscalar mediators [23]. The results obtained can be found in Figure 2.33.

Last year, a CMS combination of the  $t/\bar{t}+\text{DM}$  and  $t\bar{t}+\text{DM}$  analyses has also been published [24], combining this time only the hadronic and semileptonic channels of both analyses. The limits obtained in this case are represented in Figure 2.34, where the limits obtained by each analysis

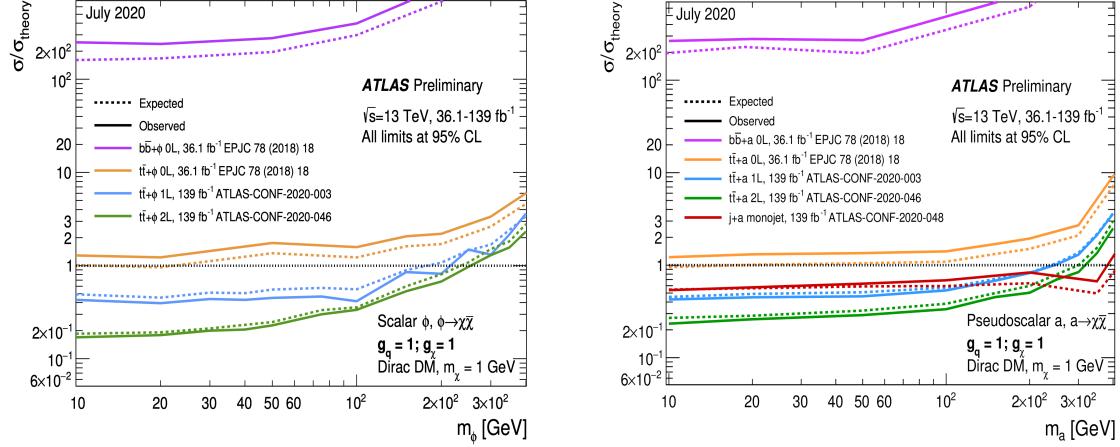


Figure 2.32: Exclusion limits at the 95% CL obtained by ATLAS considering scalar (on the left) and pseudoscalar (on the right) mediators, for a DM mass of 1 GeV [21].

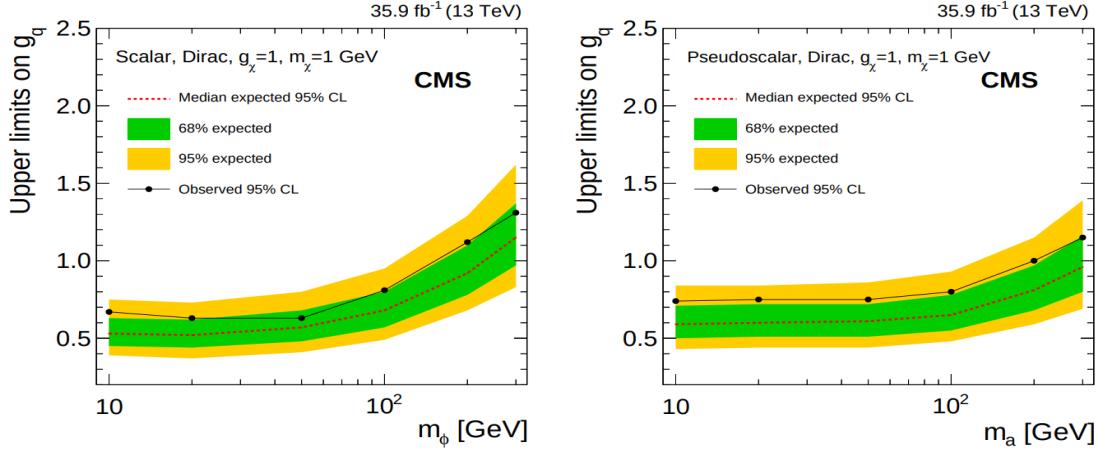


Figure 2.33: 95% CL exclusion plots on the signal strength computed as a function of the mediator and DM masses obtained by CMS considering a scalar (on the left) and a pseudoscalar (on the right) mediator for the interaction [23].

on their own along with the results of the combination, leading to a factor 2 improvement of the limits obtained, have been represented.

This combination managed to exclude the production of scalar mediators up to 290 GeV and pseudoscalar mediators up to 300 GeV, at the 95% CL for the couplings considered. This combined analysis actually leads to the most stringent exclusion limits of the LHC on the production of DM through these categories of spin-0 mediators.

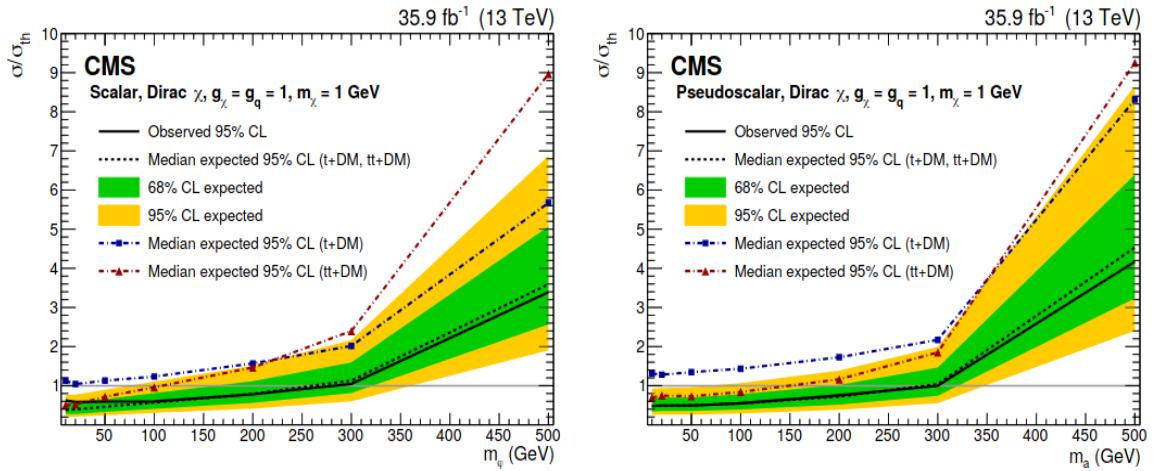


Figure 2.34: Expected and observed 95% CL limits on the DM production cross sections shown considering scalar (on the left) and pseudoscalar (on the right) mediators for the interaction [24].



---

---

# Chapter 3

---

## The experimental setup

The data analyzed throughout this work is the result of several years of proton-proton collisions produced at the LHC with a center of mass energy of 13 TeV, and the particles emerging from these collisions have been recorded with the CMS detector. In order to provide the experimental context of this thesis, the Section 3.1 of this section will be dedicated to the description of the accelerator, while the detector, in which thousands of people were involved and made of several different layers, will be described in Section 3.2.

### 3.1 The Large Hadron Collider (LHC)

The Large Hadron Collider (LHC) is a superconducting particle accelerator able to accelerate protons and lead ions up to velocities close to the speed of light ( $0.99999999c$ ). Planned since the end of the 20th century, this accelerator, a 27 kilometers ring put 100 meters underground (under France and Switzerland) to avoid part of the contamination due to the cosmic rays and located at CERN, has now been running for 10 years. The LHC is the result of the collaboration between thousands of scientists of more than 100 different nationalities and its main objective was first of all to either infer or confirm the possible existence of the Higgs boson, theoretically predicted in the 1960s [2, 3] but never observed in any experiment. The discovery of the Higgs boson, announced at CERN on the 4th of July 2012 [4, 5], was then a magnificent achievement of the accelerator, after only a few years of operation.

Now that the Higgs boson has been discovered, the priority of the LHC shifted a bit. Even though many different teams are still studying this particle in order to determine precisely its most fundamental properties such as its mass, couplings or spin, many groups of scientists are involved in different kinds of BSM physics since the LHC, whose center of mass energy has kept increasing over the years, allows us to reach a level of energy never reached before and therefore allows us to probe new parts of the phase space, searching for eventual hints of new physics. These kind of searches of course include the search for DM production as the one performed in this work.

### 3.1.1 The LHC in a nutshell

The LHC is an underground particle accelerator built in the same 27 kilometers tunnel where the Large Electron Positron collider (LEP) was previously used [89]. This machine is accelerating two beams made out of  $10^{10}$  protons or  $7 \cdot 10^7$  lead ions in each direction and is mostly made out of more than 4000 superconducting magnets, in majority dipoles and quadrupoles, allowing respectively to curve the beam to maintain a nominal circular trajectory and to focus it by compensating its natural dispersion due to the repulsion of the protons making up these beams. A dedicated small section of the accelerator is then made out of 16 radio-frequency cavities synchronized in such a way that these protons always face a negative electric charge, which is used as the driving process of the actual acceleration of such particles.

Once the nominal center of mass energy  $\sqrt{s} = 13$  TeV is reached (this concept will be described in Subsection 3.1.2), the protons are then smashed together in four different places on the LHC, where the four detectors (ATLAS, CMS, A Large Ion Collider Experiment (ALICE) and LHCb) have been placed in order to study the collisions. Both ATLAS and CMS are general detectors able to study exotic processes such as the production of DM and to make precision measurements on SM physics as well (the decision to build two separate detectors was made in order to introduce some redundancy and to check the results). ALICE on the other hand is mostly dedicated to the study of heavy ions collisions that happen  $\sim 10\%$  of the time in order to study in particular a specific state of matter, called quark-gluon plasma [90]. Finally, LHCb has been designed to study in particular the CP violation phenomena, which could be the sign of some new physics [91].

It is important to note at this point that the LHC is not a standalone accelerator in the sense that protons enter the LHC with a velocity already close to the speed of light. In order to reach such input energies ( $\sim 450$  GeV), previous smaller accelerators of CERN are still used today. A chain of accelerators is then formed: first of all, the protons, extracted from a bottle of ordinary hydrogen, are injected into the LINAC 2, a linear accelerator, before being transferred to the Proton Synchrotron (PS), the Super Proton Synchrotron (SPS) and finally the LHC itself (all this chain of acceleration can be found in Figure 3.1).

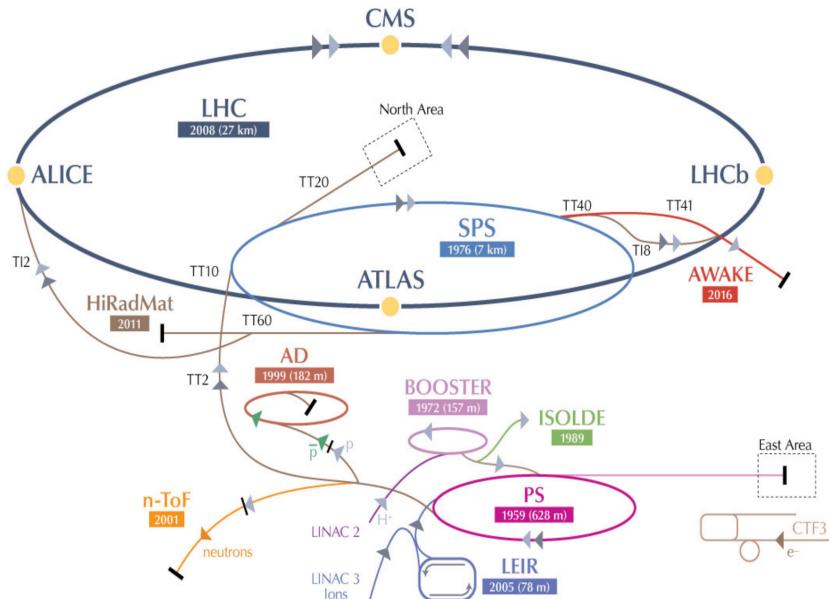


Figure 3.1: LHC injection chain and experiments performed at CERN [93].

During this phase of acceleration, the beam is separated in 2808 bunches of protons nominally separated by 25ns (giving a collision rate of 40 MHz), so that the experiments have the time to

Parameter	Run I	Run II	Run III	Design
Energy [TeV]	7 → 8	13	13	14
Bunch spacing [ns]	50	25	25	25
Intensity [ $10^{11}$ protons per beam]	1.6	1.2	Up to 1.8	1.15
Bunches	1400	2500	2800	2800
Emittance [ $\mu\text{m}$ ]	2.2	2.2	2.5	3.5
$\beta^*$ [cm]	80	30 → 25	30 → 25	55
Crossing angle [ $\mu\text{rad}$ ]	-	300 → 260	300 → 260	285
Peak luminosity [ $10^{34} \text{ cm}^{-2} \text{ s}^{-1}$ ]	0.8	2.0	2.0	1.0
Peak PU	45	60	55	25

Table 3.1: Expected and observed main parameters of  $pp$  operation of the LHC across the different eras of operation [92].

record the collision, clean the detectors and get ready for the next bunch crossing, coming just a few nanoseconds later. Each time one bunch crossing happens, around 30-35 collisions of protons happen at once, on average: this phenomena, usually referred to as Pile up (PU), has to be taken into account as well and will be described in the next section.

As previously stated, the amount of data collected is a crucial parameter for many analyses, meaning that the LHC would ideally need to run 24 hours a day, 365 days per year in order to maximize the data taking. However, this is typically not possible, as setting up a beam takes time and we cannot keep it rotating at maximal energy in the machine for a long time so in ideal conditions, around 20 hours of data taking a day are expected. The LHC is then running around 9 months per year, being usually stopped during winter for maintenance operations, and the data taking periods are defined into Runs of a few years, after which the accelerator is usually stopped for a longer period of time, a Long Shutdown (LS), in order to also have the time to perform additional upgrade operations of the machine.

The data analyzed in this work corresponds to the second phase of the Run II of operation of the LHC, from 2016 to 2018, while the Run III is now expected to start in the Spring of 2021. The summary of the main parameters of operation of the LHC across the different Runs of operations can be found in Table 3.1.

### 3.1.2 Key parameters

#### Center of mass energy

The center of mass energy is defined as a Lorentz invariant quantity under any kind of boost resulting of the collisions between two protons (defined as  $E_1, \vec{p}_1, m_1$  and  $E_2, \vec{p}_2, m_2$ ) with a  $\theta$  angle, as developed in Equation 3.1. It is a key variable of the LHC since the phase space of particles that can be explored directly depends on this value.

$$\sqrt{s} = \sqrt{(m_1)^2 + (m_2)^2 + 2(E_1 E_2 - 2|\vec{p}_1| |\vec{p}_2| \cos(\theta))} \quad (3.1)$$

The LHC started its operation in 2008 running at an energy of 7 TeV, quickly moved to 8 TeV

and kept this level of energy during the end of the Run I of operation. In 2015, for the Run II of data taking, the energy was increased until reaching 13 TeV (2 times 6.5 TeV for each beam) and an expected value of 14 TeV, the nominal energy for which the LHC was originally built, is expected to be reached in the near future.

## Luminosity

The luminosity is another extremely important variable for the operation of the LHC since it gives an indication on the number of collisions per second given by the accelerator. Increasing this parameter is then crucial to collect as much data as possible, in order to be able to isolate processes having a low production cross section and therefore an extremely low probability of creation when colliding two protons.

Mathematically, we can first of all define in Equation 3.2 the rate of production  $R$  (in number of events per second) of any given process using its cross section  $\sigma$ , equivalent to its production probability, and the instantaneous luminosity  $\mathcal{L}$ . From this rate can be extracted easily the number of expected interactions  $N$  in a certain amount of time  $T$  as well.

$$\begin{cases} R = \mathcal{L} \cdot \sigma \\ N(T) = \sigma \int_0^T \mathcal{L}(t) dt = \sigma L \end{cases} \quad (3.2)$$

This instantaneous luminosity we just introduced  $\mathcal{L}$  can be defined using Equation 3.3, assuming that the beams have a Gaussian profile, while the integrated luminosity  $L$  can be defined by simply integrating the instantaneous luminosity over time  $L = \int \mathcal{L} dt$  [94].

$$\mathcal{L} = \frac{\gamma f_{\text{rev}} k_B N_p^2}{4\pi\epsilon_n \beta^*} G, \text{ where } G = \frac{1}{\sqrt{1 + \frac{\theta_C \sigma_z}{2\sigma}}} \quad (3.3)$$

In this last equation, the following properties of the accelerator have been introduced, giving the LHC a nominal instantaneous luminosity  $\mathcal{L} = 10^{34} \text{ cm}^{-2} \text{ s}^{-1}$ :

- $\gamma$  is the usual relativistic Lorentz factor (for the LHC, when protons go at their maximum velocity of 99.9999991% of the speed of light,  $\gamma = 7460$ )
- $f_{\text{rev}}$  is the frequency of revolution (11.2 kHz)
- $k_B$  is the number of proton bunches per beam (2808 for a 25ns bunch spacing)
- $N_p$  is the number of protons per bunch ( $1.15 \cdot 10^{11}$  protons)
- $\epsilon_n$  is the transverse normalized emittance ( $3.75 \mu\text{m}$ )
- $\beta^*$  is the betatron function at the interaction point (0.55 m)
- $G$  is the geometrical factor accounting for the fact that the collisions do not happen exactly head-on, therefore reducing the effective luminosity, expressed from the full crossing angle between colliding beam  $\theta_C$  (285  $\mu\text{rad}$ ), and  $\sigma$ ,  $\sigma_z$ , the transverse and longitudinal r.m.s. sizes (respectively 16.7  $\mu\text{m}$  and 7.55 cm).

As previously stated, increasing the luminosity of the LHC is always something interesting in order to produce processes having a low production cross-section, and we can then see that many different

parameters can be tweaked in order to achieve the highest possible instantaneous luminosity, such as the number of protons per bunch, the number of bunches per beam or the beam crossing angle at the interaction point. New radio-frequency crab cavities will probably be installed during the next LS of the LHC in order to increase the value of the geometrical factor  $R$  and the instantaneous luminosity by a factor  $\sim 10$  (HL-LHC project [95]).

The total integrated luminosity taken by the LHC during its different years of operation has been summarized in Figure 3.2. The final datasets available for the Run II and analyzed in this work have an integrated luminosity of  $(35.9 \pm 0.9) \text{ fb}^{-1}$  (2016),  $(41.5 \pm 1.0) \text{ fb}^{-1}$  (2017) and  $(59.7 \pm 1.5) \text{ fb}^{-1}$  (2018), resulting in a total dataset of  $(137.1 \pm 2.0) \text{ fb}^{-1}$  recorded during the Run II of operation. This roughly means that we expect to have produced around 137 events of any process whose cross section of production would be equal to 1  $\text{fb}$ .

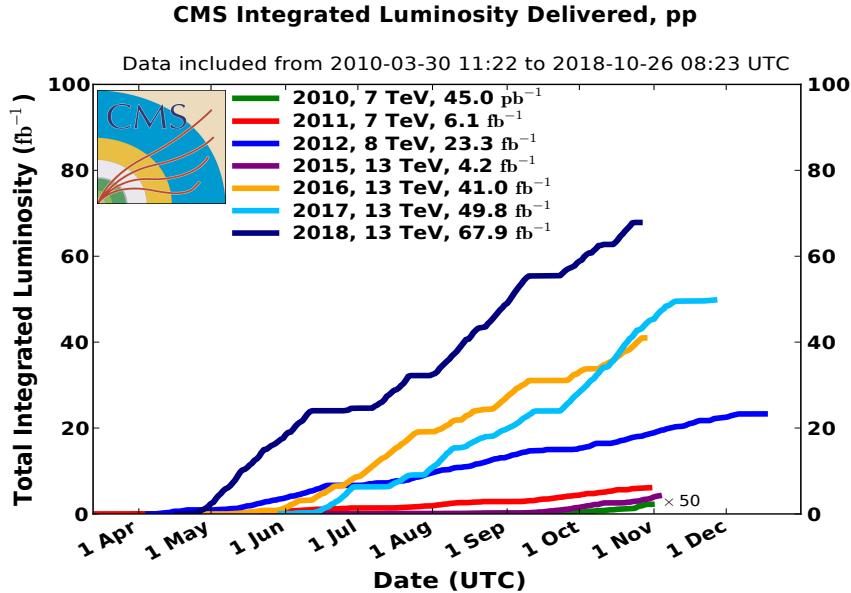


Figure 3.2: Integrated luminosity collected by CMS over its different years of operation so far.

### Number of simultaneous interactions: Pile up (PU)

The last key parameter of the LHC discussed here is the number of simultaneous interactions, also called PU. Usually, because of the high density of protons within the beams, a bunch crossing in an experiment produces around 30–35 proton collisions, as seen in Figure 3.3, defining the Primary Vertex (PV) as the most interesting one, while the other vertices are usually referred to as the PU. The tracker of CMS, which will be introduced in Section 3.2.1, therefore needs to be able to reconstruct all these events in order to define the PV of the interaction.

## 3.2 The CMS detector

The Compact Muon Solenoid (CMS) is one of the two general purpose detectors of the LHC and is installed at the access point 5 of the LHC. Its main purposes are to discover the Higgs boson, make precision measurement of its properties and also of other SM processes and to discover BSM physics, such as the possible existence of DM.

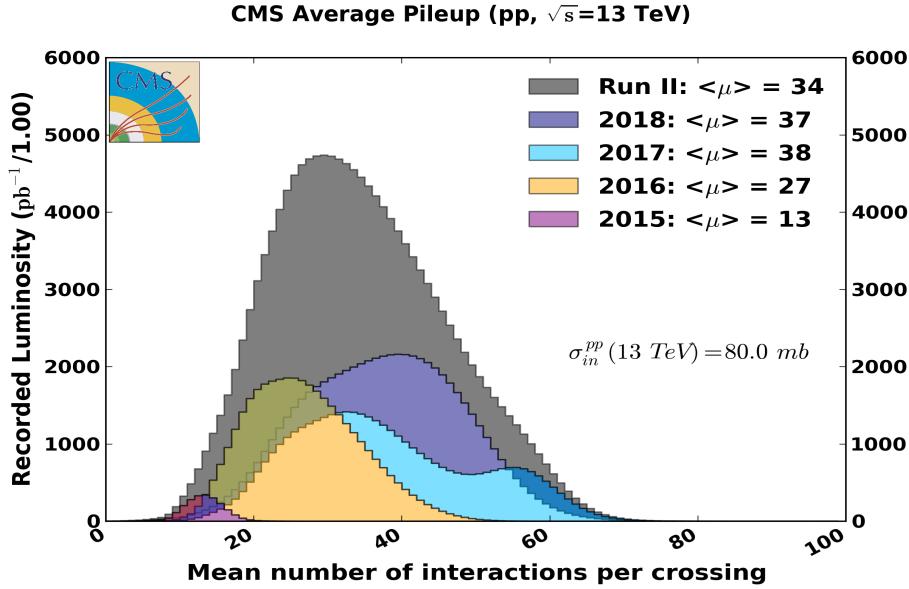


Figure 3.3: Mean PU distribution and luminosity recorded by CMS over the different years of operation of the LHC.

CMS has been carefully designed by hundreds of different physicists and engineers in order to be as hermetic as possible, covering all the possible angles around the beam pipe. It therefore counts with a central part with cylindrical shape, the barrel, and two endcaps, one on each side of the detector. With its 14.000 tons, distributed over a diameter of 15 meters and a length of 28.7 meters, the detector can be considered quite compact and was lowered into the experimental cavern after being built on the surface in 14 different moving pieces, a flexible design allowing to access its inner parts, by opening and closing the detector when needed.

### CMS subdetectors

As shown in Figure 3.4, CMS is made out of different layers corresponding to different sub-detectors, each able to provide different kinds of information about the particles created by each collision [96]. These sub-detectors will be described in detail in the following sections, but they have all been designed in order to make the reconstruction of the different events efficient, fast and precise while matching the tight conditions occurring at the LHC.

The inner part of the CMS detector is the so-called tracker, a device made out of silicon pixels and strips, described in Section 3.2.1 and responsible for the precise reconstruction of all the charged particles coming from the different interaction vertices. A bit further, the Electromagnetic Calorimeter (ECAL) can be found, made out of thousands of crystals as described in Section 3.2.2 and used to precisely measure the energy of particles able to interact electromagnetically by producing an electromagnetic shower that can be detected. Then, the Hadronic Calorimeter (HCAL) comes, described in Section 3.2.3 and whose main purpose is to identify and measure the main properties of the hadrons produced in the collisions.

The CMS name partly comes from its central piece, a huge superconducting solenoid described in Section 3.2.4 and able to produce a 3.8 T magnetic field in the detector, with a magnetic flux density increased even more by the addition of the steel return yoke layers (red parts in Figure 3.4). This magnetic field is essential in the sense that it is able to deflect the charged particles which

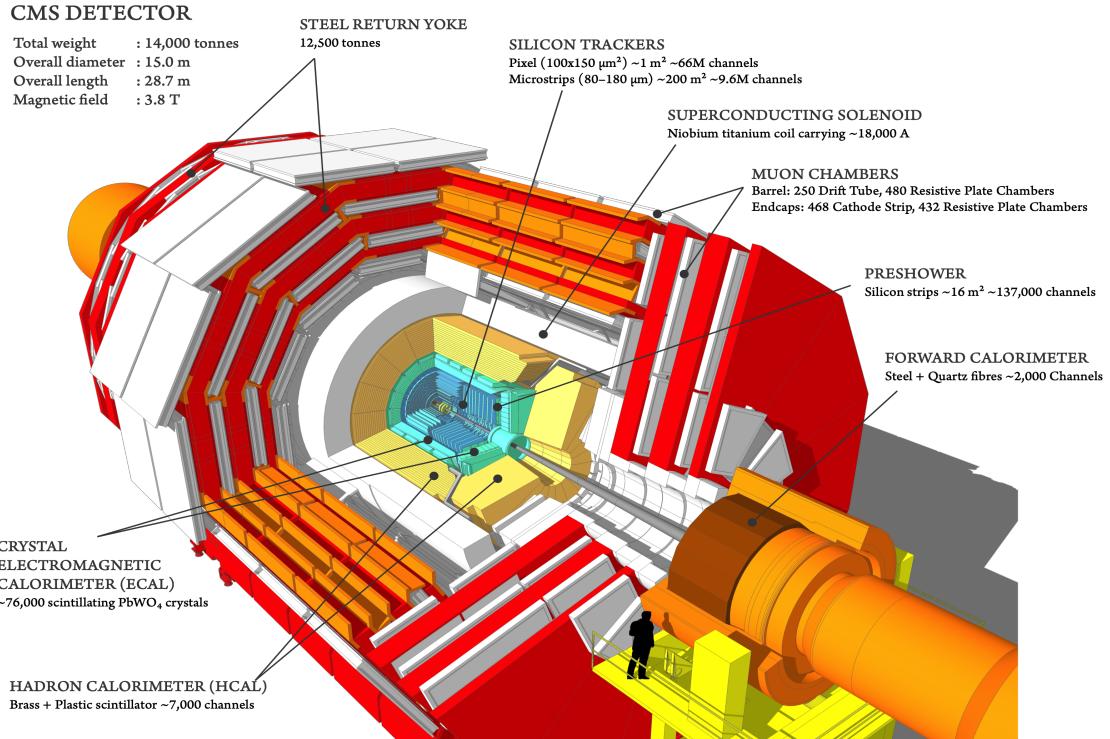


Figure 3.4: Schematic representation of the CMS detector, along with all its sub-detectors and main characteristics.

have been produced via the Lorentz interaction, therefore giving us a way to measure their charge and energy from the measurement of the curvature of the induced binding.

Finally, on the outside of the detector the complete muon system can be found and particularly performing in CMS. This sub-detector is currently made out of three main sub-systems (the Drift tubes (DTs), Cathode Strip Chambers (CSCs) and Resistive Plate Chambers (RPCs)), as explained in Section 3.2.5, and is responsible for the identification and measurement of the momentum of the muons produced by the collisions.

### Coordinates system

As a convention, it has been decided to use within the CMS collaboration a right-handed Cartesian coordinate system with the origin defined as the interaction point, and with an x-axis pointing towards the center of the ring, an y-axis pointing upwards and a z-axis pointing towards the Jura mountains (along the counterclockwise beam direction), as represented in Figure 3.5.

The  $\theta$  and  $\phi$  angles are then defined as the angles between the z and y axes and the x and y axes respectively and the pseudorapidity  $\eta$ , defined in Equation 3.4, a Lorentz invariant quantity under boosts quite often used in the different analyses since the multiplicity of high energy particles is roughly expected to be constant in  $\eta$ .

$$\eta = -\log \left( \tan \left( \frac{\theta}{2} \right) \right) \quad (3.4)$$

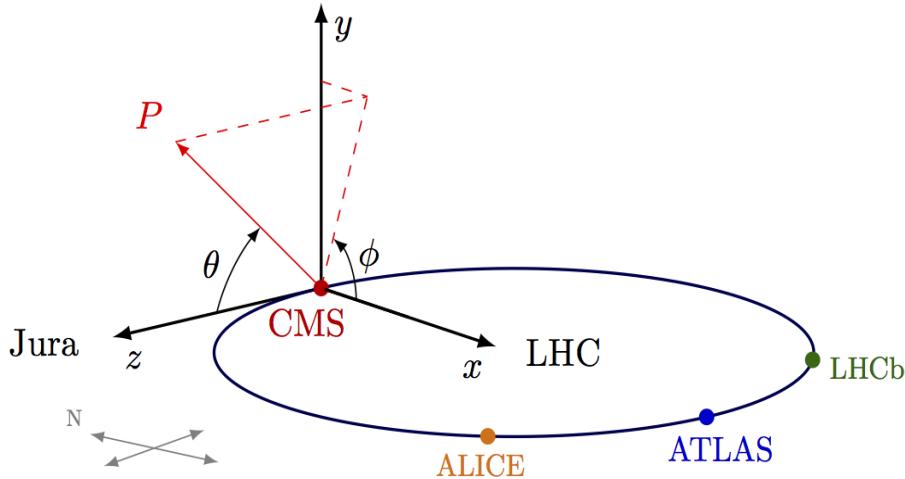


Figure 3.5: Schematic representation of the CMS coordinate system used by convention.

### 3.2.1 Tracker

The tracker is the innermost part of the CMS detector, is 5.4 meters long and has a diameter of 2.5 meters. Its main purpose is the reconstruction of the trajectories of charged particles issued from the primary and secondary interaction vertices in a quick and precise way in order to identify them and measure their individual momentum.

Many challenges were faced when designing this system mainly because of the hard conditions provided by the LHC. First of all, at its nominal instantaneous luminosity, an average of 1000 particles are created after each bunch crossing, every 25 nanoseconds. The tracking system then needs to read all its channels extremely fast in order to be ready for the next bunch crossing. However, this fast electronics then needs some cooling to work optimally, which would in return increase the size of the tracker and therefore increase the interactions between the detector and the particles created (by multiple scattering, bremsstrahlung, photon conversion and nuclear interactions). This would affect the trajectory of the particles so a compromise had to be found between the velocity and size of the tracker. Finally, this device needs to be resistant to the extreme radiation environment for its expected lifetime of at least 10 years.

This device is then made out of silicon pixels and strips mainly because of the granularity, reading velocity and radiation hardness offered by such material. It has been set up on several different layers disposed in such a way to make the detector as hermetic as possible, as shown in Figure 3.6. A charged particle crossing the tracker will then leave a hit each time it crosses one of the silicon sensors, from which the track of the particle can be reconstructed using reconstruction algorithms that will be detailed in Chapter 4.

The presence of the magnetic field due to the solenoid can then help us estimate the momentum of the particle, since the Lorentz force applied on such particle will introduce a deviation directly proportional to its momentum. The radius of curvature of particles with a high momentum ( $> 100$  GeV) is really large but the density of pixels and the algorithm can still manage to estimate the momentum of such particles, even though the uncertainty associated will then be greater.

In particular, the tracker is made out of two distinct parts: the smaller and inner pixel detector and the larger silicon strip detector:

- The **pixel detector** is made out of three barrel pixel layers and two endcap disks for hermeticity (located at radii  $r = 4.4, 7.3$  and  $10.2$  cm of the PV), one on each side. In total,

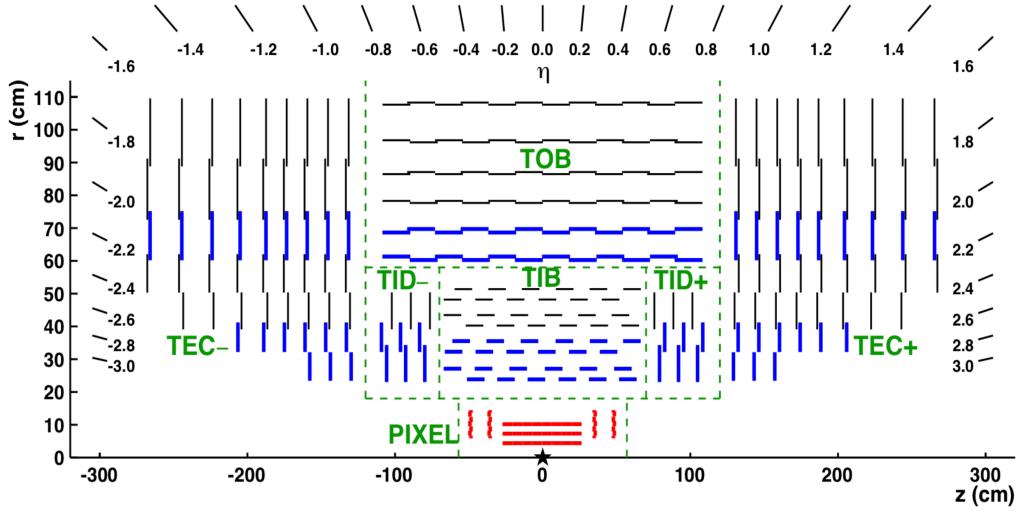


Figure 3.6: Schematic representation of the CMS tracker, for different pseudorapidity values and along the z-axis [97].

more than 60 millions pixels make up the 1440 (1856 after an upgrade in 2017) modules of this detector, covering an area corresponding to  $\sim 1 \text{ m}^2$ .

- The **silicon strip detector** on the other hand is composed of three different sub-systems, as seen in Figure 3.6 and covers a total area of  $\sim 200 \text{ m}^2$ . First of all, the Tracker Inner Barrel and Disks (TIB/TBD) add an additional 4 barrel layers and 3 endcap disks to the tracker system up to a distance  $r = 55 \text{ cm}$  to the PV of the interaction, using silicon micro-strip sensors parallel in the barrel and perpendicular to the beam axis in both endcaps. Then, the Tracker Outer Barrel (TOB) surrounds this first layer; having an outer radius of 166 cm and going up to 118 cm along the z-axis, it adds 6 layers to the inner tracking system. Finally, the Tracker EndCaps (TECs) are made out of 9 disks and complete the measurement of particles emitted along the z-axis and having a high pseudorapidity  $\eta$ .

The CMS tracker is extremely performing: one can see first of all in Figure 3.7 that for high momentum tracks of 100 GeV, the resolution is 1-2% up to  $|\eta| < 1.6$ .

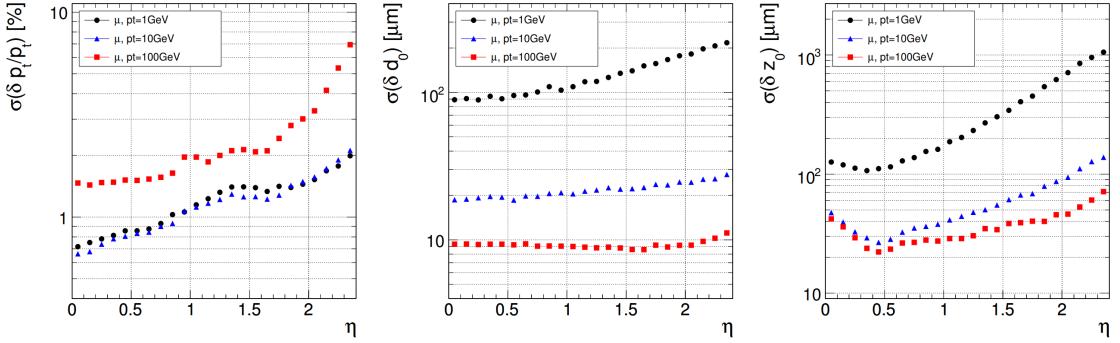


Figure 3.7: Expected resolution of muons transverse momentum (left), transverse impact parameter (middle) and longitudinal impact parameter (right), as a function of pseudorapidity and muon momentum (1, 10 and 100 GeV) [96].

Additionally, Figure 3.8 shows that muons are reconstructed with an efficiency higher than 99% for most of the pseudorapidity spectrum, even though this efficiency drops at high  $\eta$  values mainly

because of the reduced coverage provided by the pixel forward disks. The interactions between the hadrons and the tracking system is also a bit higher, which results in a lower reconstruction efficiency for such particles.

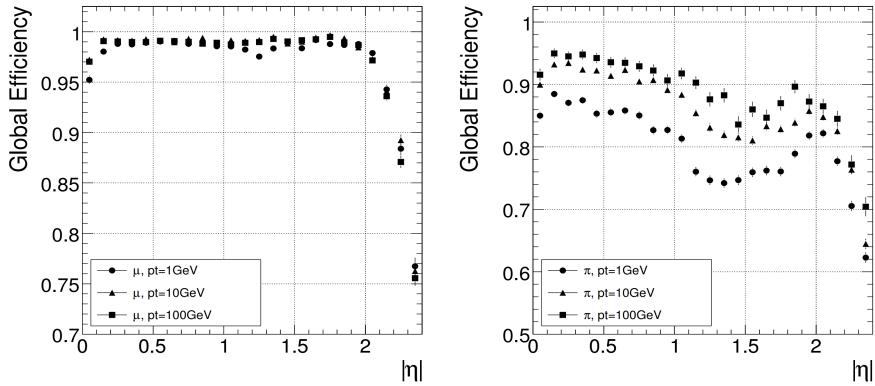


Figure 3.8: Tracker reconstruction efficiency of muons (on the left) and pions (on the right) in simulation for different pseudorapidities and particle momenta (1, 10 and 100 GeV) [96].

### 3.2.2 Electromagnetic Calorimeter (ECAL)

The ECAL of CMS is a mostly homogeneous detector inside the solenoid and enclosing the tracker system that gives information about the energy of electrons and photons, both able to interact electromagnetically with its crystals.

The ECAL can also be divided into several sections: first of all, at pseudorapidities  $|\eta| < 1.479$  is found the barrel part of the ECAL (EB), made out of 61 200 lead tungstate ( $\text{PbWO}_4$ ) crystals, located at a radius of 1.29 meters of the beam pipe. Then, two endcaps, each made out of 7 324 of those crystals, increase the coverage of the detector up to  $|\eta| < 3$ , as shown in Figure 3.9, and the preshower completes the ECAL.

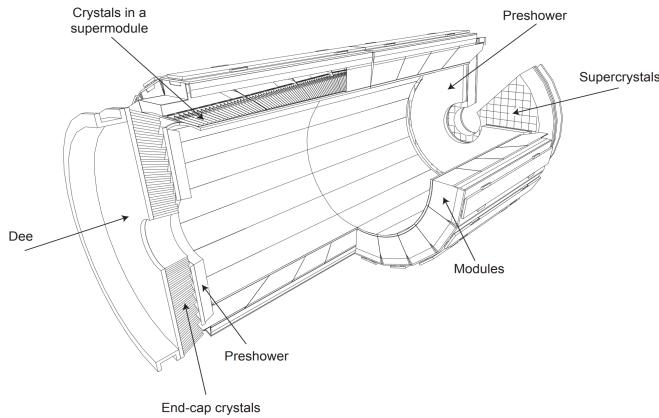


Figure 3.9: Schematic representation of the sub-systems of the CMS ECAL [96].

Its principle of action is simple and based on electromagnetic showers: when a particle such as an electron or a photon enters the ECAL, it will start to interact in different ways, depending on its nature. Photons will mainly produce pairs of electrons and anti-electrons, while the electrons themselves tend to emit additional photons by bremsstrahlung effect. This results in a chain

reaction during which the incident particle will give most of its energy to the detector itself, energy measurable using photo-detectors and photo-multipliers. This effect is known as electromagnetic shower, is represented in Figure 3.10 and is usually characterized using the so-called radiation length  $X_0$ , the mean distance over which a high-energy particle loses all but  $1/e$  of its energy, then determining the total length of interaction of a particle in the ECAL.

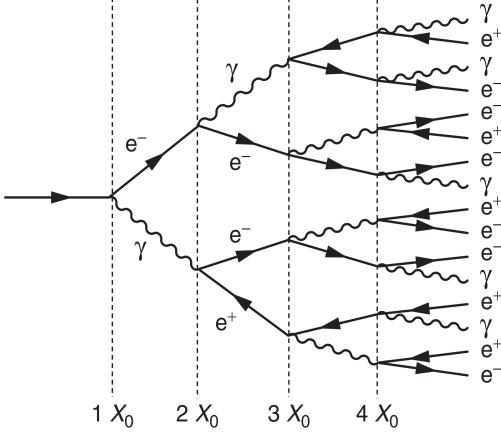


Figure 3.10: Schematic representation of a typical electromagnetic shower and the radiation length  $X_0$  concept [96].

The high density, short radiation length and low scintillation decay time (smaller than the bunch spacing of 25ns) of the PbWO<sub>4</sub> crystals make them perfect candidates towards a compact ECAL in CMS. These crystals do have some drawbacks as well, mainly their relative fragility when it comes to radiation, and the dependence on the temperature of their response. Indeed, a cooling system had to be built in order to keep the huge detector under temperature variations lower than 0.1° to avoid eventual fluctuations in the response of the crystals.

These crystals, which had to be grown individually in laboratories, measure 2.2 x 2.2 x 23 cm in the barrel and 3x3x22cm in the endcaps, cover a solid angle equal to  $(\Delta\eta, \Delta\phi) = (0.0174, 0.0174)$ , and have therefore a length corresponding to around 26 radiation lengths, more than enough to stop even the most energetic particles. Since the light output of such crystals is quite low (only around 4.5 photo-electrons per MeV of energy), they have to be connected to both avalanche photo-diodes and vacuum photo-triodes in order to multiply the signal measured. Finally, they have been mounted using a specific installation, slightly tilted with respect to both  $\phi$  and  $\eta$  in order to remove any possible gap between two adjacent crystals.

The typical energy resolution of the ECAL installed at CMS is given by Equation 3.5 [96], accounting for several different effects, such as the stochastic nature of the observed scintillation, the electronics and PU noise and the calibration and detector non-uniformity uncertainty.

$$\frac{\sigma_E}{E} = \sqrt{\left(\frac{2.8\%}{\sqrt{E}}\right)^2 + \left(\frac{0.12\%}{E}\right)^2 + (0.30\%)^2} \sim \frac{3 - 10\%}{\sqrt{E/\text{GeV}}} \quad (3.5)$$

Finally, a preshower layer has been set up in the fiducial region ( $1.653 < |\eta| < 2.6$ ) of the endcaps, where the angle between the two photons coming from the decay of neutral pions  $\pi^0 \rightarrow \gamma\gamma$  is small enough to be misidentified as individual photons. This detector has then been installed in order to reduce the possible misidentification of such events and to help with the identification of electrons against minimum ionizing particles. It is made of a lead layer able to initiate the electromagnetic shower process, followed by two layers of silicon strips for the actual measurement.

### 3.2.3 Hadronic Calorimeter (HCAL)

We know that charged hadrons lose energy in a continuous way when they traverse matter due to the ionization process and that all the hadrons strongly interact with the nuclei of any given medium. These principles are actually used in order to measure the energy of the hadrons produced by the LHC collisions using the HCAL sub-system of CMS.

In this case as well, showers of particles due to a chain reaction are expected since the primary hadronic interaction will produce several additional hadrons, themselves interacting even more with the medium while losing energy. This kind of hadronic showers is characterized by the  $\lambda$  parameter, the nuclear interaction length, defined as the mean distance between two interactions of relativistic hadrons. The nuclear interaction length  $\lambda$  is usually much larger than the radiation length  $X_0$ , resulting in a HCAL typically much larger in size than the ECAL.

A typical HCAL setup consists in alternating thick and high-density layers of absorber material, in which the showers can develop, and thin layers of active material used for the actual detection by sampling the energy deposition. This measurement is usually much less precise than the measurement provided by the ECAL, mostly since  $\pi^0$  decaying into photons can appear in these showers, leading to an electromagnetic component of the shower that cannot be measured, and because around 30% of the incident energy is usually lost due to nuclear excitation and break-up effects [94]. In this case, the energy resolution can therefore be expressed using Equation 3.6.

$$\frac{\sigma_E}{E} > \frac{50\%}{\sqrt{E/\text{GeV}}} \quad (3.6)$$

In CMS, the HCAL, represented in Figure 3.11, is also divided into a barrel (HB), radially constrained between radii values of 1.77 meters (outer radius of the ECAL) and 2.95 meters (inner radius of the solenoid), two endcaps (HE) and two symmetrical forward regions (HF) extending the pseudorapids coverage from  $|\eta| = 3$  to  $|\eta| = 5.2$  and located at a distance of 11.2 meters to the PV. A final part composing the HCAL is the so-called Hadron Outer (HO), which has to be put outside of the solenoid in order to increase the amount of shower absorber material of the HCAL and therefore the effective nuclear radiation length  $\lambda$ .



Figure 3.11: Schematic representation of the HCAL sub-system in CMS [96].

The HCAL barrel uses 36 identical azimuthal wedges as absorber, placed along the beam axis in such a way that the eventual cracks between them is smaller than 2 mm. The total absorber thickness at a  $90^\circ$  incidence angle is equal to only  $5.8\lambda$ , which explains why the HO had to be added in order to increase this value to make sure to slow down and completely stop even the most energetic hadrons. The active medium of the HB is made out of 70 000 tiles able to collect the scintillation light, using the wavelength shifting fiber concept to reduce the energy of detected photons and measure the energy of the hadrons. The HF on the other hand are using a Cerenkov-based radiation-hard technology to make the measurements required.

The barrel covers  $|\eta|$  regions up to 1.3, while the coverage up to  $|\eta| < 5.2$  is given by two endcaps on each side of the detector, placed in such a way to minimize the eventual cracks between the HB and the two HF. Finally, the HO, built in order to ensure adequate sampling depth at low pseudorapidity values, actually uses the solenoid itself as additional absorber material, since it is placed a bit outside of this coil. Its shape is constrained by the muon system and the mean fraction of recovered energy from the HO has been estimated to be equal to 0.38% for 10 GeV pions and up to 4.3% for 300 GeV pions, as shown in Figure 3.12.

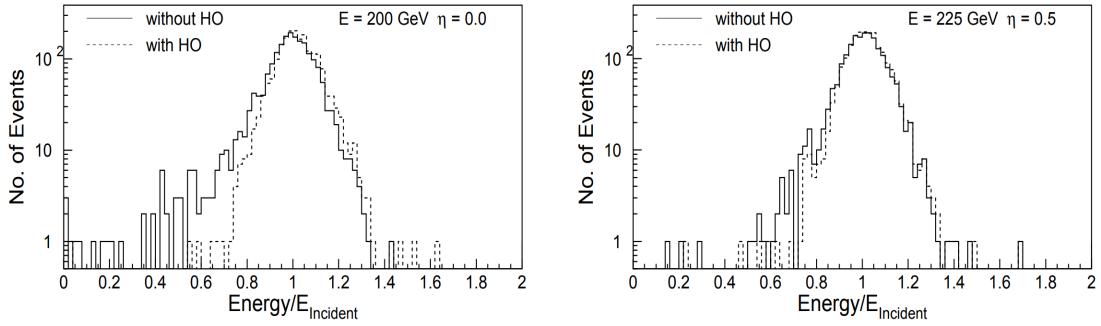


Figure 3.12: Distribution of the measured energy scaled to the incident energy for pions with incident energies of 200 GeV at  $\eta = 0$  (on the left) and 225 GeV at  $\eta = 0.5$  (on the right), with and without the inclusion of the HO in the HCAL system [96].

### 3.2.4 Solenoid

The central piece of CMS is its extremely large (12.5 meters of length and 6 meters of diameter) and heavy (220 tons) superconducting solenoid able to produce a 3.8 T magnetic field, storing when active a huge energy of 2.6GJ. It is the largest magnet of its type ever constructed, therefore allowing the tracker, ECAL and HCAL calorimeters to be placed inside the coil, resulting in a detector that is, overall, quite compact compared to detectors of similar weight.

The magnetic field produced by this coil is extremely useful since it allows to measure quite precisely the charge and the momentum of the different charged particles interacting with the detector just by measuring the curvature of their track, according to the Lorentz Equation 3.7. This solenoid has been designed to reach a momentum resolution  $\Delta p/p \sim 10\%$  at  $p = 1$  TeV.

$$\vec{F} = \frac{m \vec{v}^2}{R} = q \vec{E} + q \vec{v} \times \vec{B} = q \vec{v} \times \vec{B} \quad (3.7)$$

The hoop strain, normal stress parallel to the axis of cylindrical symmetry applied throughout this solenoid is quite large ( $\epsilon = 130\text{MPa}$ ) compared to the strain applied on other previous detectors and it had to be taken into account during the conception of this magnet. It has then been designed in such a way that a large fraction of the CMS coil has a structural function, dividing

the strain between the layers of the magnet and the support of the coil itself ( $\sim 30\%$ ). At the end, the conductor of this solenoid, made from a Rutherford-type cable combined with aluminum, is mechanically reinforced with an aluminum alloy.

The coil of the magnet is then completed with a huge steel yoke return system, as seen in Figure 3.13, made out of 6 endcap disks and 5 barrel wheels, weighting in total more than 12 000 tons and therefore accounting for most of the weight of CMS. This system is composed of many steel blocks up to 620 mm thick combined and actually also serves as the absorber plates of the HO and muon detection system that will be described in the next section.



Figure 3.13: Picture of the solenoid system of CMS being setup in the assembly hall.

Finally, a two pumping stations system has been put in place in order to setup a vacuum as strong as possible inside the  $40 \text{ m}^3$  volume of the coil cryostat and an helium refrigeration plant has been installed near the site of the detector, able to cool down the solenoid up to 4.5 K, giving a 2 K security margin with respect to the critical field of the superconducting coil. All these systems were extensively tested on the surface during the summer of 2006, before lowering down the complete solenoid in the experimental cavern where it now stands.

### 3.2.5 Muon system

The muon detection is extremely useful since many interesting processes are expected to produce such particles. Their detection and the correct measurement of their main properties such as their position and momentum is therefore crucial in most of the analyses performed. Detecting muons is at the end of the day quite easy, as we will see, and the data extracted from them is usually more reliable than the one obtained from electrons since muons are less likely to be affected by the inner parts of the detector, such as the tracker, because of their low interaction cross section.

The muon system of CMS is actually made out of three different gaseous sub-subsystems combined in order to perform a reconstruction as precise as possible of such particles over the entire kinematic range of the LHC. These different muon chambers systems do share some characteristics: they mostly have to be distributed over a cylindrical area, because of the geometrical shape of the inner systems of CMS and they have to be reliable and cheap, since they cover a total area corresponding to more than  $25\,000 \text{ m}^2$ .

Each category of muon chambers is typically used in a different pseudorapidity area, as shown in Figure 3.14, in order to form a muon system as hermetic as possible.

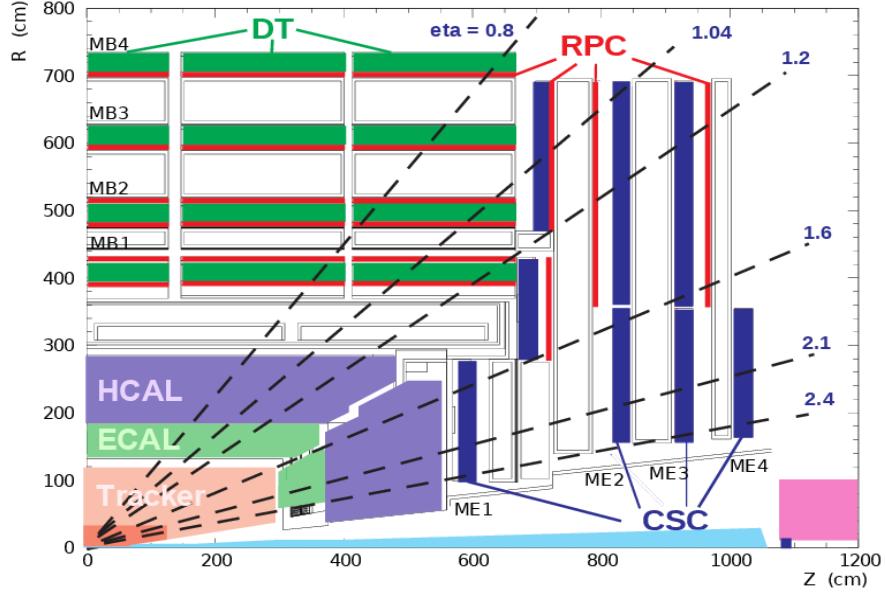


Figure 3.14: Geometrical repartition along the z-axis of the different muons chambers in CMS [98].

### Drift tubes (DTs)

First of all, in the barrel region, where the flux of muons is low and where the magnetic field is mostly uniform and low as well, the **DTs** have been installed. This system covers the  $|\eta| < 1.2$  area and has been divided into 4 different layers, each containing a number of stations optimized in order to provide a full coverage of the  $\theta$  angles, a good efficiency for the muon hits reconstruction into a single track and a good rejection of eventual background hits. This distribution of the DTs is represented in Figure 3.15.

The DT system is made out of 172 000 sensitive wires able to collect the residuals charges left by the ionization tracks of muons through the 250 chambers installed. The system has been set up in such a way that the maximal drift of any charge is lower than 21mm, corresponding to a drift time of 380 ns in the gaseous chambers made out of 85% of Ar and 15% of CO<sub>2</sub>, a value small enough to produce negligible occupancy in the different wires and to avoid the need of multi-hits electronics. Redundancy of the DTs provided by the installation of multiple layers is extremely important, mainly to reduce the backgrounds coming from eventual neutrons or photons, whose rate is actually much larger than the one obtained from prompt muons.

### Cathode Strip Chambers (CSCs)

In the two endcap regions, where the muon rates and the background levels are much larger and where the magnetic field is large and non-uniform, a different system had to be installed. First of all, the **CSCs**, multi-wire proportional chambers providing a fast response while being resistant to the radiation, are able to identify muons in a  $0.9 < |\eta| < 2.4$  region (in the  $0.9 < |\eta| < 1.2$  region, muons cross both DTs and CSCs while in the  $1.2 < |\eta| < 2.4$  area, muons typically cross between 3 and 4 CSCs only).

This sub-system is made out of 540 different chambers in total, all perpendicular to the beam pipe. The sensitive plates of this sub-system are made out of 2 million wires, cover about 5000 m<sup>2</sup> and the total gas volume included in such chambers is equal to about 50 m<sup>3</sup>.

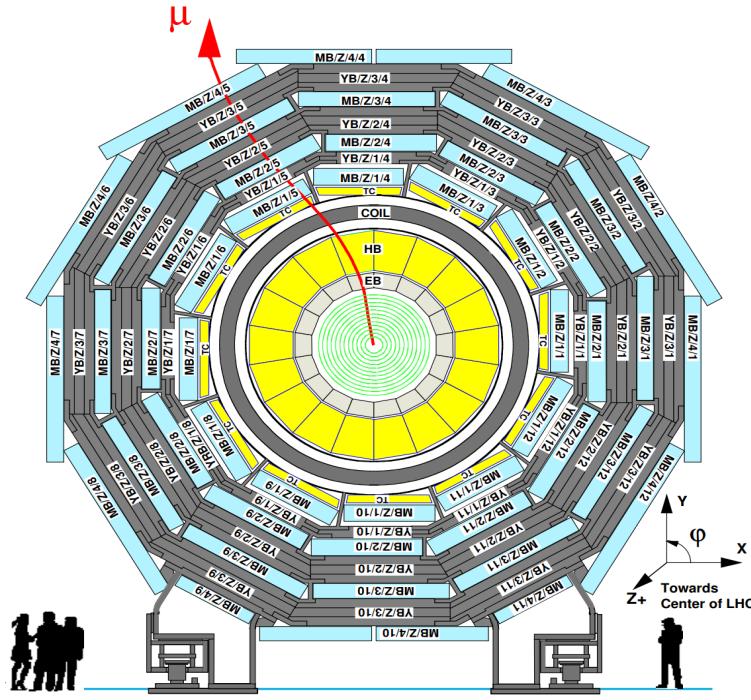


Figure 3.15: Lateral geometrical division of the different DT chambers in one of the 5 wheels of CMS [96].

### Resistive Plate Chambers (RPCs)

Finally, some **RPCs** have been added to the barrel and to the endcap regions in order to cope with the uncertainty associated with the eventual background rates and with the (in)ability of the previous muon system to identify unequivocally the correct bunch crossing when the LHC is running at full luminosity. Indeed, the time resolution of the DTs of 380 ns is way larger than the bunch spacing in the LHC, while a RPC is capable of tagging an ionizing event in less than 25ns, making it an ideal candidate to trigger the event.

The RPCs are double-gap chambers operated in avalanche mode to ensure good operation at high rates and they are able to produce a fast response with good time resolution, even though its position resolution is worse than the one obtained with DTs or CSCs. The RPCs are also useful in the sense that they can help to resolve ambiguities when attempting to construct tracks from multiple hits in a chamber.

Finally, the different features of the three muon sub-systems used by the CMS detector are summarized in Table 3.2.

Another advantage of the muon system such as the one built in CMS is that it can also directly be used by the trigger system, which will be described in Section 3.2.6, independently of the rest of the detector and in addition of being able to detect, identify and measure several properties of muons crossing it.

The muon reconstruction efficiency obtained by the muon system strongly depends on the pseudorapidity value of the muon considered, along with its transverse momentum, as shown in Figure 3.16. In this figure, we can also see that several different kinds of muons can be defined, such as the **standalone muons**, defined using only the data coming from the muon system and the **global**

Muon sub-system	DT	CSC	RPC
$ \eta $ coverage	0.0-1.2	0.9-2.4	0.0-1.9
Stations	4	4	4
Chambers	250	540	480 (barrel) 576 (endcaps)
Readout channels	172 000	266 112 (strips) 210 816 (anode channels)	68 136 (barrel) 55 296 (endcaps)
Spatial resolution	80-120 $\mu\text{m}$	40-150 $\mu\text{m}$	0.8-1.2 cm
Average efficiency (13 TeV)	97.1%	97.4%	94.2% (barrel) 96.4% (endcaps)

Table 3.2: Comparison of the three main sub-systems currently used by CMS in order to identify and measure muons [99].

**muons**, defined using both the information coming from the muon system and the tracker. This distinction will be detailed when discussing about the muons reconstruction in Section 4.3.1.

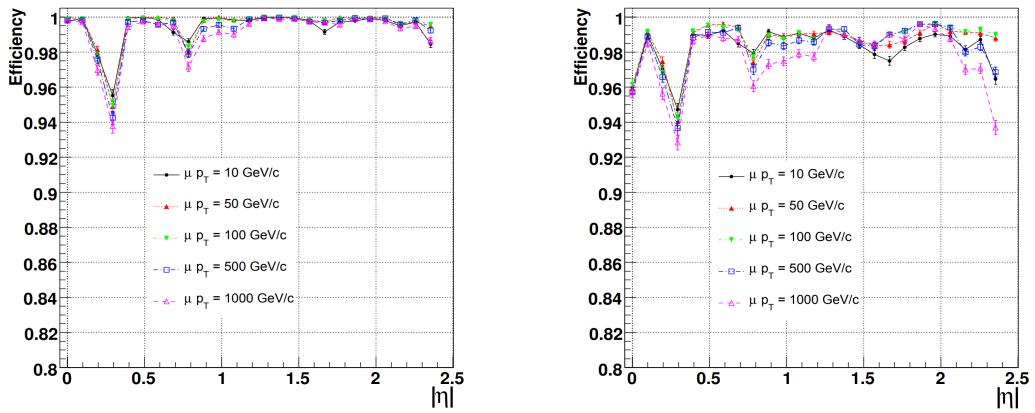


Figure 3.16: Muon reconstruction efficiency with different  $p_T$  and  $\eta$  values, considering only the muon system (on the left), and the combined information from the muon system and the tracker (on the right) [96].

Taking advantage of the LS2, a new muon system is currently being installed in the experimental cavern: the so-called **Gas Electron Multipliers (GEMs)**, placed in the endcaps, where the radiation and event rates are the highest. This new sub-detector will provide additional redundancy and measurement points to the current system, therefore allowing a better muon track identification and reconstruction and a wider coverage in the very forward region.

The first 144 chambers of the GEM sub-system, filled with a mixture of Ar and CO<sub>2</sub> and where the primary ionization due to incident muon is expected to happen, are currently being installed in the first disk of both endcaps (cf. Figure 3.17), while the rest will be set up during the next LS expected in 2024, before the phase II of operation of the LHC.

### 3.2.6 Trigger system

The CMS experiment faces a data acquisition limitation since the collision rate delivered by the LHC (one bunch crossing each 25ns, leading to an impressive rate of collisions of 40 MHz) is much

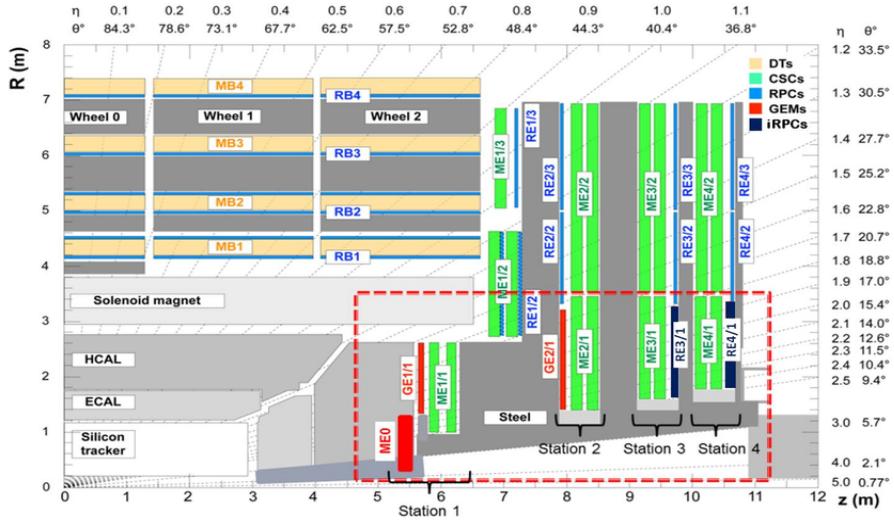


Figure 3.17: Location of the new GEM muon subsystem currently being installed in the very-forward region of CMS [96].

larger than the data acquisition rate currently achievable by nowadays electronics (around 1kHz only, more or less equivalent to 1Gb of data per second). It is therefore impossible to store and process all the collisions provided by the LHC; instead, a selection needs to be made in order to select and record only the most interesting events.

A system, called the trigger system, has therefore been put in place in order to select extremely quickly 1 kHz of interesting events out the 40 MHz. This system is based on two different levels: first of all, the Level-1 Trigger (L1), a hardware set of electronics selecting around 100 kHz of data, followed by the High-Level Trigger (HLT), a software layer improving the selection even more.

### Level-1 Trigger (L1)

The L1 is the first level of trigger [100], based directly on hardware. In order to maximize its efficiency, it is mostly implemented in different subsystems of the detector (on the calorimeters and muons system) and in the service cavern, just next to the detector, so that the electric signals do not have to travel large distances, therefore saving a few precious nanoseconds of decision time.

The L1 trigger, whose architecture is represented in Figure 3.18, has local, regional and global components. First of all, the local triggers, also referred to the Trigger Primitive Generators (TPG), are using the energy deposits measured in the calorimeter and the track segments or hit patterns in the muon chambers. Then, regional triggers combine this information and use pattern logic to determine ranked and sorted trigger objects such as electron or muon candidates in limited spatial regions. The rank is typically determined as a function of energy or momentum and quality, which reflects the level of confidence attributed to the L1 parameter measurements, based on detailed knowledge of the detectors and trigger electronics and on the amount of information available. Finally, the global calorimeter and muon triggers combine this local information to determine the highest-rank calorimeter and muon objects across the entire experiment and transfer them to the top entity of the L1 hierarchy, the global trigger. This level is in charge of taking the final decision on whether to keep or reject the event, and eventually passes it to the HLT.

This trigger gets new data each 25ns and today's electronics is not fast enough in order to deal with such a massive input of data, so an ingenious systems of buffers had to be put in place in

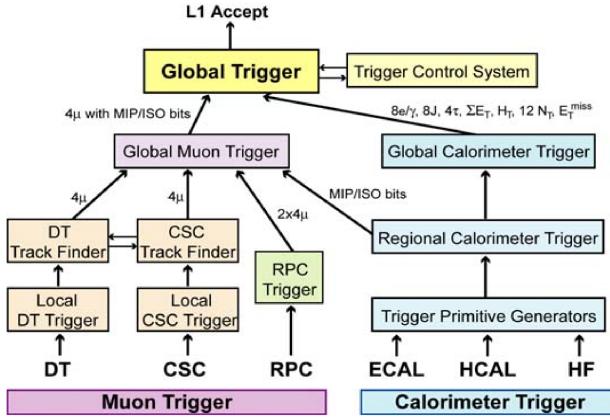


Figure 3.18: Architecture of the L1 trigger of CMS [96].

order to put in line several events before analyzing all of them at once, using only basic segmented data provided by the detector.

### High-Level Trigger (HLT)

On the other hand, the HLT [101] does get access to the complete read-out data of the detector, since the rate has already been strongly reduced by the L1 Trigger, allowing it to perform complete calculations such as the ones that will be later performed offline. Since this is a software based layer and the decision time is not as critic, it runs on a farm of computers on the surface and is in constant evolution, getting constantly improvements and updates in order to select and reconstruct in a better and more efficient way interesting data from the collisions.

The HLT uses the so-called **trigger paths** in order to select specific particle topologies in the collisions. This selection imposes some thresholds on key quantities of the different objects, such as their transverse momentum  $p_T$ , to avoid passing the bandwidth limit of 1kHz. In some cases, not all the events of a given category can be collected and only a fraction of them are therefore recorded. Trigger paths designed to behave like this are called pre-scaled trigger paths.

### 3.2.7 Data AcQuisition system (DAQ)

In summary, the CMS DAQ, whose global architecture is represented in Figure 3.19, has been designed to collect and analyze the data information at the nominal collision rate of 40MHz and is fed directly from the L1 trigger. This means that it has to be able to read a flux of data of the order of 100kHz ( $\sim 100\text{Gb/s}$ ) coming from approximately 650 different sources at once, while providing enough computing power for the HLT to be able to reduce this rate by a factor  $\sim 100$ , while keeping some resources available for other tasks.

The DAQ is indeed also in charge of performing additional tasks, such as the generation of the Data Quality Monitoring (DQM) information resulting from online event processing in the HLT, the transfer of the data from local storage at the CMS site to mass storage in the CERN data center at the Meyrin site and the operation of the Detector Control System (DCS) system, ensuring the correct operation of the detector and a high quality data taking at all times.

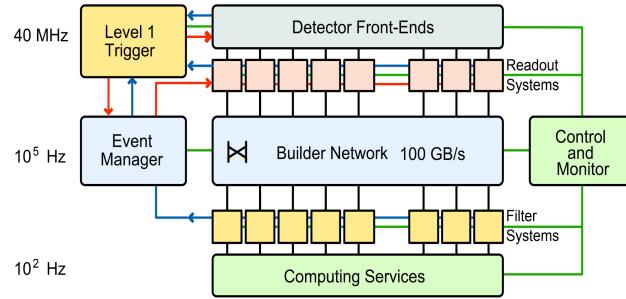


Figure 3.19: Architecture of the CMS DAQ system [96].

### 3.3 CMS main goals

At the end of the day, the CMS detector has been built in order to meet the goals of the LHC physics program, which have been summarized in [96] as:

- Good muon identification and momentum resolution over a wide range of momenta and angles, good dimuon mass resolution ( $\sim 1\%$  at 100 GeV), and the ability to determine unambiguously the charge of muons with momentum  $p < 1$  TeV.
- Good charged-particle momentum resolution and reconstruction efficiency in the inner tracker. Efficient triggering and offline tagging of taus and b-jets, requiring pixel detectors close to the interaction region.
- Good electromagnetic energy resolution, good diphoton and dielectron mass resolution ( $\sim 1\%$  at 100 GeV), wide geometric coverage,  $\pi^0$  rejection, and efficient photon and lepton isolation at high luminosities.
- Good missing-transverse-energy and dijet-mass resolution, requiring hadron calorimeters with a large hermetic geometric coverage and with fine lateral segmentation.

---

---

# Chapter 4

---

## Event reconstruction

The CMS detector is made out of different layers, each able to convert the interaction of the particle into electronic signals that can be measured and stored.

However, the recorded electronic signals contain only low-level, partial information about the particles, in such a way that an algorithmic strategy is needed in order to fully reconstruct physical information such as the number of particles and their charge, momentum and direction. Producing this kind of data is essential for all the offline analyses which usually rely on these high-level physics objects to make precision measurements or search for new physics.

The algorithm able to combine this raw data and to compute and produce useful kinematic variables and physical objects (such as leptons and jets) is the so-called Particle Flow (PF) algorithm [102], which will be first of all described in Section 4.1. Then, a particular focus will be given to the definition and reconstruction of different objects of our particular analysis, such as electrons and muons (Section 4.3), jets (Section 4.4), the MET (Section 4.5) and the top reconstruction (Section 4.6) of the different  $pp$  collisions recorded.

### 4.1 Particle Flow (PF) algorithm

The PF is an algorithm aiming to combine in the best way possible all the information coming from the different parts of the CMS detector (mostly, tracks and clusters of energy) in order to identify and reconstruct the hundreds of new particles produced by each  $pp$  collision provided by the LHC. This reconstruction can be divided into two main steps: first of all, the data coming from the different subsystems of the detector is read in order to identify and measure the properties of some basic stable objects, such as leptons, photons and hadrons. Then, more complex calculations are performed to identify eventual unstable particles, jets from the hadronization of quarks and gluons, and to compute complex variables such as the leptons isolation and the MET.

The most basic elements used by this algorithm for the reconstruction of high-level physics objects are the hits left by charged particles in the tracker and in the muon chambers, and the clusters of energy left in the two calorimeters. For this algorithm to be as efficient as possible, the detector has been carefully designed, as described in Chapter 3: a magnetic field as intense as possible

and a small calorimeter granularity are indeed crucial in order to separate efficiently charged and neutral particles, and the tracker was designed to be as efficient and small as possible to have the smallest material budget possible in front of the calorimeters. The muon system has been carefully designed as well and, in general, the whole detector is obviously as hermetic as possible.

The way the different particles produced in each collision are identified is quite easy to summarize and is represented in Figure 4.1. Basically, the different kinds of particles produced interact with different parts of the detector and the combination of this information then allows to unequivocally identify each particle. The PF algorithm starts from the collection of tracks and vertices reconstructed by the dedicated tracking and vertexing algorithms and from the local reconstruction of energy clusters in the calorimeters. The sequence follows a specific order to make this reconstruction process as efficient as possible:

1. First of all, the most energetic **Primary Vertex (PV)** is identified by taking into account the PU and by assigning the different tracks to the different  $pp$  collisions happening during a single bunch-crossing.  
All the tracks originating from less energetic PV or from a secondary vertex of interaction are typically ignored and the corresponding hits in the tracker left by such particles can therefore be ignored later on, leaving less hits available for the clustering algorithm later on, allowing for a more efficient reconstruction of the following objects.
2. Then, **muons** are the easiest particles to identify since they are at first order the only particles leaving many hits in the muon chambers placed on the outside of the detector. Depending on the version of the algorithm used, each muon identified can be associated to its track in the tracker, where all the hits matching a muon can therefore be subsequently ignored to simplify the following reconstruction steps.
3. **Electrons** do have a charge, so they are visible by the tracker, and can interact electromagnetically with the ECAL, where they are going to produce an electromagnetic shower. Identifying electrons is a bit more challenging than muons because of their associated bremsstrahlung emission of photons that need to be attached to the original electrons to measure efficiently their complete energy and avoid any double counting. All the tracker hits corresponding to electrons are also ignored for the rest of the reconstruction after identification.
4. **Charged hadrons** also leave hits in the tracker and some energy deposits in the ECAL, but mostly in the HCAL, so they are easy to identify as well as a fourth step, using the last hits available in the tracker.
5. **Photons** are on the other hand neutral particles, so they do not leave any hits in the tracker. They then appear as some energy deposits in the ECAL for which no corresponding tracker track can be associated.
6. Finally, **neutral hadrons** can be identified as particles leaving some energy mostly in the HCAL for which no corresponding tracker track has been found as well.

We will now study in more detail the reconstruction method applied in order to reconstruct the main objects of this analysis, i.e. the leptons, jets, MET and bottom/top quarks.

## 4.2 Primary vertex definition

Different kinds of vertices originating from a single  $pp$  collision can usually be defined, as shown in Figure 4.2. First of all, we have the set of **PU vertices**, corresponding to the different simultaneous

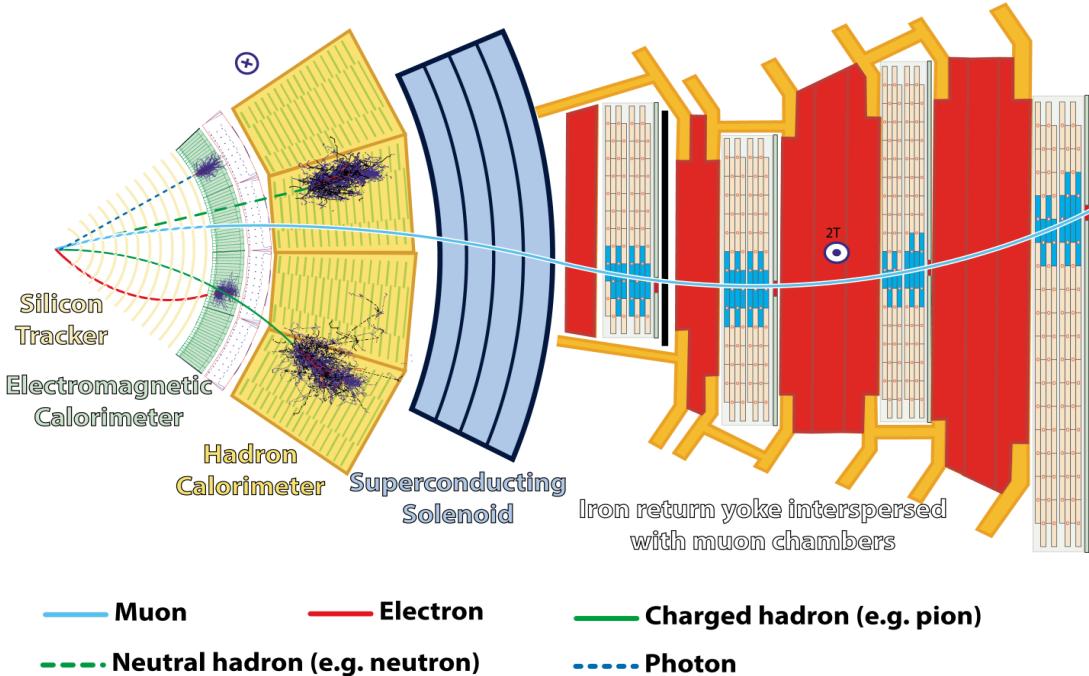


Figure 4.1: Transverse section of CMS showing the different interactions expected by different kinds of particles in the detector.

collisions of a single bunch-crossing of the LHC or eventual remnants coming from a previous collision (the out-of-time PU). The **PV** is then defined as the most energetic PU vertex and the only vertex considered in most of the physics analyses, and finally, we have the **secondary vertices**, mainly due to the eventual presence of hadrons living long enough to travel a significant distance in the detector, decay chains or jets.

The PF algorithm is not able to reconstruct the vertices itself, but it can use the information coming from a dedicated vertexing algorithm able to identify all these kinds of vertices. This is done by considering all the tracker hits observed, by clustering them together and by performing fits to determine the likelihood these tracks originating from a common vertex. The reconstructed vertex with the largest  $p_T^2$  summed over all the physics objects of the event is then assumed to be the PV, as it is considered to be the origin of the collision that actually fired the trigger, from which many different tracks are emitted.

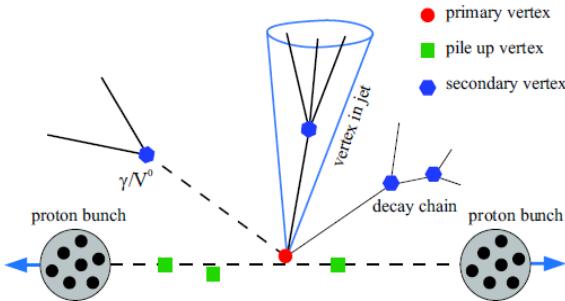


Figure 4.2: Different kinds of vertices typically observed in a  $pp$  collision in the LHC.

## 4.3 Leptons reconstruction

Different kinds of leptons are typically produced by a  $pp$  collision. The muons and the electrons and their respective anti-particles can be quite easily identified, mainly because their lifetime and velocity is high enough, meaning that they are not expected to decay inside the CMS detector, so they can be directly identified. Taus on the other hand are a bit trickier to deal with because they usually decay inside of the beam pipe itself,  $\sim 35\%$  of the time to electrons and muons and  $\sim 65\%$  of the time to hadrons. However, since our analysis does not consider taus directly but only the leptons originating from their decays, the details of their reconstruction will not be explained in this section.

### 4.3.1 Muons

Muons are the first leptons to be reconstructed by the PF algorithm since, by design and at first order, they are the only particles expected to reach the muon chambers, resulting in an easier and more efficient identification.

The typical signature of a muon consists of several hits in the silicon tracker forming a track associated with several hits in the muon chambers, electronic signals coming from the wires and strips of these chambers due to the gas ionization induced by the passage of these charged particles. Muons only deposit a negligible amount of energy within the two calorimeters since their interaction cross section is quite low for their full range of energies, going from a few hundreds MeV up to a few TeV.

The data coming from the different subsystems of CMS are then combined and fed to the PF algorithm, then able to reconstruct different kinds of muons [99].

#### Standalone muons

The standalone muons are muons reconstructed using only the hits observed in the muon system without trying to relate this data to the tracker hits. Basically, the PF algorithm looks in this case at the eventual hits left in the DTs, CSCs and RPCs and tries to reconstruct a vector of trajectory in each case using a Kalman Filter (KF) [103]. The segments obtained are then combined in the best statistical possible way in order to form a candidate track for each muon of the event, by extrapolating the innermost vectors obtained to the next chamber and by comparing it with the local track segment; the trajectory parameters are then updated and the process continues until reaching the outermost chamber. Finally, this complete process is repeated in the reverse order, from the outside chambers to the inside, to estimate the innermost track parameters as well.

It is important to note that cosmic rays reaching the Earth every day are an important source of contamination, since we estimate that around 10.000 muons per  $m^2$  and per minute coming from such processes can be observed at the sea level. Putting the detector underground removes part of this contamination but, to further limit the possibility of misidentification due to showering of cosmic rays, the tracks need to pass some quality criteria: for example, checking the extrapolation of the trajectory to the point of closest approach to the beam line allows to reduce this contamination. Additionally, at least two hits need to be measured for the fit to be performed, one of them coming from either the DTs or CSCs, in order to remove fake segments contamination due to combinatorics.

In any case, candidates reconstructed as standalone muons typically have a worse momentum reconstruction and are more sensitive to cosmic muons contamination.

### Tracker muons

The algorithm able to reconstruct the so-called tracker muons on the other hand is able to propagate tracks identified in the inner silicon tracker (having a momentum  $p > 2.5$  GeV and  $p_T > 0.5$  GeV) to the muon system itself in order to try and find corresponding segments in the different muon chambers (these tracks are therefore said to be built *inside-out*). An extrapolated track and a segment are only matched if the difference between their positions in the  $x$  coordinate is smaller than 3 cm or if the pull, the ratio of this distance to its uncertainty, is smaller than 4.

These muons are particularly efficient for less instrumented regions of the detector and for the low  $p_T$  end of the energy spectrum but they are also quite contaminated with fake muon tracks, since a single hit in any of the muon chambers is enough for the candidate to be considered a valid tracker muon, even though hadron shower remnants can for example quite easily reach the innermost muon station. The momentum assigned to such muons is the same as the one measured by the silicon tracker track itself.

### Global muons

Finally, these muons are built *outside-in* since they are obtained by matching standalone muon tracks with independently reconstructed tracks coming from the tracker itself (of course, in order to avoid any double counting, global muons and tracker muons that share the same tracker track are actually merged into a single candidate).

This category of muons presents the advantage of being less sensitive to the muon misidentification rate than tracker muons since it uses the information from more than one muon chamber. The  $p_T$  measurement in this case is also improved (especially for  $p_T > 200$  GeV) by exploiting the information from both the inner tracker and the muon system, while at low momentum, the best momentum resolution for muons is obtained from the inner silicon tracker directly.

Using this strategy, about 99% of the muons produced within the geometrical acceptance of the muon system are reconstructed either as global or tracker muons, as seen in Figure 4.3.

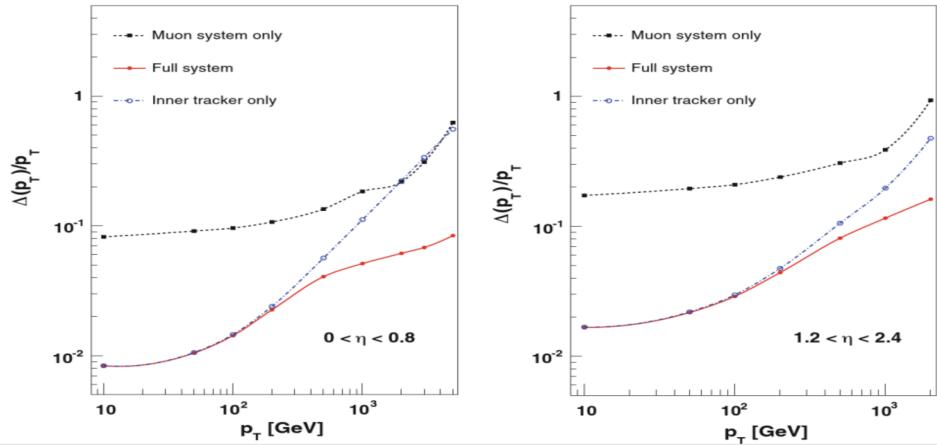


Figure 4.3: Muon  $p_T$  resolution obtained in simulation in the barrel (on the left) and endcap (on the right) for different muon reconstruction algorithms [104].

Once reconstructed, candidates are required to pass some selection criteria and are then fed to the actual PF algorithm itself to start the global reconstruction of the event. This selection consists mainly in applying identification and isolation (evaluated relative to its  $p_T$  by summing up the energy in geometrical cones of radius  $\Delta R = \sqrt{(\Delta\phi)^2 + (\Delta\eta)^2}$  surrounding the muon in the  $(\eta, \phi)$

plane, as shown in Figure 4.4) criteria in order to enhance the purity of the reconstructed prompt muons (muons directly originating from the main collision taking place in the event) by rejecting muons coming from the decay of heavy flavour quarks, typically surrounded by a large amount of hadronic activity. The calculation of the lepton isolation is typically performed considering only the PV since higher levels of PU are expected to bias this measurement by increasing the hadronic activity.

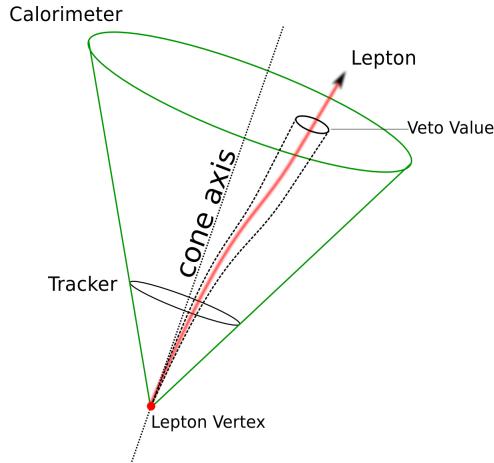


Figure 4.4: Lepton isolation cone typically used to enhance the prompt leptons purity.

Different identification Working Points (WPs) can then be defined for the offline analyses, from veto to tight, in order to reject more or less contamination from misidentified leptons, bearing in mind that a tighter selection is cleaner, but less efficient: the loose and tight WP have then been defined in order to respectively achieve 95% and 98% efficiencies, respectively. These efficiencies are tuned using simulated  $Z \rightarrow \mu^+ \mu^-$  events with a  $p_T > 20$  GeV, while the efficiency to reject muons in jets is done using simulated QCD and  $W+jets$  processes [99].

Our particular muon definition is based on the tight WP provided centrally by the muon Physics Object Group (POG), but has been slightly modified and made more robust against non-prompt leptons, according to the selection detailed in Chapter 6.

### 4.3.2 Electrons

Electrons and positrons are reconstructed by combining the tracker tracks and the several clusters of energy deposited in the ECAL by the electromagnetic showers appearing due to the interaction between the electron and the crystals composing this sub-detector.

It is usually a bit harder to reconstruct electrons than muons mainly because electrons do interact much more with the tracker and this interaction therefore needs to be modeled to understand the exact behaviour of such particles: it is for example responsible for the emission of secondary bremsstrahlung photons crashing into the ECAL but not coming from the PV. In fact, it is estimated that in CMS between 33% and 86% of the energy of an electron is actually radiated before it reaches the ECAL, depending mostly on its pseudorapidity [105]. In order to measure precisely the energy of an electron, all the photons emitted by bremsstrahlung before reaching the ECAL (usually, along the  $\phi$  plane because of the deviation implied by the solenoid) then need to be collected as well and associated to the correct electron of the event.

The actual PF reconstruction of electrons is performed in different steps:

1. A **clustering algorithm** is first of all defining the so-called **Super Cluster (SC)**. Its goal is to reconstruct the particle showers individually by identifying a seed crystal for the cluster, defined as the crystal collecting the most energy, since the energy deposited in the ECAL is usually spread into several different crystals because of the electromagnetic shower effect discussed in Section 3.2.2 and because of the bremsstrahlung emission of photons due to the interaction with the tracker. The algorithm therefore searches for eventual crystals around this seed whose energy detected would be superior to  $2\sigma$  of the electronic noise and matching some quality criteria ( $E_{\text{seed}} > 230$  MeV in the barrel,  $E_{\text{seed}} > 600$  MeV and  $E_{\text{seed}}^T > 150$  MeV in the endcaps).

The excited contiguous crystals found are then grouped into clusters, themselves considered candidates for the final global cluster, the SC, if their energy is higher than another given threshold ( $E_{\text{cluster}} > 350$  MeV in the barrel,  $E_{\text{cluster}}^T > 1$  GeV in the endcaps) [105]. The SC energy is then given by the sum of the energies of all its constituent clusters, while its position is calculated as the energy-weighted mean position of the different clusters.

2. Once the SC is identified, **electron tracker tracks** are reconstructed using a procedure a bit different than the usual KF reconstruction method for all the tracks of the silicon tracker [103] because of the large radiative losses for electrons in the tracker material.

This reconstruction is known to be very time consuming, so a good identification of potential electron seeds has to be performed as the method efficiency greatly relies on this first identification. Two different strategies can be used to perform this seeding (even though the electron seeds found using the two algorithms are usually combined afterwards):

- The **ECAL-based seeding** relies on the information obtained for the SC energy and position in order to estimate the electron trajectory to find compatible hits in the tracker. This can be done knowing that the electron is moving according to an helix in the magnetic field of the detector. This seeding is mostly optimized for isolated electrons in the  $p_T$  range relevant for the Z and W decays.
- The other way to proceed is the **tracker-based seeding**, based on tracks reconstructed using the usual KF algorithm and looking for matches within the possible reconstructed SC. This seeding is mostly suitable for low  $p_T$  electrons and also performs quite well with electrons inside jets.

Once the seeds have been identified, the identification of tracks can begin. First of all, the gathering of compatible hits from the different seeds is done using a dedicated modeling of the electron energy loss and a combinatorial KF algorithm allowing to construct possible tracks when compatible hits are found. The compatibility matching between the predicted and found hits is usually chosen to be quite loose in order to maintain a good efficiency even in case of bremsstrahlung emission.

Finally, once the hits are collected, a Gaussian Sum Filter (GSF) fit is performed to estimate the different track parameters by reconstructing the layer-to-layer propagation of electrons in the tracker. A mix of Gaussian distributions is used in this case to approximate the loss in each layer, associating a different weight and  $\chi^2$  penalty to each distribution, depending for example on the number of missing hits. This fit is also able to take into account sudden changes in the curvature radius caused by an eventual bremsstrahlung photon emission.

3. The final step consists in identifying the clusters left in the ECAL by the photons emitted by extrapolation of the GSF track and in **merging this GSF track and the ECAL SCs** previously built. This step is also designed to preserve the highest efficiency possible while keeping the misidentification probability low and ambiguities related to single electron seeds which can often lead to several reconstructed tracks are also resolved at this stage.

Finally, a loose pre-selection is applied to the electron candidates in order to reject fake electrons and the variables related to the energy and geometrical matching between the GSF

track and the ECAL cluster(s) are combined into a Multi-Variate Analysis (MVA) estimator allowing to define several electron WPs as well.

This electron reconstruction workflow has been summarized in Figure 4.5.

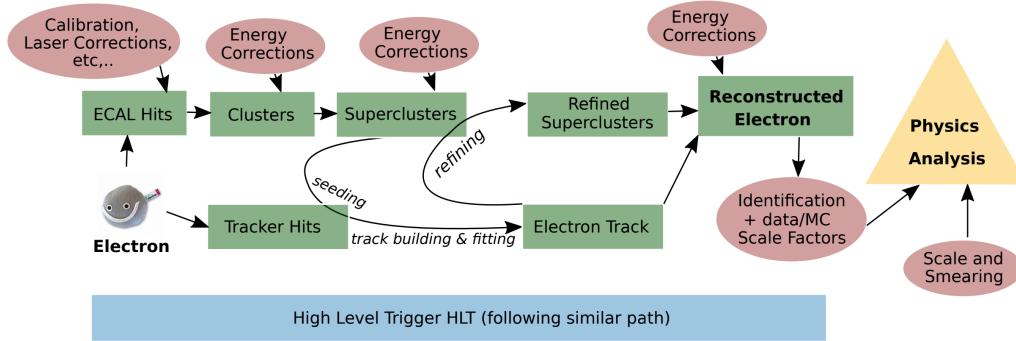


Figure 4.5: Schematic representation of the full electron reconstruction workflow in CMS [106].

## 4.4 Jets reconstruction

Eventual jets originating from quarks and/or gluons produced by a  $pp$  collision usually manifest themselves as hadronic jets in the detector because of the colour confinement principle stating that coloured particles, such as the quarks, can not be isolated and therefore be observed on their own.

This practically means that once a single quark is produced, it will start losing energy by forming new  $q\bar{q}$  pairs, themselves forming additional  $q\bar{q}$  pairs. This chain continues until the resulting pairs of quarks have such a low energy that they can start combining into colourless hadrons. This is called the *hadronization* process and the actual result of the apparition of a quark is a shower of collimated particles, usually called jet, and seen by the detector as a set of tracks and energy deposits in the calorimeters, as shown in Figure 4.6.

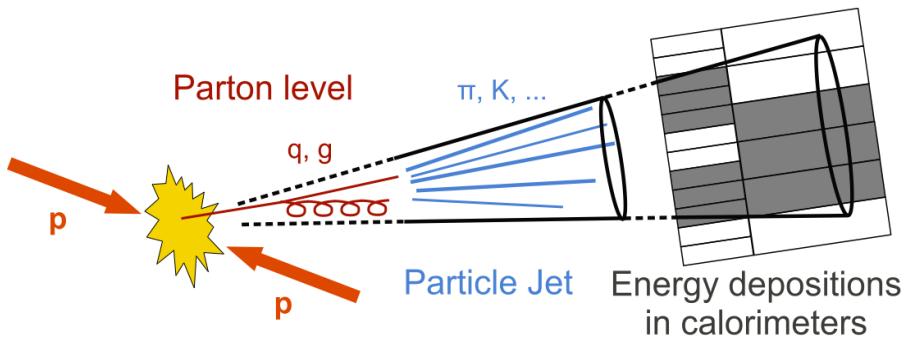


Figure 4.6: Schematic representation of the typical development of a jet within the CMS detector.

Several algorithms can be used to reconstruct the jets by linking the information coming from the tracker and the calorimeters, but the most used tool in CMS is the so-called anti- $k_T$  algorithm, able to cluster all the charged and neutral hadrons along with the eventual non-isolated photons or lepton produced and merge them into a single jet [107]. Its main objective is to compute the energy and direction of the original quark as precisely as possible. This is actually the best algorithm

developed so far to resolve jets, but the worst for studying jet substructure due to its clustering preference; in this case, other algorithms can be applied.

To perform such a job, sequential clustering algorithms such as this one rely on the value of two distances:  $d_{ij}$ , the distance between two particles  $i$  and  $j$  that need to be clustered and  $d_{iB}$ , the distance between the particle  $i$  and the beam axis  $B$ . As seen in Equation 4.1, these distances can be computed using different variables such as  $\Delta R_{ij}^2 = (\eta_i - \eta_j)^2 + (\phi_i - \phi_j)^2$ , the distance between  $i$  and  $j$  in the  $(\eta, \phi)$  space, the  $p_T^2$  of each particle and the clustering algorithm radius parameters determining the final jet size and usually set to 0.4 by the CMS collaboration. This distance parameter defines a cone in which the momenta of all the particles is summed to get the momentum of the jet itself.

$$\begin{cases} d_{ij} = \min \left( \frac{1}{p_{T,i}^2}, \frac{1}{p_{T,j}^2} \right) \Delta R_{ij}^2 \\ d_{iB} = \frac{1}{p_{T,i}^2} \end{cases} \quad (4.1)$$

The algorithm works by looking at all the  $i, j$  combinations, comparing the distances  $d_{ij}$  and  $d_{iB}$  until only jets and individual unclustered particles are present in the event:

- If  $d_{ij}$  is smaller than  $d_{iB}$ , then  $i$  and  $j$  are combined into a single particle ( $ij$ ) by summing their 4-vectors and both are removed from the list of particles to be clustered.
- If  $d_{iB}$  is smaller than  $d_{ij}$ , then  $i$  is considered to be the final jet and is therefore removed from the list of jet candidates as well.

Several corrections, mainly the so-called Jet Energy Corrections (JECs) [108], are then usually applied to the jets constructed using this algorithm in order to take into account several parameters such as the non-linearity of the response of the calorimeter, the electronic noise, the PU effects and the dependence of the reconstruction on the jet flavor. This typically introduces a source of systematic uncertainty that will be taken into account and discussed in Section 8.2.

The efficiency of the PF algorithm for jet identification and reconstruction has been checked using simulation, as shown in Figure 4.7. This study clearly shows that between 95 and 97% of the energy of the PF jet candidates can be reconstructed, compared to a 40-60% reconstruction efficiency using only the calorimeters data, and that this algorithm also leads to a gain in resolution up to a factor 3, depending on the jet  $p_T$ .

#### 4.4.1 B-tagging

Jets coming from bottom quarks are usually interesting to study in many different fields of particles physics, such as in this analysis which relies heavily on the number of b-jets produced to define the control and signal regions, as will be discussed in Chapter 6.

This specific kind of jets can be distinguished from other jets because of the relatively long lifetime of the bottom quark that produces in the detector a secondary vertex displaced by up to a few millimeters with respect to the PV, as shown in Figure 4.8; and this gives a perfect way to discriminate b-jets and jets coming from light quarks. Another consequence of the large mass of the bottom quark is that a large number of particles is typically present inside this particular kind of jets and that the decay of the bottom quark even leads to the apparition of soft leptons in the decay chain in around 20% of the cases.

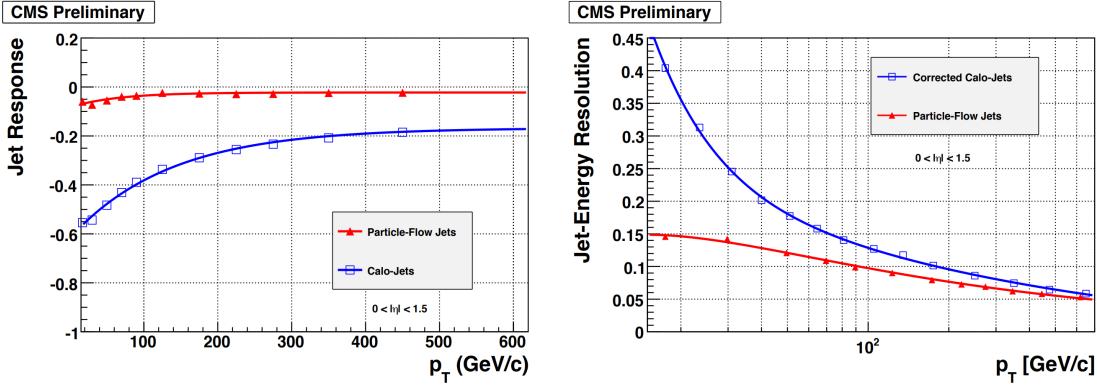


Figure 4.7: Comparison of the jet energy response (percentage of reconstructed energy, on the left) and jet energy resolution (on the right) for dijets simulated events in the barrel for jets reconstructed using only the calorimeters (in blue) and jet candidates from the PF algorithm (in red) [109].

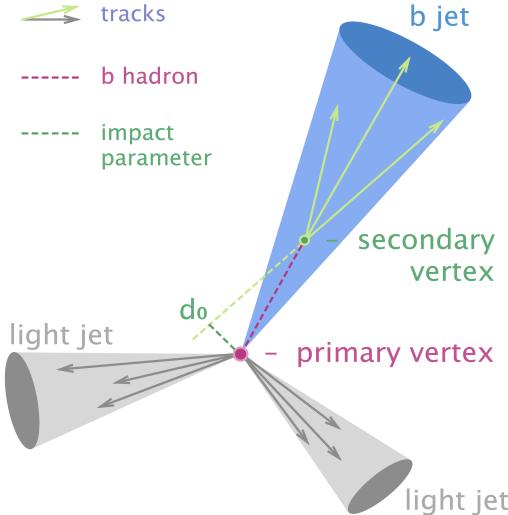


Figure 4.8: Schematic representation of the production of a b-jet originating from a slightly displaced secondary vertex.

Because of these specific properties, an algorithm can distinguish between jets coming from a bottom quark or from a lighter quark, and this will be a key point in this analysis. In our case, this discrimination is additionally optimized by using a multivariate technique able to combine all the discriminating power of the previous typical characteristics of any heavy flavour jet in the best way possible after reconstruction of all the vertices of the event. The main objective of the algorithm is to be able to identify b-jets as efficiently as possible while reducing the risk of possible misidentification of a jet.

In this analysis, the deep Combined Secondary Vertex (CSV) algorithm able to combine the information on the secondary vertex with the one on the track impact parameters and based on a Analysis Neural Network (ANN), has been used to identify such b-jets. The performance of this method can be observed in Figure 4.9, where we can see that this deep CSV algorithm is one of the best b-jets identification algorithms, depending on the phase space, while keeping a relatively low misidentification rate for light-flavor jets (u, d, s and gluons).

Different WPs are then also made available in order to obtain the desired combo b-jet identification

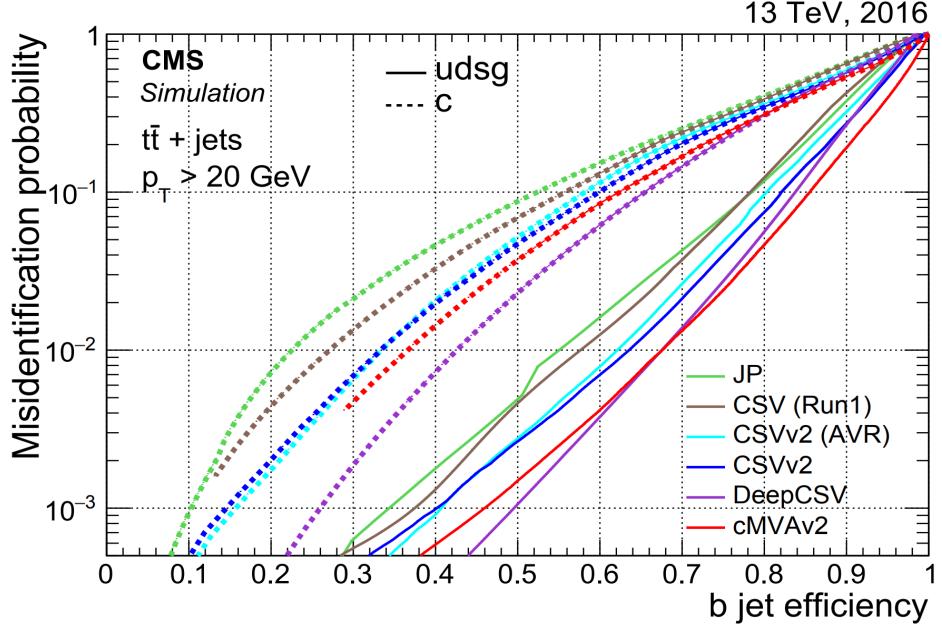


Figure 4.9: b-jets identification efficiency and misidentification rate considering different b-taggers, including the deep CSV b-tag used in this analysis [110].

efficiency/misidentification rate. The loose, medium and tight b-jets WPs have been developed in such a way to limit this misidentification rate of a light jet as a b-jet to 10%, 1% and 0.1% respectively.

## 4.5 Missing Transverse Momentum (MET)

Since the  $pp$  collisions happen mostly head-on, we know that the initial total transverse momentum of the event is exactly equal to 0 before the collision and we expect that it stays 0 afterwards because of the momentum conversation. This statement is not totally true though since we are aware of several effects that could induce an imbalance in this transverse momentum, as shown in Figure 4.10:

- Even though the CMS detector has been carefully designed, some particles could be created outside of its acceptance and therefore escape the detection (a particle can for example be created with such a boost that it could be emitted back to the beam pipe itself, making it impossible to detect it).
- Because of their extremely low interaction cross-section, SM neutrinos are expected to escape the detector with some energy while staying completely undetected.
- The finite momentum resolution of the detector can also lead to some inaccuracies in the measurement of the transverse momentum of all the particles created, leading to an instrumental MET in some cases.
- Some events also present an anomalous MET measurement which can arise because of a variety of reconstruction failures or malfunctioning detectors [111].
- Finally, the eventual exotic weakly interacting particles produced, such as Dark Matter (DM), are typically expected to leave some MET in the detector as well.

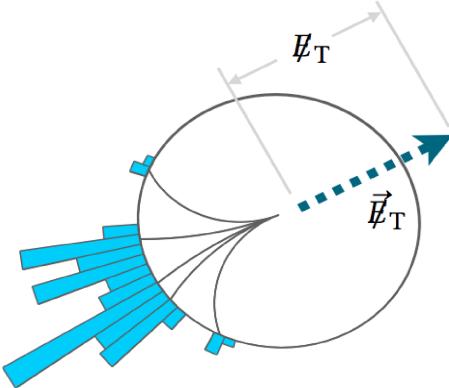


Figure 4.10: Schematic representation of the MET.

The  $E_T^{\text{miss}}$  variable, defined in Equation 4.2 as the negative sum of the transverse momentum of all the particles  $j$  of the event, accounts for this eventual imbalance in the transverse momentum and is therefore a key variable in some of the analyses searching for new BSM physics, in this particular case when such new physics is not expected to interact with the detector.

$$E_T^{\text{miss}} = \vec{p}_T^{\text{miss}} = - \sum_j \vec{p}_{T,j} \quad (4.2)$$

Different algorithms can be used in order to reconstruct this variable, the most common being [114]:

- The **particle flow MET** (PFMET), including all the information of the detector (as opposed to the calorimeter or tracker MET, for example) and only the PF reconstructed objects to estimate the MET value. This is the typical variable used in most of the analyses today, because of its simple, robust, yet very efficient estimate of the MET spectrum.
- The **Pileup Per Particle Identification (PUPPI) MET** [112], has been developed on top of the PFMET in order to further reduce the dependence on the pileup of this variable by using local shape information around each PF candidate in the event along with event PU properties and tracking information.

For this analysis, we decided to use the PFMET for simplicity after observing that using the PUPPI MET did not seem to give a better agreement between data and Monte Carlo (MC) in most of the regions defined, as shown in Figure 4.11 in a 2018 Drell-Yan (DY) inclusive region. In particular, the Type I corrected MET is used in this work, replacing the vector sum of transverse momenta of particles which can be clustered as jets with the vector sum of the transverse momenta of the jets to which JECs are applied [113].

Several corrections need to be applied to this spectrum to filter anomalous high MET events arising because of a variety of reconstruction failures induced by the detector due to several effects, such as the electronic noise and eventual dead cells in the calorimeters or the presence of an eventual beam halo particles from the LHC itself, leading to a global miscalculation of the final energy of the event. These filters, detailed in Table 4.1, are extremely important, especially in the end of the MET spectrum, as observed in Figure 4.12.

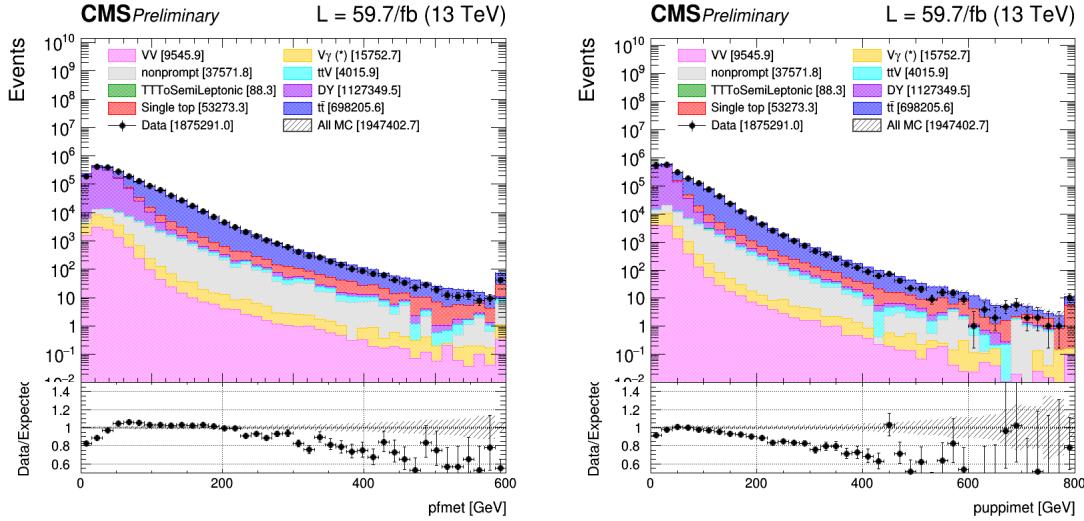


Figure 4.11: PFMET (on the left) and PUPPI MET (on the right) distributions observed in a 2018 DY inclusive control region.

Filter name	Applied to data	Applied to MC
Flag_goodVertices	✓	✓
Flag_globalSuperTightHalo2016Filter	✓	✓
Flag_HBHENoiseFilter	✓	✓
Flag_HBHENoiseIsoFilter	✓	✓
Flag_EcalDeadCellTriggerPrimitiveFilter	✓	✓
Flag_BadPFPixelFilter	✓	✓
Flag_ecalBadCalibFilterV2 <sup>†</sup>	✓	✓
Flag_eeBadScFilter	✓	—

<sup>†</sup> applied only to 2017 and 2018.

Table 4.1: MET filters applied to events selected in data and to simulated events, an hyphen (—) indicating the filter is not applied.

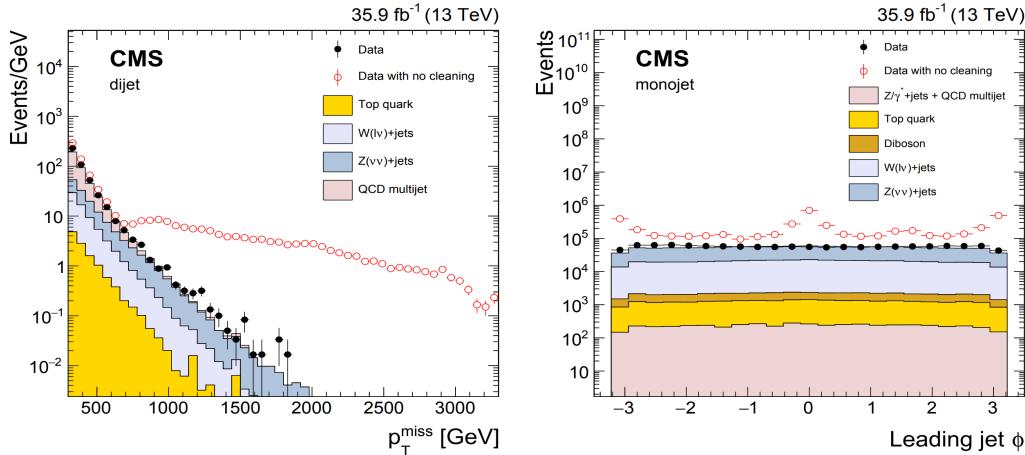


Figure 4.12: MET (on the left) and jet (on the right)  $\phi$  distributions with and without MET filters applied [114].

## 4.6 Top reconstruction

Although not formerly a part of the PF algorithm and done offline, the kinematic reconstruction of the  $t\bar{t}$  system is still an extremely important part of this analysis. Indeed, many variables allowing us to discriminate the signal from the different background processes introduced in Section 7.1 are spin correlated variables, sensitive to the nature of the DM mediator but which require knowledge of the top quark and anti-quark four-momenta, not immediately available without such complete reconstruction of the  $t\bar{t}$  system.

However, this reconstruction from channels containing two leptons is typically quite challenging, given the fact that neutrinos are not directly observed, meaning that the only observable about the neutrinos is the vectorial sum of their transverse momentum which is inferred from the total momentum imbalance of the event and which frequently has a bad resolution. Additionally, the determination of the individual momenta of each neutrino require advanced computation techniques, either numerical or analytical.

### 4.6.1 Numerical and analytical top reconstruction

Two main methods exist in order to solve this reconstruction problem:

- The **Sonnenschein numerical method** [115] first of all relies on the kinematics of the system and on the expression of the four-momenta of the different particles involved in the top quark decay chains (in this case, we consider only the chain  $t \rightarrow bW \rightarrow bl\nu$  for the two tops, as previously explained), as expressed in Equations 4.3a to 4.3c, if we assume that the MET of the event is coming only from the two neutrinos produced.

$$\begin{cases} p_x^{\text{miss}} = p_{\nu_x} + p_{\bar{\nu}_x} \\ p_y^{\text{miss}} = p_{\nu_y} + p_{\bar{\nu}_y} \end{cases} \quad (4.3a)$$

$$\begin{cases} m_{W+}^2 = (E_{l+} + E_\nu)^2 - (\vec{p}_{l+} + \vec{p}_\nu)^2 \\ m_{W-}^2 = (E_{l-} + E_{\bar{\nu}})^2 - (\vec{p}_{l-} + \vec{p}_{\bar{\nu}})^2 \end{cases} \quad (4.3b)$$

$$\begin{cases} m_t^2 = (E_b + E_{l+} + E_\nu)^2 - (\vec{p}_b + \vec{p}_{l+} + \vec{p}_\nu)^2 \\ m_{\bar{t}}^2 = (E_{\bar{b}} + E_{l-} + E_{\bar{\nu}})^2 - (\vec{p}_{\bar{b}} + \vec{p}_{l-} + \vec{p}_{\bar{\nu}})^2 \end{cases} \quad (4.3c)$$

In this case, we therefore have 6 equations to solve and exactly 6 unknowns (since the energy of the neutrinos is considered equal to their momentum because of their extremely low mass and if we assume the W boson and top quark masses to be known and fixed) corresponding to the three momentum components of each neutrino produced, a problem that can in principle be solved, leading to a quartic equation in  $p_{\nu_x}$ , analytically solvable but quite ambiguous given the variable number of solutions of such equation (typically, the solution giving the lowest invariant mass for the  $t\bar{t}$  system is then chosen).

- The **Betchart analytical method** [116] on the other hand is able to describe the decay  $t \rightarrow bl\nu$  using this time a geometric approach. This method was chosen in this analysis because it offers the following advantages:

- With this method, the invariant mass constraints from the top quark and the W boson are both exact and do not suffer the same kind of ambiguity as observed previously.  
**FIXME: Talk to Pablo about this (cf. email 21/08)**
- The solution set for each neutrino momentum in this case is an ellipse that can be described precisely and from which the solutions for the neutrino momenta can be reduced to a discrete set of values, which will be helpful when trying to estimate the  $p_T$  of the mediator of the DM signals considered, while also giving us information about the precision of this measurement.
- The results obtained here can be useful for other event topologies featuring similar kinematic constraints as well.

Basically, this methods relies on two observations constraining the geometrical shape of the W boson momentum vector. First of all, the decay of top quark constrains this vector to an ellipsoidal surface of revolution about an axis coincident with the bottom quark momentum. The decay of the W boson itself on the other hand additionally constrain this vector to another ellipsoidal surface of revolution about an axis matching the momentum of the resulting charged lepton. The W boson momentum vector will then be defined by the intersection of the surfaces given by these two constraints, resulting in an ellipse in the phase space. The neutrino momentum vector can then be expressed as a translation of this ellipse, using a parametric expression, as described in [116].

In the two neutrinos final state, it is then possible to show that the elliptical solution sets for the neutrino momenta ( $\nu_\perp, \bar{\nu}_\perp$ ) respective to the two top quarks decaying to leptons are given by Equation 4.4, where  $N_\perp$  and  $\bar{N}_\perp$  are the solution ellipses of the (anti)neutrino in the transverse plane and expressed in the laboratory coordinate system.

$$\begin{cases} \nu_\perp^T N_\perp \nu_\perp = 0 \\ \bar{\nu}_\perp^T \bar{N}_\perp \bar{\nu}_\perp = 0 \end{cases} \quad (4.4)$$

Given the fact that the measured components ( $\cancel{x}, \cancel{y}$ ) of the MET are the sum of the  $\nu_\perp$  and  $\bar{\nu}_\perp$  components, they can be related by Equation 4.5.

$$\bar{\nu}_\perp = \begin{pmatrix} -1 & 0 & \cancel{x} \\ 0 & -1 & \cancel{y} \\ 0 & 0 & 1 \end{pmatrix} \nu_\perp \equiv \Gamma \nu_\perp \quad (4.5)$$

The solutions for the momenta of the neutrinos will then be given by the intersections of these two ellipses giving either zero, two or four solution pairs ( $\mathbf{p}_\nu, \mathbf{p}_{\bar{\nu}}$ ), as shown in Figure 4.13. If the two ellipses do not intersect, a  $\chi^2$  method can be used to check the compatibility between the solution obtained and the standard  $t\bar{t}$  process, the best solution being defined as the point of closest approach between the ellipses.

So far we have seen two methods able to reconstruct the individual momentum of both neutrinos. However, both these methods usually assume that the MET is only coming from these neutrinos and this assumption is no longer verified for our signals, for which the DM particles produced will contribute by a significant amount to the global MET, resulting in a slightly lower reconstruction efficiency compared to the one obtained considering only the standard  $t\bar{t}$  process. The method used then needs to be slightly adapted to our particular case.

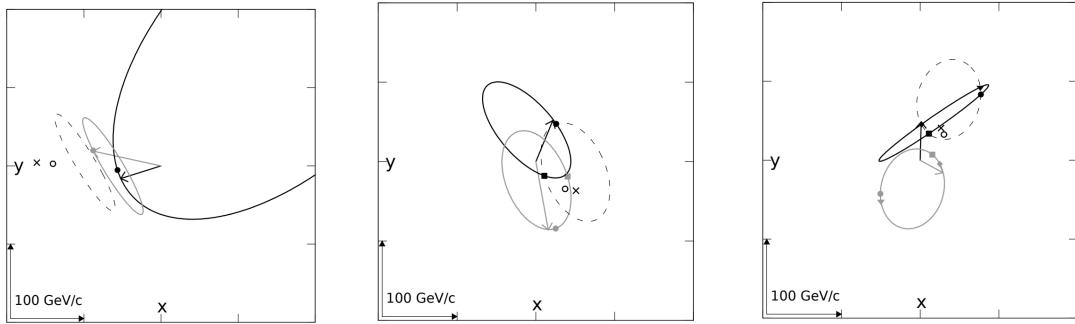


Figure 4.13: Three events constraining the neutrino and antineutrino momenta (black and grey arrows, respectively) resulting in 0 (on the left), 2 (on the center) or 4 (on the right) solutions. The dashed ellipses is obtained by using the additional constraint according to which measured MET is equal to the sum of neutrino transverse momenta [116].

#### 4.6.2 Top reconstruction with additional dark matter

In this particular case, including an additional contribution  $\phi_\perp = (\phi_x, \phi_y, 0)^T$  to the MET is equivalent to slightly modifying the Equation 4.5 to obtain Equation 4.6.

$$\bar{\nu}_\perp = \begin{pmatrix} \bar{\nu}_x \\ \bar{\nu}_y \\ 1 \end{pmatrix} = \begin{pmatrix} -1 & 0 & \cancel{x} \\ 0 & -1 & \cancel{y} \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \nu_x + \phi_x \\ \nu_y + \phi_y \\ 1 \end{pmatrix} \equiv \Gamma(\nu_\perp + \phi_\perp) \quad (4.6)$$

By substituting these values into Equation 4.4, we then get a new relation between the two ellipses for the neutrino,  $N_\perp$ , and for the antineutrino,  $\bar{N}_\perp$ , given by Equation 4.7.

$$\nu_\perp^T \Gamma^T \bar{N}_\perp \Gamma \nu_\perp + \nu_\perp^T \Gamma^T \bar{N}_\perp \Gamma \phi_\perp + \phi_\perp^T \Gamma^T \bar{N}_\perp \Gamma \nu_\perp + \phi_\perp^T \Gamma^T \bar{N}_\perp \Gamma \phi_\perp = 0 \quad (4.7)$$

This modification does make the reconstruction a bit more complicated by adding crossed terms between  $\nu_\perp$  and  $\phi$ , which modify the phase space of solutions. Even though the reconstruction gets more complex in this case, performing it is extremely important because it provides us a way to determine two excellent discriminating variables:

- First of all, the so-called **dark  $p_T$**  is a variable estimating the value of the  $p_T$  of the eventual mediator of the interaction. Indeed, if the system considered does not admit any solution, then there is no intersection between the ellipses formed in the MET phase space. However, in this case, the distance between the ellipses gives us a way to approximate by how much we miss a perfect standard  $t\bar{t}$  reconstruction, and it is actually possible to show that the distance between the center of both ellipses is related to the  $p_T$  of the mediator of the interaction [117]. This method being completely analytical, estimating the distance between the ellipses is trivial and is giving us the first background-signal discriminating variable that will be developed later in Section 7.1.
- The second interesting variable is the so-called **overlapping factor  $R$** , defined in Equation 4.8, where  $l_1$  and  $l_2$  are the two sizes of the ellipses measured along the axis joining their centers, and  $d$  is the distance between these centers.

$$R = \frac{l_1 + l_2}{d} \quad (4.8)$$

By its definition, this factor is able to take into account not only the distance between these two ellipses but also their respective sizes and this is extremely important since, at the end of the day, we want an event having small ellipses far away from each other to have a much higher weight than an event featuring two incredibly large ellipses and therefore less significant when trying to distinguish the background and signal processes. Both these extreme cases are represented in Figure 4.14.

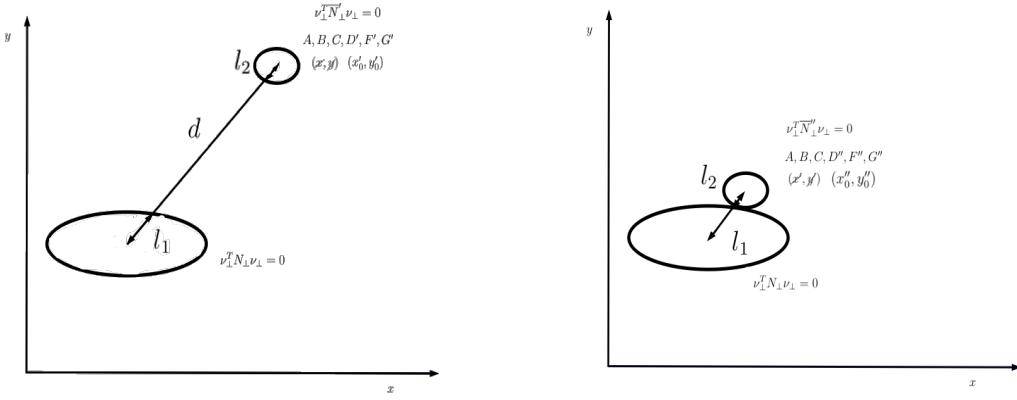


Figure 4.14: Schematic representation of the two extreme cases that can be observed when defining the overlapping factor: the reconstruction of a system with (on the left) and without (on the right) the presence of DM [117].

### 4.6.3 Top reconstruction in practice

The idea of the algorithm is quite simple: it basically takes as input the 4-momenta of the two leptons, the two b-jets, and the two components of the transverse missing momentum and performs all the calculations just described while returning the optimal momenta value for the two neutrinos. However, several complications quickly appear when solving this problem in practice.

First of all, as we already saw, when the two ellipses intersect each other, we can observe either 2 or 4 solutions, even though only one actually corresponds to the physical momenta of the neutrinos. The solution giving the lowest possible invariant mass for the  $t\bar{t}$  system is then simply chosen and taken in consideration for the analysis since this optimized the probability of making the right choice: indeed, it has been shown [118] that in 85% of the cases, the solution satisfying this requisite actually matches the true simulated kinematics.

Then, even though the process studied is always the same and is supposed to give the same final state, each recorded event typically comes with a different number of leptons, jets and b-jets, mostly due to mismeasurements and b-tagging inefficiencies. This means that this algorithm needs to be applied several times, considering all the possible leptons and (b-)jets combinations, taken two by two. In practice, the following categories are therefore defined:

- If exactly 0 b-jets are observed, the event is not considered in this analysis, according to the selection of the signal regions performed and explained later in Chapter 6.
- If 2 or more b-jets are observed, combinations between the two leptons, the first b-jet (ordered in  $p_T$ ) and the other b-jets of the event are considered.
- Finally, if exactly 1 b-jet is observed, then it is kept and used as input, while all the non b-tagged jets are considered as the second b-jet candidate. In this case as well, all the possible combinations between the two leptons and this set of jets are then tested.

When taking into account all these possible combinations, a reconstruction efficiency of  $\sim 61\%$  has been achieved when considering standard  $t\bar{t}$  MC samples.

Finally, we know any measurement made is not perfect and comes with uncertainties. The impact that imperfectly measured kinematic variables can have on the top reconstruction process can be estimated through the **smearing** method, by repeating the reconstruction 100 times for each combination previously defined, by updating in each iteration several parameters:

- The energy of the jets are updated within their respective uncertainties, using a random number drawn from a Gaussian distribution whose mean value  $\mu$  is null and whose standard deviation  $\sigma$  has been estimated to 0.3 using the jet energy corrections information.
- The energy of the leptons is also updated using a random correction factor generated from the distribution shown in Figure 4.15, computed in MC and representing the energy deviation between the generation and reco information.

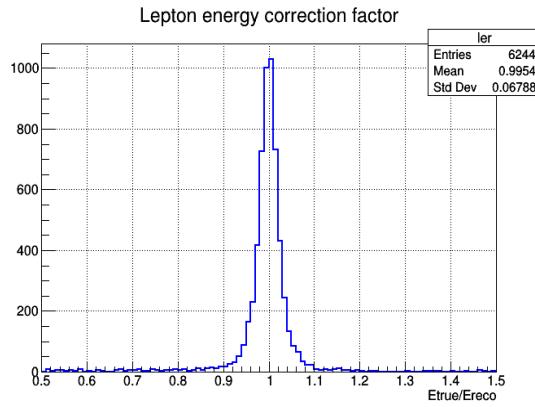


Figure 4.15: Definition of the correction factor used in the smearing of the leptons energy [119].

- The MET is also recomputed each time using both these newly calculated parameters, according to Equation 4.9.

$$E_{T_{x,y}}^{\text{updated}} = E_{T_{x,y}}^{\text{reco}} + \sum_{\text{jet}_i=1}^2 \Delta(p_{x,y}^{\text{jet}_i}) + \sum_{\text{lep}_i=1}^2 \Delta(p_{x,y}^{\text{lep}_i}) \quad (4.9)$$

- An angular smearing is also performed, using the angles  $\alpha$  and  $\omega$  shown in Figure 4.16.

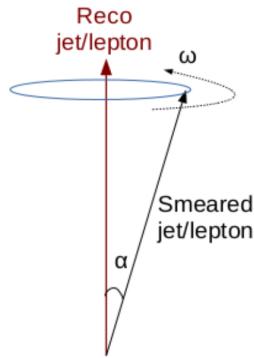


Figure 4.16: Graphical definition of the two angles  $\alpha$  and  $\omega$  used in the smearing of the leptons and jets directions [119].

The  $\omega$  angle is taken from a simple uniform distribution  $[0, 2\pi]$  while the angle  $\alpha$  is taken as a random number generated from the distributions of the angle between the particle and

detector level direction, as shown in Figure 4.17. These distributions have been obtained in MC using Equation 4.10, relating the generation and reco normalized vectors  $\hat{p}$  representing the direction of the momentum of a given object (lepton or b-jet).

$$\begin{cases} \alpha = \arccos(\hat{p}^{\text{reco}} \cdot \hat{p}^{\text{gen}}) \\ \hat{p} \equiv \frac{\vec{p}}{|\vec{p}|} = (\cos(\phi) \cdot \sin(\theta), \sin(\phi) \cdot \sin(\theta), \cos(\theta)) \end{cases} \quad (4.10)$$

- Finally, the mass of the W boson, entering the reconstruction as a parameter and not a universal constant in nature, is updated each time by generating a random number according to a Breit-Wigner distribution of mean 80.38 GeV and width 2.08 GeV.

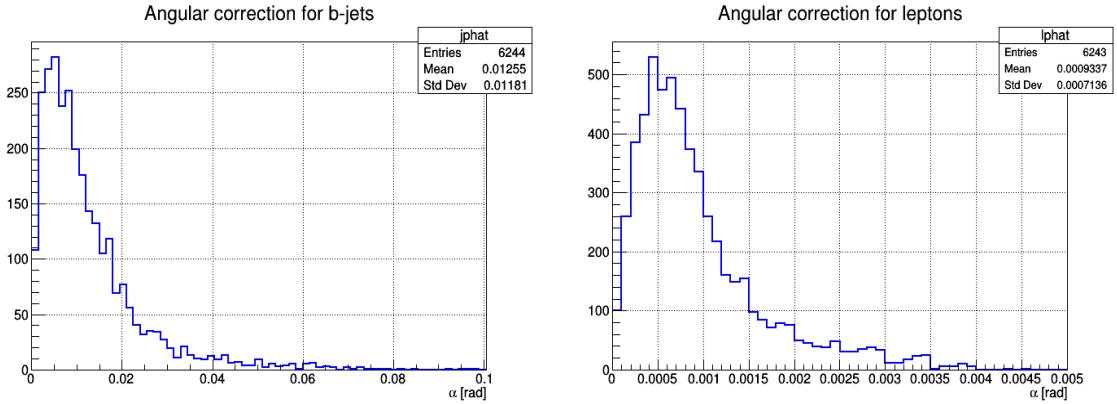


Figure 4.17: Simulated distributions of the  $\alpha$  angle between the particle and detector level direction for b-jets (on the left) and leptons (on the right).

During each iteration of the smearing process, an individual weight  $w$  is computed from the true invariant mass  $m_{lb}$  distribution previously obtained using generation. Both the W mass and true  $m_{lb}$  distributions obtained are represented in Figure 4.18.

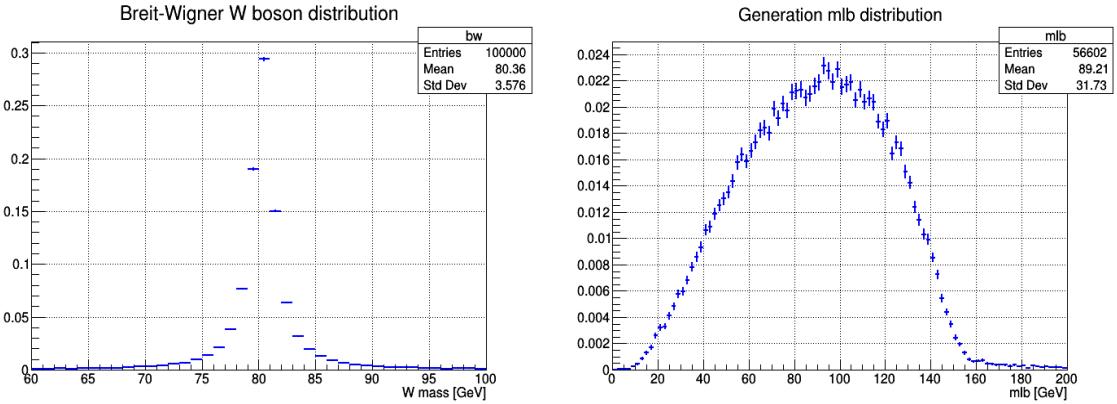


Figure 4.18: Breit-Wigner spectrum obtained when generating randomly the mass of the W boson (on the left) and true  $m_{lb}$  distribution obtained using the generation information from the standard  $t\bar{t}$  process (on the right).

A global weight  $W$  is then assigned to each lepton/b-jet combination by summing the weight  $\sum w(m_{lb}) \cdot w(m_{\bar{b}})$  obtained for each smearing iteration, and the combination with the largest weight is considered to be the solution of the event. This weight distribution will be another variable used to discriminate the signal from the backgrounds processes, since the higher this variable is, the larger the likelihood for the event to be an authentic  $t\bar{t}$  event is.

Even though this smearing increases the efficiency of the reconstruction to  $\sim 90\%$  when considering the  $t\bar{t}$  process, this kinematic reconstruction algorithm does not allow us to find a solution to all the events. If an event does not present any solution, then all its variable which do not depend on the neutrino momenta are still considered for the analysis, while the variables which do require the reconstruction information are set to non-physical default values.

---

---

# Chapter 5

---

## Data, signals and backgrounds

In order to find a possible hint for the production of DM in the LHC collisions considering our signal models of interest, briefly described in Section 2.6, the data collected needs to be compared with Monte Carlo (MC) simulations produced in a central way for each SM process. Indeed, any deviation of the data observed with respect to what we expect to see, obtained from these MC simulations, might be the sign of some BSM physics.

All of the steps needed to mathematically simulate the  $pp$  collisions of the LHC and to take into account the effect of the detector on the particles produced will first of all be introduced in Section 5.1. Then, the different formats of files available to perform the analysis and the code used will be briefly introduced in Sections 5.2 and 5.3 and the different data samples collected during the Run II of operation of the LHC will be detailed in Section 5.4, while the signal models and samples considered in this particular analysis along with the MC samples used for the simulation of the different backgrounds will be introduced in Sections 5.5 and 5.6 respectively.

### 5.1 The Monte-Carlo simulation method

As previously explained, the generation of MC simulations for the most common SM processes is a crucial step of any analysis because they are considered to be the reference to which the data collected is compared in order to try and find some discrepancies, which could be the sign of the existence of BSM physics. Searches for exotic physics therefore heavily depend on these simulations, which need to be generated with great care and to which a large uncertainty is typically associated since the collision between the partons of two protons and the interaction between the particles produced and the detector itself are extremely complex by nature.

The basic idea of the MC simulation consists in using a random number generator to simulate the randomness of nature and produce as many events as computationally possible for all the SM processes, taking into account the probability density functions of these processes. This is performed by specific softwares called **event generators** and it is important to note that since we usually don't know everything about the SM or BSM process being generated, the perfect event generator does not exist.

To make the generation of such simulations a bit easier, the description of a typical  $pp$  collision can usually be divided into several steps, as shown with the color code used in Figure 5.1. The typical approximations used to make this kind of simulation possible from the computational point of view will also be briefly introduced at this point.

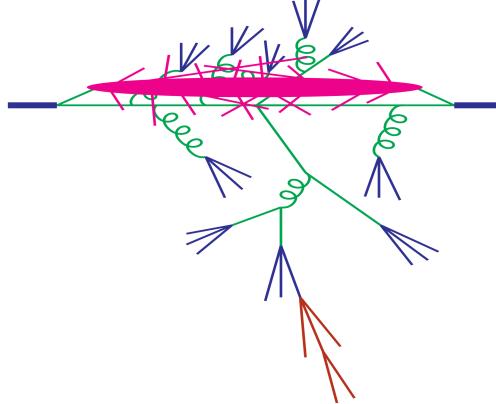


Figure 5.1: Structure of a  $pp$  collision and different steps of the MC simulation used by the event generators, such as the parton shower (in green), the Underlying Event (UE) (in pink), the hadronization (in blue) and the decay of unstable particles (in red) [120].

### Hard scattering

A typical  $pp$  collision at a center of mass energy  $\sqrt{s}$  is usually described by an event generator as the interaction between a parton  $i$  coming from one proton with a parton  $j$  coming from the other, leading to the production of a final state  $A$ , made out of  $n$  different particles. The total cross section of such process can be expressed with Equation 5.1 [121].

$$\sigma_A(s) = \sum_{i,j} \int \int dx_1 dx_2 f_i(x_1, \mu^2) f_j(x_2, \mu^2) \hat{\sigma}_{ij \rightarrow A}(\hat{s}, \mu^2) \quad (5.1)$$

In this equation, several variables have been introduced, such as:

- The artificial parameter  $\mu^2$  used as the delimitation between short and long range physics.
- The Parton Density Functions (PDFs)  $f_i(x, \mu^2)$  of both partons involved in the collision, giving the probability of finding in the proton a parton of flavor  $i$  (quark or gluon) carrying a fraction  $x$  of the proton momentum.
- The integrated parton-level cross section  $\hat{\sigma}_{ij \rightarrow A}$  describing the short range physics between the partons, taking into account the phase space and the matrix element obtained considering all the Feynman diagrams of a given process.
- The square invariant mass of the two partons  $\hat{s} = (p_i + p_j)^2$ .

Many algorithms have been developed in order to select a hard process  $ij \rightarrow A$  and determine its kinematics by solving this equation using different methods. The samples used in this work have actually been produced at different orders and by different hard scattering generators, such as MADGRAPH [123] (at LO) and POWHEG [124] and MC@NLO [125] (at NLO).

## Parton showers

The parton shower phase is then used to describe what happens to the incoming and outgoing partons after the initial collision that has just been described. The hard process induce by definition a large acceleration to the partons involved, which then tend to emit QCD radiation under the forms of gluons, just like accelerated electric charges do by emitting photons. However, the gluons emitted do have a color charge and can therefore emit further radiation until reaching such a low energy that they are able to form colourless hadrons, as discussed in Section 4.4. This process typically leads to the creation of the so-called **parton showers**, approximate higher-order real-emission corrections to the hard scattering, that need to be simulated by the event generators as well since they are an important part of the kinematics of the collision.

The parton showering then consists in simulating these showers for not only the final state particles produced by the hard scattering, but also for the particles in the initial state and for the remnants of the colliding protons, since gluons can actually be emitted by Initial State Radiation (ISR) and by these remnants themselves.

## Underlying Event (UE)

Once the hard scattering and all the possible gluon emissions simulated, the next step consists in considering the so-called **Underlying Event (UE)** arising from the parton showers just described and from the secondary collisions between partons not involved in the primary hard process, the so-called Multiple Parton Interactions (MPIs). The UE is usually responsible for the production of particles at low transverse momenta  $p_T$  that cannot be experimentally distinguished from particles produced from initial or final state radiation but still need to be accounted for and simulated.

These secondary collisions typically lead to the production of extra hadrons and therefore need to be simulated as well by events generators, usually by distributing the partons of the incoming protons in an area of  $1\text{fm}^2$ : an increased UE will be obtained when the so-called impact parameter, the distance between the parton and the centre of this area, is decreased, making the collision mostly central and almost head-on [126]. The UE is typically well simulated using softwares such as Herwig [127] and PYTHIA [128]. The spectrum for the generation of some variables in a top enriched sample can be found in Figure 5.2.

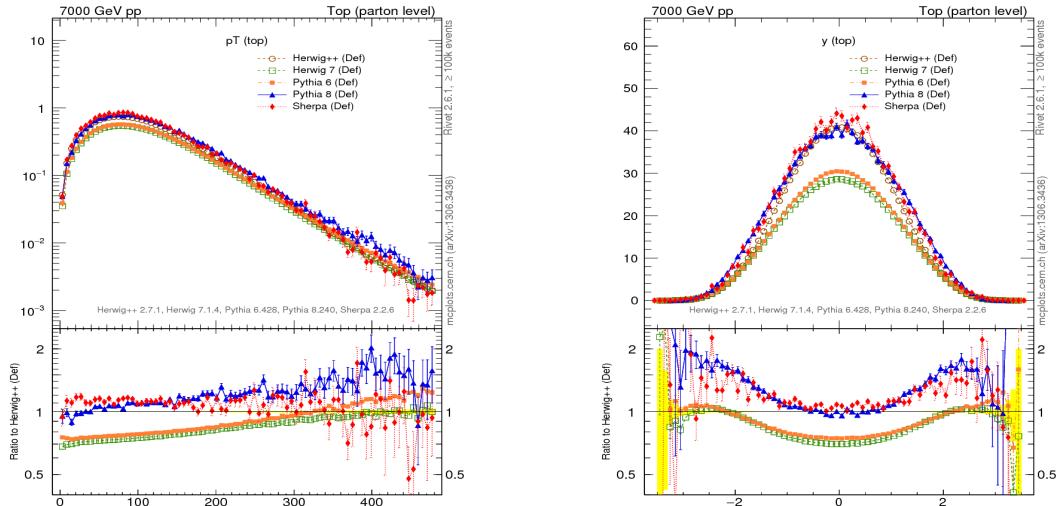


Figure 5.2: Top  $p_T$  (on the left) and rapidity (on the right) distributions obtained using different MC generators [129].

## Hadronization

Once all the primary and secondary collisions simulated, it is time for the event generators to simulate the **hadronization** and binding processes of the different coloured partons emitted into colourless hadrons, as explained in Section 4.4. This hadronization process happen at low energies, when the perturbation theory becomes invalid and the dynamics enter a non-perturbative phase, which leads to the formation of the observed final-state hadrons. Non-perturbative calculations then have to be used by the event generators in order to simulate this effect.

## Unstable particle decays

The last step of the MC generation consists in finding a model allowing the unstable hadrons created in the hadronization process to decay, and to study these decays. This is extremely important because experimental data clearly shows that a large fraction of the observed final state particles come from the decays of such excited hadronic states.

## Detector simulation

Once the event completely simulated using the event generators and the PU taken into account by reproducing the hard scattering process several times, another step is required: simulating the interaction between the "perfect" particles previously created and the "imperfect" CMS detector.

This is typically done by the GEANT4 software [130], able to model different effects, such as:

- Modeling of the interaction region;
- Modeling of the particle passage through the volumes that compose CMS detector and of the accompanying physics processes;
- Modeling of the effect of multiple interactions per beam crossing and/or the effect of events overlay (PU simulation);
- Modeling of the detector's electronic response.

This modeling accounts for all the cracks and for the disposition of the subsystems inside of the CMS detector. This software is for example able to model the interaction of the electrons with the tracker, responsible for the emission of bremsstrahlung photons, as explained in Section 4.3.2.

The results of the comparison between the output of two different versions of the GEANT4 software and prototypes of the CMS calorimeter in the test beam facility at CERN lead to comparable results, as shown in Figure 5.3.

However, the modeling of the detector is not perfect and not all the inefficiencies can be accounted for. In some cases, Scale Factors (SF) are then used to correct the MC simulations and correct some expected discrepancies between data and MC.

## 5.2 Files format

Once recorded (or simulated), the data (or MC) still needs to go under a complete post-processing in order to change its format and reduce the total size of the samples to be considered in the different

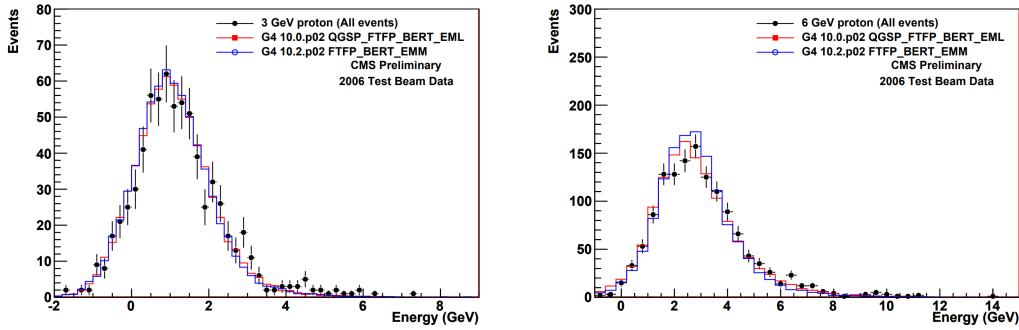


Figure 5.3: Proton energy distribution at 3 (on the left) and 6 (on the right) GeV compared for the test beam data (in black) and two different GEANT4 versions [131].

analyses. Different types of analyses are expected to need different levels of data reduction, so the data is usually accessible at different levels [132]:

- **Virgin-R<sub>A</sub>W**: used only in low rate runs with heavy ions collisions (10-15Mb/event)
- **R<sub>A</sub>W** : standard raw data event content (1Mb/event)
- **RECO**: detailed information on reconstructed physics objects (3Mb/event)
- **Analysis Object Data (AOD)**: physics objects used in analysis (400-500kB/event)

Two additional formats were introduced since the end of the Run I. First of all the miniAOD was introduced to reduce the size of the AOD by a factor 10 while retaining most of the information about all the particles that were created, without applying any further selection.

Because of the increased integrated luminosity collected by CMS over the last few years, a brand new file format featuring another reduction of the file size of a factor  $\sim 50$  was recently introduced: the nanoAOD, able to retain most of the information of each collision in around 1kB of data per event only. This reduction in size was achieved by optimizing the floating point of the variables, by not storing quantities that can be recomputed from the available information and by limiting the number of physics objects available, for example. This means that some low-level analyses cannot use this format to work, but it has been estimated that around 50-70% of the analyses performed at CMS can actually rely on such files in order for their work.

In this particular case, the 6th version of the nanoAOD, introducing a series of bug fixes and the latest jet energy corrections, was used for both the data and the MC samples (signal and backgrounds).

### 5.3 Analysis code

The code used for the event generation, simulation and reconstruction is the version 10\_4\_X of the official software of the CMS collaboration, called CMSSW [136]. This software contains the CMS Event Data Model (EDM) which is able to describe every event as a C++ object containing all the RAW and reconstructed information related to the collision. These object are stored using the ROOT file format [137], an analysis package written in C++.

Once all the different samples produced centrally up to the nanoAOD stage, another framework was put in place in order to do a post-processing of such samples, by selecting objects interesting for different dileptonic analyses, reducing therefore even more the size of the samples to be considered by selecting only events having 2 tight leptons. This selection will be detailed in Chapter 6. This *Latino* framework, written in python, is common to several different analyses and has been developed by tens of different people over the past few years, providing several tools to produce samples, read the files, apply different corrections to the MC samples and produce the histograms needed to perform a search such as this one.

## 5.4 Data samples

As already explained in Section 3.1.2, the data analyzed in this work has been taken at a center of mass energy of 13 TeV during the second part of the Run II of operation of the LHC.

During this period, an integrated luminosity of  $(35.9 \pm 0.9) \text{ fb}^{-1}$  (2016) [138],  $(41.5 \pm 1.0) \text{ fb}^{-1}$  (2017) [139] and  $(59.7 \pm 1.5) \text{ fb}^{-1}$  (2018) [140] has been collected, resulting in a total dataset of  $(137.1 \pm 2.0) \text{ fb}^{-1}$  recorded by the CMS detector and ready to be analyzed. This data has been obtained by combining a set of single and double lepton triggers that will be described in Section 6.1.1 and by taking care of avoiding any eventual double counting due to events present in different triggers. All the data samples considered for this analysis are listed in Appendix B.1.

## 5.5 Signal samples

Two different sets of MC signal samples have been produced centrally with MADGRAPH and PYTHIA8 at Leading Order (LO) for this analysis, corresponding to the  $t/\bar{t}+\text{DM}$  and to the  $t\bar{t}+\text{DM}$  signals. Different mass points were produced first privately and then centrally in both cases, considering different dark matter masses, from 1 to 51 GeV, and different scalar or pseudoscalar mediator masses depending on the model considered, ranging from 10 to 1000 GeV. In this context, 13 (17) different mass points have been produced for each mediator considered for the  $t/\bar{t}+\text{DM}$  ( $t\bar{t}+\text{DM}$ , respectively) signals, as listed in Appendix B.2.

The impact on the kinematics (in this particular case, on the spectrum of the pfMET) of these different mass points available can be observed in Figures 5.4 (scalar  $t/\bar{t}+\text{DM}$ ), 5.5 (pseudoscalar  $t/\bar{t}+\text{DM}$ ), 5.6 (scalar  $t\bar{t}+\text{DM}$ ) and 5.7 (pseudoscalar  $t\bar{t}+\text{DM}$ ). As expected from Table B.5, we can see first of all in these figures that the higher the mediator mass is, the lowest is its spectrum because of the lower cross section associated to the model. On the other hand, we can also observe that the mass of the DM itself has little to no impact regarding to the kinematics of the event, as expected given the decay of such exotic matter to a pair of invisible particles. Finally, we can also observe in these plots the difference in kinematics between the background (in blue) and the different signals, which will be the basis of the analysis performed.

## 5.6 Backgrounds prediction

Several different SM background processes have been considered for this analysis, all listed in Appendix B.3 and mostly estimated directly from MC. In this section, the main backgrounds to consider for this particular analysis will be reviewed:

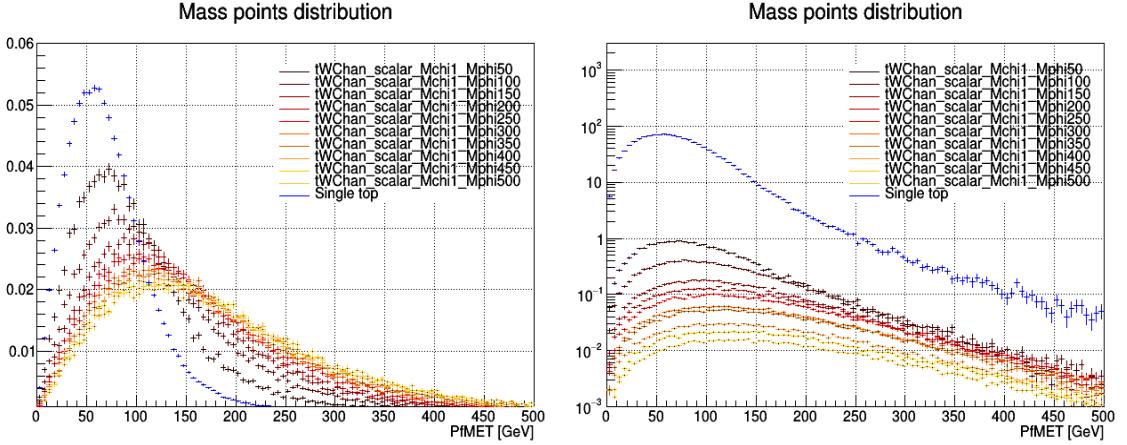


Figure 5.4: MET spectrum for several  $t/\bar{t}$ +DM **scalar** mediators, with (on the left) and without normalization (on the right).

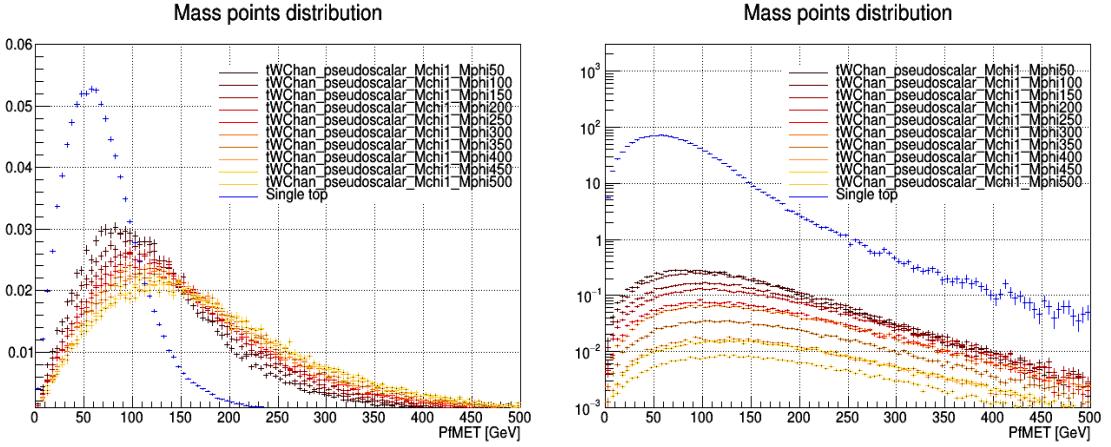


Figure 5.5: MET spectrum for several  $t/\bar{t}$ +DM **pseudoscalar** mediators, with (on the left) and without normalization (on the right).

- The major background of this analysis is the SM  $t\bar{t}$ , kinematically really close to the  $t\bar{t}$ +DM signal searched for (Section 5.6.1). The dilepton decay of the top quarks is obviously more important in this analysis, but its semi-leptonic is also considered to cover possible events in which the jet emitted is misreconstructed by the detector. The single top production is also important to consider, being kinematically close to our  $t/\bar{t}$ +DM signal (Section 5.6.1).
- Because of its huge cross section at 13 TeV, as shown in Figure 5.8, the DY process is also important to consider. However, this process is kinematically quite different to our signals and it is therefore relatively easy to reduce it in the actual Signal Regions (SRs) with some specific cuts, by removing the Z peak in the  $ee$  and  $\mu\mu$  channels and applying a cut on the dilepton stranverse mass, for example.
- The  $t\bar{t}+V$  ( $t\bar{t}+Z$  and  $t\bar{t}+W$ ) process, with a W or Z boson in the final state able to decay into neutrino(s) and therefore responsible for the production of some MET, has a typical event kinematic even closer to our signals than the other processes previously quoted, making it irreducible in most of the cases. This background is therefore extremely important in our signal regions, even though its low cross section does limit its impact.
- Finally, the non-prompt background is another important piece of some analyses mainly because of the particular data-driven method that can be used to compute its expected kinematic and yields, described in Section 5.6.4, since the MC are not fully reliable to describe

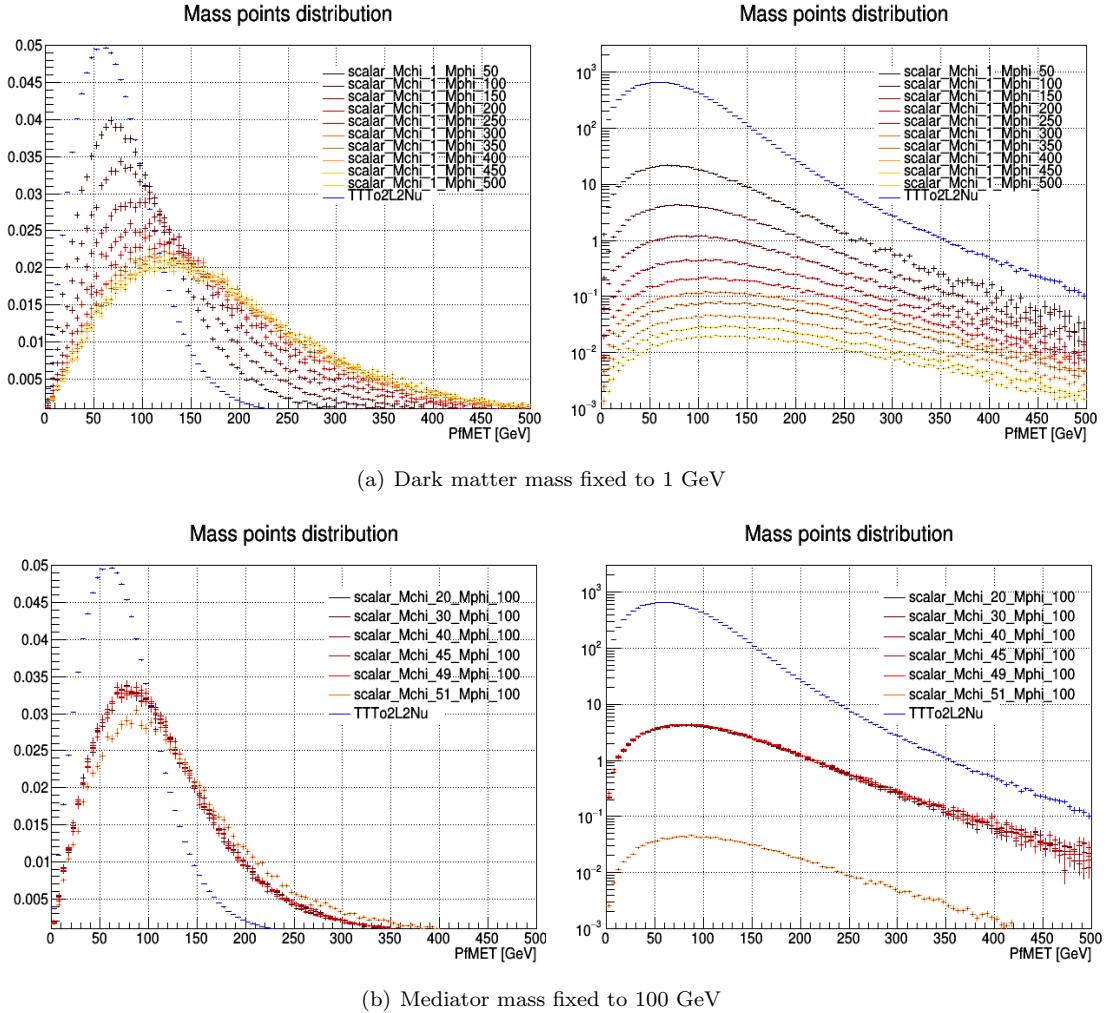


Figure 5.6: MET spectrum for several  $t\bar{t}$ +DM **scalar** mediators (on the top) and dark matter (on the bottom) masses, with (on the left) and without normalization (on the right).

such events. However, most of the non-prompt contamination in this analysis comes from semi-leptonic b-decays, expected to be well modeled. After detailed studies, we therefore decided to exclusively rely on MC simulations for this background as well.

Finally, some smaller backgrounds will be introduced in Section 5.6.5, such as the diboson and triboson production. All of the backgrounds used in this analysis are estimated directly from MC, even though some weights and corrections are typically applied to these MC samples, as will be detailed in Section 5.6.6. Ways to mitigate the impact these processes have on the signal regions will be presented in Chapter 6, and the actual impact of all these different processes on the different signal regions is shown in Table 5.1.

**FIXME:** Add corrected yields and percentages after unblinding

### 5.6.1 Top production

Because of the relatively high production cross section of top quarks at 13 TeV, the production of one or two top quarks, but without the production of associated DM, is obviously the dominant background in both searches.

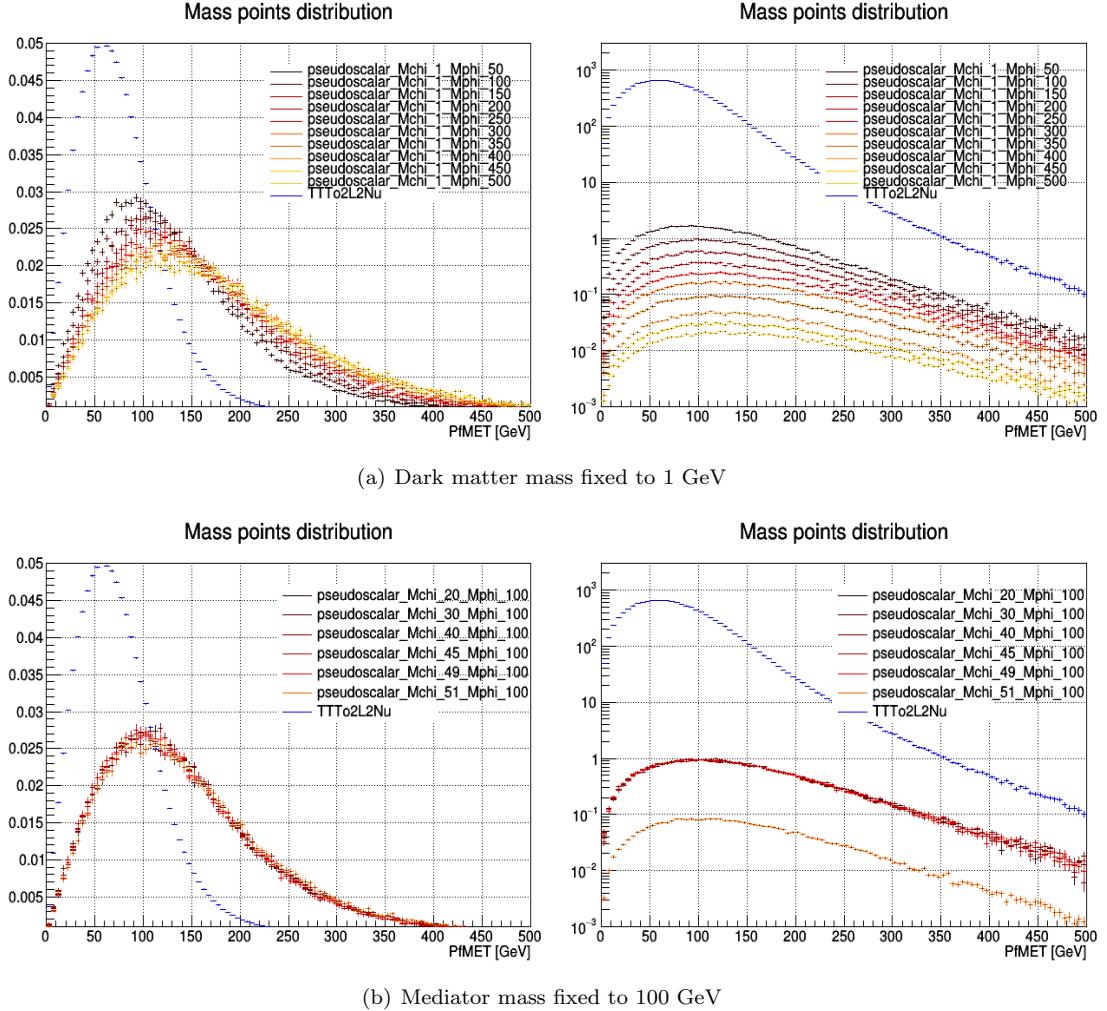


Figure 5.7: MET spectrum for several  $t\bar{t}$ +DM **pseudoscalar** mediators (on the top) and dark matter (on the bottom) masses, with (on the left) and without normalization (on the right).

Such backgrounds have kinematics quite close to the one expected for our signals. However, the additional MET expected because of the production of a pair of DM particles when considering the signals, leads to a few differences in kinematics, allowing for some discrimination between these kind of processes. To get the best discrimination possible, this work relies on the use of advanced Machine Learning (ML) techniques, as will be discussed in Chapter 7.

### The main background: $t\bar{t}$

Different Feynman diagrams contribute to this process at LO in a hadron collider, as shown in Figure 5.9. This background is estimated directly from MC in this analysis and checked in a dedicated control region. NNPDF3.0 [141] was used as the default PDF, while the default POWHEG setup with the PYTHIA8 CUETP8M2 (for 2016) and CP5 (for 2017 and 2018) tunes were used for the generation of this particular sample [142].

### Single top

Different Feynman diagrams also account for this process in the s-channel (Figure 5.10), t-channel and tW production modes (Figure 5.11), the latter being the dominant contribution in this anal-

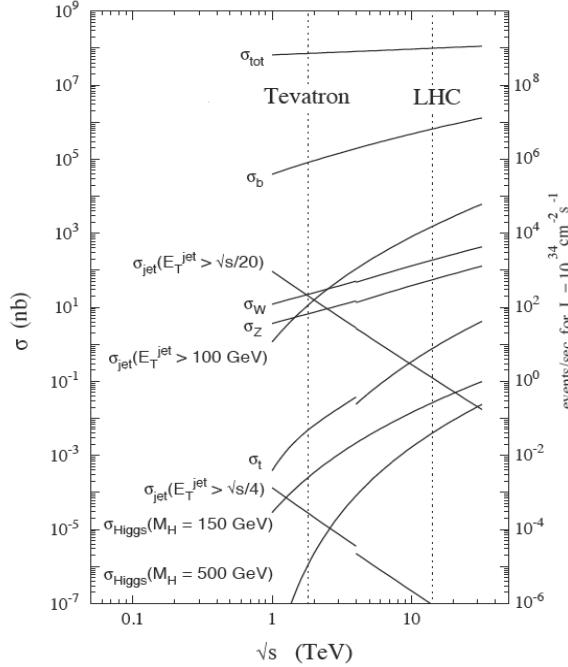


Figure 5.8: Production cross section of the most common SM processes considering different center of mass energies, such as the 13 TeV of the LHC.

Process	Scalar 100 GeV SR	Scalar 500 GeV SR	Pseudoscalar 100 GeV SR	Pseudoscalar 500 GeV SR
SM $t\bar{t}$	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)
Single top	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)
DY	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)
WW	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)
$t\bar{t} + V$	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)
Non-prompt	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)
Data	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)	XXXX (XX%)

Table 5.1: Number of yields and percentage of different processes in some of the 2018 signal regions.

ysis given its kinematics and cross-section. This process is also simulated using POWHEG (or MC@NLO, depending on the year and on the channel) with NNPDF3.0 and PYTHIA8 with the CUETP8M1 (for 2016) and CP5 (for 2017 and 2018) tunes.

### Top decay

As previously mentioned, the top is the heaviest particle of the SM and is expected to decay inside of the beam pipe itself, usually into a bottom quark, giving us a b-jet, and a W boson; this boson then decays itself into different channels even though only its leptonic decay is consider in this particular case. The decay considered for the top/anti-top produced is represented in Figure 5.12.

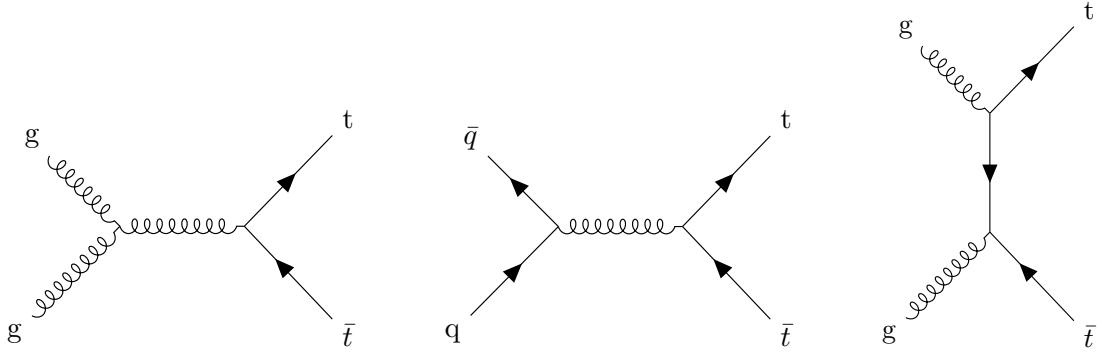


Figure 5.9: Main feynman diagrams for the production of the SM  $t\bar{t}$  process.

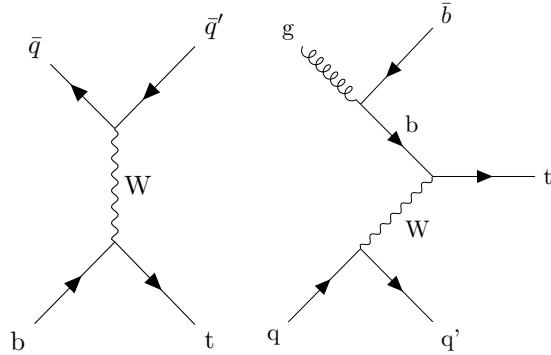


Figure 5.10: Feynman diagrams for the s-channel production mode of a single top quark.

### 5.6.2 Drell-Yan estimation

As previously mentioned, most of the DY, produced through the Feynman diagram represented in Figure 5.13, is not expected to survive the selection applied to the analysis. However, because of the huge cross section of this process, two to three orders of magnitude larger than the production of a top quark, a thorough description of such process is still extremely important. Expected to be found mostly in the  $ee$  and  $\mu\mu$  channels, it does survive even in the  $e\mu$  channels because of possible tau decays ( $Z\gamma^* \rightarrow \tau\tau \rightarrow e\mu\nu_e\nu_\mu\nu_\tau\nu_\tau$ ). These samples are generated with the MADGRAPH generator, MLM matching and interfaced to PYTHIA8 with the CUETP8M1 (for 2016) and CP5 (for 2017 and 2018) tunes for hadronization. The generation is then split into two distinct Z invariant mass ranges: 10-50 GeV and  $>50$  GeV. The HT-binned samples have been chosen over the inclusive ones in order to increase the statistics available.

The shape of this background is also estimated from MC but a correction factor to its normalization

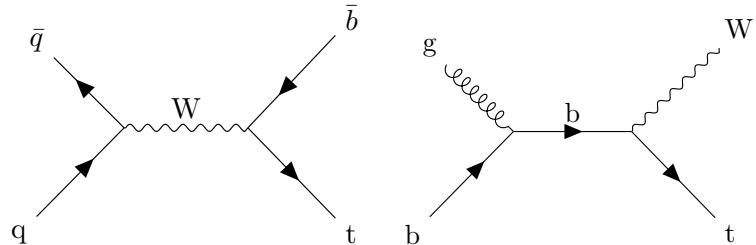


Figure 5.11: Feynman diagrams for the t-channel (on the left) and tW (on the right) production modes of a single top quark.

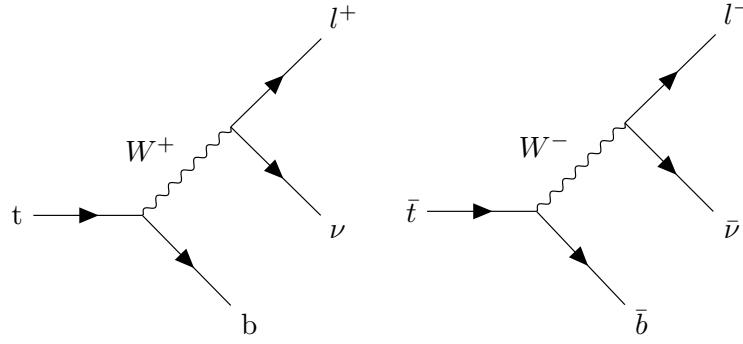


Figure 5.12: Feynman diagrams for the leptonic decay of the top (on the left) and anti-top (on the right) quarks.

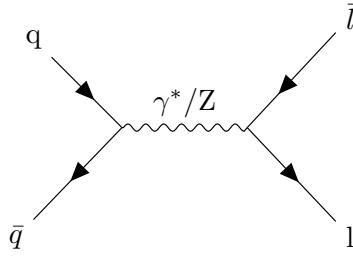


Figure 5.13: Feynman diagram for the DY process involving a virtual  $\gamma^*$  or Z boson.

is obtained using data and obtained from a general  $R_{in-out}$  method is applied to this particular process, which is additionally being checked in a specific control region as well.

### $R_{in-out}$ method

The idea behind this semi data-driven method is simple: since a 15 GeV veto in the dilepton invariant mass  $m_{ll}$  is applied to both the  $ee$  and  $\mu\mu$  channels to reduce the DY contamination in the signal region, then we could use such vetoed veto events in order to estimate the DY contribution outside of the Z mass window defined, referred to as  $M_Z$ . Mathematically, we can start by expressing in Equation 5.2 the total number of DY events  $N_{DY}^{total}$  as the sum of events observed inside  $N_{DY}^{in}$  and outside  $N_{DY}^{out}$  of  $M_Z$ .

$$N_{DY}^{total} = N_{DY}^{in} + N_{DY}^{out} \quad (5.2)$$

The parameter  $R_{out/in}$ , is then defined as the ratio between the number of MC or data events outside and inside this Z mass window. This ratio is then simply used to estimate the number of DY events outside the veto region in data from the number of observed events in the veto region, as shown in Equation 5.3.

$$N_{DY}^{out} = N_{DY,data}^{in} \cdot \left( \frac{N_{DY,MC}^{out}}{N_{DY,MC}^{in}} \right) \equiv N_{DY,data}^{in} \cdot R_{out/in, MC} \quad (5.3)$$

This previous equation relies on several assumptions. First of all, we have to assume that  $M_Z$  is dominated by actual DY events and peaking backgrounds leading to a similar  $R_{out/in}$  factor because of the prompt Z boson present in such backgrounds. Since this is typically a strong assumption to make, we actually decided to remove the contribution of non-peaking backgrounds

from the data yields to take this effect into account properly. Additionally, we also have to assume that the  $R_{out/in}$  factors measured in data and MC are similar, which is not verified when potential mismodeling of the data mass shape by the simulation appear. To account for this effect, we defined a region close to the signal region by reversing the b-jet requirement. If we assume that this mismodeling does not have a dependence in the number of b-jets, as shown in Figure 5.14, then we can simply correct the transfer factor previously obtained by a factor  $\kappa$  defined in Equation 5.4, experimentally determined to be of the order of  $\sim 5\%$ .

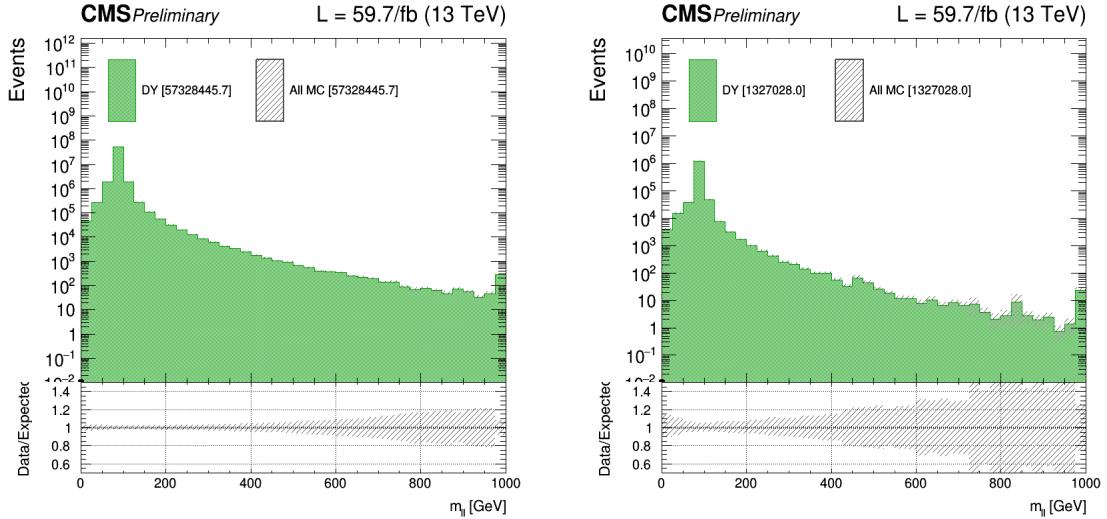


Figure 5.14: Normalized  $m_{ll}$  DY distributions obtained using 2018 MC simulations in the 0 (on the left) and 1+ b-jet (on the right) bins.

$$R_{out/in, MC}^{corr} = \kappa \cdot R_{out/in, MC} = \frac{R_{out/in, MC}^{0bj}}{R_{out/in, data}^{0bj}} \cdot R_{out/in, MC} \quad (5.4)$$

Additionally, we checked the behaviour of the transfer factor  $R_{out/in, MC}$  across a large range of MET values to make sure it is stable enough, as shown in Figure 5.15. Given the relative flatness of the distributions obtained, we decided to apply a single scale factor to the DY process, considering a systematic uncertainty simply taken as the observed maximum observed deviation between the central value and the different bins (around 20% in this case).

### 5.6.3 $t\bar{t} + W/t\bar{t} + Z$

These backgrounds are coming from an usual  $t\bar{t}$  production along with an Initial State Radiation (ISR) or Final State Radiation (FSR) production of a W or Z boson, as shown in Figure 5.16 and checked in specific control regions. The contribution of both these backgrounds is also taken directly from MC, generated using either MADGRAPH or MC@NLO, depending on the sample, with the PYTHIA8 CUETP8M1 (for 2016) and CP5 (for 2017 and 2018) tunes for the hadronization process.

The resulting cross section of such processes is a bit lower than the production of the SM  $t\bar{t}$  on its own but the kinematics of this background can be extremely close to our signal, since the W or Z boson produced can give a SM neutrino, leading to some actual MET.

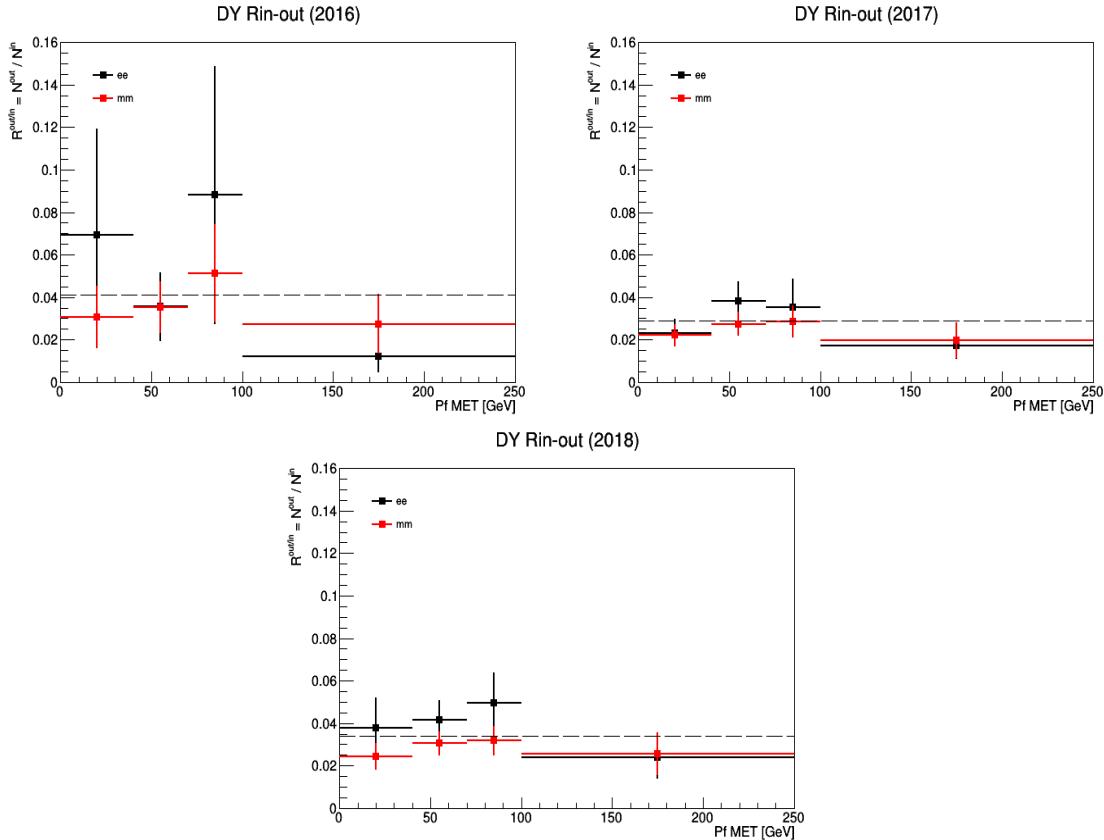


Figure 5.15:  $R_{out/in}$ , MC transfer factor obtained with the Rin-out data-driven method in bins of MET in 2016 (on the top left), 2017 (on the top right) and 2018 (on the bottom).

#### 5.6.4 Non prompt leptons contamination

Even though not extremely important in the sense that its kinematics allow us to remove most of its contributions in the signal regions, this background is interesting in the sense that it can be estimated using a data-driven method instead of being taken directly from MC.

A few definitions are first of all needed to explain the method used to compute the importance of this background in the different regions of the analysis:

- First of all, a **prompt lepton** is defined as a real lepton, in the sense that the lepton is originating from the PV of a  $pp$  collision.
- The **Prompt Rate (PR)** is defined as the number of prompt leptons passing the tight selection criteria of the analysis over the number of leptons passing the loose selection criteria.
- On the other hand, by **fake** or **non-prompt lepton**, we usually refer to truly **fake leptons**, such as jets misidentified by the detector as leptons, as shown in Figure 5.17, and real leptons coming from eventual heavy flavor decays.
- The **Fake Rate (FR)** is then defined similarly to the prompt rate but considering this time fake leptons only for the tight-to-loose ratio. This ratio therefore corresponds to the probability for a fake lepton to be considered as a real lepton in the analysis.

This background is particularly important at low  $p_T$ , where the misidentification rate is higher, and is not expected to be modeled correctly by MC because of its complexity: a general **tight-**

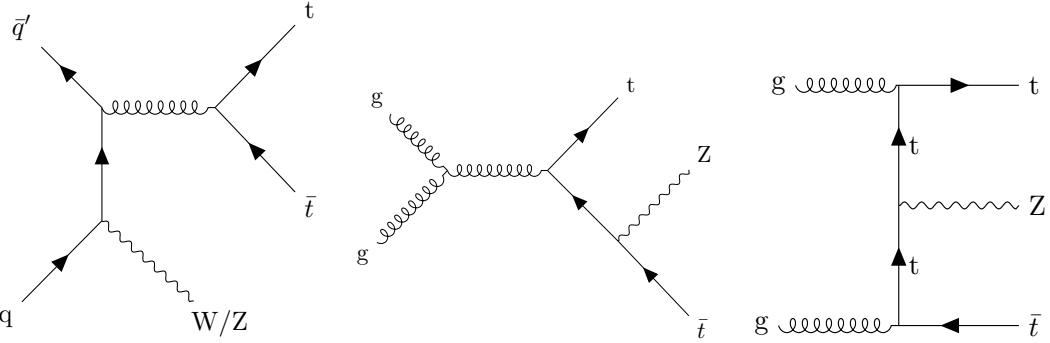


Figure 5.16: Possible Feynman diagrams for the ISR  $t\bar{t}$  with a W/Z boson (on the left) and for the production of an FSR  $ttZ$  (on the center and right).

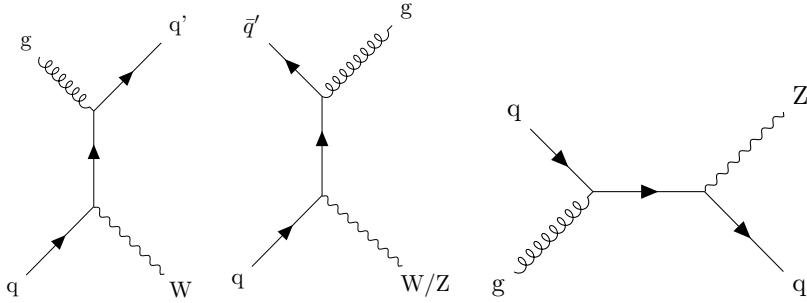


Figure 5.17: Possible Feynman diagrams for the production of a W/Z boson with a jet.

**to-lose data-driven method** is then used to compute its kinematics and final contribution in the different regions of the analysis.

In general, this method contains three main steps: the computation of the FR and PR, the extension of these rates in a region kinematically close to the SR of the analysis and the definition of a same sign control region enriched in fakes in order to perform a closure test of the yields and kinematics of this background.

### Fake Rate (FR) computation

Because of its definition, the FR is computed in a prompt lepton-free region, typically in a loose QCD enriched region, defined with the following cuts:

- Exactly 1 lepton
- $mtw1 < 20$  GeV
- $p_T > 13$  (10) GeV for  $e$  ( $\mu$ )
- $pfMET < 20$  GeV
- $|\eta| < 2.5$  (2.4) for  $e$  ( $\mu$ )
- PassJets

All the previous cuts have been designed to define a loose QCD region as pure as possible by removing most of the W+jets and Z+jets contribution. The PassJets cut is a boolean obtained by looping over all the jets of the event trying to find a jet having an  $E_T$  higher than a given threshold in order to control the average  $p_T$  of the jet that fakes the lepton (actually, different FR have been computed for different  $E_T$  thresholds, from 10 to 50 GeV).

Using this method, the jet that fakes a lepton is actually the one recoiling against ( $\Delta R > 1.0$ ) the jet used for the systematics, as shown in Figure 5.18.

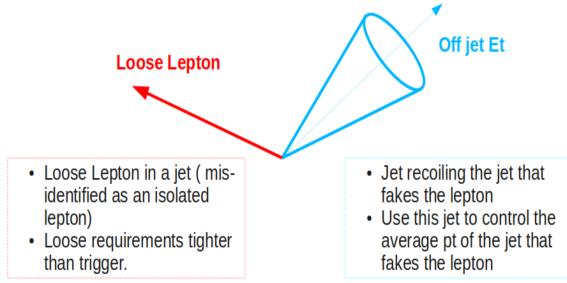


Figure 5.18: Schematic representation of the two jets used for the systematics and for the jet faking a lepton in the tight-to-loose data-driven method.

Events passing the following pre-scaled triggers are then selected in this region:

$$\text{Muon triggers} = \begin{cases} \text{HLT\_Mu8\_TrkIsoVVL (if } p_T < 20 \text{ GeV)} \\ \text{HLT\_Mu17\_TrkIsoVVL (if } p_T \geq 20 \text{ GeV)} \end{cases} \quad (5.5a)$$

$$\text{Electron triggers} = \begin{cases} \text{HLT\_Ele8\_CaloIdL\_TrackIdL\_IsoVL\_PFJet30 (if } p_T < 25 \text{ GeV)} \\ \text{HLT\_Ele23\_CaloIdL\_TrackIdL\_IsoVL\_PFJet30 (if } p_T \geq 25 \text{ GeV)} \end{cases} \quad (5.5b)$$

The remaining of the Electroweak (EWK) processes ( $W+jets$ ,  $Z+jets$ ) able to pass the previous cuts of the QCD region, are then simply subtracted: this is the so-called **EWK subtraction**.

Since both the FR and PR heavily depend on the kinematics of the event and on the Working Point (WP) chosen for the leptons of the analysis, they are computed separately depending on the flavor of the lepton and 2D histograms (accounting for the  $p_T$  and  $\eta$  of the lepton) need to be created at this stage to calculate this factor, for a given input jet  $E_T$  threshold; 1D histograms corresponding to the projections of these 2D histograms along both their axes are also defined at this point, as shown in Figure 5.19.

### Prompt Rate (PR) computation

The PR, taking into account the real lepton contamination in the region defined, is also important to calculate, even though the objects WP are usually chosen in such a way that this ratio is quite close to 1 and can therefore typically be ignored.

In our case, this rate has been calculated as well using a general tag and probe method in a  $Z+jets$  enriched sample. The main objective is to reconstruct  $Z \rightarrow ll$  events in this region and to select all the events for which the first lepton can be characterized as tight. Then, we search for the second lepton coming from the decay of  $Z$  within all the leptons detected by calculating the reconstructed mass of all the possible leptons combinations and selecting the one which is closer to the expected mass of the  $Z$  boson. We can then simply count how many times this second lepton, expected to be tight, has actually been measured as a tight lepton to estimate this PR.

The results obtained in this case have been represented in Figure 5.20.

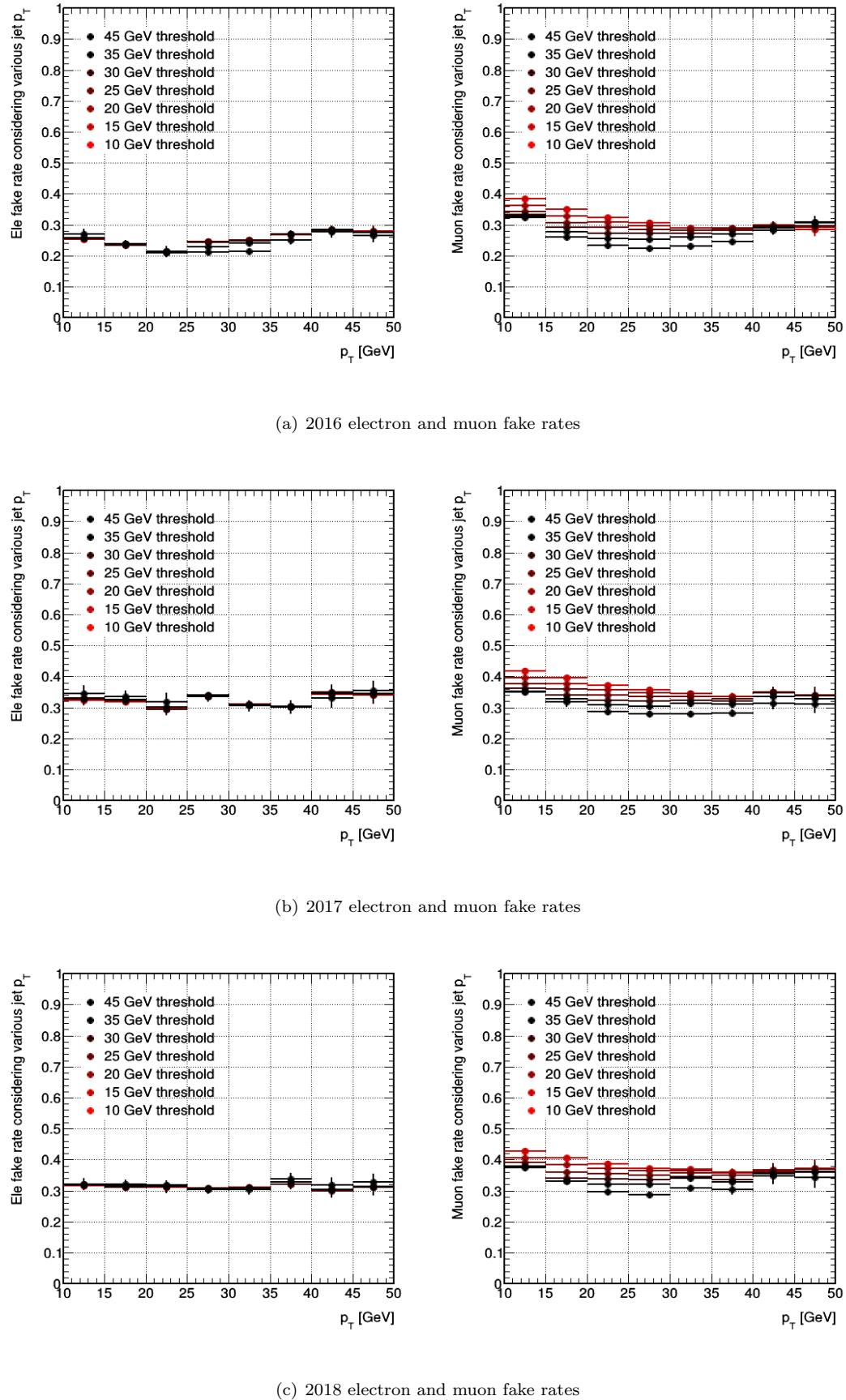


Figure 5.19: Electron (on the left) and muon (on the right) FR obtained in a QCD enriched region for different jet  $E_T$  thresholds for 2016, 2017 and 2018 with respect to the  $p_T$  of the lepton.

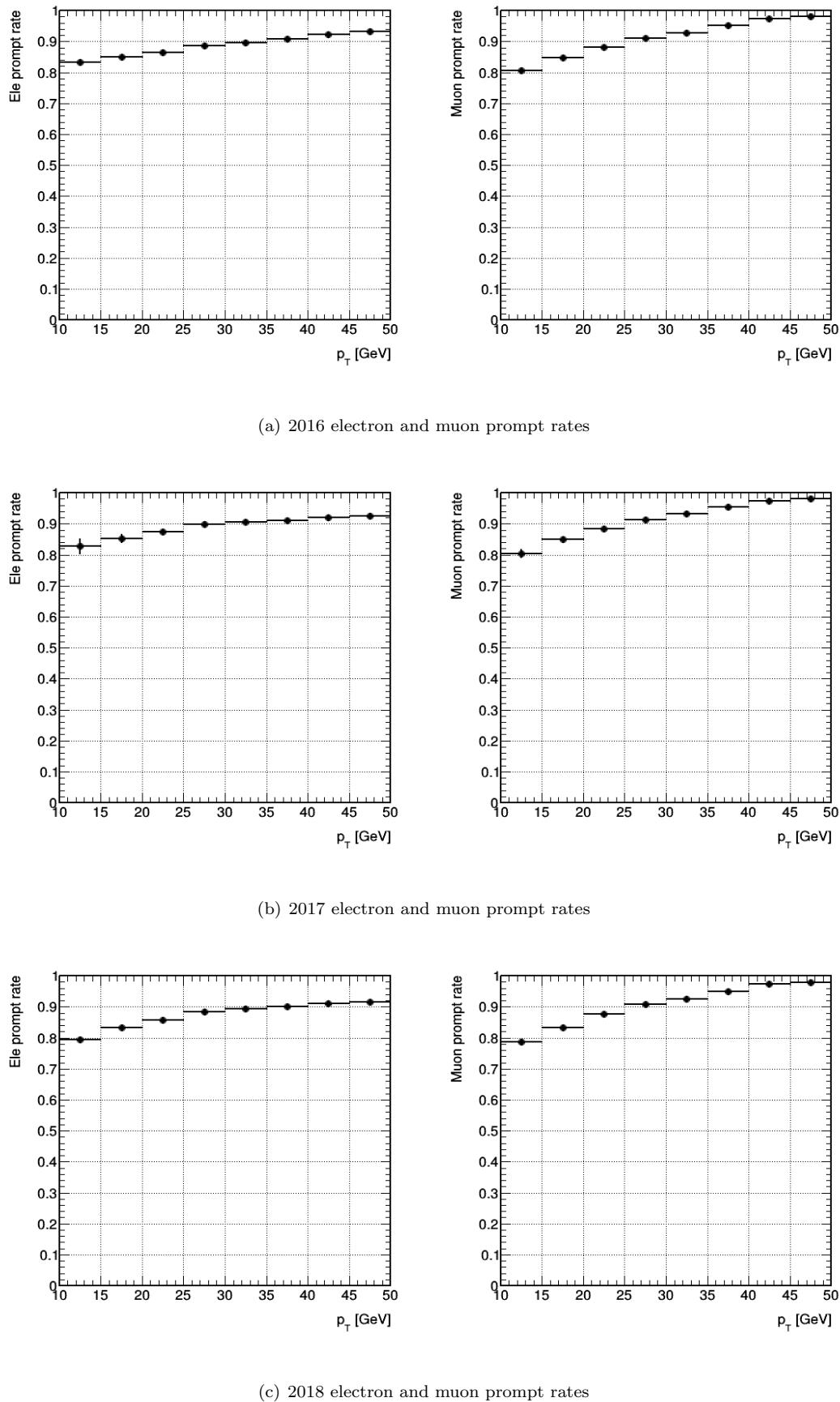


Figure 5.20: Electron (on the left) and muon (on the right) PR obtained in a  $Z+jets$  enriched region by a tag and probe method for 2016, 2017 and 2018 with respect to the  $p_T$  of the lepton.

### Fake weight calculation

Once the fake and prompt rates computed in their specific region, it is still necessary to apply them to a fake-lepton region kinematically close to the SRs of the analysis (usually, a l2loose region). For this, a simple set of equations can be used. We start by defining the following quantities:

- $N_{pp}$ , the number of events where both leptons are prompt
- $N_{fp}$ , the number of events where one lepton is prompt and the other is fake
- $N_{ff}$ , the number of events where both leptons are fake
- $N_{tx}$  ( $x = 0, 1, 2$ ), the number of events with 0, 1 or 2 leptons passing the right cuts, the **only quantity directly measurable** by the detector

It is then possible to see in Equation 5.6 that, if  $p$  is the PR and  $f$  the FR previously calculated, we can find an expression relating all these quantities.

$$\begin{cases} N_t = N_{pp} + N_{fp} + N_{ff} = N_{t2} + N_{t1} + N_{t0} \\ N_{t0} = (1-p)^2 N_{pp} + (1-p)(1-f)N_{fp} + (1-f)^2 N_{ff} \\ N_{t1} = 2p(1-p)N_{pp} + (f(1-p) + p(1-f))N_{fp} + 2f(1-f)^2 N_{ff} \\ N_{t2} = p^2 N_{pp} + pfN_{fp} + f^2 N_{ff} \end{cases} \quad (5.6)$$

These equations can be inverted in order to represent the unknowns with respect to the known variables, as shown in Equation 5.7, giving us a way to apply the weights previously calculated to this particular l2loose region.

$$\begin{pmatrix} N_{pp} \\ N_{fp} \\ N_{ff} \end{pmatrix} = \frac{f-p}{-(p-f)^3} \cdot \begin{pmatrix} f^2 & -f(1-f) & (1-f)^2 \\ -2fp & p(1-f) + f(1-p) & -2(1-p)(1-f) \\ p^2 & -p(1-p) & (1-p)^2 \end{pmatrix} \cdot \begin{pmatrix} N_{t0} \\ N_{t1} \\ N_{t2} \end{pmatrix} \quad (5.7)$$

A same sign control region enriched in fakes has also been defined in order to check this background.

The non-prompt contamination of this particular analysis however is expected to be dominated by the semi-leptonic decay of the SM  $t\bar{t}$  process, when one of the bottom quarks produced leads to a fake lepton. This particular process, contrary to the W+jets and Z+jets processes, is expected to be modeled quite well by the MC simulations, so we decided at the end of the day to rely on them for the analysis instead of using the data-driven fakes that had been produced.

#### 5.6.5 Smaller backgrounds

Even though quite negligible and reducible, some additional backgrounds still need to be considered, such as the SM  $t\bar{t}$  decaying semi-leptonically, the dibosons (WW, WZ and ZZ) and tribosons

(WWW, WWZ, WZZ, ZZZ, WWG) productions, as shown in Figure 5.21. All these smaller backgrounds are taken directly from MC and account in total for less than 1% of the total backgrounds in the different signal regions.

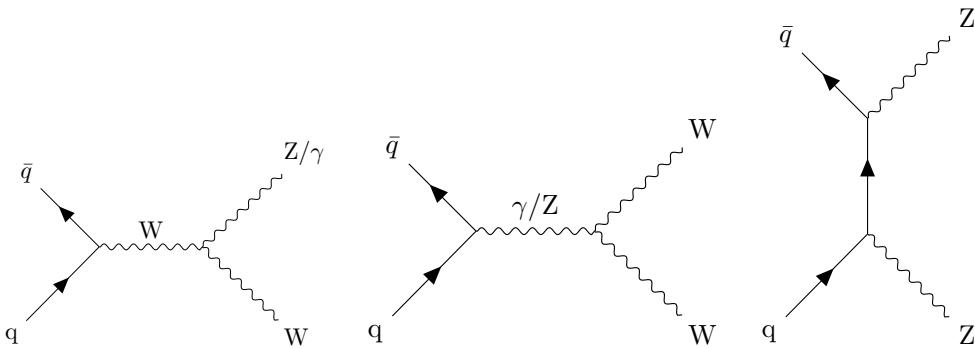


Figure 5.21: Possible Feynman diagrams for smaller backgrounds of this analysis: WW (on the left), W $\gamma$  and WZ (on the center) and ZZ (on the right).

### 5.6.6 Weights and corrections applied

Several different weights and SF usually need to be applied to the different MC processes in order to account for several effects observed in data but not accounted for during the generation of the MC simulation, such as the efficiency of the selection of the different objects, which will be discussed in Chapter 6, usually directly provided by the different Physics Object Groups (POGs) for their own default objects definitions.

The MC samples are additionally typically reweighted to match the distribution of true interactions observed in data due to multiple collisions collisions happening in the same bunch-crossings (the PU, as defined in Section 3.1.2). To illustrate the effect of these corrections, a few particular weights that are applied to the MC samples will now be detailed.

#### MET corrections

The MET distribution is expected to be independent of  $\phi$  because of the rotational symmetry of the collisions around the beam axis but we do observe that the reconstructed MET does depend on  $\phi$  because of possible anisotropic detector responses, inactive calorimeter cells or tracking regions, detector misalignment, or eventual displacements of the beam spot. XY-shift corrections have therefore been applied to reduce this modulation. In 2017, large level of noise were also observed in the data collected by the ECAL at high pseudorapidities. To mitigate this effect, especially on the MET tails, the correction consisting in completely excluding jets with raw transverse momentum  $< 50$  GeV and  $2.650 < |\eta| < 3.139$  from the calculation of the MET was considered, as recommended by the JET/MET POG [144]. The effect of these corrections is shown in Figure 5.22.

#### Top $p_T$ reweighting

Previous studies of the  $t\bar{t}$  generator typically considered in the physics analysis predict a harder top quark  $p_T$  spectrum than the one observed in data [146]. This known mismodeling of the top quark is corrected using a general reweighting recipe [147], and its effect on the  $p_T$  spectrum can be observed in Figure 5.23.

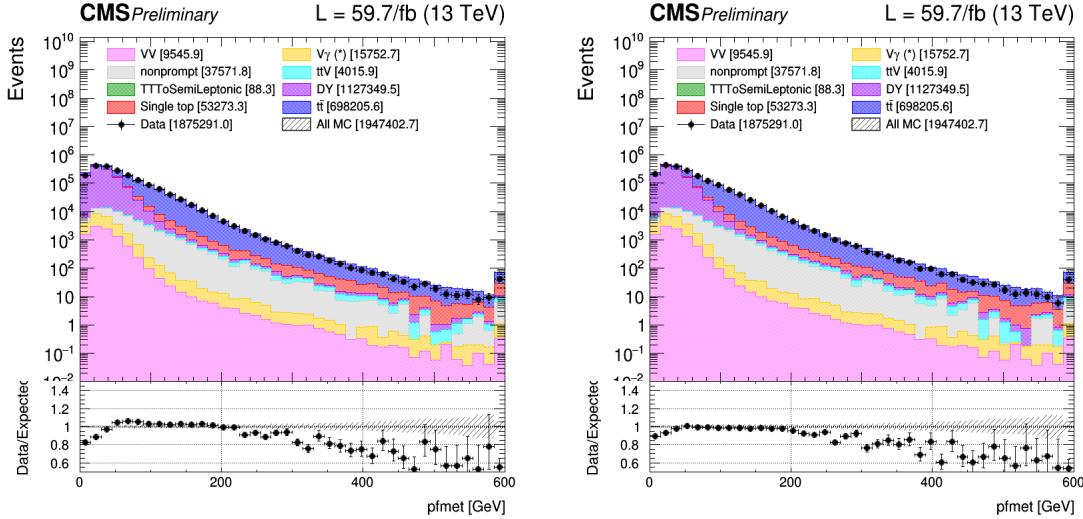


Figure 5.22: Uncorrected PFMET (on the left) and PFMET after applying the EE noise and XY-shift corrections (on the right) distributions observed in a 2018 DY inclusive control region.

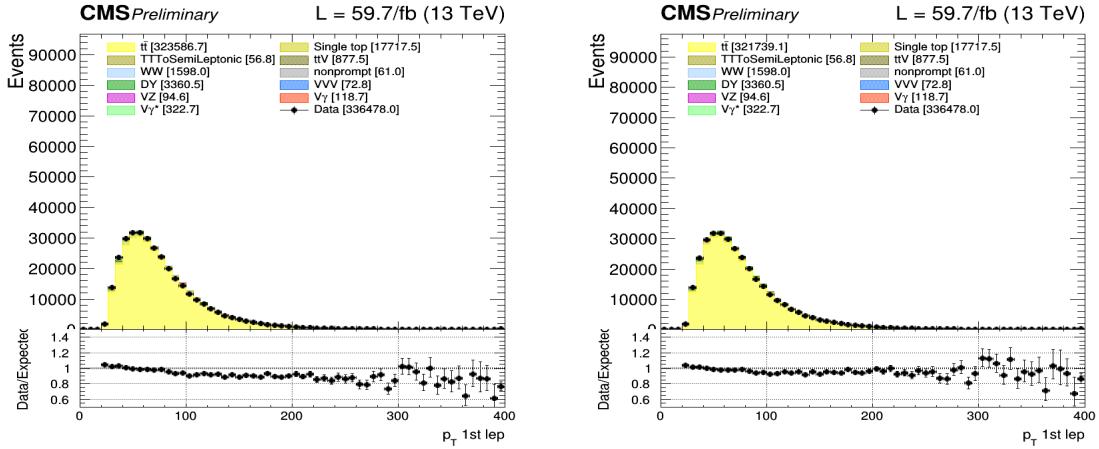


Figure 5.23: Data/MC agreement without (on the left) and with (on the right) top  $p_T$  reweighting corrections in a general 2017 top control region.

Even though we initially planned on applying this correction factor, the most recent recommendation is not to apply this factor for our kind of BSM searches. In this case, the effect of the top  $p_T$  mismodeling can be considered covered by the existing uncertainties and no additional correction or uncertainty is needed [148].

### Other factors

In 2016 and 2017 an issue, known as the **prefiring**, causing highly energetic readout from jets, photons and electrons in the ECAL endcap to be assigned by the L1 trigger to the previous bunch crossing was discovered. To make up for this difference, a weight  $(1 - x)$  is usually applied to all MC events, where  $x$  is the probability of an event to be prefired [143]. The effect this correction has on the data/MC agreement in a 2017 top enriched control region is shown in Figure 5.24.

In 2018, another issue affecting this time two endcaps of the HCAL and known as the HEM15/16 issue was reported [145], resulting in a loss of around 2% of HCAL coverage. This issue has a small

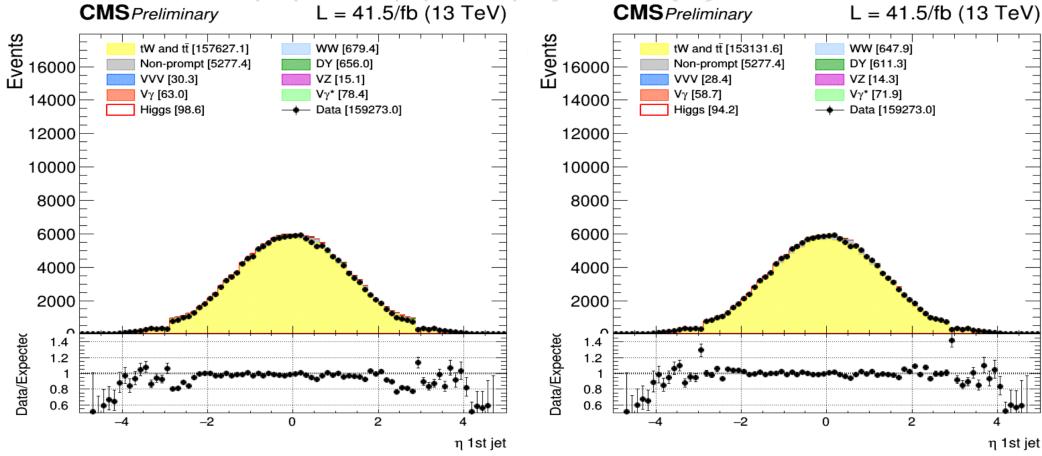


Figure 5.24: Data/MC agreement without (on the left) and with (on the right) prefiring corrections in a general 2018 top control region.

but measurable effect on the MET spectrum and is therefore taken into account and corrected by introducing another weight to the simulation.

---

---

# Chapter 6

---

## Event selection

This Chapter will be dedicated to the analysis itself, by defining first of all the different objects actually used in this case, along with the actual selection that has been applied to enhance the quality of such objects in this particular search in Section 6.1. Then, the different Signal Regions (SRs) defined in which a high purity of signal is expected are defined in Section 6.2 while all the different Control Regions (CRs) defined in order to check the behavior of the MC simulation performed for the major backgrounds on this analysis, such as the single top or SM  $t\bar{t}$  production, will be introduced in Section 6.3.

### 6.1 Objects selection

We already described what to expect from a typical  $t/\bar{t}$  or  $t\bar{t}+\text{DM}$  signal: the typical signature of such signals is made out of a certain number of b tagged jets along with two leptons (electrons and/or muons) and some MET coming from the two DM particles created along the way. It is extremely important to describe the WP chosen and the selection applied in order to select the objects of the analysis, such as the leptons and the jets used, chosen in such a way to optimize the lepton reconstruction efficiency while reducing as much as possible the possible misidentification rates of the different objects.

First of all, the different triggers used to collect the data will be detailed in Section 6.1.1. Then, the leptons used in this analysis will be introduced in Sections 6.1.2 (for electrons) and 6.1.3 (for muons). Finally, given the nature of the DM signal searched for, a complete description of the jets selected in the analysis will be necessary and performed in Section 6.1.5.

#### 6.1.1 Triggers selection

The triggers, described in Section 3.2.6, and particularly the trigger paths chosen are an important part of each analysis since they describe the kind of data that can be collected and therefore analyzed. The triggers used in this analysis for the datasets available for the years 2016, 2017 and 2018 can be found in Tables 6.1, 6.2 and 6.3 respectively.

Dataset	Run range	<b>HLT trigger path</b>
SingleMu	[297020,306462]	HLT_IsoMu27_v*
SingleEle	[297020,306462]	HLT_Ele35_WPTight_Gsf_v*
DoubleEG	[297020,306462]	HLT_Ele23_Ele12_CaloIdL_TrackIdL_IsoVL_v*
DoubleMu	[297020,299336]	HLT_Mu17_TrkIsoVVL_Mu8_TrkIsoVVL_DZ_v*
	[299337,306462]	HLT_Mu17_TrkIsoVVL_Mu8_TrkIsoVVL_DZ_Mass8_v*
MuonEG	[297020,306462]	HLT_Mu12_TrkIsoVVL_Ele23_CaloIdL_TrackIdL_IsoVL_DZ_v*
	[297020,299336]	HLT_Mu23_TrkIsoVVL_Ele12_CaloIdL_TrackIdL_IsoVL_DZ_v*
	[299337,306462]	HLT_Mu23_TrkIsoVVL_Ele12_CaloIdL_TrackIdL_IsoVL_v*

Table 6.1: 2016 trigger paths considered for this analysis.

Dataset	Run range	<b>HLT trigger path</b>
SingleMu	[297020,306462]	HLT_IsoMu27_v*
SingleEle	[297020,306462]	HLT_Ele35_WPTight_Gsf_v*
DoubleEG	[297020,306462]	HLT_Ele23_Ele12_CaloIdL_TrackIdL_IsoVL_v*
DoubleMu	[297020,299336]	HLT_Mu17_TrkIsoVVL_Mu8_TrkIsoVVL_DZ_v*
	[299337,306462]	HLT_Mu17_TrkIsoVVL_Mu8_TrkIsoVVL_DZ_Mass8_v*
MuonEG	[297020,306462]	HLT_Mu12_TrkIsoVVL_Ele23_CaloIdL_TrackIdL_IsoVL_DZ_v*
	[297020,299336]	HLT_Mu23_TrkIsoVVL_Ele12_CaloIdL_TrackIdL_IsoVL_DZ_v*
	[299337,306462]	HLT_Mu23_TrkIsoVVL_Ele12_CaloIdL_TrackIdL_IsoVL_v*

Table 6.2: 2017 trigger paths considered for this analysis.

Our analysis relies on the dilepton final state, so the single lepton trigger are only considered in order to recover some of the efficiency lost in some cases when one lepton passes the tight identification criteria while the second one does not, and does therefore not trigger the event. The logical *or* of all the trigger paths is considered. Eventual events passing several triggers is taken into account as well to make sure to avoid any double counting.

In particular, the trigger efficiency, defined for each trigger as the ratio between the number of events passing our object selection and the trigger itself in the numerator and the number of events passing our selection in the denominator has been computed using orthogonal MET datasets to avoid any bias. Studying this efficiency is important in the sense that we want to make sure that the triggers used are efficient enough in the  $p_T$  region of the leptons of the analysis to avoid any undesired effect due to the turn-on of any trigger, and because we actually use them to reweight the simulated samples. The efficiencies have in this sense been calculated using orthogonal MET triggers for the different data taking periods and for each dataset individually, as shown in Figure 6.2.

### 6.1.2 Electrons selection

Four different electron Working Points (WPs) (veto, loose, medium, tight) are then defined by the CMS EGamma POG [149] with slightly different quality cuts in the barrel or in the endcaps

Dataset	Run range	HLT trigger path
SingleMu	[315252,325175]	HLT_IsoMu24_v*
SingleEle	[315252,325175]	HLT_Ele32_WPTight_Gsf_v*
DoubleEG	[315252,325175]	HLT_Ele23_Ele12_CaloIdL_TrackIdL_IsoVL_v*
DoubleMu	[315252,325175]	HLT_Mu17_TrkIsoVVL_Mu8_TrkIsoVVL_DZ_Mass3p8_v*
MuonEG	[315252,325175]	HLT_Mu23_TrkIsoVVL_Ele12_CaloIdL_TrackIdL_IsoVL_v*
		HLT_Mu12_TrkIsoVVL_Ele23_CaloIdL_TrackIdL_IsoVL_DZ_v*

Table 6.3: 2018 trigger paths considered for this analysis.

in order to select electrons with a given efficiency while trying to limit the misidentification rate. The tight WP is then the one with the lowest electron selection efficiency (of the order of 70% for electrons with  $p_T > 20$  GeV) but the best to reject misidentified electrons from other source, to be used when the backgrounds are expected to be large. On the other hand, the veto WP corresponds to an average electron selection efficiency of the order of 95%.

For this analysis, we rely on the cut based medium identification POG WP to define our electrons [150]. Electrons originating from photon conversions are also rejected by requiring that the electron track not have missing hits in the innermost layers of the tracker, on top of the conversion veto included in the cut based ID definition.

### 6.1.3 Muons selection

The selection applied to muons is based on the Muon POG as well, providing references efficiencies for standard selection and recommendations for the tight muon selection [151].

For our analysis, we decided to use directly the medium ID WP provided by the POG [152], along with a tighter isolation criterium ( $< 0.15$ ) with  $\Delta\beta$  correction and in a cone size  $\Delta R < 0.4$ , as defined in Equation 6.1, in order to reduce the number of muons coming from the hadronization process of bottom and charm quarks.

$$\text{ISO} = \frac{\sum p_T^{\text{ch. had. (PV)}} + \max(0, \sum E_T^{\text{neut. had.}} + \sum E_T^\gamma - 0.5 \times \sum p_T^{\text{ch. had. (PU)}})}{p_T(\mu)} \quad (6.1)$$

### 6.1.4 Leptons selection

A few additional quality cuts are applied to both the electrons and muons. We ask them to have a  $p_T > 8$  GeV and a pseudo-rapidity  $|\eta| < 2.4$  to be considered valid, and candidate lepton trajectories are further required to be compatible with the primary interaction vertex by imposing constraints on their transverse ( $|d_0| < 0.05$  cm) and longitudinal ( $d_z < 0.10$  cm) impact parameters, and on the three-dimensional impact parameter significance ( $S_{3D}^d < 4$ ), computed as the ratio of the three-dimensional impact parameter and its uncertainty.

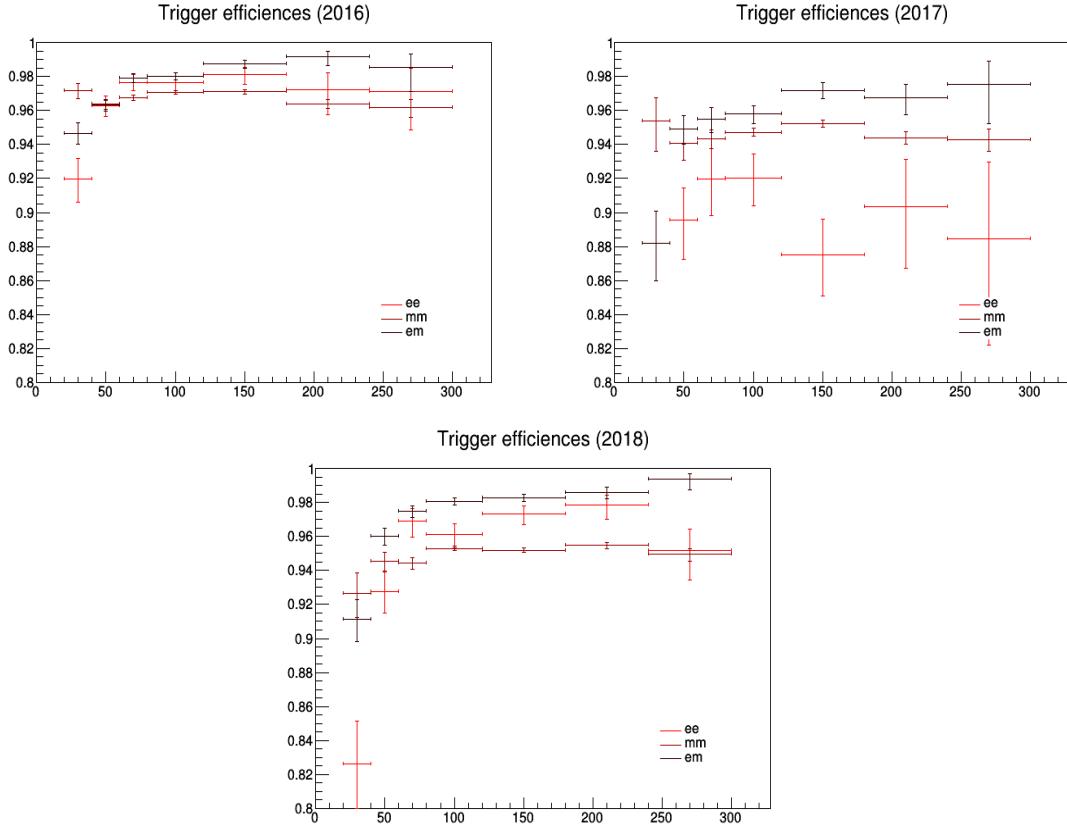


Figure 6.1: Trigger efficiencies using orthogonal MET datasets for each dataset in 2016 (on the top left), 2017 (on the top right) and 2018 (on the bottom).

### Veto WP

As we will see later on, events featuring more than two leptons are rejected, and this is achieved by defining a looser lepton selection. For electrons, the cut based veto ID definition is used, and the requirement of not having missing hits is removed while for muons, candidates satisfying the loose ID and very loose isolation criteria are selected to define the veto WP.

#### 6.1.5 Jet selection

Jets are an important part of this analysis as well given the final state searched for in this case. As explained in Section 4.4, the jets are clustered from the PF candidates using the anti- $k_T$  algorithm (with a typical distance parameter  $R = 0.4$ ).

In this analysis, jets are selected following the tight WP definition given by the CMS JET/MET POG [153], whose selection depends on the pseudorapidity of the jet and on the year of data taking [?]. The tight POG WP has been chosen since it offers an efficiency higher than 98-99% for all the jets and a background rejection higher than 98% for jets having  $|\eta| < 3.0$ . Jets are further required to have a  $p_T > 20$  GeV (30 GeV for the first jet) and a  $|\eta| < 2.4$  to be considered in this analysis, and jets overlapping with any selected lepton within a cone of radius  $R < 0.4$  are removed to prevent signal leptons clustered as jets from entering the jet counting.

Finally, a slightly different selection (tight PU jet ID) is applied to jets having a  $p_T < 50$  GeV in

order to reject jets coming from PU interactions. Following the official POG recommendations, jet energy correction and resolution are also taken into account [155, 156].

### b-jets

The b-jets are selected among the jets using the recommendations given by the B-Tagging and Vertexing POG [157]. The medium deepCSV b-tagging WP is used in this case, as explained in Section 4.4.1, for which the rate for misidentifying a light jet as a b-jet is around 10%. A jet is therefore in particular considered to be a b-jet in this analysis if it passes the jet requirements and if its deep CSV b-tagging weight is larger than 0.6321, 0.4941 or 0.4184 for the year 2016, 2017 and 2018, respectively (medium POG working point).

## 6.2 Signal regions

It is important to note at this point that a strict **blinding policy** has been followed for this search, in order to avoid optimizing the analysis based on what has already been seen. At first, the data available to be plotted in the following signal regions has therefore been limited to  $1\text{fb}^{-1}$  for each year, before the unblinding, allowing us to look at the whole Run II dataset. Additionally, only the statistical errors are represented in all the plots, unless stated otherwise.

First of all, a global pre-selection was agreed on between the groups involved, applied to all the signal regions, reducing by a large factor the contamination due to several backgrounds, such as the SM  $t\bar{t}$  and the DY processes:

- Exactly two opposite sign good leptons ( $p_T > 25$  (20) GeV for the leading (trailing) lepton, both having  $|\eta| < 2.4$ ) are required;
- Then, events with a third lepton having a  $p_T > 10$  GeV are rejected;
- Low mass resonances are removed by asking for  $m_{ll} > 20$  GeV;
- The DY is then strongly reduced by asking for a 15 GeV Z-veto in the  $ee$  and  $\mu\mu$  channels, and by asking for at least 1 jet with a  $p_T > 30$  GeV and  $|\eta| < 2.4$  in the event;
- 1 medium deep CSV b-jet is also required;
- And, finally, a cut using the stranverse mass (defined in Section 7.1)  $M_{T2}(ll) > 80$  GeV has been applied, mainly to keep our signal regions orthogonal to the  $t\bar{t}$  control regions used by the semi-leptonic channel, making it easier to combine our results while removing more background than signal anyway.

Some distributions obtained in this pre-selection region are shown in Figure 6.3.

Several different signal regions have then actually been defined, targeting each of the two signals of interest, depending on the number of jets and b-jets observed. The signal region targeting the  $t/\bar{t}+\text{DM}$  process is therefore made out of events having either exactly 1 jet, or 2 jets and exactly 1 b-jet. On the other hand, all the events passing the pre-selection but having either 2 jets and more than 1 b-jet or more than 2 jets make up the signal region targeting the  $t\bar{t}+\text{DM}$  process. Different MVA trainings were done in these two regions, as explained in Section 7. The signal regions can be even more divided, depending on the mass point that was used for the training, since the mediators with different masses typically feature different kinematic distributions.

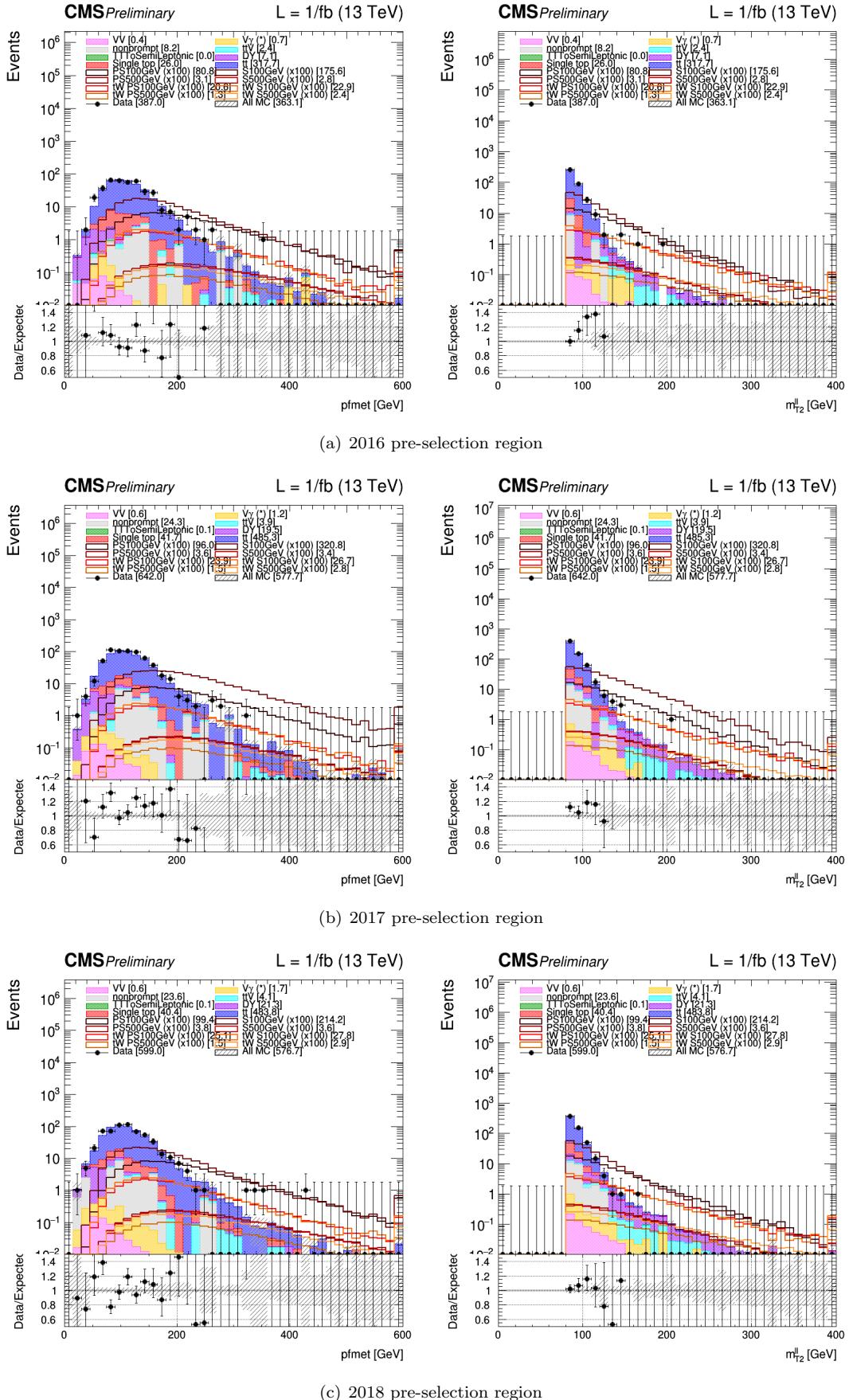


Figure 6.2: Two different variables (pfMET, on the left, and the transverse mass  $M_{T2}^{ll}$ , on the right) represented in the pre-selection region used for the training of the MVA.

Some distributions obtained in these signal regions are shown in Figures 6.4 to 6.9.

**FIXME:** Add final plots everywhere

## 6.3 Control regions

Different control regions where the data/MC agreement can be checked have been defined in order to check the validity of the MC simulations performed for the SM processes corresponding to the main backgrounds of this analysis. In this context, an inclusive control region is defined in Section 6.3.1 to spot initial issues and problems. Then, a data validation region is defined to ensure the quality of the MC prediction of both the  $t\bar{t}$  and single top processes in Section 6.3.2 and a DY enriched specific control region is defined in Section 6.3.3. Finally, the  $ttV$  and non-prompt contamination processes will be checked in Sections 6.3.4 and 6.3.5 respectively.

### 6.3.1 Inclusive control region

First, several distributions in a control region as inclusive as technically possible and mostly enriched in DY were studied, in all the different channels available. This control region is defined as a two opposite sign tight leptons region, with a veto on any events featuring three leptons, with some small additional cuts ( $m_{ll} > 20$  GeV to avoid low mass resonances, while at least 1 jet and 1 medium deep CSV b-jet are also required). The distributions obtained in all the channels in this particular control region are shown in Figure 6.10.

In all the DY enriched regions, a slight MC mismodeling can typically be observed, and a general problem in the estimation of the pfMET in DY events is observed. These features are known in CMS and mitigated anyway with the analysis cuts applied.

### 6.3.2 Top control region

As previously explained, the SM  $t\bar{t}$  control region is defined with the same pre-selection applied to both the signal regions, but selecting only events having a stranverse mass  $60 < M_{T2}^l < 80$  GeV, making this region perfectly orthogonal to the signal regions. Distributions obtained in this region are shown in Figure 6.11.

### 6.3.3 DY control region

Given the huge cross section of the DY, this process is expected to be present at almost any selection level, making it important to study in a dedicated control region. The selection applied to the signal regions does reduce it a lot though, especially by specifically asking for a Z-veto in the  $ee$  and  $\mu\mu$  channels.

At the end of the day, this control region was defined exactly the same as the signal regions, but by simply reversing the Z-veto requirement, making it perfectly orthogonal and allowing us to use the  $R_{\text{in-out}}$  method described in Section 5.6.2. Some distributions obtained in these control regions are shown in Figure 6.12. It is important to keep in mind that this last region is mostly used to compute the  $R_{\text{in-out}}$  scale factor, and that the 0-bjet correction factor  $\kappa$  described in Section 5.6.2

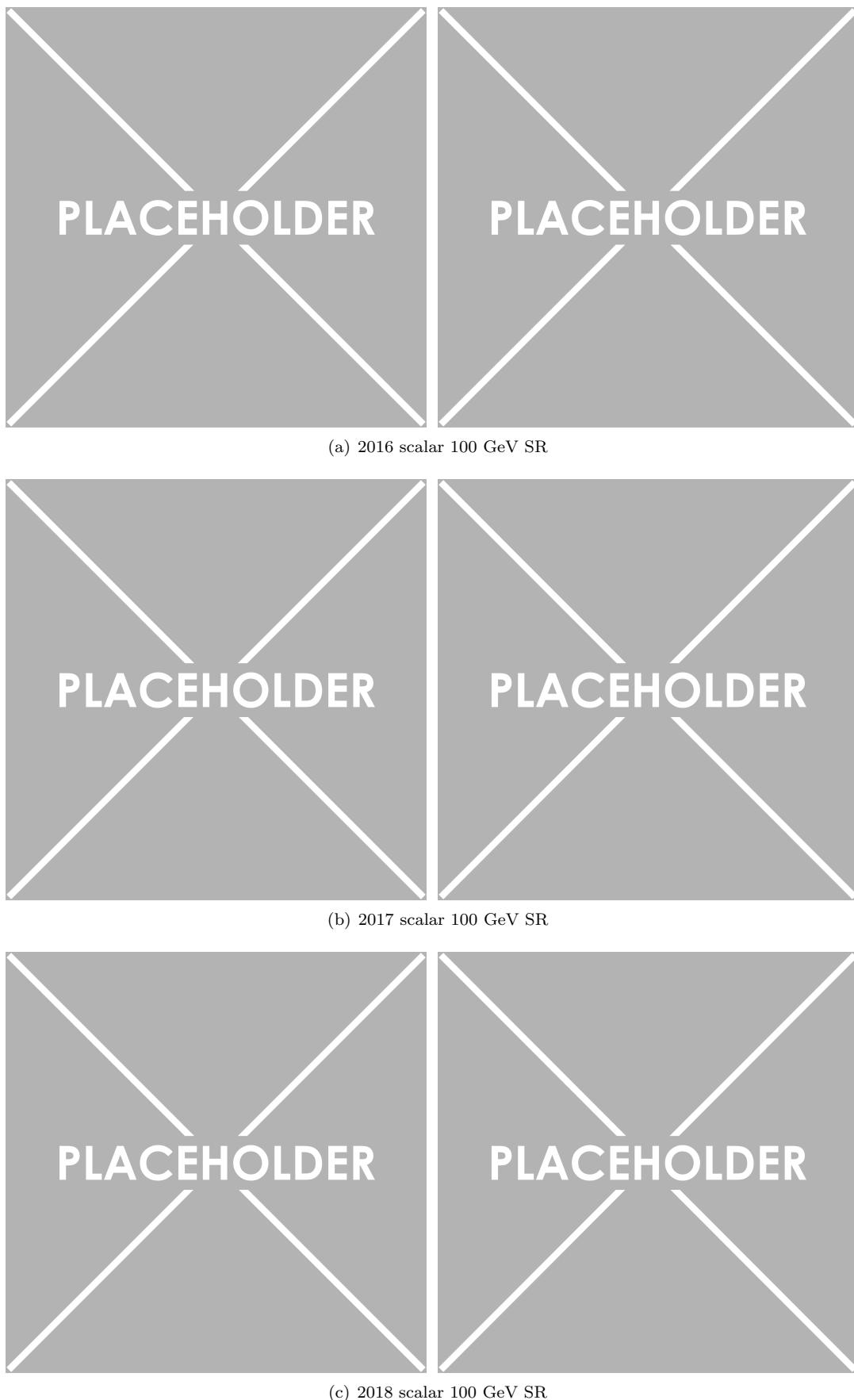


Figure 6.3: Two different variables (pfMET, on the left, and the stransverse mass  $M_{T2}^{ll}$ , on the right) represented in the  $t/\bar{t}+{\rm DM}$  signal region defined from the 100 GeV scalar training.

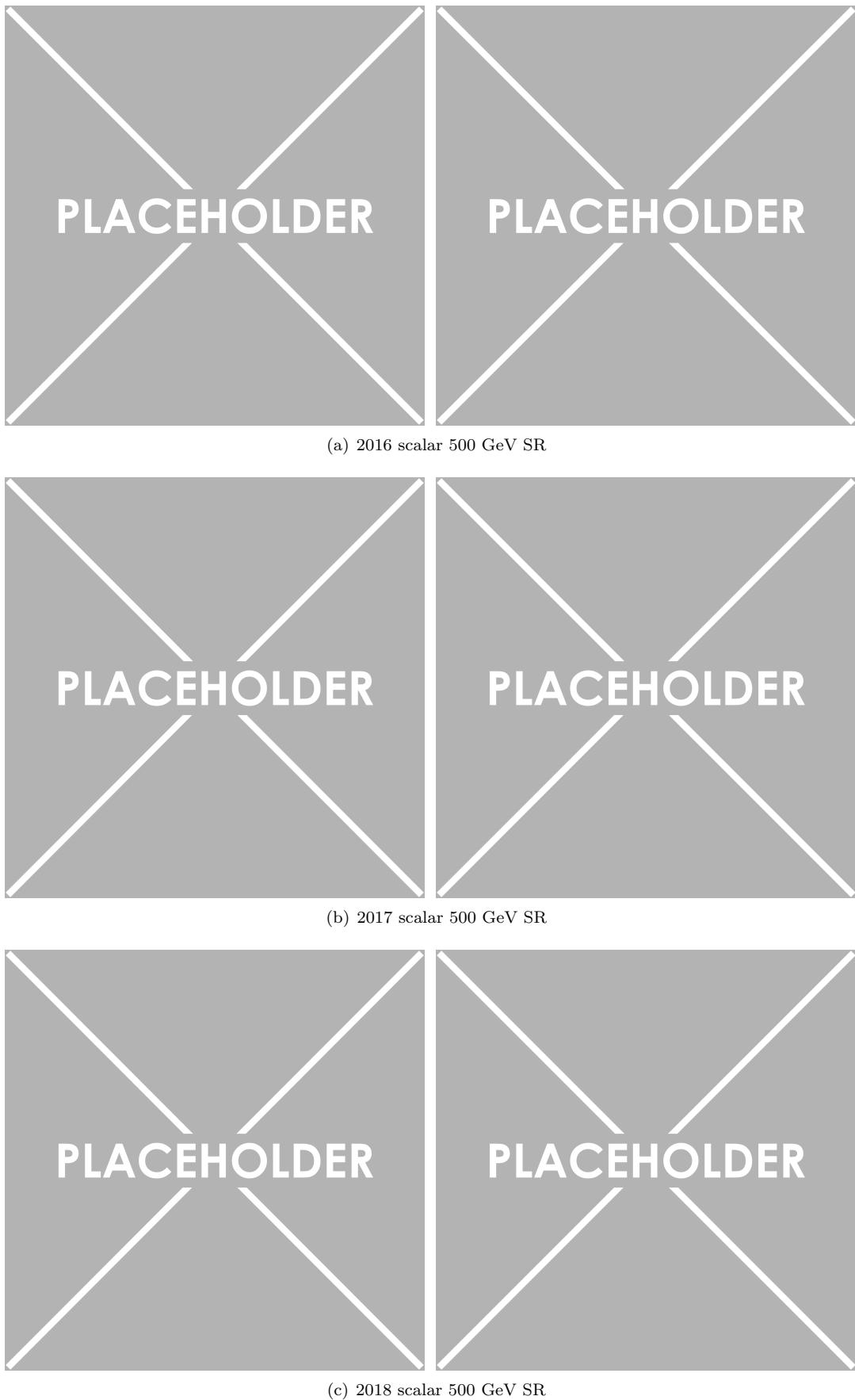


Figure 6.4: Two different variables (pfMET, on the left, and the stransverse mass  $M_{T2}^{ll}$ , on the right) represented in the  $t/\bar{t}+DM$  signal region defined from the 500 GeV scalar training.

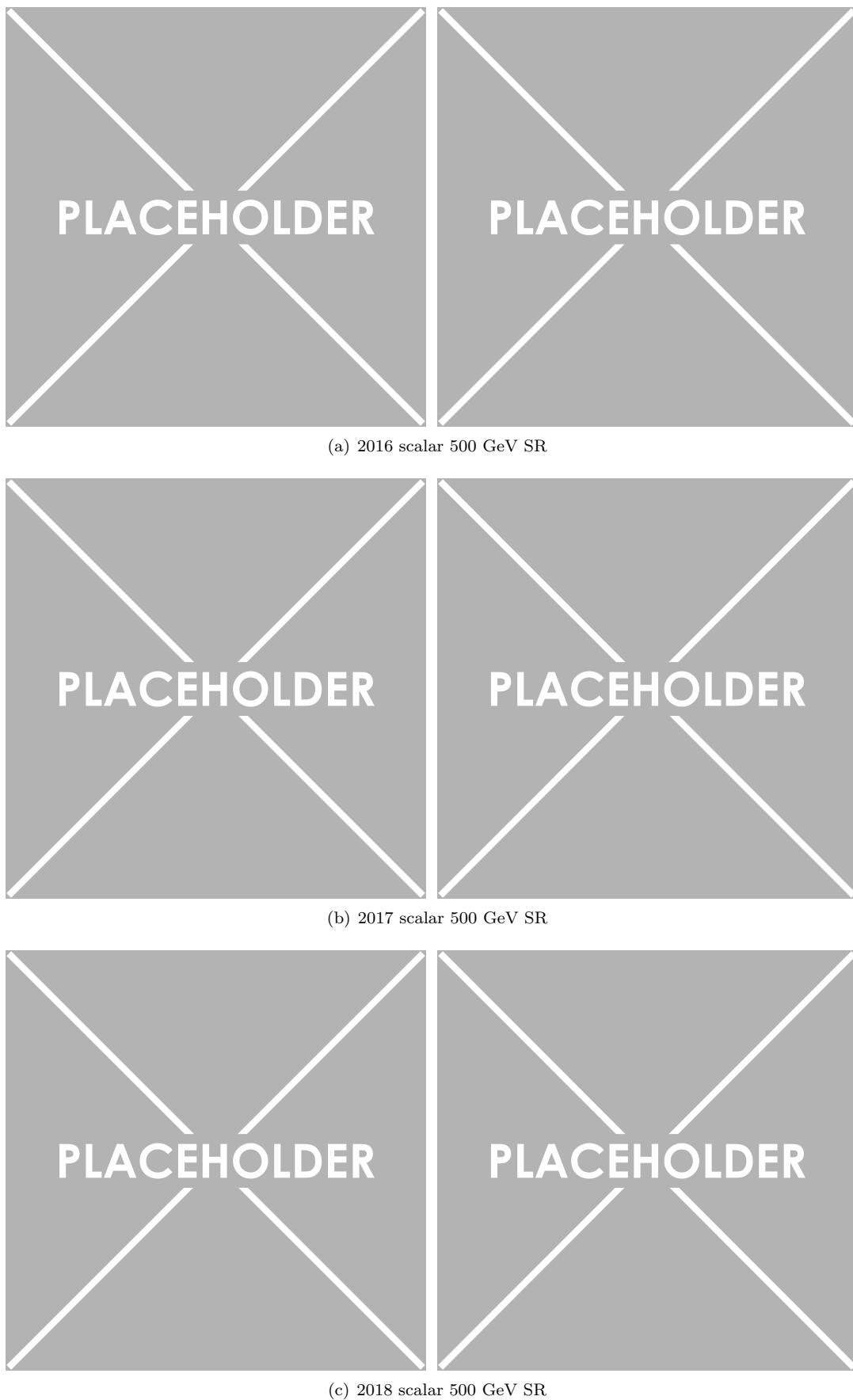


Figure 6.5: Two different variables (pfMET, on the left, and the stransverse mass  $M_{T2}^{ll}$ , on the right) represented in the  $t/\bar{t}+{\rm DM}$  signal region defined from the 100 GeV pseudoscalar training.

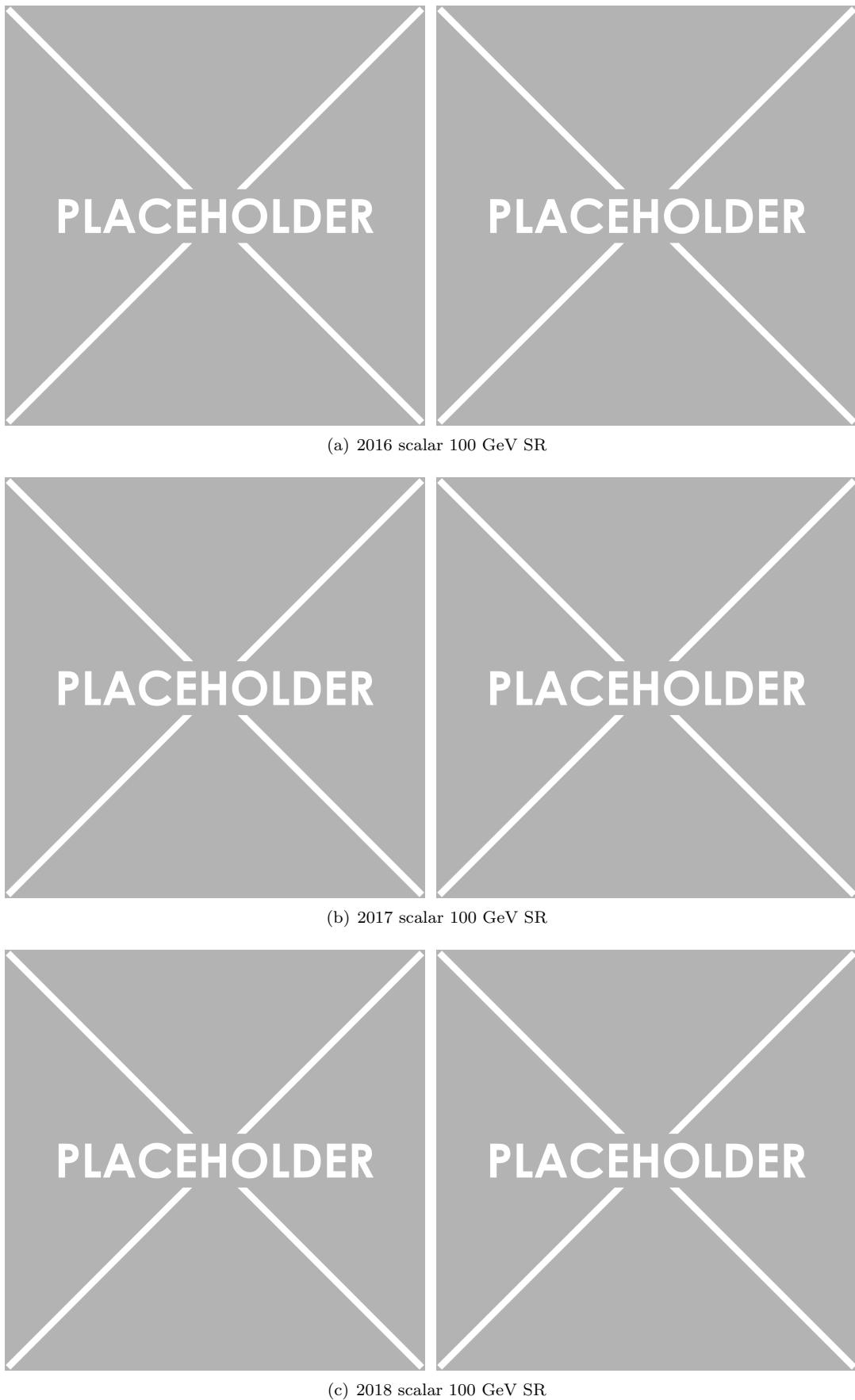


Figure 6.6: Two different variables (pfMET, on the left, and the stransverse mass  $M_{T2}^{ll}$ , on the right) represented in the  $t\bar{t}$ +DM signal region defined from the 100 GeV scalar training.

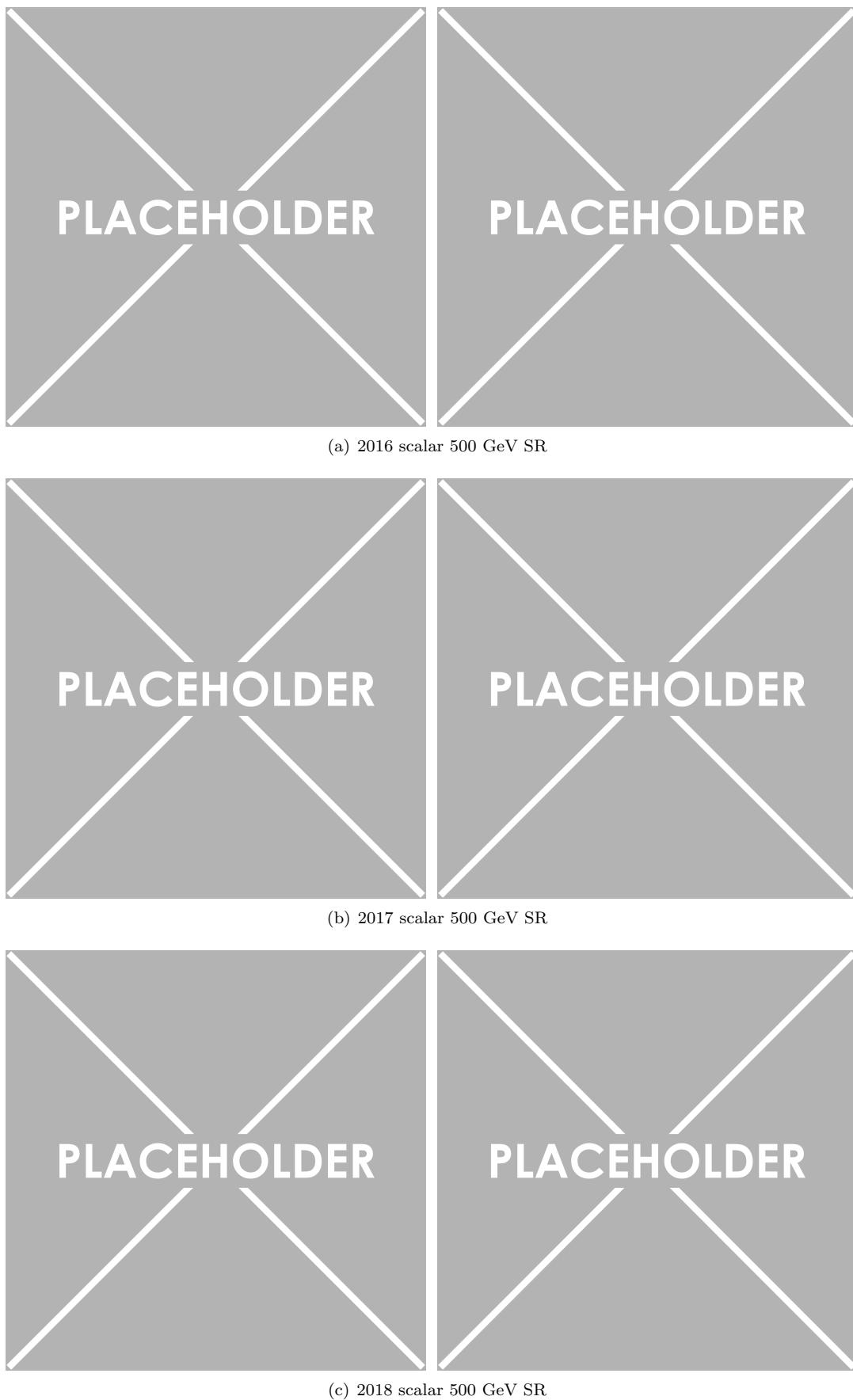


Figure 6.7: Two different variables (pfMET, on the left, and the stransverse mass  $M_{T2}^{ll}$ , on the right) represented in the  $t\bar{t}$ +DM signal region defined from the 500 GeV scalar training.

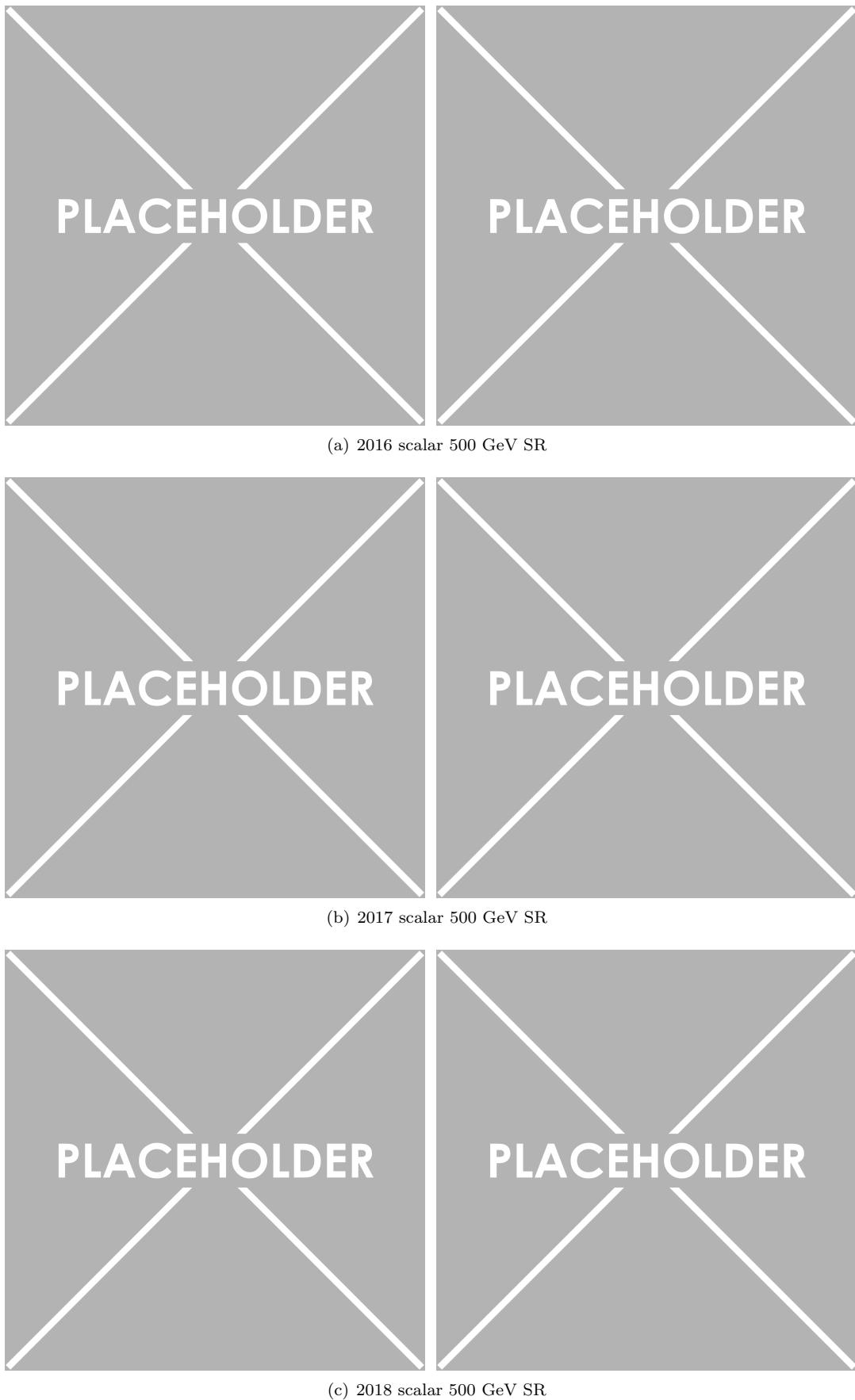


Figure 6.8: Two different variables (pfMET, on the left, and the stransverse mass  $M_{T2}^{ll}$ , on the right) represented in the  $t\bar{t}$ +DM signal region defined from the 100 GeV pseudoscalar training.

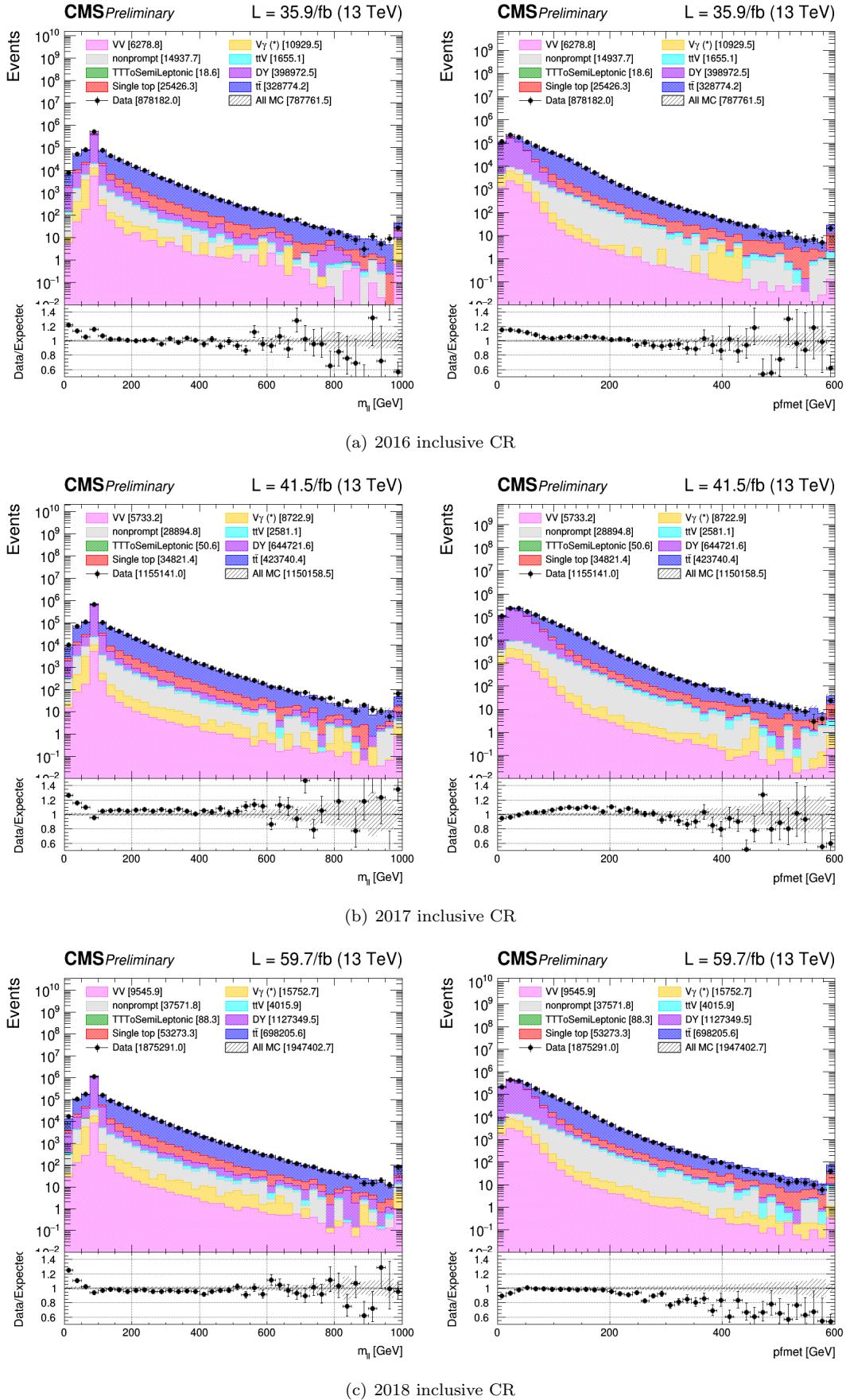


Figure 6.9: Two different variables ( $m_{ll}$ , on the left, and pfMET, on the right) represented in the inclusive control region defined.

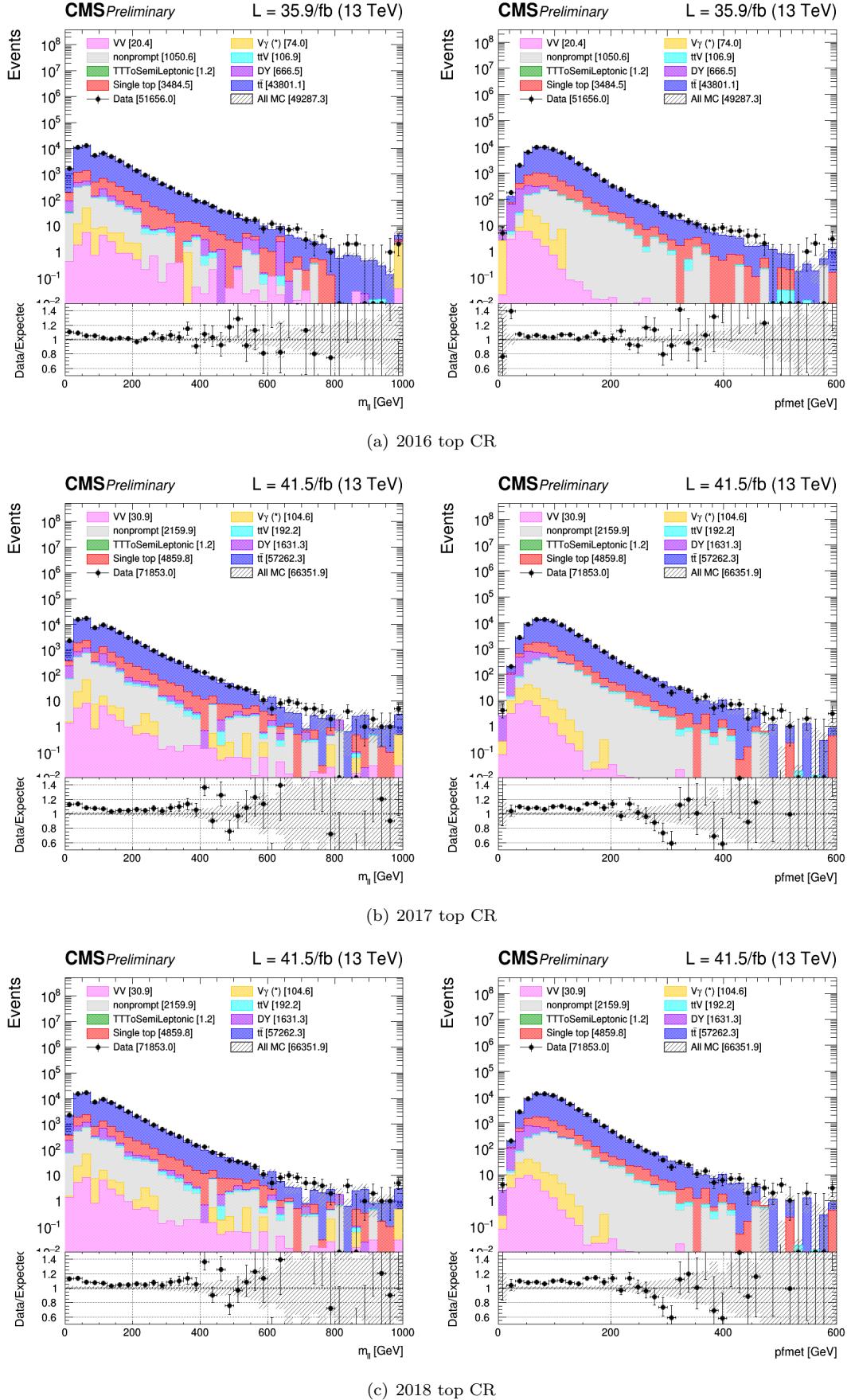


Figure 6.10: Two different variables ( $m_{ll}$ , on the left, and pfMET, on the right) represented in the  $t\bar{t}$  control region defined.

allows us to take into account the data/MC discrepancies observed in this region. Additionally, a large systematic uncertainty is associated to this particular background, strongly reduced in the signal regions thanks to the cuts applied to the analysis.

### 6.3.4 $t\bar{t} + W/t\bar{t} + Z$ control region

An additional control region targeting both the  $t\bar{t}+W$  and  $t\bar{t}+Z$  processes has also been obtained, by applying a similar selection than the pre-selection of the signal regions, but by requiring at least 3 leptons instead of 2. This leads us to a region enriched in both of these backgrounds but with a very limited statistics, given the large number of strong cuts applied. Results obtained in this case are shown in Figure 6.13.

### 6.3.5 Same sign control region

Finally, a same sign control region has also been defined in order to check the non-prompt background, calculated using a data-driven tight-to-loose method described in Section 5.6.4. This control region is defined using a similar pre-selection than the one applied to the signal regions, but by asking for two same sign leptons. Some resulting plots can be found in Figure 6.14.

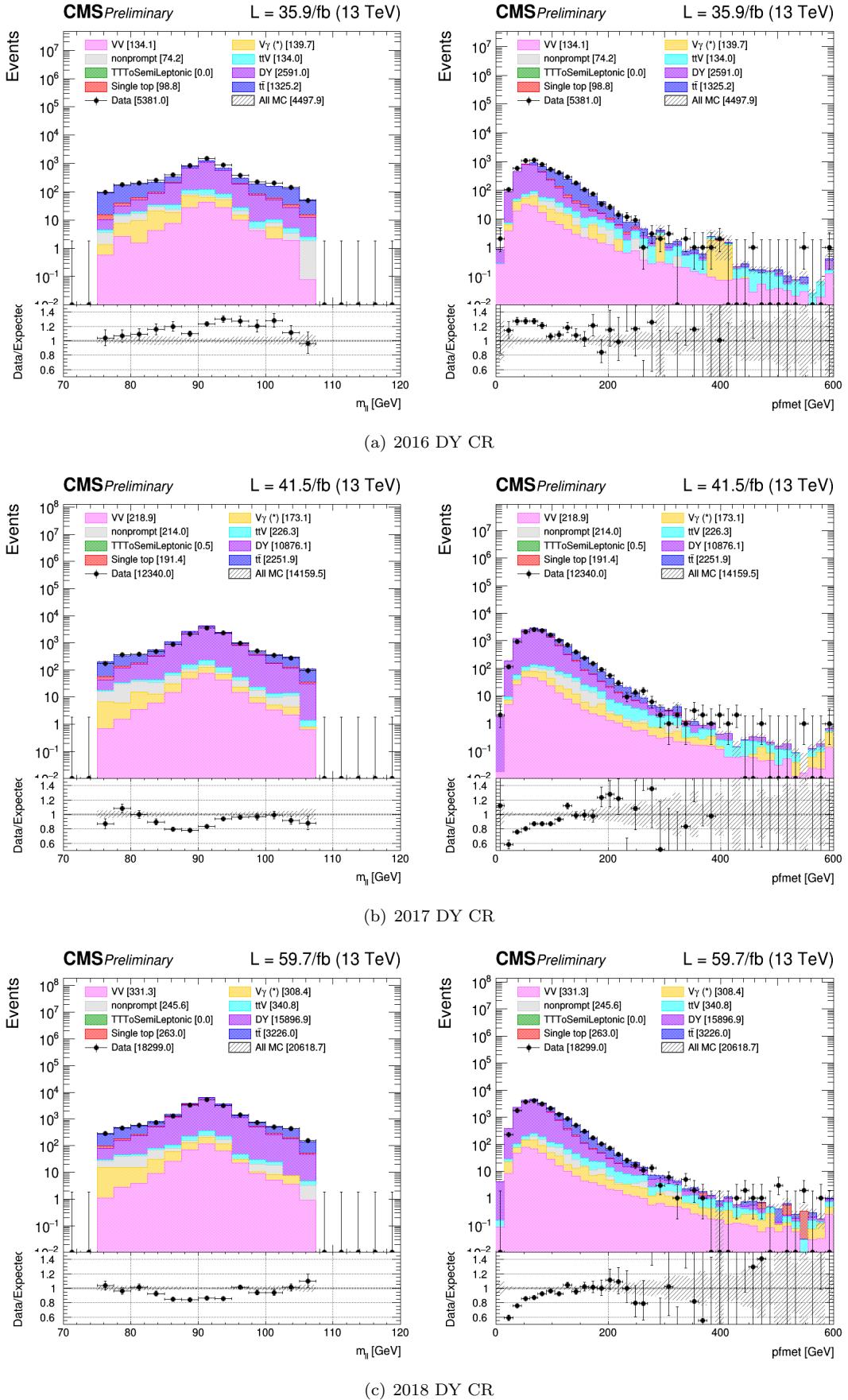


Figure 6.11: Two different variables ( $m_{ll}$ , on the left, and pfMET, on the right) represented in the DY control region defined.

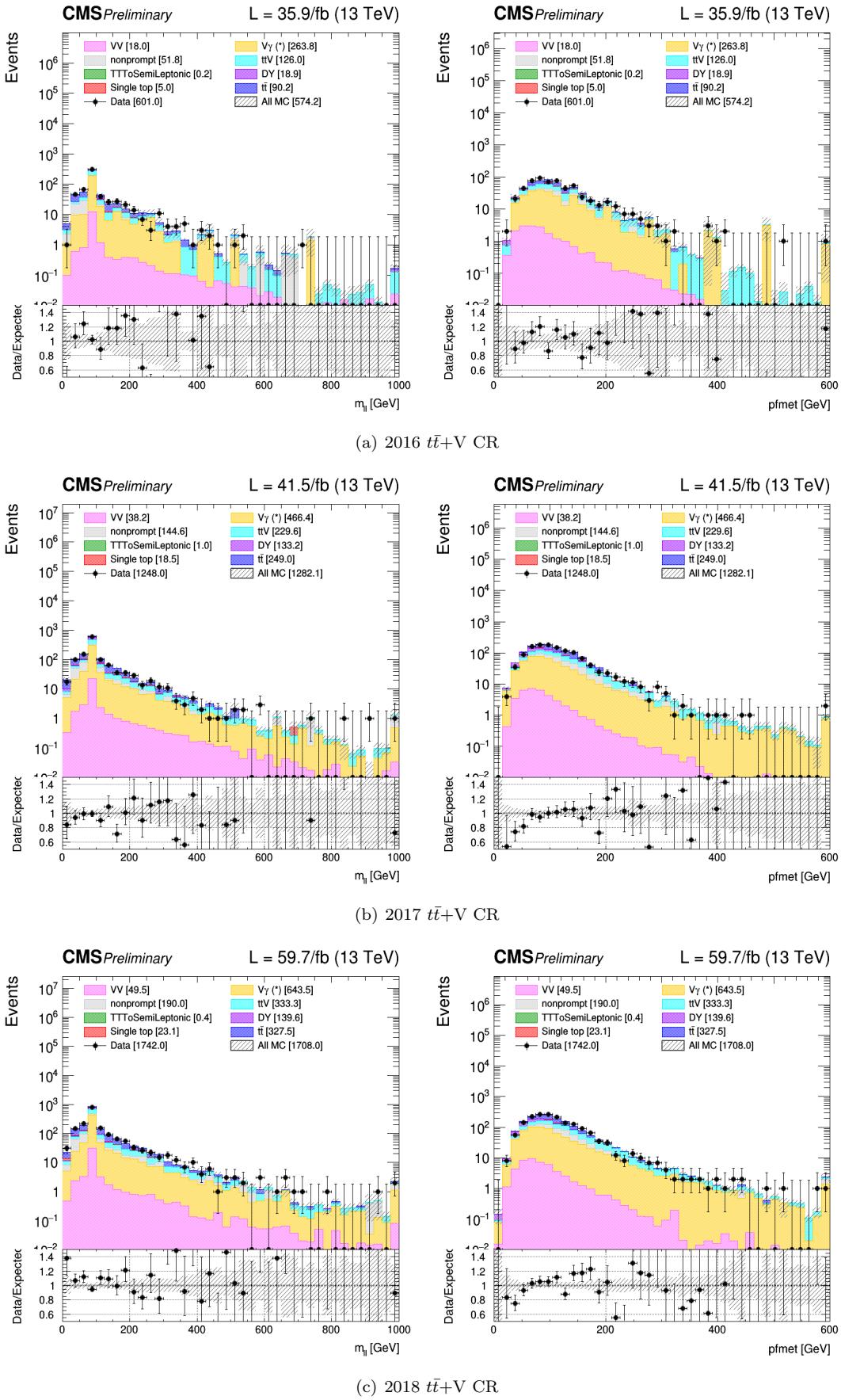


Figure 6.12: Two different variables ( $m_{ll}$ , on the left, and pFMET, on the right) represented in the  $t\bar{t}+V$  control region defined.

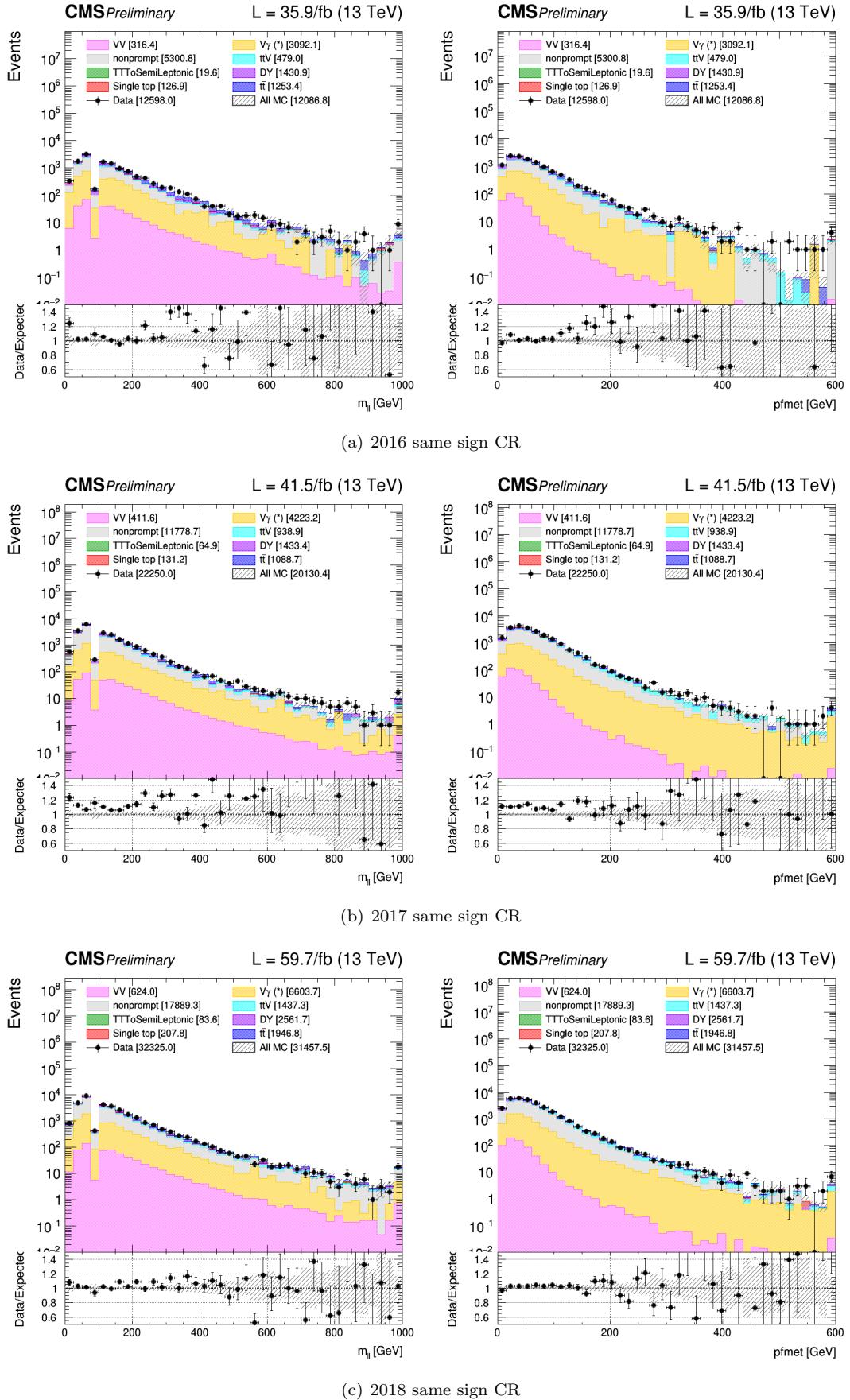


Figure 6.13: Two different variables ( $m_{ll}$ , on the left, and pfMET, on the right) represented in the same sign control region defined.



---

---

# Chapter 7

---

## Signal extraction

In this Chapter, a description about the different variables expected to naturally introduce some discrimination of the  $t/\bar{t}$  and  $t\bar{t}+{\rm DM}$  signals with respect to the different backgrounds, mainly the SM  $t\bar{t}$  and single top process, will first of all be given in Section 7.1, while a global description of the ML techniques employed in order to optimize the discriminating power of these variables in the best way possible will be detailed in Section 7.2.

### 7.1 Discriminating variables

Several variables can at principle be used in order to either separate one of the signals of interest from the different background processes. Some of the variables now presented did not feature a degree of discrimination as high as expected though and have been left behind in the actual analysis, but the ones featuring a discriminating power high enough will actually be used as input to the MVA technique developed in order to combine their discriminating power into a single variable, as described in Section 7.2.

It is also important to note that some of these variables rely on the information obtained from top reconstruction performed, which might fail in some cases. In this case, a default negative value is then assigned, making it easy to select only events which actually pass the  $t\bar{t}$  reconstruction.

#### Number of bjets and $m_{bl}^t$

Let's start with the two variables we found to be the most useful in order to separate the  $t/\bar{t}+{\rm DM}$  and  $t\bar{t}+{\rm DM}$  signals. Since both signatures have a different number of top quarks in their final state, they are also expected to lead to a different number of observed b-jets, so this quantity is an obvious choice for a good discriminating variable. However, simply rejecting events featuring two b-jets or more to define a region enriched in  $t/\bar{t}+{\rm DM}$  signal turns out in practice to be an ineffective strategy, mainly since the b-tagger efficiency for our chosen working point is of the order of 85% only, as discussed in Section 4.4.1, which can result in a large surviving  $t\bar{t}+{\rm DM}$  background, especially given the difference in cross-section between the two processes.

A slightly more effective strategy than vetoing events with multiple b-jets consists in observing that if a b-jet is produced in a top-quark decay, its invariant mass is bounded from above by  $\sqrt{m_t^2 - m_W^2} = 153$  GeV. Events compatible with two semileptonic top-quark decays can then be selected or rejected by introducing the observable  $m_{bl}^t$ , defined in Equation 7.1 [133], where the minimization is performed either over all the possible combinations of jets  $j_a, j_b$  among the b-jets of the events if three or more jets are observed, or otherwise over the b-jet(s) observed plus the non b-tagged jet having the highest b-tag weight of the event. The final distributions obtained for both variables are shown in Figure 7.1.

$$m_{bl}^t = \min (\max(m_{l_1 j_a}, m_{l_2 j_b})) \quad (7.1)$$

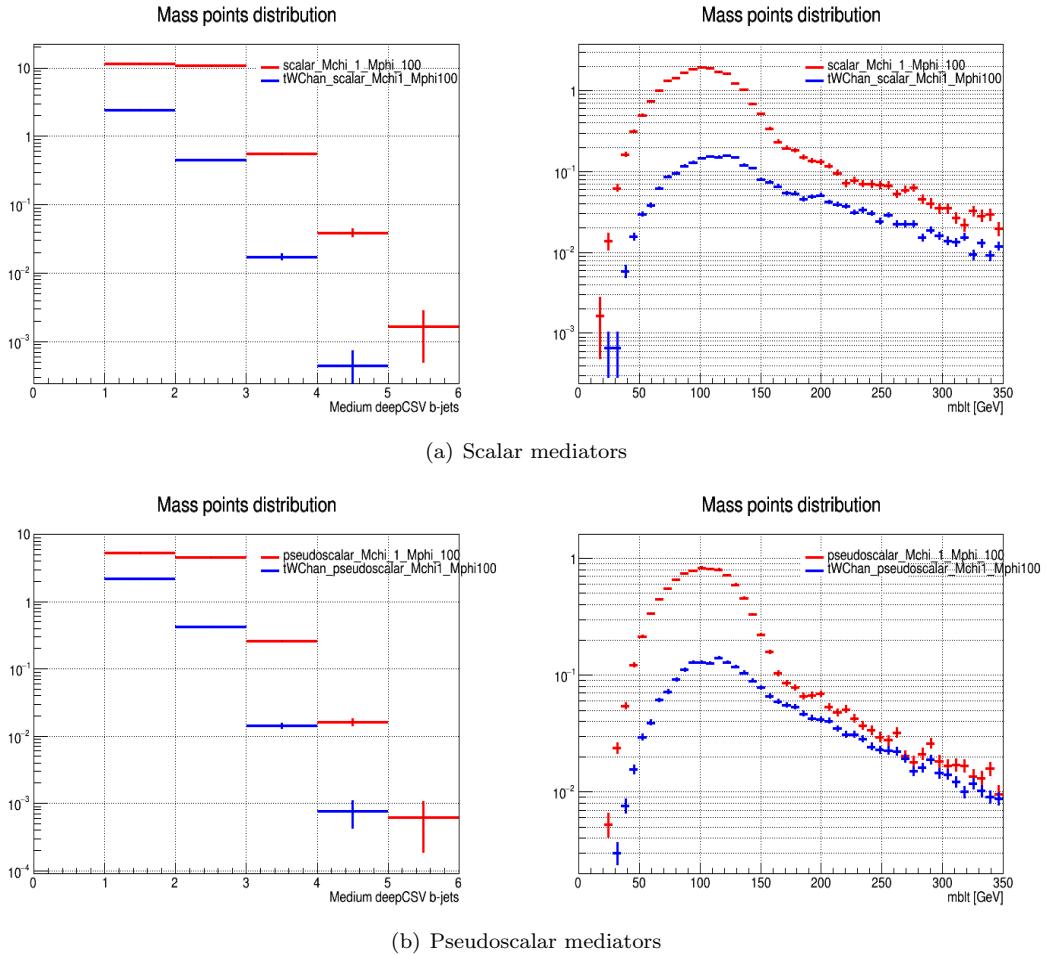


Figure 7.1: Number of b-jets (on the left) and  $m_{bl}^t$  variable (on the right), mostly used to separate our two signals in this analysis, in a control region close to the actual signal region, for different mediator categories and signal samples of the analysis.

As we can see in these plots, these variables do offer us some discrimination between the different signals considered but even by using them, we cannot easily define a signal region enriched in  $t/\bar{t}+DM$  because of the low cross-section of this process. This is one of the main points on which the MVA will be able to become useful later on.

### Missing transverse momentum

Already defined in Section 4.5, this variable corresponds to the imbalance in transverse momentum which can be left by different phenomena, such as the apparition of a SM neutrino or the existence of DM particles, able to escape the detector without being detected.

This is one of the most important variables of this analysis, expected to induce some discrimination between the signal and the backgrounds because, even though some SM processes such as the SM  $t\bar{t}$  production in the dilepton final state is also expected to produce some neutrinos and therefore some MET, both the  $t/\bar{t}+DM$  and  $t\bar{t}+DM$  signal model are expected to have mostly the same contribution to the MET from their own neutrinos, plus an additional contributions from the pair  $\chi\bar{\chi}$  produced. The MET spectrum is therefore expected to reach higher values for the signals than the backgrounds, as shown in Figure 7.2 in the pre-selection region of the analysis.

**FIXME: Put final plots for all the variables once available**

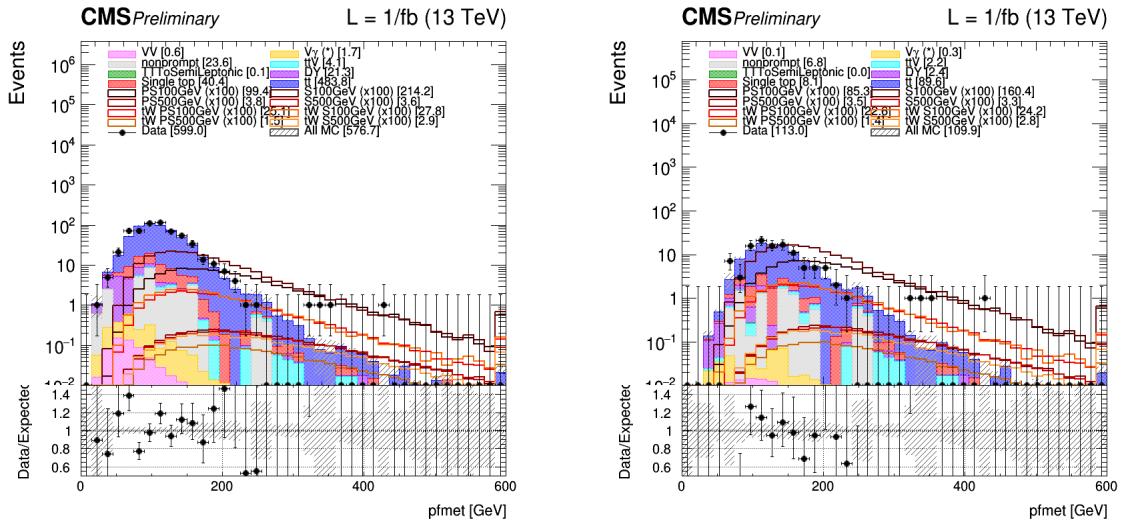


Figure 7.2: pfMET distribution in the 2018 pre-selection (on the left) and 2018 scalar 100 GeV signal (on the right) signal regions.

### Stransverse mass

The  $m_{T2}$  variable, also called **stransverse mass**, is an extension of the definition of the transverse mass  $m_T$  to cases when pairs of particles with the same flavor decay into one visible and one invisible particle, such as what happens in the double  $W \rightarrow l\nu$  decay, for example.

For the signals considered in this work, two neutrinos contribute to the presence of MET and the individual contribution of each particle ( $\not{p}_{T_1}$  and  $\not{p}_{T_2}$ ) to this missing energy cannot be inferred. The stransverse mass is then defined according to Equation 7.2, where  $\not{p}_{T_i} = \vec{p}_{T_i}$  is the (visible) transverse momentum of the particle  $i$  and  $\alpha$  is the angle between the visible and invisible  $p_T$  of the particles involved in the decay considered [134].

$$\begin{cases} M_{T2}^2 = \min_{\not{p}_{T_1} + \not{p}_{T_2} = \not{p}_{T_{\text{tot}}}} \left( \max \left( m_T^2(\not{p}_{T_1}, \not{p}_{T_1}), m_T^2(\not{p}_{T_2}, \not{p}_{T_2}) \right) \right) \\ m_T^2(\not{p}_T, \not{p}'_T) = 4 |\not{p}_T| |\not{p}'_T| \sin^2 \left( \frac{\alpha}{2} \right) \end{cases} \quad (7.2)$$

This equation can be understood in the following way: to compute the  $m_{T2}$  variable, different

combinations ( $\not{p}_{T_1}$ ,  $\not{p}_{T_2}$ ) satisfying the condition  $\not{p}_{T_1} + \not{p}_{T_2} = \not{p}_{T_{\text{tot}}}$  need to be probed, keeping only the combination which results in the lowest possible value.

In this particular analysis,  $M_{T_2}^{ll}$  is calculated from a general algorithm described in [135], since the role of the visible particles is played by the two final state leptons. This variable is expected to introduce some discrimination because the  $M_{T_2}^{ll}$  variable for a SM  $t\bar{t}$  process is expected to have an endpoint exactly at the mass of the W boson, while an eventual DM signal does not have this limitation in the  $M_{T_2}^{ll}$  spectrum because of the pair of DM particles produced, which also contributes to the total MET of the event. In practice however, we do observe a tail in this spectrum even for SM  $t\bar{t}$  without DM, because of the instrumental MET sometimes observed or the fact that some selected leptons are not actually prompt leptons but can be jets misidentified as leptons by the detector, as shown in Figure 7.3.

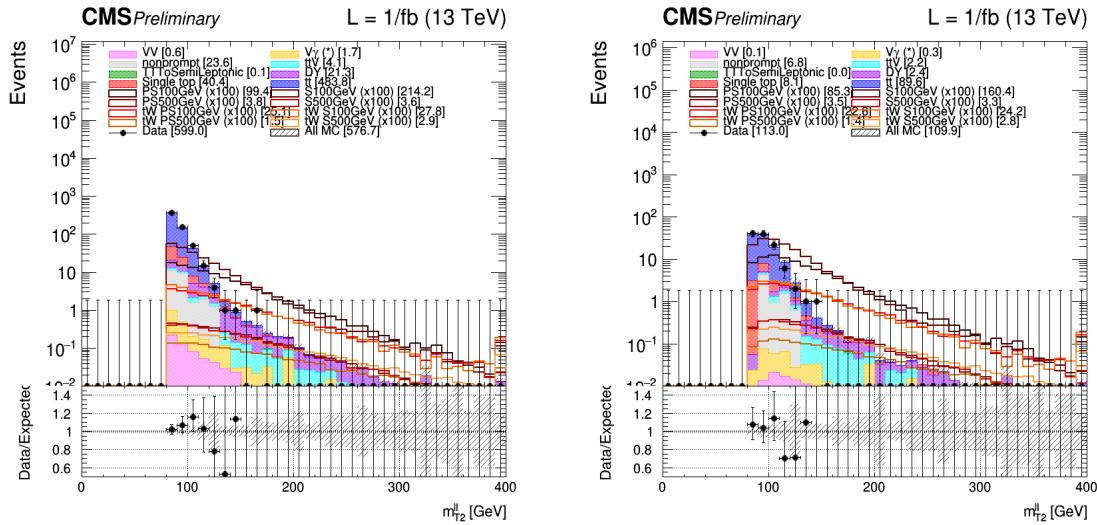


Figure 7.3:  $M_{T_2}^{ll}$  distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

A second variable,  $M_{T_2}(bl, bl)$ , the stranverse mass of the b-jet lepton pairs can also be defined. This variable is constructed in the same way, except that in this case, the lepton is paired with a b-jet. The b-jet/lepton permutation giving the smallest value of  $M_{T_2}(bl, bl)$  is kept. The distribution of this variable in our signal regions is represented in Figure 7.4.

### Dark $p_T$ and overlapping factor $R$

Already introduced with the top reconstruction method in Section 4.6.2, these two variables are by construction expected to give some discrimination between the signals and the different background processes, mainly the SM  $t\bar{t}$ .

On one hand, the dark  $p_T$  is for example expected to take slightly higher values for the signals, where we expect the ellipses to be further apart from each other due to the production of DM. On the other hand, the overlapping factor  $R$  defined in Equation 4.8 is expected to peak around 1 for the standard  $t\bar{t}$  process (since in this case, we expect to have  $d = l_1 + l_2$  for most of the events, as seen in Figure 4.14), and to take slightly lower values for the different signal mass points considered. This intuition is confirmed from the plots represented in Figure 7.5.

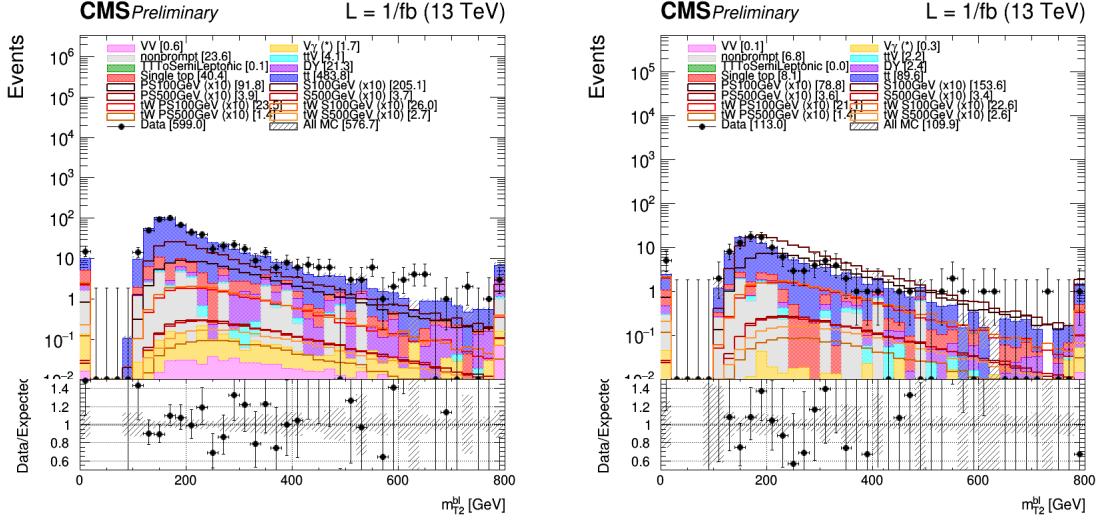


Figure 7.4:  $M_{T2}(bl, bl)$  distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

### Additional variables

A few additional variables which are expected to introduce some discrimination between the signals and the backgrounds have been considered as well. Among them, we can for example quote:

- The  $m_{bl}^t$  variable and number of b-jets, already described in the previous section and whose distributions can be found in Figures 7.6 and 7.7 respectively.
- $\Delta\Phi(E_T^{\text{miss}}, ll)$ : the distribution in  $\Phi$  of the two leptons is expected to change depending on the eventual production of DM, as shown in Figure 7.8. Some distributions obtained for this particular variable can additionally be found in Figure 7.9. In the same way, the difference in azimuthal angle between the system "b-jet and closest lepton" and the pfMET  $\Delta\Phi(lb^{\Delta R_{\min}}, E_T^{\text{miss}})$  has also been considered, as shown in Figure 7.10.
- The significance of the MET, shown in Figure 7.11, which evaluates the likelihood the the measured MET of an event is due to a resolution fluctuation because of detector-related limitations like finite measurement resolution [158]. This variable does feature a high degree of discrimination but is expected to be highly correlated to the usual MET variable.
- The kinematic reconstruction weight  $W$  obtained from the  $t\bar{t}$  reconstruction process is also expected to give us some discrimination to isolate our signals, as shown in Figure 7.12 for reasons already explained in Section 4.6.
- Two other interesting variables used by the ATLAS collaboration for their own analysis and shown in Figure 7.13 are the so-called  $r2l$ , defined as the ratio between the pfMET and the  $p_T$  of the two leptons observed and  $r2l4j$  variables, defined in a similar way but considering additionally the  $p_T$  of the first 4 jets (if they exist) in the sum in the denominator.
- We also observed that the variable  $massT$  which corresponds to the scalar sum of the transverse component of the pfMET, the two leptons and the two b-jets obtained by the top reconstruction process helps with the discrimination process, as shown in Figure 7.14.
- The spin correlation in a  $t\bar{t}$ -like event is expected to be conserved, because of the short lifetime of the top quark, and can actually be inferred from the top quark decay products, accessible to us from the top reconstruction method described in Section 4.6. These variables are

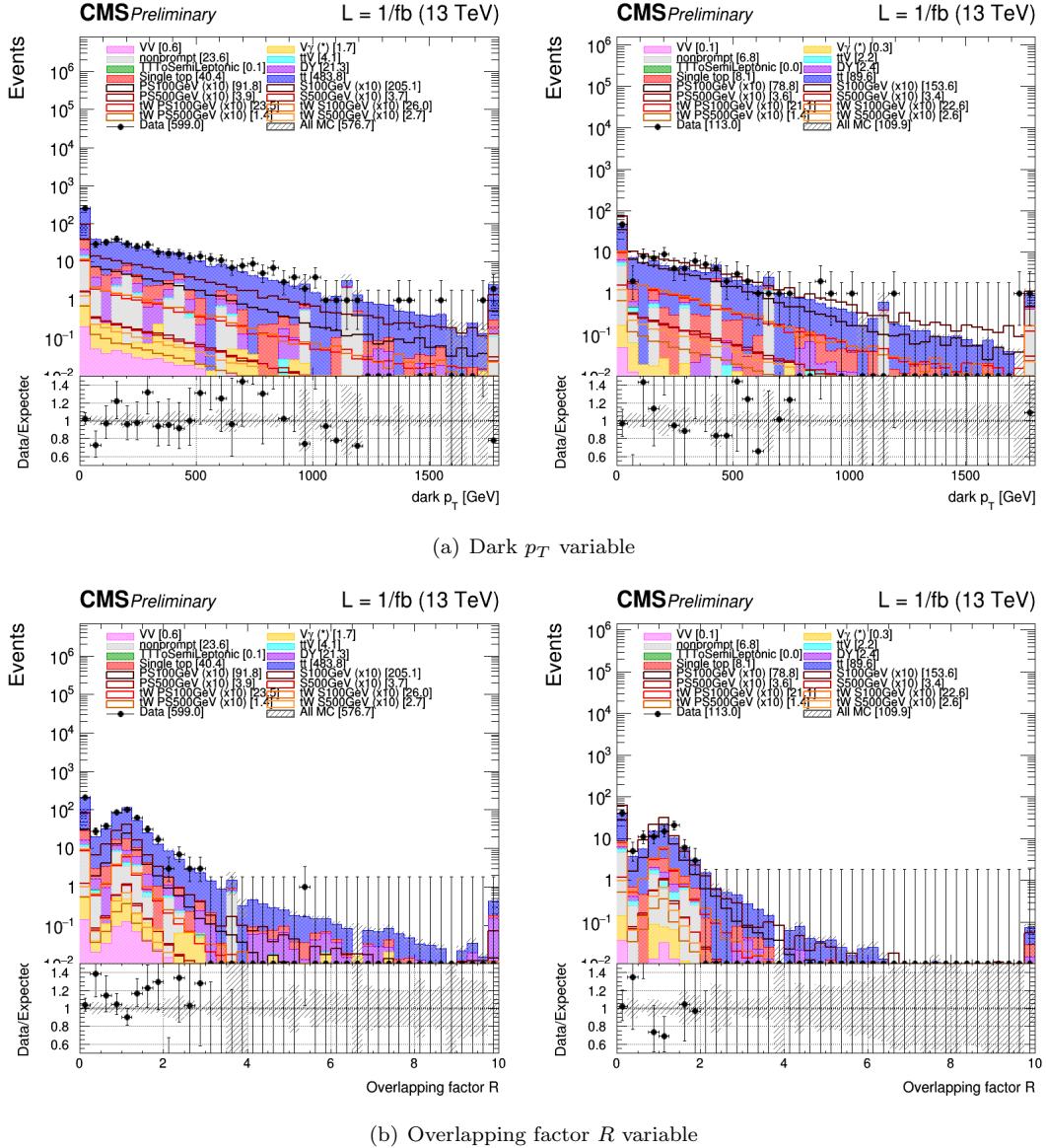


Figure 7.5: Dark  $p_T$  and overlapping factor  $R$  distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) regions (on the right) for the backgrounds and several signal samples.

interesting because the spin correlation in such events depend on the production mechanism and will be influenced by the additional coupling to a scalar or pseudoscalar mediator.

Several variables belonging to this category have been considered, such as  $\xi = \cos(\theta_l) \cos(\theta_{\bar{l}})$ , shown in Figure 7.15, and  $c_{\text{hel}}$ , the full angle between the two leptons in their parent mass frame shown in Figure 7.16.

- Finally, general variables coming from different systems have also been computed since they can feature some discrimination. Among this category, we considered the invariant and transverse masses along with the  $\Delta\phi$  angles of the different children in both the  $t\bar{t}$  and  $t\bar{b}b$  systems, where the b-jet is considered to be the b-jet with the highest  $p_T$  in the event, as shown in Figures 7.17 and 7.18 respectively.

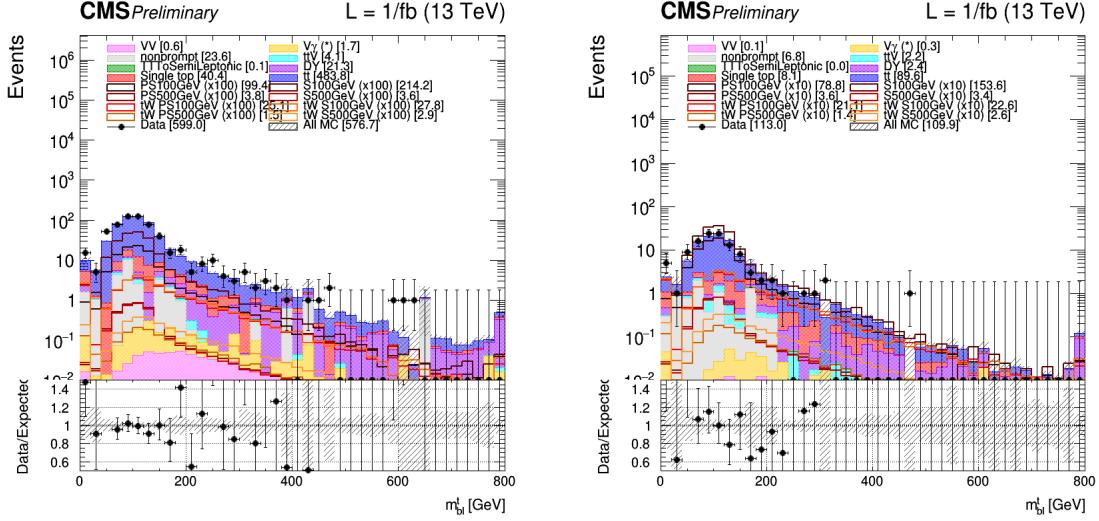


Figure 7.6:  $m_{bl}^t$  distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

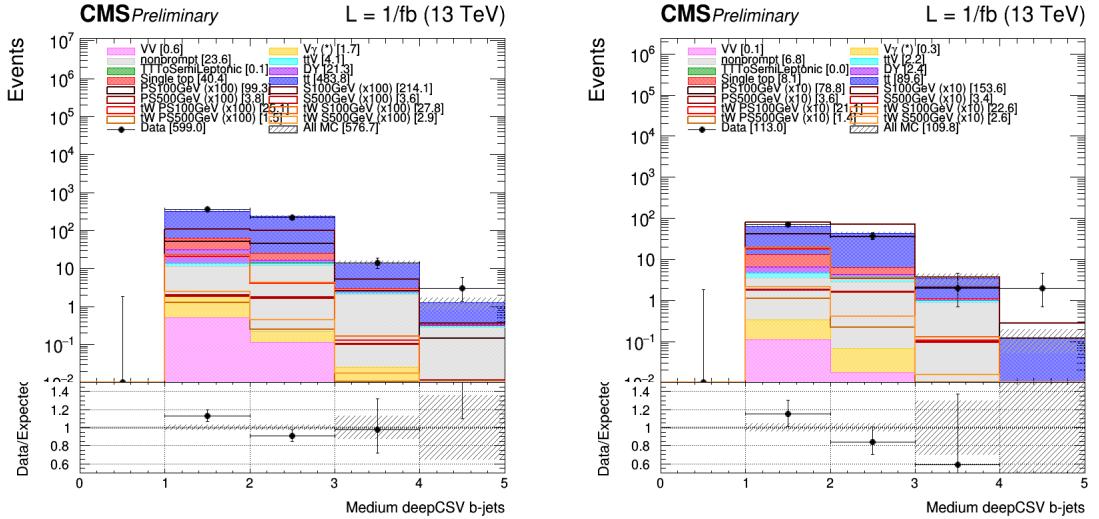


Figure 7.7: Number of b-jets distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

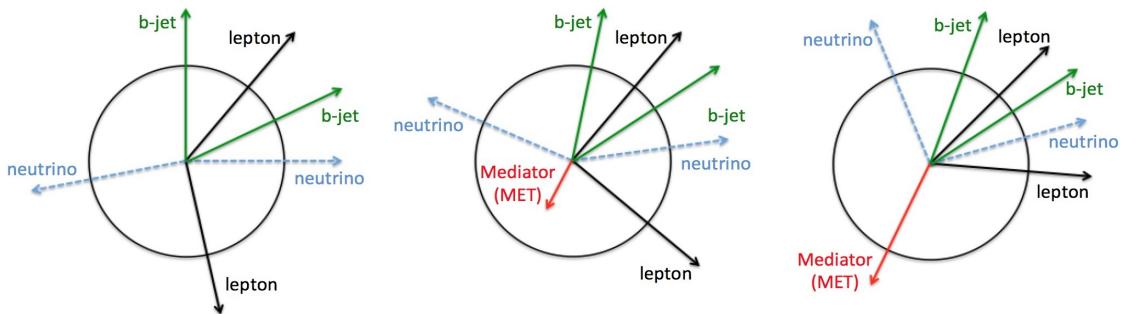


Figure 7.8: Schematic representation in the  $\Phi$  plane of the distribution of the particles for the  $t\bar{t}$  process (on the left) and for the  $t\bar{t} + \text{DM}$  (on the center and right).

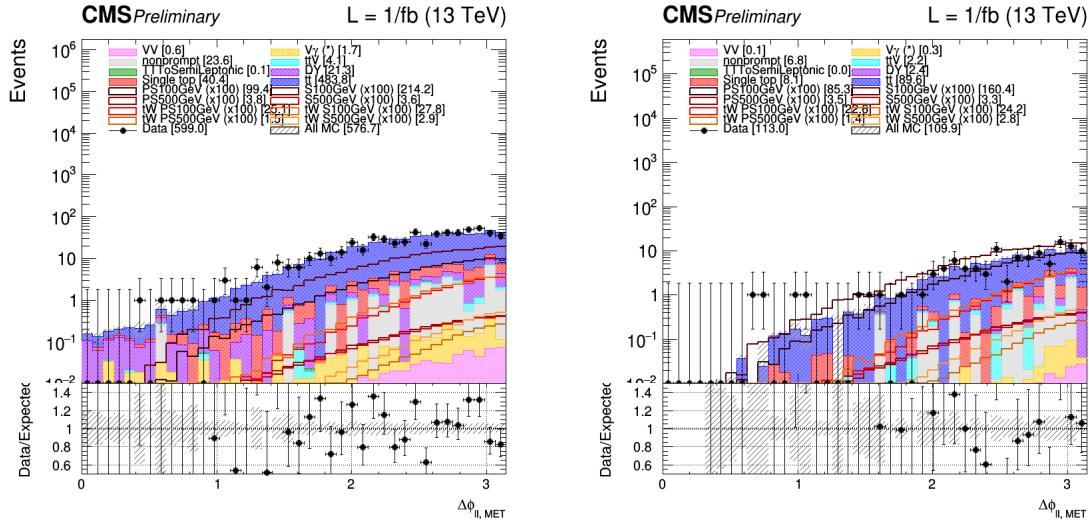


Figure 7.9:  $\Delta\Phi(E_T^{\text{miss}}, ll)$  distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

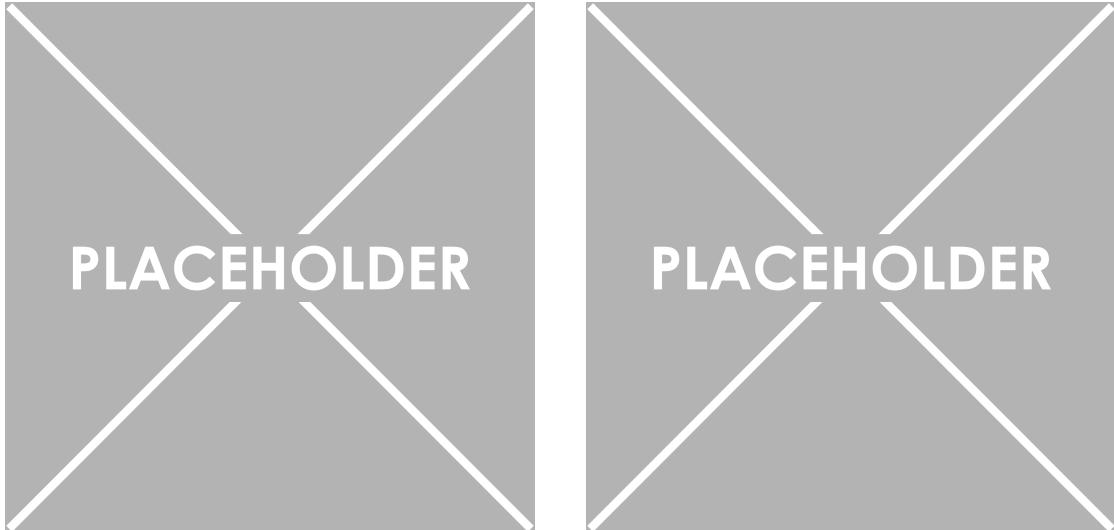


Figure 7.10:  $\Delta\Phi(lb\Delta R_{\min}, E_T^{\text{miss}})$  distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

## 7.2 Multivariate analysis

As already seen in the previous sections, we know that our signal regions are going to be contaminated with some background processes, given the fact that such backgrounds have kinematics close to the signal searched for. A general Multi-Variate Analysis (MVA) technique has therefore been used in order to perform the actual signal extraction of the analysis by combining the discriminating power of several of the variables previously presented, in order to find a way to characterize reliably any single event as either more likely to be a signal (either  $t/\bar{t}+DM$  or  $t\bar{t}+DM$ ) or a background event.

The idea behind this kind of techniques is quite simple. We have at our disposal on one hand MC simulations of the most interesting background and signal processes and on the other hand, we have unknown data collected by the detector and which needs to be characterized and labeled based on the value of several variables or features, by assigning a value of probability of belonging to any of these 3 particular categories to every single event. By definition, these MC simulations

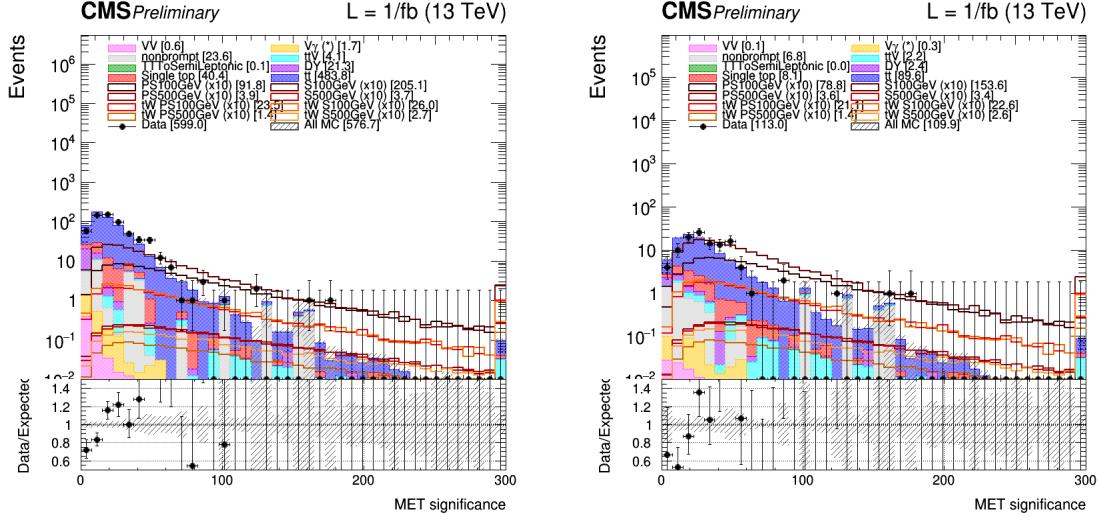


Figure 7.11: MET significance distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal regions.

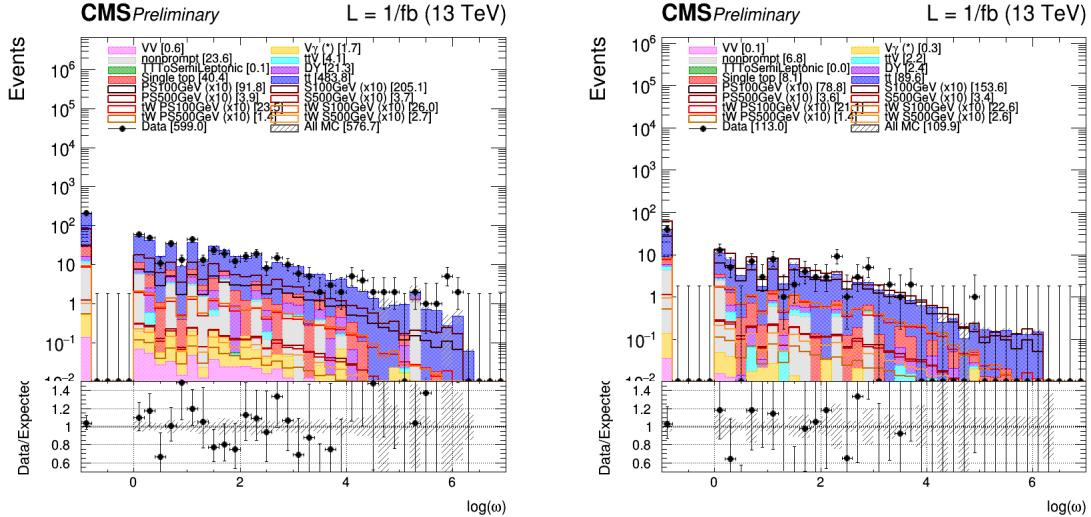


Figure 7.12: Kinematic reconstruction weight  $W$  distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

do come with a label already assigned, and this set of events can be divided into two subsets: one used for the adjustment of different internal parameters of the model used (the so-called **training process**), while the other subset is used in order to evaluate the performance of the classifier that has been defined (the so-called **testing process**).

### 7.2.1 Methods used

Two different MVA methods among the most popular nowadays have been studied in this work: a Boosted Decision Tree (BDT) and a Analysis Neural Network (ANN) have been defined and optimized separately in order to get the best discriminating power possible. Details regarding this optimization process can be found in Appendix C.

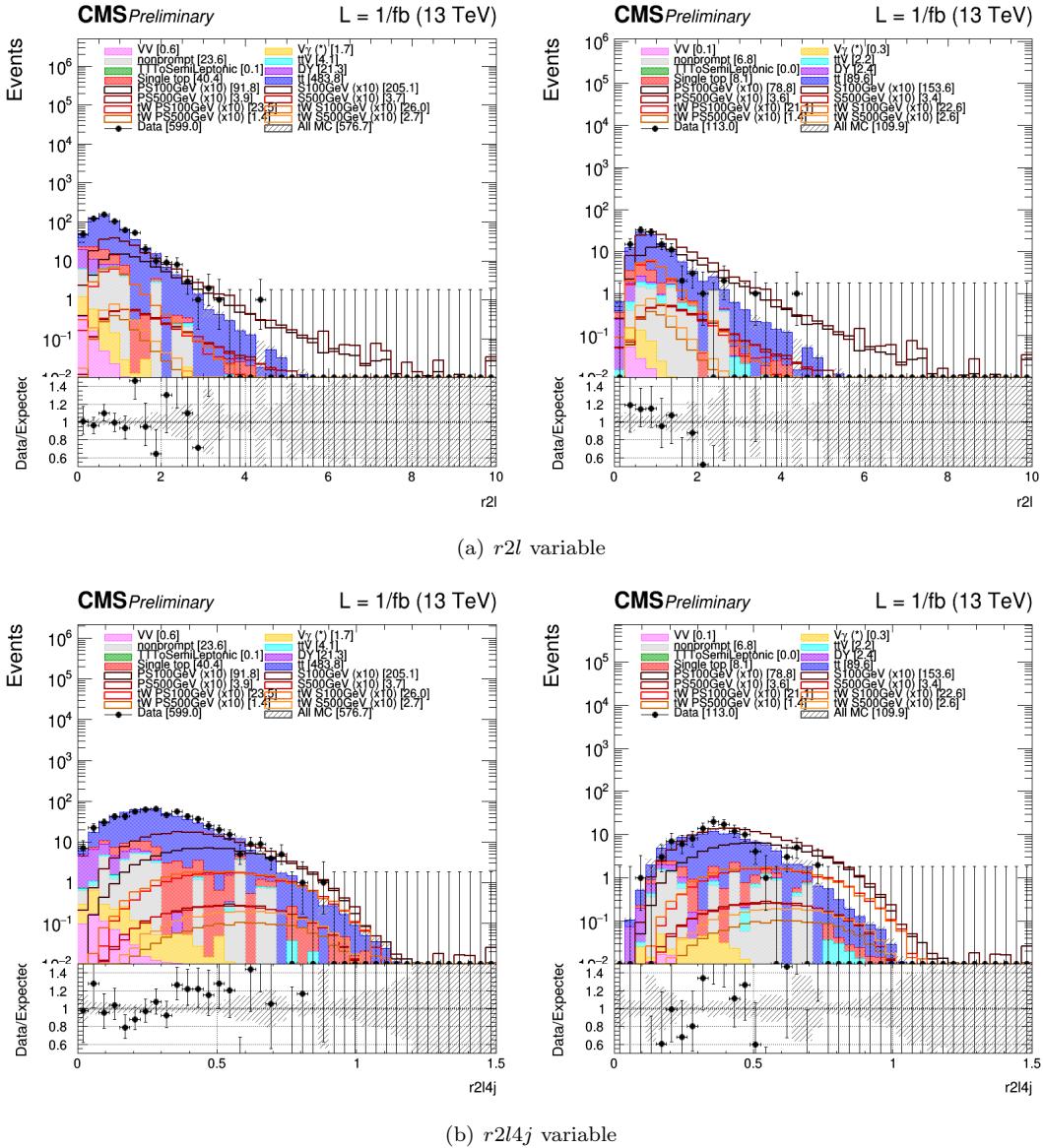


Figure 7.13:  $r2l$  (on the top) and  $r2l4j$  (on the bottom) variables considered for the signal discrimination process in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

### Boosted Decision Trees (BDTs)

This method relies on the definition of different decision trees, able to split the data recursively based on a set of input variables or **features**. A typical BDT, as shown in Figure 7.19, is then made out of several **nodes**, allowing to make a binary decision (yes/no) and split the data based on a given features until reaching a terminal condition and therefore a **leaf**, terminal node representing a class label or a probability. The split performed is obviously not chosen randomly, as a thorough training process needs to take place in order to define the optimal splitting based on the features given in such a way to maximize information gain. The **boosting** process then consists in training several trees in such a way and then to combine all these trees together into a single strong classifier, by giving a weight to each tree depending on their actual accuracy. In our case, a classification is performed, by defining a BDT able to read several input variables at once, create hundred of trees made out of different nodes, train them and apply the training to uncategorized events, labeling every single event as signal or background-like.

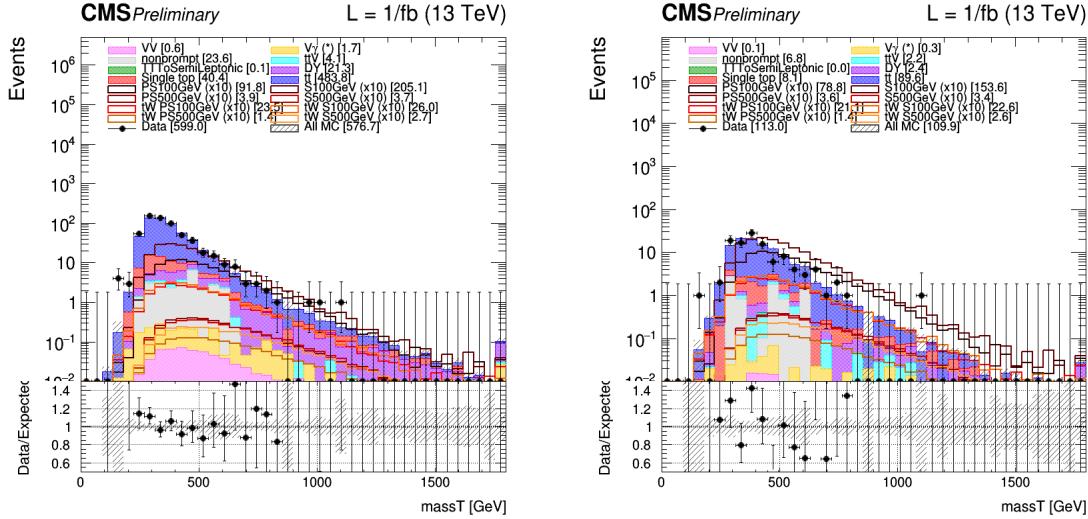


Figure 7.14: Scalar sum of the transverse energy of the particles produced in the events distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

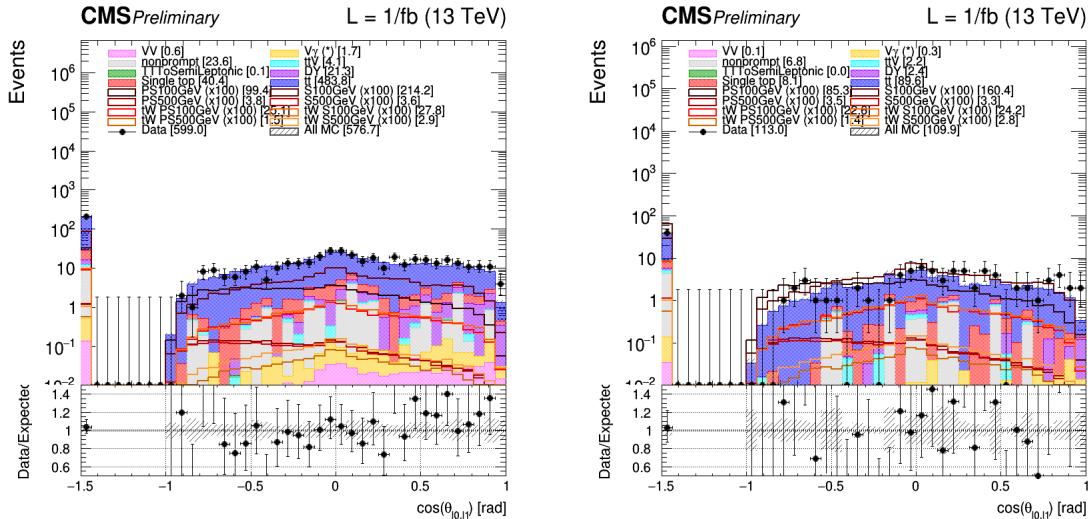


Figure 7.15:  $\xi = \cos(\theta_l) \cos(\theta_t)$  distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

The structure and complexity of a BDT usually depends on the problem considered and can be characterized from several different parameters which need to be tuned in order to get the best possible signal extraction accuracy while avoiding overfitting:

- First of all, the **features**, or input variables, allowed to be used by the BDT are extremely important, and a maximum number of features can be set when building a tree, forcing a random selection of input variables if too many features are given as input.
- The **maximum depth** is usually chosen as a small number which characterizes the maximum number of vertical levels allowed for the BDT.
- The **minimum samples per leaf** puts a requirement on the minimum number of samples required in order to create a new leaf.
- Regarding the boosting process itself, several additional parameters can be defined:
  - The **loss function** is helpful when trying the distance between the prediction made by the BDT and the actual value expected.

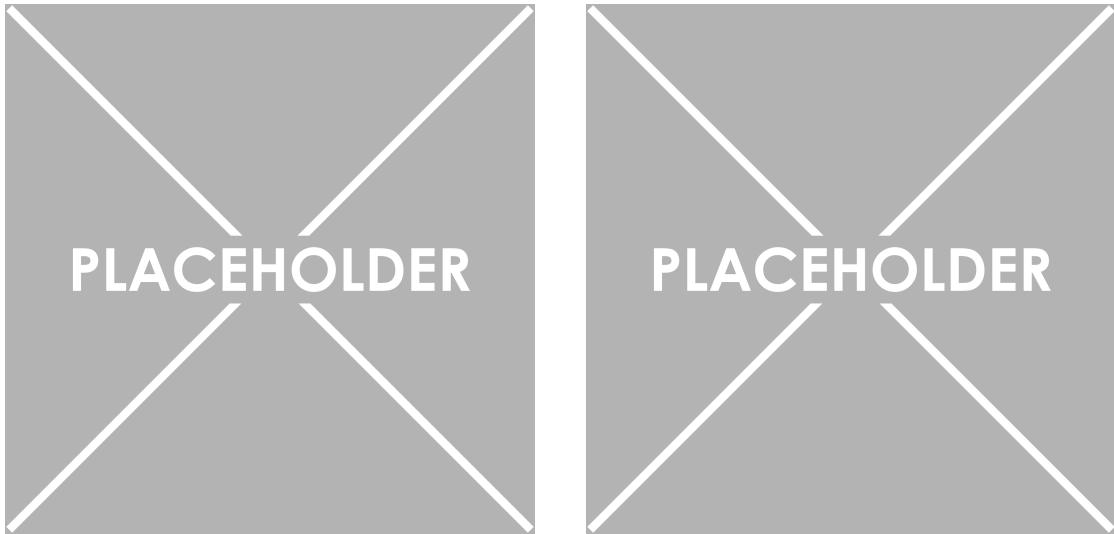


Figure 7.16:  $c_{\text{hel}}$  distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

BDT parameter	Optimized value
Maximum depth	4
Minimum samples per leaf	1%
Loss function	Quadratic
Boost algorithm	Gradient descent
Shrinkage	0.2
Grid points $n_{\text{cut}}$	250
Number of trees	100

Table 7.1: Summary of the parameters used for the training of the BDT in this analysis.

- The **learning rate**, or **shrinkage**, is another important parameter in the sense that it tells how much the weights of the tree should be adjusted after each training iteration. A small shrinkage typically demands more trees to be grown but can significantly improve the accuracy of the prediction made.
- The **number of grid points**  $n_{\text{cut}}$ , defined as the granularity used to find the optimal cut when determining the node splitting is also important.
- Finally, the actual **number of trees** used to define the forest is obviously an important parameter as well that need to be optimized.

The actual parameters used for this particular analysis obtained after a thorough optimization process can be found in Table 7.1.

### Analysis Neural Networks (ANNs)

Another large family of MVA methods rely on the so-called **neural networks**, sets of algorithms designed to be able to learn how to recognize patterns. Neural networks help us cluster any unknown dataset, allowing us to group and classify unlabeled data according to either similarities among the example inputs or based on the information available in a training dataset.

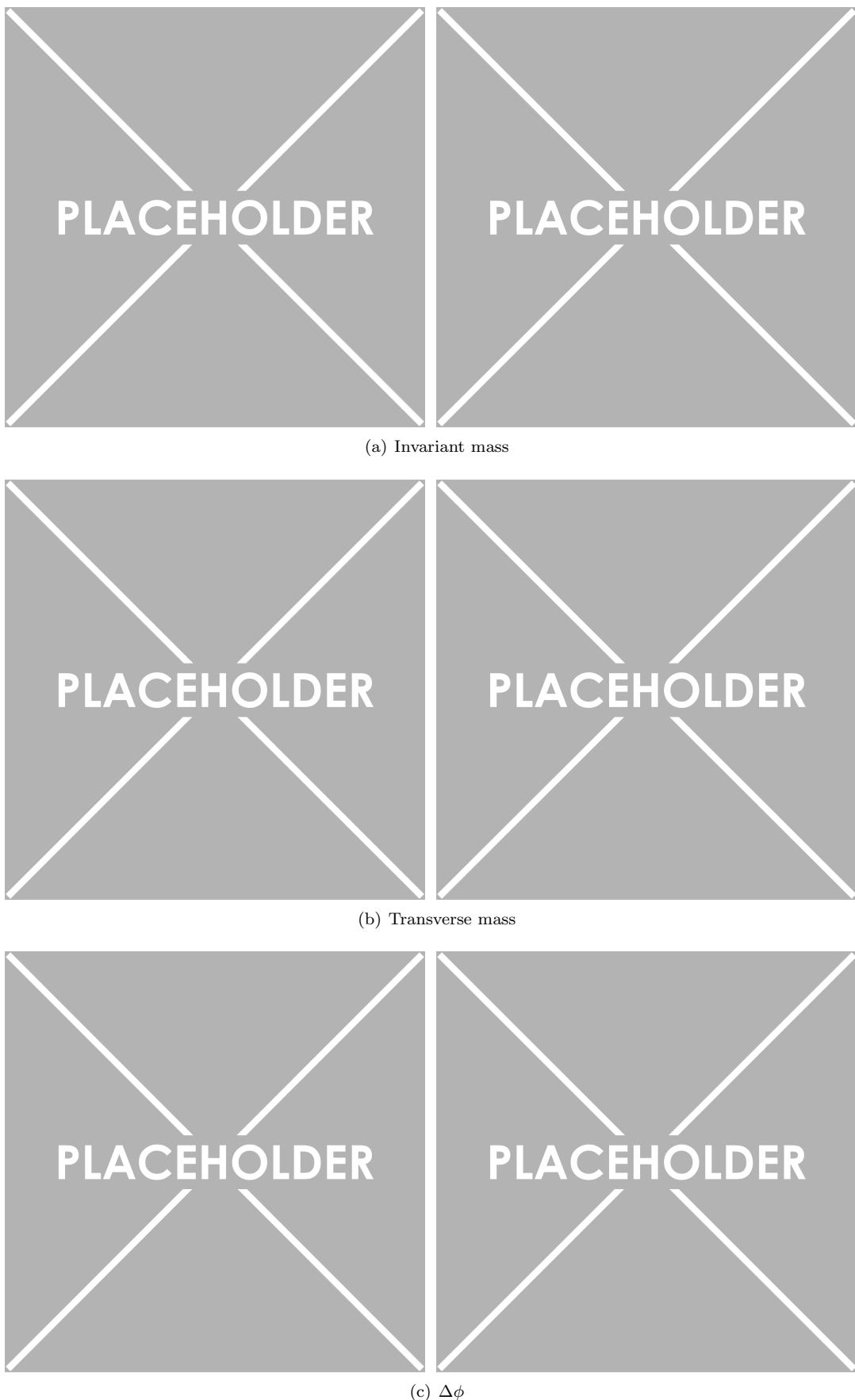


Figure 7.17: Some kinematic variables computed in the  $t\bar{t}$  system and considered for the signal discrimination process in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

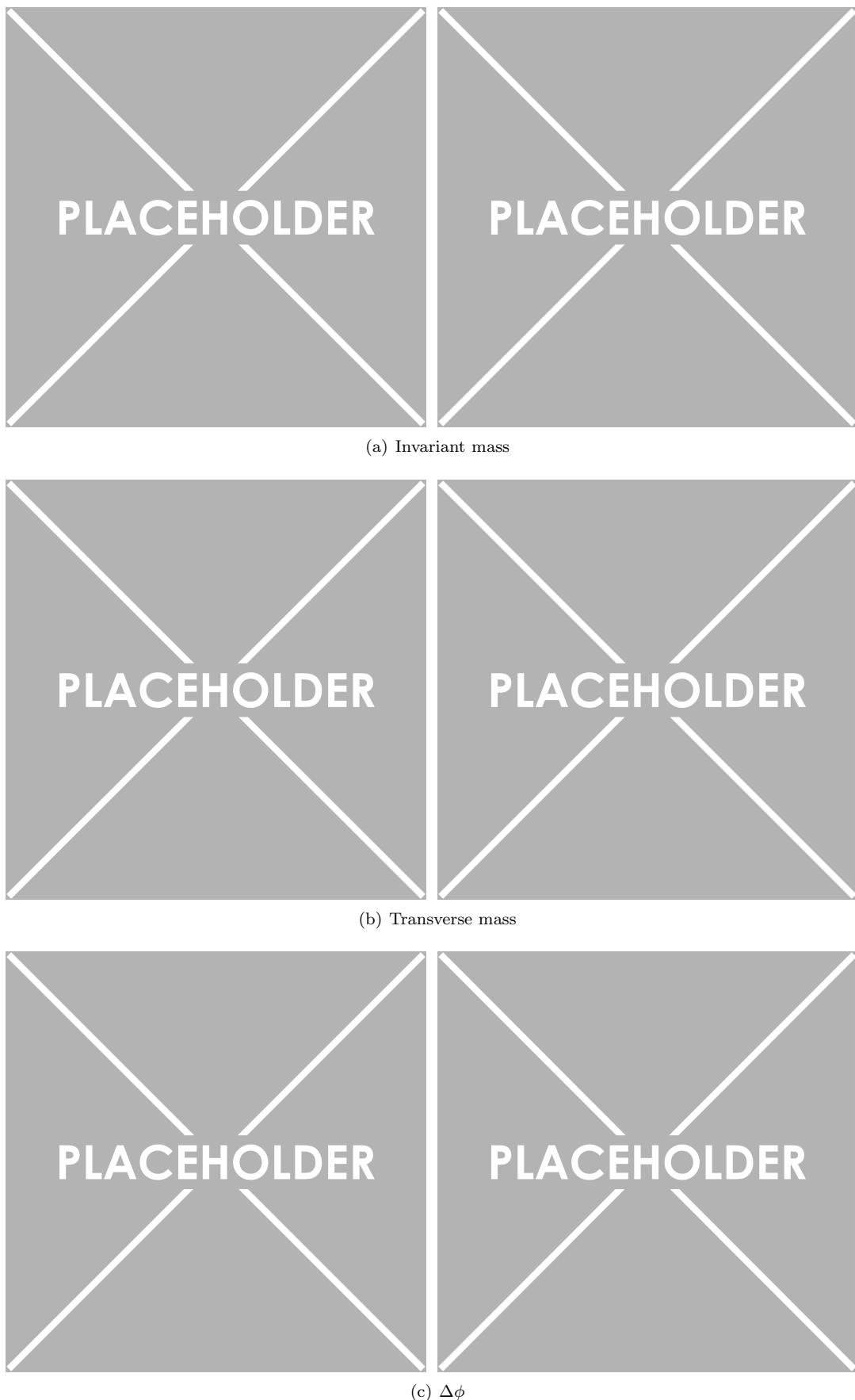


Figure 7.18: Some kinematic variables computed in the  $l\bar{b}$  system and considered for the signal discrimination process in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions.

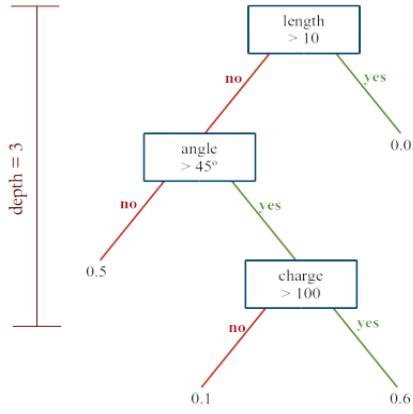


Figure 7.19: Schematic representation of a typical BDT, with its nodes (represented by the blue boxes) and leaves [159].

The basic idea for the existence of neural networks relies on the human brain itself, a complex system of thousands of billions of neurons interacting with each other and able to perform complex calculations and generate ideas. In this case, highly connected mathematical units called **artificial neurons** and grouped into layers are then defined, as shown in Figure 7.20, connected with each other with connections (also called **edges**) characterized by a **weight**, denoting its significance. At the end of the day, these neurons are simply mathematical tools able to receive an input, modify it and send the computed signal to the next layer of neurons, while the connection to the next neuron will either increase or decrease this signal based on the weight of this connection.

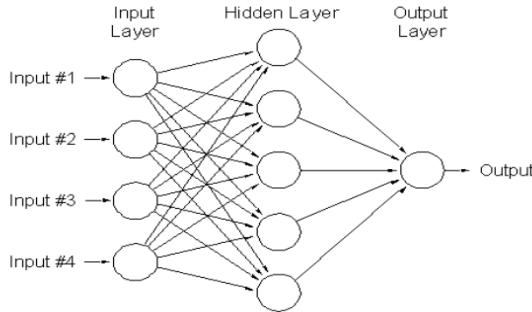


Figure 7.20: Schematic representation of a typical neural network with a single hidden layer [160].

Mathematically, each neuron  $i$  of an intermediate layer receives an input  $x_j$  from each of the  $n$  neurons of the previous layers, and combines them by multiplying them with the weight  $\omega_{ij}$  of the corresponding connection between the two neurons involved. Once done, the neuron then applies an internal non-linear **activation function**  $f$  to this previously computed value, and fires the computed final output  $y_i$  to the next layer of neurons, as shown in Equation 7.3, and this process is repeated until reaching the output layer.

$$y_i = f \left( \sum_{j=0}^n \omega_{ij} x_j \right) \quad (7.3)$$

The objective of the training process in this case is to find the optimal weights  $\omega$  for each connection that are able to deliver an output as faithful to the known event category as possible. This faithfulness of a given set of weights can be obtained from the so-called **error function**, which can take different forms but usually simply computes the normalized sum of the difference between

DNN parameter	Optimized value
Hidden layers neurons	80, 80, 40
Activation functions	Relu (x3), softmax (output)
Error function	Mean square error
Optimizer	Adam
Learning rate	0.005
Training epochs	250
Batch size	250

Table 7.2: Summary of the parameters used for the training of the ANN in this analysis.

the expected  $\hat{y}_i^p$  and obtained  $y_i^p$  values for each neuron  $i$  and each of the  $p$  events found in the training dataset, as shown in Equation 7.4.

$$E(\boldsymbol{\omega}) = \sum_{i,p} \left( y_i^p(\boldsymbol{\omega}) - \hat{y}_i^p \right)^2 \quad (7.4)$$

This error function allows us to determine the goodness of the output obtained from a given set of weights, but we still need a way to update these weights at any given training iteration to make this method useful. This can be done using different methods, such as by applying the back-propagation or performing a gradient descent method, which consists in starting from a random set of weights  $\boldsymbol{\omega}_0$  and modifying them in the direction of the gradient of the error function  $E$  at each iteration  $k$ , as shown in Equation 7.5, with a given **learning rate**  $\eta$ .

$$\boldsymbol{\omega}_{k+1} = \boldsymbol{\omega}_k - \eta \nabla_{\boldsymbol{\omega}} E(\boldsymbol{\omega}) \quad (7.5)$$

In this work, we used models with several hidden layers, therefore sometimes referring it as deep neural networks, where the **depth** parameter is equivalent to the number of hidden layers. The use of such methods with lots of hidden layers is helpful because it allows to solve the feature engineering problem, which states that to get the most of a certain dataset someone should previously extract by hand the representations/features helping to solve the problem. However, with deep learning, the features are learned on the fly, without the need of doing a previous transformation by hand. In summary, the ANNs allow to study phenomena a bit more disrupt, but are usually more sensitive to the overfitting even though some options, such as the dropout, dropping randomly some weights in the network at each training iteration, have been used in this work to mitigate this issue.

Several parameters are important to optimize during the training of a neural network, such as the number of neurons in each layer, the activation and error functions, the learning rate (giving an idea of how much the weights should be updated after each training iteration, to be optimized in order to reach the minimum of the cost function faster while avoiding getting stuck in an eventual local minimum), the number of **training epochs** and the **batch size**, allowing to reduce the time required for the training by dividing the dataset into mini-batches and updating the weights for each mini-batch instead of doing it for every single event. After a complete optimization process extensively described in Appendix C, the actual parameters used for this particular analysis can be found in Table 7.2.

After a trial and error process, the best results have actually been obtained with our ANN over the BDT, considering as input a relatively simple set of variables, all described in Section 7.1. By

Rank	Variable	Importance
1	$M_{T2}^{ll}$	$2.30 \cdot 10^{-1}$
2	$\Delta\Phi(E_T^{\text{miss}}, ll)$	$2.03 \cdot 10^{-1}$
3	pfMET	$1.90 \cdot 10^{-1}$
4	$m_{bl}^t$	$1.85 \cdot 10^{-1}$
5	$\cos(\theta_i) \cos(\theta_j)$	$1.18 \cdot 10^{-1}$
7	nbJet	$6.92 \cdot 10^{-2}$

Table 7.3: Ranking for the importance of the different variables used as input to the ANN, considering the scalar 100 GeV training.

order of importance, we therefore decided to use as input the stransverse mass  $M_{T2}^{ll}$ , the angle  $\Delta\Phi(E_T^{\text{miss}}, ll)$ , the pfMET,  $m_{bl}^t$ , one of the two spin correlated variables  $\xi = \cos(\theta_i) \cos(\theta_j)$  and the number of b-jets observed. The extensive explanation behind the choice of all these parameters can be found in Appendix C.

### 7.2.2 Training process

Given the kinematics of the signal events considered in this analysis, the SM  $t\bar{t}$  and the single top processes are the only processes considered to be part of the background for the training process. Since several mass points are inspected for each signal category, different ANNs have been trained for each mass point, and for each mediator category (scalar or pseudoscalar), to enhance the sensitivity of the analysis. As explained previously, given the limited number of training events available after applying the pre-selection, we decided to define a simple bi-class MVA, training at the same time only two categories of processes against each other: both the  $t/\bar{t}+DM$  and  $t\bar{t}+DM$  signals and the backgrounds, made out of a mix of both SM  $t\bar{t}$  and single top processes. Only the events passing the pre-selection of the analysis, as described in Section 6.2, are considered for the training.

The TMVA package, a toolkit for Multivariate Data Analysis with ROOT [161], was used in order to define both the BDT and the ANN previously defined and to perform the optimization of the different hyper-parameters of both methods. Such package is extremely helpful because it gives us plenty of information when training a specific model for a given mass point, such as a ranking of the importance of the input variables given in Tables 7.3 and 7.4, a plot representing the correlation between these input variables for the backgrounds and the signals in Figure 7.21, the Receiver Operating Characteristic (ROC) curve showing the background rejection achievable for a given selected signal efficiency in Figures 7.22 and 7.23, a plot to check for any sign of possible overtraining in Figure 7.24, and the input variable distributions in Figure 7.25. All these results have been obtained by splitting the dataset available as 50% for the training process and the remaining 50% to test the classifier obtained.

### 7.2.3 Evaluation process

Once the training performed and the weights of both MVA methods computed, we have been able to use them in order to estimate the probability of each data and MC event to be a background or a signal event. Once the weights applied, we have therefore been able to obtain a distribution of

Rank	Variable	Importance
1	$M_{T2}^{ll}$	$2.27 \cdot 10^{-1}$
2	$\Delta\Phi(E_T^{\text{miss}}, ll)$	$1.94 \cdot 10^{-1}$
2	pfMET	$1.92 \cdot 10^{-1}$
4	$m_{bl}^t$	$1.54 \cdot 10^{-1}$
5	$\cos(\theta_i) \cos(\theta_j)$	$1.53 \cdot 10^{-1}$
6	nbJet	$7.90 \cdot 10^{-2}$

Table 7.4: Ranking for the importance of the different variables used as input to the ANN, considering the scalar 500 GeV training.

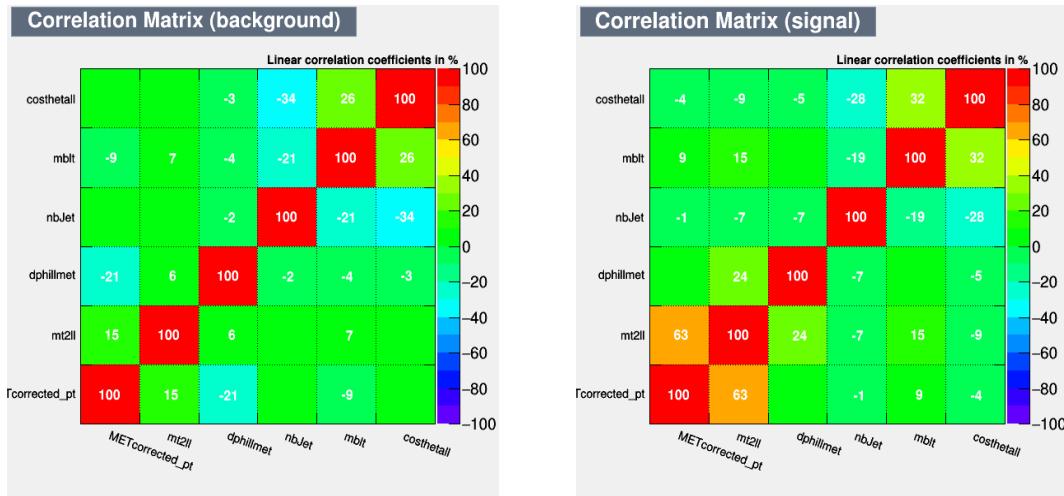


Figure 7.21: Correlation coefficient between the different variables used as input to the ANN for the backgrounds (on the left) and both signals (on the right).

the most probably category for all the processes considered in this analysis, as shown in Figure 7.26, and the actual ANN output distributions for each mediator mass point, used in a general shape analysis later on, as shown in Figure 7.27 to 7.28.

As we can see in these last plots, each bin correspond to one particular category, supposed to be enriched either in signals or backgrounds in the first and second bins respectively. We can also see that we still have a majority of background in each of the two bins, but this strategy does allow us to define separate regions enriched in the different processes anyway.

**FIXME:** add percentage table for each process? Did we end up using the control region of the MVA or simply

#### 7.2.4 Shape analysis

The output of the MVA method is used by a general shape analysis in order to extract the signals of interest. First of all, in order to perform such an analysis, the binning of the distribution was optimized in order to make sure that each bin of the distribution contains a significant number of MC events.

**FIXME:** complete this?

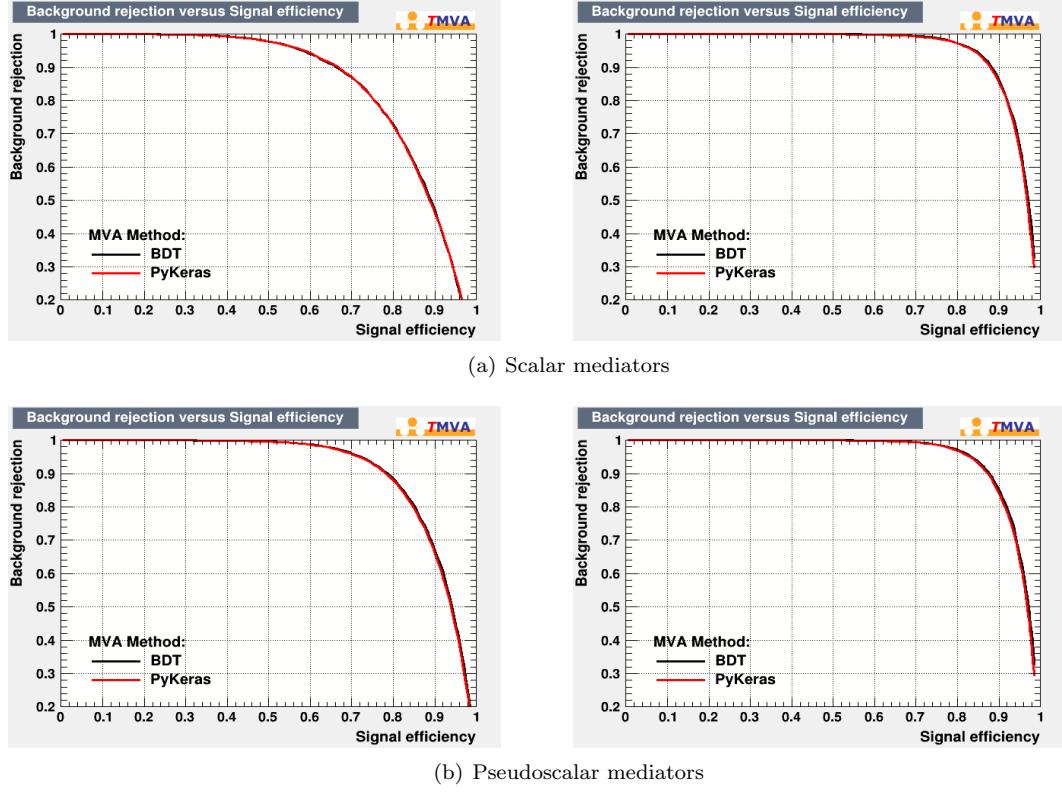


Figure 7.22: Signal versus background ROC curves obtained after the training performed, considering 100 (on the left) and 500 GeV (on the right) mediators.

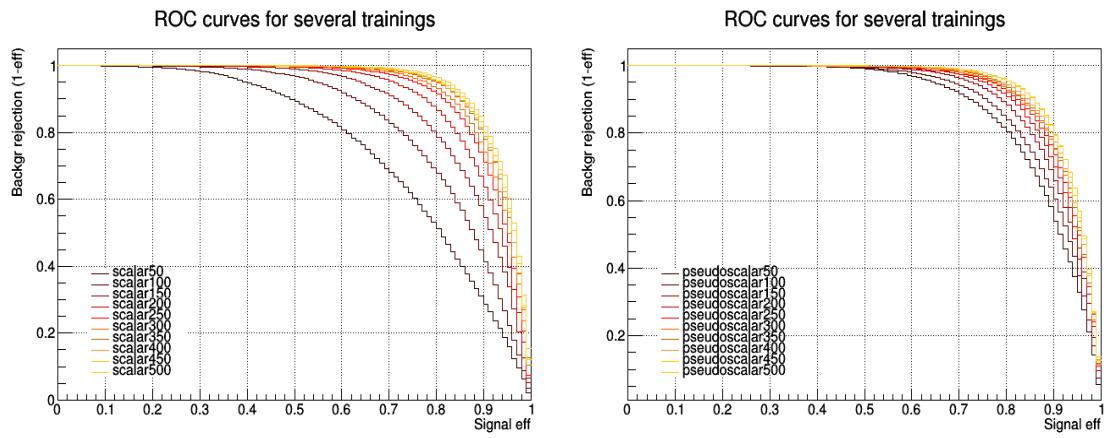


Figure 7.23: Signal versus background ROC curves obtained after the training performed with the ANN, considering all the scalar (on the left) and pseudoscalar (on the right) mediators available.

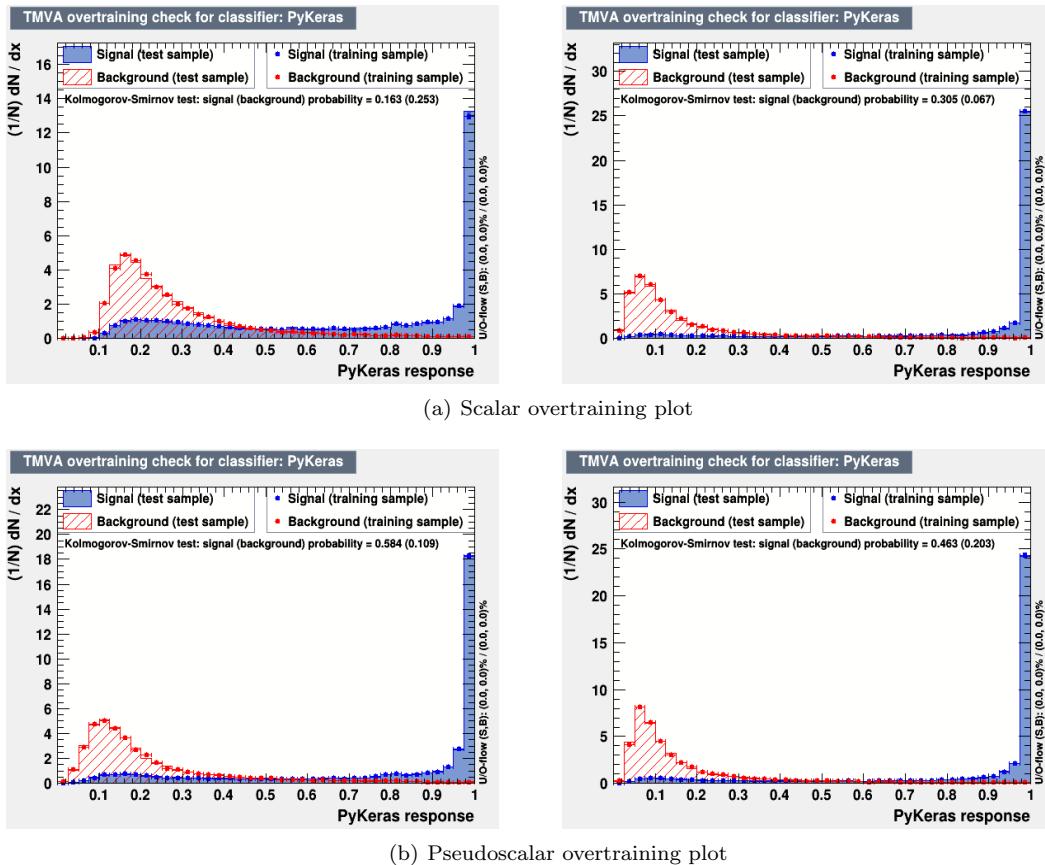


Figure 7.24: Overtraining distributions obtained for the 100 (on the left) and 500 GeV (on the right) mediators.

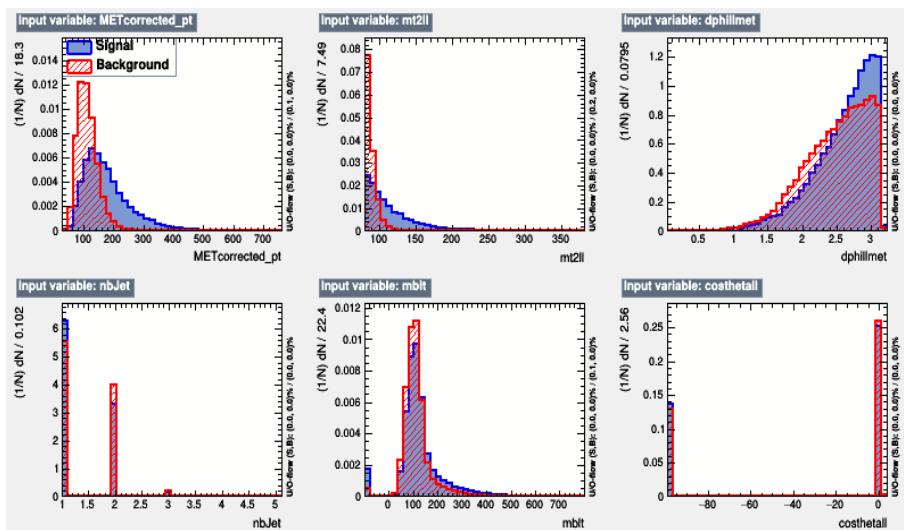


Figure 7.25: Input variables distributions used for the MVA.

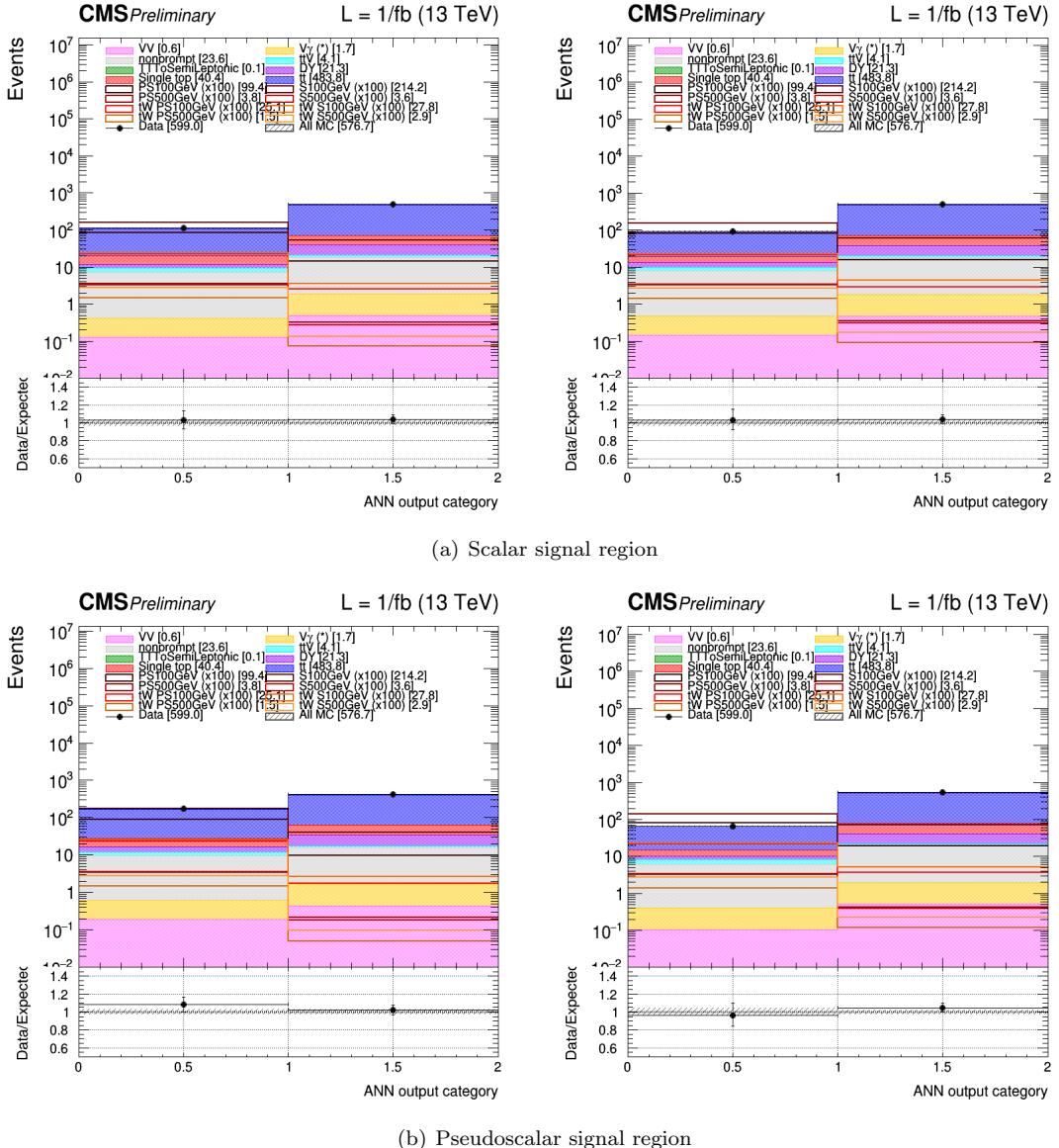


Figure 7.26: Most likely category of the events corresponding to the different processes for the 100 (on the left) and 500 GeV (on the right) mediators.

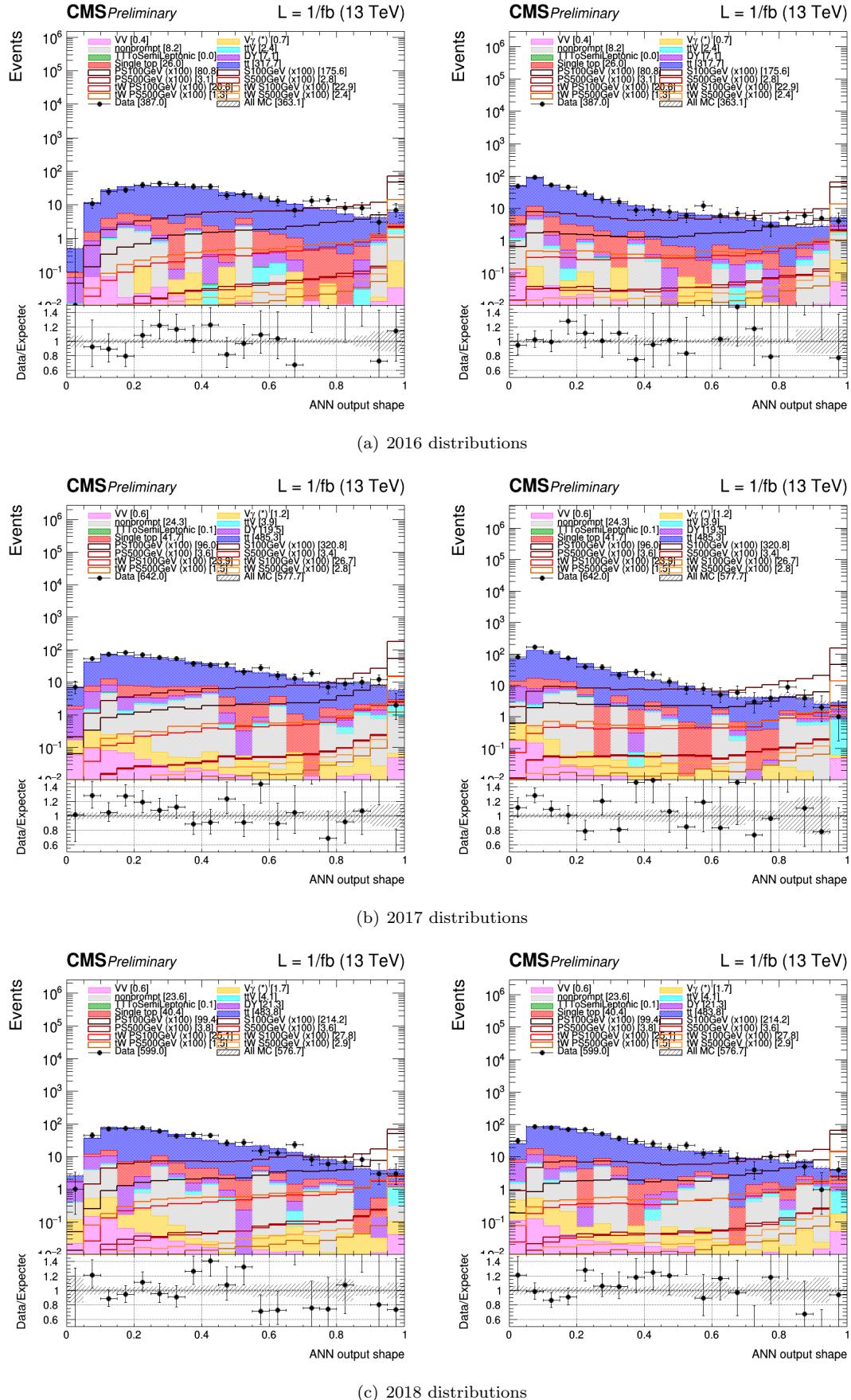


Figure 7.27: ANN output distribution used in the shape analysis performed later on, for the scalar 100 GeV (on the left) and 500 GeV (on the right) training.

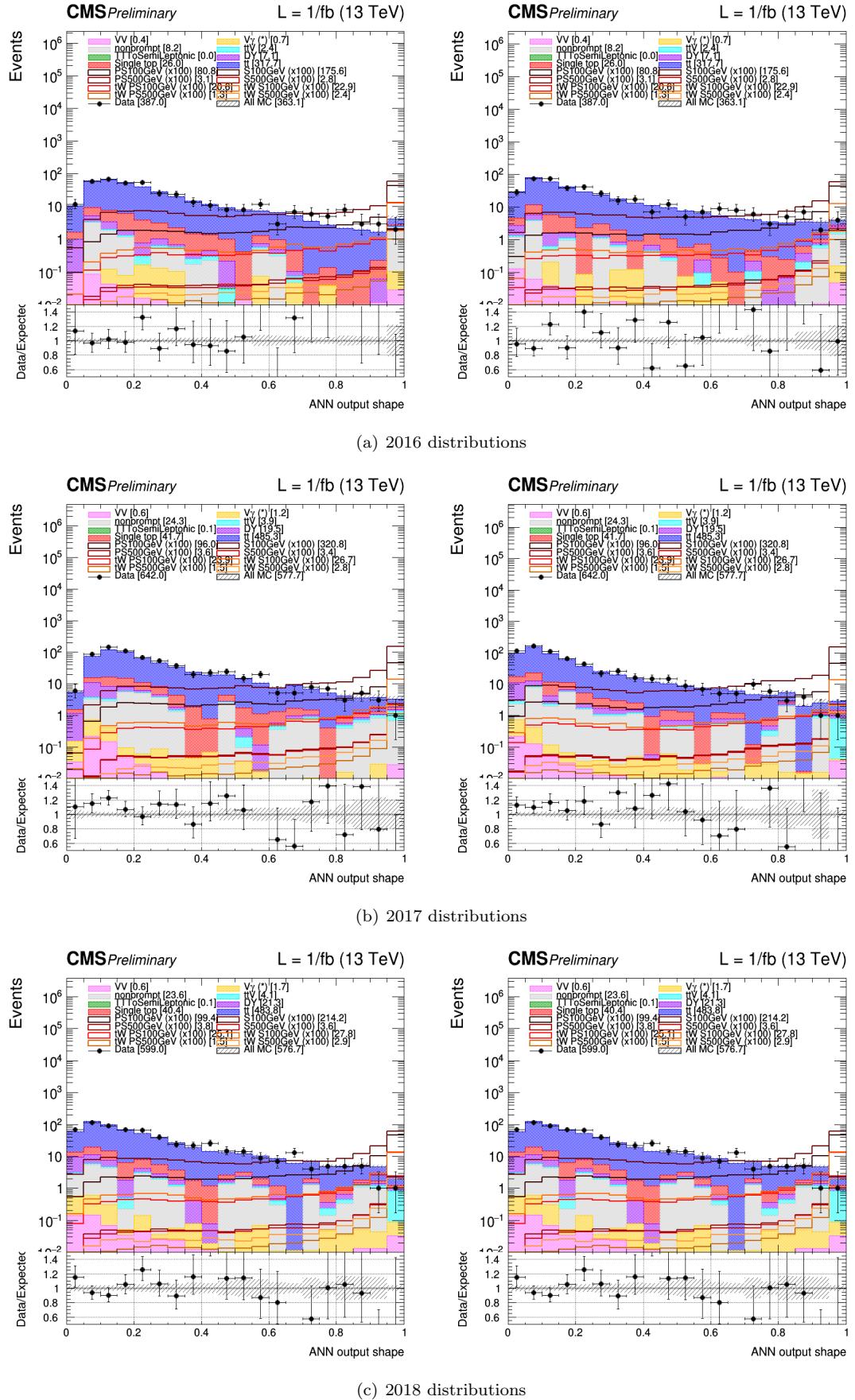


Figure 7.28: ANN output distribution used in the shape analysis performed later on, for the pseudoscalar 100 GeV (on the left) and 500 GeV (on the right) training.



---

---

# Chapter 8

---

## Results and interpretations

We now have all the ingredients needed to proceed with the interpretation of the data in terms of exclusion limits on the signal strength, for the different mediators spins, mass points and signal models considered. In this chapter, a short but necessary theoretical introduction to the statistical concepts used will be performed in Section 8.1, before discussing about the impact of the systematics and uncertainties on the final results in Section 8.2. Finally, the results and final limits obtained for this particular analysis will be shown in Section 8.3.

### 8.1 Statistical interpretation

The results presented in this chapter follow the recommendations and use the tools given directly by the LHC Higgs Combination Group [162], result of a collaboration between teams of ATLAS and CMS. The exclusion limits on the signal strength presented here have been calculated according to the asymptotic formulas described in this section, allowing us to quantify how sensitive the experiment is to the new physics searched for by calculating both the expected and observed significance obtained for a variety of signal hypotheses, corresponding to different mediator spins and mass points, as described previously.

Setting limits is typically performed using a frequentist statistical test, which rely on the definition of two hypotheses [164, 165]:

- The **null hypothesis**  $H_0$ , or  $b$ , according to which the signal searched for does not exist and does not contribute to the data observed in a given part of the phase space.
- And the **alternative hypothesis**  $H_1$  tested against  $H_0$  and stating exactly the opposite, mainly that our signal does actually exist. Since the typical search is not free of background, this hypothesis is usually called  $s + b$ .

In practice, we can usually not give a definite true or false answer to the question whether the alternative hypothesis is correct or not: we actually need to compute and quote a **CL** and a **confidence limit**, corresponding to the value of a population parameter excluded at a specified CL, for the exclusion. The objective is in this sense to quantify the level of agreement of the

observed data with a given hypothesis  $H$  by computing the so-called **p-values**, corresponding to the probability of finding data of equal or greater incompatibility with the predictions of  $H$ , under the assumption of  $H$ . This hypothesis  $H$  can then be disregarded if its p-value is observed below a specified threshold, typically set to 0.05, resulting in a CL common for the exclusion limits of most of the BSM searches performed at CERN.

In order to establish an eventual discovery or exclude certain models, this analysis then uses a quite typical frequentist significance test using a likelihood ratio, considering parameters of interest, such as the cross-section of the signal process and nuisance parameters, such as the background normalizations, as discussed in Section 8.2.

Mathematically, as shown in [163], this method can be understood by considering the  $x$  variable corresponding to any of the shape of the ANN output distributions previously obtained, containing each a certain expected number of events per bin  $E[n_i]$ , given by Equation 8.1, where  $b_i$  ( $s_i$ ) is the number of background (signal) events in the  $i$ -th bin and  $\mu$  is the so-called **signal strength**, a positive parameter characterizing the normalization of the signal under consideration and independent on the bin considered ( $\mu = 0$  therefore corresponds to the background-only hypothesis and  $\mu = 1$  is the nominal signal hypothesis).

$$E[n_i] = b_i + \mu s_i \quad (8.1)$$

The number of background and signal events in each bin can be expressed through Equation 8.2, where  $b_{\text{tot}}$  ( $s_{\text{tot}}$ ) correspond to the total number of background (signal) events, the function  $f_b(x; \boldsymbol{\theta}_b)$  ( $f_s(x; \boldsymbol{\theta}_s)$ ) are the PDFs of the variable  $x$  considered for signal and background events and  $\boldsymbol{\theta} = (\boldsymbol{\theta}_s, \boldsymbol{\theta}_b, b_{\text{tot}})$  represent all the nuisance parameters considered.

$$\begin{cases} b_i = b_{\text{tot}} \int_{\text{bin } i} f_b(x; \boldsymbol{\theta}_b) dx \\ s_i = s_{\text{tot}} \int_{\text{bin } i} f_s(x; \boldsymbol{\theta}_s) dx \end{cases} \quad (8.2)$$

To help constraining even more the nuisance parameters considered, control regions enriched in background where we expect to observe a low number of signal events if it were to exist in nature are usually defined. The expected number of events in each bin  $E[m_j]$  can then be expressed in Equation 8.3 in a similar way, where  $u_j$  is a calculable quantities depending on the set of nuisance parameters  $\boldsymbol{\theta}$ . The definition of such control region is particularly helpful when constraining the background normalization parameters introduced in Section 8.2.

$$E[m_j] = u_j(\boldsymbol{\theta}) \quad (8.3)$$

If we assume that the number of events in each bin is large enough, then we can define in Equation 8.4 the likelihood function as a product of Poisson probabilities for each of the  $N$  and  $M$  bins, by using the mean expected number of events ( $E[n_i]$  and  $E[m_j]$ ) and the actual number of events measured ( $n_i$  and  $m_j$ ) in the signal and control regions.

$$\mathcal{L}(\mu, \boldsymbol{\theta}) = \left( \prod_{i=1}^N \frac{E[n_i]^{n_i}}{n_i!} e^{-E[n_i]} \right) \left( \prod_{j=1}^M \frac{E[m_j]^{m_j}}{m_j!} e^{-E[m_j]} \right) \quad (8.4)$$

This last relation is helpful because it allows us to test any value of  $\mu$ , by considering the profile

likelihood ratio  $\lambda(\mu)$  defined in Equation 8.5, where  $\hat{\theta}_\mu$  is defined as the value of  $\theta$  that maximizes the value of the likelihood for a given  $\mu$ , while  $\hat{\mu}$  and  $\hat{\theta}$  are respectively the unconditional maximum-likelihood estimators of  $\mu$  and  $\theta$ .

$$\lambda(\mu) = \frac{\mathcal{L}(\mu, \hat{\theta}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})} \quad (8.5)$$

This equation is extremely important in the sense that it gives us a way to test the two hypotheses we have at our disposal against each other: the background-only hypothesis  $b$  (if  $\mu = 0$ ) and the signal+background hypothesis  $s + b$  (if  $\mu = 1$ ). Additionally, this equation allows us to determine which values of signal strength  $\mu$  can be discarded at a given confidence level, usually set to 95% within the CMS collaboration. At this stage, it is important to note as well that if a significant excess is observed, additional tests are required to be able to claim for a discovery, as this method only allows to reject some signal models not significant enough.

In order to finally be able to establish upper limits on the value of the parameter  $\mu$ , one final test-statistic  $q_\mu$  is usually defined according to Equation 8.6. Larger values for this parameter than represent higher incompatibilities between the data and the given signal strength. The reason why we set the estimator to 0 when  $\hat{\mu} > \mu$  is simply to make sure to have a one-sided confidence interval, because eventual upwards fluctuations of the data able to pass this condition should not be considered as evidence for a  $\mu$ -strong  $s + b$  hypothesis.

$$q_\mu = \begin{cases} -2 \ln \lambda(\mu) & \hat{\mu} \leq \mu \\ 0 & \hat{\mu} > \mu \end{cases} \quad (8.6)$$

Additionally, p-values can be used in order to quantify the level of agreement between the data and the  $\mu$ -strong signal under investigation, as shown in Equation 8.7 for an observed value  $q_{\mu,\text{obs}}$ , if  $f(q_\mu|\mu)$  is the PDF of the test-statistic  $q_\mu$  assuming a certain value of  $\mu$ .

$$p_\mu = \int_{q_{\mu,\text{obs}}}^{\infty} f(q_\mu|\mu) dq_\mu \quad (8.7)$$

The confidence level on the signal  $CL_s$  for a given signal strength is then defined in Equation 8.8, where  $p_\mu$  corresponds to the p-value obtained using the PDF of the test-statistic  $q_\mu$  assuming a value of  $\mu$ , while  $p_b$  is the PDF obtained for the same parameter but assuming  $\mu = 0$ . Most of these parameters introduced can be seen in Figure 8.1.

$$CL_s(\mu) = \frac{p_\mu}{1 - p_b} \quad (8.8)$$

This latest parameter  $CL_s$  can finally be tuned by modifying the value of  $\mu$  until reaching a given threshold  $\alpha$ , set to 0.05 in this thesis and in most of the publications of the CMS collaboration.

Performing this complete calculation with observed is giving us the so-called **observed upper limits**, but we also typically need to calculate the expected sensitivity of the analysis, defined as the mean value of  $\mu$  expected for a background-only hypothesis, computed from  $N$  simulated datasets. Placing all the  $\mu$  values obtained for a given threshold and a given dataset in a cumulative histogram then allows us to plot separate curves corresponding to both the observed and expected upper limits along with its corresponding error bars depending on the distribution of the  $\mu$  values obtained for all the  $N$  datasets. This actually allows us to make sure that the sensitivity of the

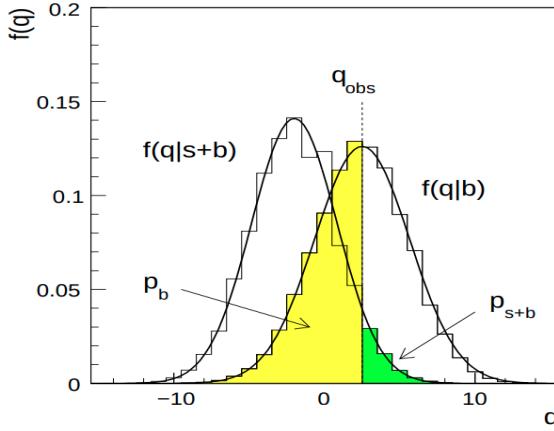


Figure 8.1: Schematic representation of the concept of PDF and p-value used in the computation of our test-statistic  $q_\mu$ , under the hypotheses  $\mu = 0$  and  $\mu = 1$ .

search is good enough to be able to effectively see a signal if it exists. Getting an expected value of  $\mu < 1$  for a given mass point is then typically the sign that if the signal did in fact exist, the analysis should be sensitive enough to see it and, if we don't see any significant deviation between the expected and observed curves, we can then rule out the possibility of such model actually existing.

Finally, one last important concept can be discussed here. Two different kind of signals are searched for in this work at once, mainly the  $t/\bar{t}+DM$  and the  $t\bar{t}+DM$  processes. Even though we can search for both processes independently, they are both able to exclude dark matter production models and a combination of both searches is therefore expected to give better results than two separate searches. Some of the nuisance parameters are common between the channels, but not all of them and, if we take the same initial signal strength value for both channels and if the channels are statistically independent, which is a condition verified in this case, then the full likelihood function is simply given by the product of the individual likelihoods in each channel  $i$ , allowing us to define a common test-statistic parameter shown in Equation 8.9.

$$\lambda(\mu) = \frac{\prod_i \mathcal{L}_i(\mu, \hat{\boldsymbol{\theta}}_{\mu,i})}{\prod_i \mathcal{L}_i(\hat{\mu}, \hat{\boldsymbol{\theta}}_i)} = \prod_i \lambda_i(\mu) \quad (8.9)$$

Because of this, it is possible to determine the values of the profile likelihood ratio separately for each channel, which simplifies greatly the task of estimating the median significance that would result from the full combination. The same considerations can also apply when combining for example the different final states given by the different decays of the  $W$  bosons, or when combining the data collected during the different years of operation of the LHC during the Run II.

## 8.2 Systematics and uncertainties

All the measurements we can make in high energy physics and in physics in general present some uncertainties, and the determination of their value is a critical point of every analysis, especially in this case since we try to detect a really low signal over a large background. Indeed, as we already saw, the nuisance parameters  $\boldsymbol{\theta}$  are extremely important when defining the upper limits on the signal strength, so a good characterization of such source of uncertainties is crucial in most

of the analyses searching for new physics. The uncertainties considered in this work belong to two major categories:

- On one hand, **statistical uncertainties** appear in any counting experiment [166] since every given measurement, made by definition of a finite set of observations, typically results in different observations when repeated, and the statistical uncertainty is then just a measure giving us an idea about the range of this kind of variations. Given the high rate of collisions happening at the LHC, we usually assume that our counting experiment can be approximated with a Poisson distribution, for which we know that the error on the number of measurements is directly given by its square root.
- On the other hand, **systematic uncertainties** are just as important, but are usually harder to estimate and are different in nature [167] since they arise directly from the theory or from the detector itself. Some systematic uncertainties are treated as individual nuisance parameters when fitting the MC observed to the data, nuisances which can either only affect the normalization of some processes, or affect the shape of the predictions across the distribution of the observables. These uncertainties share the need to be propagated through the full analysis chain all the way to the discriminating distributions.

Even though the signal enriched bins are to be found in a region of the MVA discriminant with lower number of events, where the impact of statistical uncertainties is higher, systematic uncertainties still play an important role in most of the analyses performed at CERN, including this one.

### Statistical uncertainties

Shape uncertainties due to the limited size of the simulated signal and background samples are simply included by allowing each bin of the distributions included in the signal extraction to fluctuate independently according to the statistical uncertainty on the simulation.

### Theoretical uncertainties

This category of systematic uncertainties is related directly to the theoretical models used in the analysis and is mostly related to the production of the MC simulations according to the process described in Section 5.1. Several different uncertainties belong to this particular category:

- **PDF and higher order corrections.** As accurate as Feynman diagrams can be, they do not represent the complete picture since we are computationally limited and therefore cannot consider high order perturbations, MC samples being usually limited to the Leading Order (LO) or Next to Leading Order (NLO). This limitation introduce a small systematic uncertainty that usually needs to be taken into account. Additionally, uncertainties due to the choice of the PDFs and  $\alpha_s$  parameters themselves are usually estimated by considering a hundred NNPDF3.0 [141] replicas, according to the PDF4LHC recommendations for the LHC Run II [168]. The uncertainty obtained in a given bin is then set as the standard deviation of the content of the bin obtained from all these computed replicas.
- **Underlying event and parton shower modeling.** The choice of the parton shower generator and UE tune is also to have an effect on the final results, but since we did not observe any dependency of UE variations on the number of jets, a flat 1.5% UE uncertainty is assigned to cover all the up and down variations.

- **Renormalization and factorization scales.** Uncertainties usually emerge due to the presence of the renormalization  $\mu_R$  and factorization  $\mu_F$  scales in the QCD calculations regarding the hard collisions of hadrons happening at the LHC, which are updated by a factor 0.5 and 2 and then propagated to the distributions of the analysis. This uncertainty is considered to be uncorrelated among the different background processes considered.

## Experimental uncertainties

Additional uncertainties related to the detector or to the experimental conditions are also taken into account when producing the upper limits on the signal strength.

- **Luminosity.** The actual integrated luminosity collected is not a fixed number and comes with an associated uncertainty of 2.5%, 2.3% and 2% for 2016, 2017 and 2018 respectively, based on van der Meer scans performed [138, 139, 140].
- **Pileup modeling.** Varying the inelastic cross-section used to calculate the pileup distribution in simulation results in a systematic of the order of 5% [169].
- **Lepton trigger.** An uncertainty associated to the lepton trigger efficiency is taken into account. It has an impact of the order of  $\sim 2\%$ , estimated from a general Tag and Probe method, by varying the Z window considered and the tag lepton  $p_T$  selection cut.
- **Lepton efficiency and energy scale.** Scale factors are applied to the MC processes in order to mimic the measured lepton reconstruction and selection efficiencies in data, and such scale factors come with estimated  $p_T$  and  $\eta$  uncertainties of the order of  $\sim 2\%$  for electrons and muons that have been applied.
  - In particular, for **electrons**, efficiencies are computed using a Tag and Probe method as well, involving background and signal fits, for which different fitting functions have been considered for the computation of the systematics. The tag selection has also been updated, mainly with different  $p_T$  cuts, studying the impact such change has on the final efficiencies. Uncertainties from such different sources have been added in quadrature to get the final systematic associated to electrons.
  - Then, four different parameters have been considered and updated with respect to the nominal Tag and Probe parameters to study their impact on the final **muon** efficiencies: the isolation of the tag muon, the signal fitting function used, the number of mass bins where the fit is done and finally, the Z-window considered for the calculation.
- **Jet Energy Scale (JES).** Single nuisance parameter applied to the  $p_T$  of each jets, as a function of the jet  $p_T$  and  $\eta$  values. These variations are then coherently propagated to important discriminating variables, such as the  $E_T^{\text{miss}}$  and  $M_{T2}^l$ , according to a procedure described by the JET/MET POG [153].
- **MET mismodeling.** The uncertainties on jets, electrons and muons that make up the MET are applied to the respective objects, while the MET is recalculated with the same up and down variations alongside. Additionally, any eventual unclustered energy, mostly due to ECAL or HCAL deposits not assigned to any object is shifted up and down, and the impact of such variations has been estimated on the MET.
- **b-tagging efficiency.** The b-tagging process is not perfect and its efficiency, typically different in data and MC also needs to be taken into account by applying a  $p_T$ ,  $\eta$  and flavour dependent scale factor on the MC. The uncertainties on such scale factors are measured in an independent control sample and propagated to the analysis [170].

- **Top quark  $p_T$  reweighting.** As described in Section 5.6.6, the measured  $p_T$  spectrum is expected to be softer than the one obtained using simulation, so a factor correcting this effect is usually applied to the MC. A systematic covering for such differences observed between data and MC is taken into account.

### Background related uncertainties

Finally, the systematic associated to the backgrounds normalization is typically one of the largest source of uncertainties in an analysis, and is computed by estimating the normalization of the backgrounds that are estimated on data control samples whenever possible.

- **Drell-Yan (DY) background.** To cover for the non-flatness of the  $R_{\text{in-out}}$  transfer factor and possible mismodeling of this particular background, especially at high pfMET values, a 20% systematic is associated to this process.
- **Non-prompt background.** A normalization and shape uncertainty is associated to the estimation of this background in a specific control region according to the tight-to-loose method previously described. The combined uncertainty due to statistics and the impact of a change in the input jet  $E_T$  threshold chosen has been estimated to be less than  $\sim 15\%$ , while an additional flat 30% systematic is considered for this background to cover for eventual discrepancies between data and MC in the same-sign control region defined in Section 6.3.5.

The impact plots, showing the impact the each systematic has on the final upper limits obtained, are shown in Appendix D.

**FIXME:** check all systematics, which ones we actually applied and add impact plots

## 8.3 Results

Expected and observed upper limits on the different signals production cross-section at the 95% confidence level have been obtained and plotted against all the possible dark matter mediator masses considered in this actual analysis, considering both the scalar and pseudoscalar mediators, as shown in Figure 8.2, for  $m_\chi = 1 \text{ GeV}$  and when setting all the couplings of the different models to 1. In both cases, limits obtained separately for the  $t/\bar{t}+\text{DM}$  and  $t\bar{t}+\text{DM}$  have been obtained, along with a combination of the limits obtained after combining these two signals of interest. Finally, Figure 8.3 show the limits obtained for each data taking period in the same plot, along with a combination of all the different exclusion limits, giving us the final exclusion limits for this analysis, for both the scalar and pseudoscalar mediators.

By looking at all these plots, we can make several different observations and conclusions:

- First, these plots allow us to see that the **inclusion of the  $t/\bar{t}+\text{DM}$  signal is extremely important**, especially at high scalar mediator masses, where the limits obtained are even comparable to the ones obtained by considering the  $t\bar{t}+\text{DM}$  signal.
- We can also observe that the limits obtained in 2016 are slightly worse than the other two data taking periods. This is an expected effect as expected limits are roughly expected to decrease with the square root of the luminosity, slightly lower in 2016.

- Finally, we can observe that low masses mediators are easier to exclude than higher masses simply because, even though they feature a worse global signal/background discrimination, their cross-section is also much higher. In this sense, pseudoscalar exclusion limits are also typically slightly worse than the scalar ones, mainly because of the difference in cross-section between such processes.

At the end of the day, this analysis allowed us to achieve an expected (observed) exclusion for both scalar and pseudoscalar mediators up to 200 (XXX) GeV, more than doubling the previous results obtained in 2016 for the  $t\bar{t}$ +DM model alone.

**FIXME:** check conclusions, talk about observed limits

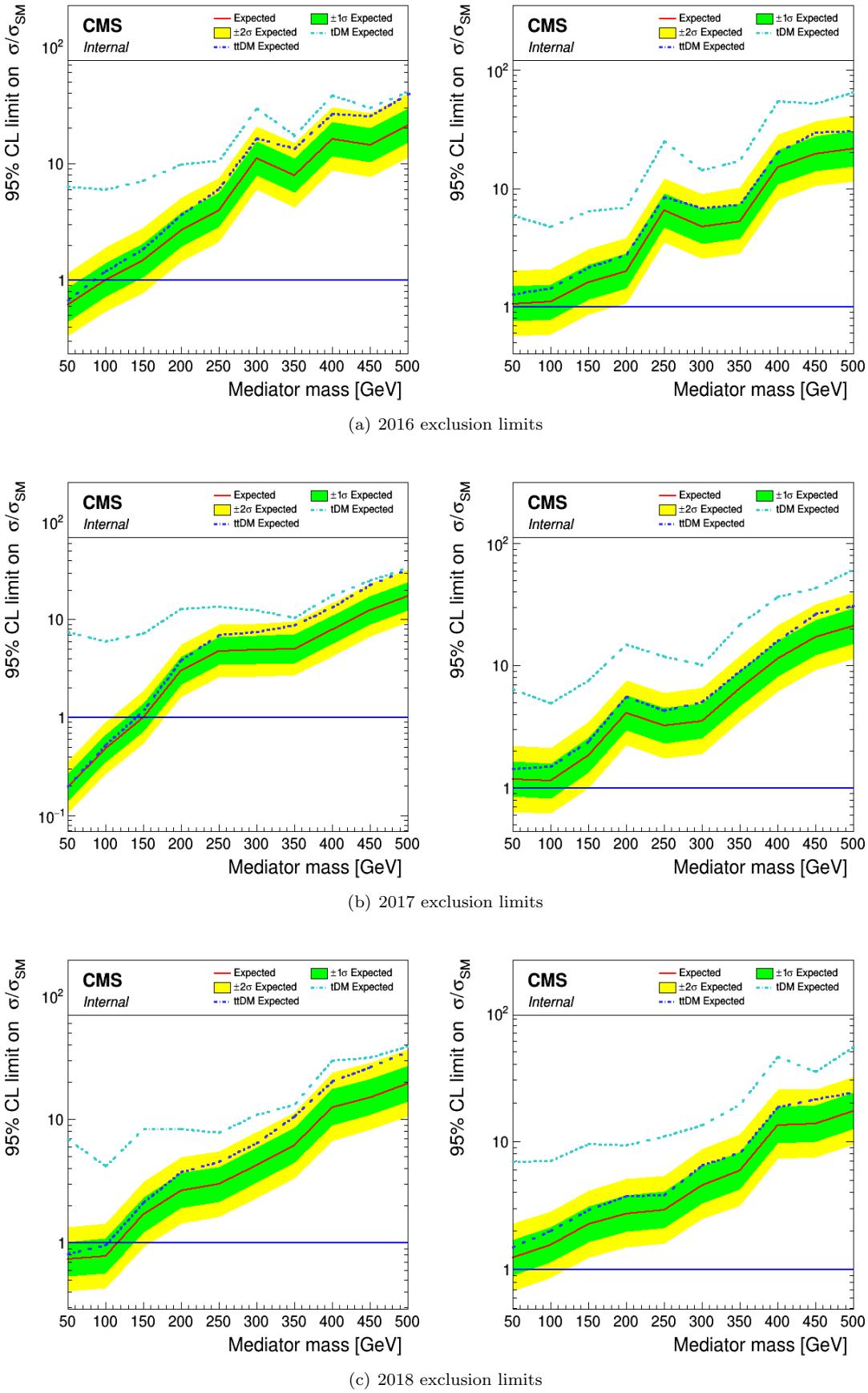


Figure 8.2: Expected and observed upper limits on the signal strength  $\mu$  for scalar (on the left) and pseudoscalar (on the right) models, considering with  $m_\chi = 1 \text{ GeV}$  and  $g_\chi = g_q = 1$ .

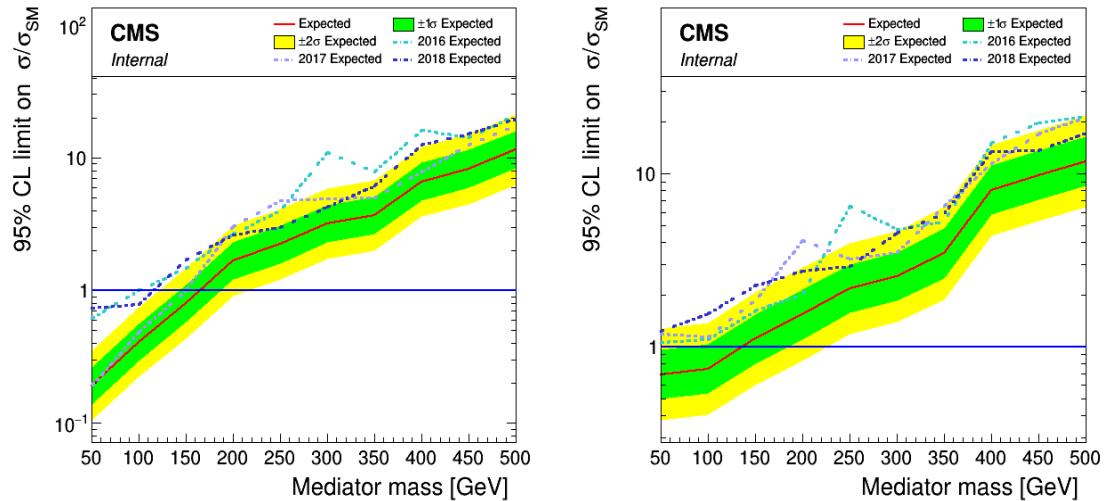


Figure 8.3: Expected and observed upper limits on the signal strength  $\mu$  for scalar (on the left) and pseudoscalar (on the right) models, considering with  $m_\chi = 1$  GeV and  $g_\chi = g_q = 1$ , after combining the different Run II data taking periods.

---

---

# Chapter 9

---

## Conclusions

In conclusion, a search for the production of dark matter in association with one or two top quarks has been performed in this work, by studying in particular the dilepton decay channel of both production modes. This analysis was done by considering the full  $(137.1 \pm 2.0) \text{ fb}^{-1}$  of proton-proton collisions data collected by the CMS detector during the Run II of operation of the LHC, at a center of mass energy of  $\sqrt{s} = 13 \text{ TeV}$ . No evidence for the existence of dark matter has been found, but upper limits on the signal strength have been obtained by considering different production models and channels. Several steps were needed to reach this goal:

- First, the triggers were chosen in such a way to collect as much interesting data as possible. Different objects were also defined and chosen at this point in coordination with other groups, such as the working point used for leptons and for the b-tag. The choice regarding the MET used in this particular analysis was a bit trickier, given the impact this choice might have on the final results, but we finally settled on the usual pf typeI correct MET, applying all the corrections recommended to this particular variable.
- Once this was done, a full top reconstruction method was developed in order to reconstruct in the best way possible the kinematic of the  $t\bar{t} \rightarrow 2l$  process. This method allowed us for example to get the information related to each top quark separately, and to define new discriminating variables, such as an estimation of the mediator  $p_T$ .
- All the signal samples were then produced, first privately and then centrally, considering all the different models that might be of interest for this analysis.
- The background were then carefully studied. Even though most of the backgrounds are estimated directly from MC, some of them did require a bit extra work, such as the DY process, for which a semi data-driven method was used for its estimation. The predictions made for most of the major backgrounds of the analysis were then checked in dedicated control regions.
- Several discriminating variables available to us were then explored, and a MVA method was set up in order to combine the discriminating power of all these variables into a single output, used to perform a general shape analysis. The optimization of all the hyper-parameters was thoroughly done, as we ended up choosing an ANN to perform the actual analysis, with a simple set of 6 different input variables.

- Finally, exclusion limits on the signal strength at the 95% confidence level were set for the different models considered.

At the end of the day, this analysis allowed us to achieve an expected (observed) exclusion for both scalar and pseudoscalar mediators up to 200 (XXX) GeV, improving by a factor of more than 2 the previous results obtained in 2016 for the  $t\bar{t}$ +DM model alone.

**FIXME:** Check this once analysis done

## 9.1 Future prospects

This is the first time that such a search combining the  $t/\bar{t}$  and  $t\bar{t}$ +DM models is performed in the dilepton final state. Even though large parts of the phase space have already been excluded, several additions could be considered to improve and complete the results obtained by this analysis:

- The remaining models to be excluded correspond to high mediator masses, which feature a large signal/background discrimination on their own but which also have a low cross-section of production. The continuous operation of the LHC is therefore expected to give more data to analyze, in turn decreasing the limits on the signal strength roughly as  $\sqrt{\frac{1}{\mathcal{L}}}$ . This channel will therefore benefit from the larger luminosity collected during the full Run III of the LHC data taking period, expected to start soon in 2021.
- A single set of couplings was studied in this work, but exploring other accessible regions of the space (by considering different  $m_\chi$  values, even though they are expected to have a minimal impact on the kinematics of the event, and mostly by changing the value of the couplings considered) would in this sense be an interesting addition to this particular analysis.
- One possible addition to this work would be to include a ttH to invisible reinterpretation of this search, given the relative similarities between this kind of model and the dark matter models considered here.

This analysis is in any case expected to gain momentum and provide us with even better exclusion limits over the course of the next years of operation of the LHC.

# Appendices



---

---

# Appendix A

---

## Resumen en español

El Modelo Estándar de la física de partículas [1] es hoy en día el modelo matemático más aceptado para describir las partículas fundamentales y algunas de las diferentes interacciones que existen entre ellas. Basado en ideas simples, este modelo ha permitido explicar la gran mayoría de los fenómenos observados en la naturaleza con un grado de precisión excelente, y ha sido capaz a lo largo de los años de hacer predicciones sobre la existencia de nuevas partículas, tal y como el postulado del mecanismo de Brout-Englert-Higgs [2, 3], seguido por el descubrimiento del bosón de Higgs mismo en el año 2012 [4, 5] por los experimentos CMS [6] y ATLAS [7] del CERN, analizando colisiones entre protones producidas por el LHC a una energía en el centro de masa  $\sqrt{s} = 7 \text{ TeV}$  y  $8 \text{ TeV}$ .

A pesar de estas predicciones, pensamos hoy en día que el Modelo Estándar no es completo y tiene unos pocos defectos ya que no es capaz de explicar algunas observaciones cosmológicas hechas a lo largo del siglo XX. Eventuales partículas exóticas que no caben dentro de este modelo podrían ser el signo de la existencia de física más allá de este Modelo Estándar. Este trabajo consiste por lo tanto en tratar de mejorar nuestro entendimiento del Universo en el cuál vivimos buscando a partículas de materia oscura que podrían ser producidas con las colisiones entre protones del Large Hadron Collider (LHC) del CERN.

### A.1 Materia oscura

La primera hipótesis seria sobre la posible existencia de la materia oscura se desarrolló en los años 1970, después de que varios astrofísicos hayan observado anomalías gravitacionales, siendo una manera sencilla de explicar la aparente falta de materia luminosa en el Universo [8]. En efecto, la masa visible dentro de las galaxias parecía ser muy baja para explicar algunos fenómenos observados, tal y como las curvas de rotación de las galaxias [9], que parecen ser incompatibles con las leyes de la gravitación de Newton. Algunas medidas adicionales del efecto lente (en el Bullet Cluster, por ejemplo [10]) y las anisotropías observadas en el fondo cósmico de microondas [11] son otras evidencias para la posible existencia de materia oscura.

Estas observaciones cosmológicas han permitido determinar que la materia bariónica ordinaria solo constituye más o menos el 5% del Universo visible, mientras la materia oscura cuenta para el 26% de

la densidad de energía total del Universo (la energía oscura, otro fenómeno completamente distinto constituye la parte faltante de esta densidad de energía). Estudiar la naturaleza y las propiedades de este nuevo tipo de partículas es por lo tanto muy importante para entender las leyes físicas que rigen nuestro Universo, con muchos científicos y experimentos alrededor del mundo dedicados a este tipo de búsquedas.

La existencia de la materia oscura está muy bien motivada, pero nunca hemos sido capaces de observar este tipo de partículas hasta la fecha y la única evidencia sobre su posible existencia viene de sus efectos gravitacionales observables a larga distancia. Se desconocen la masa, espín y propiedades básicas de este tipo de materia, pero la mayoría de las teorías contemplan un tipo particular de candidatos para constituir la materia oscura: las Weakly Interactive Massive Particles (WIMPs), ya que este tipo de partículas cumplen con propiedades generalmente asociadas a esta idea de materia oscura. En efecto, sabemos que el candidato ideal para formar esta materia exótica se debe de ser en otros insensible a la radiación electromagnética, no-bariónico, frío, estable a largo plazo y en un rango de masa entre 10 GeV y 1 TeV.

Diferentes métodos se pueden usar para buscar a este nuevo tipo de partículas, perteneciendo a tres categorías principales:

- Las **búsquedas directas**, buscando posibles interacciones entre materia oscura y partículas del Modelo Estándar, estudiando por ejemplo choques y por lo tanto cambio de energía entre estos dos sectores.
- Las **búsquedas indirectas**, tratando de encontrar a nuevas partículas del Modelo Estándar que podrían emerger de la interacción entre dos partículas de materia oscura.
- Y la **producción dentro de colisionadores de partículas**, el método de investigación principal de este trabajo, ya que se supone que chocar dos partículas del Modelo Estándar con una determinada energía podría llegar a producir un par de partículas de materia oscura. Este tipo de búsquedas resulta ser muy interesante ya que permite estudiar candidatos a materia oscura de muy baja masa.

Chocar partículas fundamentales se puede hacer hoy en día en aceleradores de partículas como el LHC, pero estudiar la posible producción de materia oscura resulta ser muy complicado por la naturaleza misma de este tipo de materia, ya que apenas interacciona con el detector dispuesto alrededor del punto de colisión. Esto significa que, incluso si logramos producirlas, se espera que estas nuevas partículas escapen el detector sin ser detectadas y este tipo de búsquedas tiene por lo tanto que buscar materia oscura como energía faltante producida en asociación con partículas del Modelo Estándar que se pueden detectar.

En este trabajo en particular, se busca materia oscura producida en asociación con uno o dos quarks top con el detector CMS a una energía en el centro de masa  $\sqrt{s} = 13$  TeV, considerando los  $\sim 137 \text{ fb}^{-1}$  de datos colectados durante el RunII de operación del LHC, a lo largo de los años 2016, 2017 y 2018. Este canal de investigación resulta ser muy interesante ya que se espera que el acoplamiento entre los sectores bariónico y oscuro sea de tipo Yukawa, y que sea por lo tanto mayor para partículas que tengan una masa mayor, tal y como el quark top. Estos quarks tienen sin embargo un inconveniente en el sentido de que decaen rápidamente, antes de llegar al detector, y solo podemos observar el resultado de este decaimiento, siendo en general formado por unos leptones y unos jets, resultado del proceso de hadronización de los quarks producidos. El número de leptones que esperamos observar depende del decaimiento del bosón W que aparece en la cadena de decaimiento del quark top. En este trabajo, nos concentramos en el decaimiento dileptónico de este bosón, que presenta la ventaja de ser un canal de desintegración bastante limpio, con pocos fondos pero con el inconveniente de tener un branching ratio bastante pequeño.

## A.2 El dispositivo experimental

### El LHC

El LHC es un colisionador de partículas ubicado 100 metros por debajo tierra en la frontera entre Francia y Suiza, resultando de la colaboració entre miles de institutos y científicos del mundo. Capaz de acelerar protones o iones de plomo a una velocidad muy cerca a la velocidad de la luz, este dispositivo, circular y de un tamaño de 27 kilómetros, empezó a funcionar en el año 2010 a una energía en el centro de masa de 7 TeV. A lo largo de los años, esta energía subió hasta llegar a los 13 TeV al empezar el Run II de funcionamiento en el 2016, después de una parada técnica de unos años. El LHC es una máquina ideal para estudiar física más allá del Modelo Estándar por los niveles de energías alcanzables y por su alto número de colisiones producidos por segundo (su luminosidad instantánea), lo que permite estudiar procesos con baja tasa de producción.

### El detector CMS

En el anillo principal de la cadena de aceleradores formando el LHC, 4 detectores han sido construidos alrededor del *beam pipe* para estudiar las partículas producidas después de las colisiones entre los protones. Estos detectores, CMS, ATLAS, ALICE y LHCb, han sido desarrollados usando equipos y métodos distintos, para lograr varios objetivos.

En particular, en este trabajo se estudian los datos colectados por el detector CMS, un detector de forma cilíndrica, de propósito general y diseñado para ser lo más hermético y compacto posible. Este detector pesa unos 12 500 toneladas, mide unos 14 metros de diámetro y unos 22 metros de altura, y está compuesto de diferentes capas, cada una diseñada con un objetivo distinto:

- Primero, en el interior del detector y muy cerca del punto de colisión, un sistema perfeccionado de detección de trazas cargadas conocido como el *tracker* ha sido instalado, para identificar el punto de origen de la colisión y ayudar con la identificación y la medida del momento de las partículas creadas. Este detector, compuesto por unas capas de píxeles de silicio rodeada por un detector de tiras de silicio de 10 capas y por sus tapas correspondientes para que esta capa sea la más hermética posible, tiene que ser muy resistente a la radiación y rápido para medir todas las colisiones, que se producen cada 25 nanosegundos.
- Luego viene los dos calorímetros: el ECAL y el HCAL. Estos dispositivos permiten medir la energía de las partículas interaccionando de manera electromagnéticas o hadrónicas respectivamente de manera muy eficiente. El ECAL está hecho de una parte central compuesta por 61 200 cristales de tungsteno de plomo PbWO<sub>4</sub> y completada por dos barriles de 7 324 cristales adicionales. El HCAL por el otro lado está hecho por unas capas de centelladores y de material que permite desarrollar cascadas hadrónicas, lo que permite medir con precisión la energía de los eventuales hadrones producidos.
- La parte central del detector es un imán capaz de crear un campo magnético de 3.8T, lo que permite determinar el momento de las partículas gracias a la relación de Lorentz.
- Fuera de esta parte central se encuentra el sistema de detección de muones, ya que este tipo de partículas interacciona poco con la materia y puede escapar del detector. Este último sistema de detección se basa en varias tecnologías: su zona central, donde el campo magnético es uniforme, viene por ejemplo formada por las cámaras de derivas (las DTs). Estas cámaras están divididas en 4 capas y están hechas por más de 172 000 cables que permiten medir la posición exacta de los muones gracias al fenómeno de ionización. Las tapas de este sistema,

en la zona donde el campo magnético no es uniforme y donde la tasa de muones y de fondos es más alta, las cámaras de tiras catódicas CSCs están instaladas y, por fin, el sistema está completado por las cámaras de tiras resistivas RPCs, mucho más rápidas que las CSCs.

Debido a limitaciones computacionales, es imposible guardar todas las colisiones que se producen en el LHC y, de todas maneras, la gran mayoría de las colisiones son muy bien conocidas y poco interesantes. Por lo tanto, un sistema de *trigger* de dos niveles ha sido desarrollado a lo largo de los años para guardar solamente unos 1000 eventos por segundo que parecen ser interesantes desde el punto de vista de la física para un análisis ulterior.

### A.3 Reconstrucción de objetos

Como lo acabamos de ver, el detector CMS está hecho por varias capas dedicadas a las medidas de diferentes propiedades de las miles de partículas que se generan después de cada colisión entre dos protones en el *beam pipe*. Cada capa del detector atravesada por una de estas partículas genera señales eléctricas que se colectan y que se almacenan. Esta señal solamente contiene datos brutos poco interesantes al principio, así que un algoritmo conocido como el *Particle Flow* ha sido desarrollado para analizar las señales que provienen de cada parte del detector, para combinarlas y extraer información física sobre la colisión estudiada.

Para la búsqueda de materia oscura llevada a cabo en este trabajo, se seleccionan eventos que tienen dos leptones y dos jets y, por lo tanto, un sistema eficiente de detección, reconstrucción e identificación de este tipo de objetos es imprescindible. Los muones se reconstruyen primero con este algoritmo, ya que son las únicas partículas cargadas capaces de llegar a las cámaras de muones en el exterior del detector. Luego, se pueden reconstruir los electrones, indentificados como una cascada electromagnética en el ECAL asociada con trazas compatibles en el tracker y, por fin, se identifican los hadrones cargados y neutros asociando los hits que quedan en el tracker con la información del HCAL. Los eventuales quarks que aparecen en una colisión se hadronizan y se manifiestan como jets de partículas en el detector, más difíciles de reconstruir en general. Se desarrolló para llevar a cabo esta tarea el algoritmo *anti- $kT$* , capaz de agrupar los hadrones, fotones y leptones que forman un solo jet. Este algoritmo es también capaz de distinguir un jet procedente de un bottom quark, ya que estos quarks tienen un tiempo de vida un poco más alto y que pueden por lo tanto viajar unos poco milímetros antes de decaer. Este pequeño desplazamiento basta para distinguir estos jets, llamados b-jets, que tienen un papel importante en este análisis.

El *Particle Flow* es también capaz de reconstruir un objeto llamado MET y que representa la energía transversal faltante de una colisión, siendo definido como la suma del momento transverso de todas las partículas creadas y medidas por el detector. Como sabemos que estas colisiones se producen de manera frontal, el momento transverso total de una colisión es exactamente 0 al principio y esperamos que se quede así después del choque por el fenómeno de conservación del momento. Sin embargo, no siempre es el caso y a veces se puede observar un momento transverso total no nulo más que nada por imperfecciones en el detector, o debido a la producción de partículas capaces de escapar el detector sin dejar ninguna señal, tal y como los neutrinos del Modelo Estándar o eventuales partículas exóticas. Una reconstrucción completa y precisa de esta variable es por lo tanto muy importante en este análisis buscando partículas de materia oscura.

Por fin, hemos desarrollado para esta análisis en particular un método completo de reconstrucción *offline* de los sistemas  $t\bar{t}$ . Este algoritmo es muy importante ya que se espera que muchas de las variables que nos permiten distinguir entre la señal buscada y los fondos dependen del 4-momento de los quarks top, típicamente no disponible sin este proceso de reconstrucción.

## A.4 Análisis de datos

El objetivo principal de este tipo de búsqueda de física nueva consiste principalmente en observar los datos colectados por el detector CMS en este caso, y compararlos con simulaciones matemáticas de tipo Monte-Carlo que permiten simular la naturaleza aleatoria de los procesos cuánticos involucrados en la física de partículas. Calquiera desviación que se podría ver entre los datos y simulaciones podrían ser el signo de la existencia de física nueva y, por lo tanto, el proceso de generación de muestras de Monte-Carlo tiene que ser muy preciso y suele ser costoso computacionalmente.

Para este análisis, se estudian los  $(137.1 \pm 2.0) \text{ fb}^{-1}$  de datos colectados por el detector CMS a lo largo de los años 2016, 2017 y 2018, con el objetivo de intentar descubrir algunos eventos en los cuales se observa materia oscura producida en asociación con uno o dos quarks top. Dado lo poco que se conoce sobre la materia oscura hoy en día, diferentes modelos se consideran para este tipo de búsquedas, considerando diferentes canales de producción, masas (de 10 GeV a 1 TeV) y espines (0 o 1) de mediadores y distintas masas para la materia oscura misma (de 1 GeV a 51 GeV), aunque no se espera que cambie mucho la cinemática del evento con esta parámetro.

Dos modelos principales de producción de materia oscura se consideran en este análisis: con un solo quark top, llevando al análisis llamado  $t/\bar{t}+\text{DM}$ , o bien con dos quarks top, caracterizado como  $t\bar{t}+\text{DM}$ . En ambos casos, en este trabajo se estudia solamente el estado final dileptonico, aunque se puedan observar 0, 1 o 2 leptones en el estado final, que vienen acompañados por un número de b-jets y una cantidad de energía faltante que puede variar. Todos estos modelos se diseñaron siguiendo las recomendaciones del ATLAS-CMS Dark Matter Forum [75].

### Estimación de fondos

Por otro lado, la estimación de los distintos fondos del análisis también es muy importante, ya que sirve de referencia. La mayoría de los fondos están bastante bien entendidos y se pueden por lo tanto estimar usando simulaciones de Monte-Carlo directamente, mientras que zonas de control se suelen definir para comprobar el acuerdo datos/MC de los fondos más importantes del análisis, tal y como el Drell-Yan (más que nada por su alta sección eficaz de producción, aunque reducir este fondo suele ser bastante fácil ya que su cinemática resulta ser muy distinta a la de nuestras señales de interés), el  $t\bar{t}$  del Modelo Estándar y la producción de un simple top.

La mayoría de los fondos se estiman directamente a partir de simulaciones matemáticas de tipo Monte-Carlo, y se pueden controlar en zonas de control enriquecidas en fondos en particular. Aún así, dos fondos tienen un tratamiento un poco peculiar:

- En el caso del Drell-Yan, se usan los datos mismos para calcular un factor de corrección que permite estimar el número de eventos de este proceso fuera del pico del Z a partir de los eventos observados dentro de este pico.
- Por otro lado, algunos procesos como el W+jets que solo tiene un lepton en el estado final también puede llevar a contaminar la zona de señal por culpa de los *fakes*, que suelen ser objetos como jets misidentificados por el detector como leptones, o leptones que no provienen del Primary Vertex de la colisión protón-protón. La estimación de este fondo se hace mediante un método data-driven ya que se no se espera poder simular este tipo de fondo de manera óptima usando simulaciones de Monte-Carlo, y se controla este fondo peculiar en una zona enriquecida en fakes, definida con dos leptones que tengan el mismo signo de carga.

## Definición de objetos

Los objetos que se quieren usar en el análisis mismo se tienen que definir de antemano:

- Los *triggers*, que permiten reducir el flujo de datos que se guardan, seleccionando eventos que tengan leptones aislados (o no) en este caso, son de esta manera los primeros objetos que hemos definido para los periodos de toma de datos 2016, 2017 y 2018. Se tomaron precauciones a la hora de seleccionar estos *triggers*, más que nada para evitar cualquier problema de *double counting* que podría producirse si un evento tiene la posibilidad de entrar en dos *triggers* distintos a la vez.
- Luego, otra parte fundamental del análisis consiste en definir los leptones que se quieren usar. De esta manera, definimos electrones basandónos en la definición usual de un electrón *tight* dada por el EGamma POG [149], con unos cortes adicionales en los parámetros de impacto  $d_{xy}$  y  $d_z$  para reducir la contaminación debida a los fakes en nuestras zonas de señal.
- Seguimos una estrategia similar a la hora de definir nuestros muones, basandónos esta vez en el *tight working point* del muon POG [151], con unos cortes adicionales también en los parámetros de impacto también.
- Luego se requiere que los jets tengan un  $p_T > 30$  GeV y un  $|\eta| < 2.4$  (límite de la extensión del tracker a partir de la cuál no se espera poder medir los jets con un alto grado de precisión) para formar parte del análisis. Además, el *tight working point* dado por el JET/MET POG [153] se usa, porque ofrece una eficiencia muy alta de rejección del ruido. Los jets de *pile-up* también se quitan aplicando además la definición de loose PU para jets de alto  $p_T$  (2017, 2018).
- Por fin, dada la presencia de quark tops que tienen un tiempo de vida pequeño y que decaen rápidamente a bottom quarks, los b-jets son muy importantes y se definen a partir del *medium working point* del POG de b-tagging y vertexing [157], para el cuál la tasa de misidentificación de un jet ligero como un b-jet es del 10%.

## Extracción de señal

Las señales de interés en este análisis están muy cerca de algunos fondos del Modelo Estándar desde el punto de vista de la cinética: se espera por ejemplo que sea complicado separar la señal de  $t\bar{t}+DM$  del proceso  $t\bar{t}$  sin que este sea acompañado por materia oscura. Se espera aún así que algunas variables presenten algo de discriminación entre estos tipos de procesos muy similares: por ejemplo, la MET debería ser estadísticamente un poquito más alta cuando consideramos la señal ya que, aunque el proceso  $t\bar{t}$  y la señal suelen presentar neutrinos con forma de energía transversa faltante en el detector, una contribución adicional a la MET se espera en el caso de la señal, por la aparición de un par de partículas de materia oscura.

En este trabajo en concreto, un método multivariado avanzado ha sido desarrollado con el objetivo de poder combinar el poder de discriminación de unas 6 variables que se esperan que presenten algo de discriminación entre las señales y los fondos del análisis, para poder definir zonas de señal enriquecidas en las señales de interés. Con este objetivo de discriminación en mente, se desarrollaron en paralelo dos métodos de *machine learning* distintos: una BDT, y una red neuronal ANN, ambas entrenadas con un conjunto etiquetado de simulaciones de Monte-Carlo para que sean capaces de asignar una probabilidad de pertenencia a cada una de las 2 categorías de interés: las dos señales ( $t\bar{t}+DM$  y  $t/\bar{t}+DM$ ) o fondo. Un largo proceso de optimización de las redes tuvo por lo tanto lugar, con el objetivo de definir redes tan eficientes como sea posible y, una vez el training hecho

y el proceso de validación terminado, pasamos los datos colectados a través de ambas redes para asignar un label a cada evento. De esta manera, definimos una única variable de salida que nos permite discriminar estos procesos y con la cuál hemos podido definir nuestras zonas de señal del análisis. Esta variable de salida se usa además para hacer un *shape analysis*, que consiste en aplicar métodos estadísticos avanzados para estudiar la forma de esta variable para los diferentes procesos del análisis, con el objetivo de tratar de descubrir una posible señal de materia oscura.

## A.5 Resultados obtenidos

El resultado más importante de este trabajo consiste en pintar los límites de exclusión de la sección eficaz de producción de las señales de interés, para las dos señales y para los 3 años de toma de datos por separado, y luego combinando todos estos parámetros. Esto se hace considerando diferentes fuentes de errores estadísticos y sistemáticos, ambos teóricos (tal y como el sistemático relacionado con la estimación de las PDF de los diferentes procesos) y experimentales (tal y como el error asociado a la medida de la luminosidad colectada por el detector CMS).



---

---

# Appendix B

---

## Samples used

### B.1 Data samples

All the data samples considered for this analysis are listed in Tables B.1, B.2 and B.3. The luminosity of each dataset has been computed using the Brilcalc tool provided by CMS [173], while the number of generated events has been obtained using the CERN official Data Aggregation System (DAS).

### B.2 Signal samples

The MC signal samples have been produced centrally, considering different dark matter and mediator masses and different production channels, for both signals of this analysis. All the mass points considered along with their respective cross sections can be found in Tables B.4 ( $t/\bar{t}$ +DM signal) and B.5 ( $t\bar{t}$ +DM signal). The mass points generated are the same for 2016, 2017 and 2018.

### B.3 Backgrounds samples

All the background MC samples considered for this analysis are listed in Tables B.6, B.7 and B.8 for 2016, 2017 and 2018 respectively.

Dataset	Events (size)	$\mathcal{L}$ [fb $^{-1}$ ]
<b>Run 2016B</b>		
/DoubleEG/Run2016B_ver2-Nano02Apr2020_ver2-v1/NANOAOD	143073268 (99.4Gb)	
/DoubleMuon/Run2016B_ver2-Nano02Apr2020_ver2-v1/NANOAOD	82535526 (53.2Gb)	
/MuonEG/Run2016B_ver2-Nano02Apr2020_ver2-v1/NANOAOD	32727796 (26.8Gb)	5.8
/SingleElectron/Run2016B_ver2-Nano02Apr2020_ver2-v1/NANOAOD	246440440 (167.8Gb)	
/SingleMuon/Run2016B_ver2-Nano02Apr2020_ver2-v1/NANOAOD	158145722 (96.4Gb)	
<b>Run 2016C</b>		
/DoubleEG/Run2016C-Nano02Apr2020-v1/NANOAOD	47677856 (35.3Gb)	
/DoubleMuon/Run2016C-Nano02Apr2020-v1/NANOAOD	27934629 (19.7Gb)	
/MuonEG/Run2016C-Nano02Apr2020-v1/NANOAOD	15405678 (12.8Gb)	2.6
/SingleElectron/Run2016C-Nano02Apr2020-v1/NANOAOD	97259854 (69.3Gb)	
/SingleMuon/Run2016C-Nano02Apr2020-v1/NANOAOD	67441308 (42.4Gb)	
<b>Run 2016D</b>		
/DoubleEG/Run2016D-Nano02Apr2020-v1/NANOAOD	53324960 (39.6Gb)	
/DoubleMuon/Run2016D-Nano02Apr2020-v1/NANOAOD	33861745 (24.1Gb)	
/MuonEG/Run2016D-Nano02Apr2020-v1/NANOAOD	23482352 (19.4Gb)	4.2
/SingleElectron/Run2016D-Nano02Apr2020-v1/NANOAOD	148167727 (104.4Gb)	
/SingleMuon/Run2016D-Nano02Apr2020-v1/NANOAOD	98017996 (61.3Gb)	
<b>Run 2016E</b>		
/DoubleEG/Run2016E-Nano02Apr2020-v1/NANOAOD	49877710 (37.9Gb)	
/DoubleMuon/Run2016E-Nano02Apr2020-v1/NANOAOD	28246946 (20.8Gb)	
/MuonEG/Run2016E-Nano02Apr2020-v2/NANOAOD	22519303 (19.0Gb)	4.0
/SingleElectron/Run2016E-Nano02Apr2020-v1/NANOAOD	117321545 (86.5Gb)	
/SingleMuon/Run2016E-Nano02Apr2020-v1/NANOAOD	90984718 (58.7Gb)	
<b>Run 2016F</b>		
/DoubleEG/Run2016F-Nano02Apr2020-v1/NANOAOD	34577629 (26.9Gb)	
/DoubleMuon/Run2016F-Nano02Apr2020-v1/NANOAOD	20329921 (15.3Gb)	
/MuonEG/Run2016F-Nano02Apr2020-v1/NANOAOD	16002165 (13.6Gb)	3.1
/SingleElectron/Run2016F-Nano02Apr2020-v1/NANOAOD	70593532 (51.4Gb)	
/SingleMuon/Run2016F-Nano02Apr2020-v1/NANOAOD	65489554 (42.4Gb)	
<b>Run 2016G</b>		
/DoubleEG/Run2016G-Nano02Apr2020-v1/NANOAOD	78797031 (61.6Gb)	
/DoubleMuon/Run2016G-Nano02Apr2020-v1/NANOAOD	45235604 (34.2Gb)	
/MuonEG/Run2016G-Nano02Apr2020-v1/NANOAOD	33854612 (29.0Gb)	7.6
/SingleElectron/Run2016G-Nano02Apr2020-v1/NANOAOD	153363109 (109.2Gb)	
/SingleMuon/Run2016G-Nano02Apr2020-v1/NANOAOD	149912248 (94.6Gb)	
<b>Run 2016H</b>		
/DoubleEG/Run2016H-Nano02Apr2020-v1/NANOAOD	85388734 (67.7Gb)	
/DoubleMuon/Run2016H-Nano02Apr2020-v1/NANOAOD	48912812 (37.3Gb)	
/MuonEG/Run2016H-Nano02Apr2020-v1/NANOAOD	29236516 (26.0Gb)	8.6
/SingleElectron/Run2016H-Nano02Apr2020-v1/NANOAOD	128854598 (93.8Gb)	
/SingleMuon/Run2016H-Nano02Apr2020-v1/NANOAOD	174035164 (110.2Gb)	

Table B.1: Datasets collected in 2016 and considered for this analysis.

Dataset	Events (size)	$\mathcal{L}$ [fb $^{-1}$ ]
<b>Run 2017B</b>		
/DoubleEG/Run2017B-Nano02Apr2020-v1/NANOAOD	58088760 (46.6Gb)	
/DoubleMuon/Run2017B-Nano02Apr2020-v1/NANOAOD	14501767 (10.8Gb)	
/SingleElectron/Run2017B-Nano02Apr2020-v1/NANOAOD	60537490 (42.2Gb)	4.8
/SingleMuon/Run2017B-Nano02Apr2020-v1/NANOAOD	136300266 (86.2Gb)	
/MuonEG/Run2017B-Nano02Apr2020-v1/NANOAOD	4453465 (4.1Gb)	
<b>Run 2017C</b>		
/DoubleEG/Run2017C-Nano02Apr2020-v1/NANOAOD	65181125 (53.8Gb)	
/DoubleMuon/Run2017C-Nano02Apr2020-v1/NANOAOD	49636525 (39.5Gb)	
/SingleElectron/Run2017C-Nano02Apr2020-v1/NANOAOD	136637888 (102.5Gb)	9.7
/SingleMuon/Run2017C-Nano02Apr2020-v1/NANOAOD	165652756 (109.5Gb)	
/MuonEG/Run2017C-Nano02Apr2020-v1/NANOAOD	15595214 (15.0Gb)	
<b>Run 2017D</b>		
/DoubleEG/Run2017D-Nano02Apr2020-v1/NANOAOD	25911432 (21.6Gb)	
/DoubleMuon/Run2017D-Nano02Apr2020-v1/NANOAOD	23075733 (18.6Gb)	
/SingleElectron/Run2017D-Nano02Apr2020-v1/NANOAOD	51526710 (38.5Gb)	4.2
/SingleMuon/Run2017D-Nano02Apr2020-v1/NANOAOD	70361660 (47.2Gb)	
/MuonEG/Run2017D-Nano02Apr2020-v1/NANOAOD	9164365 (8.9Gb)	
<b>Run 2017E</b>		
/DoubleEG/Run2017E-Nano02Apr2020-v1/NANOAOD	56233597 (49.8Gb)	
/DoubleMuon/Run2017E-Nano02Apr2020-v1/NANOAOD	51589091 (44.4Gb)	
/SingleElectron/Run2017E-Nano02Apr2020-v1/NANOAOD	102121689 (81.3Gb)	9.3
/SingleMuon/Run2017E-Nano02Apr2020-v1/NANOAOD	154630534 (111.0Gb)	
/MuonEG/Run2017E-Nano02Apr2020-v1/NANOAOD	19043421 (19.2Gb)	
<b>Run 2017F</b>		
/DoubleEG/Run2017F-Nano02Apr2020-v1/NANOAOD	74307066 (67.1Gb)	
/DoubleMuon/Run2017F-Nano02Apr2020-v1/NANOAOD	79756560 (68.0Gb)	
/SingleElectron/Run2017F-Nano02Apr2020-v1/NANOAOD	128467223 (105.2Gb)	13.5
/SingleMuon/Run2017F-Nano02Apr2020-v1/NANOAOD	242135500 (178.3Gb)	
/MuonEG/Run2017F-Nano02Apr2020-v1/NANOAOD	25776363 (26.3Gb)	

Table B.2: Datasets collected in 2017 and considered for this analysis.

Dataset	Events (size)	$\mathcal{L}$ [fb $^{-1}$ ]
<b>Run 2018A</b>		13.5
/DoubleMuon/Run2018A-Nano02Apr2020-v1/NANO AOD	75499908 (62.6Gb)	
/EGamma/Run2018A-Nano02Apr2020-v1/NANO AOD	327843843 (261.8Gb)	
/SingleMuon/Run2018A-Nano02Apr2020-v1/NANO AOD	241608232 (167.7Gb)	
/MuonEG/Run2018A-Nano02Apr2020-v1/NANO AOD	32958503 (32.3Gb)	
<b>Run 2018B</b>		6.8
/DoubleMuon/Run2018B-Nano02Apr2020-v1/NANO AOD	35057758 (28.3Gb)	
/EGamma/Run2018B-Nano02Apr2020-v1/NANO AOD	153822427 (123.1Gb)	
/SingleMuon/Run2018B-Nano02Apr2020-v1/NANO AOD	119918017 (82.3Gb)	
/MuonEG/Run2018B-Nano02Apr2020-v1/NANO AOD	16211567 (15.8Gb)	
<b>Run 2018C</b>		6.6
/DoubleMuon/Run2018C-Nano02Apr2020-v1/NANO AOD	34565869 (27.6Gb)	
/EGamma/Run2018C-Nano02Apr2020-v1/NANO AOD	147827904 (119.2Gb)	
/SingleMuon/Run2018C-Nano02Apr2020-v1/NANO AOD	110032072 (75.7Gb)	
/MuonEG/Run2018C-Nano02Apr2020-v1/NANO AOD	15652198 (15.3Gb)	
<b>Run 2018D</b>		32.0
/DoubleMuon/Run2018D-Nano02Apr2020_ver2-v1/NANO AOD	168605834 (128.6Gb)	
/EGamma/Run2018D-Nano02Apr2020-v1/NANO AOD	751348648 (583.6Gb)	
/SingleMuon/Run2018D-Nano02Apr2020-v1/NANO AOD	513867253 (344.5Gb)	
/MuonEG/Run2018D-Nano02Apr2020_ver2-v1/NANO AOD	71961587 (68.6Gb)	

Table B.3: Datasets collected in 2018 and considered for this analysis.

Mass point	Cross-section [pb]
<b>Scalar mediators</b>	
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi10	$4.959 \cdot 10^{-2}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi20	$3.235 \cdot 10^{-2}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi50	$1.323 \cdot 10^{-2}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi100	$5.633 \cdot 10^{-3}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi150	$3.397 \cdot 10^{-3}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi200	$2.359 \cdot 10^{-3}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi250	$1.720 \cdot 10^{-3}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi300	$1.328 \cdot 10^{-3}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi350	$1.018 \cdot 10^{-3}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi400	$6.717 \cdot 10^{-4}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi450	$4.535 \cdot 10^{-4}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi500	$3.206 \cdot 10^{-4}$
DMscalar_Dilepton_top_tWChan_Mchi1_Mphi1000	$3.045 \cdot 10^{-5}$
<b>Pseudoscalar mediators</b>	
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi10	$6.151 \cdot 10^{-3}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi20	$5.869 \cdot 10^{-3}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi50	$4.946 \cdot 10^{-3}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi100	$3.658 \cdot 10^{-3}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi150	$2.754 \cdot 10^{-3}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi200	$2.097 \cdot 10^{-3}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi250	$1.616 \cdot 10^{-3}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi300	$1.253 \cdot 10^{-3}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi350	$7.851 \cdot 10^{-4}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi400	$4.371 \cdot 10^{-4}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi450	$3.095 \cdot 10^{-4}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi500	$2.321 \cdot 10^{-4}$
DMpseudoscalar_Dilepton_top_tWChan_Mchi1_Mphi1000	$2.791 \cdot 10^{-5}$

Table B.4: Signal samples mass points and LO dileptonic cross-sections considered for the  $t/\bar{t}+DM$  signal used in this analysis.

Mass point	Cross-section [pb]
<b>Scalar mediators</b>	
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_50	$3.405 \cdot 10^{-1}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_100	$8.027 \cdot 10^{-2}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_150	$2.673 \cdot 10^{-2}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_200	$1.158 \cdot 10^{-2}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_250	$6.020 \cdot 10^{-3}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_300	$3.579 \cdot 10^{-3}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_350	$2.376 \cdot 10^{-3}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_400	$1.443 \cdot 10^{-3}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_450	$9.025 \cdot 10^{-4}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_500	$6.204 \cdot 10^{-4}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_20_mPhi_100	$7.993 \cdot 10^{-2}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_30_mPhi_100	$8.052 \cdot 10^{-2}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_40_mPhi_100	$8.147 \cdot 10^{-2}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_45_mPhi_100	$8.319 \cdot 10^{-2}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_49_mPhi_100	$8.304 \cdot 10^{-2}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_51_mPhi_100	$9.735 \cdot 10^{-4}$
TTbarDMJets_Dilepton_scalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_55_mPhi_100	$4.835 \cdot 10^{-4}$
<b>Pseudoscalar mediators</b>	
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_50	$3.440 \cdot 10^{-2}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_100	$2.164 \cdot 10^{-2}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_150	$1.414 \cdot 10^{-2}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_200	$9.773 \cdot 10^{-3}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_250	$6.753 \cdot 10^{-3}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_300	$4.808 \cdot 10^{-3}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_350	$2.742 \cdot 10^{-3}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_400	$1.409 \cdot 10^{-3}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_450	$9.302 \cdot 10^{-4}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_1_mPhi_500	$6.618 \cdot 10^{-4}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_20_mPhi_100	$2.166 \cdot 10^{-2}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_30_mPhi_100	$2.164 \cdot 10^{-2}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_40_mPhi_100	$2.162 \cdot 10^{-2}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_45_mPhi_100	$2.180 \cdot 10^{-2}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_49_mPhi_100	$2.151 \cdot 10^{-2}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_51_mPhi_100	$1.993 \cdot 10^{-3}$
TTbarDMJets_Dilepton_pseudoscalar_LO_TuneCP5_13TeV-madgraph-mcatnlo-pythia8_mChi_55_mPhi_100	$7.750 \cdot 10^{-4}$

Table B.5: Signal samples mass points and LO dileptonic cross-sections considered for the  $t\bar{t}+DM$  signal used in this analysis.

Process	Sample	Cross section [pb]
Drell-Yan	DYJetsToLL_M-10to50_TuneCUETP8M1_13TeV-madgraphMLM-pythia8 ( $H_T < 70$ GeV)	18610.0
	DYJetsToLL_M-5to50_HT-70to100_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	303.8
	DYJetsToLL_M-5to50_HT-100to200_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	224.2
	DYJetsToLL_M-5to50_HT-200to400_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	37.2
	DYJetsToLL_M-5to50_HT-400to600_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	3.581
	DYJetsToLL_M-5to50_HT-600toInf_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	1.124
	DYJetsToLL_M-50_TuneCUETP8M1_13TeV-madgraphMLM-pythia8 ( $H_T < 70$ GeV)	6025.20
	DYJetsToLL_M-50_HT-70to100_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	169.9
	DYJetsToLL_M-50_HT-100to200_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	147.4
	DYJetsToLL_M-50_HT-200to400_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	40.99
	DYJetsToLL_M-50_HT-400to600_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	5.678
	DYJetsToLL_M-50_HT-600to800_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	1.367
	DYJetsToLL_M-50_HT-800to1200_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	0.6304
	DYJetsToLL_M-50_HT-1200to2500_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	0.1514
	DYJetsToLL_M-50_HT-2500toInf_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	0.003565
TTTo2L2Nu	TTTo2L2Nu_TuneCUETP8M2_ttHtranche3_13TeV-powheg-pythia8	87.310
Single top	ST_tW_antitop_5f_inclusiveDecays_13TeV-powheg-pythia8_TuneCUETP8M1	35.60
	ST_tW_top_5f_inclusiveDecays_13TeV-powheg-pythia8_TuneCUETP8M1	35.60
TTToSemiLeptonic	TToSemilepton_TuneCUETP8M2_ttHtranche3_13TeV-powheg-pythia8	364.35
ttV	TTZToLLNuNu_M-10_TuneCP5_PSweights_13TeV-amcatnlo-pythia8	0.2529
	TTZToQQ_TuneCUETP8M1_13TeV-amcatnlo-pythia8	0.5297
	TTWJetsToLNu_TuneCUETP8M1_13TeV-amcatnloFXFX-madspin-pythia8	0.2043
	TTWJetsToQQ_TuneCUETP8M1_13TeV-amcatnloFXFX-madspin-pythia8	0.4062
VZ	WWTo2L2Nu_13TeV-powheg	12.178
	WZTo3LNu_TuneCUETP8M1_13TeV-powheg-pythia8	4.42965
	WZTo2L2Q_13TeV_amcatnloFXFX_madspin_pythia8	5.595
	ZZTo2L2Nu_13TeV_powheg_pythia8	0.5640
	ZZTo2L2Q_13TeV_powheg_pythia8	3.22
Others	WWW, WWZ, WZZ, ZZZ, WWG	//

Table B.6: Main 2016 MC simulations for the different background processes considered for this analysis and their respective cross sections.

Process	Sample	Cross section [pb]
Drell-Yan	DYJetsToLL_M-10to50_TuneCP5_13TeV-madgraphMLM-pythia8 ( $H_T < 100$ GeV)	18610
	DYJetsToLL_M-4to50_HT-100to200_TuneCP5_13TeV-madgraphMLM-pythia8	204.0
	DYJetsToLL_M-4to50_HT-200to400_TuneCP5_13TeV-madgraphMLM-pythia8	54.39
	DYJetsToLL_M-4to50_HT-400to600_TuneCP5_13TeV-madgraphMLM-pythia8	5.697
	DYJetsToLL_M-4to50_HT-600toInf_TuneCP5_13TeV-madgraphMLM-pythia8	1.85
	DYJetsToLL_M-50_TuneCP5_13TeV-madgraphMLM-pythia8 ( $H_T < 70$ GeV)	6025.20
	DYJetsToLL_M-50_HT-70to100_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	169.9
	DYJetsToLL_M-50_HT-100to200_TuneCP5_13TeV-madgraphMLM-pythia8	161.1
	DYJetsToLL_M-50_HT-200to400_TuneCP5_13TeV-madgraphMLM-pythia8	48.66
	DYJetsToLL_M-50_HT-400to600_TuneCP5_13TeV-madgraphMLM-pythia8	6.968
	DYJetsToLL_M-50_HT-600to800_TuneCP5_13TeV-madgraphMLM-pythia8	1.743
	DYJetsToLL_M-50_HT-800to1200_TuneCP5_13TeV-madgraphMLM-pythia8	0.8052
	DYJetsToLL_M-50_HT-1200to2500_TuneCP5_13TeV-madgraphMLM-pythia8	0.1933
	DYJetsToLL_M-50_HT-2500toInf_TuneCP5_13TeV-madgraphMLM-pythia8	0.003468
TTTo2L2Nu	TTTo2L2Nu_TuneCP5_13TeV-powheg-pythia8	87.310
Single top	ST_tW_antitop_5f_inclusiveDecays_TuneCP5_13TeV-powheg-pythia8	35.60
	ST_tW_top_5f_inclusiveDecays_TuneCP5_13TeV-powheg-pythia8	35.60
TTToSemiLeptonic	TTToSemiLeptonic_TuneCP5_13TeV-powheg-pythia8	364.35
ttV	TTZToLLNuNu_M-10_TuneCP5_PSweights_13TeV-amcatnlo-pythia8	0.2529
	TTZToQQ_TuneCUETP8M1_13TeV-amcatnlo-pythia8	0.5297
	TTWJetsToLNu_TuneCP5_PSweights_13TeV-amcatnloFXFX-madspin-pythia8	0.2043
	TTWJetsToQQ_TuneCUETP8M1_13TeV-amcatnloFXFX-madspin-pythia8	0.4062
VZ	WWTo2L2Nu_NNPDF31_TuneCP5_PSweights_13TeV-powheg-pythia8	12.178
	WZTo3LNu_TuneCUETP8M1_13TeV-powheg-pythia8	4.42965
	WZTo2L2Q_13TeV_amcatnloFXFX_madspin_pythia8	5.595
	ZZTo2L2Nu_13TeV_powheg_pythia8	0.5640
	ZZTo2L2Q_13TeV_amcatnloFXFX_madspin_pythia8	3.22
	WWW, WWZ, WZZ, ZZZ, WWG	//

Table B.7: Main 2017 MC simulations for the different background processes considered for this analysis and their respective cross sections.

Process	Sample	Cross section [pb]
Drell-Yan	DYJetsToLL_M-10to50_TuneCP5_13TeV-madgraphMLM-pythia8 ( $H_T < 100$ GeV)	18610.0
	DYJetsToLL_M-4to50_HT-100to200_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	204.0
	DYJetsToLL_M-4to50_HT-200to400_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	54.39
	DYJetsToLL_M-4to50_HT-400to600_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	5.697
	DYJetsToLL_M-4to50_HT-600toInf_TuneCP5_PSWeights_13TeV-madgraphMLM-pythia8	1.85
	DYJetsToLL_M-50_TuneCP5_13TeV-madgraphMLM-pythia8 ( $H_T < 70$ GeV)	6025.20
	DYJetsToLL_M-50_HT-70to100_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	169.9
	DYJetsToLL_M-50_HT-100to200_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	161.1
	DYJetsToLL_M-50_HT-200to400_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	48.66
	DYJetsToLL_M-50_HT-400to600_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	6.968
	DYJetsToLL_M-50_HT-600to800_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	1.743
	DYJetsToLL_M-50_HT-800to1200_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	0.8052
	DYJetsToLL_M-50_HT-1200to2500_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	0.1933
	DYJetsToLL_M-50_HT-2500toInf_TuneCP5_PSweights_13TeV-madgraphMLM-pythia8	0.003468
TTTo2L2Nu	TTTo2L2Nu_TuneCP5_13TeV-powheg-pythia8	87.310
Single top	ST_tW_antitop_5f_inclusiveDecays_TuneCP5_13TeV-powheg-pythia8	35.60
	ST_tW_top_5f_inclusiveDecays_TuneCP5_13TeV-powheg-pythia8	35.60
TTToSemiLeptonic	TTToSemiLeptonic_TuneCP5_13TeV-powheg-pythia8	364.35
ttV	TTZToLLNuNu_M-10_TuneCP5_PSweights_13TeV-amcatnlo-pythia8	0.2529
	TTZToQQ_TuneCUETP8M1_13TeV-amcatnlo-pythia8	0.5297
	TTWJetsToLNu_TuneCP5_PSweights_13TeV-amcatnloFXFX-madspin-pythia8	0.2043
	TTWJetsToQQ_TuneCUETP8M1_13TeV-amcatnloFXFX-madspin-pythia8	0.4062
VZ	WWTo2L2Nu_NNPDF31_TuneCP5_13TeV-powheg-pythia8	12.178
	WZTo3LNu_TuneCP5_13TeV-amcatnloFXFX-pythia8	4.42965
	WZTo2L2Q_13TeV_amcatnloFXFX_madspin_pythia8	5.595
	ZZTo2L2Nu_TuneCP5_13TeV_powheg_pythia8	0.5640
	ZZTo2L2Q_13TeV_amcatnloFXFX_madspin_pythia8	3.22
Others	WWW, WWZ, WZZ, ZZZ, WWG	//

Table B.8: Main 2018 MC simulations for the different background processes considered for this analysis and their respective cross sections.



---

---

# Appendix C

---

## MVA optimization

The study of the different hyper-parameters defining both the ANN and the BDT considered in this work was an important part of the analysis, to make sure to use optimized parameters when defining our MVA and to get the best possible discrimination between the backgrounds and the two signals of interest, the  $t/\bar{t}$ +DM and the  $t\bar{t}$ +DM. In general, such optimization is a complex task to perform because of the high number of hyper-parameters that usually influence the quality of the final results. As a starting point, the optimal parameters found by DESY (for the BDT) and IFCA (for the ANN) during the 2016 analysis were used.

The global strategy used for the actual optimization process is quite straightforward: even though some of the hyper-parameters considered are actually related to each other, we decided for the sake of simplicity to treat them as independent and to optimize the network by optimizing each one of these parameters individually. Among all the candidate values for a given parameter, the one leading to the best possible ROC curve while keeping the time needed to train the network reasonably low was then simply chosen for the analysis. The training optimization has been performed for the 100 GeV scalar mediators for both signal processes, an intermediate mass point that should feature some of the lowest discrimination between the signals and the backgrounds, given their kinematics similarities and a similar number of background and signal events was used to train the different methods.

### Boosted Decisions Trees

Starting with the BDT, the first obvious hyper-parameter among all the parameters typical of a BDT and described in Section 7.2.1 that we decided to optimize was the number of trees, whose impact on the final discrimination obtained between the processes can be seen in the ROC curves shown in Figure C.1.

Another important parameter of the BDT and which was optimized is the maximal depth of the trees that we define, as shown in Figure C.2. However, given the relative small number of input variables used in this analysis and given to the BDT, this parameter does not seem to be highly relevant since the results obtained are similar in all the cases.

The  $n_{\text{cut}}$  parameter, corresponding to the granularity used when scanning over the variable range to find the optimal splitting criterion has also been studied, as shown in Figure C.3. Finally,

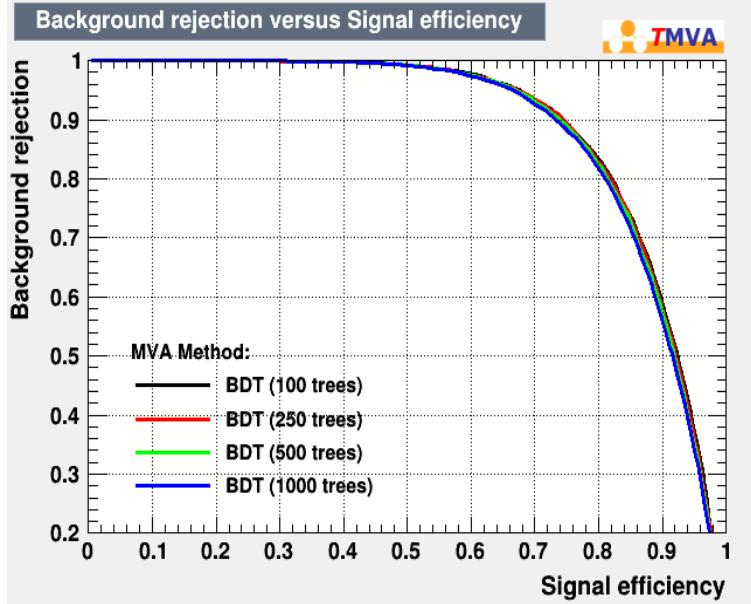


Figure C.1: ROC curves obtained for the BDT by trying out different number of trees.

the impact that the shrinkage, or learning rate for the GradBoost algorithm, has on the final discrimination results has also been studied and is shown in Figure C.4.

For most of the hyper-parameters described so far, we usually expect that increasing (or decreasing) their values will result in better discrimination results but unfortunately, this usually also increases the time needed to train the BDT so a compromise had to be found for each one of the hyper-parameters described here, settling at the end for the set of parameters described in Table 7.1.

### Analysis Neural Network

The same process has been repeated by considering this time the ANN of the analysis. In this case, we decided to first of all focus the optimization of the neural network on its actual architecture, trying out different number of hidden layers and different number of neurons in each layer, while trying to keep the validation loss obtained to a minimum. The results for the optimization of this first parameter can be found in Figure C.5. We can see in these plots that with the 50%/50% test/training splitting used, which gave us around 45.000 total events for the training, all the networks seem to give similar results. At the end of the day, the best combination between the background rejection and signal efficiency was obtained for the architecture with 3 hidden layers made out of 80/80/40 neurons.

After this first step of optimization, we kept the architecture which seemed to give the best results and started a subsequent optimization process, evaluating this time the learning rate of the Adam algorithm used, as shown in Figure C.6. Usually, a larger learning rate allow us to avoid any eventual local minimum in the curve corresponding to the error function defined and to reach the minimum global faster, while a lower value is typically a bit more precise.

The last parameter that has been optimized in this case is the batch size, allowing us to gain some time by avoiding updating the weight of the network each time an input event is considered but by instead grouping them in batches. The results obtained in this case are shown in Figure C.7. In this case, time was an important factor when choosing the parameter to use for the analysis, given the fact that it takes around two to three more times to train a ANN with a batch size 20 than a higher value of 100.

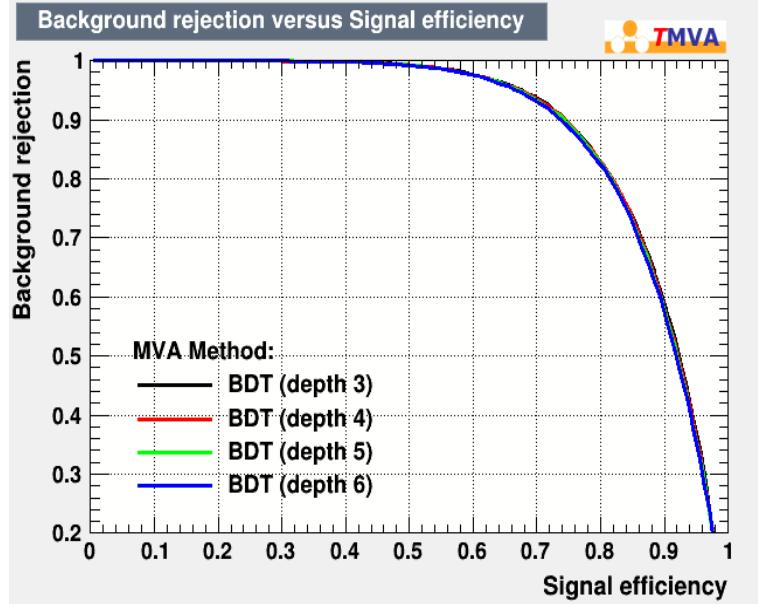


Figure C.2: ROC curves obtained for the BDT by trying out different maximum depths.

For all the previous tests, we took care of defining a number of epochs high enough by using a simple trial and error process, making sure to run at least 200 epochs and then making sure that the validation loss obtained did not decrease for at least 30 consecutive training epochs.

### Input variables chosen

Finally, the input variables given to the MVA is also an extremely important parameter to study. It has not been mentioned so far because it actually received a different treatment than the others, given its importance. Instead of simply changing this parameter to study the background rejection curve obtained, we decided to train several networks using different sets of input variables, to apply the weights for each network to the actual data and MC and to repeat the calculation of the upper limits in each case. The objective of this process was to find the optimal set of input variables that leads directly to the lowest possible upper limits for the different models considered in a more precise way. We therefore decided to train several times the ANN and the BDT considering each time a different set of input variables, all defined in Section 7, as shown in Table C.1.

As an example, the scalar  $t\bar{t}$ +DM upper limits obtained considering different sets of variables and a global training (for the 2016 data taking period only, without considering any systematics and for a single 100 GeV scalar training) are shown in Figures C.8 and C.9. Obviously, it is important to note that the optimization of the MVA has been performed at an early stage of the analysis, and this has several consequences that can be observed on these plots: i) the blinding policy was still in place during this process, meaning that the limits shown have been obtained with a single  $\text{fb}^{-1}$  and, more importantly, ii) the actual values obtained for the upper limits plots are not accurate and should not be considered to set exclusion limits on some DM production models, because the systematics were not considered when the optimization was done, for example. However, even though the actual value of the upper limit should not be trusted, the differences observed between the trainings considering the different sets of input variables is able to bring us valuable information in order to choose the input variables for the actual analysis.

According to these last results, we decided to use the ANN for the analysis instead of the BDT, trained with the simplest possible set of input variables, plus one of the two spin correlated variables from the set 3 (mainly because the two variables were found to be correlated at 99%), since this

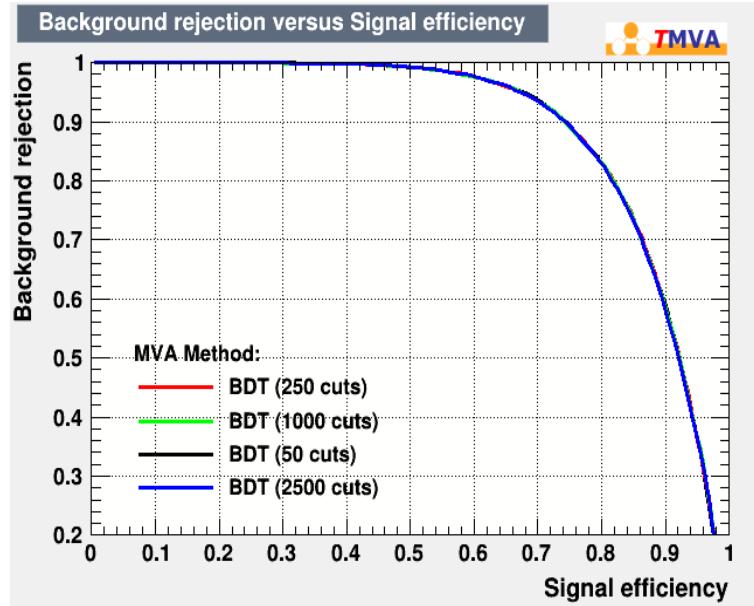


Figure C.3: ROC curves obtained for the BDT by trying out different  $n_{\text{cut}}$  values.

seems to be the set of variables able to give us the best possible upper limits around 100 GeV, the mass point chosen for the training process.

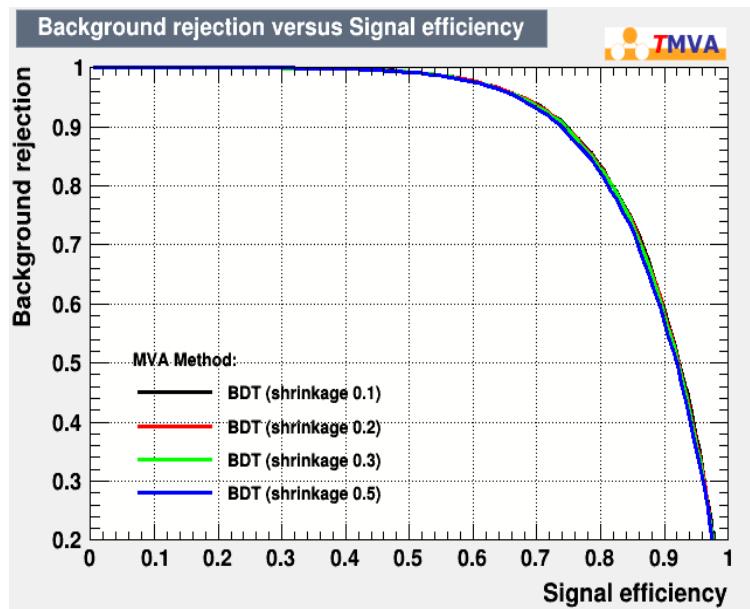


Figure C.4: ROC curves obtained for the BDT by trying out different shrinkage values.

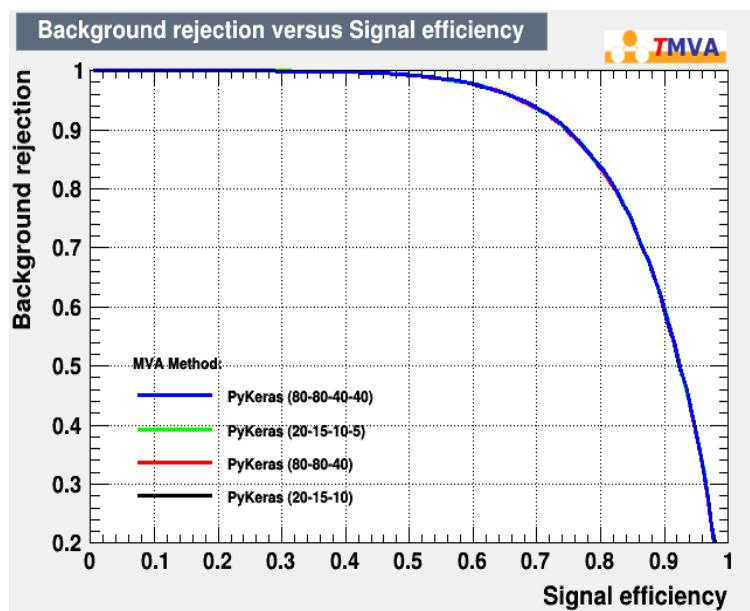


Figure C.5: ROC curves obtained for the ANN by trying out different architectures.

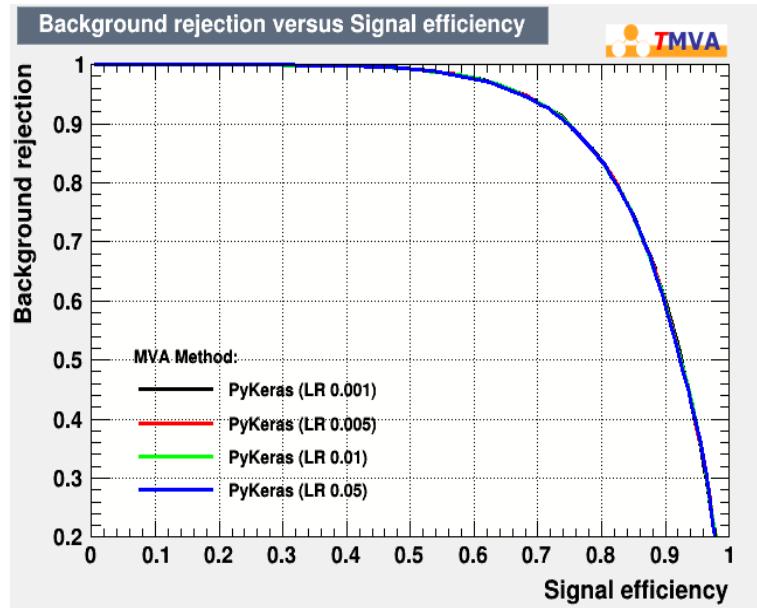


Figure C.6: ROC curves obtained for the ANN by trying out different learning rates.

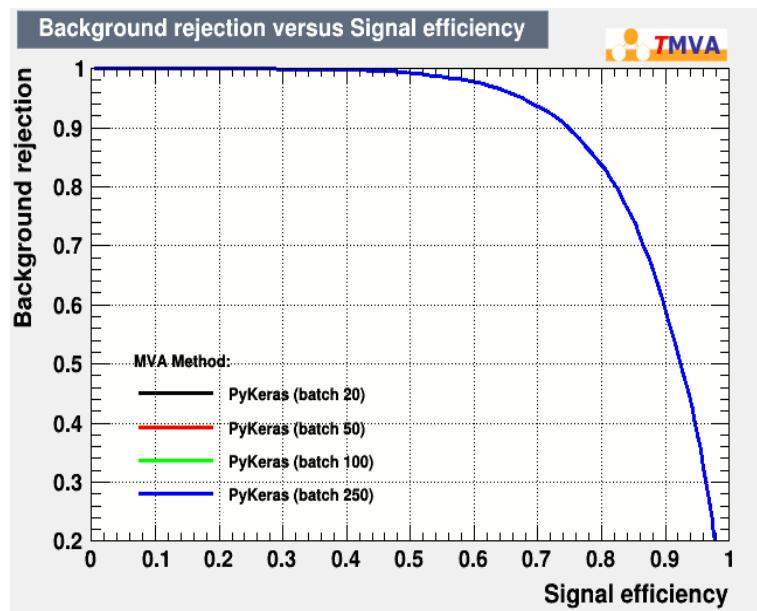


Figure C.7: ROC curves obtained for the ANN by trying out different batch sizes.

	Variables
Set 1	Number of b-jets $m_{bl}^t$ pfMET $\Delta\Phi(E_T^{\text{miss}}, ll)$ $M_{T2}^{ll}$
Set 2	Set 1 and $m_{T2}(bl, bl)$ massT
Set 3	Set 2 and Spin correlated variables
Set 4	Set 3 and $r2l$ $r2l4j$
Set 5	Set 4 and reco weight $w$
Set 6	Set 5 and Dark $p_T$ Overlapping factor $R$

Table C.1: Sets of input variables considered when optimizing the MVA method.

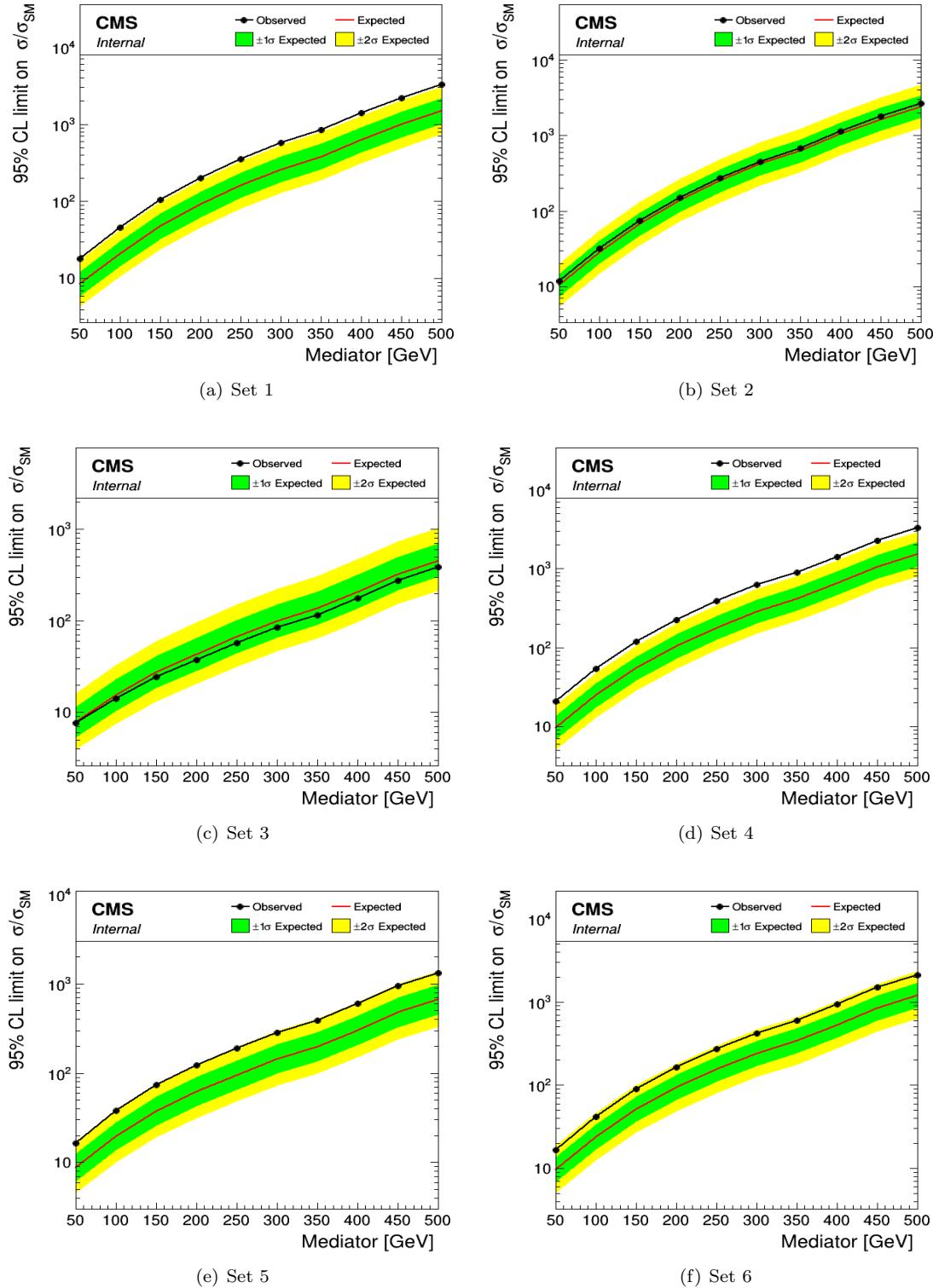


Figure C.8: Upper limits obtained for the scalar  $t/\bar{t}+DM$  signal considering the optimal BDT and different sets of variables used as input.

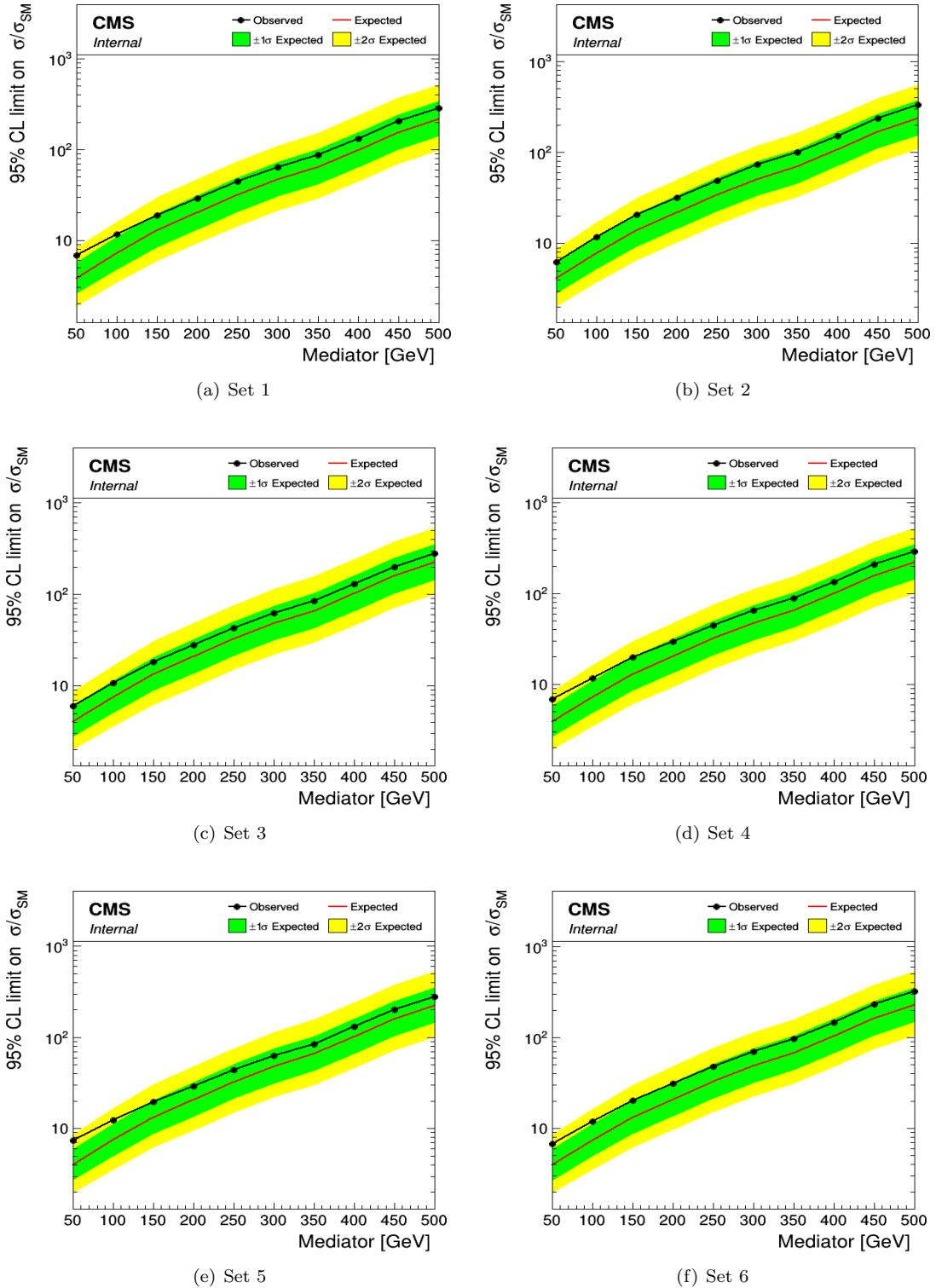


Figure C.9: Upper limits obtained for the scalar  $t/\bar{t}+DM$  signal considering the optimal ANN and different sets of variables used as input.



---

---

## Appendix D

---

### Pulls and impacts plots

Given the large number of systematic uncertainties considered in this analysis, as described in Section 8.2, it is usually useful to try and understand the impact each one of them has on the final upper limits on the signal strengths obtained. This is where the so-called impact plots become useful, since they allow us to study the quality of the systematic estimation, expressed as the difference between the maximum likelihood estimator and our own estimation, along with the importance of every systematic of the analysis, according to the effect each systematic has on the signal strength  $\mu$ , by rerunning the fit performed with each nuisance parameters fixed at their  $\pm 1\sigma$  values, to understand the impact they might have.

The impact plots obtained for different mediators are shown in Figures D.1 and D.2, for several scalar and pseudoscalar mediators, respectively.

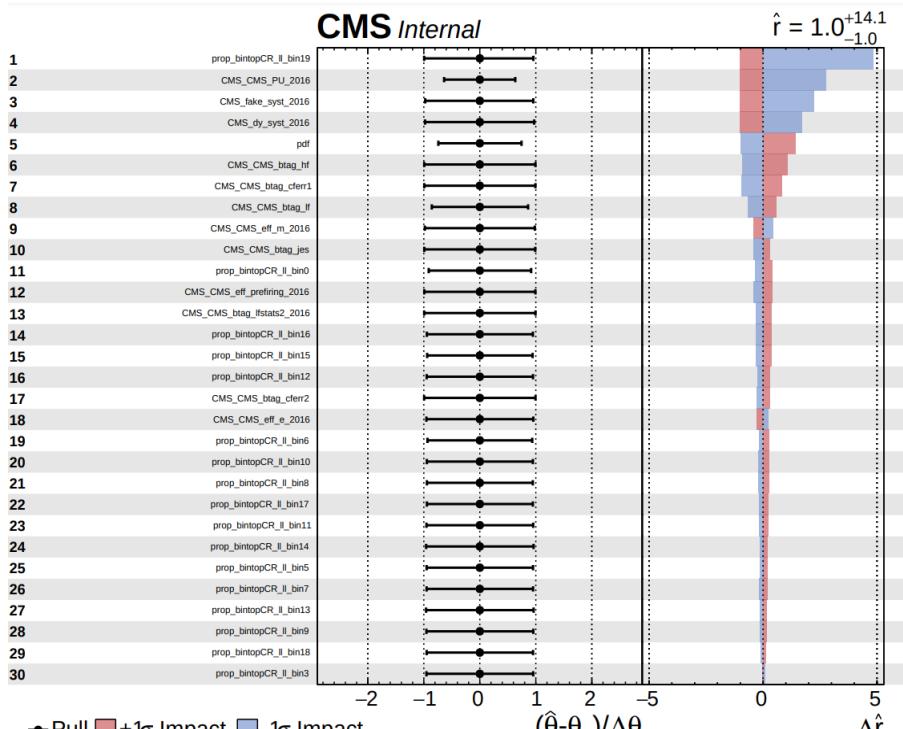
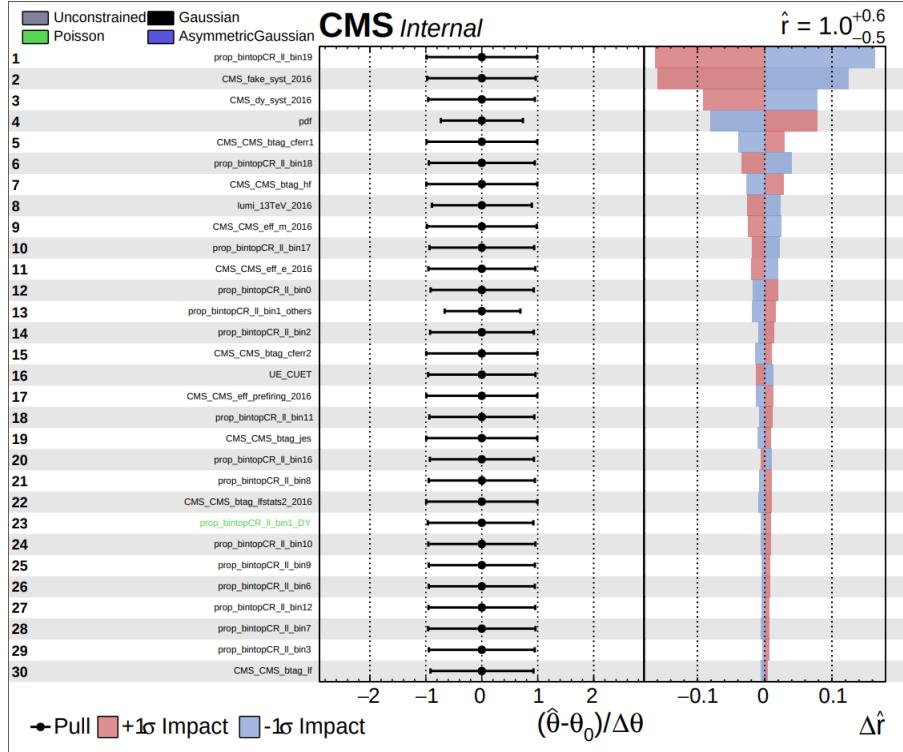


Figure D.1: 2016 impacts and pulls obtained for the 30 most important systematics of the analysis, considering the 100 and 500 GeV scalar mediators.

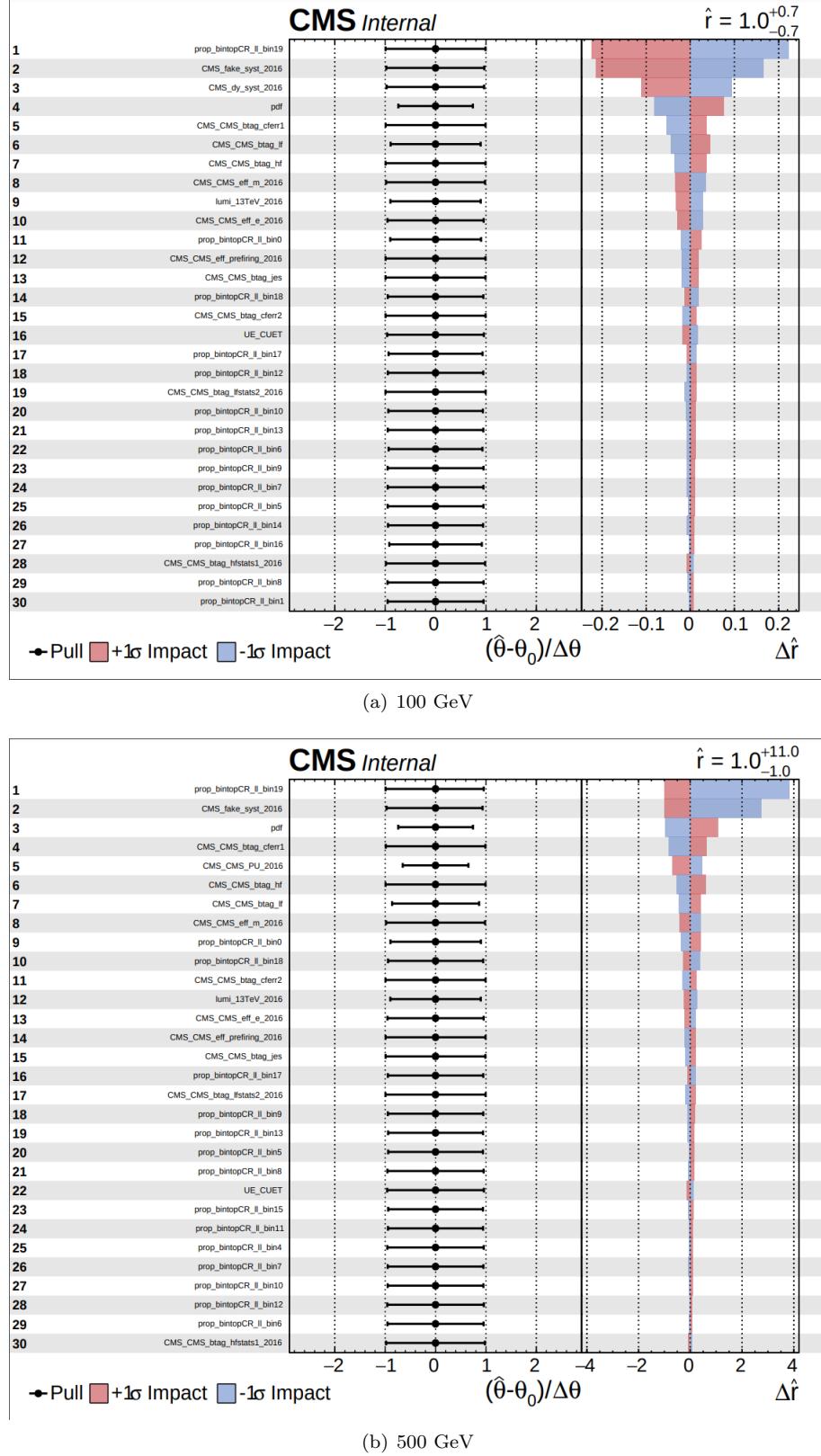


Figure D.2: 2016 impacts and pulls obtained for the 30 most important systematics of the analysis, considering the 100 and 500 GeV pseudoscalar mediators.



---

## List of figures

2.1	Representation of the 12 fermions of the SM [25] along with the main force carriers and the Higgs boson, discovered in 2012 and completing the SM. . . . .	6
2.2	Expected and observed rotation curves of the galaxy NGC 6503 [9]. The black dots correspond to the data and the <i>luminous</i> line corresponds to the rotation curve decreasing as $r^{-1/2}$ expected from Newtonian dynamics. . . . .	8
2.3	Anisotropies at the $10^{-5}$ level in the temperature of the CMB, as observed by the Planck satellite in 2018 [35]. . . . .	9
2.4	Power spectrum of the CMB obtained by Planck, representing the fluctuations of the temperature of the radiation with respect to the angular angle of observation [37].	10
2.5	Mass distributions obtained by the Magellan telescope in the visible (on the left) and the Chandra telescope on the X-rays spectrum (on the right) telescopes of the Bullet Cluser. Being shifted compared to each other, this is yet another clear evidence for the existence of DM [39]. . . . .	11
2.6	Computer simulations for cold (on the left) and warm (on the right) DM scenarios and their impact on a galactic halo at 0 redshift [42]. . . . .	12
2.7	Schematic representation of the freeze-out process, representing the abundance of a 500 GeV DM as $Y_\chi$ with respect to the time and the impact of increasing cross-section annihilation values on this freeze-out abundance [43]. . . . .	13
2.8	Halo fraction upper limits at the 95% CL compared to the mass of the lensing object for different MACHO models considered by the EROS (on the top) and the MACHO (on the bottom) collaborations [47]. . . . .	14
2.9	Neutrino cross section of interaction from the charged current as measured by different experiments over a large range of energies, for both neutrinos $\nu$ and antineutrinos $\bar{\nu}$ [48]. . . . .	15
2.10	3.57 keV emission line detected with a $4.5\sigma$ CL by the XMM-Newton telescope in 2014, which could be a hint of the presence of DM [52]. . . . .	16
2.11	Axions exclusion summary plot and projected coverage of axion searches experiments, such as ADMX, CAST and IAXO [58]. . . . .	17
2.12	Schematic view of the three main DM detection strategies: direct, indirect and collider production searches [61]. . . . .	18
2.13	Nuclear recoil spectra induced in different materials for a given DM WIMP of 100 GeV, assuming a WIMP-nucleon SI cross section [62]. . . . .	19

---

2.14 Schematic representation of the annual modulation of the WIMP wind introduced by the motion of rotation of the Earth around the Sun [64]. . . . .	20
2.15 Impact of different experimental parameters on the final limits depending on the cross section and WIMP mass (on the left) or sensitivity and exposure (on the right), with respect to the expected limits (black curve) [65]. . . . .	21
2.16 Schematic representation of the three main strategies to detect directly the interaction between DM particles and an ordinary nucleus [65]. . . . .	21
2.17 Exclusion limits obtained by various direct detection experiments considering a SI interaction cross section for low WIMP (on the left) or high WIMP masses (on the right) [65]. . . . .	22
2.18 Observed and expected annual modulation in single hits events in the 2-6 keV energy range by the DAMA experiment [66]. . . . .	22
2.19 Upper limits on the DM annihilation cross section considering $b\bar{b}$ (on the left) and $\mu^+\mu^-$ (on the right) final states as a function of the WIMP mass, for different clusters studied [69]. . . . .	24
2.20 Neutrino spectra for a scalar DM candidate of 1 TeV for different indirect detection experiments and the corresponding background level expected [71]. . . . .	25
2.21 Limits of the decay width of the interaction with respect to the DM mass obtained by both IceCube and AMS [72]. . . . .	26
2.22 Schematic representation of a typical EFT modelization of an LHC event with an ISR object used to trigger the event [61]. . . . .	27
2.23 Schematic representation of a typical collider DM production through s-channel and t-channel processes [74]. . . . .	28
2.24 Observed and expected 95% exclusion limits obtained by different searches of the CMS collaboration as a function the spin-0 scalar (on the left) or pseudoscalar (on the right) mediator. . . . .	29
2.25 Observed and expected 95% exclusion limits obtained by different searches of the CMS collaboration as a function of the spin-1 mediator, considering axial-vector (on the left) and axial (on the right) interactions. . . . .	30
2.26 CMS 90% exclusion limits compared to the most famous direct detection experiments for the SD (on the left) and SI (on the right) scenarios, obtained using similar couplings. . . . .	30
2.27 Feynman diagrams involving the production of DM with a single top quark according to its t-channel W boson (on the left), or tW (on the center and on the right) production modes. . . . .	31
2.28 Schematic representation of a typical $t\bar{t}$ +DM event. . . . .	31
2.29 Limits on the DM and mediator masses obtained by ATLAS using $13.3 \text{ fb}^{-1}$ of 13 TeV data, considering scalar (on the left) and pseudoscalar (on the right) mediators [17]. . . . .	32
2.30 Exclusion limits at the 95% CL obtained by ATLAS considering scalar (on the left) and pseudoscalar (on the right) mediators, for a DM mass of 1 GeV [20]. . . . .	33

---

2.31	Exclusion limits at the 95% CL obtained by ATLAS considering scalar (on the left) and pseudoscalar (on the right) mediators, for a DM mass of 1 GeV, considering the dilepton final state of such model and the full Run II dataset [20]. . . . .	33
2.32	Exclusion limits at the 95% CL obtained by ATLAS considering scalar (on the left) and pseudoscalar (on the right) mediators, for a DM mass of 1 GeV [21]. . . . .	34
2.33	95% CL exclusion plots on the signal strength computed as a function of the mediator and DM masses obtained by CMS considering a scalar (on the left) and a pseudoscalar (on the right) mediator for the interaction [23]. . . . .	34
2.34	Expected and observed 95% CL limits on the DM production cross sections shown considering scalar (on the left) and pseudoscalar (on the right) mediators for the interaction [24]. . . . .	35
3.1	LHC injection chain and experiments performed at CERN [93]. . . . .	38
3.2	Integrated luminosity collected by CMS over its different years of operation so far. . . . .	41
3.3	Mean PU distribution and luminosity recorded by CMS over the different years of operation of the LHC. . . . .	42
3.4	Schematic representation of the CMS detector, along with all its sub-detectors and main characteristics. . . . .	43
3.5	Schematic representation of the CMS coordinate system used by convention. . . . .	44
3.6	Schematic representation of the CMS tracker, for different pseudorapidity values and along the z-axis [97]. . . . .	45
3.7	Expected resolution of muons transverse momentum (left), transverse impact parameter (middle) and longitudinal impact parameter (right), as a function of pseudorapidity and muon momentum (1, 10 and 100 GeV) [96]. . . . .	45
3.8	Tracker reconstruction efficiency of muons (on the left) and pions (on the right) in simulation for different pseudorapidities and particle momenta (1, 10 and 100 GeV) [96]. . . . .	46
3.9	Schematic representation of the sub-systems of the CMS ECAL [96]. . . . .	46
3.10	Schematic representation of a typical electromagnetic shower and the radiation length $X_0$ concept [96]. . . . .	47
3.11	Schematic representation of the HCAL sub-system in CMS [96]. . . . .	48
3.12	Distribution of the measured energy scaled to the incident energy for pions with incident energies of 200 GeV at $\eta = 0$ (on the left) and 225 GeV at $\eta = 0.5$ (on the right), with and without the inclusion of the HO in the HCAL system [96]. . . . .	49
3.13	Picture of the solenoid system of CMS being setup in the assembly hall. . . . .	50
3.14	Geometrical repartition along the z-axis of the different muons chambers in CMS [98].	51
3.15	Lateral geometrical division of the different DT chambers in one of the 5 wheels of CMS [96]. . . . .	52
3.16	Muon reconstruction efficiency with different $p_T$ and $\eta$ values, considering only the muon system (on the left), and the combined information from the muon system and the tracker (on the right) [96]. . . . .	53

---

3.17	Location of the new GEM muon subsystem currently being installed in the very-forward region of CMS [96]. . . . .	54
3.18	Architecture of the L1 trigger of CMS [96]. . . . .	55
3.19	Architecture of the CMS DAQ system [96]. . . . .	56
4.1	Transverse section of CMS showing the different interactions expected by different kinds of particles in the detector. . . . .	59
4.2	Different kinds of vertices typically observed in a $pp$ collision in the LHC. . . . .	59
4.3	Muon $p_T$ resolution obtained in simulation in the barrel (on the left) and endcap (on the right) for different muon reconstruction algorithms [104]. . . . .	61
4.4	Lepton isolation cone typically used to enhance the prompt leptons purity. . . . .	62
4.5	Schematic representation of the full electron reconstruction workflow in CMS [106].	64
4.6	Schematic representation of the typical development of a jet within the CMS detector.	64
4.7	Comparison of the jet energy response (percentage of reconstructed energy, on the left) and jet energy resolution (on the right) for dijets simulated events in the barrel for jets reconstructed using only the calorimeters (in blue) and jet candidates from the PF algorithm (in red) [109]. . . . .	66
4.8	Schematic representation of the production of a b-jet originating from a slightly displaced secondary vertex. . . . .	66
4.9	b-jets identification efficiency and misidentification rate considering different b-taggers, including the deep CSV b-tag used in this analysis [110]. . . . .	67
4.10	Schematic representation of the MET. . . . .	68
4.11	PFMET (on the left) and PUPPI MET (on the right) distributions observed in a 2018 DY inclusive control region. . . . .	69
4.12	MET (on the left) and jet (on the right) $\phi$ distributions with and without MET filters applied [114]. . . . .	69
4.13	Three events constraining the neutrino and antineutrino momenta (black and grey arrows, respectively) resulting in 0 (on the left), 2 (on the center) or 4 (on the right) solutions. The dashed ellipses is obtained by using the additional constraint according to which measured MET is equal to the sum of neutrino transverse momenta [116]. . . . .	72
4.14	Schematic representation of the two extreme cases that can be observed when defining the overlapping factor: the reconstruction of a system with (on the left) and without (on the right) the presence of DM [117]. . . . .	73
4.15	Definition of the correction factor used in the smearing of the leptons energy [119].	74
4.16	Graphical definition of the two angles $\alpha$ and $\omega$ used in the smearing of the leptons and jets directions [119]. . . . .	74
4.17	Simulated distributions of the $\alpha$ angle between the particle and detector level direction for b-jets (on the left) and leptons (on the right). . . . .	75

---

4.18 Breit-Wigner spectrum obtained when generating randomly the mass of the W boson (on the left) and true $m_W$ distribution obtained using the generation information from the standard $t\bar{t}$ process (on the right). . . . .	75
5.1 Structure of a $pp$ collision and different steps of the MC simulation used by the event generators, such as the parton shower (in green), the UE (in pink), the hadronization (in blue) and the decay of unstable particles (in red) [120]. . . . .	78
5.2 Top $p_T$ (on the left) and rapidity (on the right) distributions obtained using different MC generators [129]. . . . .	79
5.3 Proton energy distribution at 3 (on the left) and 6 (on the right) GeV compared for the test beam data (in black) and two different GEANT4 versions [131]. . . . .	81
5.4 MET spectrum for several $t/\bar{t}+DM$ <b>scalar</b> mediators, with (on the left) and without normalization (on the right). . . . .	83
5.5 MET spectrum for several $t/\bar{t}+DM$ <b>pseudoscalar</b> mediators, with (on the left) and without normalization (on the right). . . . .	83
5.6 MET spectrum for several $t\bar{t}+DM$ <b>scalar</b> mediators (on the top) and dark matter (on the bottom) masses, with (on the left) and without normalization (on the right). . . . .	84
5.7 MET spectrum for several $t\bar{t}+DM$ <b>pseudoscalar</b> mediators (on the top) and dark matter (on the bottom) masses, with (on the left) and without normalization (on the right). . . . .	85
5.8 Production cross section of the most common SM processes considering different center of mass energies, such as the 13 TeV of the LHC. . . . .	86
5.9 Main feynman diagrams for the production of the SM $t\bar{t}$ process. . . . .	87
5.10 Feynman diagrams for the s-channel production mode of a single top quark. . . . .	87
5.11 Feynman diagrams for the t-channel (on the left) and tW (on the right) production modes of a single top quark. . . . .	87
5.12 Feynman diagrams for the leptonic decay of the top (on the left) and anti-top (on the right) quarks. . . . .	88
5.13 Feynman diagram for the DY process involving a virtual $\gamma^*$ or Z boson. . . . .	88
5.14 Normalized $m_{ll}$ DY distributions obtained using 2018 MC simulations in the 0 (on the left) and 1+ b-jet (on the right) bins. . . . .	89
5.15 $R_{out/in}$ , MC transfer factor obtained with the Rin-out data-driven method in bins of MET in 2016 (on the top left), 2017 (on the top right) and 2018 (on the bottom). . . . .	90
5.16 Possible Feynman diagrams for the ISR $t\bar{t}$ with a W/Z boson (on the left) and for the production of an FSR ttZ (on the center and right). . . . .	91
5.17 Possible Feynman diagrams for the production of a W/Z boson with a jet. . . . .	91
5.18 Schematic representation of the two jets used for the systematics and for the jet faking a lepton in the tight-to-loose data-driven method. . . . .	92
5.19 Electron (on the left) and muon (on the right) FR obtained in a QCD enriched region for different jet $E_T$ thresholds for 2016, 2017 and 2018 with respect to the $p_T$ of the lepton. . . . .	93

---

5.20 Electron (on the left) and muon (on the right) PR obtained in a Z+jets enriched region by a tag and probe method for 2016, 2017 and 2018 with respect to the $p_T$ of the lepton. . . . .	94
5.21 Possible Feynman diagrams for smaller backgrounds of this analysis: WW (on the left), W $\gamma$ and WZ (on the center) and ZZ (on the right). . . . .	96
5.22 Uncorrected PFMET (on the left) and PFMET after applying the EE noise and XY-shift corrections (on the right) distributions observed in a 2018 DY inclusive control region. . . . .	97
5.23 Data/MC agreement without (on the left) and with (on the right) top $p_T$ reweighting corrections in a general 2017 top control region. . . . .	97
5.24 Data/MC agreement without (on the left) and with (on the right) prefire corrections in a general 2018 top control region. . . . .	98
 6.1 DoubleEG trigger efficiencies with respect to the $p_T$ (on the left) and $\eta$ (on the right), computed using a tag and probe method, for the 2017 data taking period. .	101
6.2 Trigger efficiencies using orthogonal MET datasets for each dataset in 2016 (on the top left), 2017 (on the top right) and 2018 (on the bottom). . . . .	102
6.3 Two different variables (pfMET, on the left, and the stransverse mass $M_{T2}^{ll}$ , on the right) represented in the pre-selection region used for the training of the MVA. . .	104
6.4 Two different variables (pfMET, on the left, and the stransverse mass $M_{T2}^{ll}$ , on the right) represented in the $t/\bar{t}$ +DM signal region defined from the 100 GeV scalar training. . . . .	106
6.5 Two different variables (pfMET, on the left, and the stransverse mass $M_{T2}^{ll}$ , on the right) represented in the $t/\bar{t}$ +DM signal region defined from the 500 GeV scalar training. . . . .	107
6.6 Two different variables (pfMET, on the left, and the stransverse mass $M_{T2}^{ll}$ , on the right) represented in the $t/\bar{t}$ +DM signal region defined from the 100 GeV pseudoscalar training. . . . .	108
6.7 Two different variables (pfMET, on the left, and the stransverse mass $M_{T2}^{ll}$ , on the right) represented in the $t\bar{t}$ +DM signal region defined from the 100 GeV scalar training. . . . .	109
6.8 Two different variables (pfMET, on the left, and the stransverse mass $M_{T2}^{ll}$ , on the right) represented in the $t\bar{t}$ +DM signal region defined from the 500 GeV scalar training. . . . .	110
6.9 Two different variables (pfMET, on the left, and the stransverse mass $M_{T2}^{ll}$ , on the right) represented in the $t\bar{t}$ +DM signal region defined from the 100 GeV pseudoscalar training. . . . .	111
6.10 Two different variables ( $m_{ll}$ , on the left, and pfMET, on the right) represented in the inclusive control region defined. . . . .	112
6.11 Two different variables ( $m_{ll}$ , on the left, and pfMET, on the right) represented in the $t\bar{t}$ control region defined. . . . .	113
6.12 Two different variables ( $m_{ll}$ , on the left, and pfMET, on the right) represented in the DY control region defined. . . . .	115

---

6.13	Two different variables ( $m_{ll}$ , on the left, and pfMET, on the right) represented in the $t\bar{t}+V$ control region defined. . . . .	116
6.14	Two different variables ( $m_{ll}$ , on the left, and pfMET, on the right) represented in the same sign control region defined. . . . .	117
7.1	Number of b-jets (on the left) and $m_{bl}^t$ variable (on the right), mostly used to separate our two signals in this analysis, in a control region close to the actual signal region, for different mediator categories and signal samples of the analysis. . . . .	120
7.2	pfMET distribution in the 2018 pre-selection (on the left) and 2018 scalar 100 GeV signal (on the right) signal regions. . . . .	121
7.3	$M_{T2}^{ll}$ distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	122
7.4	$M_{T2}(bl, bl)$ distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	123
7.5	Dark $p_T$ and overlapping factor $R$ distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) regions (on the right) for the backgrounds and several signal samples. . . . .	124
7.6	$m_{bl}^t$ distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	125
7.7	Number of b-jets distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	125
7.8	Schematic representation in the $\Phi$ plane of the distribution of the particles for the $t\bar{t}$ process (on the left) and the for the $t\bar{t} + \text{DM}$ (on the center and right). . . . .	125
7.9	$\Delta\Phi(E_T^{\text{miss}}, ll)$ distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	126
7.10	$\Delta\Phi(lb^{\Delta R_{\min}}, E_T^{\text{miss}})$ distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	126
7.11	MET significance distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	127
7.12	Kinematic reconstruction weight $W$ distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	127
7.13	$r2l$ (on the top) and $r2l4j$ (on the bottom) variables considered for the signal discrimination process in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	128
7.14	Scalar sum of the transverse energy of the particles produced in the events distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	129
7.15	$\xi = \cos(\theta_l) \cos(\theta_{\bar{l}})$ distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	129
7.16	$c_{\text{hel}}$ distribution in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	130

---

7.17 Some kinematic variables computed in the $t\bar{t}$ system and considered for the signal discrimination process in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	131
7.18 Some kinematic variables computed in the $l\bar{l}b$ system and considered for the signal discrimination process in the 2018 pre-selection (on the left) and scalar 100 GeV signal (on the right) signal regions. . . . .	132
7.19 Schematic representation of a typical BDT, with its nodes (represented by the blue boxes) and leaves [159]. . . . .	133
7.20 Schematic representation of a typical neural network with a single hidden layer [160].	133
7.21 Correlation coefficient between the different variables used as input to the ANN for the backgrounds (on the left) and both signals (on the right). . . . .	136
7.22 Signal versus background ROC curves obtained after the training performed, considering 100 (on the left) and 500 GeV (on the right) mediators. . . . .	137
7.23 Signal versus background ROC curves obtained after the training performed with the ANN, considering all the scalar (on the left) and pseudoscalar (on the right) mediators available. . . . .	137
7.24 Overtraining distributions obtained for the 100 (on the left) and 500 GeV (on the right) mediators. . . . .	138
7.25 Input variables distributions used for the MVA. . . . .	138
7.26 Most likely category of the events corresponding to the different processes for the 100 (on the left) and 500 GeV (on the right) mediators. . . . .	139
7.27 ANN output distribution used in the shape analysis performed later on, for the scalar 100 GeV (on the left) and 500 GeV (on the right) training. . . . .	140
7.28 ANN output distribution used in the shape analysis performed later on, for the pseudoscalar 100 GeV (on the left) and 500 GeV (on the right) training. . . . .	141
 8.1 Schematic representation of the concept of PDF and p-value used in the computation of our test-statistic $q_\mu$ , under the hypotheses $\mu = 0$ and $\mu = 1$ . . . . .	146
8.2 Expected and observed upper limits on the signal strength $\mu$ for scalar (on the left) and pseudoscalar (on the right) models, considering with $m_\chi = 1$ GeV and $g_\chi = g_q = 1$ . . . . .	151
8.3 Expected and observed upper limits on the signal strength $\mu$ for scalar (on the left) and pseudoscalar (on the right) models, considering with $m_\chi = 1$ GeV and $g_\chi = g_q = 1$ , after combining the different Run II data taking periods. . . . .	152
 C.1 ROC curves obtained for the BDT by trying out different number of trees. . . . .	176
C.2 ROC curves obtained for the BDT by trying out different maximum depths. . . . .	177
C.3 ROC curves obtained for the BDT by trying out different $n_{cut}$ values. . . . .	178
C.4 ROC curves obtained for the BDT by trying out different shrinkage values. . . . .	179
C.5 ROC curves obtained for the ANN by trying out different architectures. . . . .	179
C.6 ROC curves obtained for the ANN by trying out different learning rates. . . . .	180

C.7	ROC curves obtained for the ANN by trying out different batch sizes. . . . .	180
C.8	Upper limits obtained for the scalar $t/\bar{t}$ +DM signal considering the optimal BDT and different sets of variables used as input. . . . .	182
C.9	Upper limits obtained for the scalar $t/\bar{t}$ +DM signal considering the optimal ANN and different sets of variables used as input. . . . .	183
D.1	2016 impacts and pulls obtained for the 30 most important systematics of the analysis, considering the 100 and 500 GeV scalar mediators. . . . .	186
D.2	2016 impacts and pulls obtained for the 30 most important systematics of the analysis, considering the 100 and 500 GeV pseudoscalar mediators. . . . .	187



---

## List of tables

3.1	Expected and observed main parameters of $pp$ operation of the LHC across the different eras of operation [92]. . . . .	39
3.2	Comparison of the three main sub-systems currently used by CMS in order to identify and measure muons [99]. . . . .	53
4.1	MET filters applied to events selected in data and to simulated events, an hyphen (-) indicating the filter is not applied. . . . .	69
5.1	Number of yields and percentage of different processes in some of the 2018 signal regions. . . . .	86
6.1	2016 trigger paths considered for this analysis. . . . .	100
6.2	2017 trigger paths considered for this analysis. . . . .	100
6.3	2018 trigger paths considered for this analysis. . . . .	101
7.1	Summary of the parameters used for the training of the BDT in this analysis. . . .	130
7.2	Summary of the parameters used for the training of the ANN in this analysis. . . .	134
7.3	Ranking for the importance of the different variables used as input to the ANN, considering the scalar 100 GeV training. . . . .	135
7.4	Ranking for the importance of the different variables used as input to the ANN, considering the scalar 500 GeV training. . . . .	136
B.1	Datasets collected in 2016 and considered for this analysis. . . . .	166
B.2	Datasets collected in 2017 and considered for this analysis. . . . .	167
B.3	Datasets collected in 2018 and considered for this analysis. . . . .	168
B.4	Signal samples mass points and LO dileptonic cross-sections considered for the $t/\bar{t}+DM$ signal used in this analysis. . . . .	169
B.5	Signal samples mass points and LO dileptonic cross-sections considered for the $t\bar{t}+DM$ signal used in this analysis. . . . .	170
B.6	Main 2016 MC simulations for the different background processes considered for this analysis and their respective cross sections. . . . .	171

B.7 Main 2017 MC simulations for the different background processes considered for this analysis and their respective cross sections. . . . .	172
B.8 Main 2018 MC simulations for the different background processes considered for this analysis and their respective cross sections. . . . .	173
C.1 Sets of input variables considered when optimizing the MVA method. . . . .	181

---

## Bibliography

- [1] G. Altarelli, "The Standard Model of Particle Physics", CERN-PH-TH/2005-206, 2005
- [2] F. Englert and R. Brout, "Broken symmetry and the mass of gauge vector mesons", Phys. Rev. Lett. 13, pp. 321-323, 1964
- [3] P. W. Higgs, "Broken symmetries and the masses of gauge bosons", Phys. Rev. Lett. 13, pp. 508-509, 1964
- [4] S. Chatrchyan et al., "Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC", Phys. Lett. B716, pp. 30-61, 2012 [arXiv: 1207.7235]
- [5] G. Aad et al., "Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC", Phys. Lett. B716, pp. 1-29, 2012 [arXiv: 1207.7214]
- [6] CMS Collaboration, "The CMS Experiment at the CERN LHC", JINST 3 S08004, 2008
- [7] ATLAS Collaboration, "The ATLAS Experiment at the CERN Large Hadron Collider", JINST 3 S08003, 2008
- [8] V.C. Rubin, W.K. Ford and N. Thonnard, "Rotational properties of 21 SC galaxies with a large range of luminosities and radii, from NGC 4605 (R=4kpc) to UGC 2885 (R=122kpc)", Astrophysical Journal 238, pp. 471-487, 1980
- [9] K.G. Begeman, A.H. Broeils and R.H. Sanders, "Extended rotation curves of spiral galaxies - Dark haloes and modified dynamics", Monthly Notices of the Royal Astronomical Society, vol. 249, issue 3, ISSN 0035-8711, 1991
- [10] A. Robertson, R. Massey and V. Eke, "What does the Bullet Cluster tell us about self-interacting dark matter?", Monthly Notices of the Royal Astronomical Society, vol. 465, issue 1, 2017 [arXiv: 1605.04307]
- [11] J.B. Muñoz, C. Dvorkin and A. Loeb, "21-cm Fluctuations from Charged Dark Matter", Phys. Rev. Lett. 121, 121301 (2018) [arXiv: 1804.01092]
- [12] A. Natarajan, "A closer look at CMB constraints on WIMP dark matter", Phys. Rev. D85, 2012 [arXiv:1201.3939 ]
- [13] G. D'Ambrosio G.F. Giudice, G. Isidori and A. Strumia, "Minimal Flavour Violation: an effective field theory approach", Nucl.Phys. 645, pp 155-187, 2002 [arXiv:0207.036 ]
- [14] CMS Collaboration, "Search for the production of dark matter in association with top-quark pairs in the single-lepton final state in proton-proton collisions at  $\sqrt{s} = 8$  TeV", JHEP, vol. 6 121, 2015
- [15] CMS Collaboration,, "Search for the Production of Dark Matter in Association with Top Quark Pairs in the Di-lepton Final State in pp collisions at  $\sqrt{s} = 8$  TeV", CMS-PAS-B2G-13-004, 2014

- [16] "Search for dark matter in events with heavy quarks and missing transverse momentum in pp collisions with the ATLAS detector", *Eur. Phys. J. C* (2015) 75:92
- [17] ATLAS Collaboration, Search for the Supersymmetric Partner of the Top Quark in the Jets+Emiss Final State at  $\sqrt{s} = 13$  TeV", *ATLAS-CONF-2016-077*
- [18] ATLAS Collaboration, "Search for top squarks in final states with one isolated lepton, jets, and missing transverse momentum in  $\sqrt{s} = 13$  TeV pp collisions with the ATLAS detector", *ATLAS-CONF-2016-050*, 2016
- [19] ATLAS Collaboration, "Search for direct top squark pair production and dark matter production in final states with two leptons in  $\sqrt{s} = 13$  TeV pp collisions using  $13.3 \text{ fb}^{-1}$  of ATLAS data", *ATLAS-CONF-2016-076*, 2016
- [20] ATLAS Collaboration, "Search for dark matter produced in association with bottom or top quarks in  $\sqrt{s} = 13$  TeV pp collisions with the ATLAS detector", *Eur. Phys. J. C* 78 (2018) 18 [arXiv: 1710.11412]
- [21] ATLAS Collaboration, "Search for new phenomena in events with two opposite-charge leptons, jets and missing transverse momentum in pp collisions at  $\sqrt{s} = 13$  TeV with the ATLAS detector", *ATLAS-CONF-2020-046*, July 2020
- [22] CMS Collaboration, Search for dark matter produced in association with heavy-flavor quark pairs in proton-proton collisions at  $\sqrt{s} = 13$  TeV", *Eur. Phys. J. C* (2017) 77: 845
- [23] CMS Collaboration, "Search for dark matter particles produced in association with a top quark pair at  $\sqrt{s} = 13$  TeV", *Phys. Rev. Lett.* 122, 011803 (2019) [arXiv: 1807.06522]
- [24] CMS Collaboration, "Search for dark matter produced in association with a single top quark or a top quark pair in proton-proton collisions at  $\sqrt{s} = 13$  TeV", *JHEP*, vol. 03 141, 2019 [arXiv: 1901.01553]
- [25] S. Manzoni, "The Standard Model and the Higgs Boson", *Physics with Photons Using the ATLAS Run 2 Data*, Springer Theses, 2019
- [26] A.B. Balantekin, A. Gouvea and B.Kayser, "Addressing the Majorana vs. Dirac Question with Neutrino Decays", *FERMILAB-PUB-18-418-T, NUHEP-TH/18-09* [arXiv: 1808.10518]
- [27] J. Woithe, G.J. Wiener and F. Van der Vecken, "Let's have a coffee with the Standard Model of particle physics!", *Physics education* 52, number 3, 2017
- [28] Y. Shadmi, "Introduction to Supersymmetry", *CERN Yellow Report CERN 2016-003*, pp. 95-123
- [29] H. Poincare, "The Milky Way and the Theory of Gases", *Popular Astronomy*, vol. 14, pp.475-488, 1906
- [30] F. Zwicky, "Die Rotverschiebung von extragalaktischen Nebeln", *Helvetica Physica Acta* , vol. 6, pp. 110-127, 1933
- [31] S. Van den Bergh, *Phys Rev D* "The early history of dark matter", Dominion Astrophysical Observatory, 1999
- [32] V.C. Rubin, W.K. Ford, "Rotation of the Andromeda Nebula from a Spectroscopic Survey of Emission Regions", *Astrophysical Journal* 159, p. 379, 1970
- [33] A. A. Penzias, R.W. Wilson, "A Measurement of Excess Antenna Temperature at 4080 Mc/s", *Astrophysical Journal* 142, pp. 419-421, 1965

- [34] D.J. Fixsen, "The temperature of the cosmic microwave background", *Astrophysical Journal*, 2009
- [35] Planck Collaboration, "Planck 2018 results. I. Overview and the cosmological legacy of Planck", 2018 [arXiv: 1807.06205]
- [36] R. Tojeiro, "Understanding the Cosmic Microwave Background Temperature Power Spectrum", 2006
- [37] Planck Collaboration, "Planck 2018 results. VI. Cosmological parameters", 2018 [arXiv: 1807.06209]
- [38] "Astrophysical Constants and Parameters", 2019
- [39] D. Clowe et all., "A Direct Empirical Proof of the Existence of Dark Matter", *Astrophysical Journal Letters* 648, 2006
- [40] L. Heurtier, H. Partouche, "Spontaneous Freeze Out of Dark Matter From an Early Thermal Phase Transition", CPHT-RR065.112019 [arXiv: 1912.02828]
- [41] K.R. Dienes, J. Fennick, J. Kumar, B. Thomas "Dynamical Dark Matter from Thermal Freeze-Out", *Phys. Rev. D* 97, 063522 (2018) [arXiv: 1712.09919]
- [42] C.S. Frenk, S.D.M. White, "Dark matter and cosmic structure", *Annalen der Physik*, p. 22 , 2012 [arXiv: 1210.0544]
- [43] R. Kirk, "Dark matter genesis", PhD Royal Holloway, U. of London, 2017
- [44] M. Drewes et all., "A White Paper on keV Sterile Neutrino Dark Matter", 2016 [arXiv: 1602.04816]
- [45] C. Alcock et all., "The MACHO Project: Microlensing Results from 5.7 Years of LMC Observations", *Astrophys.J.* 542 (2000) 281-307
- [46] P. Tisserand et all., "Limits on the Macho content of the Galactic Halo from the EROS-2 Survey of the Magellanic Clouds", *A & A* 469, pp. 387-404 (2007)
- [47] EROS and MACHO collaborations, "EROS and MACHO Combined Limits on Planetary Mass Dark Matter in the Galactic Halo", 1998
- [48] Particle Data Group, "Neutrino Cross Section Measurements", PDG 2019
- [49] K. McFarland, "Neutrino Interactions", 2008 [arXiv: 0804.3899]
- [50] E. Morgante, "Aspects of WIMP Dark Matter Searches at Colliders and Other Probes", Springer theses, 2016
- [51] F. Couchot et all., "Cosmological constraints on the neutrino mass including systematic uncertainties", *A & A* 606, A104 (2017)
- [52] E. Bulbul et all., "Detection of An Unidentified Emission Line in the Stacked X-ray spectrum of Galaxy Clusters", 2014 [arXiv: 1402.2301]
- [53] A. Boyarsky et all., "An unidentified line in X-ray spectra of the Andromeda galaxy and Perseus galaxy cluster", *Phys. Rev. Lett.* 113, 251301 (2014) [arXiv: 1402.4119]
- [54] A. Boyarsky et all., "Checking the dark matter origin of 3.53 keV line with the Milky Way center", *Phys. Rev. Lett.* 115, 161301 (2015) [arXiv: 1408.2503]
- [55] T. Jeltema1 and S.Profumo, "Deep XMM Observations of Draco rule out at the 99% Confidence Level a Dark Matter Decay Origin for the 3.5 keV Line", 2015 [arXiv: 1512.01239]

- [56] D. Wu, "A Brief Introduction to the Strong CP Problem", Superconducting Super Collider Laboratory, 1991
- [57] R.D. Peccei, H.R. Quinn, "CP Conservation in the Presence of Pseudoparticles", Phys. Rev. Lett. 38, 1440, 1977
- [58] P.W. Graham et all., "Experimental Searches for the Axion and Axion-like Particles", Annual Review of Nuclear and Particle Science 65, 2015 [arXiv: 1602.00039]
- [59] CAST collaboration, "New CAST limit on the axion-photon interaction", Nature Physics 13, pp. 584-590 (2017)
- [60] S.K. Vempati, "Introduction to MSSM", 2012 [arXiv: 1201.0334]
- [61] B. Penning, "The Pursuit of Dark Matter at Colliders - An Overview", 2017 [arXiv: 1712.01391]
- [62] M. Schumann, "Direct Detection of WIMP Dark Matter: Concepts and Status", J. Phys. G46 (2019) no.10, 103003 [arXiv: 1903.03026]
- [63] S.C. Martin et all., "The RAVE survey: constraining the local Galactic escape speed", Mon.Not.Roy.Astron.Soc.379:755-772, 2007
- [64] K. Freese, M. Lisanti, C. Savage, "Annual Modulation of Dark Matter: A Review", [arXiv: 1209.3339]
- [65] T.M. Undagoitia and L. Rauch, "Dark matter direct-detection experiments", J. Phys. G43 (2016) no.1, 013001 [arXiv: 1509.08767]
- [66] R. Bernabei et all., "First results from DAMA/LIBRA and the combined results with DAMA/NaI", Eur.Phys.J.C56:333-355, 2008 [arXiv: 0804.2741]
- [67] J.M. Gaskins, "A review of indirect searches for particle dark matter", Contemporary Physics, 2016 [arXiv: 1604.00014]
- [68] F.S. Queiroz, "Dark Matter Overview: Collider, Direct and Indirect Detection Searches", Max-Planck Institute of Physics
- [69] LAT collaboration, "Constraints on Dark Matter Annihilation in Clusters of Galaxies with the Fermi Large Area Telescope", JCAP 05(2010)025 [arXiv: 1002.2239]
- [70] A.A. Moiseev et all., "Dark Matter Search Perspectives with GAMMA-400", 2013 [arXiv: 1307.2345]
- [71] L. Covi et all., "Neutrino Signals from Dark Matter Decay", JCAP 1004:017, 2010 [arXiv: 0912.3521]
- [72] B. Lu and H. Zong, "Limits on the Dark Matter from AMS-02 antiproton and positron fraction data", Phys. Rev. D 93, 103517 (2016) [arXiv: 1510.04032]
- [73] J. Abdallah et all., "Simplified Models for Dark Matter Searches at the LHC", Phys. Dark Univ. 9-10 (2015) 8-23 [arXiv: 1506.03116]
- [74] H. An, L. Wang, H. Zhang, "Dark matter with t-channel mediator: a simple step beyond contact interaction", Phys. Rev. D 89, 115014 (2014) [arXiv: 1308.0592]
- [75] D. Abercrombie et all, "Dark Matter Benchmark Models for Early LHC Run-2 Searches: Report of the ATLAS/CMS Dark Matter Forum", Phys. Dark Univ. 26 (2019) 100371 [arXiv: 1507.00966]

- [76] ATLAS Collaboration, "Search for dark matter and other new phenomena in events with an energetic jet and large missing transverse momentum using the ATLAS detector", JHEP 01 (2018) 126 [arXiv: 1711.03301]
- [77] CMS Collaboration, "Search for new physics in the monophoton final state in proton-proton collisions at  $\sqrt{s} = 13$  TeV", J. High Energy Phys. 10 (2017) 073 [arXiv: 1706.03794]
- [78] CMS Collaboration, "Search for dark matter produced with an energetic jet or a hadronically decaying W or Z boson at  $\sqrt{s} = 13$  TeV", JHEP 07 (2017) 014 [arXiv: 1703.01651]
- [79] CMS Collaboration, "Search for new physics in final states with an energetic jet or a hadronically decaying W or Z boson and transverse momentum imbalance at  $\sqrt{s} = 13$  TeV", Phys. Rev. D 97, 092005 (2018) [arXiv: 1712.02345]
- [80] ATLAS Collaboration, "Search for dark matter in association with a Higgs boson decaying to two photons at  $\sqrt{s} = 13$  TeV with the ATLAS detector", Phys. Rev. D 96 (2017) 112004 [arXiv: 1706.03948]
- [81] CMS Collaboration, "Search for associated production of dark matter with a Higgs boson decaying to  $b\bar{b}$  or  $\gamma\gamma$  at  $\sqrt{s} = 13$  TeV", JHEP 10 (2017) 180 [arXiv: 1703.05236]
- [82] Atlas Collaboration, "Search for new phenomena in dijet events using 37 fb<sup>-1</sup> of pp collision data collected at  $\sqrt{s} = 13$  TeV with the ATLAS detector", Phys. Rev. D 96, 052004 (2017) [arXiv: 1703.09127]
- [83] CMS Collaboration, "Search for narrow and broad dijet resonances in proton-proton collisions at  $\sqrt{s} = 13$  TeV and constraints on dark matter mediators and other new particles", JHEP 08 (2018) 130 [arXiv: 1806.00843]
- [84] C. Munoz, "Models of Supersymmetry for Dark Matter", FTUAM 17/2, IFT-UAM/CSIC-17-005, 2017 [arXiv: 1701.05259]
- [85] CMS Collaboration, "Searches for invisible decays of the Higgs boson in pp collisions at  $\sqrt{s} = 7, 8,$  and  $13$  TeV", JHEP 02 (2017) 135 [arXiv: 1610.09218]
- [86] J. Alimena et all., "Searching for long-lived particles beyond the Standard Model at the Large Hadron Collider", 2019 [arXiv: 1903.04497]
- [87] A. Albert et all., "Recommendations of the LHC Dark Matter Working Group: Comparing LHC searches for heavy mediators of dark matter production in visible and invisible decay channels", 2017 [arXiv: 1703.05703]
- [88] M. Tanabashi et al., Particle Data Group, Phys. Rev. D98, 030001 (2018)
- [89] G. Giacomelli and R. Giacomelli, "The LEP legacy", 2005 [arXiv:hep-ex/0503050]
- [90] R. Schicker, "The ALICE detector at LHC", 2005 [arXiv:hep-ex/0509259]
- [91] LHCb Collaboration, "LHCb Detector Performance", Int. J. Mod. Phys. A 30, 1530022 (2015) [arXiv: 1412.6352]
- [92] J.T. Boyd, "LHC Run-2 and Future Prospects", 2020
- [93] E. Gschwendtner, "AWAKE, A Particle-driven Plasma Wakefield Acceleration Experiment", CERN Yellow Report CERN 2016-001, pp.271-288 [arXiv: 1705.10573]
- [94] M. Thomson, "Modern Particle Physics", Cambridge University Press, 2013
- [95] G. Apollinari et all., "High Luminosity Large Hadron Collider HL-LHC", CERN Yellow Report CERN-2015-005, pp.1-19 [arXiv: 1705.08830]

- [96] CMS Collaboration, "The CMS experiment at the CERN LHC", JINST 3 (2008) S08004
- [97] CMS Collaboration, "Precision measurement of the structure of the CMS inner tracking system using nuclear interactions", JINST 13 (2018) P10034 [arXiv: 1807.03289]
- [98] M.S. Kim, "CMS reconstruction improvement for the muon tracking by the RPC chambers", 2013 JINST 8 T03001 [arXiv: 1209.2646]
- [99] CMS Collaboration, "Performance of the CMS muon detector and muon reconstruction with proton-proton collisions at  $\sqrt{s} = 13$  TeV", JINST 13 (2018) P06015 [arXiv: 1804.04528]
- [100] CMS Collaboration, "CMS TriDAS project : Technical Design Report, Volume 1: The Trigger Systems", CERN-LHCC-2000-038, 2000
- [101] CMS Collaboration, "CMS The TriDAS Project : Technical Design Report, Volume 2: Data Acquisition and High-Level Trigger", CERN-LHCC-2002-026, 2002
- [102] CMS Collaboration, "Particle-Flow Event Reconstruction in CMS and Performance for Jets, Taus, and MET", CMS-PAS-PFT-09-001, 2009
- [103] CMS Collaboration, "Description and performance of track and primary-vertex reconstruction with the CMS tracker", JINST 9 (2014) P10009 [arXiv: 1405.6569]
- [104] V. Knunz, "Measurement of Quarkonium Polarization to Probe QCD at the LHC", Springer theses, 2015
- [105] CMS Collaboration, "Performance of electron reconstruction and selection with the CMS detector in proton-proton collisions at  $\sqrt{s} = 8$  TeV", JINST 10 (2015) P06005 [arXiv: 1502.02701]
- [106] J. Rembser, "CMS Electron and Photon Performance at 13 TeV", J. Phys. Conf. Ser. 1162 012008, 2019
- [107] P.L.S. Connor, "Review of jet reconstruction algorithms", Ryan Atkin J. Phys. Conf. Ser. 645 012008, 2015
- [108] Jet Energy Resolutions and Corrections Twiki, "Introduction to Jet Energy Corrections at CMS", <https://twiki.cern.ch/twiki/bin/view/CMS/IntroToJEC>, 2016
- [109] F. Beaudette, "The CMS Particle Flow Algorithm", 2014 [arXiv: 1401.8155]
- [110] CMS Collaboration, "Identification of heavy-flavour jets with the CMS detector in pp collisions at 13 TeV", JINST 13 (2018) P05011 [arXiv: 1712.07158]
- [111] CMS Collaboration, "Performance of missing transverse momentum reconstruction in proton-proton collisions at  $\sqrt{s} = 13$  TeV using the CMS detector", JINST 14 (2019) P07004 [arXiv: 1903.06078]
- [112] D. Bertolini et all., "Pileup Per Particle Identification", JHEP 1410 (2014) 59 [arXiv: 1407.6013]
- [113] CMS Collaboration, "MET Corrections and Uncertainties for Run-II", <https://twiki.cern.ch/twiki/bin/viewauth/CMS/MissingETRun2Corrections>, 2020
- [114] CMS Collaboration, "Performance of missing transverse momentum reconstruction in proton-proton collisions at  $\sqrt{s} = 13$  TeV using the CMS detector", JINST 14 (2019) P07004 [arXiv: 1903.06078]
- [115] L. Sonnenschein, "Analytical solution of ttbar dilepton equations", Phys.Rev.D73:054015, 2016

- [116] B.A. Betchart, R. Demina and A. Harel, "Analytic solutions for neutrino momenta in decay of top quarks", 2013 [arXiv: 1305.1878]
- [117] A. Lantero Barreda, "Improvement of the signal-to-background discrimination in search for Dark Matter produced in association with a pair of top-antitop quarks", UNICAN repository, 2018
- [118] CDF Collaboration, "Measurement of the Top Quark Mass using Template Methods on Dilepton Events in Proton-Antiproton Collisions at  $\sqrt{s} = 1.96$  TeV", Phys.Rev.D73:112006, 2006
- [119] CMS Collaboration, "Search for dark matter production in association with dileptonically decaying top quarks at 13 TeV", AN2016-240
- [120] M.H. Seymour and M. Marx, "Monte Carlo Event Generators", MCnet-13-05, 2013 [arXiv:1304.6677]
- [121] B. Cabouat, J.R. Gaunt and K. Ostrolenk, "A Monte-Carlo Simulation of Double Parton Scattering", JHEP11(2019)061 [arXiv: 1906.04669]
- [122] R. Placakyte, "Parton Distribution Functions", 2011 [arXiv:1111.5452]
- [123] J. Alwall et all., "MadGraph 5 : Going Beyond", 2011 [arXiv:1106.0522]
- [124] C. Oleari, "The POWHEG-BOX", Nucl.Phys.Proc.Suppl.205-206:36-41 [arXiv:1007.3893]
- [125] S. Frixione et all., "The MC@NLO 4.0 Event Generator", CERN-TH/2010-216 [arXiv: 1010.0819]
- [126] B. Webber, "Parton shower Monte Carlo event generators", Scholarpedia
- [127] M. Bahr et all., "Herwig++ Physics and Manual", Eur.Phys.J.C58:639-707, 2008 [arXiv: 0803.0883]
- [128] T. Sjostrand, "A Brief Introduction to PYTHIA 8.1" Comput.Phys.Commun.178:852-867, 2008 [arXiv: 0710.3820]
- [129] A. Karneyeu et all., "MCPlots: a particle physics resource based on volunteer computing", European Physical Journal C 74 (2014) [arXiv: 1306.3436]
- [130] V. Lefebure and S. Banerjee, "CMS Simulation Software Using Geant4", CMS-NOTE-1999-072, 1999
- [131] S. Banerjee, "Validation of Geant4 Physics Models Using Collision Data from the LHC", J. Phys.: Conf. Ser. 898 042005
- [132] A. Rizzi, G. Petrucciani and M. Peruzzi, "A further reduction in CMS event data for analysis: the NANOAOOD format", J. Phys.: Conf. Ser. 214 06021
- [133] U. Aisch and G. Polesello, "Searching for production of dark matter in association with top quarks at the LHC", High Energy Physics - Phenomenology, 2018 [arXiv: 1812.00694]
- [134] C.G. Lester and D.J. Summers, "Measuring masses of semi-invisibly decaying particles pair produced at hadron colliders", Phys.Lett.B463:99-103, 1999
- [135] C.G. Lester and B. Nachman, "Bisection-based asymmetric MT2 computation: a higher precision calculator than existing symmetric methods", JHEP03(2015)100, 1999 [arXiv: 1411.4312]
- [136] K. Bloom, "CMS software and computing for LHC Run 2", ICHEP 2016 [arXiv: 1611.03215]

- [137] W. Tanenbaum, "A ROOT/IO Based Software Framework for CMS", ECONF C0303241:TU KT010, 2003
- [138] CMS Collaboration, "CMS Luminosity Measurements for the 2016 Data Taking Period", CMS-PAS-LUM-17-001, 2017
- [139] CMS Collaboration, "CMS Luminosity Measurements for the 2017 Data Taking Period", CMS-PAS-LUM-17-001, 2018
- [140] CMS Collaboration, "CMS Luminosity Measurements for the 2018 Data Taking Period", CMS-PAS-LUM-17-001, 2019
- [141] NNPDF Collaboration, "Parton distributions for the LHC Run II", 10.1007/JHEP04(2015)040 [arXiv: 1410.8849]
- [142] CMS Collaboration, "Extraction and validation of a new set of CMS PYTHIA8 tunes from underlying-event measurements", Eur. Phys. J. C 80 (2020) 4 [arXiv: 1903.12179]
- [143] CMS Twiki, "Reweighting recipe to emulate Level 1 ECAL prefiring", <https://twiki.cern.ch/twiki/bin/view/CMS/L1ECALPrefiringWeightRecipe>, as seen in May 2020
- [144] CMS Twiki, "Official Prescription for calculating corrections and uncertainties on Missing Transverse Energy (MET)", [https://twiki.cern.ch/twiki/bin/viewauth/CMS/MissingETUncertaintyPrescription#Instructions\\_for\\_2017\\_data\\_with](https://twiki.cern.ch/twiki/bin/viewauth/CMS/MissingETUncertaintyPrescription#Instructions_for_2017_data_with), as seen in May 2021
- [145] J. Hirschauer, "CMS status report", LHCC Open Session, September 2018
- [146] CMS Collaboration, "Measurement of the differential cross sections for top quark pair production as a function of kinematic event variables in pp collisions at  $\sqrt{s} = 7$  and 8 TeV", Phys. Rev. D 94 (2016) 052006 [arXiv: 1607.00837]
- [147] CMS Twiki, "pt(top-quark) based reweighting of ttbar MC", <https://twiki.cern.ch/twiki/bin/viewauth/CMS/TopPtReweighting>, as seen in May 2020
- [148] CMS Twiki, "The modeling of the top quark  $p_T$ ", [https://twiki.cern.ch/twiki/bin/viewauth/CMS/TopPtReweighting#Case\\_3\\_1\\_Analyses\\_with\\_SM\\_tt\\_as](https://twiki.cern.ch/twiki/bin/viewauth/CMS/TopPtReweighting#Case_3_1_Analyses_with_SM_tt_as), as seen in February 2021
- [149] CMS Twiki, "Electron and Photon Physics Object Offline Guide", <https://twiki.cern.ch/twiki/bin/view/CMSPublic/SWGuideEgamma>, as seen in May 2020
- [150] CMS Twiki, "Cut Based Electron ID for Run 2", <https://twiki.cern.ch/twiki/bin/viewauth/CMS/CutBasedElectronIdentificationRun2>, as seen in May 2021
- [151] CMS Twiki, "Baseline muon selections for Run-II", <https://twiki.cern.ch/twiki/bin/viewauth/CMS/SWGuideMuonIdRun2>, as seen in May 2020
- [152] CMS Twiki, "Baseline muon selections for Run-II", <https://twiki.cern.ch/twiki/bin/viewauth/CMS/SWGuideMuonIdRun2>, as seen in May 2021
- [153] CMS Twiki, "Jet Identification", <https://twiki.cern.ch/twiki/bin/view/CMS/JetID>, as seen in May 2020
- [154] CMS Twiki, "Jet Identification", [https://twiki.cern.ch/twiki/bin/view/CMS/JetID#Recommendations\\_for\\_13\\_TeV\\_data](https://twiki.cern.ch/twiki/bin/view/CMS/JetID#Recommendations_for_13_TeV_data), as seen in May 2021
- [155] CMS Twiki, "Recommended Jet Energy Corrections and Uncertainties For Data and MC", <https://twiki.cern.ch/twiki/bin/view/CMS/JECDATAmc>, as seen in May 2021

- [156] CMS Twiki, "Jet Energy Resolution", [https://twiki.cern.ch/twiki/bin/view/CMS/JetResolution#JER\\_Scaling\\_factors\\_and\\_Uncertai](https://twiki.cern.ch/twiki/bin/view/CMS/JetResolution#JER_Scaling_factors_and_Uncertai), as seen in May 2021
- [157] CMS Twiki, "Heavy flavor identification at CMS with deep neural networks", <https://twiki.cern.ch/twiki/bin/view/CMSPublic/BTV13TeVDPDeepCSV>, as seen in May 2020
- [158] CMS Twiki, "MET significance", <https://twiki.cern.ch/twiki/bin/view/CMSPublic/SWGuideMETSignificance>, as seen in August 2020
- [159] K. Woodruff, "Introduction to boosted decision trees", Machine Learning Group Meeting, September 2017
- [160] K. Albertsson et all, "TMVA 4: Toolkit for Multivariate Data Analysis with ROOT, users guide", CERN-OPEN-2007-007
- [161] M. Jachowski, "Multivariate Analysis, TMVA, and Artificial Neural Networks", Michigan REU Final Presentations, August 2006
- [162] CMS Twiki, "LHC Higgs Combination Group (LHC-HCG)", <https://twiki.cern.ch/twiki/bin/view/LHCPhysics/HiggsCombination>, as seen in August 2020
- [163] G. Cowan et all, "Asymptotic formulae for likelihood-based tests of new physics", *Eur.Phys.J.C*71:1554, 2011 [arXiv: 1007.1727]
- [164] A.L. Read, "Presentation of search results: the CLs technique", *J. Phys. G: Nucl. Part. Phys.* 28 (2002) 2693-2704
- [165] A.L. Read, "1st Workshop on Confidence Limits", CERN-2000-005, pp81-101"
- [166] T. Junk, "Statistical Methods for Experimental Particle Physics", TRIUMF Summer Institute, July 2009
- [167] P.K. Sinervo, "Definition and Treatment of Systematic Uncertainties in High Energy Physics and Astrophysics", PHYSTAT2003, Stanford Linear Accelerator Center, September 2003
- [168] J. Buttwerworth et all., "PDF4LHC recommendations for LHC Run II", *J. Phys. G: Nucl. Part. Phys.* 43 023001, 2016 [arXiv: 1510.03865]
- [169] CMS collaboration., "Measurement of the inelastic proton-proton cross section at  $\sqrt{s} = 13$  TeV", *JHEP* 07 (2018) 161 [arXiv: 1802.02613]
- [170] CMS collaboration., "Identification of heavy-flavour jets with the CMS detector in pp collisions  $\sqrt{s} = 13$  TeV", *JINST* 13 (2018) P05011 [arXiv: 1712.07158]
- [171] CERN Twiki, "ATLAS-CMS recommended predictions for top-quark-pair cross sections using the Top++v2.0 program", <https://twiki.cern.ch/twiki/bin/view/LHCPhysics/TtbarNNLO>, as seen in August 2020
- [172] CERN Twiki, "ATLAS-CMS recommended predictions for single-top cross sections using the Hathor v2.1 program", <https://twiki.cern.ch/twiki/bin/view/LHCPhysics/SingleTopRefXsec>, as seen in August 2020
- [173] CMS Collaboration, "BRIL Work Suite", <https://cms-service-lumi.web.cern.ch/cms-service-lumi;brilwsdoc.html>, as seen in June 2020