

1 SEARCH FOR PRODUCTION OF A HIGGS BOSON AND A SINGLE TOP
2 QUARK IN MULTILEPTON FINAL STATES IN pp COLLISIONS AT $\sqrt{s} = 13$
3 TeV.

4 by

5 Jose Andres Monroy Montañez

A DISSERTATION

7 Presented to the Faculty of

8 The Graduate College at the University of Nebraska

9 In Partial Fulfilment of Requirements

10 For the Degree of Doctor of Philosophy

11 Major: Physics and Astronomy

12 Under the Supervision of Kenneth Bloom and Aaron Dominguez

13 Lincoln, Nebraska

14 July, 2018

15 SEARCH FOR PRODUCTION OF A HIGGS BOSON AND A SINGLE TOP
16 QUARK IN MULTILEPTON FINAL STATES IN pp COLLISIONS AT $\sqrt{s} = 13$
17 TeV.

18 Jose Andres Monroy Montañez, Ph.D.

19 University of Nebraska, 2018

20 Adviser: Kenneth Bloom and Aaron Dominguez

²¹ Table of Contents

²²	Table of Contents	iii
²³	List of Figures	vii
²⁴	List of Tables	xi
²⁵	1 Theoretical approach	1
²⁶	1.1 Introduction	1
²⁷	1.2 Standard model of particle physics	2
²⁸	1.2.1 Fermions	4
²⁹	1.2.1.1 Leptons	5
³⁰	1.2.1.2 Quarks	7
³¹	1.2.2 Fundamental interactions	11
³²	1.2.3 Gauge invariance.	15
³³	1.2.4 Gauge bosons	17
³⁴	1.3 Electroweak unification and the Higgs mechanism	18
³⁵	1.3.1 Spontaneous symmetry breaking (SSB)	26
³⁶	1.3.2 Higgs mechanism	30
³⁷	1.3.3 Masses of the gauge bosons	33
³⁸	1.3.4 Masses of the fermions	34

39	1.3.5	The Higgs field	35
40	1.3.6	Production of Higgs bosons at LHC	36
41	1.3.7	Higgs boson decay channels	40
42	1.4	Experimental status of the anomalous Higgs-fermion coupling	42
43	1.5	Associated production of a Higgs boson and a single top quark	44
44	1.6	CP-mixing in tH processes	49
45	2	The CMS experiment at the LHC	54
46	2.1	Introduction	54
47	2.2	The LHC	55
48	2.3	The CMS experiment	65
49	2.3.1	CMS coordinate system	68
50	2.3.2	Tracking system	70
51	2.3.3	Silicon strip tracker	73
52	2.3.4	Electromagnetic calorimeter	74
53	2.3.5	Hadronic calorimeter	76
54	2.3.6	Superconducting solenoid magnet	77
55	2.3.7	Muon system	79
56	2.3.8	CMS trigger system	80
57	2.3.9	CMS computing	81
58	3	Event generation, simulation and reconstruction	86
59	3.1	Event generation	87
60	3.2	Monte Carlo Event Generators.	90
61	3.3	CMS detector simulation.	92
62	3.4	Event reconstruction.	94
63	3.4.1	Particle-Flow Algorithm.	94

64	3.4.1.1	Missing transverse energy.	108
65	3.4.2	Event reconstruction examples	109
66	5	Statistical methods	112
67	5.1	Multivariate analysis	112
68	5.1.1	Decision trees	116
69	5.1.2	Boosted decision trees (BDT)	119
70	5.1.3	Overtraining	121
71	5.1.4	Variable ranking	122
72	5.1.5	BDT output example	122
73	5.2	Statistical inference	123
74	5.2.1	Nuisance parameters	124
75	5.2.2	Maximum likelihood estimation method	125
76	5.3	Upper limits	126
77	5.4	Asymptotic limits	130
78	6	Search for production of a Higgs boson and a single top quark in multilepton final states in pp collisions at $\sqrt{s} = 13$ TeV	132
80	6.1	Introduction	132
81	6.2	<i>tHq</i> signature	136
82	6.3	Background processes	137
83	6.4	Data and MC Samples	139
84	6.4.1	Full 2016 data set	139
85	6.4.2	Triggers	140
86	6.4.3	Signal modeling and MC samples	143
87	6.5	Object Identification	145
88	6.5.1	Lepton reconstruction and identification	146

89	6.5.2	Lepton selection efficiency	152
90	6.5.3	Jets and b -jet tagging	156
91	6.5.4	Missing Energy MET	157
92	6.6	Event selection	158
93	6.7	Background modeling and predictions	160
94	6.7.1	$t\bar{t}V$ and diboson backgrounds	161
95	6.7.2	Non-prompt and charge mis-ID backgrounds	163
96	6.8	Signal discrimination	168
97	6.8.1	MVA classifiers evaluation	168
98	6.8.2	Discriminating variables	169
99	6.8.3	BDTG classifiers response	174
100	6.8.4	Additional discriminating variables	176
101	6.8.5	Signal extraction procedure	177
102	6.8.6	Binning and selection optimization	179
103	6.9	Signal model	182
104	A Datasets and triggers		186
105	B BDTG aditional plots		190
106	B.1	BDTG input variables distributions for $2lss$ channel	190
107	B.2	Input variables distributions from BDTG classifiers	192
108	C Other binning strategies		196
109	D BDTG output variation with κ_V/κ_t		199
110	Bibliography		199
111	References		200

¹¹² List of Figures

113	1.1	Standard Model of particle physics.	3
114	1.2	Transformations between quarks	11
115	1.3	Fundamental interactions in nature.	12
116	1.4	SM interactions diagrams	13
117	1.5	Neutral current processes	19
118	1.6	Spontaneous symmetry breaking mechanism	27
119	1.7	SSB Potential form	28
120	1.8	Potential for complex scalar field	29
121	1.9	SSB mechanism for complex scalar field	30
122	1.10	Proton-Proton collision	36
123	1.11	Proton PDFs	37
124	1.12	Higgs boson production mechanism Feynman diagrams	38
125	1.13	Higgs boson production cross section and decay branching ratios	39
126	1.14	κ_t - κ_V plot of the coupling modifiers. ATLAS and CMS combination.	42
127	1.15	Higgs boson production in association with a top quark	45
128	1.16	Cross section for tHq process as a function of κ_t	48
129	1.17	Cross section for tHW process as a function of κ_{Htt}	48
130	1.18	NLO cross section for tX_0 and $t\bar{t}X_0$.	52

131	1.19 NLO cross section for $tWX_0, t\bar{t}X_0$.	53
132	2.1 CERN accelerator complex	55
133	2.2 LHC protons source. First acceleration stage.	56
134	2.3 The LINAC2 accelerating system at CERN.	57
135	2.4 LHC layout and RF cavities module.	58
136	2.5 LHC dipole magnet.	60
137	2.6 Integrated luminosity delivered by LHC and recorded by CMS during 2016	62
138	2.7 LHC interaction points	63
139	2.8 Multiple pp collision bunch crossing at CMS.	65
140	2.9 Layout of the CMS detector	66
141	2.10 CMS detector transverse slice	67
142	2.11 CMS detector coordinate system	69
143	2.12 CMS tracking system schematic view.	70
144	2.13 CMS pixel detector	71
145	2.14 SST Schematic view.	73
146	2.15 CMS ECAL schematic view	75
147	2.16 CMS HCAL schematic view	77
148	2.17 CMS solenoid magnet	78
149	2.18 CMS Muon system schematic view	79
150	2.19 CMS Level-1 trigger architecture	81
151	2.20 WLCG structure	82
152	2.21 Data flow from CMS detector through hardware Tiers	84
153	3.1 Event generation process.	87
154	3.2 Particle flow algorithm.	95
155	3.3 Stable cones identification	102

156	3.4	Jet reconstruction.	104
157	3.5	Jet energy corrections.	106
158	3.6	Secondary vertex in a b-hadron decay.	107
159	3.7	HIG-13-004 Event 1 reconstruction.	109
160	3.8	$e\mu$ event reconstruction.	110
161	3.9	Recorded event reconstruction.	111
162	5.1	Scatter plots-MVA event classification.	114
163	5.2	Scalar test statistical.	115
164	5.3	Decision tree.	116
165	5.4	Decision tree output example.	119
166	5.5	BDT output example.	122
167	5.6	t_r p.d.f. assuming each H_0 and H_1	128
168	5.7	Illustration of the CL_s limit.	129
169	5.8	Example of Brazilian flag plot	130
170	6.1	Analysis strategy workflow	135
171	6.2	tHq event signature	136
172	6.3	$t\bar{t}$ and $t\bar{t}W$ events signature	139
173	6.4	Trigger efficiency for the same-sign $\mu\mu$ category	141
174	6.5	Trigger efficiency for the $e\mu$ category	142
175	6.6	Trigger efficiency for the $3l$ category	143
176	6.7	tHq and tHW cross section in the κ_t - κ_V phase space	144
177	6.8	Tight vs loose lepton selection efficiencies in the $2lss$ channel.	154
178	6.9	Tight vs loose lepton selection efficiencies in the $3l$ channel.	155
179	6.10	Kinematic distributions in the diboson control region.	163
180	6.11	Input variables to the BDT for signal discrimination normalized.	166

181	6.12 MVA classifiers performance.	169
182	6.13 BDTG classifier Input variables distributions.	171
183	6.14 BDT input variables. Discrimination against $t\bar{t}$ and $t\bar{t}V$ in $3l$ channel. . .	172
184	6.15 Correlation matrices for the BDT input variables.	173
185	6.16 BDTG classifier response. Default parameters.	174
186	6.17 BDTG classifier output.	175
187	6.18 Additional discriminating variables distributions.	176
188	6.19 2D BDT classifier output planes	178
189	6.20 Binning overlaid on the S/B ratio map on the plane of classifier outputs. .	179
190	6.21 Binning combination scheme.	180
191	B.1 Distributions of input variables to the BDT for signal discrimination, two lepton same sign channel.	191
192	B.2 BDT input variables. Discrimination against $t\bar{t}$ in $2lss$ channel.	192
193	B.3 BDT input variables. Discrimination against $t\bar{t}V$ in $2lss$ channel.	193
194	B.4 BDT input variables. Discrimination against $t\bar{t}$ in $3l$ channel.	194
195	B.5 BDT input variables. Discrimination against $t\bar{t}V$ in $3l$ channel.	195
196	C.1 Binning by S/B regions for $2lss$ (left) and $3l$ (right).	196
197	C.2 Final bins (corresponding to S/B regions in the 2D plane)	197
198	C.3 Binning into geometric regions using a k -means algorithm.	198
199	C.4 Final bins using a k -means algorithm.	198
201	D.1 BDTG output variation with κ_V/κ_t	199

²⁰² List of Tables

203	1.1	Fermions of the SM.	4
204	1.2	Fermion masses.	5
205	1.3	Lepton properties.	7
206	1.4	Quark properties.	8
207	1.5	Fermion weak isospin and weak hypercharge multiplets.	9
208	1.6	Fundamental interactions features.	14
209	1.7	SM gauge bosons.	18
210	1.8	Higgs boson properties.	36
211	1.9	Predicted branching ratios for a SM Higgs boson with $m_H = 125 \text{ GeV}/c^2$	41
212	1.10	Predicted SM cross sections for tH production at $\sqrt{s} = 13 \text{ TeV}$	46
213	1.11	Predicted enhancement of the tHq and tHW cross sections at LHC	49
214	6.1	Trigger efficiency scale factors and associated uncertainties.	142
215	6.2	MC signal samples.	144
216	6.3	Effective areas, for electrons and muons.	150
217	6.4	Requirements on each of the three muon selections.	152
218	6.5	Criteria for each of the three electron selections.	153
219	6.6	Summary of event pre-selection.	158

220	6.7	Expected and observed yields for 35.9fb^{-1} after the selection in all final states. Uncertainties are statistical only.	159
221	6.8	Signal yields split by decay channels of the Higgs boson.	161
223	6.9	BDTG input variables.	170
224	6.10	Configuration used in the final BDTG training.	174
225	6.11	Input variables ranking for BDTG classifiers	175
226	6.12	ROC-integral for all the testing cases.	177
227	6.13	Selection cuts optimization.	180
228	6.14	Limit variation as a function of bin size.	181
229	6.15	Limit variation as a function of bin size in the same-sign dimuon channel. (In bold: the final bin borders used in the $2lss$ channel.)	181
231	6.16	The 33 distinct values of κ_t/κ_V and f_t as mapped by the 51 κ_t and κ_V points.	184
233	A.1	Full 2016 dataset.	186
234	A.2	HLT	187
235	A.3	κ_V and κ_t combinations.	188
236	A.4	List of background samples used in this analysis (CMSSW 80X).	189

²³⁷ **Chapter 1**

²³⁸ **Theoretical approach**

²³⁹ **1.1 Introduction**

240 The physical description of the universe is a challenge that physicists have faced by
241 making theories that refine existing principles and proposing new ones in an attempt
242 to embrace emerging facts and phenomena.

243 At the end of 1940s Julian Schwinger [1] and Richard P. Feynman [2], based on
244 the work of Sin-Itiro Tomonaga [3], developed an electromagnetic theory consistent
245 with special relativity and quantum mechanics that describes how matter and light
246 interact; the so-called *quantum electrodynamics* (QED) was born.

247 QED has become the blueprint for developing theories that describe the universe.
248 It was the first example of a quantum field theory (QFT), which is the theoretical
249 framework for building quantum mechanical models that describes particles and their
250 interactions. QFT is composed of a set of mathematical tools that combines classical
251 fields, special relativity and quantum mechanics, while keeping the quantum point
252 particles and locality ideas.

253 This chapter gives an overview of the standard model of particle physics, starting

254 with a description of the particles and their interactions, followed by a description of
 255 the electroweak interaction, the Higgs boson and the associated production of Higgs
 256 boson and a single top quark (tH). The description contained in this chapter is based
 257 on References [4–6].

258 1.2 Standard model of particle physics

259 The *standard model of particle physics (SM)* describes particle physics at the funda-
 260 mental level in terms of a collection of interacting particles and fields. The full picture
 261 of the SM is composed of three fields¹ whose excitations are interpreted as particles
 262 called mediators or force-carriers, a set of fields whose excitations are interpreted as
 263 elementary particles interacting through the exchange of those mediators, and a field
 264 that gives the mass to elementary particles. Figure 1.1 shows a scheme of the SM
 265 particles’ organization. In addition, for each of the particles in the scheme there exists
 266 an antiparticle with the same mass and opposite quantum numbers. The existence of
 267 antiparticles is a prediction of the relativistic quantum mechanics from the solution
 268 of the Dirac equation for which a negative energy solution is also possible. In some
 269 cases a particle is its own anti-particle, like photon or Higgs boson.

270 The mathematical formulation of the SM is based on group theory and the use of
 271 Noether’s theorem [8] which states that for a physical system modeled by a Lagrangian
 272 that is invariant under a group of transformations a conservation law is expected. For
 273 instance, a system described by a time-independent Lagrangian is invariant (symmet-
 274 ric) under time changes (transformations) with the total energy conservation law as
 275 the expected conservation law. In QED, the charge operator (Q) is the generator of

¹ The formal and complete treatment of the SM is out of the scope of this document, however a plenty of textbooks describing it at several levels are available in the literature. The treatment in References [5, 6] is quite comprehensive and detailed. Note that gravitational field is not included in the standard model formulation

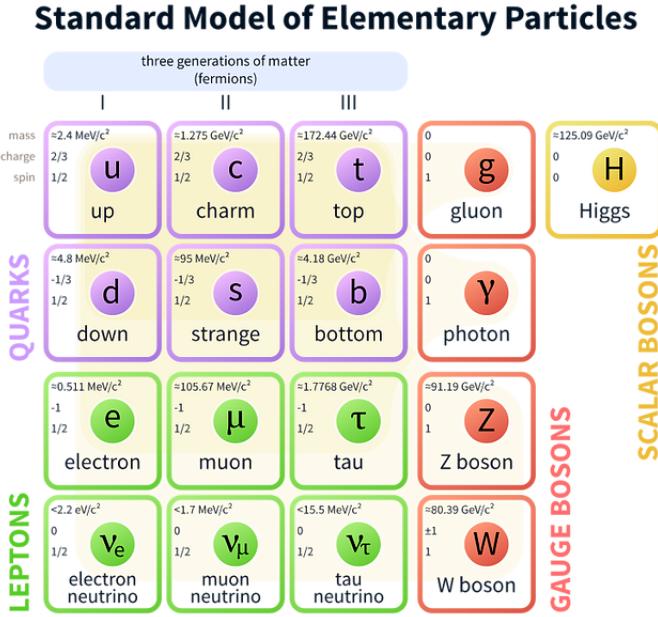


Figure 1.1: Schematic representation of the Standard Model of particle physics. The SM is a theoretical model intended to describe three of the four fundamental forces of the universe in terms of a set of particles and their interactions. [7].

the U(1) symmetry which according to the Noether's theorem means that there is a conserved quantity; this conserved quantity is the electric charge and thus the law of conservation of electric charge is established.

In the SM, the symmetry group $SU(3)_C \otimes SU(2)_L \otimes U(1)_Y$ describes three of the four fundamental interactions in nature (see Section 1.2.2): strong interaction (SI), weak interaction (WI) and electromagnetic interactions (EI) in terms of symmetries associated to physical quantities:

- 283 • Strong: $SU(3)_C$ associated to color charge

284 • Weak: $SU(2)_L$ associated to weak isospin and chirality

285 • Electromagnetic: $U(1)_Y$ associated to weak hypercharge and electric charge

286 It will be shown that the electromagnetic and weak interactions are combined in

287 the so-called electroweak interaction where chirality, hypercharge, weak isospin and
 288 electric charge are the central concepts.

289 **1.2.1 Fermions**

290 The basic constituents of the ordinary matter at the lowest level, which form the set
 291 of elementary particles in the SM formulation, are quarks and leptons. All of them
 292 have spin 1/2, therefore they are classified as fermions since they obey Fermi-Dirac
 293 statistics. There are six *flavors* of quarks and three of leptons organized in three
 294 generations, or families, as shown in Table 1.1.

		Generation		
		1st	2nd	3rd
Leptons	Type	Charged	Electron (e)	Moun(μ)
	Neutral	Neutral	Electron neutrino (ν_e)	Muon neutrino (ν_μ)
Quarks	Up-type	Up (u)	Charm (c)	Top (t)
	Down-type	Down (d)	Strange (s)	Bottom (b)

Table 1.1: Fermions of the SM. There are six flavors of quarks and three of leptons, organized in three generations, or families, composed of two pairs of closely related particles. The close relationship is motivated by the fact that each pair of particles is a member of an $SU(2)_L$ doublet that has an associated invariance under isospin transformations. WI between leptons is limited to the members of the same generation; WI between quarks is not limited but greatly favored, to same generation members.

295

296 There is a mass hierarchy between generations (see Table 1.2), where the higher
 297 generation particles decays to the lower one, which can explain why the ordinary
 298 matter is made of particles from the first generation. In the SM, neutrinos are modeled
 299 as massless particles so they are not subject to this mass hierarchy; however, today it
 300 is known that neutrinos are massive so the hierarchy could be restated. The reason
 301 behind this mass hierarchy is one of the most important open questions in particle
 302 physics, and it becomes more puzzling when noticing that the mass difference between

303 first and second generation fermions is small compared to the mass difference with
 304 respect to the third generation.

Lepton	Mass (MeV/c ²)	Quark	Mass (MeV/c ²)
e	0.51	u	2.2
μ	105.65	c	1.28×10^3
τ	1776.86	t	173.1×10^3
ν_e	Unknown	d	4.7
ν_μ	Unknown	s	96
τ_μ	Unknown	b	4.18×10^3

Table 1.2: Fermion masses [9]. Generations differ by mass in a way that has been interpreted as a mass hierarchy. Approximate values with no uncertainties are used, for comparison purpose.

305

306 Usually, the second and third generation fermions are produced in high energy
 307 processes, like the ones recreated in particle accelerators.

308 **1.2.1.1 Leptons**

309 A lepton is an elementary particle that is not subject to the SI. As seen in Table 1.1,
 310 there are two types of leptons, the charged ones (electron, muon and tau) and the
 311 neutral ones (the three neutrinos). The electric charge (Q) is the property that gives
 312 leptons the ability to participate in the EI. From the classical point of view, Q plays
 313 a central role determining, among others, the strength of the electric field through
 314 which the electromagnetic force is exerted. It is clear that neutrinos are not affected
 315 by EI because they don't carry electric charge.

316 Another feature of the leptons that is fundamental in the mathematical description
 317 of the SM is the chirality, which is closely related to spin and helicity. Helicity
 318 defines the handedness of a particle by relating its spin and momentum such that
 319 if they are parallel then the particle is right-handed; if spin and momentum are

320 antiparallel the particle is said to be left-handed. The study of parity conservation
 321 (or violation) in β -decay has shown that only left-handed electrons/neutrinos or right-
 322 handed positrons/anti-neutrinos are created [10]; the inclusion of that feature in the
 323 theory was achieved by using projection operators for helicity, however, helicity is
 324 frame dependent for massive particles which makes it not Lorentz invariant and then
 325 another related attribute has to be used: *chirality*.

326 Chirality is a purely quantum attribute which makes it not so easy to describe in
 327 graphical terms but it defines how the wave function of a particle transforms under
 328 certain rotations. As with helicity, there are two chiral states, left-handed chiral (L)
 329 and right-handed chiral (R). In the highly relativistic limit where $E \approx p \gg m$ helicity
 330 and chirality converge, becoming exactly the same for massless particles.

331 In the following, when referring to left-handed (right-handed) it will mean left-
 332 handed chiral (right-handed chiral). The fundamental fact about chirality is that
 333 while EI and SI are not sensitive to chirality, in WI left-handed and right-handed
 334 fermions are treated asymmetrically, such that only left-handed fermions and right-
 335 handed anti-fermions are allowed to couple to WI mediators, which is a violation of
 336 parity. The way to translate this statement in a formal mathematical formulation is
 337 based on the isospin symmetry group $SU(2)_L$.

338 Each generation of leptons is seen as a weak isospin doublet.² The left-handed
 339 charged lepton and its associated left-handed neutrino are arranged in doublets of
 340 weak isospin $T=1/2$ while their right-handed partners are singlets:

$$\begin{pmatrix} \nu_l \\ l \end{pmatrix}_L, l_R := \begin{pmatrix} \nu_e \\ e \end{pmatrix}_L, \begin{pmatrix} \nu_\mu \\ \mu \end{pmatrix}_L, \begin{pmatrix} \nu_\tau \\ \tau \end{pmatrix}_L, e_R, \mu_R, \tau_R, \nu_{eR}, \nu_{\mu R}, \nu_{\tau R} \quad (1.1)$$

341 The isospin third component refers to the eigenvalues of the weak isospin operator

² The weak isospin is an analogy of the isospin symmetry in strong interaction where neutron and proton are affected equally by strong force but differ in their charge.

342 which for doublets is $T_3 = \pm 1/2$, while for singlets it is $T_3 = 0$. The physical meaning
 343 of this doublet-singlet arrangement falls in that the WI couples the two particles in
 344 the doublet by exchanging the interaction mediator while the singlet member is not
 345 involved in WI. The main properties of the leptons are summarized in Table 1.3.

346 Although all three flavor neutrinos have been observed, their masses remain un-
 347 known and only some estimations have been made [11]. The main reason is that
 348 the flavor eigenstates are not the same as the mass eigenstates which implies that
 349 when a neutrino is created its mass state is a linear combination of the three mass
 350 eigenstates and experiments can only probe the squared difference of the masses. The
 351 Pontecorvo-Maki-Nakagawa-Sakata (PMNS) mixing matrix encodes the relationship
 352 between flavor and mass eigenstates.

Lepton	$Q(e)$	T_3	L_e	L_μ	L_τ	Lifetime (s)
Electron (e)	-1	-1/2	1	0	0	Stable
Electron neutrino(ν_e)	0	1/2	1	0	0	Unknown
Muon (μ)	-1	-1/2	0	1	0	2.19×10^{-6}
Muon neutrino (ν_μ)	0	1/2	0	1	0	Unknown
Tau (τ)	-1	-1/2	0	0	1	290.3×10^{-15}
Tau neutrino (ν_τ)	0	1/2	0	0	1	Unknown

Table 1.3: Lepton properties [9]. Q: electric charge, T_3 : weak isospin. Only left-handed leptons and right-handed anti-leptons participate in the WI. Anti-particles with inverted T_3 , Q and lepton number complete the leptons set but are not listed. Right-handed leptons and left-handed anti-leptons, neither listed, form weak isospin singlets with $T_3 = 0$ and do not take part in the weak interaction.

353

354 1.2.1.2 Quarks

355 Quarks are the basic constituents of protons and neutrons. The way quarks join to
 356 form bound states, called *hadrons*, is through the SI. Quarks are affected by all the
 357 fundamental interactions which means that they carry all the four types of charges:
 358 color, electric charge, weak isospin and mass.

Flavor	$Q(e)$	I_3	T_3	B	C	S	T	B'	Y	Color
Up (u)	2/3	1/2	1/2	1/3	0	0	0	0	1/3	r,b,g
Charm (c)	2/3	0	1/2	1/3	1	0	0	0	4/3	r,b,g
Top(t)	2/3	0	1/2	1/3	0	0	1	0	4/3	r,b,g
Down(d)	-1/3	-1/2	-1/2	1/3	0	0	0	0	1/3	r,b,g
Strange(s)	-1/3	0	-1/2	1/3	0	-1	0	0	-2/3	r,b,g
Bottom(b)	-1/3	0	-1/2	1/3	0	0	0	-1	-2/3	r,b,g

Table 1.4: Quark properties [9]. Q: electric charge, I_3 : isospin, T_3 : weak isospin, B: baryon number, C: charmness, S: strangeness, T: topness, B' : bottomness, Y: hypercharge. Anti-quarks posses the same mass and spin as quarks but all charges (color, flavor numbers) have opposite sign.

359

360 Table 1.4 summarizes the features of quarks, among which the most remarkable
 361 is their fractional electric charge. Note that fractional charge is not a problem, given
 362 that quarks are not found isolated, but serves to explain how composed particles are
 363 formed out of two or more valence quarks³.

364 Color charge is responsible for the SI between quarks and is the symmetry ($SU(3)_C$)
 365 that defines the formalism to describe SI. There are three colors: red (r), blue (b)
 366 and green (g) and their corresponding three anti-colors; thus each quark carries one
 367 color unit while anti-quarks carries one anti-color unit. As explained in Section 1.2.2,
 368 quarks are not allowed to be isolated due to the color confinement effect, hence, their
 369 features have been studied indirectly by observing their bound states created when

- 370 • one quark with a color charge is attracted by an anti-quark with the correspond-
 371 ing anti-color charge forming a colorless particle called a *meson*.

 372 • three quarks (anti-quarks) with different color (anti-color) charges are attracted
 373 among them forming a colorless particle called a *baryon (anti-baryon)*.

³ Hadrons can contain an indefinite number of virtual quarks and gluons, known as the quark and gluon sea, but only the valence quarks determine hadrons' quantum numbers.

374 In practice, when a quark is left alone isolated a process called *hadronization* occurs
 375 where the quark emits gluons (see Section 1.2.4) which eventually will generate new
 376 quark-antiquark pairs and so on; those quarks will recombine to form hadrons that
 377 will decay into leptons. This proliferation of particles looks like a *jet* coming from
 378 the isolated quark. More details about the hadronization process and jet structure
 379 will be given in chapter3.

380 In the first version of the quark model (1964), M. Gell-Mann [12] and G. Zweig
 381 [13, 14] developed a consistent way to classify hadrons according to their properties.
 382 Only three quarks (u, d, s) were involved in a scheme in which all baryons have baryon
 383 number $B=1$ and therefore quarks have $B=1/3$; non-baryons have $B=0$. Baryon
 384 number is conserved in SI and EI which means that single quarks cannot be created
 385 but in pairs $q - \bar{q}$.

386 The scheme organizes baryons in a two-dimensional space (I_3 - Y); Y (hyper-
 387 charge) and I_3 (isospin) are quantum numbers related by the Gell-Mann-Nishijima
 388 formula [15, 16]:

$$Q = I_3 + \frac{Y}{2} \quad (1.2)$$

389 where $Y = B + S + C + T + B'$ are the quantum numbers listed in Table 1.4.

390 There are six quark flavors organized in three generations (see Table 1.1) fol-
 391 lowing a mass hierarchy which, again, implies that higher generations decay to first
 392 generation quarks.

	Quarks			T_3	Y_W	Leptons			T_3	Y_W
Doublets	$(\frac{u}{d'})_L$	$(\frac{c}{s'})_L$	$(\frac{t}{b'})_L$	$(\frac{1/2}{-1/2})$	$1/3$	$(\nu_e)_L$	$(\nu_\mu)_L$	$(\nu_\tau)_L$	$(\frac{1/2}{-1/2})$	-1
Singlets	u_R	c_R	t_R	0	$4/3$	ν_{eR}	$\nu_{\mu R}$	$\nu_{\tau R}$		
	d'_R	s'_R	b'_R	0	$-2/3$	e_R	μ_R	τ_R	0	-2

Table 1.5: Fermion weak isospin and weak hypercharge multiplets. Weak hypercharge is calculated through the Gell-Mann-Nishijima formula 1.2 but using the weak isospin and charge for quarks.

393

394 Isospin doublets of quarks are also defined (see Table 1.5), and same as for neutrinos,
 395 the WI eigenstates are not the same as the mass eigenstates which means that
 396 members of different quark generations are connected by the WI mediator; thus, up-
 397 type quarks are coupled not to down-type quarks (the mass eigenstates) directly but
 398 to a superposition of down-type quarks (q'_d ; *the weak eigenstates*) via WI according
 399 to:

400

$$\begin{aligned} q'_d &= V_{CKM} q_d \\ \begin{pmatrix} d' \\ s' \\ b' \end{pmatrix} &= \begin{pmatrix} V_{ud} & V_{us} & V_{ub} \\ V_{cd} & V_{cs} & V_{cb} \\ V_{td} & V_{ts} & V_{tb} \end{pmatrix} \begin{pmatrix} d \\ s \\ b \end{pmatrix} \end{aligned} \quad (1.3)$$

401 where V_{CKM} is known as Cabibbo-Kobayashi-Maskawa (CKM) mixing matrix [17,18]
 402 given by

$$\begin{pmatrix} |V_{ud}| & |V_{us}| & |V_{ub}| \\ |V_{cd}| & |V_{cs}| & |V_{cb}| \\ |V_{td}| & |V_{ts}| & |V_{tb}| \end{pmatrix} = \begin{pmatrix} 0.97427 \pm 0.00015 & 0.22534 \pm 0.00065 & 0.00351^{+0.00015}_{-0.00014} \\ 0.22520 \pm 0.00065 & 0.97344 \pm 0.00016 & 0.0412^{+0.0011}_{-0.0005} \\ 0.00867^{+0.00029}_{-0.00031} & 0.0404^{+0.0011}_{-0.0005} & 0.999146^{+0.000021}_{-0.000046} \end{pmatrix}. \quad (1.4)$$

403 The weak decays of quarks are represented in the diagram of Figure 1.2; again
 404 the CKM matrix plays a central role since it contains the probabilities for the differ-
 405 ent quark decay channels, in particular, note that quark decays are greatly favored
 406 between generation members.

407 CKM matrix is a 3×3 unitary matrix parametrized by three mixing angles and
 408 the *CP-mixing phase*; the latter is the parameter responsible for the Charge-Parity

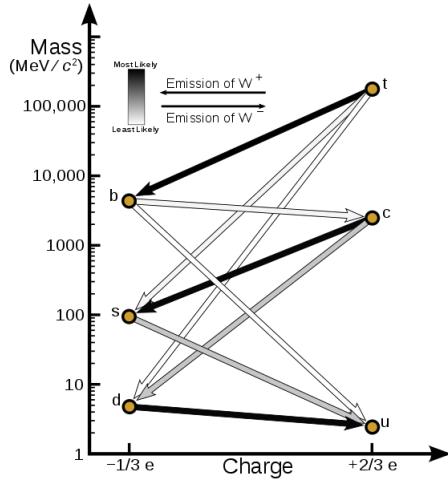


Figure 1.2: Transformations between quarks through the exchange of a WI. Higher generations quarks decay to first generation quarks by emitting a W boson. The arrow color indicates the likelihood of the transition according to the grey scale in the top left side which represent the CKM matrix parameters [19].

409 symmetry violation (CP-violation) in the SM. The fact that the top quark decays
 410 almost all the time to a bottom quark is exploited in this thesis when making the
 411 selection of the signal events by requiring the presence of a jet tagged as a jet coming
 412 from a b quark in the final state.

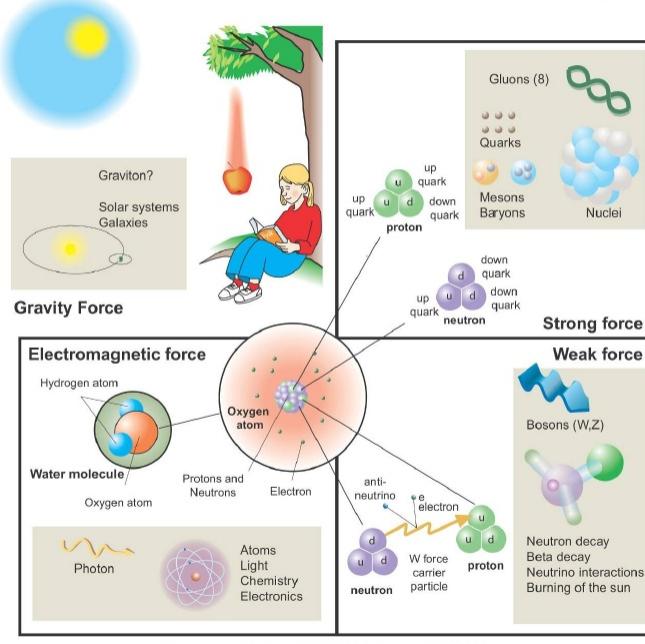
413 1.2.2 Fundamental interactions

414 Even though there are many manifestations of force in nature, like the ones repre-
 415 sented in Figure 1.3, we can classify all of them in four fundamental interactions:

- 416 • *Electromagnetic interaction (EI)* affects particles that are *electrically charged*,
 417 like electrons and protons. Figure 1.4a. shows a graphical representation, known
 418 as *Feynman diagram*, of electron-electron scattering.
- 419 • *Strong interaction (SI)* described by Quantum Chromodynamics (QCD). Hadrons
 420 like the proton and the neutron have internal structure given that they are com-

Fundamental interactions.

Illustration: Typoform



Summer School KPI 15 August 2009

Figure 1.3: Fundamental interactions in nature. Despite the many manifestations of forces in nature, we can track all of them back to one of the fundamental interactions. The most common forces are gravity and electromagnetic given that all of us are subject and experience them in everyday life.

421 posed of two or more valence quarks⁴. Quarks have fractional electric charge
 422 which means that they are subject to electromagnetic interaction and in the case
 423 of the proton they should break apart due to electrostatic repulsion; however,
 424 quarks are held together inside the hadrons against their electrostatic repulsion
 425 by the *Strong Force* through the exchange of *gluons*. The analog to the electric
 426 charge is the *color charge*. Electrons and photons are elementary particles as
 427 quarks but they don't carry color charge, therefore they are not subject to SI. A
 428 Feynman diagram for gluon exchange between quarks is shown in Figure 1.4b.

429 • *Weak interaction (WI)* described by the weak theory (WT), is responsible, for
 430 instance, for the radioactive decay in atoms and the deuterium production

⁴ Particles made of four and five quarks are exotic states not so common.

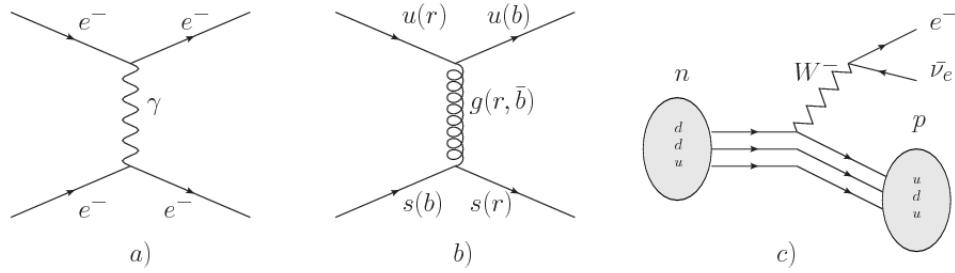


Figure 1.4: Feynman diagrams representing the interactions in SM; a) EI: e - e scattering; b) SI: gluon exchange between quarks ; c) WI: β -decay

within the sun. Quarks and leptons are the particles affected by the weak interaction; they possess a property called *flavor charge* (see 1.2.1) which can be changed by emitting or absorbing one weak force mediator. There are three mediators of the *weak force* known as Z boson in the case of electrically neutral flavor changes and W^\pm bosons in the case of electrically charged flavor changes. The *weak isospin* is the WI analog to electric charge in EI, and color charge in SI, and defines how quarks and leptons are affected by the weak force. Figure 1.4c. shows the Feynman diagram of β -decay where a neutron (n) is transformed in a proton (p) by emitting a W^- particle.

- *Gravitational interaction (GI)* described by General Theory of Relativity (GR). It is responsible for the structure of galaxies and black holes as well as the expansion of the universe. As a classical theory, in the sense that it can be formulated without even appeal to the concept of quantization, it implies that the space-time is a continuum and predictions can be made without limitation to the precision of the measurement tools. The latter represents a direct contradiction of the quantum mechanics principles. Gravity is deterministic while quantum mechanics is probabilistic; despite that, efforts to develop a quantum theory of gravity have predicted the *graviton* as mediator of the gravitational

449 force⁵.

Interaction	Acts on	Relative strength	Range (m)	Mediators
Electromagnetic (QED)	Electrically charged particles	10^{-2}	Infinite	Photon
Strong (QCD)	Quarks and gluons	1	10^{-15}	Gluon
Weak (WI)	Leptons and quarks	10^{-6}	10^{-18}	W^\pm, Z
Gravitational (GI)	Massive particles	10^{-39}	Infinite	Graviton

Table 1.6: Fundamental interactions features [20].

450

451 Table 1.6 summarizes the main features of the fundamental interactions. The
 452 strength of the interactions is represented by the coupling constants which depend
 453 on the energy scale at which the interaction is evaluated, therefore, it is the relative
 454 strength of the fundamental forces that reveals the meaning of strong and weak; in
 455 a context where the relative strength of the SI is 1, the EI is about hundred times
 456 weaker and WI is about million times weaker than the SI. A good description on how
 457 the relative strength and range of the fundamental interactions are calculated can
 458 be found in References [20, 21]. In the everyday life, only EI and GI are explicitly
 459 experienced due to the range of these interactions; i.e., at the human scale distances
 460 only EI and GI have appreciable effects, in contrast to SI which at distances greater
 461 than 10^{-15} m become negligible. Is it important to clarify that the weakness of the
 462 WI is attributed to the fact that its mediators are highly massive which affects the
 463 propagators of the interaction, as a result, the effect of the coupling constant is
 464 reduced.

⁵ Actually a wide variety of theories have been developed in an attempt to describe gravity; some famous examples are string theory and supergravity.

465 **1.2.3 Gauge invariance.**

466 QED was built successfully on the basis of the classical electrodynamics theory (CED)
 467 of Maxwell and Lorentz, following theoretical and experimental requirements imposed
 468 by

- 469 • Lorentz invariance: independence on the reference frame.
- 470 • Locality: interacting fields are evaluated at the same space-time point to avoid
 action at a distance.
- 472 • Renormalizability: physical predictions are finite and well defined.
- 473 • Particle spectrum, symmetries and conservation laws already known must emerge
 from the theory.
- 475 • Local gauge invariance.

476 The gauge invariance requirement reflects the fact that the fundamental fields
 477 cannot be directly measured but associated fields which are the observables. Electric
 478 (**E**) and magnetic (**B**) fields in CED are associated with the electric scalar potential
 479 V and the vector potential **A**. In particular, **E** can be obtained by measuring the
 480 change in the space of the scalar potential (ΔV); however, two scalar potentials
 481 differing by a constant f correspond to the same electric field. The same happens
 482 in the case of the vector potential **A**; thus, different configurations of the associated
 483 fields result in the same set of values of the observables. The freedom in choosing one
 484 particular configuration is known as *gauge freedom*; the transformation law connecting
 485 two configurations is known as *gauge transformation* and the fact that the observables
 486 are not affected by a gauge transformation is called *gauge invariance*.

487 When the gauge transformation:

$$\begin{aligned} \mathbf{A} &\rightarrow \mathbf{A} - \Delta f \\ V &\rightarrow V - \frac{\partial f}{\partial t} \end{aligned} \tag{1.5}$$

488 is applied to Maxwell equations, they are still satisfied and the fields remain invariant.
 489 Thus, CED is invariant under gauge transformations and is called a *gauge theory*.
 490 The set of all gauge transformations form the *symmetry group* of the theory, which
 491 according to the group theory, has a set of *group generators*. The number of group
 492 generators determine the number of *gauge fields* of the theory.

493 As mentioned in the first lines of Section 1.2, QED has one symmetry group ($U(1)$)
 494 with one group generator (the Q operator) and one gauge field (the electromagnetic
 495 field A^μ). In CED there is not a clear definition, beyond the historical convention,
 496 of which fields are the fundamental and which are the associated, but in QED the
 497 fundamental field is A^μ . When a gauge theory is quantized, the gauge fields are
 498 quantized and their quanta are called *gauge bosons*. The word boson characterizes
 499 particles with integer spin which obey Bose-Einstein statistics.

500 As will be detailed in Section 1.3, interactions between particles in a system can
 501 be obtained by considering first the Lagrangian density of free particles in the sys-
 502 tem, which of course is incomplete because the interaction terms have been left out,
 503 and demanding global phase transformation invariance. Global phase transforma-
 504 tion means that a gauge transformation is performed identically to every point
 505 in the space⁶ and the Lagrangian remains invariant. Then, the global transforma-
 506 tion is promoted to a local phase transformation (this time the gauge transforma-
 507 tion depends on the position in space) and again invariance is required.

⁶ Here space corresponds to the 4-dimensional space i.e. space-time.

508 Due to the space dependence of the local transformation, the Lagrangian density is
 509 not invariant anymore. In order to reinstate the gauge invariance, the gauge covariant
 510 derivative is introduced in the Lagrangian and with it the gauge field responsible for
 511 the interaction between particles in the system. The new Lagrangian density is gauge
 512 invariant, includes the interaction terms needed to account for the interactions and
 513 provides a way to explain the interaction between particles through the exchange of
 514 the gauge boson.

515 This recipe was used to build QED and the theories that aim to explain the
 516 fundamental interactions.

517 1.2.4 Gauge bosons

518 The importance of the gauge bosons comes from the fact that they are the force
 519 mediators or force carriers. The features of the gauge bosons reflect those of the fields
 520 they represent and they are extracted from the Lagrangian density used to describe
 521 the interactions. In Section 1.3, it will be shown how the gauge bosons of the EI and
 522 WI emerge from the electroweak Lagrangian. The SI gauge bosons features are also
 523 extracted from the SI Lagrangian but it is not detailed in this document. The main
 524 features of the SM gauge bosons will be briefly presented below and summarized in
 525 Table 1.7.

- 526 • **Photon.** EI occurs when the photon couples to (is exchanged between) parti-
 527 cles carrying electric charge; however, The photon itself does not carry electric
 528 charge, therefore, there is no coupling between photons. Given that the photon
 529 is massless the EI is of infinite range, i.e., electrically charged particles interact
 530 even if they are located far away one from each other; this also implies that
 531 photons always move with the speed of light.

- 532 • **Gluon.** SI is mediated by gluons which just as photons are massless. They
 533 carry one unit of color charge and one unit of anticolor charge, hence, gluons
 534 can couple to other gluons. As a result, the range of the SI is not infinite
 535 but very short due to the attraction between gluons, giving rise to the *color*
 536 *confinement* which explains why color charged particles cannot be isolated but
 537 live within composite particles, like quarks inside protons.
- 538 • **W, Z.** W^\pm and Z, are massive which explains their short-range. Given that
 539 the WI is the only interaction that can change the flavor of the interacting
 540 particles, the W boson is the responsible for the nuclear transmutation where
 541 a neutron is converted into a proton or vice versa with the involvement of an
 542 electron and a neutrino (see Figure 1.4c). The Z boson is the responsible for the
 543 neutral weak processes like neutrino elastic scattering where no electric charge
 544 but momentum transference is involved. WI gauge bosons carry isospin charge
 545 which makes interaction between them possible.

Interaction	Mediator	Electric charge (e)	Color charge	Weak Isospin	mass (GeV/c ²)
Electromagnetic	Photon (γ)	0	No	0	0
Strong	Gluon (g)	0	Yes -octet	No	0
Weak	W^\pm	± 1	No	± 1	80.385 ± 0.015
	Z	0	No	0	91.188 ± 0.002

Table 1.7: SM gauge bosons main features [9].

546

547 1.3 Electroweak unification and the Higgs 548 mechanism

549 Physicists dream of building a theory that contains all the interactions in one single
 550 interaction, i.e., showing that at some scale in energy all the four fundamental inter-

actions are unified and only one interaction emerges in a *Theory of everything*. The first sign of the feasibility of such unification came from success in the construction of the CED. Einstein spent years trying to reach that full unification, which by 1920 only involved electromagnetism and gravity, with no success; however, a new partial unification was achieved in the 1960's, when S.Glashow [22], A.Salam [23] and S.Weinberg [24] independently proposed that electromagnetic and weak interactions are two manifestations of a more general interaction called *electroweak interaction* (EWI). EWI was developed by following the useful prescription provided by QED and the gauge invariance principles.

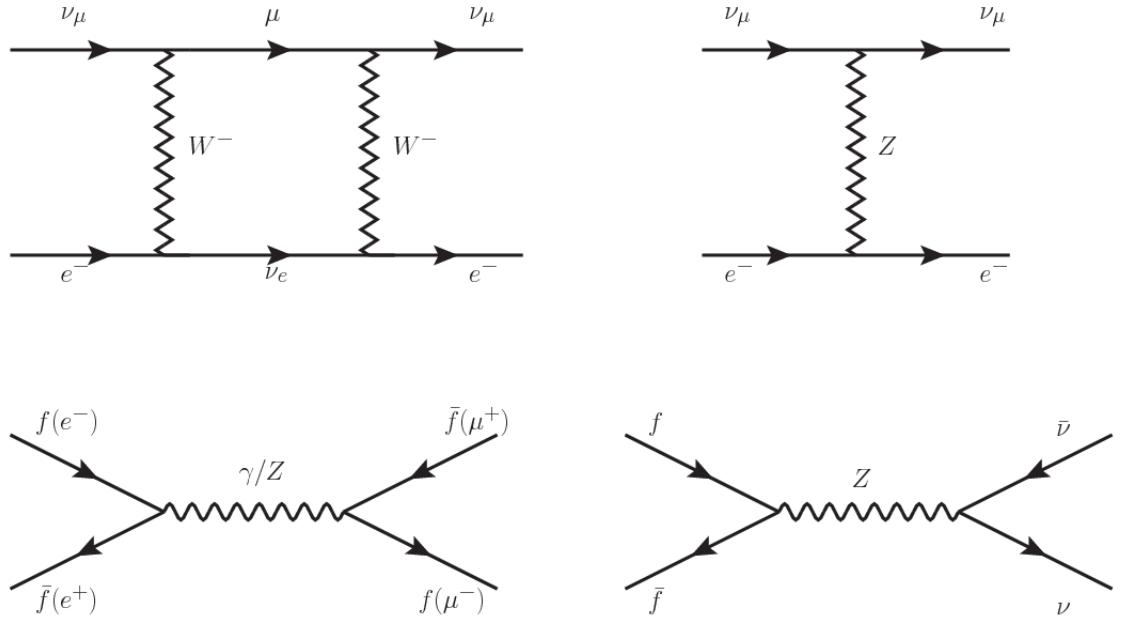


Figure 1.5: Top: $\nu_\mu - e^-$ scattering going through charged currents (left) and neutral currents (right). Bottom: neutral current processes for charged fermions (left) and involving neutrinos (right). While neutral current processes involving only charged fermions can proceed through EI or WI, those involving neutrinos can only proceed via WI.

The *classic* weak theory developed by Fermi, did not have the concept of the W boson but instead it was treated as a point interaction with the dimensionful constant G_F associated with it. It works really well at low energies very far off the W mass

563 shell. When going up in energy, the theory of weak interactions involving the W
 564 boson is capable of explaining the β -decay and in general the processes mediated by
 565 W^\pm bosons. However, there were some processes like the $\nu_\mu - e$ scattering which
 566 would require the exchange of two W bosons (see Figure 1.5 top diagrams) giving
 567 rise to divergent loop integrals and then non-finite predictions. The EWI theory, by
 568 including neutral currents involving fermions via the exchange of a neutral bosons Z,
 569 overcomes those divergences and the predictions become realistic.

570 Neutral weak interaction vertices conserve flavor in the same way as the electro-
 571 magnetic vertices do, but additionally, the Z boson can couple to neutrinos which
 572 implies that processes involving charged fermions can proceed through EI or WI but
 573 processes involving neutrinos can proceed only through WI.

574 The prescription to build a gauge theory of the WI consists of proposing a free
 575 field Lagrangian density that includes the particles involved; next, by requesting
 576 invariance under global phase transformations first and generalizing to local phase
 577 transformations invariance later, the conserved currents are identified and interactions
 578 are generated by introducing gauge fields. Given that the goal is to include the EI
 579 and WI in a single theory, the group symmetry considered should be a combination of
 580 $SU(2)_L$ and $U(1)_{em}$, however the latter cannot be used directly because the EI treats
 581 left and right-handed particles indistinctly in contrast to the former. Fortunately, the
 582 weak hypercharge, which is a combination of the weak isospin and the electric charge
 583 (Eqn. 1.2) is suitable to be used since it is conserved by the EI and WI. Thus, the
 584 symmetry group to be considered is

$$G \equiv SU(2)_L \otimes U(1)_Y \quad (1.6)$$

585 The following treatment applies to any of the fermion generations, but for sim-

586 plicity the first generation of leptons will be considered [5, 6, 25, 26].

587 Given the first generation of leptons

$$\psi_1 = \begin{pmatrix} \nu_e \\ e^- \end{pmatrix}_L, \quad \psi_2 = \nu_{eR}, \quad \psi_3 = e_R^- \quad (1.7)$$

588 the charged fermionic currents are given by

$$J_\mu \equiv J_\mu^+ = \bar{\nu}_{eL} \gamma_\mu e_L, \quad J_\mu^\dagger \equiv J_\mu^- = \bar{e}_L \gamma_\mu \nu_{eL} \quad (1.8)$$

589 and the free Lagrangian is given by

$$\mathcal{L}_0 = \sum_{j=1}^3 i\bar{\psi}_j(x) \gamma^\mu \partial_\mu \psi_j(x). \quad (1.9)$$

590 Mass terms are included directly in the QED free Lagrangians since they preserve
 591 the invariance under the symmetry transformations involved which treat left and right
 592 handed particles similarly, however mass terms of the form

$$m_W^2 W_\mu^\dagger(x) W^\mu(x) + \frac{1}{2} m_Z^2 Z_\mu(x) Z^\mu(x) - m_e \bar{\psi}_e(x) \psi_e(x) \quad (1.10)$$

593 which represent the mass of W^\pm , Z and electrons, are not invariant under G trans-
 594 formations, therefore the gauge fields described by the EWI are in principle massless.

595 Experiments have shown that the EWI gauge fields are not massless [27–30];
 596 however, they have to acquire mass through a mechanism compatible with the gauge
 597 invariance; that mechanism is known as the *Higgs mechanism* and will be considered
 598 later in this Section. The global transformations in the combined symmetry group G
 599 can be written as

$$\begin{aligned}
\psi_1(x) &\xrightarrow{G} \psi'_1(x) \equiv U_Y U_L \psi_1(x), \\
\psi_2(x) &\xrightarrow{G} \psi'_2(x) \equiv U_Y \psi_2(x), \\
\psi_3(x) &\xrightarrow{G} \psi'_3(x) \equiv U_Y \psi_3(x)
\end{aligned} \tag{1.11}$$

600 where U_L represent the $SU(2)_L$ transformation acting only on the weak isospin dou-
601 blet and U_Y represent the $U(1)_Y$ transformation acting on all the weak isospin mul-
602 tiplets. Explicitly

$$U_L \equiv \exp\left(i \frac{\sigma_i}{2} \alpha^i\right), \quad U_Y \equiv \exp(i y_i \beta) \quad (i = 1, 2, 3) \tag{1.12}$$

603 with σ_i the Pauli matrices and y_i the weak hypercharges. In order to promote the
604 transformations from global to local while keeping the invariance, it is required that
605 $\alpha^i = \alpha^i(x)$, $\beta = \beta(x)$ and the replacement of the ordinary derivatives by the covariant
606 derivatives

$$\begin{aligned}
D_\mu \psi_1(x) &\equiv \left[\partial_\mu + ig\sigma_i W_\mu^i(x)/2 + ig'y_1 B_\mu(x) \right] \psi_1(x) \\
D_\mu \psi_2(x) &\equiv \left[\partial_\mu + ig'y_2 B_\mu(x) \right] \psi_2(x) \\
D_\mu \psi_3(x) &\equiv \left[\partial_\mu + ig'y_3 B_\mu(x) \right] \psi_3(x)
\end{aligned} \tag{1.13}$$

607 introducing in this way four gauge fields, $W_\mu^i(x)$ and $B_\mu(x)$, in the process. The
608 covariant derivatives (Eqn. 1.13) are required to transform in the same way as fermion
609 fields $\psi_i(x)$ themselves, therefore, the gauge fields transform as:

$$\begin{aligned} B_\mu(x) &\xrightarrow{G} B'_\mu(x) \equiv B_\mu(x) - \frac{1}{g'} \partial_\mu \beta(x) \\ W_\mu^i(x) &\xrightarrow{G} W_\mu^{i\prime}(x) \equiv W_\mu^i(x) - \frac{i}{g} \partial_\mu \alpha_i(x) - \varepsilon_{ijk} \alpha_i(x) W_\mu^j(x). \end{aligned} \quad (1.14)$$

610 The G invariant version of the Lagrangian density 1.9 can be written as

$$\mathcal{L}_0 = \sum_{j=1}^3 i \bar{\psi}_j(x) \gamma^\mu D_\mu \psi_j(x) \quad (1.15)$$

611 where free massless fermion and gauge fields and fermion-gauge boson interactions
 612 are included. The EWI Lagrangian density must additionally include kinetic terms
 613 for the gauge fields (\mathcal{L}_G) which are built from the field strengths, according to

$$B_{\mu\nu}(x) \equiv \partial_\mu B_\nu - \partial_\nu B_\mu \quad (1.16)$$

$$W_{\mu\nu}^i(x) \equiv \partial_\mu W_\nu^i(x) - \partial_\nu W_\mu^i(x) - g \varepsilon^{ijk} W_\mu^j W_\nu^k \quad (1.17)$$

614 the last term in Eqn. 1.17 is added in order to hold the gauge invariance; therefore,

$$\mathcal{L}_G = -\frac{1}{4} B_{\mu\nu}(x) B^{\mu\nu}(x) - \frac{1}{4} W_{\mu\nu}^i(x) W_i^{\mu\nu}(x) \quad (1.18)$$

615 which contains not only the free gauge fields contributions, but also the gauge fields
 616 self-interactions and interactions among them.

617 The three weak isospin conserved currents resulting from the $SU(2)_L$ symmetry
 618 are given by

$$J_\mu^i(x) = \frac{1}{2} \bar{\psi}_1(x) \gamma_\mu \sigma^i \psi_1(x) \quad (1.19)$$

619 while the weak hypercharge conserved current resulting from the $U(1)_Y$ symmetry is
 620 given by

$$J_\mu^Y = \sum_{j=1}^3 \bar{\psi}_j(x) \gamma_\mu y_j \psi_j(x) \quad (1.20)$$

621 In order to evaluate the electroweak interactions modeled by an isos triplet field
 622 W_μ^i that couples to isospin currents J_μ^i with strength g and additionally the singlet
 623 field B_μ which couples to the weak hypercharge current J_μ^Y with strength $g'/2$. The
 624 interaction Lagrangian density to be considered is

$$\mathcal{L}_I = -g J^{i\mu}(x) W_\mu^i(x) - \frac{g'}{2} J^{Y\mu}(x) B_\mu(x) \quad (1.21)$$

625 Note that the weak isospin currents are not the same as the charged fermionic cur-
 626 rents that were used to describe the WI (Eqn. 1.8), since the weak isospin eigenstates
 627 are not the same as the mass eigenstates, but they are closely related

$$J_\mu = \frac{1}{2}(J_\mu^1 + iJ_\mu^2), \quad J_\mu^\dagger = \frac{1}{2}(J_\mu^1 - iJ_\mu^2). \quad (1.22)$$

628 The same happens with the gauge fields W_μ^i which are related to the mass eigen-
 629 states W^\pm by

$$W_\mu^+ = \frac{1}{\sqrt{2}}(W_\mu^1 - iW_\mu^2), \quad W_\mu^- = \frac{1}{\sqrt{2}}(W_\mu^1 + iW_\mu^2). \quad (1.23)$$

630 The fact that there are three weak isospin conserved currents is an indication that
 631 in addition to the charged fermionic currents, which couple charged to neutral leptons,
 632 there should be a neutral fermionic current that does not involve electric charge
 633 exchange; therefore, it couples neutral fermions or fermions of the same electric charge.
 634 The third weak isospin current contains a term that is similar to the electromagnetic

635 current (j_μ^{em}), indicating that there is a relation between them and resembling the
 636 Gell-Mann-Nishijima formula 1.2 adapted to electroweak interactions

$$Q = T_3 + \frac{Y_W}{2}. \quad (1.24)$$

637 Just as Q generates the $U(1)_{em}$ symmetry, the weak hypercharge generates the
 638 $U(1)_Y$ symmetry as said before. It is possible to write the relationship in terms of
 639 the currents as

$$j_\mu^{em} = J_\mu^3 + \frac{1}{2} J_\mu^Y. \quad (1.25)$$

640 The neutral gauge fields W_μ^3 and B_μ cannot be directly identified with the Z
 641 and the photon fields since the photon interacts similarly with left and right-handed
 642 fermions; however, they are related through a linear combination given by

$$A_\mu = B_\mu \cos \theta_W + W_\mu^3 \sin \theta_W \quad (1.26)$$

$$Z_\mu = -B_\mu \sin \theta_W + W_\mu^3 \cos \theta_W$$

where θ_W is known as the *Weinberg angle*. The interaction Lagrangian is now given by

$$\begin{aligned} \mathcal{L}_I = -\frac{g}{\sqrt{2}}(J^\mu W_\mu^+ + J^{\mu\dagger} W_\mu^-) - & \left(g \sin \theta_W J_\mu^3 + g' \cos \theta_W \frac{J_\mu^Y}{2} \right) A^\mu \\ & - \left(g \cos \theta_W J_\mu^3 - g' \sin \theta_W \frac{J_\mu^Y}{2} \right) Z^\mu \end{aligned} \quad (1.27)$$

643 the first term is the weak charged current interaction, while the second term is the

644 electromagnetic interaction under the condition

$$g \sin \theta_W = g' \cos \theta_W = e, \quad \frac{g'}{g} = \tan \theta_W \quad (1.28)$$

645 contained in the Eqn.1.25; the third term is the neutral weak current.

646

647 Note that the neutral fields transformation given by the Eqn. 1.26 can be written
 648 in terms of the coupling constants g and g' as:

$$A_\mu = \frac{g' W_\mu^3 + g B_\mu}{\sqrt{g^2 + g'^2}}, \quad Z_\mu = \frac{g W_\mu^3 - g' B_\mu}{\sqrt{g^2 + g'^2}}. \quad (1.29)$$

649 So far, the Lagrangian density describing the non-massive EWI is:

$$\mathcal{L}_{nmEWI} = \mathcal{L}_0 + \mathcal{L}_G \quad (1.30)$$

650 where fermion and gauge fields have been considered massless because their regular
 651 mass terms are manifestly non invariant under G transformations; therefore, masses
 652 have to be generated in a gauge invariant way. The mechanism by which this goal is
 653 achieved is known as the *Higgs mechanism* and is closely connected to the concept of
 654 *spontaneous symmetry breaking*.

655 1.3.1 Spontaneous symmetry breaking (SSB)

656 Figure 1.6 left shows a steel nail (top) which is subject to an external force; the form
 657 of the potential energy is also shown (bottom).

658 Before reaching the critical force value, the system has rotational symmetry with
 659 respect to the nail axis; however, after the critical force value is reached the nail buck-

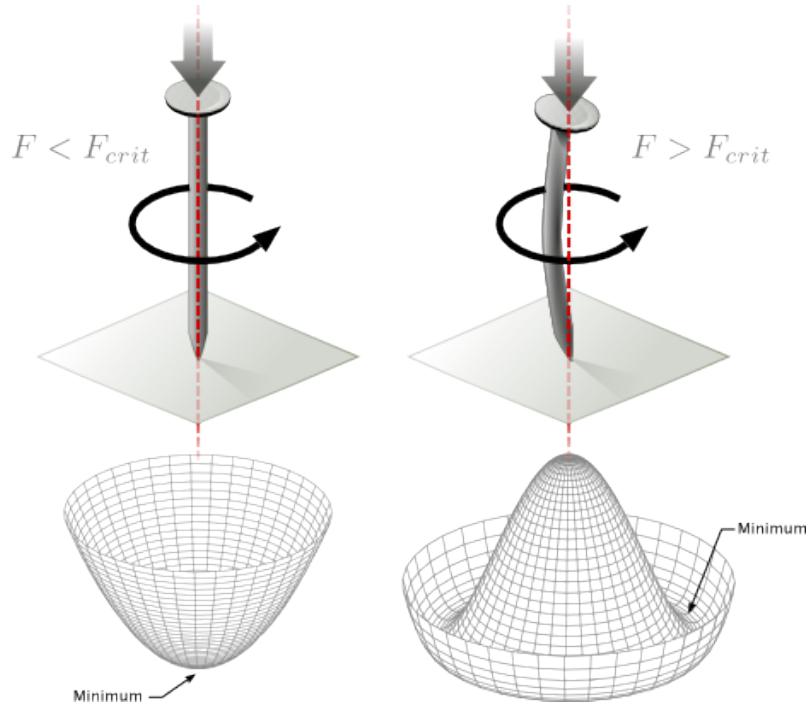


Figure 1.6: Spontaneous symmetry breaking mechanism. The steel nail, subject to an external force (top left), has rotational symmetry with respect to its axis. When the external force overcomes a critical value the nail buckles (top right) choosing a minimal energy state (ground state) and thus *breaking spontaneously the rotational symmetry*. The potential energy (bottom) changes but holds the rotational symmetry; however, an infinite number of asymmetric ground states are generated and circularly distributed in the bottom of the potential [31].

660 les (top right). The form of the potential energy (bottom right) changes appearing a
 661 set of infinity minima but preserving its rotational symmetry. Right before the nail
 662 buckles there is no indication of the direction the nail will bend because any of the
 663 directions are equivalent, but once the nail bends, choosing a direction, an arbitrary
 664 minimal energy state (ground state) is selected and it does not share the system's
 665 rotational symmetry. This mechanism for reaching an asymmetric ground state is
 666 known as *spontaneous symmetry breaking*.

667 The lesson from this analysis is that the way to introduce the SSB mechanism
 668 into a system is by adding the appropriate potential to it.

669 Figure 1.7 shows a plot of the potential $V(\phi)$ in the case of a scalar field ϕ

$$V(\phi) = \mu^2 \phi^\dagger \phi + \lambda (\phi^\dagger \phi)^2 \quad (1.31)$$

670 If $\mu^2 > 0$ the potential has only one minimum at $\phi = 0$ and describes a scalar field
 671 with mass μ . If $\mu^2 < 0$ the potential has a local maximum at $\phi = 0$ and two minima
 672 at $\phi = \pm\sqrt{-\mu^2/\lambda}$ which enables the SSB mechanism to work.

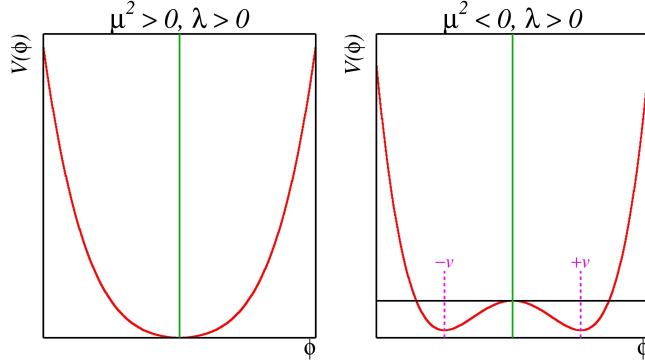


Figure 1.7: Shape of the potential $V(\phi)$ for $\lambda > 0$ and: $\mu^2 > 0$ (left) and $\mu^2 < 0$ (right). The case $\mu^2 < 0$ corresponds to the potential suitable for introducing the SSB mechanism by choosing one of the two ground states which are connected via reflection symmetry. [31].

673 In the case of a complex scalar field $\phi(x)$

$$\phi(x) = \frac{1}{\sqrt{2}}(\phi_1 + i\phi_2) \quad (1.32)$$

674 the Lagrangian (invariant under global $U(1)$ transformations) is given by

$$\mathcal{L} = (\partial_\mu \phi)^\dagger (\partial^\mu \phi) - V(\phi), \quad V(\phi) = \mu^2 \phi^\dagger \phi + \lambda (\phi^\dagger \phi)^2 \quad (1.33)$$

675 where an appropriate potential has been added in order to introduce the SSB.

676 As seen in Figure 1.8, the potential has now an infinite number of minima circularly
 677 distributed along the ξ -direction which makes possible the occurrence of the SSB by
 678 choosing an arbitrary ground state; for instance, $\xi = 0$, i.e. $\phi_1 = v, \phi_2 = 0$

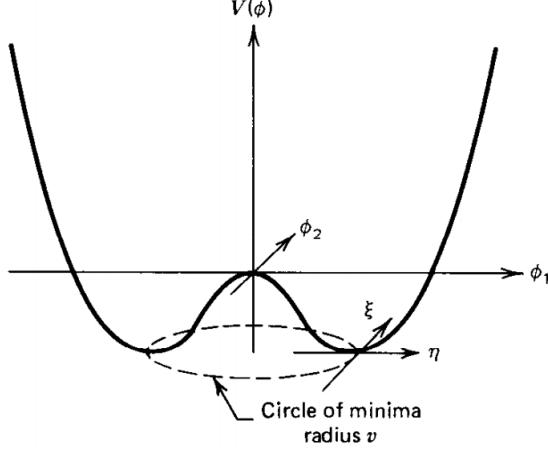


Figure 1.8: Potential for complex scalar field. There is a circle of minima of radius v along the ξ -direction [6].

$$\phi_0 = \frac{v}{\sqrt{2}} \exp(i\xi) \xrightarrow{\text{SSB}} \phi_0 = \frac{v}{\sqrt{2}} \quad (1.34)$$

679 As usual, excitations over the ground state are studied by making an expansion
 680 about it; thus, the excitations can be parametrized as:

$$\phi(x) = \frac{1}{\sqrt{2}}(v + \eta(x) + i\xi(x)) \quad (1.35)$$

681 which when substituted into Eqn. 1.33 produces a Lagrangian in terms of the new
 682 fields η and ξ

$$\mathcal{L}' = \frac{1}{2}(\partial_\mu \xi)^2 + \frac{1}{2}(\partial_\mu \eta)^2 + \mu^2 \eta^2 - V(\phi_0) - \lambda v \eta (\eta^2 + \xi^2) - \frac{\lambda}{4}(\eta^2 + \xi^2)^2 \quad (1.36)$$

683 where the last two terms represent the interactions and self-interaction between the
 684 two fields η and ξ . The particular feature of the SSB mechanism is revealed when
 685 looking to the first three terms of \mathcal{L}' . Before the SSB, only the massless ϕ field is

686 present in the system; after the SSB there are two fields of which the η -field has
 687 acquired mass $m_\eta = \sqrt{-2\mu^2}$ while the ξ -field is still massless (see Figure 1.9).

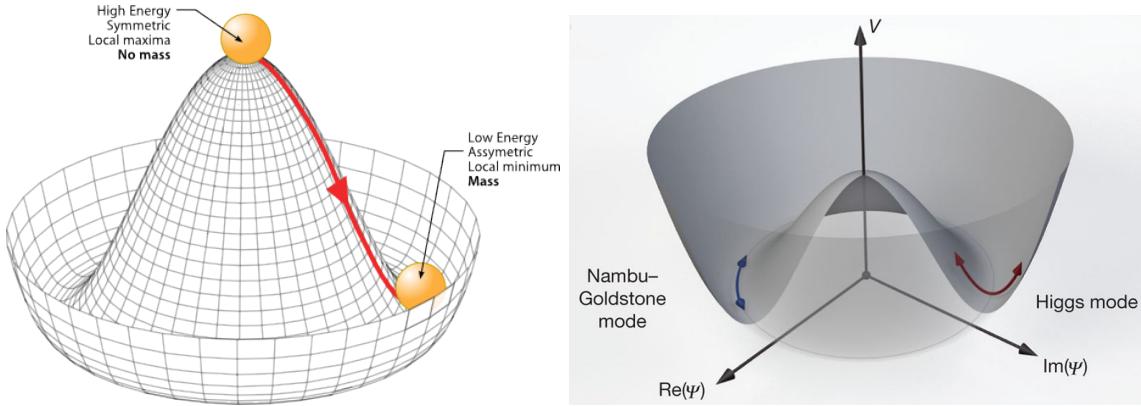


Figure 1.9: SSB mechanism for a complex scalar field [31, 32].

688 Thus, the SSB mechanism serves as a method to generate mass but as a side ef-
 689 fect a massless field is introduced in the system. This fact is known as the Goldstone
 690 theorem and states that a massless scalar field appears in the system for each con-
 691 tinuous symmetry spontaneously broken. Another version of the Goldstone theorem
 692 states that “if a Lagrangian is invariant under a continuous symmetry group G , but
 693 the vacuum is only invariant under a subgroup $H \subset G$, then there must exist as many
 694 massless spin-0 particles (Nambu-Goldstone bosons) as broken generators.” [26] The
 695 Nambu-Goldstone boson can be understood considering that the potential in the ξ -
 696 direction is flat so excitations in that direction are not energy consuming and thus
 697 represent a massless state.

698 1.3.2 Higgs mechanism

699 When the SSB mechanism is introduced in the formulation of the EWI in an attempt
 700 to generate the mass of the so far massless gauge bosons and fermions, an interesting
 701 effect is revealed. In order to keep the G symmetry group invariance and generate

702 the mass of the EW gauge bosons, a G invariant Lagrangian density (\mathcal{L}_S) has to be
 703 added to the non massive EWI Lagrangian (Eqn. 1.30)

$$\mathcal{L}_S = (D_\mu \phi)^\dagger (D^\mu \phi) - \mu^2 \phi^\dagger \phi - \lambda (\phi^\dagger \phi)^2, \quad \lambda > 0, \mu^2 < 0 \quad (1.37)$$

$$D_\mu \phi = \left(i\partial_\mu - g \frac{\sigma_i}{2} W_\mu^i - g' \frac{Y}{2} B_\mu \right) \phi \quad (1.38)$$

704 ϕ has to be an isospin doublet of complex scalar fields so it preserves the G invariance;
 705 thus ϕ can be defined as:

$$\phi = \begin{pmatrix} \phi^+ \\ \phi^0 \end{pmatrix} \equiv \frac{1}{\sqrt{2}} \begin{pmatrix} \phi_1 + i\phi_2 \\ \phi_3 + i\phi_4 \end{pmatrix}. \quad (1.39)$$

706 The minima of the potential are defined by

$$\phi^\dagger \phi = \frac{1}{2} (\phi_1^2 + \phi_2^2 + \phi_3^2 + \phi_4^2) = -\frac{\mu^2}{2\lambda}. \quad (1.40)$$

707 The choice of the ground state is critical. By choosing a ground state, invariant
 708 under $U(1)_{em}$ gauge symmetry, the photon will remain massless and the W^\pm and Z
 709 bosons masses will be generated which is exactly what is needed. In that sense, the
 710 best choice corresponds to a weak isospin doublet with $T_3 = -1/2$, $Y_W = 1$ and $Q = 0$
 711 which defines a ground state with $\phi_1 = \phi_2 = \phi_4$ and $\phi_3 = v$:

$$\phi_0 \equiv \frac{1}{\sqrt{2}} \begin{pmatrix} 0 \\ v \end{pmatrix}, \quad v^2 \equiv -\frac{\mu^2}{\lambda}. \quad (1.41)$$

712 where the vacuum expectation value v is fixed by the Fermi coupling G_F according
 713 to $v = (\sqrt{2}G_F)^{1/2} \approx 246$ GeV.

714 The G symmetry has been broken and three Nambu-Goldstone bosons will appear.

715 The next step is to expand ϕ about the chosen ground state as:

$$\phi(x) = \frac{1}{\sqrt{2}} \exp\left(\frac{i}{v} \sigma_i \theta^i(x)\right) \begin{pmatrix} 0 \\ v + H(x) \end{pmatrix} \approx \frac{1}{\sqrt{2}} \begin{pmatrix} \theta_1(x) + i\theta_2(x) \\ v + H(x) - i\theta_3(x) \end{pmatrix} \quad (1.42)$$

716 to describe fluctuations from the ground state ϕ_0 . The fields $\theta_i(x)$ represent the
 717 Nambu-Goldstone bosons while $H(x)$ is known as *Higgs field*. The fundamental fea-
 718 ture of the parametrization used is that the dependence on the $\theta_i(x)$ fields is factored
 719 out in a global phase that can be eliminated by taking the physical *unitary gauge*
 720 $\theta_i(x) = 0$. Therefore the expansion about the ground state is given by:

$$\phi(x) \frac{1}{\sqrt{2}} \begin{pmatrix} 0 \\ v + H(x) \end{pmatrix} \quad (1.43)$$

721 which when substituted into \mathcal{L}_S (Eqn. 1.37) results in a Lagrangian containing the
 722 now massive three gauge bosons W^\pm, Z , one massless gauge boson (photon) and
 723 the new Higgs field (H). The three degrees of freedom corresponding to the Nambu-
 724 Goldstone bosons are now integrated into the massive gauge bosons as their lon-
 725 gitudinal polarizations which were not available when they were massless particles.
 726 The effect by which vector boson fields acquire mass after an spontaneous symmetry
 727 breaking, but without an explicit gauge invariance breaking is known as the *Higgs*
 728 *mechanism*.

729 The mechanism was proposed by three independent groups: F.Englert and R.Brout
 730 in August 1964 [33], P.Higgs in October 1964 [34] and G.Guralnik, C.Hagen and
 731 T.Kibble in November 1964 [35]; however, its importance was not realized until
 732 S.Glashow [22], A.Salam [23] and S.Weinberg [24], independently, proposed that elec-
 733 tromagnetic and weak interactions are two manifestations of a more general interac-
 734 tion called *electroweak interaction* in 1967.

735 **1.3.3 Masses of the gauge bosons**

736 The masses of the gauge bosons are extracted by evaluating the kinetic part of La-
 737 grangian \mathcal{L}_S in the ground state (known also as the vacuum expectation value), i.e.,

$$\left| \left(\partial_\mu - ig \frac{\sigma_i}{2} W_\mu^i - i \frac{g'}{2} B_\mu \right) \phi_0 \right|^2 = \left(\frac{1}{2} v g \right)^2 W_\mu^+ W^{-\mu} + \frac{1}{8} v^2 (W_\mu^3, B_\mu) \begin{pmatrix} g^2 & -gg' \\ -gg' & g'^2 \end{pmatrix} \begin{pmatrix} W^{3\mu} \\ B^\mu \end{pmatrix} \quad (1.44)$$

738 comparing with the typical mass term for a charged boson $M_W^2 W^+ W^-$

$$M_W = \frac{1}{2} v g. \quad (1.45)$$

The second term in the right side of the Eqn.1.44 comprises the masses of the neutral bosons, but it needs to be written in terms of the gauge fields Z_μ and A_μ in order to be compared to the typical mass terms for neutral bosons, therefore using Eqn. 1.29

$$\begin{aligned} \frac{1}{8} v^2 [g^2 (W_\mu^3)^2 - 2gg' W_\mu^3 B^\mu + g'^2 B_\mu^2] &= \frac{1}{8} v^2 [g W_\mu^3 - g' B_\mu]^2 + 0[g' W_\mu^3 + g B_\mu]^2 \quad (1.46) \\ &= \frac{1}{8} v^2 [\sqrt{g^2 + g'^2} Z_\mu]^2 + 0[\sqrt{g^2 + g'^2} A_\mu]^2 \end{aligned}$$

739 and then

$$M_Z = \frac{1}{2} v \sqrt{g^2 + g'^2}, \quad M_A = 0 \quad (1.47)$$

740 **1.3.4 Masses of the fermions**

741 The lepton mass terms can be generated by introducing a gauge invariant Lagrangian
 742 term describing the Yukawa coupling between the lepton field and the Higgs field

$$\mathcal{L}_{Yl} = -G_l \left[(\bar{\nu}_l, \bar{l})_L \begin{pmatrix} \phi^+ \\ \phi^0 \end{pmatrix} l_R + \bar{l}_R (\phi^-, \bar{\phi}^0) \begin{pmatrix} \nu_l \\ l \end{pmatrix} \right], \quad l = e, \mu, \tau. \quad (1.48)$$

743 After the SSB and replacing the usual field expansion about the ground state
 744 (Eqn.1.41) into \mathcal{L}_{Yl} , the mass term arises

$$\mathcal{L}_{Yl} = -m_l (\bar{l}_L l_R + \bar{l}_R l_L) - \frac{m_l}{v} (\bar{l}_L l_R + \bar{l}_R l_L) H = -m_l \bar{l} l \left(1 + \frac{H}{v} \right) \quad (1.49)$$

745

$$m_l = \frac{G_l}{\sqrt{2}} v \quad (1.50)$$

746 where the additional term represents the lepton-Higgs interaction. The quark masses
 747 are generated in a similar way as lepton masses but for the upper member of the
 748 quark doublet a different Higgs doublet is needed:

$$\phi_c = -i\sigma_2 \phi^* = \begin{pmatrix} -\bar{\phi}^0 \\ \phi^- \end{pmatrix}. \quad (1.51)$$

749 Additionally, given that the quark isospin doublets are not constructed in terms
 750 of the mass eigenstates but in terms of the flavor eigenstates, as shown in Table 1.5,
 751 the coupling parameters will be related to the CKM matrix elements; thus, the quark
 752 Lagrangian is given by:

$$\mathcal{L}_{Yq} = -G_d^{i,j} (\bar{u}_i, \bar{d}'_i)_L \begin{pmatrix} \phi^+ \\ \phi^0 \end{pmatrix} d_{jR} - G_u^{i,j} (\bar{u}_i, \bar{d}'_i)_L \begin{pmatrix} -\bar{\phi}^0 \\ \phi^- \end{pmatrix} u_{jR} + h.c. \quad (1.52)$$

753 with $i, j = 1, 2, 3$. After SSB and expansion about the ground state, the diagonal form

754 of \mathcal{L}_{Yq} is:

$$\mathcal{L}_{Yq} = -m_d^i \bar{d}_i d_i \left(1 + \frac{H}{v}\right) - m_u^i \bar{u}_i u_i \left(1 + \frac{H}{v}\right) \quad (1.53)$$

755 Fermion masses depend on arbitrary couplings G_l and $G_{u,d}$ and are not predicted
756 by the theory.

757 1.3.5 The Higgs field

758 After the characterization of the fermions and gauge bosons as well as their interac-
759 tions, it is necessary to characterize the Higgs field itself. The Lagrangian \mathcal{L}_S in Eqn.
760 1.37 written in terms of the gauge bosons is given by

$$\mathcal{L}_S = \frac{1}{4} \lambda v^4 + \mathcal{L}_H + \mathcal{L}_{HV} \quad (1.54)$$

761

$$\mathcal{L}_H = \frac{1}{2} \partial_\mu H \partial^\mu H - \frac{1}{2} m_H^2 H^2 - \frac{1}{2v} m_H^2 H^3 - \frac{1}{8v^2} m_H^2 H^4 \quad (1.55)$$

762

$$\mathcal{L}_{HV} = m_H^2 W_\mu^+ W^{\mu-} \left(1 + \frac{2}{v} H + \frac{2}{v^2} H^2\right) + \frac{1}{2} m_Z^2 Z_\mu Z^\mu \left(1 + \frac{2}{v} H + \frac{2}{v^2} H^2\right) \quad (1.56)$$

763 The mass of the Higgs boson is deduced as usual from the mass term in the Lagrangian
764 resulting in:

$$m_H = \sqrt{-2\mu^2} = \sqrt{2\lambda}v \quad (1.57)$$

765 however, it is not predicted by the theory either. The experimental measurement of
766 the Higgs boson mass have been performed by the *Compact Muon Solenoid (CMS)*
767 experiment and the *A Toroidal LHC Appartus (ATLAS)* experiments at the *Large
768 Hadron Collider (LHC)*, [36–38], and is presented in Table 1.8.

769

Property	Value
Electric charge	0
Color charge	0
Spin	0
Weak isospin	-1/2
Weak hypercharge	1
Parity	1
Mass (GeV/c ²)	125.09±0.21 (stat.)±0.11 (syst.)

Table 1.8: Higgs boson properties. Higgs mass is not predicted by the theory and the value here corresponds to the experimental measurement.

770 1.3.6 Production of Higgs bosons at LHC

771 At the LHC, Higgs bosons are produced as a result of the collision of two counter-
 772 rotating protons beams. A detailed description of the LHC machine will be presented
 in chapter 2.

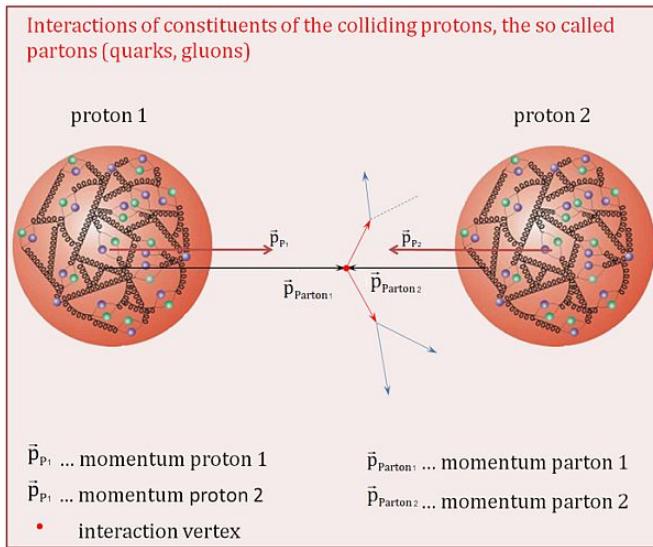


Figure 1.10: Proton-proton collision. Protons are composed of 3 valence quarks, a sea of quarks and gluons; therefore in a proton-proton collision, quarks and gluons are those who collide. [39].

773

774 Protons are composed of quarks and these quarks are bound by gluons; however,
775 what is commonly called the quark content of the proton makes reference to the
776 valence quarks. In fact, a proton is not just a rigid entity with three balls in it all
777 tied up with springs, but the gluons exchanged by the valence quarks tend to split

778 spontaneously into quark-antiquark pairs or more gluons, creating *sea of quarks and*
 779 *gluons* as represented in Figure 1.10.

780 In a proton-proton (pp) collision, the proton's constituents, quarks and gluons, are
 781 those that collide. The pp cross section depends on the momentum of the colliding
 782 particles, reason for which it is needed to know how the momentum is distributed
 783 inside the proton. Quarks and gluons are known as partons, hence, the functions
 784 that describe how the proton momentum is distributed among partons inside it are
 785 called *parton distribution functions (PDFs)*; PDFs are determined from experimental
 786 data obtained in experiments where the internal structure of hadrons is tested, and
 787 depend on the momentum transfer Q and the fraction of momentum x carried by an
 788 specific parton. Figure 1.11 shows the proton PDFs ($xf(x, Q^2)$) for two values of Q .

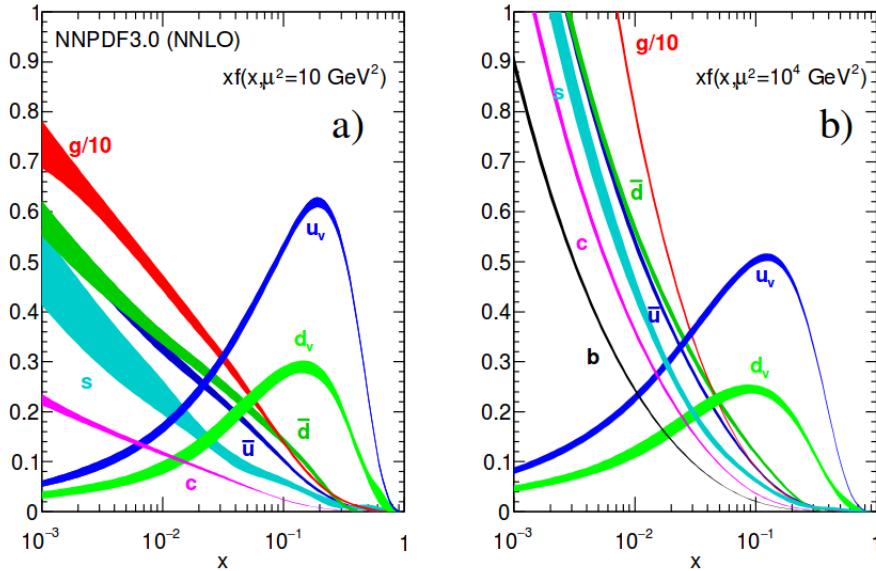


Figure 1.11: Proton PDFs for two values of Q^2 : left. $\mu^2 = Q^2 = 10 \text{ GeV}^2$, right. $\mu^2 = Q^2 = 10^4 \text{ GeV}^2$. u_v and d_v correspond to the u and d valence quarks, $s, c, b, \bar{u}, \bar{d}$ correspond to sea quarks, and g corresponds to gluons. Note that gluons carry a high fraction of the proton's momentum. [9]

789 In physics, a common approach to study complex systems consists of starting
 790 with a simpler version of them, for which a well known description is available, and

adding an additional *perturbation* which represents a small deviation from the known behavior. If the perturbation is small enough, the physical quantities associated with the perturbed system are expressed as a series of corrections to those of the simpler system. The perturbation series corresponds to an expansion in power series of a small parameter, therefore, the more terms are considered in the series (the higher order in the perturbation series), the more precise is the description of the complex system. If the perturbation does not get progressively smaller, the strategy cannot be applied and new methods have to be employed.

High energy systems, like the Higgs production at LHC explored in this thesis, usually can be treated perturbatively with the expansion made in terms of the coupling constants. The overview presented here will be oriented specifically to the Higgs boson production mechanisms in pp collisions at LHC.

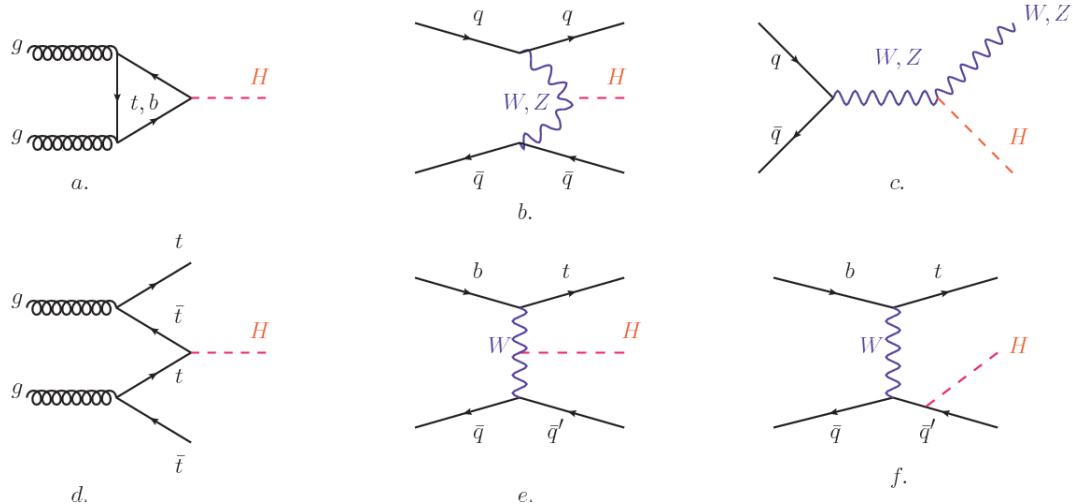


Figure 1.12: Main Higgs boson production mechanism Feynman diagrams. a. gluon-gluon fusion, b. vector boson fusion (VBF), c. Higgs-strahlung, d. Associated production with a top or bottom quark pair, e-f. associated production with a single top quark.

Figure 1.12 shows the Feynman diagrams for the leading order (first order) Higgs production processes at LHC; note that in these diagrams the incoming particles are not the protons themselves but the partons from the protons that actually participate

in the interaction, hence, theorists typically calculate the cross section for the parton interaction, and then convolute that cross section with the information from the PDFs to get a production cross section that is actually measured in experiments. The cross section for Higgs production as a function of the center of mass-energy (\sqrt{s}) for pp collisions is showed in Figure 1.13 left. The tags NLO (next to leading order), NNLO (next to next to leading order) and N3LO (next to next to next to leading order) make reference to the order at which the perturbation series have been considered while the tags QCD and EW correspond to the strong and electroweak coupling constants respectively.

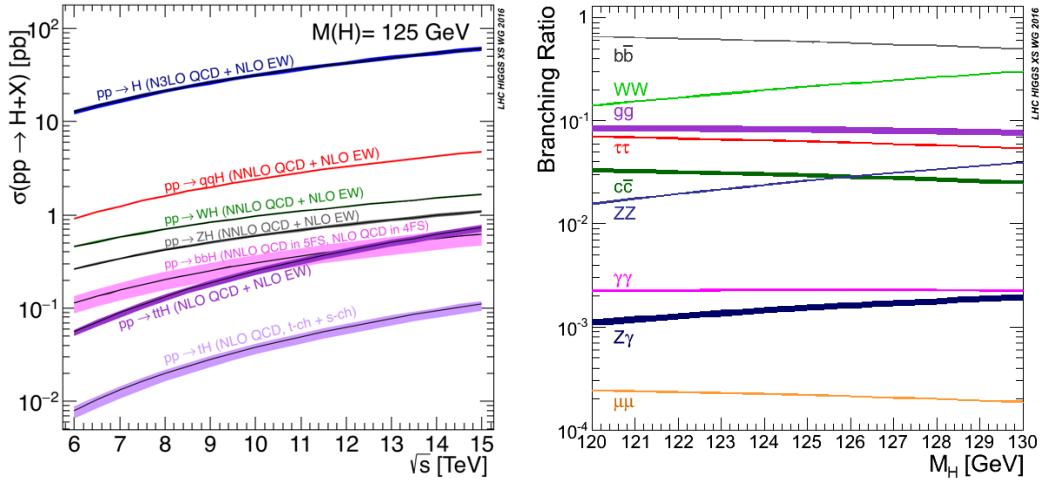


Figure 1.13: Higgs boson production cross sections (left) and decay branching ratios (right) for the main mechanisms. The VBF is indicated as qqH [40].

The main production mechanism is the gluon fusion (Figure 1.12a and $pp \rightarrow H$ in Figure 1.13) given that gluons carry the highest fraction of momentum of the protons in pp colliders (as shown in Figure ??). Since the Higgs boson does not couple to gluons, the mechanism proceeds through the exchange of a virtual top-quark loop. Note that in this process the Higgs boson is produced alone, turning out to be problematic for some Higgs decays, because such absence of anything produced in

821 association with the Higgs represent a trouble for triggering, however, this mechanism
 822 is experimentally clean when combined with the two-photon or the four-lepton decay
 823 channels (see Section 1.3.7).

824 Vector boson fusion (Figure 1.12b and $pp \rightarrow qqH$ in Figure 1.13) has the second
 825 largest production cross section. The scattering of two fermions is mediated by a weak
 826 gauge boson which later emits a Higgs boson. In the final state, the two fermions tend
 827 to be located in the central region of the detector; this kind of features are generally
 828 used as a signature when analyzing the datasets provided by the experiments⁷.

829 In the Higgs-strahlung mechanism (Figure 1.12c and $pp \rightarrow WH, pp \rightarrow ZH$ in
 830 Figure 1.13) two fermions annihilate to form a weak gauge boson. If the initial
 831 fermions have enough energy, the emergent boson might emit a Higgs boson.

832 The associated production with a top or bottom quark pair and the associated
 833 production with a single top quark (Figure 1.12d-f and $pp \rightarrow bbH, pp \rightarrow t\bar{t}H, pp \rightarrow tH$
 834 in Figure 1.13) have a smaller cross section than the main three mechanisms above,
 835 but they provide a good opportunity to test the Higgs-top coupling. The analysis
 836 reported in this thesis is developed using these production mechanisms. A detailed
 837 description of the tH mechanism will be given in Section 1.5.

838 1.3.7 Higgs boson decay channels

839 When a particle can decay through several modes, also known as channels, the prob-
 840 ability of decaying through a given channel is quantified by the *branching ratio (BR)*
 841 of the decay channel; thus, the BR is defined as the ratio of number of decays go-
 842 ing through that given channel to the total number of decays. In regard to the
 843 Higgs boson decay, the BR can be predicted with accuracy once the Higgs mass is
 844 known [41, 42]. In Figure 1.13 right, a plot of the BR as a function of the Higgs mass

⁷ More details about how to identify events of interest in this analysis will be given in chapter 6.

845 is presented; the largest predicted BR corresponds to the $b\bar{b}$ pair decay channel (see
 846 Table 1.9) given that it is the heaviest particle pair whose on-shell⁸ production is
 847 kinematically allowed in the decay.

Decay channel	Branching ratio	Rel. uncertainty
$H \rightarrow b\bar{b}$	5.84×10^{-1}	+3.2% – 3.3%
$H \rightarrow W^+W^-$	2.14×10^{-1}	+4.3% – 4.2%
$H \rightarrow \tau^+\tau^-$	6.27×10^{-2}	+5.7% – 5.7%
$H \rightarrow ZZ$	2.62×10^{-2}	+4.3% – 4.1%
$H \rightarrow \gamma\gamma$	2.27×10^{-3}	+5.0% – 4.9%
$H \rightarrow Z\gamma$	1.53×10^{-3}	+9.0% – 8.9%
$H \rightarrow \mu^+\mu^-$	2.18×10^{-4}	+6.0% – 5.9%

Table 1.9: Predicted branching ratios and the relative uncertainty for some decay channels of a SM Higgs boson with $m_H = 125\text{GeV}/c^2$ [9]; the uncertainties are driven by theoretical uncertainties for the different Higgs boson partial widths and by parametric uncertainties associated to the strong coupling and the masses of the quarks which are the input parameters. Further details on these calculations can be found in Reference [43]

848

849 Decays to other lepton and quark pairs, like electron, strange, up, and down
 850 quark pairs not listed in the table, are also possible but their likelihood is too small
 851 to measure since they are very lightweight, hence, their interaction with the Higgs
 852 boson is very weak. On other hand, the decay to top quark pairs is heavily suppressed
 853 due to the top quark mass ($\approx 173\text{ GeV}/c^2$).

854 Decays to gluons proceed indirectly through a virtual top quark loop while the
 855 decays to photons proceed through a virtual W boson loop, therefore, their branching
 856 ratio is smaller compared to direct interaction decays. Same is true for the decay to
 857 a photon and a Z boson.

⁸ In general, on-shell or real particles are those which satisfy the energy-momentum relation ($E^2 - |\vec{p}|^2 c^2 = m^2 c^4$); off-shell or virtual particles does not satisfy it which is possible under the uncertainty principle of quantum mechanics. Usually, virtual particles correspond to internal propagators in Feynman diagrams.

858 In the case of decays to pairs of W and Z bosons, the decay proceed with one of
 859 the bosons being on-shell and the other being off-shell. The likelihood of the process
 860 diminish depending on how far off-shell are the virtual particles involved, hence, the
 861 branching ratio for W boson pairs is bigger than for Z boson pairs since Z boson mass
 862 is bigger than W boson mass.

863 Note that the decay to a pair of virtual top quarks is possible, but the probability
 864 is way too small.

865 **1.4 Experimental status of the anomalous
 866 Higgs-fermion coupling**

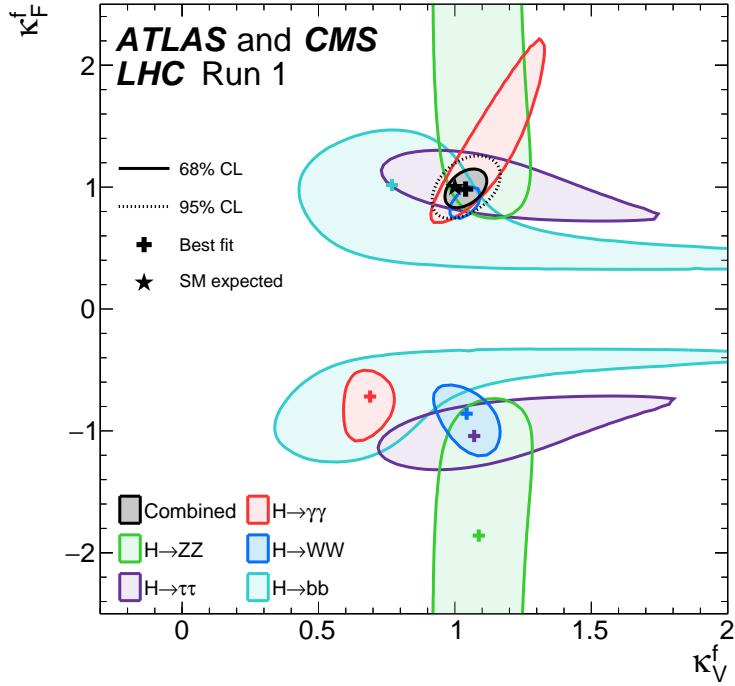


Figure 1.14: Combination of the ATLAS and CMS fits for coupling modifiers $\kappa_t - \kappa_V$; also shown the individual decay channels combination and their global combination. No assumptions have been made on the sign of the coupling modifiers [44].

867 ATLAS and CMS have performed analyses of the anomalous Higgs-fermion cou-
 868 pling by making likelihood scans for the two coupling modifiers, κ_f and κ_V , under
 869 the assumption that $\kappa_Z = \kappa_W \equiv \kappa_V$ and $\kappa_t = \kappa_\tau = \kappa_b \equiv \kappa_f$. Figure 1.14 shows the
 870 result of the combination of ATLAS and CMS fits; also the individual decay channels
 871 combination and the global combination results are shown. Note that from this plot
 872 there is limited information on the sign of the coupling since the only information
 873 available about the sign of the coupling comes from decays rather than production.

874 While all the channels are compatible for positive values of the modifiers, for
 875 negative values of κ_f there is no compatibility. The best fit for individual channels
 876 is compatible with negative values of κ_f except for the $H \rightarrow bb$ channel. The best
 877 fit for the combination yields $\kappa_f \geq 0$, in contrast to the yields from the individual
 878 channels; the reason of this yield resides in the $H \rightarrow \gamma\gamma$ coupling. $H \rightarrow \gamma\gamma$ decay
 879 proceeds through a loop of either top quarks or W bosons, hence, this channel is
 880 sensitive to κ_t thanks to the interference of these two amplitude contributions; under
 881 the assumption that no beyond SM particles take part in the loops, a flipped sign
 882 of κ_t will increase the $H \rightarrow \gamma\gamma$ branching fraction by a factor of ~ 2.4 which is not
 883 supported by measurements; thus, this large asymmetry between the positive and
 884 negative coupling ratios in the $H \rightarrow \gamma\gamma$ channel drives the yield of the global fit and
 885 would mean that the anomalous H-t coupling is excluded as stated in Reference [44],
 886 but there is a caveat, this exclusion holds only if no new particles contribute to the
 887 loop in the main diagram for that decay.

888 Although the $H \rightarrow bb$ channel is expected to be the most sensitive channel and
 889 its best fit value of κ_t is positive, and then the global fit yield is still supported,
 890 the contributions from all the other decay channels, small compared to the $H \rightarrow bb$,
 891 indicate that the anomalous H-t coupling cannot be excluded completely, motivating
 892 to look at tH processes which can help with both, the limited information on the sign

893 of the H-t coupling and the access to information from the Higgs boson production
 894 rather than from its decays.

895 It will be shown in Section 1.5 that the same interference effect enhance the
 896 tH production rate and could reveal evidence of direct production of heavy new par-
 897 ticles as predicted in composite and little Higgs models [45], or new physics related
 898 to Higgs boson mediated flavor changing neutral currents [46] as well as probes the
 899 CP-violating phase of the H-t coupling [47, 48].

900 **1.5 Associated production of a Higgs boson and a 901 single top quark**

902 The production of Higgs boson in association with a top quark has been extensively
 903 studied [47, 49–52]. While measurements of the main Higgs production mechanisms
 904 rates are sensitive to the strength of the Higgs coupling to W boson or top quark,
 905 they are not sensitive to the relative sign between the two couplings. In this thesis,
 906 the Higgs boson production mechanism explored is the associated production with a
 907 single top quark (tH) which offers sensitivity to the relative sign of the Higgs couplings
 908 to W boson and to top quark. The description given here is based on Reference [51]

A process where two incoming particles interact and produce a final state with two
 particles can proceed in three called channels (see, for instance, Figure 1.15 omitting
 the red line). The t-channel represents processes where an intermediate particle is
 emitted by one of the incoming particles and absorbed by the other. The s-channel
 represents processes where the two incoming particles merge into an intermediate par-
 ticle which eventually will split into the particles in the final state. The third channel,
 u-channel, is similar to the t-channel but the two outgoing particles interchange their

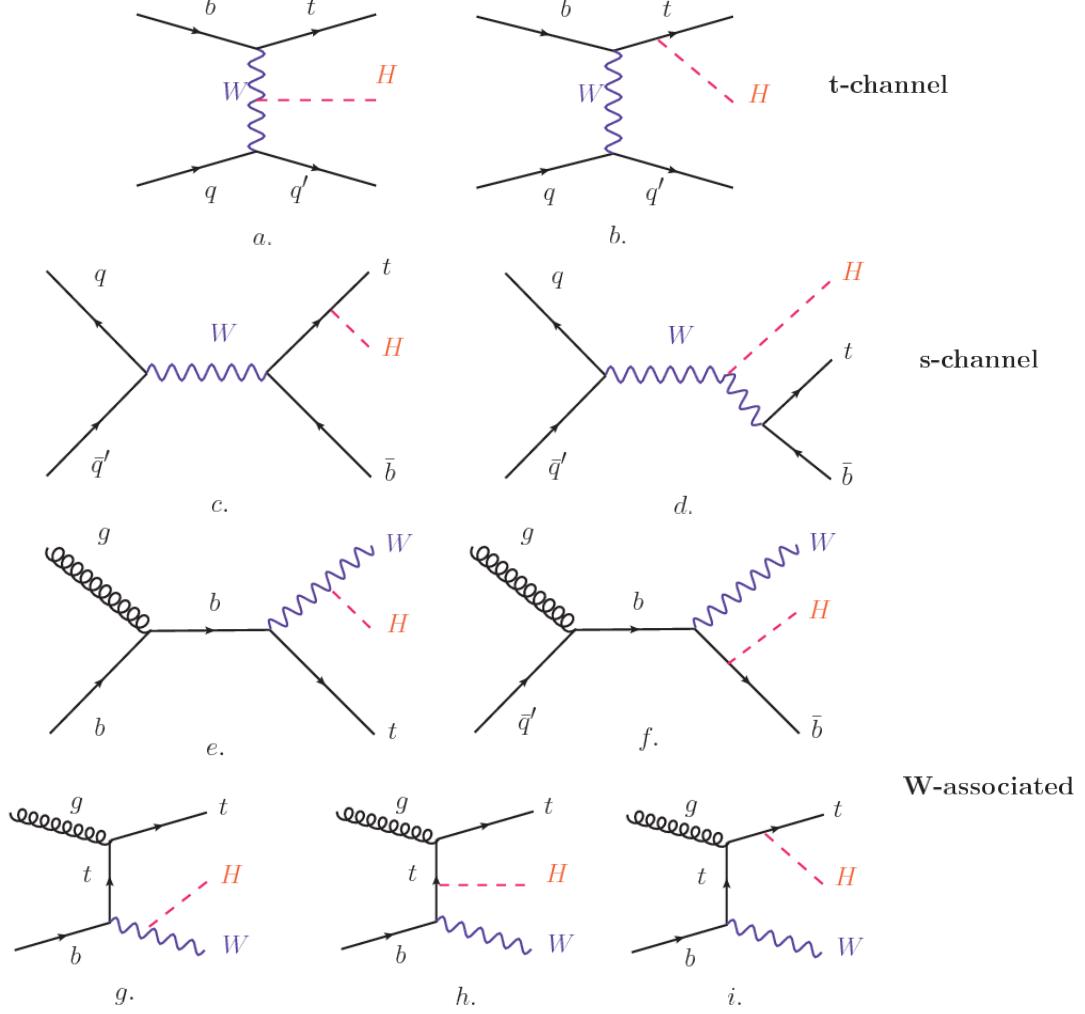


Figure 1.15: Associated Higgs boson production with a top quark mechanism Feynman diagrams. a.,b. t-channel (tHq), c.,d. s-channel (tHb), e-i. W-associated.

roles. These three channels are connected to the so-called Mandelstam variables

$$s = (p_1 + p_2)^2 = (p'_1 + p'_2)^2 \rightarrow \text{square of the center mass-energy.} \quad (1.58)$$

$$t = (p_1 - p'_1)^2 = (p'_2 - p_2)^2 \rightarrow \text{square of the four-momentum transfer.} \quad (1.59)$$

$$u = (p_1 - p'_2)^2 = (p'_1 - p_2)^2 \rightarrow \text{square of the crossed four-momentum transfer.} \quad (1.60)$$

$$s + t + u = m_1^2 + m_2^2 + m'_1^2 + m'_2^2 \quad (1.61)$$

which relate the momentum, energy and the angles of the incoming and outgoing particles in an scattering process of two particles to two particles. The importance of the Mandelstam variables reside in that they form a minimum set of variables needed to describe the kinematics of this scattering process; they are Lorentz invariant which makes them very useful when doing calculations.

The tH production, where Higgs boson can be radiated either from the top quark or from the W boson, is represented by the leading order Feynman diagrams in Figure 1.15. The cross section for the tH process is calculated, as usual, summing over the contributions from the different Feynman diagrams; therefore it depends on the interference between the contributions. In the SM, the interference for t-channel (tHq process) and W-associated (tHW process) production is destructive [49] resulting in the small cross sections presented in Table 1.10.

tH production channel	Cross section (fb)
t-channel ($pp \rightarrow tHq$)	$70.79^{+2.99}_{-4.80}$
W-associated ($pp \rightarrow tHW$)	$15.61^{+0.83}_{-1.04}$
s-channel($pp \rightarrow tHb$)	$2.87^{+0.09}_{-0.08}$

Table 1.10: Predicted SM cross sections for tH production at $\sqrt{s} = 13$ TeV [53, 54].

The s-channel contribution can be neglected. It will be shown that a deviation from the SM destructive interference would result in an enhancement of the tH cross section compared to that in SM, which could be used to get information about the sign of the Higgs-top coupling [51, 52]. In order to describe tH production processes, Feynman diagram 1.15b will be considered; there, the W boson is radiated by a quark in the proton and eventually it will interact with the b quark. In the high energy regime, the effective W approximation [55] is used to describe the process as the

929 emission of an approximately on-shell W and its hard scattering with the b quark;
 930 i.e. $Wb \rightarrow th$. The scattering amplitude for the process is given by

$$\mathcal{A} = \frac{g}{\sqrt{2}} \left[(\kappa_t - \kappa_V) \frac{m_t \sqrt{s}}{m_W v} A \left(\frac{t}{s}, \varphi; \xi_t, \xi_b \right) + \left(\kappa_V \frac{2m_W s}{v t} + (2\kappa_t - \kappa_V) \frac{m_t^2}{m_W v} \right) B \left(\frac{t}{s}, \varphi; \xi_t, \xi_b \right) \right], \quad (1.62)$$

931 where $\kappa_V \equiv g_{HVV}/g_{HVV}^{SM}$ and $\kappa_t \equiv g_{Ht}/g_{Ht}^{SM} = y_t/y_t^{SM}$ are scaling factors that quantify
 932 possible deviations of the couplings from the SM values, Higgs-Vector boson (H-
 933 W) and Higgs-top (H-t) respectively, from the SM couplings; $s = (p_W + p_b)^2$, $t =$
 934 $(p_W - p_H)^2$, φ is the Higgs azimuthal angle around the z axis taken parallel to the
 935 direction of motion of the incoming W; A and B are functions describing the weak
 936 interaction in terms of the chiral states (ξ_t, ξ_b) of the quarks b and t . Terms that
 937 vanish in the high energy limit have been neglected as well as the Higgs and b quark
 938 masses⁹.

939 The scattering amplitude grows with energy like \sqrt{s} for $\kappa_V \neq \kappa_t$, in contrast to
 940 the SM ($\kappa_t = \kappa_V = 1$), where the first term in 1.62 cancels out and the amplitude
 941 is constant for large s ; therefore, a deviation from the SM predictions represents an
 942 enhancement in the tHq cross section. In particular, for a SM H-W coupling and a
 943 H-t coupling of inverted sign with respect to the SM ($\kappa_V = -\kappa_t = 1$) the tHq cross
 944 section is enhanced by a factor greater 10 as seen in the Figure 1.16 taken from
 945 Reference [51]; Reference [56] has reported similar enhancement results.

946 A similar analysis is valid for the W-associated channel but, in that case, the in-
 947 terference is more complicated since there are more than two contributions and an ad-
 948 ditional interference with the production of Higgs boson and a top pair process($t\bar{t}H$).
 949 The calculations are made using the so-called Diagram Removal (DR) technique where

⁹ A detailed explanation of the structure and approximations used to derive \mathcal{A} can be found in Reference [51]

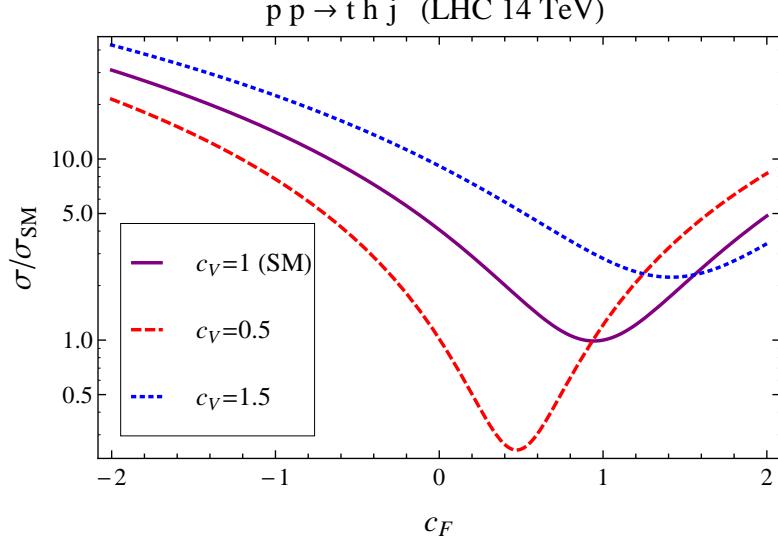


Figure 1.16: Cross section for tHq process as a function of κ_t , normalized to the SM, for three values of κ_V . In the plot c_f refers to the Higgs-fermion coupling which is dominated by the H-t coupling and represented here by κ_t . Solid, dashed and dotted lines correspond to $c_V \rightarrow \kappa_V = 1, 0.5, 1.5$ respectively. Note that for the SM ($\kappa_V = \kappa_t = 1$), the destructive effect of the interference is maximal.

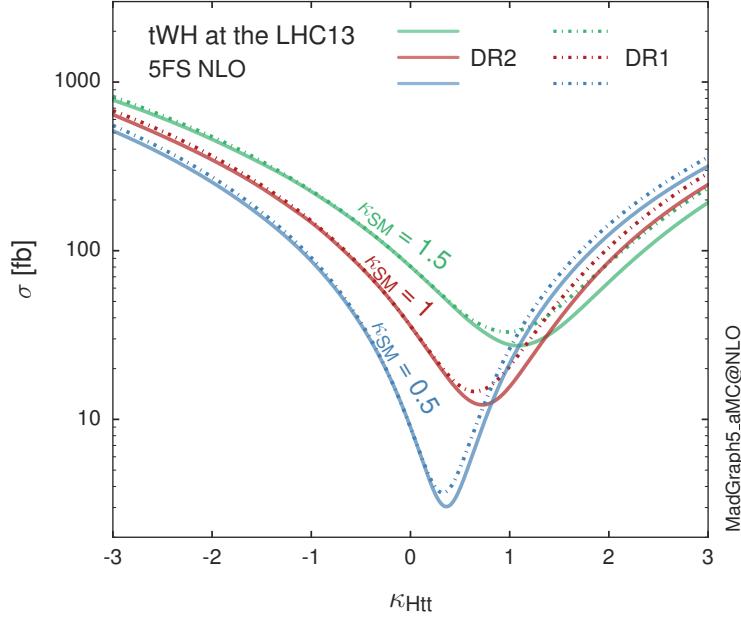


Figure 1.17: Cross section for tHW process as a function of κ_{Htt} , for three values of κ_{SM} at $\sqrt{s} = 13$ TeV. $\kappa_{Htt}^2 = \sigma_{Htt}/\sigma_{Htt}^{SM}$ is a simple re-scaling of the SM Higgs interactions.

interfering diagrams are removed (or added) from the calculations in order to evaluate the impact of the removed contributions. DR1 was defined to neglect $t\bar{t}H$ interference while DR2 was defined to take $t\bar{t}H$ interference into account [48]. As shown in Figure 1.17, the tHW cross section is enhanced from about 15 fb (SM: $\kappa_{Htt} = 1$) to about 150 fb ($\kappa_{Htt} = -1$). Differences between curves for DR1 and DR2 help to gauge the impact of the interference with $t\bar{t}H$. Results of the calculations of the tHq and tHW cross sections at $\sqrt{s} = 13$ TeV can be found in Reference [57] and a summary of the results is presented in Table 1.11.

	\sqrt{s} TeV	$\kappa_t = 1$	$\kappa_t = -1$
$\sigma^{LO}(tHq)(fb)$ [51]	8	≈ 17.4	≈ 252.7
	14	≈ 80.4	≈ 1042
$\sigma^{NLO}(tHq)(fb)$ [51]	8	$18.28^{+0.42}_{-0.38}$	$233.8^{+4.6}_{-0.0}$
	14	$88.2^{+1.7}_{-0.0}$	$982.8^{+28}_{-0.0}$
$\sigma^{LO}(tHq)(fb)$ [56]	14	≈ 71.8	≈ 893
$\sigma^{LO}(tHW)(fb)$ [56]	14	≈ 16.0	≈ 139
$\sigma^{NLO}(tHq)(fb)$ [57]	8	$18.69^{+8.62\%}_{-17.13\%}$	-
	13	$74.25^{+7.48\%}_{-15.35\%}$	$848^{+7.37\%}_{-13.70\%}$
	14	$90.10^{+7.34\%}_{-15.13\%}$	$1011^{+7.24\%}_{-13.39\%}$
$\sigma^{LO}(tHW)(fb)$ [48]	13	$15.77^{+15.91\%}_{-15.76\%}$	-
$\sigma^{NLO}DR1(tHW)(fb)$ [48]	13	$21.72^{+6.52\%}_{-5.24\%}$	≈ 150
$\sigma^{NLO}DR2(tHW)(fb)$ [48]	13	$16.28^{+7.34\%}_{-15.13\%}$	≈ 150

Table 1.11: Predicted enhancement of the tHq and tHW cross sections at LHC for $\kappa_V = 1$ and $\kappa_t = \pm 1$ at LO and NLO; the cross section enhancement of more than a factor of 10 is due to the flipping in the sign of the H-t coupling with respect to the SM one.

958

1.6 CP-mixing in tH processes

In addition to the sensitivity to sign of the H-t coupling, the tHq and tHW processes have been proposed as a tool to investigate the possibility of a H-t coupling that does

962 not conserve CP [47, 48, 58].

963 In this thesis, the sensitivity of tH processes to CP-mixing is also studied on the
 964 basis of References [47, 48] using the effective field theory framework where a generic
 965 particle (X_0) of spin-0 and a general CP violating interaction with the top quark
 966 (Htt coupling), can couple to scalar and pseudo-scalar fermionic densities. The H-W
 967 interaction is assumed to be SM-like. The Lagrangian modeling the H-t interaction
 968 is given by

$$\mathcal{L}_0^t = -\bar{\psi}_t (c_\alpha \kappa_{Htt} g_{Htt} + i s_\alpha \kappa_{Att} g_{Att} \gamma_5) \psi_t X_0, \quad (1.63)$$

969 where α is the CP-mixing phase, $c_\alpha \equiv \cos \alpha$ and $s_\alpha \equiv \sin \alpha$, κ_{Htt} and κ_{Att} are real
 970 dimensionless re-scaling parameters¹⁰ used to parametrize the magnitude of the CP-
 971 violating and CP-conserving parts of the amplitude. The model defines $g_{Htt} =$
 972 $g_{Att} = m_t/v = y_t/\sqrt{2}$ with $v \sim 246$ GeV the Higgs vacuum expectation value. In
 973 this parametrization, three special cases can be recovered

974 • CP-even coupling $\rightarrow \alpha = 0^\circ$

975 • CP-odd coupling $\rightarrow \alpha = 90^\circ$

976 • SM coupling $\rightarrow \alpha = 0^\circ$ and $\kappa_{Htt} = 1$

977 The loop induced X_0 coupling to gluons can also be described in terms of the
 978 parametrization above, according to

$$\mathcal{L}_0^g = -\frac{1}{4} \left(c_\alpha \kappa_{Hgg} g_{Hgg} G_{\mu\nu}^a G^{a,\mu\nu} + s_\alpha \kappa_{Agg} g_{Agg} G_{\mu\nu}^a \tilde{G}^{a,\mu\nu} \right) X_0. \quad (1.64)$$

979 where $g_{Hgg} = -\alpha_s/3\pi v$ and $g_{Agg} = \alpha_s/2\pi v$ and $G_{\mu\nu}$ is the gluon field strength tensors.

980 Under the assumption that the top quark dominates the gluon-fusion process at LHC

¹⁰ analog to κ_t and κ_V

981 energies, $\kappa_{Hgg} \rightarrow \kappa_{Htt}$ and $\kappa_{Agg} \rightarrow \kappa_{Att}$, so that the ratio between the gluon-gluon
 982 fusion cross section for X_0 and for the SM Higgs prediction can be written as

$$\frac{\sigma_{NLO}^{gg \rightarrow X_0}}{\sigma_{NLO,SM}^{gg \rightarrow H}} = c_\alpha^2 \kappa_{Htt}^2 + s_\alpha^2 \left(\kappa_{Att} \frac{g_{Agg}}{g_{Hgg}} \right)^2. \quad (1.65)$$

983 If the re-scaling parameters are set to

$$\kappa_{Htt} = 1, \quad \kappa_{Att} = \left| \frac{g_{Hgg}}{g_{Agg}} \right| = \frac{2}{3}. \quad (1.66)$$

984 the gluon-fusion SM cross section is reproduced for every value of the CP-mixing
 985 angle α ; therefore, by imposing that condition to the Lagrangian density 1.63, the
 986 CP-mixing angle is not constrained by current data. Figure 1.18 shows the NLO cross
 987 sections for t-channel tX_0 (blue) and $t\bar{t}X_0$ (red) associated production processes as a
 988 function of the CP-mixing angle α . X_0 is a generic spin-0 particle with top quark
 989 CP-violating coupling. Re-scaling factors κ_{Htt} and κ_{Att} have been set to reproduce
 990 the SM gluon-fusion cross sections.

991 It is interesting to notice that the tX_0 cross section is enhanced, by a factor of
 992 about 10, when a continuous rotation in the scalar-pseudoscalar plane is applied; this
 993 enhancement is similar to the enhancement produced when the H-t coupling is flipped
 994 in sign with respect to the SM ($y_t = -y_{t,SM}$ in the plot), as showed in Section 1.5. In
 995 contrast, the degeneracy in the $t\bar{t}X_0$ cross section is still present given that it depends
 996 quadratically on the H-t coupling, but more interesting is to notice that $t\bar{t}X_0$ cross
 997 section is exceeded by tX_0 cross section after $\alpha \sim 60^\circ$.

998 A similar parametrization can be used to investigate the tHW process sensitivity
 999 to CP-violating H-t coupling. As said in 1.5, the interference in the W-associated
 1000 channel is more complicated because there are more than two contributions and also
 1001 there is interference with the $t\bar{t}H$ production process.

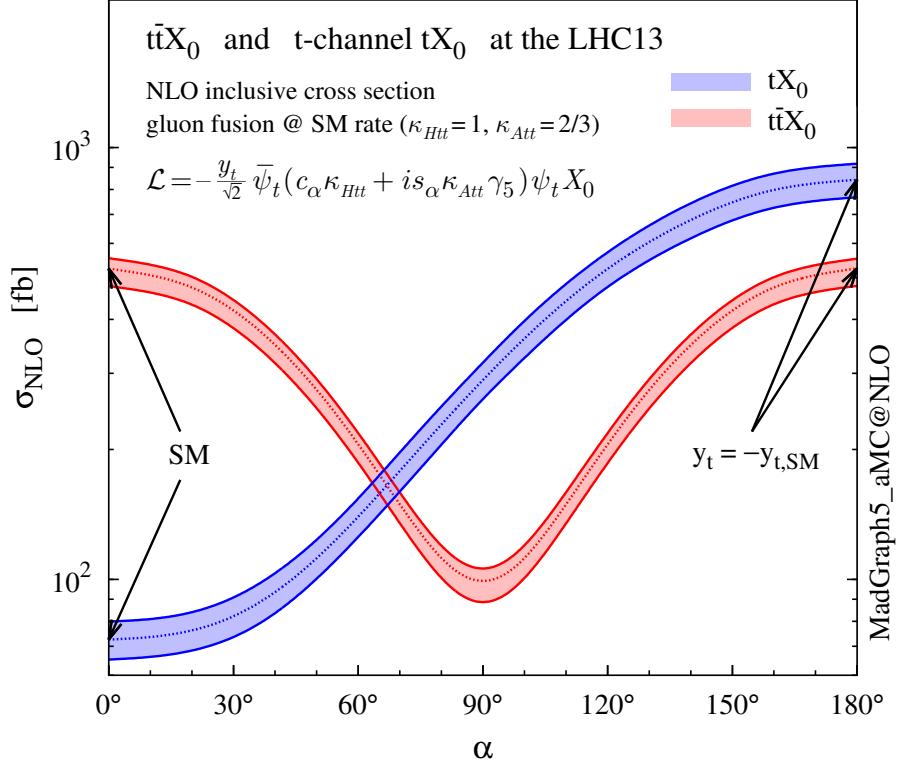


Figure 1.18: NLO cross sections for t-channel tX_0 (blue) and $t\bar{t}X_0$ (red) associated production processes as a function of the CP-mixing angle α . X_0 is a generic spin-0 particle with top quark CP-violating coupling [47].

1002 Figure 1.19 shows the NLO cross sections for t-channel tX_0 (blue), $t\bar{t}X_0$ (red)
 1003 associated production and for the combined $tWX_0+t\bar{t}X_0+interference$ (orange) as
 1004 a function of the CP-mixing angle. It is clear that the effect of the interference in the
 1005 combined case is the lifting of the degeneracy present in the $t\bar{t}X_0$ production. The
 1006 constructive interference enhances the cross section from about 500 fb at SM ($\alpha = 0$)
 1007 to about 600 fb ($\alpha = 180^\circ \rightarrow y_t = -y_{t,SM}$).

1008 An analysis combining tHq and tHW processes will be made in this thesis taking
 1009 advantage of the sensitivity improvement.

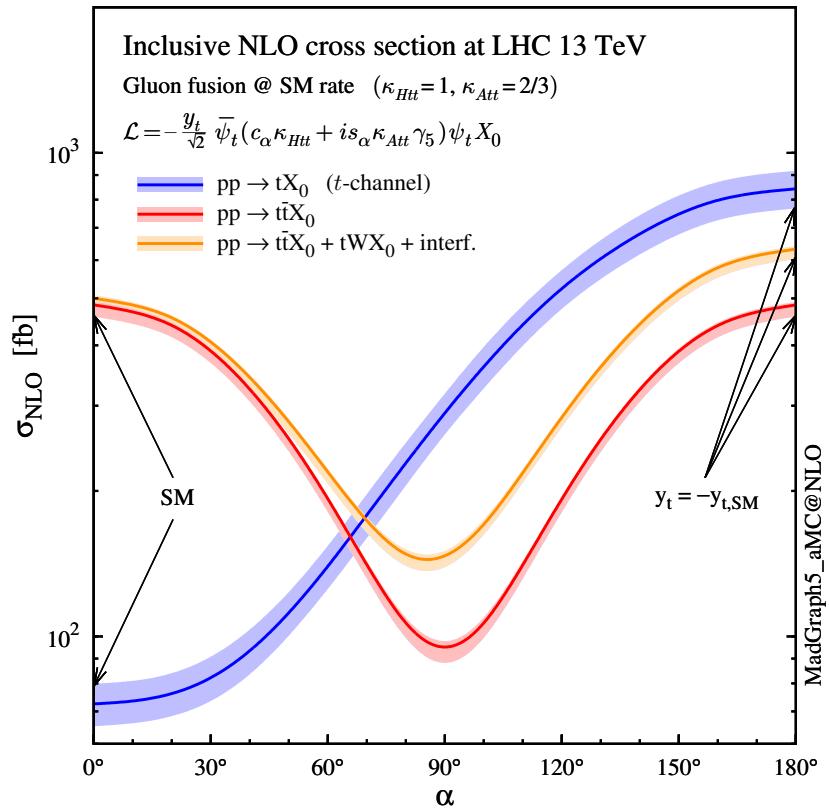


Figure 1.19: NLO cross sections for t-channel tX_0 (blue), $t\bar{t}X_0$ (red) associated production processes and combined $tWX_0 + t\bar{t}X_0$ (including interference) production as a function of the CP-mixing angle α [47].

1010 **Chapter 2**

1011 **The CMS experiment at the LHC**

1012 **2.1 Introduction**

1013 Located on the Swiss-French border, the European Council for Nuclear Research
1014 (CERN) is the largest scientific organization leading particle physics research. About
1015 13000 people in a broad range of roles including users, students, scientists, engineers,
1016 among others, contribute to the data taking and analysis, with the goal of unveiling
1017 the secrets of nature and revealing the fundamental structure of the universe. CERN
1018 is also the home of the Large Hadron Collider (LHC), the largest particle accelerator
1019 around the world, where protons (or heavy ions) traveling close to the speed of light,
1020 are made to collide. These collisions open a window to investigate how particles (and
1021 their constituents if they are composite) interact with each other, providing clues
1022 about the laws of nature. This chapter presents an overview of the LHC structure
1023 and operation. A detailed description of the CMS detector is offered, given that the
1024 data used in this thesis have been taken with this detector.

1025 2.2 The LHC

1026 With 27 km of circumference, the LHC is currently the most powerful circular accelerator
 1027 in the world. It is installed in the same tunnel where the Large Electron-Positron
 1028 (LEP) collider was located, taking advantage of the existing infrastructure. The LHC
 1029 is part of the CERN's accelerator complex composed of several successive accelerat-
 1030 ing stages before the particles are injected into the LHC ring where they reach their
 1031 maximum energy (see Figure 2.1).

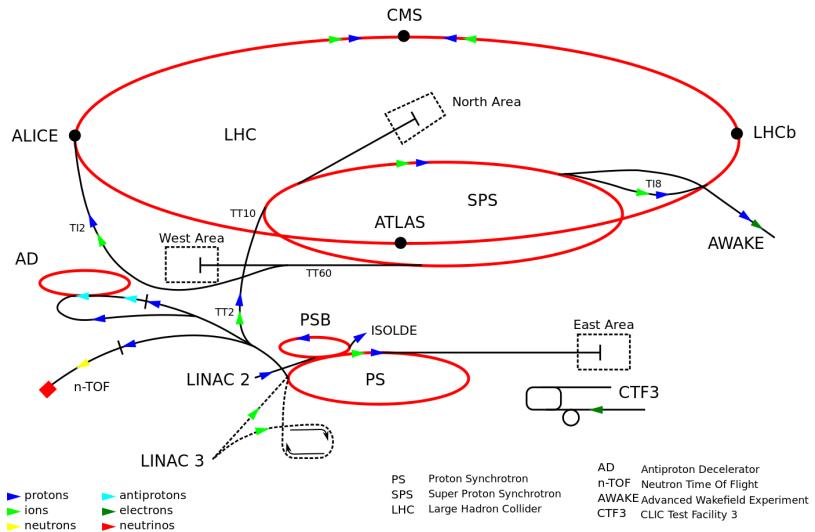


Figure 2.1: CERN accelerator complex. Blue arrows show the path followed by protons along the acceleration process [61].

1032 The LHC runs in three collision modes depending on the particles being acceler-
 1033 ated

- 1034 • Proton-Proton collisions (pp) for multiple physics experiments.
- 1035 • Lead-Lead collisions ($Pb-Pb$) for heavy ion experiments.
- 1036 • Proton-Lead collisions ($p-Pb$) for quark-gluon plasma experiments.

1037 In this thesis only pp collisions will be considered.

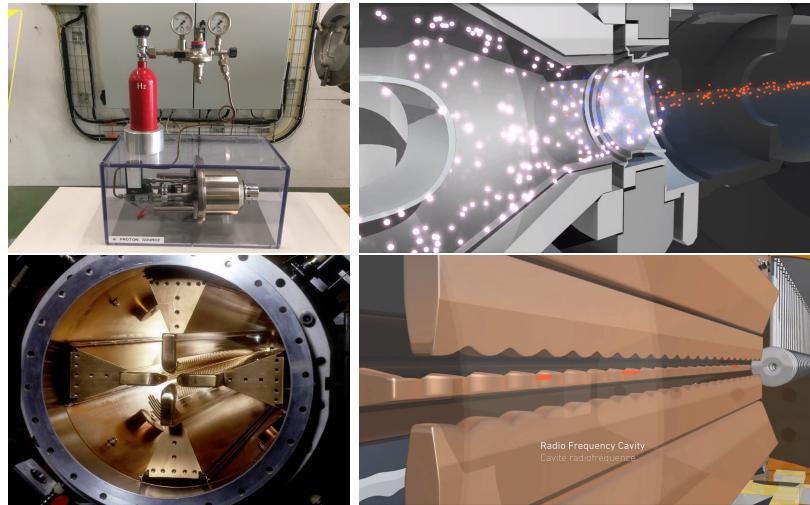


Figure 2.2: LHC protons source and the first acceleration stage. Top: the bottle contains hydrogen gas which is injected into the metal cylinder (white dots) to be broken down into electrons (blue dots) and protons (red dots); Bottom: the obtained protons are directed towards the radio frequency quadrupole which perform the first acceleration, focus the beam and create the bunches of protons. Left images are real pictures while right images are drawings [65, 66].

1038 Collection of protons starts with hydrogen atoms taken from a bottle, containing
 1039 hydrogen gas, and injecting them in a metal cylinder; hydrogen atoms are broken
 1040 down into electrons and protons by an intense electric field (see Figure 2.2 top).
 1041 The resulting protons leave the metal cylinder towards a radio frequency quadrupole
 1042 (RFQ) that focus the beam, accelerates the protons and creates the packets of protons
 1043 called bunches. In the RFQ, an electric field is generated by a RF wave at a frequency
 1044 that matches the resonance frequency of the cavity where the electrodes are contained.
 1045 The beam of protons traveling on the RFQ axis experiences an alternating electric
 1046 field gradient that generates the focusing forces.

1047 In order to accelerate the protons, a longitudinal time-varying electric field com-
 1048 ponent is added to the system; it is done by giving the electrodes a sine-like profile as
 1049 shown in Figure 2.2 bottom. By matching the speed and phase of the protons with
 1050 the longitudinal electric field the bunching is performed; protons synchronized with

1051 the RFQ (synchronous protons) do not feel an accelerating force, but those protons in
 1052 the beam that have more (or less) energy than the synchronous proton (asynchronous
 1053 protons) will feel a decelerating (accelerating) force; therefore, asynchronous protons
 1054 will oscillate around the synchronous ones forming bunches of protons [63]. From the
 1055 RFQ protons emerge with energy 750 keV in bunches of about 1.15×10^{11} protons [64].

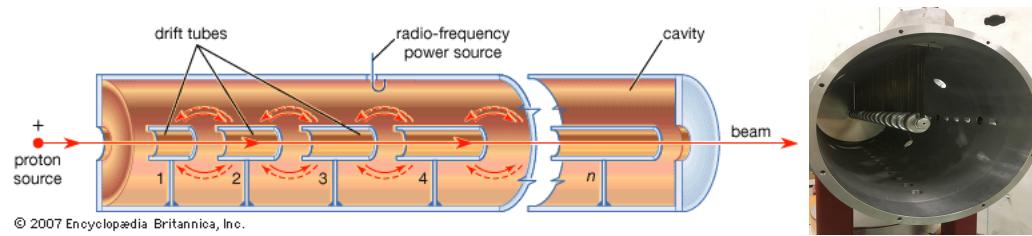


Figure 2.3: Left: drawing of the LINAC2 accelerating system at CERN. Electric fields generated by radio frequency (RF) create acceleration and deceleration zones inside the cavity; deceleration zones are blocked by drift tubes where quadrupole magnets focus the proton beam. Right: picture of a real RF cavity [67].

1056 Proton bunches coming from the RFQ go to the linear accelerator 2 (LINAC2)
 1057 where they are accelerated to reach 50 MeV energy. In the LINAC2 stage, acceleration
 1058 is performed using electric fields generated by radio frequency which create zones
 1059 of acceleration and deceleration as shown in Figure 2.3. In the deceleration zones,
 1060 the electric field is blocked using drift tubes where protons are free to drift while
 1061 quadrupole magnets focus the beam.

1062 The beam coming from LINAC2 is injected into the proton synchrotron booster
 1063 (PSB) to reach 1.4 GeV in energy. The next acceleration is provided at the proton
 1064 synchrotron (PS) up to 26 GeV, followed by the injection into the super proton
 1065 synchrotron (SPS) where protons are accelerated to 450 GeV. Finally, protons are
 1066 injected into the LHC where they are accelerated to the target energy of 6.5 TeV.

1067 PSB, PS, SPS and LHC accelerate protons using the same RF acceleration tech-
 1068 nique described before.

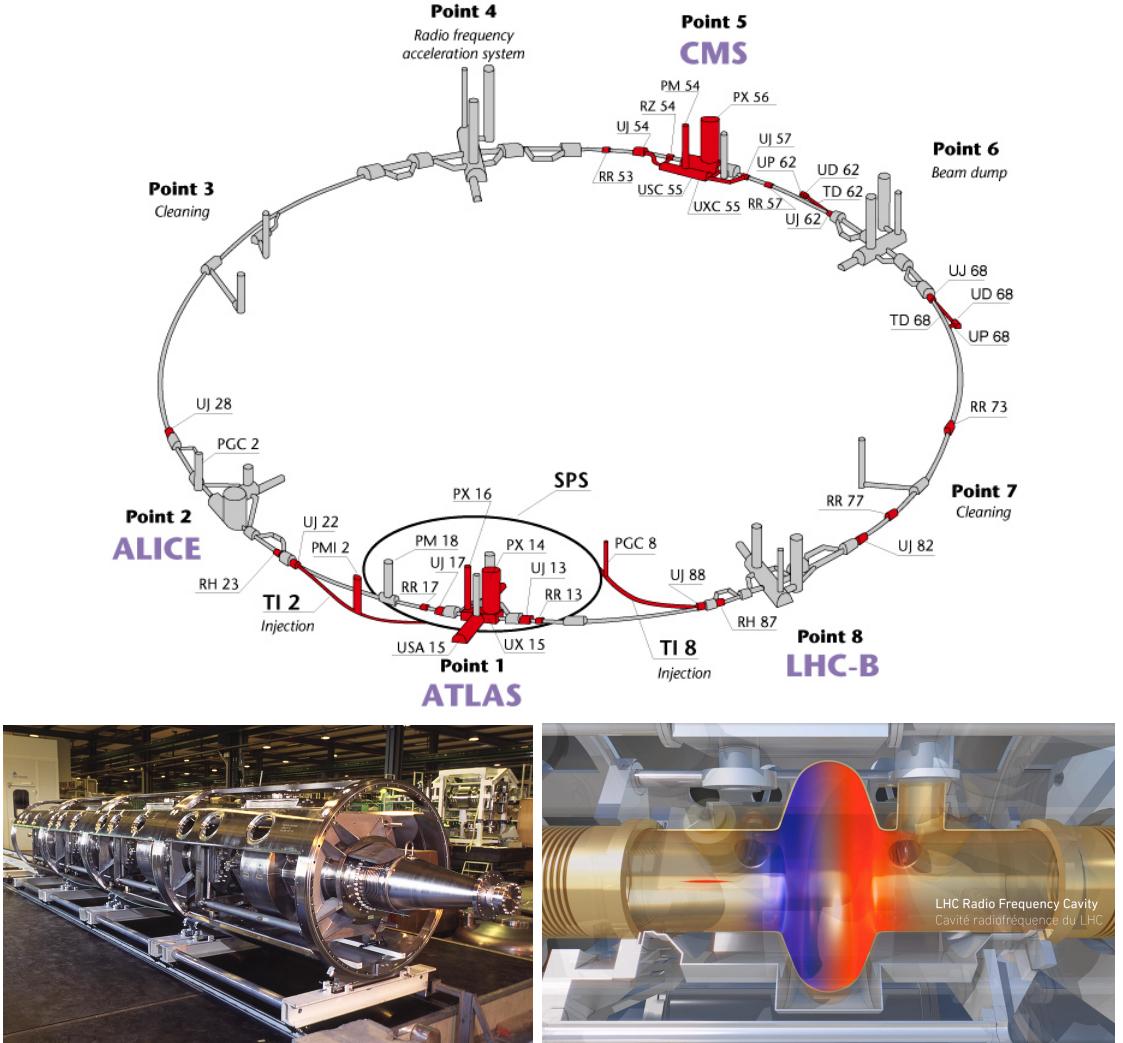


Figure 2.4: Top: LHC layout. The red zones indicate the infrastructure additions to the LEP installations, built to accommodate the ATLAS and CMS experiments which exceed the size of the former experiments located there [62]. Bottom: LHC RF cavities. A module (left picture) accommodates 4 cavities, each like the one in the right drawing, that accelerate protons and preserve the bunch structure of the beam. The color gradient pattern represents the strength of the electric field inside the cavity; red zone corresponds to the maximum of the oscillation electric field while the blue zone corresponds to the minimum. [66, 68]

1069 The LHC has a system of 16 RF cavities located in the so-called point 4, as
 1070 shown in Figure 2.4 top, tuned at a frequency of 400 MHz. The bottom side of
 1071 Figure 2.4 shows a picture of a RF module composed of 4 RF cavities working in a
 1072 superconducting state at 4.5 K; also, a representation of the accelerating electric field

1073 that accelerates the protons in the bunch is shown. The maximum of the oscillating
 1074 electric field (red region) picks the proton bunches at the entrance of the cavity
 1075 and keeps accelerating them through the whole cavity. The protons are carefully
 1076 timed so that in addition to the acceleration effect the bunch structure of the beam
 1077 is preserved.

1078 While protons are accelerated in one section of the LHC ring, where the RF cavities
 1079 are located, in the rest of their path they have to be kept in the curved trajectory
 1080 defined by the LHC ring. Technically, LHC is not a perfect circle; RF, injection, beam
 1081 dumping, beam cleaning and sections before and after the experimental points where
 1082 protons collide are all straight sections. In total, there are 8 arcs 2.45 km long each
 1083 and 8 straight sections 545 m long each. In order to curve the proton's trajectory in
 1084 the arc sections, superconducting dipole magnets are used.

1085 Inside the LHC ring, there are two proton beams traveling in opposite directions
 1086 in two separated beam pipes; the beam pipes are kept at ultra-high vacuum ($\sim 10^{-9}$
 1087 Pa) to ensure that there are no particles that interact with the proton beams. The
 1088 superconducting dipole magnets used in LHC are made of a NbTi alloy, capable of
 1089 transporting currents of about 12000 A when cooled at a temperature below 2K using
 1090 liquid helium (see Figure 2.5).

1091 Protons in the arc sections of LHC feel a centripetal force exerted by the dipole
 1092 magnets; the magnitude of magnetic field needed to keep the protons in the LHC
 1093 curved trayectomy can be found assuming that protons travel at $v \approx c$, using the
 1094 standard values for proton mass and charge and the LHC radius, as

$$F_m = \frac{mv^2}{r} = qBv \rightarrow B = 8.33T \quad (2.1)$$

1095 which is about 100000 times the Earth's magnetic field. A representation of the

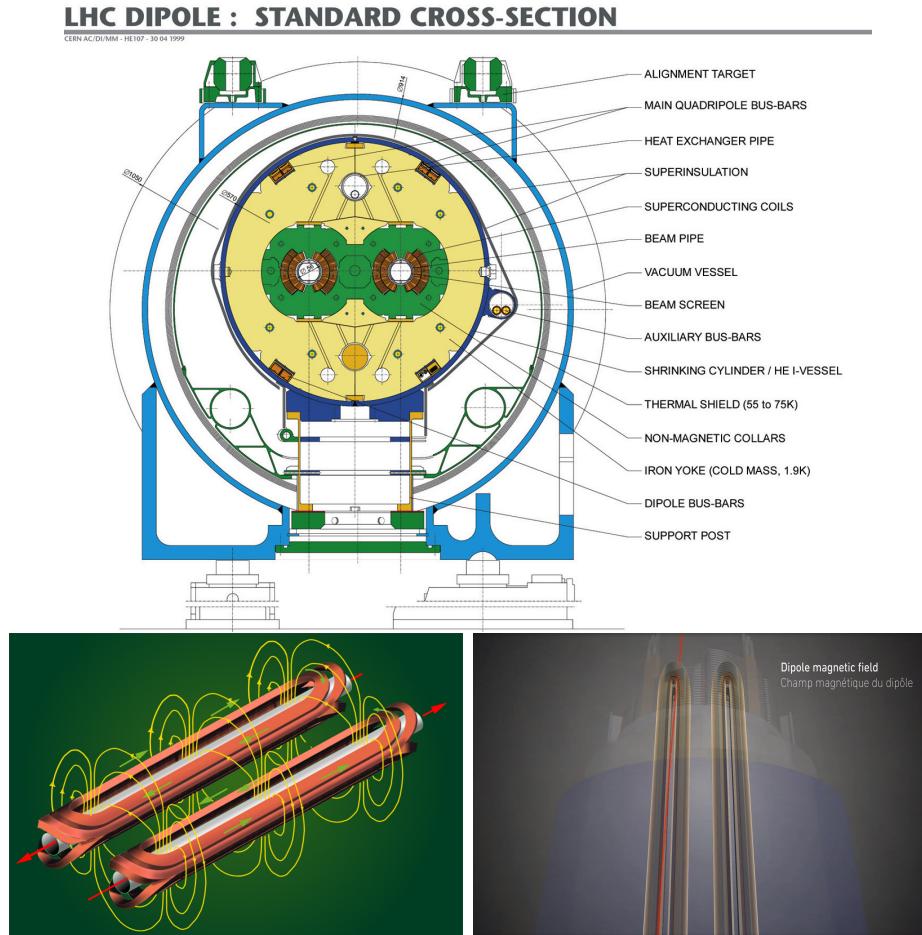


Figure 2.5: Top: LHC dipole magnet transverse view; cooling, shielding and mechanical support are indicated. Bottom left: Magnetic field generated by the dipole magnets; note that the direction of the field inside one beam pipe is opposite with respect to the other beam pipe which guarantee that both proton beams are curved in the same direction towards the center of the ring. The effect of the dipole magnetic field on the proton beam is represented on the bottom right side [66, 69, 70].

1096 magnetic field generated by the dipole magnets is shown on the bottom left side of
 1097 Figure 2.5. The bending effect of the magnetic field on the proton beam is shown on
 1098 the bottom right side of Figure 2.5. Note that the dipole magnets are not curved;
 1099 the arc section of the LHC ring is composed of straight dipole magnets of about 15
 1100 m. In total there are 1232 dipole magnets along the LHC ring.

1101 In addition to the bending of the beam trajectory, the beam has to be focused. The

focusing is performed by quadrupole magnets installed in a different straight section; in total 858 quadrupole magnets are installed along the LHC ring. Other effects like electromagnetic interaction among bunches, interaction with electron clouds from the beam pipe, the gravitational force on the protons, differences in energy among protons in the same bunch, among others, are corrected using sextupole and other magnetic multipoles.

The two proton beams inside the LHC ring are made of bunches with a cylindrical shape of about 7.5 cm long and about 1 mm in diameter; when bunches are close to the interaction point (IP), the beam is focused up to a diameter of about 16 μm in order to maximize the probability of collisions between protons. The number of collisions per second is proportional to the cross section of the bunches with the *luminosity* (L) as the proportionality factor, thus, the luminosity can be calculated using

$$L = fn \frac{N_1 N_2}{4\pi\sigma_x\sigma_y} \quad (2.2)$$

where f is the revolution frequency, n is the number of bunches per beam, N_1 and N_2 are the numbers of protons per bunch, σ_x and σ_y are the gaussian transverse sizes of the bunches. With the expected parameters, the LHC expected luminosity is about

$$f = \frac{v}{2\pi r_{LHC}} \approx \frac{3 \times 10^8 \text{ m/s}}{27 \text{ km}} \approx 11.1 \text{ kHz},$$

$$n = 2808$$

$$N_1 = N_2 \sim 1.5 \times 10^{11}$$

$$\sigma_x = \sigma_y = 16 \mu\text{m}$$

$$L = 1.28 \times 10^{34} \text{ cm}^{-2} \text{ s}^{-1} = 1.28 \times 10^{-5} \text{ fb}^{-1} \text{ s}^{-1} \quad (2.3)$$

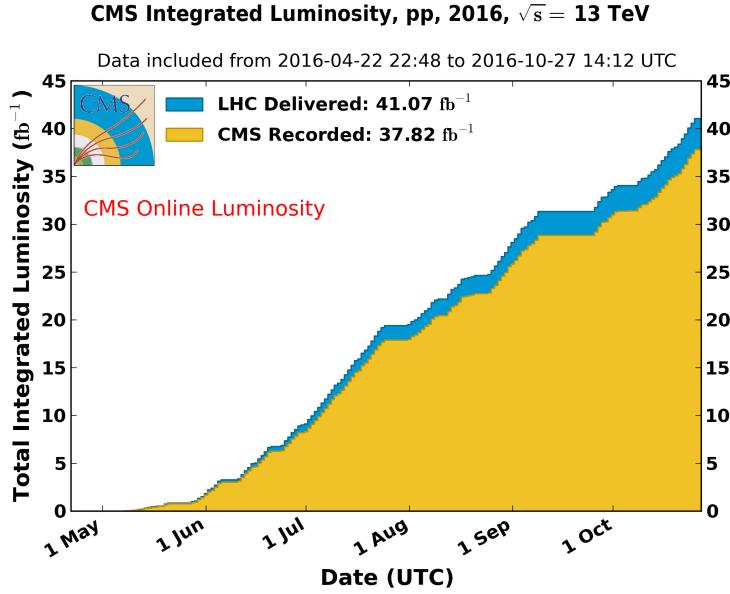


Figure 2.6: Integrated luminosity delivered by LHC and recorded by CMS during 2016. The difference between the delivered and the recorded luminosity is due to fails and issues occurred during the data taking in the CMS experiment [71].

1118 Luminosity is a fundamental aspect of LHC given that the bigger luminosity the
 1119 bigger number of collisions, which means that for processes with a very small cross
 1120 section the number of expected occurrences is increased and so the chances of being
 1121 detected. The integrated luminosity, collected by the CMS experiment during 2016
 1122 is shown in Figure 2.6; the data analyzed in this thesis corresponds to an integrated
 1123 luminosity of 35.9 fb^{-1} at a center of mass-energy $\sqrt{s} = 13 \text{ TeV}$.

1124 One way to increase L is increasing the number of bunches in the beam. Cur-
 1125 rently, the separation between two consecutive bunches in the beam is 7.5 m which
 1126 corresponds to a time separation of 25 ns. In the full LHC ring the allowed number
 1127 of bunches is $n = 27\text{km}/7.5\text{m} = 3600$; however, there are some gaps in the bunch pat-
 1128 tern intended for preparing the dumping and injection of the beam, thus, the proton
 1129 beams are composed of 2808 bunches.

1130 Once the proton beams reach the desired energy, they are brought to cross each

other producing pp collisions. The bunch crossing happens in precise places where the four LHC experiments are located, as seen in the top of Figure 2.7. In 2008 pp collisions of $\sqrt{s} = 7$ TeV were performed; the energy was increased to 8 TeV in 2012 and to 13 TeV in 2015.

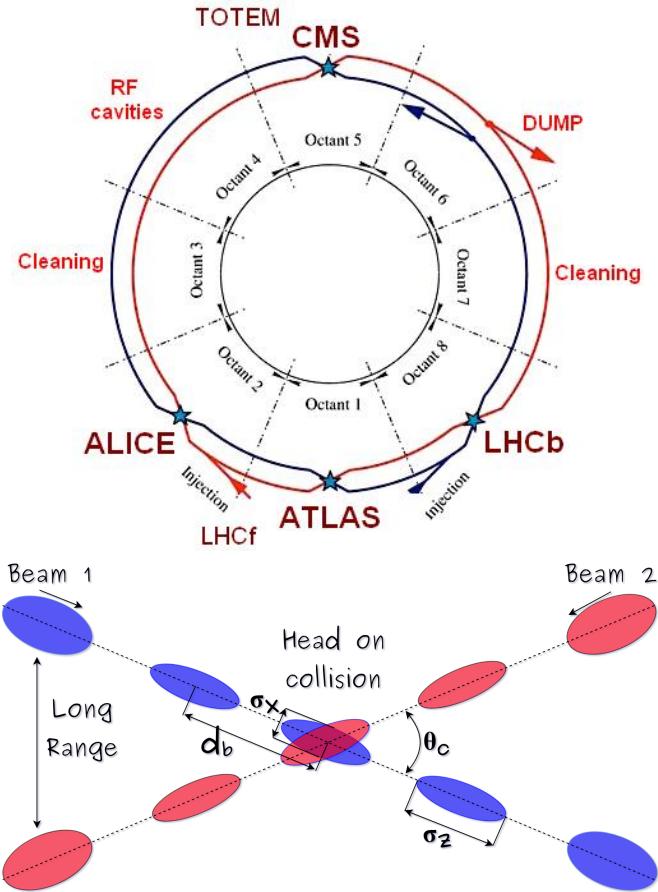


Figure 2.7: Top: LHC interaction points. Bunch crossing occurs where the LHC experiments are located [72]. Sections indicated as cleaning are dedicated to collimate the beam in order to protect the LHC ring from collisions with protons in very spreaded bunches. Bottom: bunch crossing scheme. Since the bunch crossing is not perfectly head-on, the luminosity is reduced in a factor of 17%; adapted from Reference [86].

The CMS and ATLAS experiments are multi-purpose experiments, hence, they are enabled to explore physics in any of the LHC collision modes. LHCb experiment is optimized to explore bottom quark physics, while ALICE is optimized for heavy

1138 ion collisions searches; TOTEM and LHCf are dedicated to forward physics studies;
 1139 MoEDAL (not indicated in the Figure) is intended for monopole or massive pseudo
 1140 stable particles searches.

1141 At the IP there are two interesting details that need to be addressed. The first
 1142 one is that the bunch crossing does not occur head-on but at a small crossing angle θ_c
 1143 (280 μ rad in CMS and ATLAS) as shown in the bottom side of Figure 2.7, affecting
 1144 the overlapping between bunches; the consequence is a reduction of about 17% in
 1145 the luminosity (represented by a factor not included in eqn. 2.2). The second one
 1146 is the occurrence of multiple pp collisions in the same bunch crossing; this effect is
 1147 called *pileup* (PU). A fairly simple estimation of the PU follows from estimating the
 1148 probability of collision between two protons, one from each of the bunches in the
 1149 course of collision; it depends roughly on the ratio of proton size and the cross section
 1150 of the bunch in the IP, i.e.,

$$P(pp\text{-}collision) \sim \frac{d_{proton}^2}{\sigma_x \sigma_y} = \frac{(1\text{fm})^2}{(16\mu\text{m})^2} \sim 4 \times 10^{-21} \quad (2.4)$$

1151 however, there are $N = 1.15 \times 10^{11}$ protons in a bunch, thus the estimated number of
 1152 collisions in a bunch crossing is

$$PU = N^2 * P(pp\text{-}collision) \sim 50pp \text{ collision per bunch crossing}, \quad (2.5)$$

1153 about 20 of which are inelastic. A multiple pp collision event in a bunch crossing at
 1154 CMS is shown in Figure 2.8.

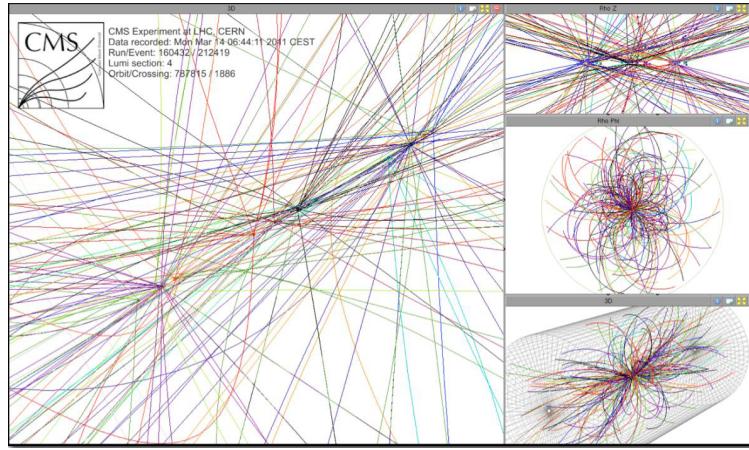


Figure 2.8: Multiple pp collision bunch crossing at CMS. [73].

1155 2.3 The CMS experiment

1156 CMS is a general-purpose detector designed to conduct research in a wide range
 1157 of physics from the standard model to new physics like extra dimensions and dark
 1158 matter. Located at Point 5 in the LHC layout as shown in Figure 2.4, CMS is
 1159 composed of several detection systems distributed in a cylindrical structure; in total,
 1160 CMS weights about 12500 tons in a very compact 21.6 m long and 14.6 m diameter
 1161 cylinder. It was built in 15 separate sections at the ground level and lowered to the
 1162 cavern individually to be assembled. A complete and detailed description of the CMS
 1163 detector and its components is given in Reference [74] on which this section is based.
 1164 Figure 2.9 shows the layout of the CMS detector. The detection system is composed
 1165 of (from the innermost to the outermost)

- 1166 • Pixel detector.
- 1167 • Silicon strip tracker.
- 1168 • Preshower detector.
- 1169 • Electromagnetic calorimeter.

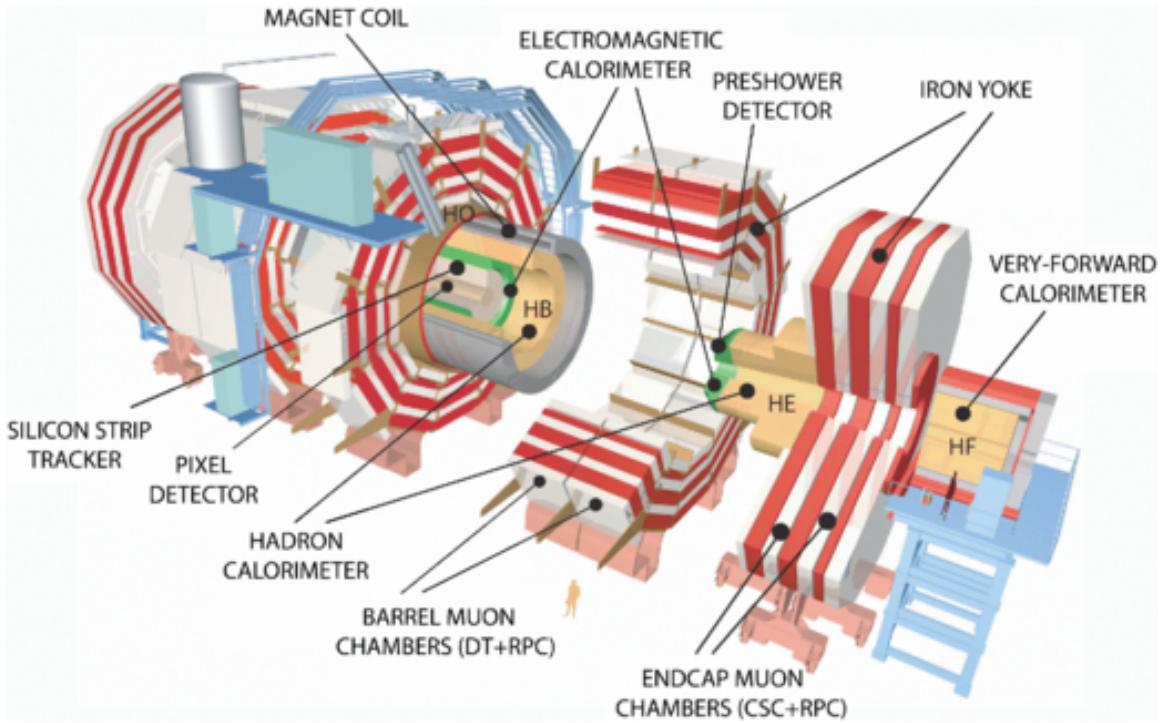


Figure 2.9: Layout of the CMS detector. The several subdetectors are indicated. The central region of the detector is referred as the barrel section while the endcaps are referred as the forward sections. [75].

1170 • Hadronic calorimeter.

1171 • Muon chambers (barrel and endcap)

1172 The central region of the detector is commonly referred as the barrel section while
 1173 the endcaps are referred as the forward sections of the detector; thus, each subdetector
 1174 is composed of a barrel section and a forward section.

1175 When a pp collision happens inside the CMS detector, many different particles are
 1176 produced, but only some of them live long enough to be detected; they are electrons,
 1177 photons, pions, kaons, protons, neutrons and muons; neutrinos are not detected by
 1178 the CMS detector. Thus, the CMS detector was designed to detect those particles and
 1179 measure their properties. Figure 2.10 shows a transverse slice of the CMS detector.
 1180 The silicon tracker (pixel detector + strip tracker) is capable to register the track of

1181 the charged particles traversing it, while calorimeters (electromagnetic and hadronic)
 1182 measure the energy of the particles that are absorbed by their materials. Considering
 1183 the detectable particles, mentioned above, emerging from the IP, a basic description
 1184 of the detection process is as follows.

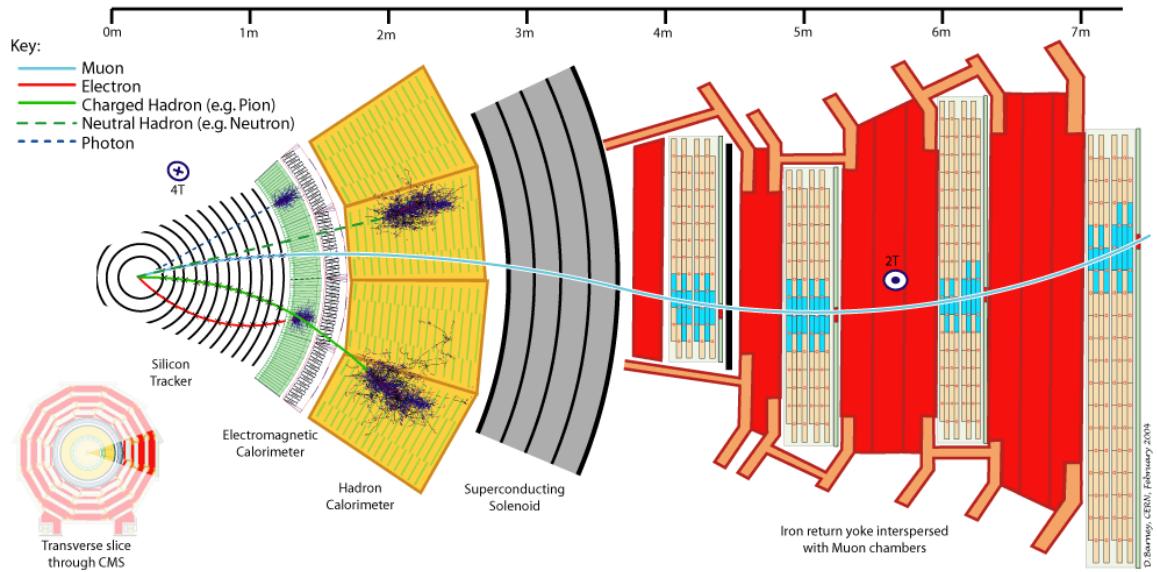


Figure 2.10: CMS detector transverse slice [76].

1185 A muon emerging from the IP, will create a track on the silicon tracker and on
 1186 the muon chambers. The design of the CMS detector is driven by the requirements
 1187 on the identification, momentum resolution and unambiguous charge determination
 1188 of the muons; therefore, a large bending power is provided by the solenoid magnet
 1189 made of superconducting cable capable of generating a 3.8 T magnetic field. The
 1190 muon track is bent twice since the magnetic field inside the solenoid is directed along
 1191 the z -direction but outside its direction is reversed. Muons interact very weakly with
 1192 the calorimeters, therefore, it is not absorbed but escape away from the detector.

1193 An electron emerging from the IP will create a track along the tracker which will
 1194 be bent due to the presence of the magnetic field, later, it will be absorbed in the
 1195 electromagnetic calorimeter where its energy is measured.

1196 A photon will not leave a track because it is neutral, but it will be absorbed in
 1197 the electromagnetic calorimeter.

1198 A neutral hadron, like the neutron, will not leave a track either but it will lose a
 1199 small amount of its energy during its passage through the electromagnetic calorimeter
 1200 and then it will be absorbed in the hadron calorimeter depositing the rest of its energy.

1201 A charged hadron, like the proton or π^\pm , will leave a curved track on the silicon
 1202 tracker, some of its energy in the electromagnetic calorimeter and finally will be
 1203 absorbed in the hadronic calorimeter.

1204 A more detailed description of each detection system will be presented in the
 1205 following sections.

1206 2.3.1 CMS coordinate system

1207 The coordinate system used by CMS is centered on the geometrical center of the
 1208 detector which is the nominal IP as shown in Figure 2.11¹. The z -axis is parallel
 1209 to the beam direction, while the Y -axis pointing vertically upward, and the X -axis
 1210 pointing radially inward toward the center of the LHC.

1211 In addition to the common cartesian and cylindrical coordinate systems, two co-
 1212 ordinates are of particular utility in particle physics: rapidity (y) and pseudorapidity
 1213 (η), defined in connection to the polar angle θ , energy and longitudinal momentum
 1214 component (momentum along the z -axis) according to

$$y = \frac{1}{2} \ln \frac{E + p_z}{E - p_z} \quad \eta = -\ln \left(\tan \frac{\theta}{2} \right) \quad (2.6)$$

1215 Rapidity is related to the angle between the XY -plane and the direction in which
 1216 the products of a collision are emitted; it has the nice property that the difference

¹ Not all the pp interaction occur at the nominal IP because of the bunch lenght, therefore, each pp collision has its own IP location

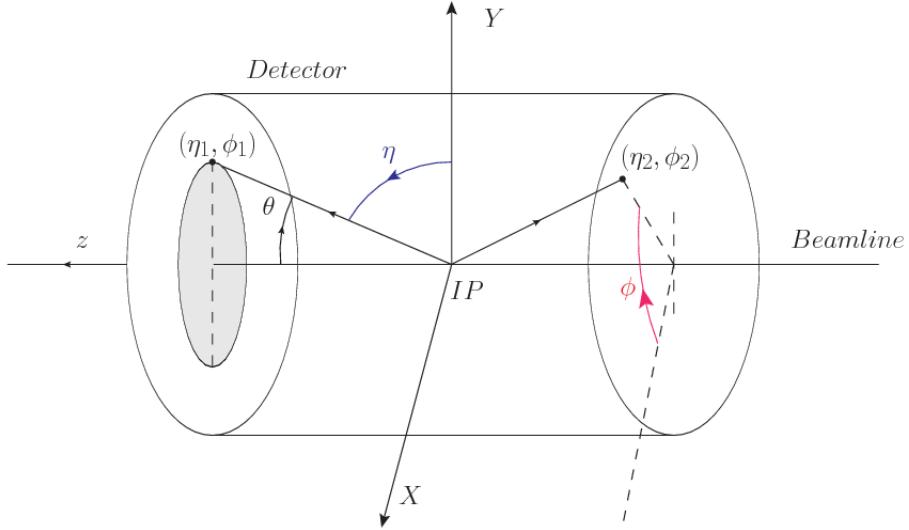


Figure 2.11: CMS detector coordinate system.

1217 between the rapidities of two particles is invariant with respect to Lorentz boosts
 1218 along the z -axis, hence, data analysis becomes more simple when based on rapid-
 1219 ity; however, it is not simple to measure the rapidity of highly relativistic particles,
 1220 as those produced after pp collisions. Under the highly relativistic motion approxi-
 1221 mation, y can be rewritten in terms of the polar angle, concluding that rapidity is
 1222 approximately equal to the pseudorapidity defined above, i.e., $y \approx \eta$. Note that η
 1223 is easier to measure than y given the direct relationship between the former and the
 1224 polar angle.

1225 The angular distance between two objects in the detector (ΔR) is commonly used
 1226 to judge the isolation of those object; it is defined in terms of their coordinates (η_1, ϕ_1) ,
 1227 (η_2, ϕ_2) as

$$\Delta R = \sqrt{(\Delta\eta)^2 - (\Delta\phi)^2} \quad (2.7)$$

1228 2.3.2 Tracking system

1229 The CMS tracking system is designed to provide a precise measurement of the trajectories (*track*) followed by the charged particles created after the pp collisions; also, the
 1230 precise reconstruction of the primary and secondary origins of the tracks (*vertices*) is
 1231 expected in an environment where, each 25 ns, the bunch crossing produces about 20
 1232 inelastic collisions and about 1000 particles.
 1233

1234 Physics requirements guiding the tracking system performance include the precise
 1235 characterization of events involving gauge bosons, W and Z, and their leptonic
 1236 decays for which an efficient isolated lepton and photon reconstruction is of capital
 1237 importance, given that isolation is required to suppress background events to a level
 1238 that allows observations of interesting processes like Higgs boson decays or beyond
 1239 SM events.

1240 The ability to identify and reconstruct *b*-jets and B-hadrons within these jets is also
 1241 a fundamental requirement, achieved through the ability to reconstruct accurately
 1242 displaced vertices, given that *b*-jets are part of the signature of top quark physics, like
 1243 the one treated in this thesis.

1244 An schematic view of the CMS tracking system is shown in Figure 2.12

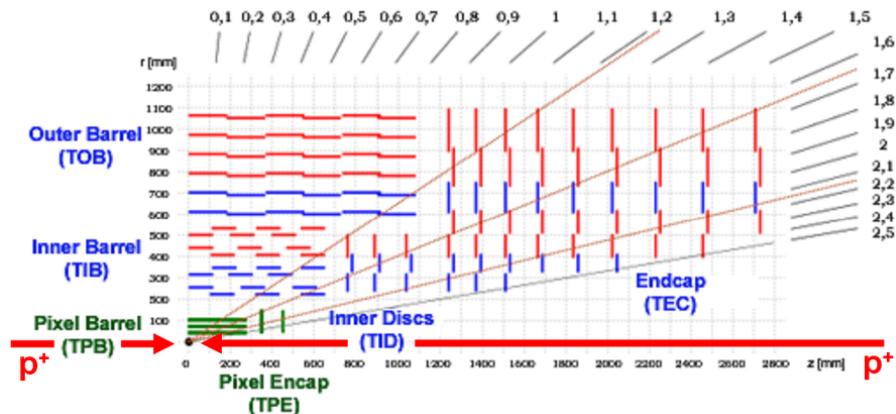


Figure 2.12: CMS tracking system schematic view [78].

1245 In order to satisfy these performance requirements, the tracking system uses two
 1246 different detector subsystems arranged in concentric cylindrical volumes, the pixel
 1247 detector and the silicon strip tracker; the pixel detector is located in the high particle
 1248 density region ($r < 20\text{cm}$) while the silicon strip tracker is located in the medium and
 1249 lower particle density regions $20\text{cm} < r < 116\text{cm}$.

1250 **Pixel detector**

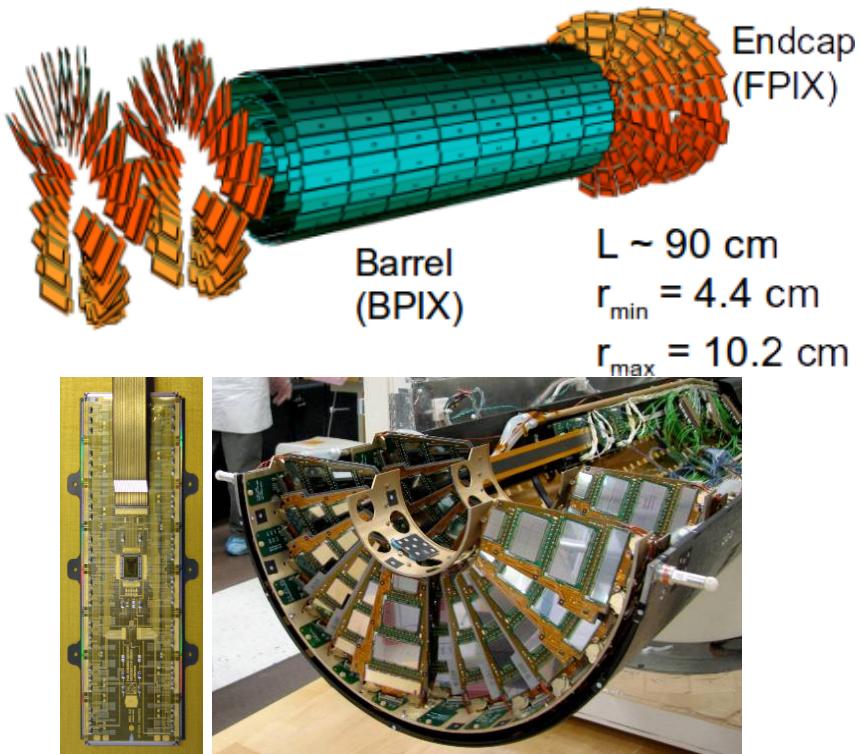


Figure 2.13: CMS pixel detector. Top: schematic view; Bottom: pictures of a barrel(BPIX) module(left) and forward modules(right) [74].

1251 The pixel detector was replaced during the 2016-2017 extended year-end technical
 1252 stop, due to the increasingly challenging operating conditions like the higher particle
 1253 flux and more radiation harsh environment, among others. The new one is responding
 1254 as expected, reinforcing its crucial role in the successful way to fulfill the new LHC

physics objectives after the discovery of the Higgs boson. Since the data sets used in this thesis were produced using the previous version of the pixel detector, it will be the subject of the description in this section. The last chapter of this thesis is dedicated to describe my contribution to the *Forward Pixel Phase 1 upgrade*.

The pixel detector was composed of 1440 silicon pixel detector modules organized in three-barrel layers in the central region (BPix) and two disks in the forward region (FPix) as shown in the top side of Figure 2.13; it was designed to record efficiently and with high precision, up to $20\mu\text{m}$ in the XY -plane and $20\mu\text{m}$ in the z -direction, the first three space-points (*hits*) nearest to the IP region in the range $|\eta| \leq 2.5$. The first barrel layer was located at a radius of 44 mm from the beamline, while the third layer was located at a radius of 102 mm, closer to the strip tracker inner barrel layer (see Section 2.3.3) in order to reduce the rate of fake tracks. The high granularity of the detector is represented in its about 66 Mpixels, each of size $100 \times 150\mu\text{m}^2$. The transverse momentum resolution of tracks can be measured with a resolution of 1-2% for muons of $p_T = 100$ GeV.

A charged particle passing through the pixel sensors produce ionization in them, giving energy for electrons to be removed from the silicon atoms, hence, creating electron-hole pairs. The collection of charges in the pixels generates an electrical signal that is read out by an electronic readout chip (ROC); each pixel has its own electronics which amplifies the signal. Combining the signal from the pixels activated by a traversing particle in the several layers of the detector allows one to reconstruct the particle's trajectory in 3D.

Commonly, the charge produced by traversing of a particle is collected by and shared among several pixels; by interpolating between pixels, the spatial resolution is improved. In the barrel section the charge sharing in the $r\phi$ -plane is due to the Lorentz effect. In the forward pixels the charge sharing is enhanced by arranging the

1281 blades in the turbine-like layout as shown in Figure 2.13 bottom left.

1282 **2.3.3 Silicon strip tracker**

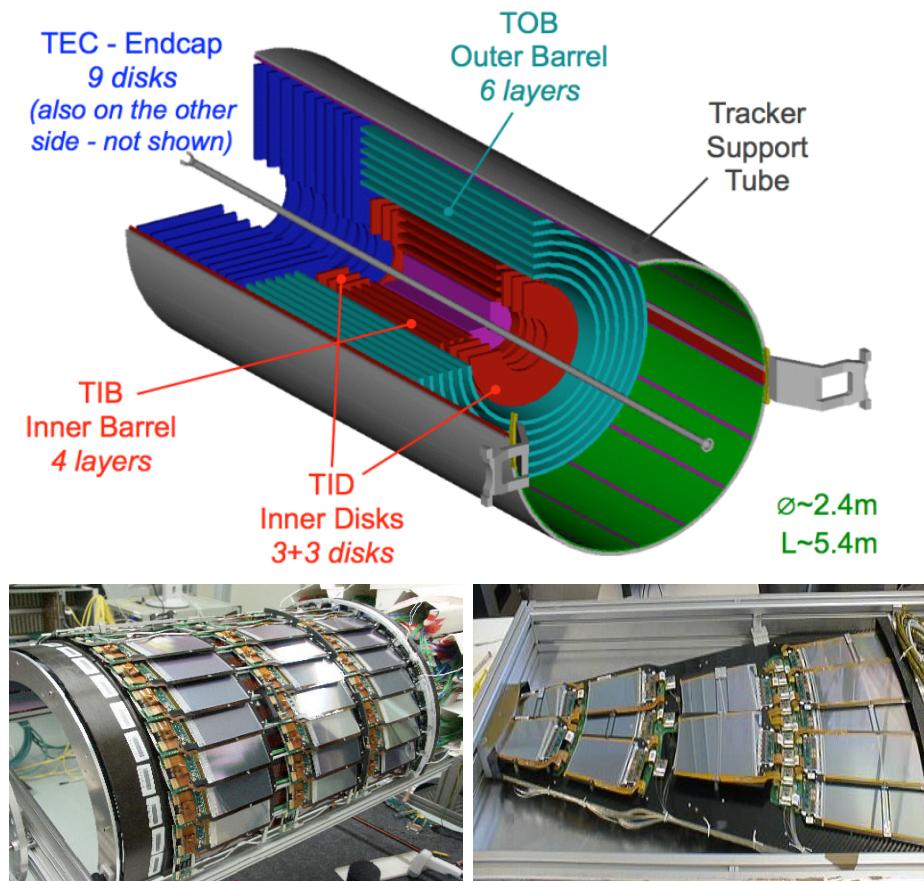


Figure 2.14: Top: CMS Silicon Strip Tracker (SST) schematic view. The SST is composed of the tracker inner barrel (TIB), the tracker inner disks (TID), the tracker outer barrel (TOB) and the tracker endcaps (TEC). Each part is made of silicon strip modules; the modules in blue represent two modules mounted back-to-back and rotated in the plane of the module by a stereo angle of 120 mrad in order to provide a 3-D reconstruction of the hit positions. Bottom: pictures of the TIB (left) and TEC (right) modules [80–82].

1283 The silicon strip tracker (SST) is the second stage in the CMS tracking system.

1284 The top side of Figure 2.14 shows a schematic of the SST. The inner tracker region is

1285 composed of the tracker inner barrel (TIB) and the tracker inner disks (TID) covering

1286 the region $r < 55$ cm and $|z| < 118$ cm. The TIB is composed of 4 layers while the TID

1287 is composed of 3 disks at each end. The silicon sensors in the inner tracker are 320
 1288 μm thick, providing a resolution of about 13-38 μm in the $r\phi$ position measurement.

1289 The modules indicated in blue in the schematic view of Figure 2.14 are two mod-
 1290 ules mounted back-to-back and rotated in the plane of the module by a *stereo* angle
 1291 of 100 mrad; the hits from these two modules, known as *stereo hits*, are combined to
 1292 provide a measurement of the second coordinate (z in the barrel and r on the disks)
 1293 allowing the reconstruction of hit positions in 3-D.

1294 The outer tracker region is composed of the tracker outer barrel (TOB) and the
 1295 tracker endcaps (TEC). The six layers of the TOB offer coverage in the region $r > 55$
 1296 cm and $|z| < 118$ cm, while the 9 disks of the TEC cover the region $124 < |z| < 282$
 1297 cm. The resolution offered by the outer tracker is about 13-38 μm in the $r\phi$ position
 1298 measurement. The inner four TEC disks use silicon sensors 320 μm thick; those in
 1299 the TOB and the outer three TEC disks use silicon sensors of 500 μm thickness. The
 1300 silicon strips run parallel to the z -axis and the distance between strips varies from 80
 1301 μm in the inner TIB layers to 183 μm in the inner TOB layers; in the endcaps the
 1302 wedge-shaped sensors with radial strips, whose pitch range between 81 μm at small
 1303 radii and 205 μm at large radii.

1304 The whole SST has 15148 silicon modules, 9.3 million silicon strips and cover a
 1305 total active area of about 198 m^2 .

1306 2.3.4 Electromagnetic calorimeter

1307 The CMS electromagnetic calorimeter (ECAL) is designed to measure the energy of
 1308 electrons and photons. It is composed of 75848 lead tungstate crystals which have a
 1309 short radiation length (0.89 cm) and fast response, since 80% of the light is emitted
 1310 within 25 ns; however, they are combined with Avalanche photodiodes (APDs) as

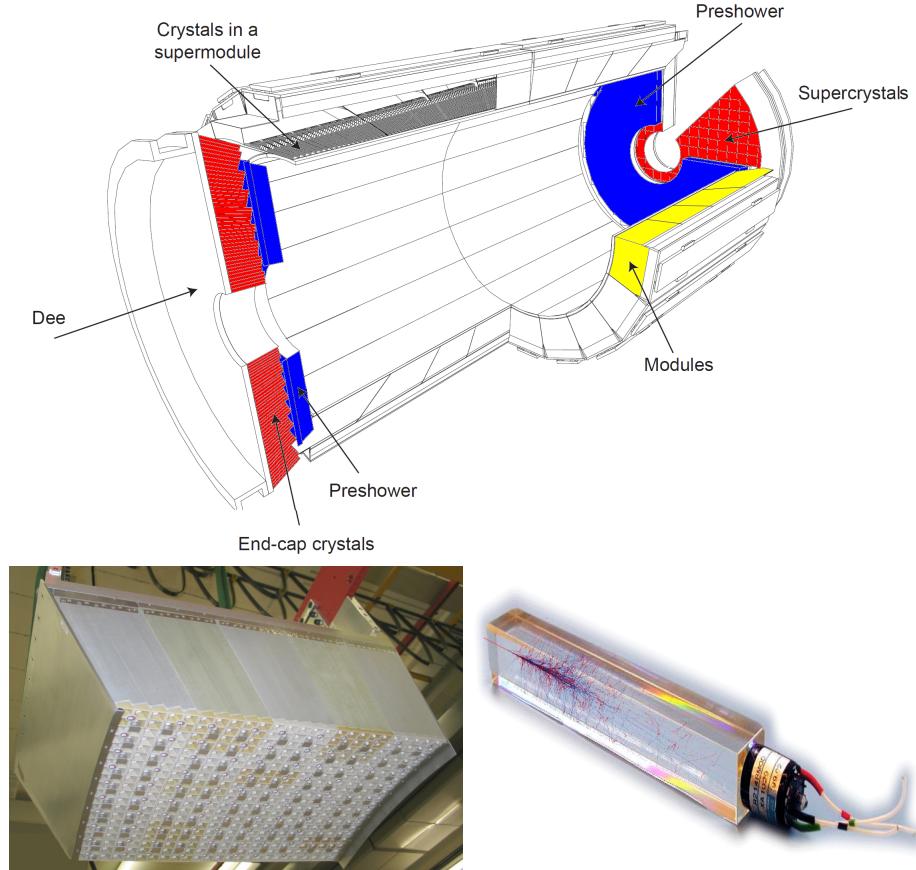


Figure 2.15: Top: CMS ECAL schematic view. Bottom: Module equipped with the crystals (left); ECAL crystal(right) with an artistic representation of an electromagnetic shower [74].

1311 photodetectors given that crystals themselves have a low light yield ($30\gamma/\text{MeV}$). A
 1312 schematic view of the ECAL is shown in Figure 2.15.

1313 Energy is measured when electrons and photons are absorbed by the crystals
 1314 which generates an electromagnetic *shower*, as seen in bottom right picture of the
 1315 Figure 2.15; the shower is seen as a *cluster* of energy which depending on the amount
 1316 of energy deposited can involve several crystals. The ECAL barrel (EB) covers the
 1317 region $|\eta| < 1.479$, using crystals of depth of 23 cm and $2.2 \times 2.2 \text{ cm}^2$ transverse
 1318 section; the ECAL endcap (EE) covers the region $1.479 < |\eta| < 3.0$ using crystals of
 1319 depth 22 cm and transverse section of $2.86 \times 2.86 \text{ cm}^2$; the photodetectors used are

vacuum phototriodes (VPTs). Each EE is divided in two structures called *Dees*.

The preshower detector (ES) is installed in front of the EE and covers the region $1.653 < |\eta| < 2.6$. The ES provides a precise measurement of the position of electromagnetic showers, which allows to distinguish electrons and photon signals from π^0 decay signals. The ES is composed of a layer of lead radiators followed by a layer of silicon strip sensors. The lead radiators initiate electromagnetic showers when reached by photons and electrons, then, the strip sensors measure the deposited energy and the transverse shower profiles. The full ES thickness is 20 cm.

2.3.5 Hadronic calorimeter

Hadrons are not absorbed by the ECAL² but by the hadron calorimeter (HCAL), which is made of a combination of alternating brass absorber layers and silicon photomultiplier(SiPM) layers; therefore, particles passing through the scintillator material produce showers, as in the ECAL, as a result of the inelastic scattering of the hadrons with the detector material. Since the particles are not absorbed in the scintillator, their energy is sampled; therefore the total energy is not measured but estimated from the energy clusters, which reduces the resolution of the detector. Brass was chosen as the absorber material due to its short interaction length ($\lambda_I = 16.42\text{cm}$) and its non-magnetivity. Figure 2.16 shows a schematic view of the CMS HCAL.

The HCAL is divided into four sections; the Hadron Barrel (HB), the Hadron Outer (HO), the Hadron Endcap (HE) and the Hadron Forward (HF) sections. The HB covers the region $0 < |\eta| < 1.4$, while the HE covers the region $1.3 < |\eta| < 3.0$. The HF, made of quartz fiber scintillator and steel as absorption material, covers the forward region $3.0 < |\eta| < 5.2$. Both the HB and HF are located inside the solenoid. The HO is placed outside the magnet as an additional layer of scintillators with the

² Most hadrons are not absorbed, but few low-energy ones might be.

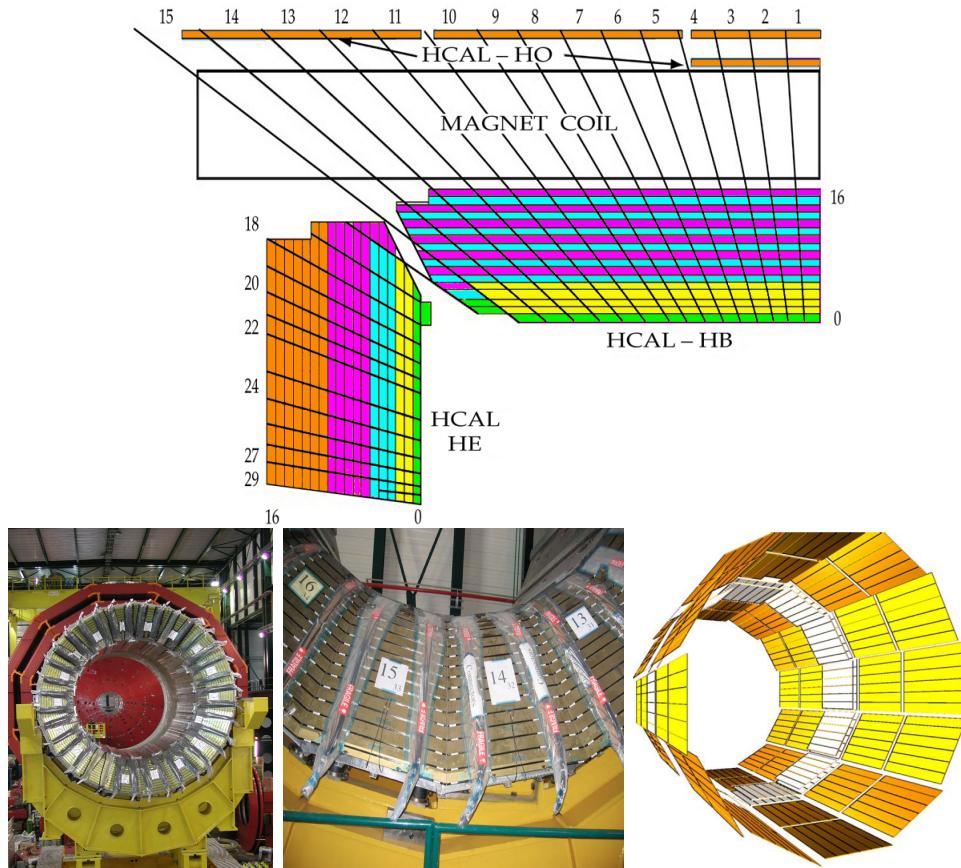


Figure 2.16: Top: CMS HCAL schematic view, the colors indicate the layers that are grouped into the same readout channels. Bottom: picture of a section of the HB; the absorber material is the golden region and scintillators are placed in between the absorber material (left and center). Schematic view of the HO (right). [83,84]

1344 purpose of measure the energy tails of particles passing through the HB and the
 1345 magnet (see Figure 2.16 top and bottom right).

1346 **2.3.6 Superconducting solenoid magnet**

1347 The superconducting magnet installed in the CMS detector is designed to provide
 1348 an intense and highly uniform magnetic field in the central part of the detector.
 1349 In fact, the tracking system takes advantage of the bending power of the magnetic
 1350 field to measure with precision the momentum of the particles that traverse it; the

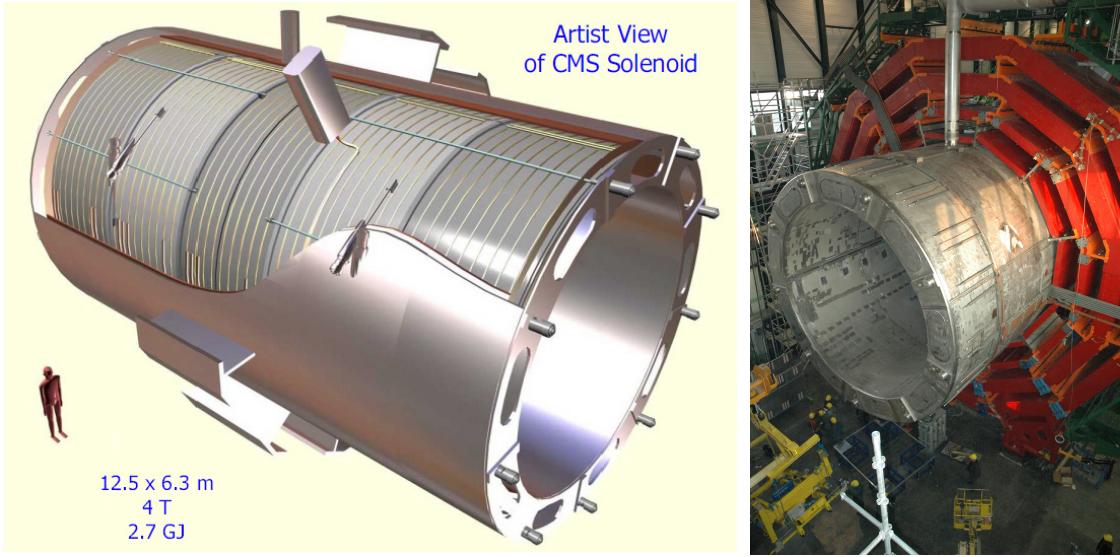


Figure 2.17: Artistic representation of the CMS solenoid magnet(left). The magnet is supported on an iron yoke (right) which also serves as the house of the muon detector and as mechanical support for the whole CMS detector [77].

1351 unambiguous determination of the sign for high momentum muons was a driving
 1352 principle during the design of the magnet. The magnet has a diameter of 6.3 m, a
 1353 length of 12.5 m and a cold mass of 220 t; the generated magnetic field reaches a
 1354 strength of 3.8T. Since it is made of Ni-Tb superconducting cable it has to operate at
 1355 a temperature of 4.7 K by using a helium cryogenic system; the current circulating in
 1356 the cables reaches 18800 A under normal running conditions. The left side of Figure
 1357 2.17 shows an artistic view of the CMS magnet, while the right side shows a transverse
 1358 view of the cold mass where the winding structure is visible.

1359 The yoke (see Figure 2.17), composed of 5 barrel wheels and 6 endcap disks made
 1360 of iron, serves not only as the media for magnetic flux return but also provides housing
 1361 for the muon detector system and structural stability to the full detector.

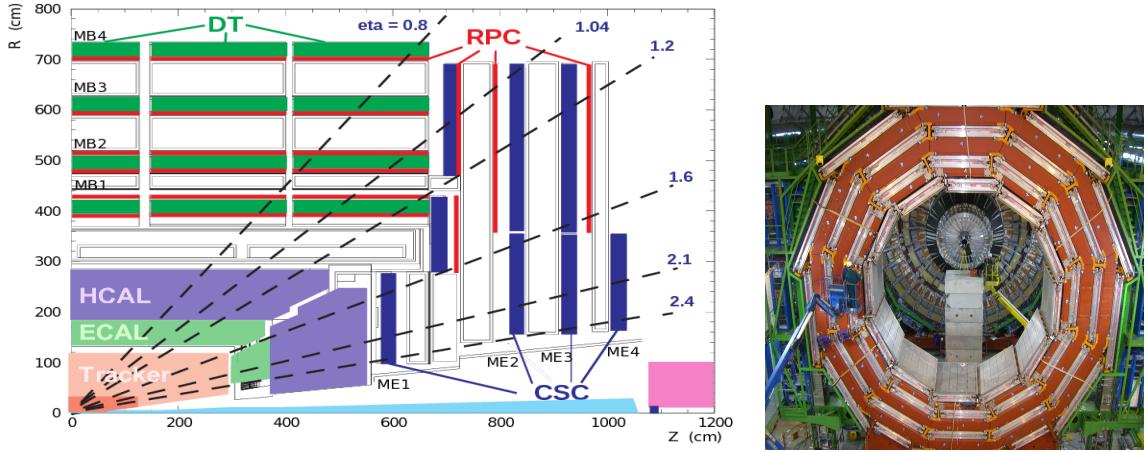


Figure 2.18: Left: CMS muon system schematic view; Right: one of the yoke rings with the muon DTs and RPCs installed; in the back it is possible to see the muon endcap [85].

1362 2.3.7 Muon system

1363 Muons are the only charged particles able to pass through all the CMS detector due
 1364 to their low ionization energy loss; thus, muons can be separated easily from the
 1365 high amount of particles produced in a pp collision. Also, muons are expected to be
 1366 produced in the decay of several new particles; therefore, good detection of muons
 1367 was one of the leading principles when designing the CMS detector.

1368 The CMS muon detection system (muon spectrometer) is embedded in the return
 1369 yoke as seen in Figure 2.18. It is composed of three different detector types, the drift
 1370 tube chambers (DT), cathode strip chambers (CSC), and resistive plate chambers
 1371 (RPC); DT are located in the central region $\eta < 1.2$ arranged in four layers of drift
 1372 chambers filled with an Ar/CO₂ gas mixture.

1373 The muon endcaps are made of CSCs covering the region $\eta < 2.4$ and filled with
 1374 a mixture of Ar/CO₂/CF₄. The reason behind using a different detector type lies on
 1375 the different conditions in the forward region like the higher muon rate and higher
 1376 residual magnetic field compared to the central region.

1377 The third type of detector used in the muon system is a set of four disks of RPCs

1378 working in avalanche mode. The RPCs provide good spatial and time resolutions. The
 1379 track of high- p_T muon candidates is built combining information from the tracking
 1380 system and the signal from up to six RPCs and four DT chambers.

1381 The muon tracks are reconstructed from the hits in the several layers of the muon
 1382 system.

1383 2.3.8 CMS trigger system

1384 CMS expects pp collisions every 25 ns, i.e., an interaction rate of 40 MHz for which
 1385 it is not possible to store the recorded data in full. In order to handle this high event
 1386 rate data, an online event selection, known as triggering, is performed; triggering
 1387 reduces the event rate to 100 Hz for storage and further offline analysis.

1388 The trigger system starts with a reduction of the event rate to 100 kHz in the
 1389 so-called *the level 1 trigger* (L1). L1 is based on dedicated programmable hardware
 1390 like Field Programmable Gate Arrays (FPGAs) and Application Specific Integrated
 1391 Circuits (ASICs), partly located in the detector itself; another portion is located in
 1392 the CMS underground cavern. Hit pattern information from the muon chambers
 1393 and the energy deposits in the calorimeter are used to decide if an event is accepted
 1394 or rejected, according to selection requirements previously defined, which reflect the
 1395 interesting physics processes. Figure 2.19 shows the L1 trigger architecture.

1396 The second stage in the trigger system is called *the high-level trigger* (HLT); events
 1397 accepted by L1 are passed to HLT in order to make an initial reconstruction of them.
 1398 HLT is software based and runs on a dedicated server farm, using selection algorithms
 1399 and high-level object definitions; the event rate at HLT is reduced to 100 Hz. The
 1400 first HLT stage takes information from the muon detectors and the calorimeters to
 1401 make the initial object reconstruction; in the next HLT stage, information from the

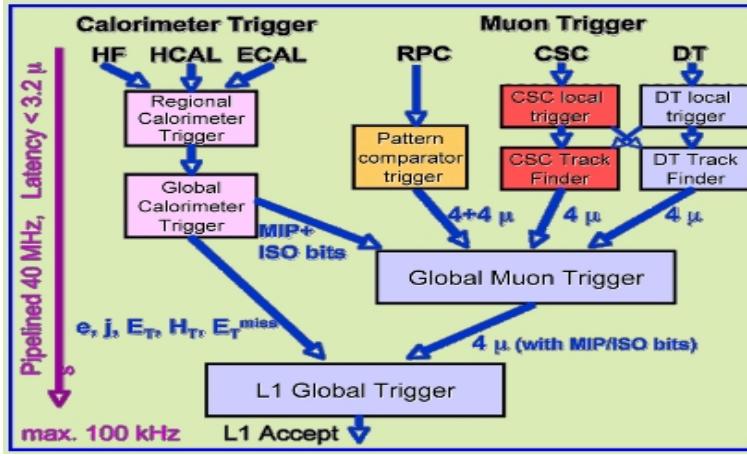


Figure 2.19: CMS Level-1 trigger architecture [86].

pixel and strip detectors is used to do first fast tracking and then full tracking online.
 This initial object reconstruction is used in further steps of the trigger system.

Events and preliminary reconstructed physics objects from HLT are sent to be fully reconstructed at the CERN computing center known also as Tier-0 facility. Again, the pixel detector information provides high-quality seeds for the track reconstruction algorithm offline, primary vertex reconstruction, electron and photon identification, muon reconstruction, τ identification, and b-tagging. After full reconstruction, data sets are made available for offline analyses.

2.3.9 CMS computing

Data coming from the experiment have to be stored and made available for further analysis; in order to cope with all the tasks implied in the offline data processing, like transfer, simulation, reconstruction and reprocessing, among others, a large computing power is required. The CMS computing system is based on the distributed architecture concept, where users of the system and physical computer centers are distributed worldwide and interconnected by high-speed networks.

The worldwide LHC computing grid (WLCG) is the mechanism used to provide

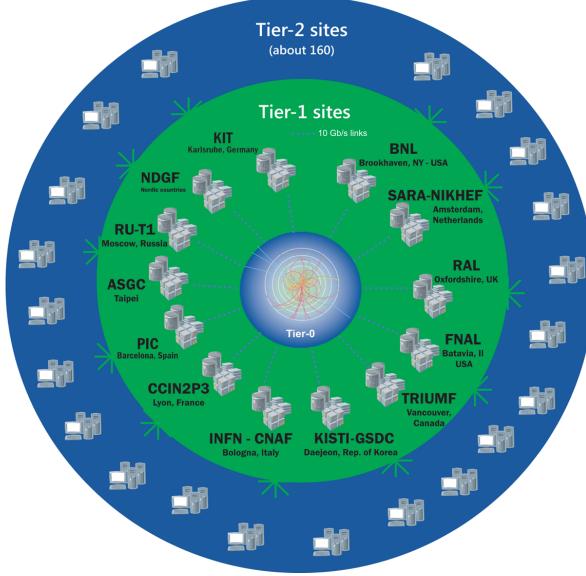


Figure 2.20: WLCG structure. The primary computer centers (Tier-0) are located at CERN (data center) and at the Wigner datacenter in Budapest. Tier-1 is composed of 13 centers and Tier-2 is composed of about 160 centers. [87].

1418 that distributed environment. WLCG is a tiered structure connecting computing
 1419 centers around the world, which provides the necessary storage and computing facil-
 1420 ties. The primary computing centers of the WLCG are located at the CERN and
 1421 the Wigner datacenter in Budapest and are known as Tier-0 as shown in Figure 2.20.
 1422 The main responsibilities for each tier level are [87]

- 1423 • **Tier-0:** initial reconstruction of recorded events and storage of the resulting
 1424 datasets, the distribution of raw data to the Tier-1 centers.
- 1425 • **Tier-1:** provide storage capacity, support for the Grid, safe-keeping of a pro-
 1426 portional share of raw and reconstructed data, large-scale reprocessing and safe-
 1427 keeping of corresponding output, generation of simulated events, distribution
 1428 of data to Tier 2s, safe-keeping of a share of simulated data produced at these
 1429 Tier 2s.
- 1430 • **Tier-2:** store sufficient data and provide adequate computing power for spe-

1431 cific analysis tasks and proportional share of simulated event production and
1432 reconstruction.

1433 Aside from the general computing strategy to manage the huge amount of data
1434 produced by experiments, CMS uses a software framework to perform a variety of
1435 processing, selection and analysis tasks. The central concept of the CMS data model
1436 referred to as *event data model* (EDM) is the *Event*; therefore, an event is the unit
1437 that contains the information from a single bunch crossing, any data derived from
1438 that information like the reconstructed objects, and the details of the derivation.

1439 Events are passed as the input to the *physics modules* that obtain information
1440 from them and create new information; for instance, *event data producers* add new
1441 data into the events, *analyzers* produce an information summary from an event set,
1442 *filters* perform selection and triggering.

1443 CMS uses several event formats with different levels of detail and precision

1444 • **Raw format:** events in this format contain the full recorded information from
1445 the detector as well as trigger decision and other metadata. An extended version
1446 of raw data is used to store information from the CMS Monte Carlo simulation
1447 tools (see Chapter 3). Raw data are stored permanently, occupying about
1448 2MB/event

1449 • **RECO format:** events in this format correspond to raw data that have been
1450 submitted to reconstruction algorithms like primary and secondary vertex re-
1451 construction, particle ID, and track finding. RECO events contain physics ob-
1452 jects and all the information used to reconstruct them; average size is about 0.5
1453 MB/event.

- 1454 • **AOD format:** Analysis Object Data (AOD) is the data format used in the
 1455 physics analyses given that it contains the parameters describing the high-level
 1456 physics objects in addition to enough information to allow a kinematic refitting if
 1457 needed. AOD events are filtered versions of the RECO events to which skimming
 1458 or other filtering have been applied, hence AOD events are subsets of RECO
 1459 events. Requires about 100 kB/event.
- 1460 • **Non-event data** are data needed to interpret and reconstruct events. Some
 1461 of the non-event data used by CMS contains information about the detector
 1462 contraction and condition data like calibrations, alignment, and detector status.

1463 Figure 2.21 shows the data flow scheme between CMS detector and tiers.

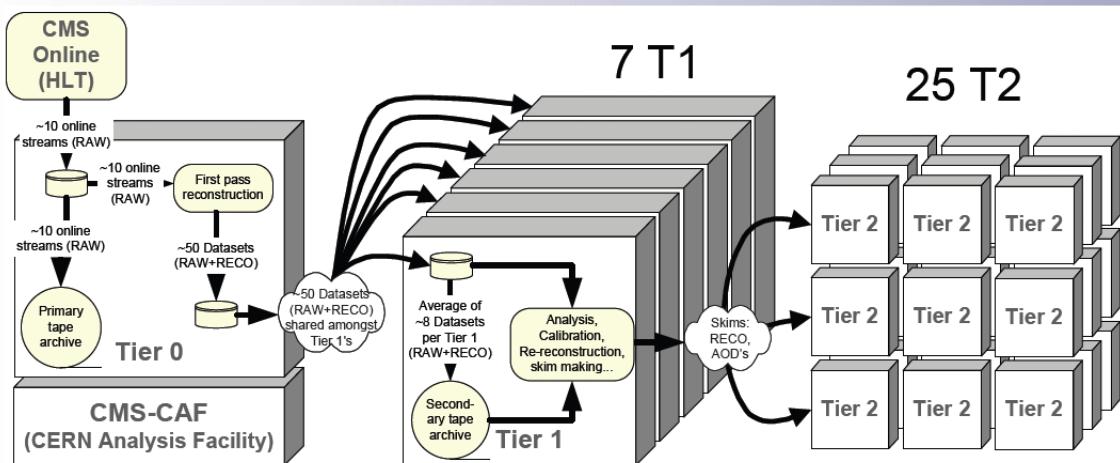


Figure 2.21: Data flow from CMS detector through tiers.

1464 The whole collection of software built as a framework is referred to as *CMSSW*. This
 1465 framework provides the services needed by the simulation, calibration and alignment,
 1466 and reconstruction modules that process event data, so that physicists can perform
 1467 analysis. The CMSSW event processing model is composed of one executable, called
 1468 `cmsRun`, and several plug-in modules which contains all the tools (calibration, recon-

1469 struction algorithms) needed to process an event. The same executable is used for
1470 both detector data and Monte Carlo simulations [88].

¹⁴⁷¹ **Chapter 3**

¹⁴⁷² **Event generation, simulation and
reconstruction**

¹⁴⁷⁴ The process of analyzing data recorded by the CMS experiment involves several stages
¹⁴⁷⁵ where the data are processed in order to interpret the information provided by all
¹⁴⁷⁶ the detection systems; in those stages, the particles produced after the pp collision
¹⁴⁷⁷ are identified by reconstructing their trajectories and measuring their features. In
¹⁴⁷⁸ addition, the SM provides a set of predictions that have to be compared with the
¹⁴⁷⁹ experimental results; however, in most of the cases, theoretical predictions are not
¹⁴⁸⁰ directly comparable to experimental results due to the diverse source of uncertainties
¹⁴⁸¹ introduced by the experimental setup and theoretical approximations, among others.

¹⁴⁸² The strategy to face these conditions consists in using statistical methods imple-
¹⁴⁸³ mented in computational algorithms to produce numerical results that can be con-
¹⁴⁸⁴ trasted with the experimental results. These computational algorithms are commonly
¹⁴⁸⁵ known as Monte Carlo (MC) methods and, in the case of particle physics, they are
¹⁴⁸⁶ designed to apply the SM rules and produce predictions about the physical observ-
¹⁴⁸⁷ ables measured in the experiments. Since particle physics is governed by quantum

mechanics principles, predictions are not allowed from single events; therefore, a high number of events are *generated* and predictions are produced in the form of statistical distributions for the observables. Effects of the detector presence are included in the predictions by introducing simulations of the detector itself.

This chapter presents a description of the event generation strategy and the tools used to perform the detector simulation and physics objects reconstruction. A comprehensive review of event generators for LHC physics can be found in Reference [89] on which this chapter is based.

3.1 Event generation

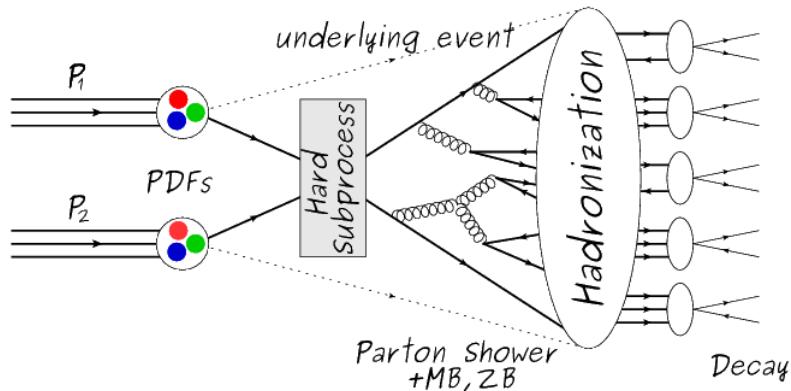


Figure 3.1: Event generation process. The actual interaction is generated in the hard subprocess. The parton shower describes the evolution of the partons from the hard subprocess. Modified from Reference [90].

The event generation is intended to create events that mimic the behavior of actual events produced in collisions; they obey a sequence of steps from the particles collision hard process to the decay process into the final state. Figure 3.1 shows a schematic view of the event generation process; the fact that the full process can be treated as several independent steps is motivated by the QCD factorization theorem.

1502 Generation starts by taking into account the PDFs of the incoming particles.
 1503 Event generators offer the option to chose from several PDF sets depending on the
 1504 particular process under simulation¹; in the following, pp collisions will be consid-
 1505 ered. The *hard subprocess* describes the actual interaction between partons from the
 1506 incoming protons; it is represented by the matrix element connecting the initial and
 1507 final states of the interaction. Normally, the matrix element can be written as a
 1508 sum over Feynman diagrams and consider interferences between terms in the sum-
 1509 mation. During the generation of the hard subprocess, the production cross section
 1510 is calculated.

1511 The order to which the cross section is calculated depends on the order of the Feyn-
 1512 man diagrams involved in the calculation; therefore, radiative corrections are included
 1513 by considering a higher order Feynman diagrams where QCD radiation dominates.
 1514 Currently, cross sections calculated to LO do not offer a satisfactory description of the
 1515 processes, i.e., the results are only reliable for the shape of distributions; therefore,
 1516 NLO calculations have to be performed with the implication that the computing time
 1517 needed is highly increased.

1518 The final parton content of the hard subprocess is subjected to the *parton shower*
 1519 which generates the gluon radiation. Parton shower evolves the partons, i.e., glouns
 1520 split into quark-antiquark pairs and quarks with enough energy radiate gluons giv-
 1521 ing rise to further parton multiplication, following the DGLAP (Dokshitzer-Gribov-
 1522 Lipatov-Altarelli-Parisi) equations. Showering continues until the energy scale is low
 1523 enough to reach the non-perturbative limit.

1524 In the simulation of LHC processes that involve b quarks, like the single top quark
 1525 or Higgs associated production, it is needed to consider that the b quark is heavier
 1526 than the proton; hence, the QCD interaction description is made in two different

¹ Tool in Reference [91] allows to plot different PDF sets under customizable conditions.

1527 schemes [95]

- 1528 • four-flavor (4F) scheme. b quarks appear only in the final state because they
 1529 are heavier than the proton and therefore they can be produced only from the
 1530 splitting of a gluon into pairs or singly in association with a t quark in high
 1531 energy-scale interactions; furthermore, during the simulation, the b -PDFs are set
 1532 to zero. Calculations in this scheme are more complicated due to the presence
 1533 of the second b quark but the full kinematics is considered already at LO and
 1534 therefore the accuracy of the description is better.

- 1535 • five-flavor (5F) scheme. b quarks are considered massless, therefore they can
 1536 appear in both initial and final states since they can now be part of the proton;
 1537 thus, during the simulation b -PDFs are not set to zero. In this scheme, calcula-
 1538 tions are simpler than in the 4F scheme and possible logarithmic divergences
 1539 are absorbed by the PDFs through the DGLAP evolution.

1540 In this thesis, the tHq events are generated using the 4F scheme in order to reduce
 1541 uncertainties, while the tHW events are generated using the 5F scheme to eliminate
 1542 LO interference with $t\bar{t}H$ process [48].

1543 Partons involved in the pp collision are the focus of the simulation, however, the
 1544 rest of the partons inside the incoming protons are also affected because the remnants
 1545 are colored objects; also, multiple parton interactions can occur. The hadronization
 1546 of the remnants and multiple parton interactions are known as *underlying event* and
 1547 it has to be included in the simulation. In addition, multiple pp collisions in the same
 1548 bunch crossing (pile-up mentioned in 2.2) occurs, actually in two forms

- 1549 • *in-time PU* which refers to multiple pp collision in the bunch crossing but that
 1550 are not considered as primary vertices.

1551 • *Out-of-time PU* which refers to overlapping pp collisions from consecutive bunch
 1552 crossings; this can occur due to the time-delays in the detection systems where
 1553 information from one bunch crossing is assigned to the next or previous one.

1554 While the underlying event effects are included in generation using generator-
 1555 specific tools, PU effects are added to the generation by overlaying Minimum-bias (MB)
 1556 and Zero-bias (ZB) events to the generated events. MB events are inelastic events
 1557 selected by using a loose trigger with as little bias as possible, therefore accepting a
 1558 large fraction of the overall inelastic event; ZB events correspond to random events
 1559 recorded by the detector when collisions are likely. MB models in-time PU and ZB
 1560 models out-of-time PU.

1561 The next step in the generation process is called *hadronization*. Since particles
 1562 with a net color charge are not allowed to exits isolated, they have to recombine
 1563 to form bound states. This is precisely the process by which the partons resulting
 1564 from the parton shower arrange themselves as color singlets to form hadrons. At
 1565 this step, the energy-scale is low and the strong coupling constant is large, therefore
 1566 hadronization process is non-perturbative and the evolution of the partons is described
 1567 using phenomenological models. Most of the baryons and mesons produced in the
 1568 hadronization are unstable and hence they will decay in the detector.

1569 The last step in the generation process corresponds to the decay of the unstable
 1570 particles generated during hadronization; it is also simulated in the hadronization
 1571 step, based on the known branching ratios.

1572 **3.2 Monte Carlo Event Generators.**

1573 The event generation described in the previous section has been implemented in
 1574 several software packages for which a brief description is given.

- 1575 • **PYTHIA 8.** It is a program designed to perform the generation of high energy
 1576 physics events which describes the collisions between particles such as electrons
 1577 and protons. Several theories and models are implemented in it, in order to
 1578 describe physical aspects like hard and soft interaction, parton distributions,
 1579 initial and final-state parton showers, multiple parton interactions, beam rem-
 1580 nants, hadronization² and particle decay. Thanks to extensive testing, several
 1581 optimized parametrizations, known as *tunings*, have been defined in order to
 1582 improve the description of actual collisions to a high degree of precision; for
 1583 analysis at $\sqrt{s} = 13$ TeV, the underline event CUETP8M1 tune is employed [97].
 1584 The calculation of the matrix element is performed at LO which is not enough
 1585 for the current required level of precision; therefore, pythia is often used for
 1586 parton shower, hadronization and decays, while other event generators are used
 1587 to generate the matrix element at NLO.
- 1588 • **MadGraph5_aMC@NLO.** MadGraph is a matrix element generator which
 1589 calculates the amplitudes for all contributing Feynman diagrams of a given
 1590 process but does not provide a parton shower while MC@NLO incorporates
 1591 NLO QCD matrix elements consistently into a parton shower framework; thus,
 1592 MadGraph5_aMC@NLO, as a merger of the two event generators MadGraph5
 1593 and aMC@NLO, is an event generator capable to calculate tree-level and NLO
 1594 cross sections and perform the matching of those with the parton shower. It is
 1595 one of the most frequently used matrix element generators; however, it has the
 1596 particular feature of the presence of negative event weights which reduce the
 1597 number of events used to reproduce the properties of the objects generated [98].
- 1598 • **POWHEG.** It is an NLO matrix element generator where the hardest emis-

² based in the Lund string model [96]

1599 sion of color charged particles is generated in such a way that the negative event
 1600 weights issue of MadGraph5_aMC@NLO is overcome; however, the method re-
 1601 quires an interface with p_T -ordered parton shower or a parton shower generator
 1602 where this highest emission can be vetoed in order to avoid double counting of
 1603 this highest-energetic emission. PYTHIA is a commonly matched to POWHEG
 1604 event generator [100].

1605 Events resulting from the whole generation process are known as MC events.

1606 **3.3 CMS detector simulation.**

1607 After generation, MC events contain the physics of the collisions but they are not
 1608 ready to be compared to the events recorded by the experiment since these recorded
 1609 events correspond to the response of the detection systems to the interaction with
 1610 the particles traversing them. The simulation of the CMS detector has to be applied
 1611 on top of the event generation; it is simulated with a MC toolkit for the simulation
 1612 of particles passing through matter called Geant4 which is also able to simulate the
 1613 electronic signals that would be measured by all detectors inside CMS.

1614 The simulation takes the generated particles contained in the MC events as input,
 1615 makes them pass through the simulated geometry, and models physics processes that
 1616 particles experience during their passage through matter. The full set of results from
 1617 particle-matter interactions corresponds to the simulated hit which contains informa-
 1618 tion about the energy loss, momentum and position. Particles of the input event are
 1619 called *primary*, while the particles originating from GEANT4-modeled interactions of
 1620 a primary particle with matter are called a *secondary*. Simulated hits are the input
 1621 of subsequent modules that emulate the response of the detector readout system and

1622 triggers. The output from the emulated detection systems and triggers is known as
 1623 digitization [101, 102].

1624 The modeling of the CMS detector corresponds to the accurate modeling of the
 1625 interaction among particles, the detector material, and the magnetic field. This
 1626 simulation procedure includes the following standard steps

1627 • Modeling of the Interaction Region.

1628 • Modeling of the particle passage through the hierarchy of volumes that compose
 1629 CMS detector and of the accompanying physics processes.

1630 • Modeling of the effect of multiple interactions per beam crossing and/or the
 1631 effect of events overlay (Pile-Up simulation).

1632 • Modeling of the detector's electronics response, signal shape, noise, calibration
 1633 constants (digitization).

1634 In addition to the full simulation, i.e., a detailed detector simulation, a faster
 1635 simulation (FastSim) have been developed, that may be used where much larger
 1636 statistics are required. In FastSim, detector material effects are parametrized and
 1637 included in the hits; those hits are used as input of the same higher-level algorithms³
 1638 used to analyze the recorded events. In this way, comparisons between fast and full
 1639 simulations can be performed [104].

1640 After the full detector simulation, the output events can be directly compared
 1641 to events actually recorded in the CMS detector. The collection of MC events that
 1642 reproduces the expected physics for a given process is known as MC sample.

³ track fitting, calorimeter clustering, b tagging, electron identification, jet reconstruction and calibration, trigger algorithms which will be considered in the next sections

1643 **3.4 Event reconstruction.**

1644 The CMS experiment use the *particle-flow event reconstruction algorithm (PF)* to do
1645 the reconstruction of particles produced in pp collisions. Next sections will present
1646 a basic description of the *Elements* used by PF (tracker tracks, energy clusters, and
1647 muon tracks), based in the References [105, 106] where more detailed descriptions can
1648 be found.

1649 **3.4.1 Particle-Flow Algorithm.**

1650 Each of the several sub detection systems of the CMS detector is dedicated to identify
1651 an specific type of particles, i.e., photons and electrons are absorbed by the ECAL
1652 and their reconstruction is based on ECAL information; hadrons are reconstructed
1653 from clusters in the HCAL while muons are reconstructed from hits in the muon
1654 chambers. PF is designed to correlate signals from all the detector layers (tracks and
1655 energy clusters) in order to reconstruct and identify each final state particle and its
1656 properties as sketched in Figure 3.2.

1657 For instance, a charged hadron is identified by a geometrical connection, known
1658 as *link*, between one or more calorimeter clusters and a track in the tracker, provided
1659 there are no hits in the muon system; combining several measurements allows a better
1660 determination of the energy and charge sign of the charged hadron.

1661 **Charged-particle track reconstruction.**

1662 The strategy used by PF in order to reconstruct tracks is called *Iterative Tracking*
1663 which occurs in four steps

- 1664 • Seed generation where initial track candidates are found by looking for a combi-
- 1665 nation of hits in the pixel detector, strip tracker, and muon chambers. In total

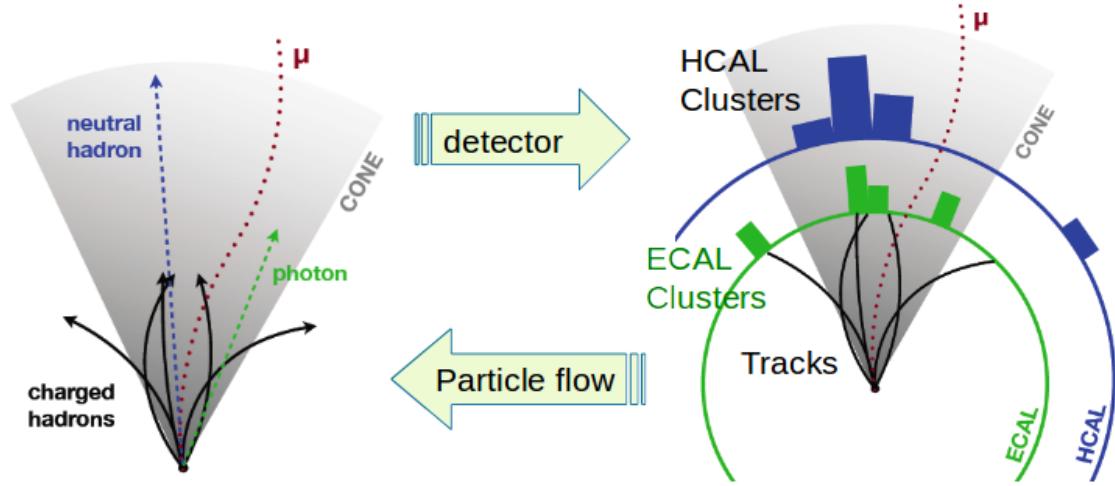


Figure 3.2: Particle flow algorithm. Information from the several CMS detection systems is provided as input to the algorithm which then combine it to identify and reconstruct all the particles in the final state and their properties. Reconstruction of simulated events is also performed by providing information from MC samples, detector and trigger simulation [107].

1666 ten iterations are performed, each one with a different seeding requirement.
 1667 Seeds are used to estimate the trajectory parameters and uncertainties at the
 1668 time of the full track reconstruction. Seeds are also considered track candidates.

- 1669 • Track finding using a tracking software known as Combinatorial Track Finder
 1670 (CTF) [108]. The seed trajectories are extrapolated along the expected flight
 1671 path of a charged particle, in agreement to the trajectory parameters obtained
 1672 in the first step, in an attempt to find additional hits that can be assigned to
 1673 the track candidates.

- 1674 • Track-fitting where the found tracks are passed as input to a module which
 1675 provides the best estimate of the parameters of each trajectory.

- 1676 • Track selection where track candidates are submitted to a selection which dis-
 1677 cards those that fail a set of defined quality criteria.

1678 Iterations differ in the seeding configuration and the final track selection as elab-

1679 orated in References [105, 106]. In the first iteration, high p_T tracks and tracks pro-
 1680 duced near to the interaction region are identified and those hits are masked thereby
 1681 reducing the combinatorial complexity. Next, iterations search for more complicated
 1682 tracks, like low p_T tracks and tracks from b hadron decays, which tend to be displaced
 1683 from the interaction region.

1684 **Vertex reconstruction.**

1685 During the track reconstruction, an extrapolation toward to the calorimeters is per-
 1686 formed in order to match energy deposits; that extrapolation is performed also toward
 1687 the beamline in order to find the origin of the track known as *vertex*. The vertex re-
 1688 construction is performed by selecting from the available reconstructed tracks, those
 1689 that are consistent with being originated in the interaction region where pp collisions
 1690 are produced. The selection involves a requirement on the number of tracker (pixel
 1691 and strip) hits and the goodness of the track fit.

1692 Selected tracks are clustered using a *deterministic annealing algorithm* (DA)⁴. A
 1693 set of candidate vertices and their associated tracks, resulting from the DA, are then
 1694 fitted with an *adaptive vertex fitter* (AVF) to produce the best estimate of the vertices
 1695 locations.

1696 The p_T of the tracks associated to a reconstructed vertex is added, squared and
 1697 used to organize the vertices; the vertex with the highest squared sum is designated
 1698 as the *primary vertex* (*PV*) while the rest are designated as PU vertices.

1699 **Calorimeter clustering.**

1700 After traversing the CMS tracker system, electrons, photons and hadrons deposit their
 1701 energy in the ECAL and HCAL cells. The PF clustering algorithm aims to provide

⁴ DA algorithm and AVF are described in detail in References [110, 111]

1702 a high detection efficiency even for low-energy particles and an efficient distinction
 1703 between close energy deposits. The clustering runs independently in the ECAL barrel
 1704 and endcaps, HCAL barrel and endcaps, and the two preshower layers, following two
 1705 steps

- 1706 • cells with an energy larger than a given seed threshold and larger than the energy
 1707 of the neighboring cells are identified as cluster seeds. The neighbor cells are
 1708 those that either share a side with the cluster seed candidate, or the eight closest
 1709 cells including cells that only share a corner with the seed candidate.
- 1710 • cells with at least a corner in common with a cell already in the cluster seed
 1711 and with an energy above a cell threshold are grouped into topological clusters.

1712 Clusters formed in this way are known as *particle-flow clusters*. With this cluster-
 1713 ing strategy, it is possible to detect and measure the energy and direction of photons
 1714 and neutral hadrons as well as differentiate these neutral particles from the charged
 1715 hadron energy deposits. In cases involving charged hadrons for which the track pa-
 1716 rameters are not determined accurately, for instance, low-quality and high- p_T tracks,
 1717 clustering helps in the energy measurements.

1718 Electron track reconstruction.

1719 Although the charged-particle track reconstruction described above works for elec-
 1720 trons, they lose a significant fraction of their energy via bremsstrahlung photon radi-
 1721 ation before reaching the ECAL; thus, the reconstruction performance depends on the
 1722 ability to measure also the radiated energy. The reconstruction strategy, in this case,
 1723 requires information from the tracking system and from the ECAL. Bremsstrahlung
 1724 photons are emitted at similar η values to that of the electron but at different values
 1725 of ϕ ; therefore, the radiated energy can be recovered by grouping ECAL clusters in a

1726 η window over a range of ϕ around the electron direction. The group is called ECAL
 1727 supercluster.

1728 Electron candidates from the track-seeding and ECAL super clustering are merged
 1729 into a single collection which is submitted to a full electron tracking fit with a
 1730 Gaussian-sum filter (GSF) [109]. The electron track and its associated ECAL su-
 1731 percluster form a *particle-flow electron*.

1732 **Muon track reconstruction.**

1733 Given that the CMS detector is equipped with a muon spectrometer capable to iden-
 1734 tify and measure the momentum of the muons traversing it, the muon reconstruction
 1735 is not specific to PF; therefore, three different muon types are defined

- 1736 • *Standalone muon.* A clustering on the DTs or CSCs hits is performed to form
 1737 track segments; those segments are used as seeds for the reconstruction in the
 1738 muon spectrometer. All DTs, CSCs, and RPCs hits along the muon trajectory
 1739 are combined and fitted to form the full track. The fitting output is called a
 1740 *standalone-muon track*.
- 1741 • *Tracker muon.* Each track in the inner tracker with p_T larger than 0.5 GeV and
 1742 a total momentum p larger than 2.5 GeV is extrapolated to the muon system.
 1743 A *tracker muon track* corresponds to a extrapolated track that matches at least
 1744 one muon segment.
- 1745 • *Global muon.* When tracks in the inner tracker (inner tracks) and standalone-
 1746 muon tracks are matched and turn out being compatibles, their hits are com-
 1747 bined and fitted to form a *global-muon track*.

1748 Global muons sharing the same inner track with tracker muons are merged into
 1749 a single candidate. PF muon identification uses the muon energy deposits in ECAL,
 1750 HCAL, and HO associated with the muon track to improve the muon identification.

1751 **Particle identification and reconstruction.**

1752 PF elements are connected by a linker algorithm that tests the connection between any
 1753 pair of elements; if they are found to be linked, a geometrical distance that quantifies
 1754 the quality of the link is assigned. Two elements may be linked indirectly through
 1755 common elements. Linked elements form *PF blocks* and each PF block may contain
 1756 elements originating in one or more particles. Links can be established between
 1757 tracks, between calorimeter clusters, and between tracks and calorimeter clusters.
 1758 The identification and reconstruction start with a PF block and proceed as follows

- 1759 • Muons. An *isolated global muon* is identified by evaluating the presence of
 inner track and energy deposits close to the global muon track in the (η, ϕ)
 plane, i.e., in a particular point of the global muon track, inner tracks and
 energy deposits are sought within a radius of $\Delta R = 0.3$ (see eqn. 2.7) from the
 muon track; if they exist and the p_T of the found track added to the E_T of the
 found energy deposit does not exceed 10% of the muon p_T then the global muon
 is an isolated global muon. This isolation condition is stringent enough to reject
 hadrons misidentified as muons.

1767 *Non-isolated global muons* are identified using additional selection requirements
 1768 on the number of track segments in the muon system and energy deposits along
 1769 the muon track. Muons inside jets are identified with more stringent criteria
 1770 in isolation and momentum as described in Reference [112]. The PF elements
 1771 associated with an identified muon are masked from the PF block.

- 1772 ● Electrons are identified and reconstructed as described above plus some addi-
 1773 tional requirements on fourteen variables like the amount of energy radiated,
 1774 the distance between the extrapolated track position at the ECAL and the po-
 1775 sition of the associated ECAL supercluster, among others, which are combined
 1776 in an specialized multivariate analysis strategy that improves the electron iden-
 1777 tification. Tracks and clusters used to identify and reconstruct electrons are
 1778 masked in the PF block.
- 1779 ● Isolated photons are identified from ECAL superclusters with E_T larger than 10
 1780 GeV, for which the energy deposited at a distance of 0.15, from the supercluster
 1781 position on the (η, ϕ) plane, does not exceed 10% of the supercluster energy;
 1782 note that this is an isolation requirement. In addition, there must not be links
 1783 to tracks. Clusters involved in the identification and reconstruction are masked
 1784 in the PF block.
- 1785 ● Bremsstrahlung photons and prompt photons tend to convert to electron-positron
 1786 pairs inside the tracker, therefore, a dedicated finder algorithm is used to link
 1787 tracks that seem to originate from a photon conversion; in case those two tracks
 1788 are compatible with the direction of a bremsstrahlung photon, they are also
 1789 linked to the original electron track. Photon conversion tracks are also masked
 1790 in the PF block.
- 1791 ● The remaining elements in the PF block are used to identify hadrons. In the
 1792 region $|\eta| \leq 2.5$, neutral hadrons are identified with HCAL clusters not linked
 1793 to any track while photons from neutral pion decays are identified with ECAL
 1794 clusters without links to tracks. In the region $|\eta| > 2.5$ ECAL clusters linked to
 1795 HCAL clusters are identified with a charged or neutral hadron shower; ECAL
 1796 clusters with no links are identified with photons. HCAL clusters not used yet,

1797 are linked to one or more unlinked tracks and to an unlinked ECAL in order to
 1798 reconstruct charged-hadrons or a combination of photons and neutral hadrons
 1799 according to certain conditions on the calibrated calorimetric energy.

- 1800 • Charged-particle tracks may be liked together when they converge to a *sec-*
 1801 *ondary vertex (SV)* displaced from the IP where the PV and PU vertices are
 1802 reconstructed; at least three tracks are needed in that case, of which at most
 1803 one has to be an incoming track with hits in tracker region between a PV and
 1804 the SV.

1805 The linker algorithm, as well as the whole PF algorithm, has been validated and
 1806 commissioned; results from that validation are presented in the Reference [105].

1807 **Jet reconstruction.**

1808 Quarks and gluons may be produced in the pp collisions, therefore, their hadronization
 1809 will be seen in the detector as a shower of hadrons and their decay products in the
 1810 form of a *jet*. Two classes of clustering algorithms have been developed based in
 1811 their jet definition [113]:

- 1812 • Iterative cone algorithms (IC). Jets are defined in terms of circles of fixed radius
 1813 R in the $\eta\text{-}\phi$ plane, known as *stable cones*, for which the sum of the momenta
 1814 of all the particles within the cone points in the same direction as the center
 1815 of the circle. The seed of the iteration is the hardest non-isolated particle in
 1816 the event, then, the resulting momentum direction is assigned as the new cone
 1817 direction and a new iteration starts; iteration process stops when the cone if
 1818 found to be stable.

1819 • Sequential recombination algorithms. The distance between non-isolated par-
 1820 ticles is calculated; if that distance is below a threshold, these particles are
 1821 recombined into a new object. The sequence is repeated until the separation
 1822 between the recombined object and any other particle is above certain thresh-
 1823 old; the recombined object is called a jet and the algorithm starts again with
 1824 the remaining particles.

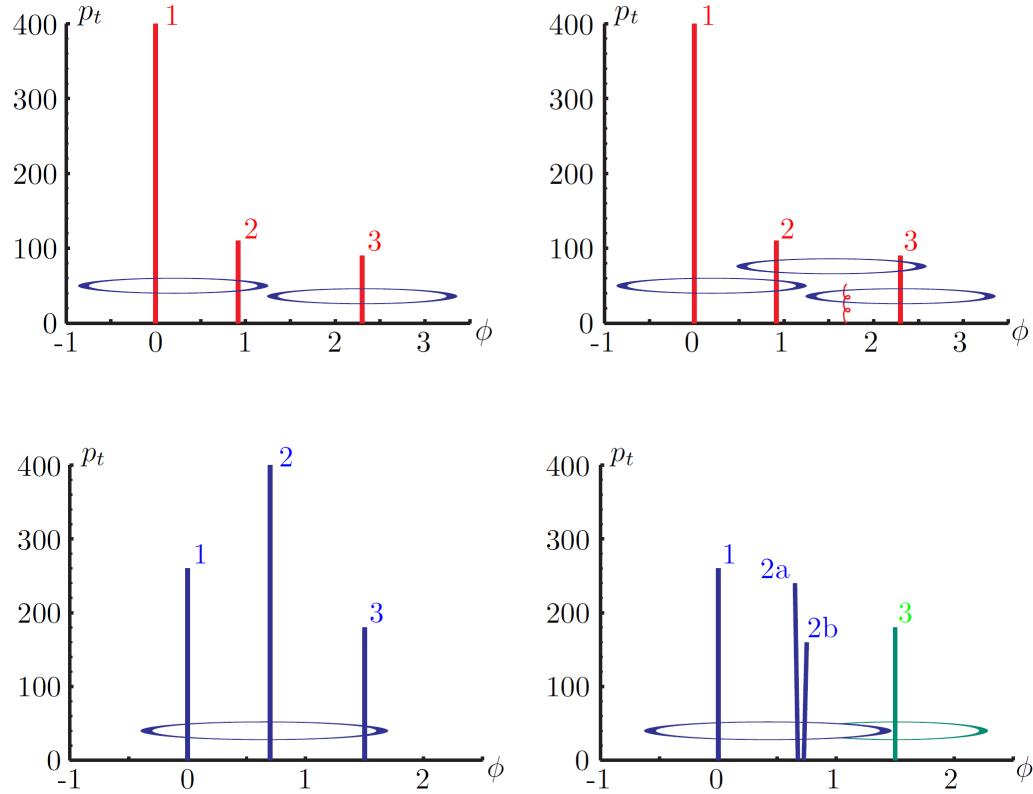


Figure 3.3: Stable cones identification using IC algorithms [113].

1825 Two conditions are of particular importance for the clustering algorithms, *infrared*
 1826 and *collinear (IRC) safety*. In order to explain the concept of infrared (IR) safety,
 1827 consider an event with three hard particles as shown in the top left side of Figure 3.3,
 1828 two stable cones are found and then two jets are identified; if a soft gluon is added, as
 1829 shown in the top right side of Figure 3.3, three stable cones are found and the three

1830 hard particles are now clustered into a single jet. If the addition of soft particles
 1831 change the outcome of the clustering, then it is said that the algorithm is IR unsafe.
 1832 Soft radiation is highly likely in perturbative QCD, which dominates the physics of
 1833 the jets, and then IR unsafe effect leads to divergences [113].

1834 The concept of collinear safety can also be explained considering a three hard
 1835 particles event, as shown in the bottom left side of Figure 3.3, where one stable cone
 1836 containing all three particles is found and one jet is identified; if the hardest particle
 1837 is split into two collinear particles (2a and 2b) in the bottom right side of Figure 3.3,
 1838 then the clustering results in a different jet identification and the algorithm is said
 1839 to be collinear unsafe. The collinear unsafe effect leads to divergences in jet cross
 1840 section calculations [114].

1841 It has been determined that IC algorithms are IRC unsafe, and therefore, they
 1842 have to be replaced by algorithms that not only provide the finite perturbative results
 1843 from theoretical computations, but also that are not highly dependent on underlying
 1844 event and pileup effects which leads to significant corrections [113].

1845 The sequential recombination algorithms arise as the IRC safe alternative used by
 1846 the CMS experiment; in particular the anti- k_t algorithm which is a generalization of
 1847 the previously existing k_t [115] and Cambridge/Aachen [116] jet clustering algorithms.

1848 The anti- k_t algorithm is used to perform the jet reconstruction by clustering those
 1849 PF particles within a cone (see Figure 3.4); previously, isolated electrons, isolated
 1850 muons, and charged particles associated with other interaction vertices are excluded
 1851 from the clustering.

1852 The anti- k_t algorithm proceeds in a sequential recombination of PF particles; the
 1853 distance between particles i and j (d_{ij}) and the distance between particles and the
 1854 beam are defined as

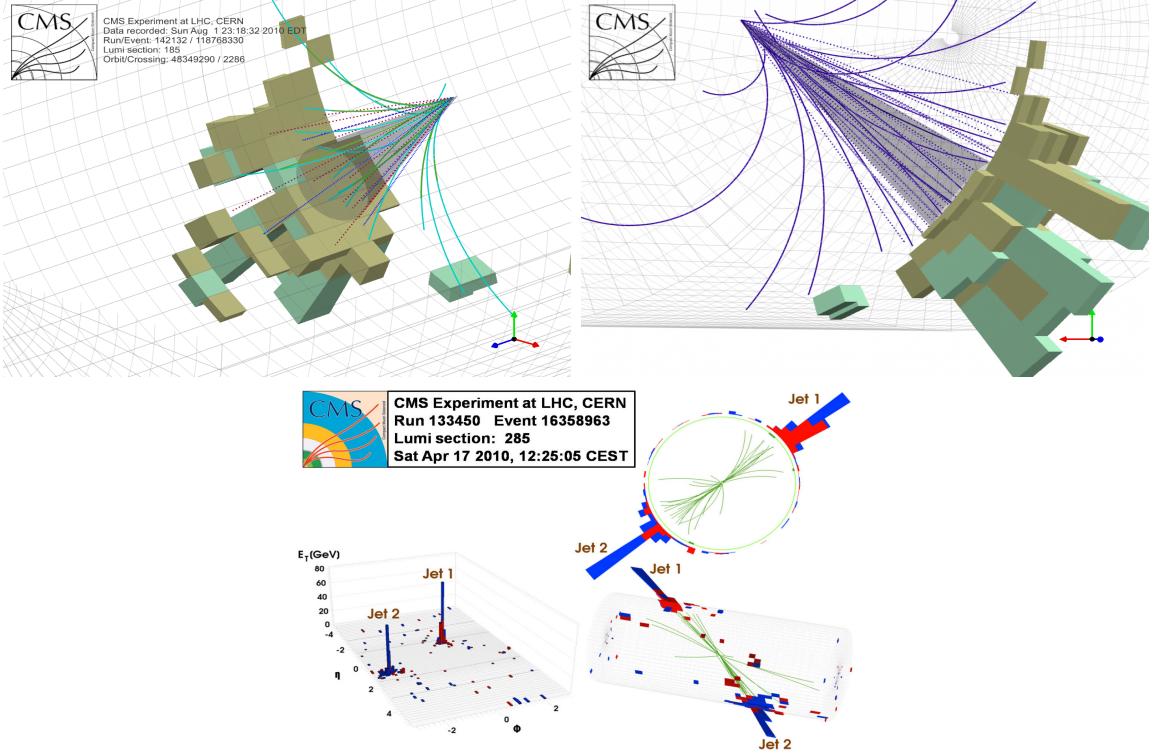


Figure 3.4: Jet reconstruction performed by the anti- k_t algorithm. Top: Two different views of a CMS recorded event are presented. Continuous lines correspond to tracks left by charged particles in the tracker while dotted lines are the imaginary paths followed by neutral particles. The green cubes represent the ECAL cells while the blue ones represent the HCAL cells; in both cases, the height of the cube represent the amount of energy deposited in the cells [117]. Bottom: Reconstruction of a recorded event with two jets [118].

$$d_{ij} = \min\left(\frac{1}{k_{ti}^2}, \frac{1}{k_{tj}^2}\right) \frac{\Delta_{ij}^2}{R^2}$$

$$d_{iB} = \frac{1}{k_{ti}^2} \quad (3.1)$$

1855 where $\Delta_{ij}^2 = (y_i - y_j)^2 + (\phi_i - \phi_j)^2$, k_{ti} , y_i and ϕ_i are the transverse momentum, 1856 rapidity and azimuth of particle i respectively and R is the called jet radius. For all 1857 the remaining PF particles, after removing the isolated ones, d_{ij} and d_{iB} are calcu-

lated⁵ and the smallest is identified; if it is a d_{ij} , particles i and j are replaced with
 a new object whose momentum is the vectorial sum of the combined particles. If the
 smallest distance is a d_{iB} the clustering process ends, the object i (which at this stage
 should be a combination of several PF particles) is declared as a *Particle-flow-jet* (PF
 jet) and all the associated PF particles are removed from the detector. The clustering
 process is repeated until no PF particles remain. R is a free parameter that can be
 adjusted according to the specific analysis conditions; usually, two values are used,
 $R=0.4$ and $R=0.5$, giving the name to the so-called AK4-jet and AK5-jet respectively.

An advantage of the anti- k_t algorithm over other clustering algorithms is the reg-
 ularity of the boundaries of the resulting jets. For all known IRC safe algorithms,
 soft radiation can introduce irregularities in the boundaries of the final jets; however,
 anti- k_t algorithm is soft-resilient, meaning that jets shape is not affected by soft radi-
 ation, which is a valuable property considering that knowing the typical shape of jets
 makes experimental calibration of jets more simple. In addition, that soft-resilience
 is expected to simplify certain theoretical calculations and reduce the momentum-
 resolution loss caused by underlying-event (UE) and pileup contamination [114].

The effect of the UE and pileup contamination over a jet identification, can be
 seen as if soft events are added to the jet; for instance, if a soft event representing UE
 or pileup is added to an event for which a set of jets J have been identified, and the
 clustering is rerun on that new extended event, the outcome will be different in two
 aspects: jets will contain some additional soft energy and the distribution of particles
 in jets may have change; that effect is called *back-reaction*. The back-reaction effect in
 the anti- k_t algorithm is suppressed not by the amount of momentum added to the jet
 but by the jet transverse momentum $p_{T,J}$, which means that this strong suppression
 leads to a smaller correction due to EU and pileup effect [114].

⁵ Notice that this is a combinatorial calculation.

1883 Jet energy Corrections

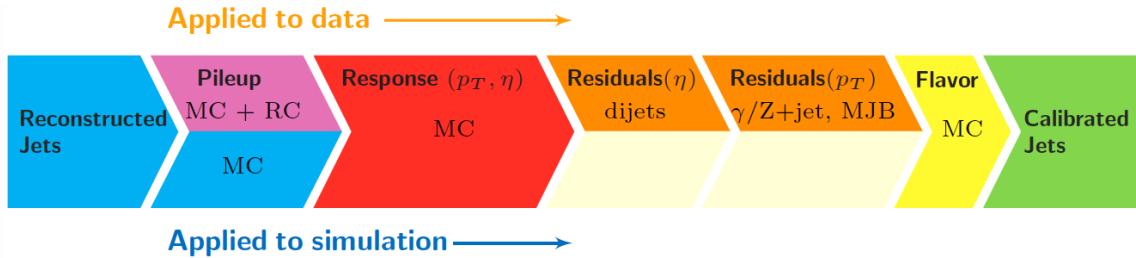


Figure 3.5: Jet energy correction diagram. Correction levels are applied sequentially in the indicated fixed order [120].

1884 Even though jets can be reconstructed efficiently, there are some effects that are
 1885 not included in the reconstruction and that lead to discrepancies between the re-
 1886 constructed results and the predicted results; in order to overcome these discrep-
 1887 ancies, a factorized model has been designed in the form of jet energy corrections
 1888 (JEC) [119, 120] applied sequentially as shown in the diagram of Figure 3.5.

1889 At each level, the jet four-momentum is multiplied by a scaling factor based on
 1890 jet properties, i.e., η , flavor, etc.

- 1891 • Level 1 correction removes the energy coming from pile-up. The scale factor is
 1892 determined using a MC sample of QCD dijet (2 jets) events with and without
 1893 pileup overlay; it is parametrized in terms of the offset energy density ρ , jet
 1894 area A , jet η and jet p_T . Different corrections are applied to data and MC due
 1895 to the detector simulation.
- 1896 • MC-truth correction accounts for differences between the reconstructed jet en-
 1897 ergy and the MC particle-level energy. The correction is determined on a QCD
 1898 dijet MC sample and is parametrized in terms of the jet p_T and η .
- 1899 • Residuals correct remaining small differences within jet response in data and
 1900 MC. The Residuals η -dependent correction compares jets of similar p_T in the

1901 barrel reference region. The Residuals p_T -dependent correct the jet absolute
 1902 scale (JES vs p_T).

- 1903 • Jet-flavor corrections are derived in the same way as MC-truth corrections but
 1904 using QCD pure flavor samples.

1905 ***b*-tagging of jets.**

1906 A particular feature of the hadrons containing bottom quarks (*b*-hadrons) is that
 1907 their lifetime is long enough to travel some distance before decaying, but it is not as
 1908 long as those of light quark hadrons; therefore, when looking at the hadrons produced
 1909 in pp collisions, *b*-hadrons decay typically inside the tracker rather than reaching the
 1910 calorimeters as some light-hadrons do. As a result, a *b*-hadron decay gives rise to a
 1911 displaced vertex (secondary vertex) with respect to the primary vertex as shown in
 1912 Figure 3.6; the SV displacement is in the order of a few millimeters. A jet resulting
 1913 from the decay of a *b*-hadron is called *b* jet; other jets are called light jets.

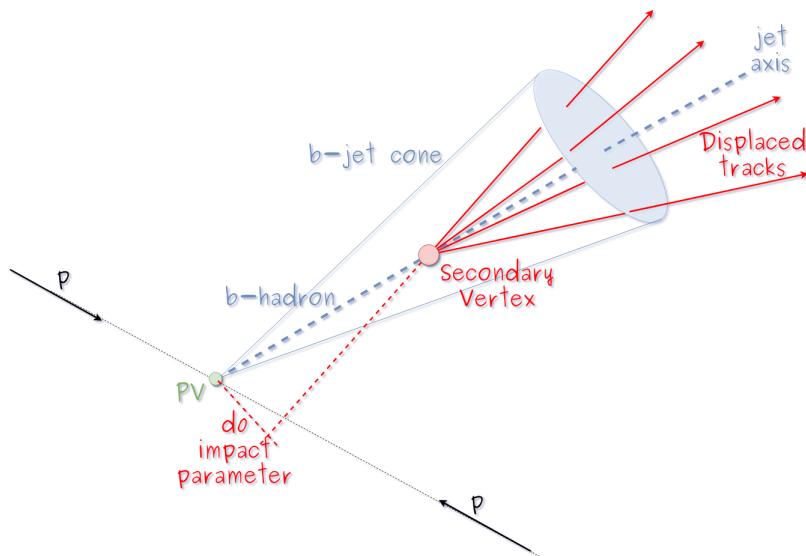


Figure 3.6: Secondary vertex in a *b*-hadron decay.

1914 Several methods to identify *b*-jets (*b*-tagging) have been developed; the method

1915 used in this thesis is known as *Combined Secondary Vertex* algorithm in its second
 1916 version (CSVv2) [121]. By using information of the impact parameter, the recon-
 1917 structed secondary vertices, and the jet kinematics as input in a multivariate analysis
 1918 that combines the discrimination power of each variable in one global discrimina-
 1919 tor variable, three working points (references): loose, medium and tight, are defined
 1920 which quantify the probabilities of mistag jets from light quarks as jets from b quarks;
 1921 10, 1 and 0.1 % respectively. Although the mistagging probability decreases with the
 1922 working point strength, the efficiency to correctly tag b -jets also decreases as 83, 69
 1923 and 49 % for the respective working point; therefore, a balance needs to be achieved
 1924 according to the specific requirements of the analysis.

1925 3.4.1.1 Missing transverse energy.

1926 The fact that proton bunches carry momentum along the z -axis implies that for
 1927 each event it is expected that the momentum in the transverse plane is balanced.
 1928 Imbalances are quantified by the missing transverse energy (MET) and are attributed
 1929 to several sources including particles escaping undetected through the beam pipe,
 1930 neutrinos produced in weak interactions processes which do not interact with the
 1931 detector and thus escaping without leaving a sign, or even undiscovered particles
 1932 predicted by models beyond the SM.

1933 The PF algorithm assigns the negative sum of the momenta of all reconstructed
 1934 PF particles to the *particle-flow MET* according to

$$\vec{E}_T = - \sum_i \vec{p}_{T,i} \quad (3.2)$$

1935 JEC are propagated to the calculation of the \vec{E}_T as described in the Reference [122].

1936 3.4.2 Event reconstruction examples

1937 Figures 3.7-3.9 show the results of the reconstruction performed on 3 recorded events.

1938 Descriptions are taken directly from the source.

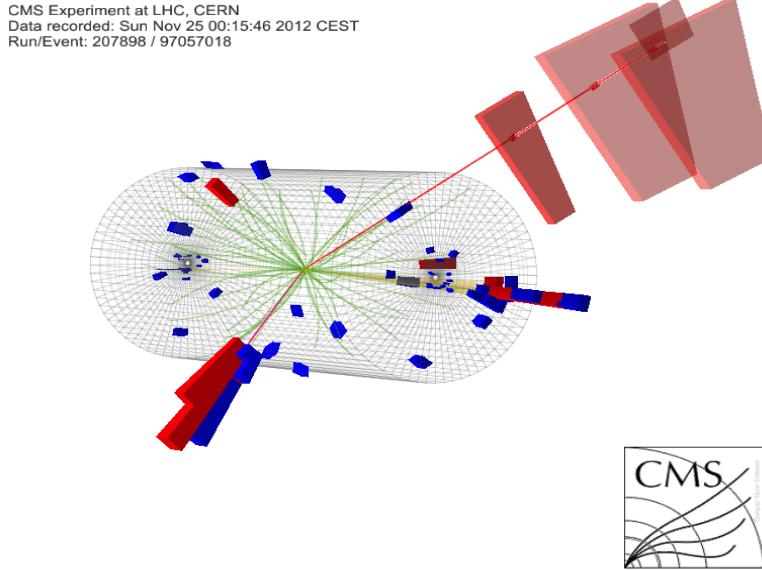


Figure 3.7: HIG-13-004 Event 1 reconstruction results; “HIG-13-004 Event 1: Event recorded with the CMS detector in 2012 at a proton-proton center-of-mass energy of 8 TeV. The event shows characteristics expected from the decay of the SM Higgs boson to a pair of τ leptons. Such an event is characterized by the production of two forward-going jets, seen here in opposite endcaps. One of the τ decays to a muon (red lines on the right) and neutrinos, while the other τ decays into a charged hadron and a neutrino.” [123].

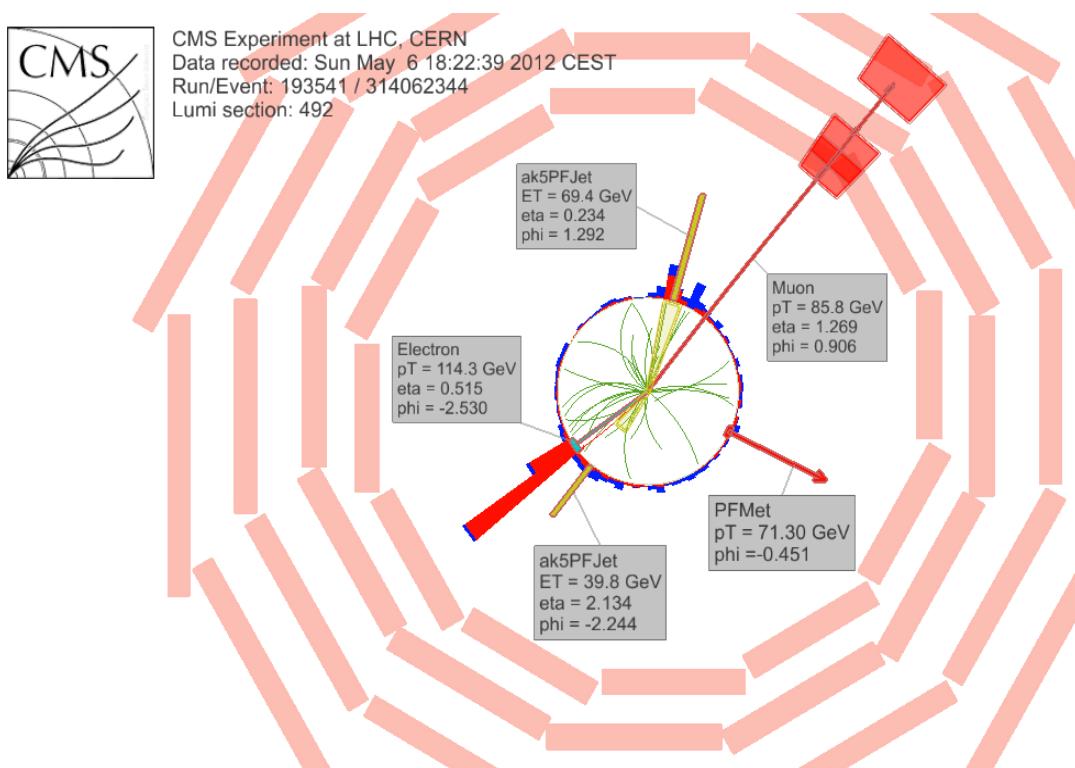


Figure 3.8: $e\mu$ event reconstruction results; “An $e\mu$ event candidate selected in 8 TeV data, as seen from the direction of the proton beams. The kinematics of the main objects used in the event selection are highlighted: two isolated leptons and two particle-flow jets. The reconstructed missing transverse energy is also displayed for reference” [124].

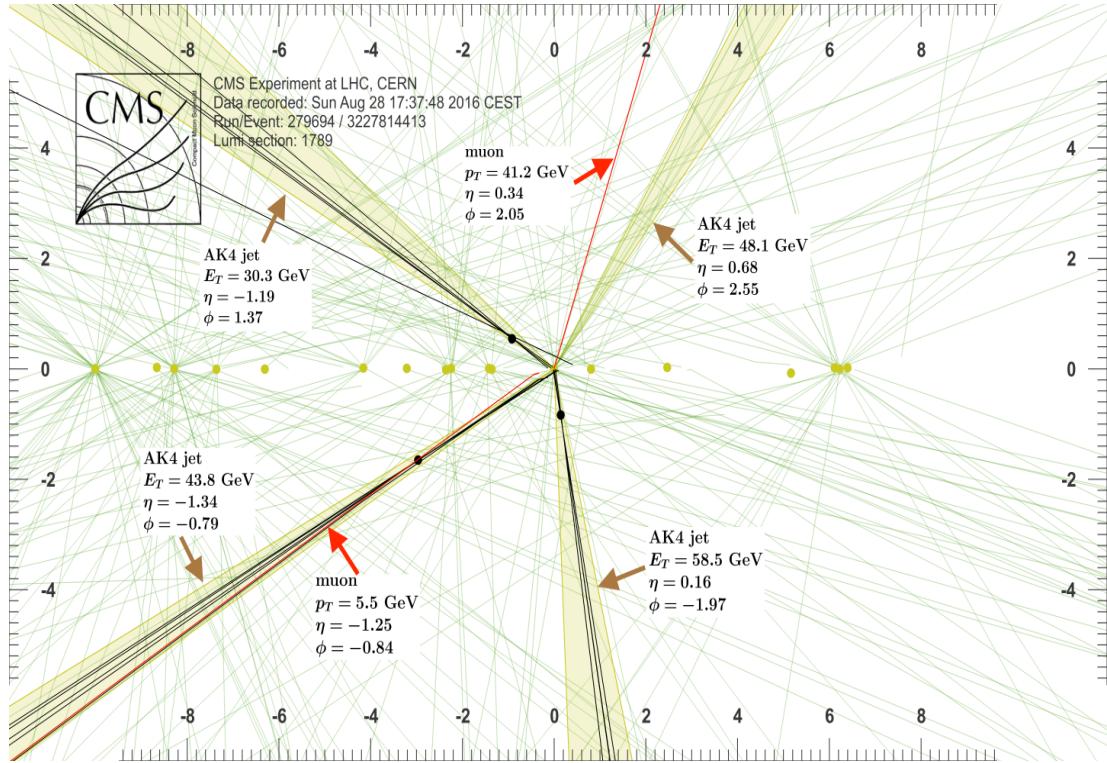


Figure 3.9: Recorded event reconstruction results; “Recorded event (ρ - z projection) with three jets with $p_T > 30$ GeV with one displaced muon track in 2016 data collected at 13 TeV. Each of the three jets has a displaced reconstructed vertex. The jet with $p_T(j) = 43.8$ GeV, $\eta(j) = -1.34$, $\phi(j) = -0.79$ contains muon with $p_T(\mu) = 5.5$ GeV, $\eta(\mu) = -1.25$, $\phi(\mu) = -0.84$. Event contains reconstructed isolated muon with $p_T(\mu) = 41.2$ GeV, $\eta(\mu) = 0.34$, $\phi(\mu) = 2.05$ and MET with $p_T = 72.5$ GeV, $\phi = -0.32$. Jet candidates for a b -jet from top quark leptonic and hadronic decays are tagged by CSVv2T algorithm. One of the other two jets is tagged by CharmT algorithm. Tracks with $p_T > 0.5$ GeV are shown. The number of reconstructed primary vertices is 18. Reconstructed $m_T(W)$ is 101.8 GeV. Beam spot position correction is applied. Reconstructed primary vertices are shown in yellow color, while reconstructed displaced vertices and associated tracks are presented in black color. Dimensions are given in cm” [125].

₁₉₃₉ **Chapter 5**

₁₉₄₀ **Statistical methods**

₁₉₄₁ In the course of analyzing the data sets provided by the CMS experiment and used in
₁₉₄₂ this thesis, several statistical tools have been employed; in this chapter, a description
₁₉₄₃ of these tools will be presented, starting with the general statement of the multivariate
₁₉₄₄ analysis methods, followed by the particularities of the Boosted Decision Trees (BDT)
₁₉₄₅ method and its application to the classification problem. Statistical inference methods
₁₉₄₆ used will also be presented. This chapter is based mainly on References [126–128].

₁₉₄₇ **5.1 Multivariate analysis**

₁₉₄₈ Multivariate data analysis (MVA) makes use of the statistical techniques developed to
₁₉₄₉ analyze more than one variable at once, taking into account all the correlations among
₁₉₅₀ variables. MVA is employed in a variety of fields like consumer and market research,
₁₉₅₁ quality control and process optimization. Using MVA it is possible to identify the
₁₉₅₂ dominant patterns in a data sample, like groups, outliers and trends, and determine
₁₉₅₃ to which group a set of values belong; in the particle physics context, MVA methods
₁₉₅₄ are used to perform the selection of certain type of events from a large data set.

1955 Processes with small cross section, such as the tHq process ($\sigma_{SM}(\sqrt{s} = 13\text{TeV}) =$
 1956 70.96 fb), are hard to detect in the presence of the processes with larger cross sections,
 1957 $\sigma_{SM}^{t\bar{t}}(\sqrt{s} = 13\text{TeV}) = 823.44$ fb for instance; therefore, only a small fraction of the data
 1958 contains events of interest (signal), the major part is signal-like events, which mimic
 1959 signal characteristics but belong to different processes, so they are a background to
 1960 the process of interest. This implies that it is not possible to say with certainty
 1961 that a given event is a signal or a background and statistical methods should be
 1962 involved. In that sense, the challenge can be formulated as one where a set of events
 1963 have to be classified according to certain special features; these features correspond
 1964 to the measurements of several parameters like energy or momentum, organized in a
 1965 set of *input variables*. The measurements for each event can be written in a vector
 1966 $\mathbf{x} = (x_1, \dots, x_n)$ for which

- 1967 • $f(\mathbf{x}|s)$ is the probability density (*likelihood function*) that \mathbf{x} is the set of mea-
 1968 sured values given that the event is a signal event (signal hypothesis).
- 1969 • $f(\mathbf{x}|b)$ is the probability density (*likelihood function*) that \mathbf{x} is the set of mea-
 1970 sured values given that the event is a background event (background hypothe-
 1971 sis).

1972 Figure 5.1 shows three ways to perform a classification of events for which mea-
 1973 surements of two properties, i.e., two input variables x_1 and x_2 , have been performed;
 1974 blue circles represent signal events while red triangles represent background events.
 1975 The classification on the left is *cut-based* requiring $x_1 < c_1$ and $x_2 < c_2$; usually the
 1976 cut values (c_1 and c_2) are chosen according to some knowledge about the event pro-
 1977 cess. In the middle plot, the classification is performed using a linear function of
 1978 the input variables, hence the boundary is a straight line, while in the right plot the

1979 the relationship between input variables is not linear thus the boundary is not linear
 1980 either.

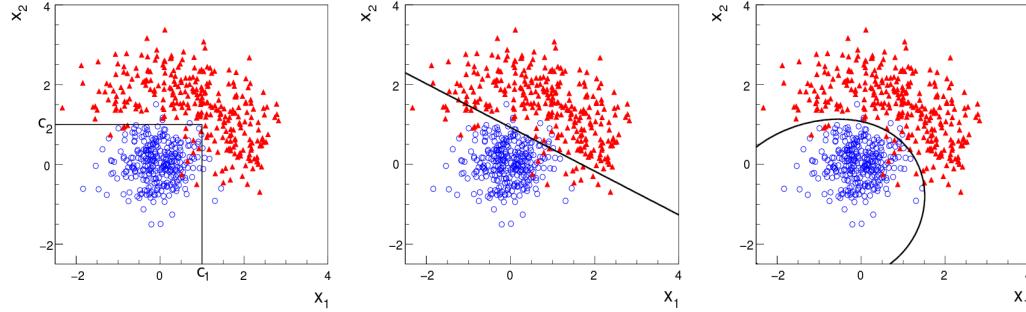


Figure 5.1: Scatter plots-MVA event classification. Distribution of two input variables x_1 and x_2 measured for a set of events; blue circles represent signal events and red triangles represent background events. The classification is based on cuts (left), linear boundary (center), and nonlinear boundary (right) [126]

1981 In general, the boundary can be parametrized in terms of the input variables such
 1982 that the cut is set on the parametrization instead of on the variables, i.e., $y(\mathbf{x}) = y_{cut}$
 1983 with y_{cut} being a constant; thus, the acceptance or rejection of an event is based on
 1984 which side of the boundary the event is located. If $y(\mathbf{x})$, usually called *test statistic*,
 1985 has functional form, it can be used to determine the probability distribution functions
 1986 $p(y|s)$ and $p(y|b)$ and then perform a test statistic with a single cut on the scalar
 1987 variable y .

1988 Figure 5.2 shows an example of what would be the probability distribution func-
 1989 tions under the signal and background hypotheses for a scalar test statistic with a cut
 1990 on the classifier y . Note that the tails of the distributions indicate that some signal
 1991 events fall in the rejection region and some background events fall on the acceptance
 1992 region; therefore, it is convenient to define the *efficiency* with which events of a given
 1993 type are accepted. The signal and background efficiencies are given by

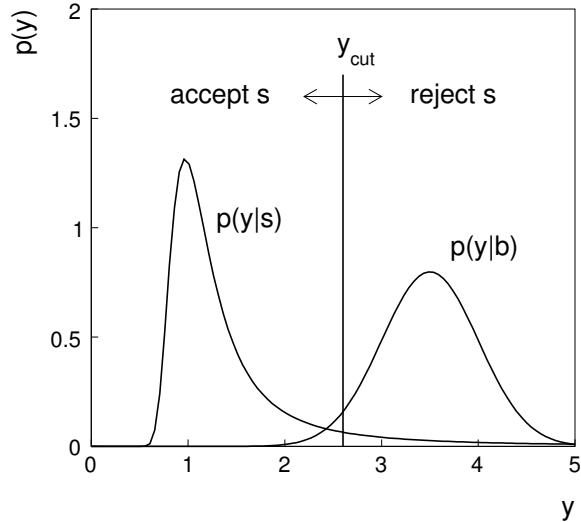


Figure 5.2: Distributions of the scalar test statistic $y(\mathbf{x})$ under the signal and background hypotheses. [126]

$$\varepsilon_s = P(\text{accept event}|s) = \int_A f(\mathbf{x}|s) d\mathbf{x} = \int_{-\infty}^{y_{cut}} p(y|s) dy , \quad (5.1)$$

$$\varepsilon_b = P(\text{accept event}|b) = \int_A f(\mathbf{x}|b) d\mathbf{x} = \int_{-\infty}^{y_{cut}} p(y|b) dy , \quad (5.2)$$

1994 where A is the acceptance region. If the background hypothesis is the *null hypothesis*
 1995 (H_0), the signal hypothesis would be *alternative hypothesis* (H_1); in this context, the
 1996 background efficiency corresponds to the significance level of the test (α) and describes
 1997 the misidentification probability, while the signal efficiency corresponds to the power
 1998 of the test ($1-\beta$)¹ and describes the probability of rejecting the background hypothesis
 1999 if the signal hypothesis is true. What is sought in an analysis is to maximize the power
 2000 of the test relative to the significance level, i.e., set a selection with the largest possible
 2001 selection efficiency and the smallest possible misidentification probability.

¹ β is the fraction of signal events that fall out of the acceptance region

2002 **5.1.1 Decision trees**

2003 For this thesis, the implementation of the MVA strategy, described above, is per-
 2004 formed through decision trees by using the TMVA software package [127] included
 2005 in the ROOT analysis framework [129]. In a simple picture, a decision tree classifies
 2006 events according to their input variables values by setting a cut on each input variable
 2007 and checking which events are on which side of the cut, just as proposed in the MVA
 2008 strategy, but in addition, as a machine learning algorithm, decision trees offer the
 2009 possibility to be trained and then perform the classification efficiently.

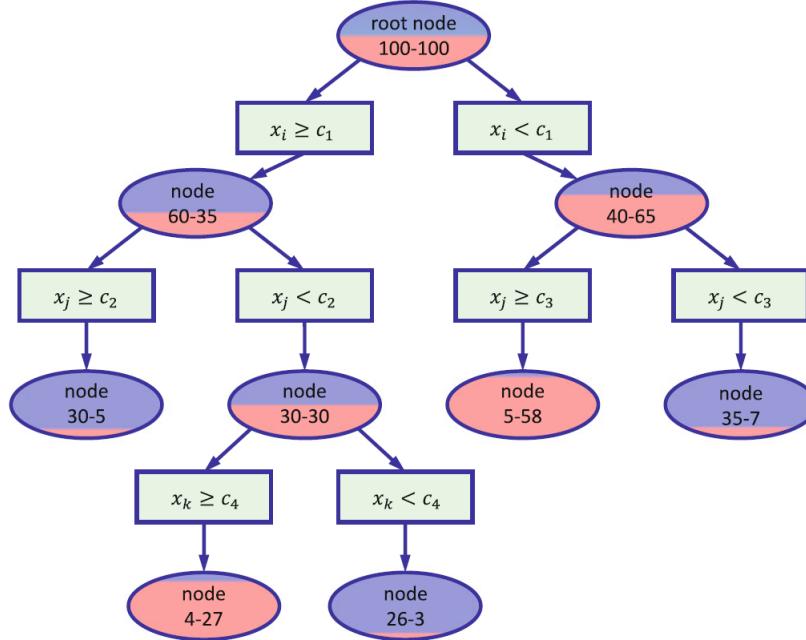


Figure 5.3: Example of a decision tree. Each node is fed with a MC sample mixing signal and background events (left-right numbers); nodes colors represent the relative number of signal/background events [128].

2010 The training or growing of a decision tree is the process where the rules for clas-
 2011 sifying events are defined; this process is represented in Figure 5.3 and consists of
 2012 several steps:

- 2013 • take MC samples of signal and background events and split them into two parts

2014 each; the first parts will be used in the decision tree training, while the second
 2015 parts will be used for testing the final classifier obtained from the training.
 2016 Each event has associated a set of input variables $\mathbf{x} = (x_1, \dots, x_n)$ which serve
 2017 to distinguish between signal and background events. The training sample is
 2018 taken in at the *root node*.

- 2019 • Pick one variable, say x_i .
- 2020 • Pick one value of x_i , each event has its own value of x_i , and split the training
 2021 sample into two subsamples B_1 and B_2 ; B_1 contains events for which $x_i < c_1$
 2022 while B_2 contains the rest of the training events;
- 2023 • scan all possible values of x_i and find the splitting value that provides the *best*
 2024 classification², i.e., B_1 is mostly made of signal events while B_2 is mostly made
 2025 of background events.
- 2026 • It is possible that variables other than the picked one produce a better classi-
 2027 fication, hence, all the variables have to be evaluated. Pick the next variable,
 2028 say x_j , and repeat the scan over its possible values.
- 2029 • At the end, all the variables and their values will have been scanned, the *best*
 2030 variable and splitting value will have been identified, say x_1, c_1 , and there will
 2031 be two nodes fed with the subsamples B_1 and B_2 .

2032 Nodes are further split by repeating the decision process until a given number of
 2033 final nodes is obtained, nodes are largely dominated by either signal or background
 2034 events, or nodes have too few events to continue. Final nodes are called *leaves* and
 2035 they are classified as signal or background leaves according to the class of the majority
 2036 of events in them. Each *branch* in the tree corresponds to a sequence of cuts.

² Quality of the classification will be treated in the next paragraph.

2037 The quality of the classification at each node is evaluated through a separation
 2038 criteria; there are several of them but the *Gini Index* (G) is the one used in the
 2039 decision trees trained for the analysis in this thesis. G is written in terms of the
 2040 purity (P), i.e., the fraction of signal events in the samples after the separation is
 2041 made; it is given by

$$G = P(1 - P) \quad (5.3)$$

2042 note that $P=0.5$ at the root node while $G=0$ for pure leaves. For a node A split into
 2043 two nodes B_1 and B_2 the G gain is

$$\Delta G = G(A) - G(B_1) - G(B_2). \quad (5.4)$$

2044 The *best* classification corresponds to that for which the gain of G is maximized;
 2045 hence, the scanning over all the variables in an event and their values is of great
 2046 importance.

2047 In order to provide a numerical output for the classification, events in a sig-
 2048 nal(background) leaf are assigned a score of 1(-1) each, defining in this way the
 2049 decision tree *classifier/weak learner* as

$$f(\mathbf{x}) = \begin{cases} 1 & \mathbf{x} \text{ in signal region,} \\ -1 & \mathbf{x} \text{ in background region.} \end{cases}$$

2050 Figure 5.4 shows an example of the classification of a sample of events, containing
 2051 two variables, performed by a decision tree.

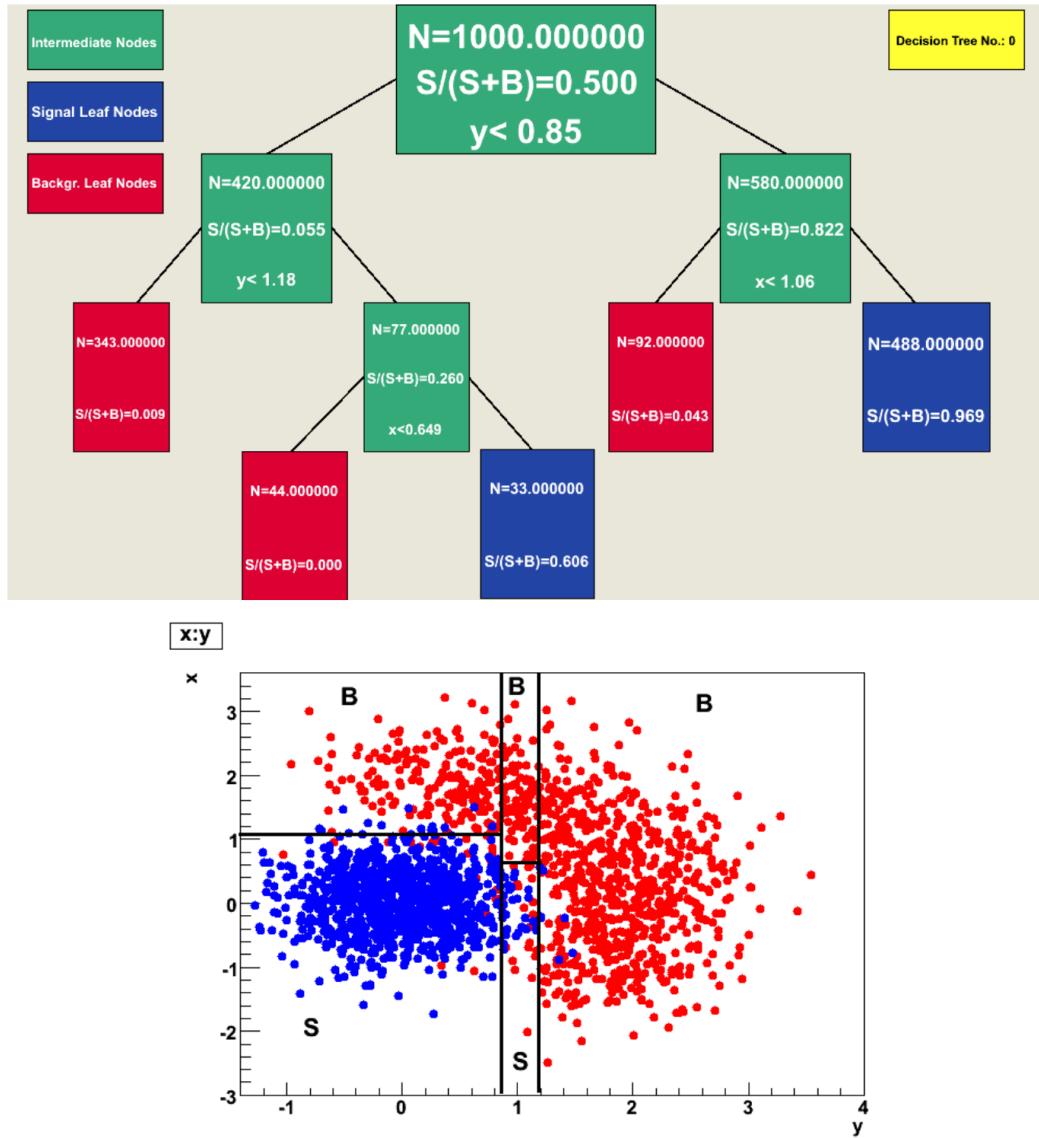


Figure 5.4: Example of a decision tree output. Each leaf, blue for signal events and red for background events, is represented by a region in the variables phase space [130].

2052 5.1.2 Boosted decision trees (BDT).

2053 Event misclassification occurs when a training event ends up in the wrong leaf, i.e., a
 2054 signal event ends up in a background leaf or a background event ends up in a signal
 2055 leaf. A way to correct it is to assign a weight to the misclassified events and train
 2056 a second tree using the reweighted events; the event reweighting is performed by a

2057 boosting algorithm in such a way that when used in the training of a new decision
 2058 tree the *boosted events* get correctly classified. The process is repeated iteratively
 2059 adding a new tree to the forest and creating a set of classifiers, which are combined
 2060 to create the next classifier; the final classifier offers more stability³ and has a smaller
 2061 misclassification rate than any individual ones. The resulting tree collection is known
 2062 as a *boosted decision tree (BDT)*.

2063 Thus, purity of the sample is generalized to

$$P = \frac{\sum_s w_s}{\sum_s w_s + \sum_b w_b} \quad (5.5)$$

2064 where w_s and w_b are the weights of the signal and background events respectively;
 2065 the Gini index is also generalized

$$G = \left(\sum_i^n w_i \right) P(1 - P) \quad (5.6)$$

2066 with n the number of events in the node. The final score of an event, after pass-
 2067 ing through the forest, is calculated as the renormalized sum of all the individual
 2068 (possibly weighted) scores; thus, high(low) score implies that the event is most likely
 2069 signal(background).

2070 The boosting procedure, implemented in the *Gradient boosting* algorithm used in
 2071 this thesis, produces a classifier $F(\mathbf{x})$ which is the weighted sum of the individual
 2072 classifiers obtained after each iteration, i.e.,

$$F(\mathbf{x}) = \sum_{m=1}^M \beta_m f(\mathbf{x}; a_m) \quad (5.7)$$

2073 where M is the number of trees in the forest. The *loss function* $L(F, y)$ represents the

³ Decision trees suffer from sensitivity to statistical fluctuations in the training sample which may lead to very different results with a small change in the training samples.

2074 deviation between the classifier $F(\mathbf{x})$ response and the true value y obtained from the
 2075 training sample (1 for signal events and -1 for background event), according to

$$L(F, y) = \ln(1 + e^{-2F(\mathbf{x})y}) \quad (5.8)$$

2076 thus, the reweighting is employed to ensure the minimization of the loss function; a
 2077 more detailed description of the minimization procedure can be found in Reference
 2078 [131]. The final classifier output is later used as a final discrimination variable, labeled
 2079 as *BDT output/response*.

2080 5.1.3 Overtraining

2081 Decision trees offer the possibility to have as many nodes as desired in order to
 2082 reduce the misclassification to zero (in theory); however, when a classifier is too much
 2083 adjusted to a particular training sample, the classifier's response to a slightly different
 2084 sample may leads to a completely different classification results; this effect is known
 2085 as *overtraining*.

2086 An alternative to reduce the overtraining in BDTs consists in pruning the tree
 2087 by removing statistically insignificant nodes after the tree growing is completed but
 2088 this option is not available for BDTs with gradient boosting in the TMVA-toolkit,
 2089 therefore, the overtraining has to be reduced by tuning the algorithm, number of
 2090 nodes, minimum number of events in the leaves, etc. The overtraining can be evaluated
 2091 by comparing the responses of the classifier when running over the training and
 2092 test samples.

2093 5.1.4 Variable ranking

2094 BDTs have a couple of particular advantages related to the input variables; they are
 2095 relatively insensitive to the number of input variables used in the vector \mathbf{x} . The
 2096 ranking of the BDT input variables is determined by counting the number of times a
 2097 variable is used to split decision tree nodes; in addition, the separation gain-squared
 2098 achieved in the splitting and the number of events in the node are accounted by
 2099 applying a weighting to that number. Thus, those variables with small or no power
 2100 to separate signal and background events are rarely chosen to split the nodes, i.e., are
 2101 effectively ignored.

2102 In addition, variables correlations play an important role for some MVA methods
 2103 like the Fisher discriminant algorithm in which the first step consist of performing a
 2104 linear transformation to a phase space where the correlations between variables are
 2105 removed; in the case of BDT algorithm, correlations do not affect the performance.

2106 5.1.5 BDT output example

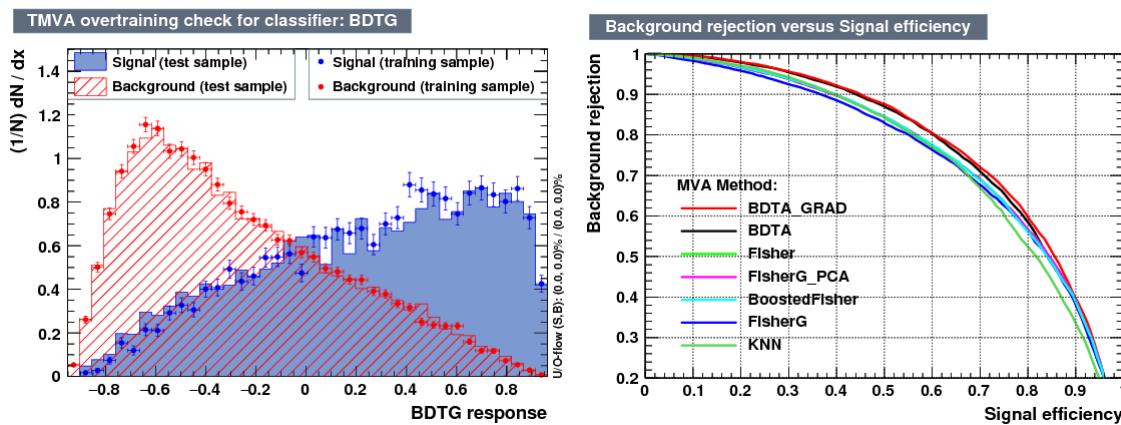


Figure 5.5: Left: Output distributions for the gradient boosted decision tree (BDTG) classifier using a sample of signal ($pp \rightarrow tHq$) and background ($pp \rightarrow tt$) events. Right: Background rejection vs signal efficiency (ROC curves) for various MVA classifiers running over the same sample used to produce the plot on the left.

2107 The left side of figure 5.5 shows the BDT output distributions for signal ($pp \rightarrow$
 2108 tHq) and background ($pp \rightarrow t\bar{t}$) events; this plot is the equivalent to the one showed
 2109 in Figure 5.2. A forest with 800 trees, maximum depth per tree = 3, and gradient
 2110 boosting have been used as training parameters. The BDTG classifier offers a good
 2111 separation power. There is a small overtraining in the signal distribution, while the
 2112 background distribution is very well predicted which might indicate that the sample
 2113 is composed of more background than signal events.

2114 The right side of figure 5.5 shows the background rejection vs signal efficiency
 2115 curves for several combinations of MVA classifiers-boosting algorithms running over
 2116 the same MC sample; these curves are known as ROC curves and give an indication
 2117 of the performance of the classifier. In this particular example, the best performance
 2118 is achieved with the BDTG classifier (BDTA_GRAD), which motivate its use in this
 2119 thesis.

2120 **5.2 Statistical inference**

2121 Once events are classified, the next step consists of finding the parameters that define
 2122 the likelihood functions $f(\mathbf{x}|s)$, $f(\mathbf{x}|b)$ for signal and background events respectively.
 2123 In general, likelihood functions depend not only on the measurements but also on
 2124 parameters (θ_m) that define their shapes; the process of estimating these *unknown*
 2125 *parameters* and their uncertainties from the experimental data is called *inference*.

2126 The statistical inference tools used in this analysis are implemented in the RooFit
 2127 toolkit [132] and COMBINE package [133] included in the CMSSW software frame-
 2128 work.

2129 **5.2.1 Nuisance parameters**

2130 The unknown parameter vector θ is made of two types of parameters: those pa-
 2131 rameters that provide information about the physical observables of interest for the
 2132 experiment or *parameters of interest*, and the *nuisance parameters* that are not of
 2133 direct interest for the experiment but that need to be included in the analysis in
 2134 order to achieve a satisfactory description of the data; they represent effects of the
 2135 detector response like the finite resolutions of the detection systems, miscalibrations,
 2136 and in general any source of uncertainty introduced in the analysis.

2137 Nuisance parameters can be estimated from experimental data; for instance, data
 2138 samples from a test beam are usually employed for calibration purposes. In cases
 2139 where experimental samples are not availables, the estimation of nuisance parameters
 2140 makes use of dedicated simulation programs to provide the required samples.

2141 The estimation of the unknown parameters involves certain deviations from their
 2142 true values, hence, the measurement of the nuisance parameter is written in terms
 2143 of an estimated value, also called central value, $\hat{\theta}$ and its uncertainty $\delta\theta$ using the
 2144 notation

$$\theta = \hat{\theta} \pm \delta\theta \quad (5.9)$$

2145 where the interval $[\hat{\theta} - \delta\theta, \hat{\theta} + \delta\theta]$ is called *confidence interval*; it is usually interpreted,
 2146 in the limit of infinite number of experiments, as the interval where the true value
 2147 of the unknown parameter θ is contained with a probability of 0.6827 (if no other
 2148 convention is stated); this interval represents the area under a Gaussian distribution
 2149 in the interval $\pm 1\sigma$.

2150 The uncertainties associated with nuisance parameters produce *systematic uncer-*
 2151 *tainties* in the final measurement, while the uncertainties related only to fluctuations

2152 in data and that affect the determination of parameters of interest produce *statistical*
 2153 *uncertainties*.

2154 **5.2.2 Maximum likelihood estimation method**

2155 The estimation of the unknown parameters that are in best agreement with the ob-
 2156 served data is performed through a function of the data sample that returns the
 2157 estimate of those parameters; that function is called an *estimator*. Estimators are
 2158 usually constructed using mathematical expressions encoded in algorithms.

2159 In this thesis, the estimator used is the likelihood function $f(\mathbf{x}|\boldsymbol{\theta})$ ⁴ which depends
 2160 on a set of measured variables \mathbf{x} and a set of unknown parameters $\boldsymbol{\theta}$. The likelihood
 2161 function for N events in a sample is the combination of all the individual likelihood
 2162 functions, i.e.,

$$L(\boldsymbol{\theta}) = \prod_{i=1}^N f(\mathbf{x}^i|\boldsymbol{\theta}) = \prod_{i=1}^N f(x_1^i, \dots, x_n^i; \theta_1, \dots, \theta_m) \quad (5.10)$$

2163 and the estimation method used is the *Maximum Likelihood Estimation* method
 2164 (MLE); it is based on the combined likelihood function defined by eqn. 5.10 and
 2165 the procedure seeks for the parameter set that corresponds to the maximum value of
 2166 the combined likelihood function, i.e., the *maximum likelihood estimator* of the un-
 2167 known parameter vector $\boldsymbol{\theta}$ is the function that produces the vector of *best estimators*
 2168 $\hat{\boldsymbol{\theta}}$ for which the likelihood function $L(\boldsymbol{\theta})$ evaluated at the measured \mathbf{x} is maximum.

2169 Usually, the logarithm of the likelihood function is used in numerical algorithm
 2170 implementations in order to avoid underflow the numerical precision of the computers
 2171 due to the product of low likelihoods. In addition, it is common to minimize the
 2172 negative logarithm of the likelihood function, therefore, the negative log-likelihood

⁴ analogue to the likelihood functions described in previous sections

2173 function is

$$F(\boldsymbol{\theta}) = -\ln L(\boldsymbol{\theta}) = -\sum_{i=1}^N f(\mathbf{x}^i | \boldsymbol{\theta}). \quad (5.11)$$

2174 The minimization process is performed by the software MINUIT [134] implemented in the ROOT analysis framework. In case of data samples with large number 2175 of measurements, the computational resources necessary to calculate the likelihood 2176 function are too big; therefore, the parameter estimation is performed using binned 2177 distributions of the variables of interest for which the *binned likelihood function* is 2178 given by

$$L(\mathbf{x}|r, \boldsymbol{\theta}) = \prod_{i=1} \frac{(r \cdot s_i(\boldsymbol{\theta}) + b_i(\boldsymbol{\theta}))^{n_i}}{n_i!} e^{-r \cdot s_i(\boldsymbol{\theta}) - b_i(\boldsymbol{\theta})} \prod_{j=1} \frac{1}{\sqrt{2\pi}\sigma_{\theta_j}^2} e^{-(\theta_j - \theta_{0,j})^2/2\sigma_{\theta_j}^2}, \quad (5.12)$$

2180 with s_i and b_i the expected number of signal and background yields for the bin i , n_i 2181 is the observed number of events in the bin i and $r = \sigma/\sigma_{SM}$ is the signal strength. 2182 Note that the number of entries per bin follows a Poisson distribution. The effect 2183 of the nuisance parameters have been included in the likelihood function through 2184 the multiplication by a Gaussian distribution that models the nuisance. The three 2185 parameters, r , s_i and b_i are jointly fitted to estimate the value of r .

2186 5.3 Upper limits

2187 In this analysis, two hypotheses are considered; the background only hypothesis 2188 ($H_0(b)$) and the signal plus background hypothesis ($H_1(s+b)$), i.e., the sample of 2189 events is composed of background only events ($r=0$) or it is a mixture of signal plus 2190 background events ($r=1$). The exclusion of one hypothesis against the other means 2191 that the observed data sample better agrees with H_0 or rather with H_1 . In order 2192 to discriminate these hypotheses, a test statistic is constructed on the basis of the

2193 likelihood function evaluated for each of the hypothesis.

2194 The *Neyman-Pearson* lemma [135] states that the test statistic that provides the
 2195 maximum power for H_1 for a given significance level (background misidentification
 2196 probability α), is given by the ratio of the likelihood functions $L(\mathbf{x}|H_1)$ and $L(\mathbf{x}|H_0)$;
 2197 however, in order to use that definition it is necessary to know the true likelihood
 2198 functions, which in practice is not always possible. Approximate functions obtained
 2199 by numerical methods, like the BDT method described above, have to be used, so
 2200 that the *profile likelihood* test statistic is defined by

$$\lambda(\mathbf{r}) = \frac{L(\mathbf{x}|r, \hat{\boldsymbol{\theta}}(r))}{L(\mathbf{x}|\hat{r}, \hat{\boldsymbol{\theta}})}, \quad (5.13)$$

2201 where, \hat{r} and $\hat{\boldsymbol{\theta}}$ maximize the likelihood function, and $\hat{\boldsymbol{\theta}}$ maximizes the likelihood
 2202 function for a given value of the signal strength modifier r . In practice, the test
 2203 statistic t_r

$$t_r = -2\ln\lambda(r) \quad (5.14)$$

2204 is used to evaluate the presence of signal in the sample, since the minimum of t_r at
 2205 $r = \hat{r}$ suggests the presence of signal with signal strength \hat{r} . The uncertainty interval
 2206 for r is determined by the values of r for which $t_r = +1$.

2207 The expected probability density function (p.d.f) $f(t_r|r, \boldsymbol{\theta})$ of the test statistic t_r
 2208 can be obtained numerically by generating MC samples where one hypothesis, $H_0(b)$
 2209 or $H_1(s+b)$, is assumed; thus, MC samples contain the possible values of t_r obtained
 2210 from *pseudo-experiments* as shown in Figure 5.6. The probability that t_r takes a value
 2211 equal or greater than the observed value ($t_{r,obs}$) when a signal with a signal modifier
 2212 r is present in the data sample, is called the *p-value* of the observation; it can be
 2213 calculated using

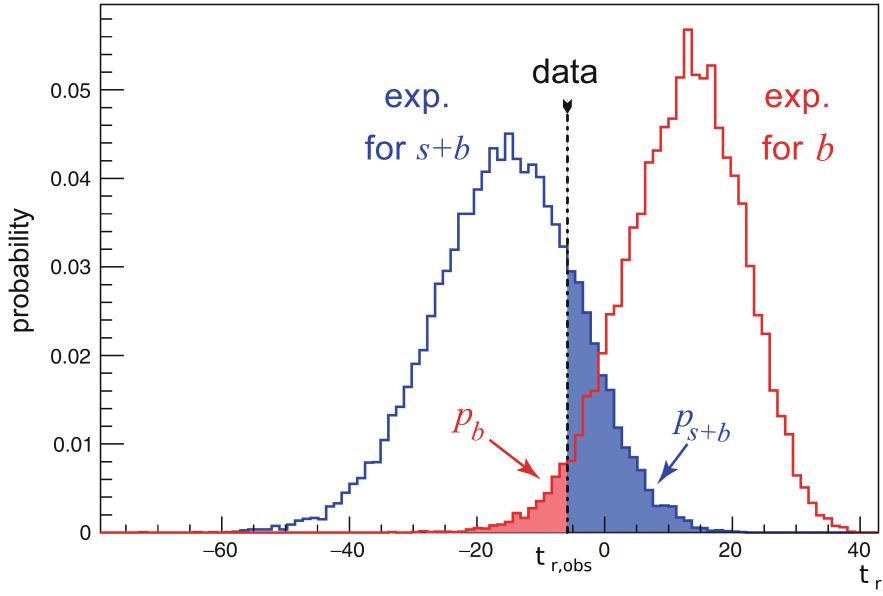


Figure 5.6: t_r p.d.f. from MC pseudo experiments assuming H_0 (red) and H_1 (blue). The black dashed line shows the value of the test statistic as measured from data. Adapted from Reference [128].

$$p_r = \int_{t_{r,obs}}^{\infty} f(t'_r | r, \boldsymbol{\theta}) dt'_r, \quad (5.15)$$

thus, $p_r < 0.05$ means that, for that particular value of r , H_1 could be excluded at 95% Confidence Level (CL). The corresponding background-only p-value is given by

$$1 - p_b = \int_{t_{r,obs}}^{\infty} f(t'_r | 0, \boldsymbol{\theta}) dt'_r, \quad (5.16)$$

If the t_r p.d.f.s for both hypotheses are well separated, as shown in the top side of Figure 5.7, the experiment is sensitive to the presence of signal in the sample. If the signal presence is small, both p.d.f.s will be largely overlapped (bottom of Figure ??) and either the signal hypothesis could be rejected with not enough justification because the experiment is not sensitive to the signal or a fluctuation of the background could be misinterpreted as presence of signal with the corresponding rejection of the

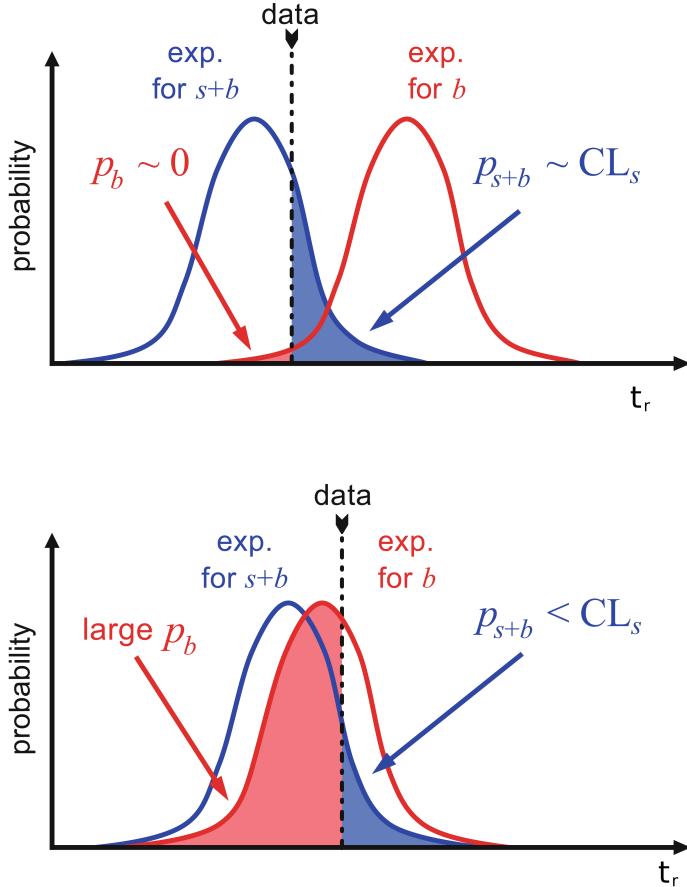


Figure 5.7: CL_s limit illustration. When the test statistic p.d.f. for the two hypotheses H_0 and H_1 are well separated (top) and when they are largely overlapped (bottom). Adapted from Reference [128].

background-only hypothesis. These issues are corrected by using the modified p-value [136]

$$p'_r = \frac{p_r}{1 - p_b} \equiv CL_s. \quad (5.17)$$

If H_1 is true, then p_b is small, $CL_s \simeq p_r$ and H_0 is rejected; if there is large overlap and a statistical fluctuation cause that p_b is large, then both numerator and denominator in Eqn. 5.17 become small but CL_s would allow the rejection of H_1 even if there is poor sensitivity to signal.

2228 The upper limit of the parameter of interest r^{up} is determined by excluding the
 2229 range of values of r for which $CL_s(r, \theta)$ is lower than the confidence level desired,
 2230 normally 90% or 95%, e.g, scanning over r and finding the value for which $p_r'^{up} =$
 2231 0.05. The expected upper limit can be calculated using pseudo-experiments based on
 2232 the background-only hypothesis and obtaining a distribution for r_{ps}^{up} ; the median of
 2233 that distribution corresponds to the expected upper limit, while the $\pm 1\sigma$ and $\pm 2\sigma$
 2234 deviations correspond to the values of the distribution that defines the 68% and 95%
 2235 of the area under the distribution centered in the median. It is usual to present all
 2236 the information about the expected and observed limits in the so-called *Brazilian-flag*
 2237 plot as the one showed in Figure 5.8. The solid line represent the observed CL_s

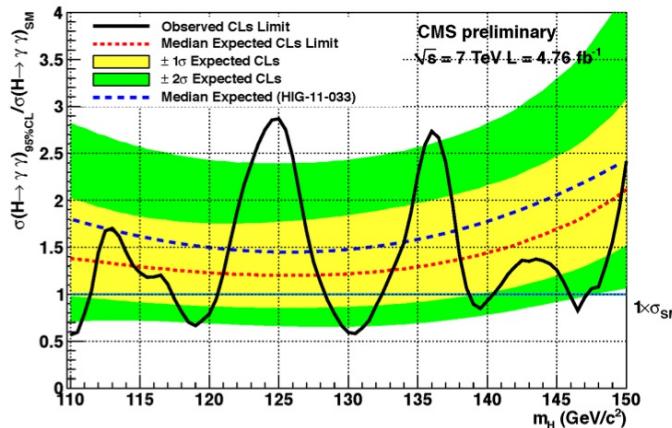


Figure 5.8: Brazilian flag plot of CMS experiment limits for Higgs boson decaying to photons [137].

2238 5.4 Asymptotic limits

2239 As said before, the complexity of the likelihood functions, the construction of test
 2240 statistics, and the calculation of the limits and their uncertainties is not always man-
 2241 ageable and requires extensive computational resources; in order to overcome those
 2242 issues, asymptotic approximations for likelihood-based test statistics, like the ones

described in previous sections, have been developed [138, 139] using Wilks' theorem.
Asymptotic approximations replace the construction of the test statistics p.d.f.s using
MC pseudo-experiments, with the approximate calculation of the test statistics p.d.f.s
by employing the so-called *Asimov dataset*.

The Asimov dataset is defined as the dataset that produce the true values of the
nuisance parameters when it is used to evaluate the estimators for all the parameters;
it is obtained by setting the values of the variables in the dataset to their expected
values [139].

Limits calculated by using the asymptotic approximation and the Asimov dataset
are know as *asymptotic limits*.

2253 **Chapter 6**

2254 **Search for production of a Higgs**

2255 **boson and a single top quark in**

2256 **multilepton final states in pp**

2257 **collisions at $\sqrt{s} = 13$ TeV**

2258 **6.1 Introduction**

2259 The Higgs boson discovery, supported on experimental observations and theoretical
2260 predictions made about the SM, gives the clue of the way in that elementary particles
2261 acquire mass through the Higgs mechanism; therefore, knowing the Higgs mass, the
2262 Higgs-vector boson and Higgs-fermion couplings can be determined. In order to test
2263 the Higgs-top coupling, several measurements have been performed, as stated in the
2264 chapter 1, but they are limited in sensitivity to measure the square of the coupling.
2265 The production of a Higgs boson in association with a single top quark (tH) not
2266 only offers access to the sign of the coupling, but also, to the CP phase of the Higgs

2267 couplings.

2268 This chapter presents the search for the associated production of a Higgs boson
 2269 and a single top quark (tHq) events, focusing on leptonic signatures provided by the
 2270 Higgs decay modes to WW , ZZ , and $\tau\tau$; the 13 TeV dataset produced in 2016, which
 2271 corresponds to an integrated luminosity of 35.9fb^{-1} , is used.

2272 As shown in Section 1.5, the SM cross section of tHq process is driven by a
 2273 destructive interference between two contributions (see Figure 1.15), where the Higgs
 2274 couples to either the W boson or the top quark; however, if the sign of the Higgs-
 2275 top coupling is flipped with respect to the SM prediction, a large enhancement of
 2276 the cross section occurs, making this analysis sensitive to such deviation. A second
 2277 process, where the Higgs boson and top quark are accompanied by a W boson (tHW)
 2278 has similar behavior, albeit with a weaker interference pattern and lower contribution
 2279 to the cross section, therefore, a combination of both processes would increase the
 2280 sensitivity to the sign of the coupling; in this analysis both contributions are combined
 2281 and referred as tH channel. A third contribution comes from $t\bar{t}H$ process. The purpose
 2282 of this analysis is to investigate the exclusion of the presence of the $tH + t\bar{t}H$ processes
 2283 under the assumption of the anomalous Higgs-top coupling modifier ($\kappa_t = -1$). The
 2284 analysis exploits signatures with two leptons of the same sign ($2lss$) channel and three
 2285 leptons ($3l$) channel in the final state.

2286 Constraints on the sign of the Higgs-top coupling (y_t) have been derived from the
 2287 decay rate of Higgs boson to photon pairs [50] and from the cross section for associated
 2288 production of Higgs and Z bosons via gluon fusion [141], with recent results disfavoring
 2289 negative signs of the coupling [44, 59, 142], although the negative sign coupling have
 2290 not been completely excluded.

2291 The analysis presented here, expands previous analyses performed at 8 TeV [143,
 2292 144] and searches for associated production of $t\bar{t}$ pair and a Higgs boson in the multi-

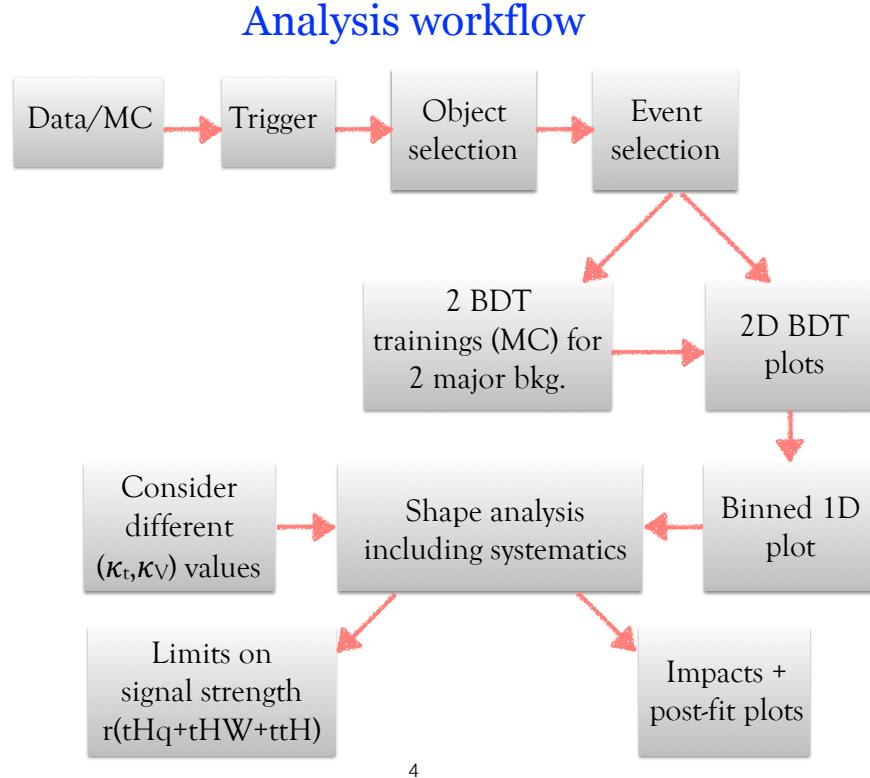
2293 lepton final state channel [145]; it also complements searches in other decay channels
 2294 targeting $H \rightarrow b\bar{b}$ [146].

2295 The first sections present the characteristic tHq signature as well as the expected
 2296 backgrounds. The MC samples, data sets, and the physics object definitions are
 2297 then defined. Following, the background predictions, the signal extraction, and the
 2298 statistical treatment of the selected events as well as the systematic uncertainties are
 2299 described. The final section present the results for the exclusion limits as a function
 2300 of the ratio of κ_t and the dimensionless modifier of the Higgs-vector boson coupling
 2301 κ_V .

2302 The analysis is designed to efficiently identify and select prompt leptons from on-
 2303 shell W and Z boson decays and to reject non-prompt leptons from b quark decays
 2304 and spurious lepton signatures from hadronic jets. Events are then selected in the
 2305 $2lss$ and $3l$ channels, and are required to contain hadronic jets, some of which must
 2306 be consistent with b quark hadronization. Finally, the signal yield is extracted by
 2307 simultaneously fitting the output of two dedicated multivariate discriminants, trained
 2308 to separate the tHq signal from the two dominant backgrounds, in all categories. The
 2309 fit result is then used to set an upper limit on the combined $t\bar{t}H + tH$ production
 2310 cross section, as a function of the relative coupling strengths of Higgs-top quark and
 2311 Higgs-Vector boson. Figure 6.1 shows an schematic overview of the analysis strategy
 2312 workflow.

2313 With respect to the 8 TeV analysis, the object selections have been adjusted for
 2314 the updated LHC running conditions at 13 TeV, the lepton identification has been
 2315 improved, and more powerful multivariate analysis techniques are used for the signal
 2316 extraction.

2317 The analysis has been made public by CMS as a Physics Analysis Summary [147]
 2318 combining the result for the three lepton and two lepton same-sign channels; the



4

Figure 6.1: A schematic overview of the analysis workflow. Based on sets of optimized physics object definitions and selection criteria, signal and background events in a data sample are discriminated. The discrimination is performed by a BDT, previously trained using MC samples of the dominant backgrounds, using discriminant variables based on the b -jet multiplicity, the activity in the forward region of the detector, and the kinematic properties of leptons. The CL_s limits on the combined $t\bar{t}H + tH$ production cross section, as a function of the relative coupling strengths are calculated.

2319 content present in this chapter is based on that document and on References [145,149]
 2320 unless other Reference is stated. Currently, an effort to turn the analysis into a paper
 2321 combining the multilepton and $H \rightarrow b\bar{b}$ is ongoing.

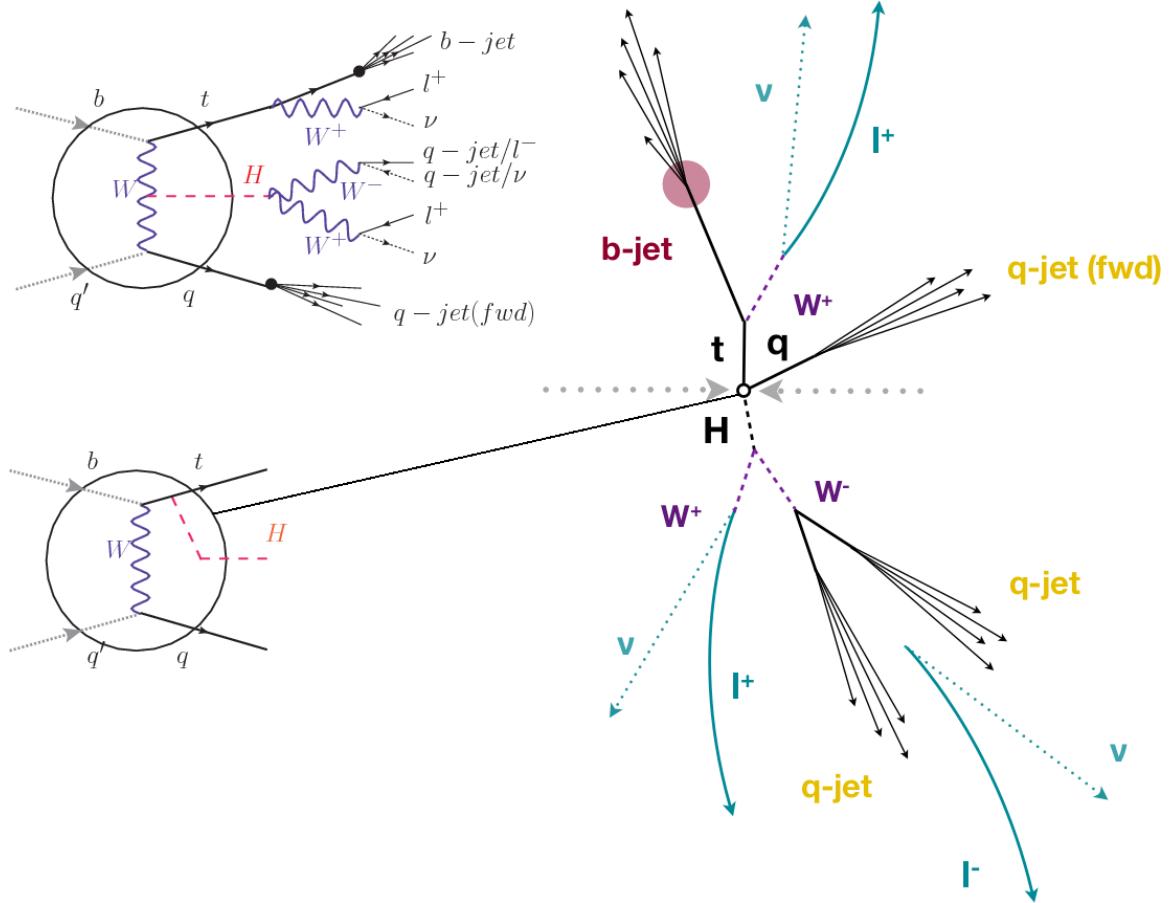


Figure 6.2: tHq event signature. Left: Feynman diagram including the whole evolution up to the final state for the case of the Higgs boson emitted by the W boson (top); Feynman diagram for the case where the Higgs boson is emitted by the top quark. Right: Schematic view as it would be seen in the detector; the circle in the Feynman diagrams on the left corresponds to the circle in the center of the schematic view as indicated by the line connecting them. In the $2lss$ channel, one of the W bosons from the Higgs boson decays to two light-quark jets while in the $3l$ channel both W bosons decay to leptons.

2322 6.2 tHq signature

2323 In order to select events of tHq process, its features are translated into a set of
 2324 selection rules; Figure 6.2 shows the Feynman diagram and an schematic view of the
 2325 tHq process from the pp collision to the final state configuration. A single top quark
 2326 is produced accompanied by a light quark, denoted as q ; this light quark is produced

2327 predominantly in the forward region of the detector. The Higgs boson can be either
 2328 emitted by the exchanged W boson or directly by the singly produced top quark.

2329 Due to their high masses/short lifetimes, top quark and Higgs boson decay after
 2330 their production within the detector. The Higgs boson is required to decay into a W
 2331 boson pair¹. The top quark almost always decays into a bottom quark and a W boson,
 2332 as encoded in the CMK matrix. The W bosons are required to decay leptonically
 2333 either all the three in the $3l$ channel or the pair with equal electrical charge in the
 2334 $2lss$ channel case; τ leptons are not reconstructed separately and only their leptonic
 2335 decays into either electrons or muons are considered in this analysis.

2336 In summary, the signal process is characterized by a the final state with

2337 • one light-flavored forward jet,

2338 • one central b-jet,

2339 • $2lss$ channel → two leptons of the same sign, two neutrinos and two light (often
 2340 soft) jets,

2341 • $3l$ channel → three leptons, three neutrinos and no central light-flavored jets,

2342 The presence of neutrinos is inferred from the presence of MET.

2343 6.3 Background processes

2344 The background processes are those that can mimic the signal signature or at least
 2345 can be reconstructed as that as a result of certain circumstances. The backgrounds
 2346 can be classified as

¹ ZZ and $\tau\tau$ decays are also include in the analysis but they are not separately reconstructed

- 2347 • irreducible backgrounds: where genuine prompt leptons are produced in on-
 2348 shell W and Z boson decays; they can be reliably estimated directly from MC
 2349 simulated events, using higher-order cross sections or data control regions for
 2350 the overall normalization.
- 2351 • reducible backgrounds: where at least one of the leptons is *non-prompt*, i.e.,
 2352 produced within a hadronic jet; genuine leptons from heavy flavor decays and
 2353 misreconstructed jets, also known as *mis-ID leptons* or *fake leptons*, are consid-
 2354 ered non-prompt leptons. These non-prompt leptons leave tracks and hits in
 2355 the detection systems as would a prompt lepton, but correlating those hits with
 2356 nearby jets could be a way of removing them. The misassignment of electron
 2357 charge in processes like $t\bar{t}$ or Drell-Yan, represent an additional source of back-
 2358 ground, but it is relevant only for the $2lss$ channel. Reducible backgrounds are
 2359 not well predicted by simulation, hence, they are estimated using data-driven
 2360 methods.

2361 The main sources of background events for tHq process are $t\bar{t}$ process and $t\bar{t} +$
 2362 $X(X = W, Z, \gamma)$ processes, the latter regarded together as $t\bar{t}V$ process. Figure 6.3
 2363 shows the signature for $t\bar{t}$ and $t\bar{t}W$ processes.

2364 The largest contribution to irreducible backgrounds comes from $t\bar{t}W$ and $t\bar{t}Z$ processes
 2365 for which the number of ($b-$)jets (($b-$)jet multiplicity) is higher than that of the sig-
 2366 nal events, while for other contributing background events, WZ , ZZ , and rare SM
 2367 processes like $W^\pm W^\pm qq$, $t\bar{t}t\bar{t}$, tZq , tZW , WWW , WWZ , WZZ , ZZZ , the ($b-$)jet
 2368 multiplicity is lower compared to that of the signal events. None of the irreducible
 2369 backgrounds present activity in the forward region of the detector.

2370 On the side of the reducible backgrounds, the largest contribution comes from the
 2371 $t\bar{t}$ events which have a very similar signature to the signal events but does no present

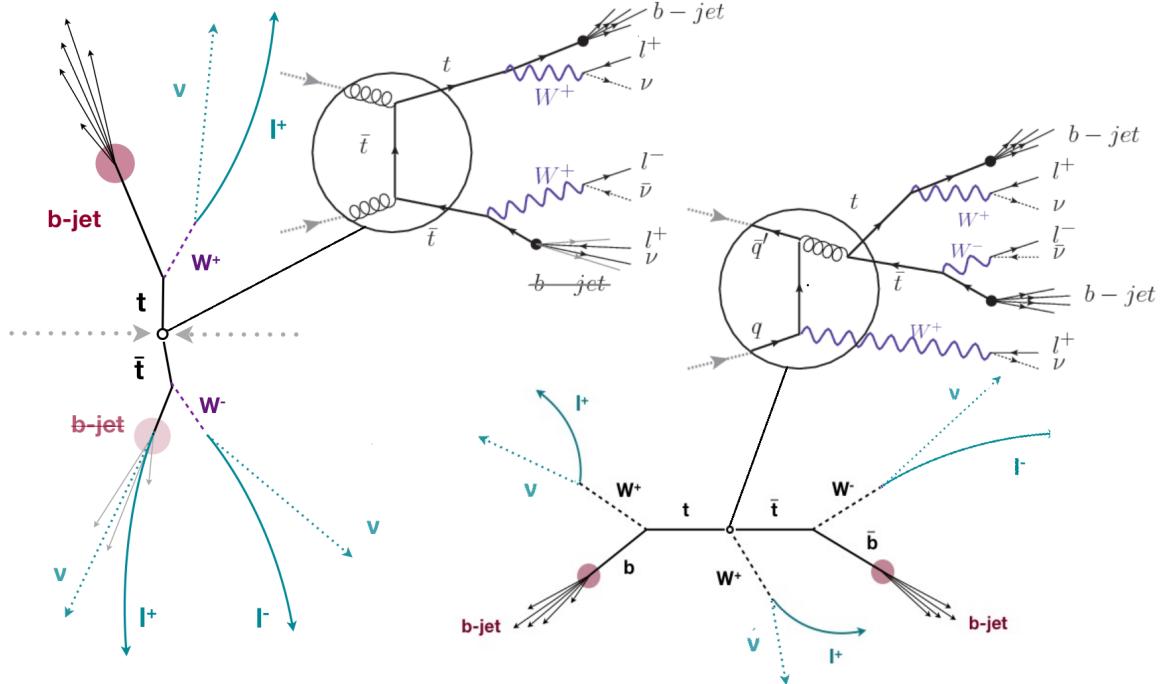


Figure 6.3: $t\bar{t}$ (left) and $t\bar{t}W$ (right) events signature as they would be seen in the detector; the Feynman diagrams including the whole evolution up to the final state are also showed. The $t\bar{t}$ process signature is very similar to that of the signal process with one fake lepton and no forward activity. The $t\bar{t}W$ process present a higher b-jet multiplicity compared to the signal process, a prompt lepton and no forward activity.

2372 activity in the forward region of the detector either; A particular feature of the $t\bar{t}$
 2373 events is their charge-symmetry, which is also a difference with respect to the signal
 2374 events.

2375 6.4 Data and MC Samples

2376 6.4.1 Full 2016 data set

2377 The data set used in this analysis was collected by the CMS experiment during 2016
 2378 at while running at $\sqrt{s} = 13\text{TeV}$ and corresponds to a total integrated luminosity
 2379 of 35.9fb^{-1} . Only periods when the CMS magnet was on were considered when

2380 selecting the data samples; that corresponds to the 23Sep2016 (Run B to G) and
 2381 **PromptReco** (Run H) versions of the datasets.

2382 Multilepton final states with either two same-sign leptons or three leptons tar-
 2383 get the case where the Higgs boson decays to a pair of W bosons, τ leptons, or Z
 2384 bosons, and where the top quark decays leptonically, hence, the **SingleElectron**,
 2385 **SingleMuon**, **DoubleEG**, **MuonEG**, **DoubleMuon** dataset (see Table A.1) compose the
 2386 full dataset. The certified luminosity sections are selected using the golden JSON file
 2387 defined by the CMS experiment [148].

2388 6.4.2 Triggers

2389 The events considered are those online-reconstructed events triggered by one, two, or
 2390 three leptons. Single-lepton triggers are included in order to boost the acceptance
 2391 of events where the p_T of the sub-leading lepton falls below the threshold of the
 2392 double-lepton triggers. The trigger efficiency is increased by including double-lepton
 2393 triggers in the $3l$ category, and single-lepton triggers in all categories; it is possible
 2394 given the logical “or” of the trigger decisions of all the individual triggers in a given
 2395 category. Table A.2 shows the lowest-threshold non-prescaled triggers present in the
 2396 High-Level Trigger (HLT) menus for both Monte-Carlo and data in 2016.

2397 Trigger efficiency scale factors

2398 Trigger efficiency describes the ability of events to pass the trigger requirements. It
 2399 is measured in simulated events using generator information given that there is no
 2400 trigger bias with the MC sample. Measuring the trigger efficiency in data requires a
 2401 more elaborated procedure; first, select a set of events collected by a trigger that is
 2402 uncorrelated with the lepton triggers such that the selected events form an unbiased

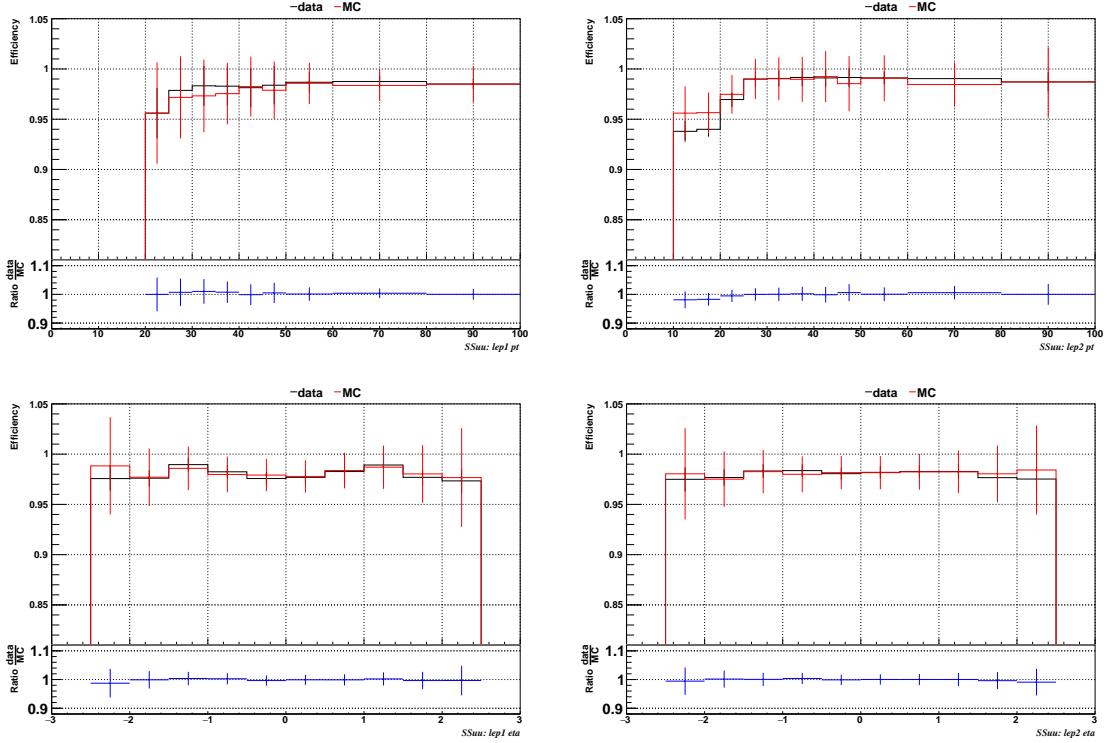


Figure 6.4: Comparison between data and MC trigger efficiencies in the same-sign $\mu\mu$ category, as a function of the p_T and η of the leading lepton (left) and the sub-leading lepton (right) [149].

sample. In this analysis, that uncorrelated trigger is a MET trigger. Second step is looking for candidate events with exactly two good leptons (exactly three good leptons for the $3l$ channel). Finally, measure the efficiency for the candidate events to pass the logical “or” of triggers being considered in a given event category as defined in Table A.2.

Comparisons between the data and MC efficiencies for each category, showed in Figures 6.4, 6.5, and 6.6, reveal that they are in good agreement; the difference is corrected by applying scale factors derived from the ratio between both efficiencies.

Applied flat scale factors in each category are shown in Table 6.1; they have been inherited from Reference [149].

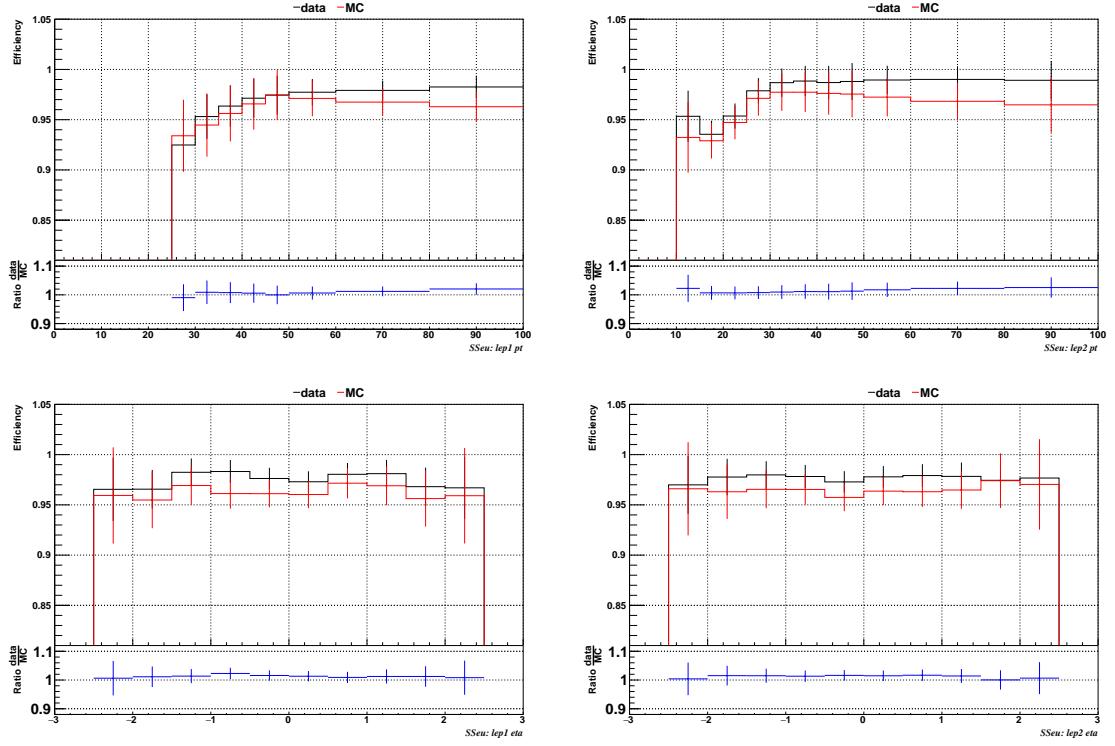


Figure 6.5: Comparison between data and MC trigger efficiencies in the same-sign $e\mu$ category as a function of the p_T and η of the leading lepton (left) and the sub-leading lepton (right) [149].

Category	Scale Factor
ee	1.01 ± 0.02
$e\mu$	1.01 ± 0.01
$\mu\mu$	1.00 ± 0.01
3l	1.00 ± 0.03

Table 6.1: Trigger efficiency scale factors and associated uncertainties, shown here rounded to the nearest percent.

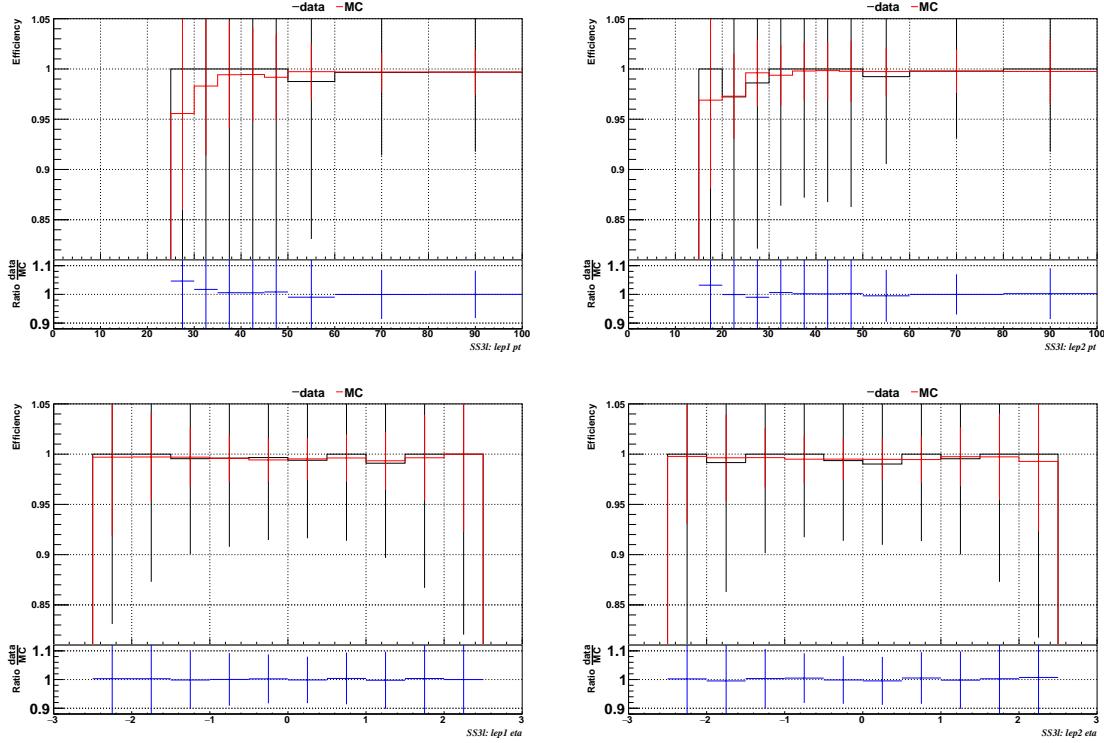


Figure 6.6: Comparison between data and MC trigger efficiencies in the $3l$ category, as a function of the p_T and η of the leading lepton (left) and the sub-leading lepton (right) [149].

2413 6.4.3 Signal modeling and MC samples

2414 Current event generators allow for adjusting the kinematics of the generated events,
 2415 based on an event-wise reweighting; in this way, several generation parameters phase
 2416 spaces can be explored according to the experimental interests. The signal samples
 2417 used in this analysis were generated in such a way that not only the case $\kappa_t = -1$, but
 2418 an extended range of κ_t and κ_V values may be investigated.

2419 tHq and tHW cross section in the κ_t - κ_V phase space are shown in Figure 6.7. As
 2420 said in section 3.1, the tHq sample was generated using the 4F scheme which provides
 2421 a better description of the additional b quark from the initial gluon splitting, while the
 2422 tHW sample was generated using the 5F scheme in order to remove its interference
 2423 with $t\bar{t}H$ at LO.

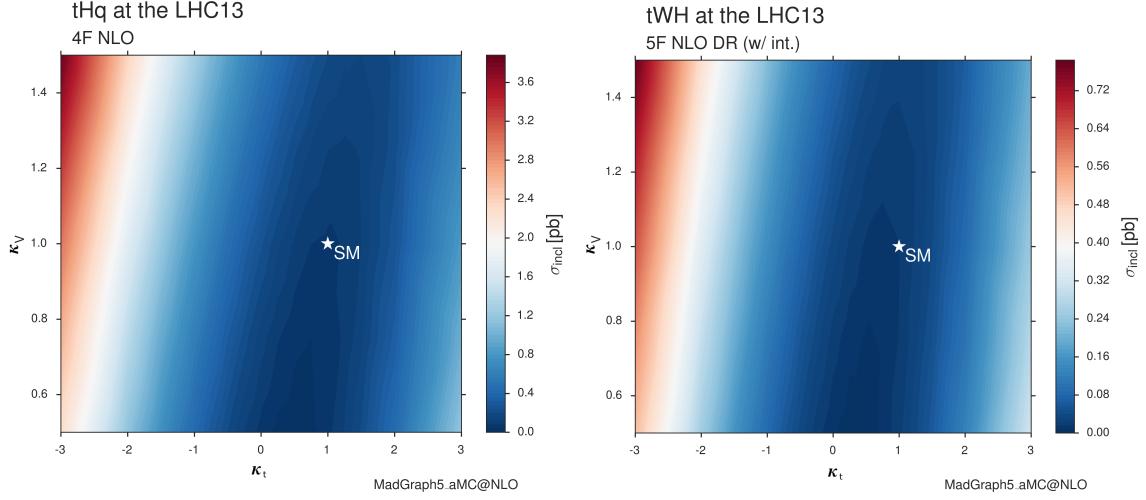


Figure 6.7: tHq and tHW cross section in the κ_t - κ_V phase space [150].

2424 MC signal samples

2425 The two signal samples, tHq and tHW , correspond to the `RunIISummer16MiniAODv2`
 2426 campaign produced with `CMSSW_80X`; they were produced with `MG5_aMC@NLO`
 2427 (version 5.2.2.3), in LO order mode at $\sqrt{s} = 13$ TeV, and are normalized to NLO cross
 2428 sections (see Table 6.2). The Higgs boson is assumed to be SM-like except for the
 2429 values of its couplings to the top quark and W boson. Each sample was generated
 2430 with a set of event weights corresponding to 51 different values of (κ_t, κ_V) couplings,
 2431 accessible in terms of LHE event weights as shown in Table A.3; however, the main
 2432 interest is the $(\kappa_t = -1, \kappa_V = 1)$ case.

Sample	σ [pb]	BF
<code>/THQ_Hincl_13TeV-madgraph-pythia8_TuneCUETP8M1/</code>	0.7927	0.324
<code>/THW_Hincl_13TeV-madgraph-pythia8_TuneCUETP8M1/</code>	0.1472	1.0
<code>/tthJetToNonbb_M125_13TeV_amcatnloFXFX_madspin_pythia8_mWCutfix/</code>	0.2151	1.0

Table 6.2: MC signal samples used in this analysis; cross section and branching fraction are also listed [150].

2433 The $t\bar{t}H$ sample was produced using `AMC@NLO` interfaced to `PYTHIA 8` for
 2434 the parton shower, and is scaled to NLO cross sections. The $t\bar{t}H$ cross section depends

2435 quadratically on κ_t ; however, in contrast to the tHq and tHW samples, the scaling
 2436 is not performed during the sample generation process but in the analysis code since
 2437 it was decided to include the $t\bar{t}H$ process as part of the signal in the course of the
 2438 analysis.

2439 **MC background samples**

2440 Several MC generators were used to generate the samples of the background processes.
 2441 The dominant background sources ($t\bar{t}$, $t\bar{t}W$, $t\bar{t}Z$) were produced using AMC@NLO
 2442 interfaced to PYTHIA8, and are scaled to NLO cross sections. Other minor back-
 2443 ground processes are simulated using POWHEG interfaced to PYTHIA, or bare
 2444 PYTHIA as stated in the sample names in Table A.4. Pileup interactions are in-
 2445 cluded in the simulation in order to reflect the observed multiplicity in data; the
 2446 simulated events are weighted according to the actual pileup in data, estimated from
 2447 the measured bunch-to-bunch instantaneous luminosity and the total inelastic cross
 2448 section, 69.2 mb. All events are finally passed through a full simulation of the CMS
 2449 detector based on GEANT4, and reconstructed using the same algorithms as used for
 2450 the data.

2451 **6.5 Object Identification**

2452 In this section, the specific definitions of the physical objects in terms of the numerical
 2453 values assigned to the reconstruction parameters are presented; thus, the provided
 2454 details summarize and complement the descriptions presented in previous chapters.
 2455 The object reconstruction and selection strategy used in this thesis is inherited from
 2456 the analyses in References [145, 149], thus, the information in this section is extracted
 2457 from those documents unless other References are stated.

2458 **6.5.1 Lepton reconstruction and identification**

2459 Two types of leptons are defined in this analysis: *signal leptons* are those coming from
 2460 W, Z and τ decays which usually are isolated from other particles; *background leptons*
 2461 are defined as leptons produced in b -jet hadron decays, light-jets misidentification,
 2462 and photon conversions.

2463 The process of reconstruction and identification of electron and muon candidates
 2464 was described in chapter3, hence, the identification variables used in order to retain
 2465 the highest possible efficiency for signal leptons while maximizing the rejection of
 2466 background leptons are listed and described in the following sections ².

2467 The identification variables include not only observables related directly to the re-
 2468 constructed leptons themselves, but also to the clustered energy deposits and charged
 2469 particles in a cone around the lepton direction (jet-related variables); an initial loose
 2470 preselection of leptons candidates is performed and then an MVA discriminator, re-
 2471 ferred to as *lepton MVA* discriminator, is used to distinguish signal leptons from
 2472 background leptons.

2473 **Muons**

2474 The Physics Objects Groups (POG) at CMS, are in charge of studying and defining
 2475 the set of selection criteria applied on the course of reconstruction and identification
 2476 of particles. These selection criteria are implemented in the CMS framework in the
 2477 form of several object identification working points according to the strength of the
 2478 requirements.

2479 The muon candidates are reconstructed by combining information from the tracker
 2480 system and the muon detection system of CMS detector and the POG defined three

² the studies performed to optimize the identification are far from the scope of this thesis, therefore, only general descriptions are provided

2481 working points for muon identification *MuonID* [153];

- 2482 • *POG Loose Muon ID* is a particle identified as a muon by the PF event re-
 2483 construction and also reconstructed either as a global-muon or as an arbitrated
 2484 tracker-muon. This identification criteria is designed to be highly efficient for
 2485 prompt muons and for muons from heavy and light quark decays; it can be com-
 2486 plemented by applying impact parameter cuts in analyses with prompt muon
 2487 signals.
- 2488 • *POG Medium Muon ID* is a Loose muon with additional track-quality and
 2489 muon-quality (spatial matching between the individual measurements in the
 2490 tracker and the muon system) requirements. This identification criteria is de-
 2491 signed to be highly efficient in the separation of the muons coming from decay
 2492 in flight of heavy quarks and muons coming from B meson decays as well as
 2493 prompt muons. An additional category *MVA Prompt ID* is defined in this iden-
 2494 tification criteria directed to discriminated muons from B mesons and prompt
 2495 muons (from W,Z and τ decays). The Medium ID provides the same fake rate as
 2496 the Tight Muon ID but a higher efficiency on prompt and B-decays muons. [154]
- 2497 • *POG Tight Muon ID* is a global muon with additional muon-quality require-
 2498 ments Tight Muon ID selects a subset of the PF muons.

2499 Only muons within the muon system acceptance $|\eta| < 2.4$ and minimum p_T of 5
 2500 GeV are considered.

2501 **Electrons**

2502 Electrons are reconstructed using information from the tracker and from the electro-
 2503 magnetic calorimeter and identified by an MVA algorithm (*MVA eID* discriminant)

2504 using the shape of the calorimetric shower variables like the shape in η and ϕ , the clus-
 2505 ter circularity, widths along η and ϕ ; track-cluster matching variables like E_{tot}/p_{in} ,
 2506 E_{Ele}/p_{out} , $\Delta\eta_{in}$, $\Delta\eta_{out}$, $\Delta\phi_{in}$, $1/E - 1/p$; and track quality variables like χ^2 of the
 2507 GSF tracks, the number of hits used by the GSF filter [155].

2508 A loose selection based on η -dependent cuts on this discriminant is used to prese-
 2509 lect electron candidates, the full shape of the discriminant is used in the lepton MVA
 2510 selection to separate signal leptons from background leptons (described in Section
 2511 6.5.1).

2512 In order to reject electrons from photon conversions, electron candidates with
 2513 missing hits in the pixel tracker layers or matched to a conversion secondary vertex
 2514 are discarded. Electrons are selected for the analysis if they have $p_T > 7$ GeV and
 2515 are located within the tracker system acceptance region ($|\eta| < 2.5$).

2516 Lepton vertexing and pile-up rejection

2517 The impact parameter in the transverse plane d_0 , impact parameter along the z -
 2518 axis d_z , and the impact parameter significance in the detector space SIP_{3D} , are
 2519 considered to perform the identification and rejection of pile-up, misreconstructed
 2520 tracks, and background leptons from b-hadron decays; pile-up and misreconstructed
 2521 track mitigation is achieved by imposing loose cuts on the impact parameter variables.
 2522 The full shape of the those variables is used in a lepton MVA classifier to achieve the
 2523 best separation between the signal and the background leptons.

2524 Lepton isolation

2525 PF is able to recognize leptons from two different sources: on one side, leptons from
 2526 the decays of heavy particles, such as W and Z bosons, which are normally isolated
 2527 in space from the hadronic activity in the event; on the other side, leptons from the

2528 decays of hadrons and jets misidentified as leptons, which are not isolated as the
 2529 former. For highly boosted systems, like the lepton and the b -jet generated in the
 2530 semileptonic decay of a boosted top, the decay products tend to be more closer and
 2531 sometimes they even overlap; thus, the PF standard definition of isolation in terms of
 2532 the separation between the lepton candidates and other PF objects in the η - ϕ plane,

$$\Delta R = \sqrt{(\eta^l - \eta^i)^2 + (\phi^l - \phi^i)^2} < 0.3 \quad (6.1)$$

2533 which considers all the neutral, charged hadrons and photons in a cone around the
 2534 leptons, is refocused to the local isolation of the leptons through the mini-isolation
 2535 I_{mini} [156] defined as the sum of particle flow candidates p_T within a cone around
 2536 the lepton, corrected for the effects of pileup and divided by the lepton p_T

$$I_{mini} = \frac{\sum_R p_T(h^\pm) - \max\left(0, \sum_R p_T(h^0) + p_T(\gamma) - \rho \mathcal{A}\left(\frac{R}{0.3}\right)^2\right)}{p_T(l)} \quad (6.2)$$

2537 where ρ is the pileup energy density, h^\pm, h^0, γ, l , represent the charged hadron, neutral
 2538 hadrons, photons, and the lepton, respectively. The radius R of the cone depends on
 2539 the p_T of the lepton according to

$$R = \frac{10\text{GeV}}{\min(\max(p_T(l), 50\text{GeV}), 200\text{GeV})}, \quad (6.3)$$

2540 The p_T dependence of the cone size allows for greater signal efficiency. Setting a
 2541 cut on I_{mini} below a given threshold ensures that the lepton is locally isolated, even
 2542 in boosted systems. The effect of pileup is mitigated using the so-called effective area
 2543 correction \mathcal{A} listed in Table 6.3.

2544 A loose cut on I_{mini} is applied to pre-select the muon and electron candidates;

$ \eta $ range	$\mathcal{A}(e)$ neutral/charged	$A(\mu)$ neutral/charged
0.0 - 0.8	0.1607 / 0.0188	0.1322 / 0.0191
0.8 - 1.3	0.1579 / 0.0188	0.1137 / 0.0170
1.3 - 2.0	0.1120 / 0.0135	0.0883 / 0.0146
2.0 - 2.2	0.1228 / 0.0135	0.0865 / 0.0111
2.2 - 2.5	0.2156 / 0.0105	0.1214 / 0.0091

Table 6.3: Effective areas, for electrons and muons used to mitigate the effect of pileup by using the so-called effective area correction.

however, the full shape is used in the lepton MVA discriminator when performing the signal lepton selection.

Jet-related variables

In order to reject misidentified leptons from b -jets, mostly coming from $t\bar{t}$ +jets, Drell-Yan+jets, and W +jets events, the vertexing and isolation described in previous sections are complemented with additional variables related to the closest reconstructed jet to the lepton, i.e., the PF jets reconstructed³ around the leptons with $\Delta R = \sqrt{(\eta^l - \eta^{jet})^2 + (\phi^l - \phi^{jet})^2} < 0.5$. The identification variables used in the MVA discriminator are the ratio p_T^l/p_T^{jet} , the CSV b-tagging discriminator value of the jet, the number of charged tracks of the jet, and the relative p_T given by

$$p_T^{rel} = \frac{(\vec{p}_{jet} - \vec{p}_l) \cdot \vec{p}_l}{||\vec{p}_{jet} - \vec{p}_l||}. \quad (6.4)$$

LeptonMVA discriminator

Electrons and muons passing the basic selection process described above are referred to as *loose leptons*. Additional discrimination between signal leptons and background leptons is crucial considering that the rate of $t\bar{t}$ production is much larger than the signal, hence, an overwhelming background from $t\bar{t}$ production. To maximally ex-

³ charged hadrons from PU vertices are not removed prior to the jet clustering.

2560 exploit the available information in each event to that end, the dedicated lepton MVA
 2561 discriminator, based on a boosted decision tree (BDT) algorithm, has been built so
 2562 that all the identification variables can be used together.

2563 The lepton MVA discriminator training is performed using simulated signal Loose
 2564 leptons from the $t\bar{t}H$ MC sample and fake leptons from the $t\bar{t}$ + jets MC sample,
 2565 separately for muons and electrons. The input variables used include vertexing, iso-
 2566 lation and jet-related variables, the p_T and η of the lepton, the electron MVA eID
 2567 discriminator and the muon segment-compatibility variables. An additional require-
 2568 ment known as *tight-charge* requirement, is imposed by comparing two independent
 2569 measurement of the charge, one from the ECAL supercluster and the other from the
 2570 tracker; thus, the consistency in the measurements of the electron charge is ensured
 2571 so that events with a wrong electron charge assignment are rejected; this variable is
 2572 particularly used in the $2lss$ channel to suppress opposite-sign events for which the
 2573 charge of one of the leptons has been mismeasured. The tight-charge requirement for
 2574 muons is represented by the requirement of a consistently well measured track trans-
 2575 verse momentum given by $\Delta p_T/p_T < 0.2$. Leptons are selected for the final analysis
 2576 if they pass a given threshold of the BDT output, and are referred to as *tight leptons*
 2577 in the following.

2578 The validation of the lepton MVA algorithm and the lepton identification variables
 2579 is performed using data in various control regions; the details about that validation
 2580 are not discussed here but can be found in Reference [149].

2581 Selection definitions

2582 Electron and muon object identification is defined in three different sets of selections
 2583 criteria; the *Loose*, *Fakeable Object*, and *Tight* selection. These three levels of selection
 2584 are designed to serve for event level vetoes, the fake rate estimation application region

(see Section 6.7.2), and the final signal selection, respectively. The p_T of fakeable objects is defined as $0.85 \times p_T(\text{jet})$, where the jet is the one associated to the lepton object. This mitigates the dependence of the fake rate on the momentum of the fakeable object and thereby improves the precision of the method.

Tables 6.4 and 6.5 list the full criteria for the different selections of muons and electrons.

Cut	Loose	Fakeable object	Tight
$ \eta < 2.4$	✓	✓	✓
p_T	$> 5\text{GeV}$	$> 15\text{GeV}$	$> 15\text{GeV}$
$ d_{xy} < 0.05$ (cm)	✓	✓	✓
$ d_z < 0.1$ (cm)	✓	✓	✓
$\text{SIP}_{3D} < 8$	✓	✓	✓
$I_{\text{mini}} < 0.4$	✓	✓	✓
is Loose Muon	✓	✓	✓
jet CSV	–	< 0.8484	< 0.8484
is Medium Muon	–	–	✓
tight-charge	–	–	✓
lepMVA > 0.90	–	–	✓

Table 6.4: Requirements on each of the three muon selections. In the cases where the cut values change between the selections, those values are listed in the table. Otherwise, whether the cut is applied is indicated.

In addition to the previously defined requirements for jets, they are required to be separated from any lepton candidates passing the fakeable object selections by $\Delta R > 0.4$.

6.5.2 Lepton selection efficiency

Efficiencies of reconstruction and selecting loose leptons are measured both for muons and electrons using a tag and probe method on both data and MC, using $Z \rightarrow \ell^+ \ell^-$ [157]. The scale factors are derived from the ratio of efficiencies $\varepsilon_i(p_T, \eta)$ measured

Cut	Loose	Fakeable Object	Tight
$ \eta < 2.5$	✓	✓	✓
p_T	$> 7\text{GeV}$	$> 15\text{GeV}$	$> 15\text{GeV}$
$ d_{xy} < 0.05 \text{ (cm)}$	✓	✓	✓
$ d_z < 0.1 \text{ (cm)}$	✓	✓	✓
$\text{SIP}_{3D} < 8$	✓	✓	✓
$I_{\text{mini}} < 0.4$	✓	✓	✓
MVA eID $> (0.0, 0.0, 0.7)$	✓	✓	✓
$\sigma_{in\eta} < (0.011, 0.011, 0.030)$	—	✓	✓
$\text{H/E} < (0.10, 0.10, 0.07)$	—	✓	✓
$\Delta\eta_{in} < (0.01, 0.01, 0.008)$	—	✓	✓
$\Delta\phi_{in} < (0.04, 0.04, 0.07)$	—	✓	✓
$-0.05 < 1/E - 1/p < (0.010, 0.010, 0.005)$	—	✓	✓
p_T^{ratio}	—	$> 0.5^\dagger / -$	—
jet CSV	—	$< 0.3^\dagger / < 0.8484$	< 0.8484
tight-charge	—	—	✓
conversion rejection	—	—	✓
Number of missing hits	< 2	$== 0$	$== 0$
lepton MVA > 0.90	—	—	✓

Table 6.5: Criteria for each of the three electron selections. In cases where the cut values change between selections, those values are listed in the table. Otherwise, whether the cut is applied is indicated. In some cases, the cut values change for different η ranges. These ranges are $0 < |\eta| < 0.8$, $0.8 < |\eta| < 1.479$, and $1.479 < |\eta| < 2.5$ and the respective cut values are given in the form (value₁, value₂, value₃). For the two p_T^{ratio} and CSV rows, the cuts marked with a \dagger are applied to leptons that fail the lepton MVA cut, while the loose cut value is applied to those that pass the lepton MVA cut.

2598 for a given lepton in data/MC, according to

$$\rho(p_T, \eta) = \frac{\varepsilon_{\text{data}}(p_T, \eta)}{\varepsilon_{\text{MC}}(p_T, \eta)}. \quad (6.5)$$

2599 The scale factor for each event is used to correct the weight of the event in the
 2600 full sample; therefore, the full simulation correction is given by the product of all
 2601 the individual scale factors. The scale factors used in this thesis are inherited from
 2602 Reference [149] which in turns inherited them from leptonic SUSY analyses using
 2603 equivalent lepton selections.

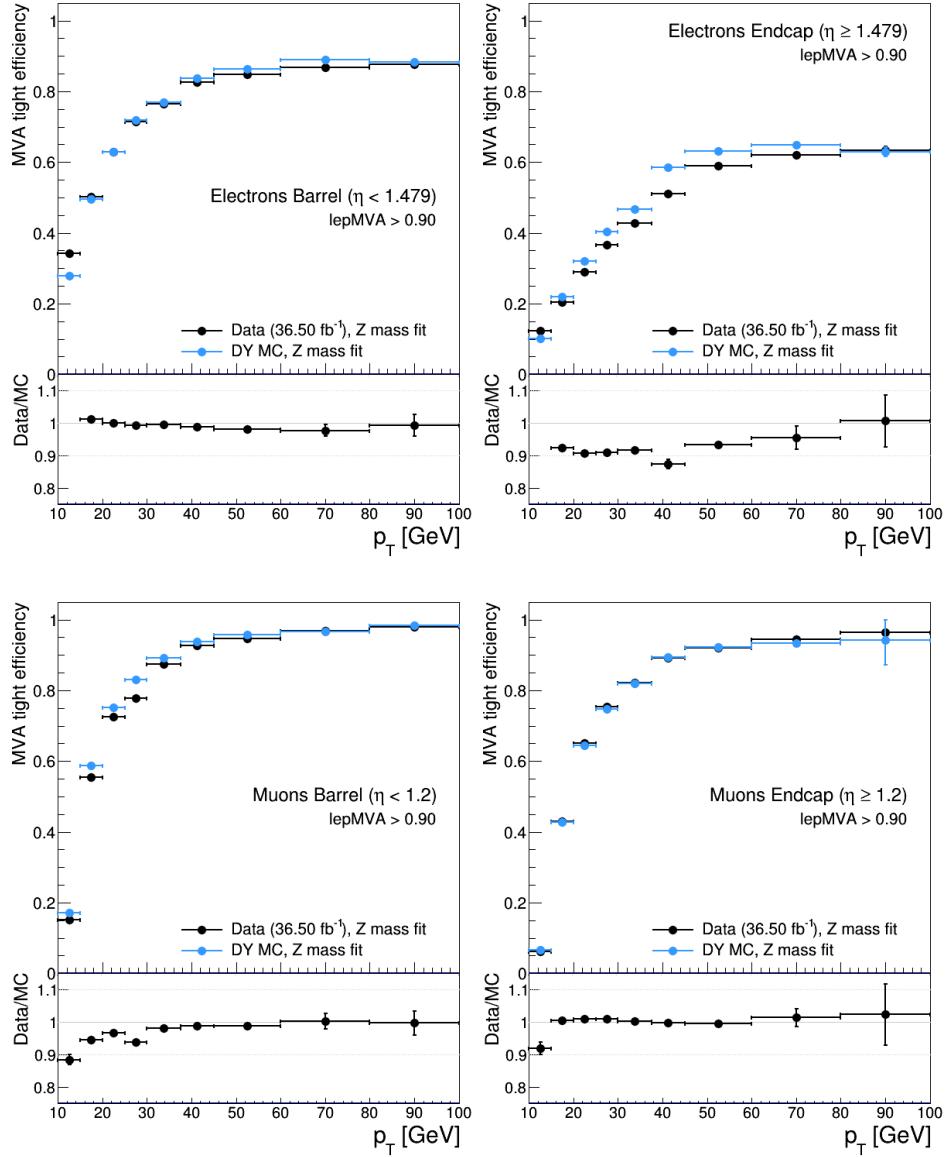


Figure 6.8: Tight vs loose selection efficiencies for electrons (top), and muons (bottom), for the $2lss$ definition, i.e., including the tight-charge requirement.

2604 The efficiency of applying the tight selection as defined in Tables 6.4 and 6.5, on the
 2605 loose leptons are determined by using a tag and probe method on a sample of Drell-
 2606 Yan enriched events. Figures 6.8 and 6.9 show the efficiencies for the $2lss$ channel and
 2607 $3l$ channel respectively. Efficiencies in the $2lss$ channel have been produced including
 2608 the tight-charge requirement, while for the $3l$ channel it is not included. Number

of passed and failed probes are determined from a fit to the invariant mass of the dilepton system. Simulation is corrected using these scale factors; note that they depends on η and p_T .

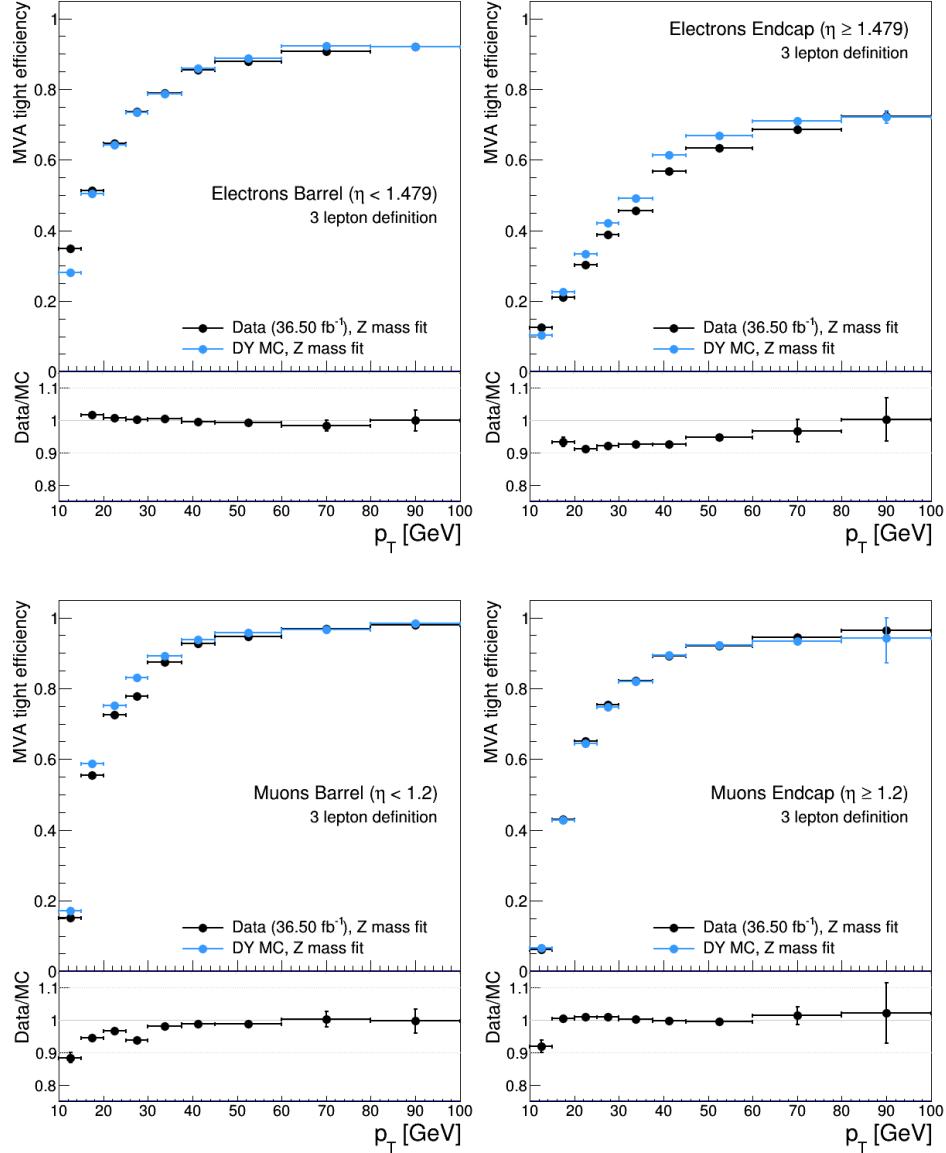


Figure 6.9: Tight vs loose selection efficiencies for electrons (top), and muons (bottom), for the $3l$ channel not including the tight-charge requirement.

²⁶¹² **6.5.3 Jets and b -jet tagging**

²⁶¹³ In this analysis, jets are reconstructed by clustering PF candidates using the anti- k_t
²⁶¹⁴ algorithm with parameter distance $\Delta R = 0.4$; those charged hadrons that are not
²⁶¹⁵ consistent with the selected primary vertex are discarded from the clustering. The
²⁶¹⁶ jet energy is then corrected for the varying response of the detector as a function
²⁶¹⁷ of transverse momentum p_T and pseudorapidity η . Jets are selected for use in the
²⁶¹⁸ analysis only if they have $p_T > 25$ GeV and are separated from any selected leptons
²⁶¹⁹ by $\Delta R > 0.4$.

²⁶²⁰ Jets coming from the primary vertex and jets coming from pile-up vertices are
²⁶²¹ distinguished using a MVA discriminator based on the differences in the jet shapes,
²⁶²² in the relative multiplicity of charged and neutral components, and in the different
²⁶²³ fraction of transverse momentum which is carried by the hardest components. Jet
²⁶²⁴ tracks are also required to be compatible with the primary vertex.

²⁶²⁵ Jets originated from the hadronization of a b quark are selected using a MVA
²⁶²⁶ likelihood discriminant which uses track-based lifetime information and reconstructed
²⁶²⁷ secondary vertices (CSV algorithm). Only jets within the CMS tracker acceptance
²⁶²⁸ ($\eta < 2.4$) are identified with this tool. Data samples are used to measure the efficiency
²⁶²⁹ of the b -jet tagging and the probability to misidentify jets from light quarks or gluons;
²⁶³⁰ in both cases the measurements are parametrized as a function of the jet p_T and η
²⁶³¹ and later used to correct differences between the data and MC simulation in the b
²⁶³² tagging performance, by applying per-jet weights to the simulation, dependent on
²⁶³³ the jet p_T , η , b tagging discriminator, and flavor (from simulation truth) [151]. The
²⁶³⁴ per-event weight is taken as the product of the per-jet weights, including those of the
²⁶³⁵ jets associated to the leptons. The weights are derived on $t\bar{t}$ and Z+jets events.

²⁶³⁶ Two working points are defined, based on the CSV algorithm output: ‘*loose*’ work-

ing point ($\text{CSV} > 0.46$) with a b signal tagging efficiency of about 83% and a mistagging rate of about 8%; and *medium* working point ($\text{CSV} > 0.80$) with b -tagging efficiency of about 69% and mistagging rate of order 1% [152]. Tagging of jets from charm quarks have efficiencies of about 40% and 18% for loose and medium working points respectively. Separate scale factors are applied to jets originating from bottom/charm quarks and from light quarks in simulated events to match the tagging efficiencies measured in the data.

6.5.4 Missing Energy MET

As stated in Section 3.4.1.1, the MET vector is calculated as the negative of the vector sum of transverse momenta of all PF candidates in the event and its magnitude is referred to as E_T^{miss} . Due to pile-up interactions, the performance in determining MET is degraded; in order to correct for that, the energy from the selected jets and leptons that compose the event is assigned to the variable H_T^{miss} . It is calculated in the same way as E_T^{miss} and although it has worse resolution than E_T^{miss} , it is more robust in the sense that it does not rely on the soft part of the event. The event selection uses a linear discriminator based on the two variables given by

$$E_T^{\text{miss}} \text{LD} = 0.00397 * E_T^{\text{miss}} + 0.00265 * H_T^{\text{miss}} \quad (6.6)$$

taking advantage of the fact that the correlation between E_T^{miss} and H_T^{miss} is less for events with instrumental missing energy than for events with real missing energy. The working point $E_T^{\text{miss}} \text{LD} > 0.2$ was chosen to ensure a good signal efficiency while keeping a good background rejection.

2657 6.6 Event selection

2658 Events are selected considering the features of the signal process and the decay sig-
 2659 nature as described in Section 6.2. At the trigger level, events are selected to contain
 2660 either one, two, or three leptons with minimal p_T thresholds:

- 2661 • single-lepton trigger → 24 GeV for muons and at 27 GeV for electrons
- 2662 • double-lepton triggers → leading and sub-leading leptons: 17 and 8 GeV for
 2663 muons and 23 and 12 GeV for electrons.
- 2664 • three-lepton triggers → threshold on the third hardest lepton in the event: 5
 2665 and 9 GeV for muons and electrons, respectively.

2666 The offline event selection level targets the specific topology of the tHq signal
 2667 with $H \rightarrow WW$ and $t \rightarrow Wb \rightarrow l\nu b$; therefore, the resulting state is composed of three
 2668 W bosons, one b quark, and a light spectator quark at high rapidity. The selection
 2669 criteria for the two channels exploited in this analysis are summarized in Table 6.6.
 2670 This selection includes contributions from $H \rightarrow \tau\tau$ and $H \rightarrow ZZ$ as well.

Same-sign $\ell\ell$ channel	$\ell\ell\ell$ channel
have fired one of the corresponding trigger paths	
No loose leptons with $m_{\ell\ell} < 12\text{GeV}$	
One or more b tagged jets (CSV medium) $ \eta < 2.4$	
One or more non-tagged jets: central → $p_T > 25\text{ GeV}$, $\eta < 2.4$ forward → $p_T > 40\text{ GeV}$, $\eta > 2.4$	
	$E_T^{\text{miss}} \text{LD} > 0.2$
Exactly two tight same-sign leptons	Exactly three tight leptons
Lepton $p_T > 25/15\text{GeV}$	Lepton $p_T > 25/15/15\text{GeV}$
Electrons are triple-charge consistent.	No OSSF lepton pair with $ m_{\ell\ell} - m_Z < 15\text{GeV}$
Muon p_T resolution: $\Delta p_T/p_T < 0.2$.	
No ee pair with $ m_{ee} - m_Z < 10\text{GeV}$	

Table 6.6: Summary of event pre-selection.

2671 The dominant background contribution is expected to arise from top quark pro-
 2672 duction processes, either $t\bar{t}$ pair production or in $t\bar{t}$ associated production with a
 2673 W/Z . Processes with production of single top quarks also contribute, mainly in the
 2674 associated production with a Z boson (tZq) or when produced with both a W and a
 2675 Z boson (tZW). Background contamination from diboson processes is strongly sup-
 2676 pressed by imposing the Z -veto, vetoing additional leptons and requiring b -jets in the
 2677 event.

	3ℓ	$\mu^\pm\mu^\pm$	$e\mu$	ee
$t\bar{t}W$	22.50 ± 0.35	68.03 ± 0.61	97.00 ± 0.71	29.63 ± 0.39
$t\bar{t}Z/\gamma^*$	32.80 ± 1.79	25.89 ± 1.12	64.82 ± 2.42	28.74 ± 1.70
WZ	8.22 ± 0.86	15.07 ± 1.19	26.25 ± 1.57	9.31 ± 0.93
ZZ	1.62 ± 0.33	1.16 ± 0.29	2.86 ± 0.45	1.09 ± 0.27
$W^\pm W^\pm qq$	–	3.96 ± 0.52	6.99 ± 0.69	2.19 ± 0.37
$W^\pm W^\pm(\text{DPS})$	–	2.48 ± 0.42	4.17 ± 0.54	0.81 ± 0.24
VVV	0.42 ± 0.16	2.99 ± 0.34	4.85 ± 0.43	1.19 ± 0.21
ttt	1.84 ± 0.44	2.32 ± 0.45	4.06 ± 0.57	0.89 ± 0.31
tZq	3.92 ± 1.48	5.77 ± 2.24	10.73 ± 3.03	7.56 ± 1.72
tZW	1.70 ± 0.12	2.13 ± 0.13	3.91 ± 0.18	1.13 ± 0.10
γ conversions	7.43 ± 1.94	–	23.81 ± 6.04	9.87 ± 4.17
Non-prompt	25.61 ± 1.26	80.94 ± 2.02	135.34 ± 2.83	47.72 ± 1.79
Charge mis-ID	–	–	58.50 ± 0.31	44.52 ± 0.31
All backgrounds	106.05 ± 3.45	210.74 ± 3.61	443.30 ± 8.01	184.65 ± 5.29
tHq ($\kappa_t = -1.0$)	7.48 ± 0.14	18.48 ± 0.22	27.41 ± 0.27	8.47 ± 0.15
tHW ($\kappa_V = -1.0$)	7.38 ± 0.16	7.72 ± 0.17	11.23 ± 0.20	3.66 ± 0.11
$t\bar{t}H$	18.29 ± 0.41	24.18 ± 0.48	35.21 ± 0.58	11.07 ± 0.32
Data ($35.9 fb^{-1}$)	127	280	525	208

Table 6.7: Expected and observed yields for $35.9 fb^{-1}$ after the selection in all final states. Uncertainties are statistical only.

2678 In the $2lss$ channel, events with additional tight leptons are vetoed as well as those
 2679 for which a loose lepton pair has an invariant mass below 12 GeV. A threshold in p_T of
 2680 the leading and sub-leading leptons is also required. Events where the two electrons
 2681 have invariant mass within 10 GeV of the Z boson mass (Z -veto) are discarded in
 2682 order to reject events from DY+jets production with charge misidentified electrons.

2683 In addition, contribution from the associated production of two W bosons of equal
 2684 charge and two light jets $W^\pm W^\pm qq$ and from same-sign W boson pairs can also be
 2685 produced in double parton scattering (DPS) processes, where each of the colliding
 2686 protons gives two partons, resulting in two hard interactions.

2687 In the $3l$ lepton channel, leptons are required to have respectively $p_T > 25\text{GeV}$, $>$
 2688 15 GeV , and $> 15\text{ GeV}$. Events with an opposite-sign same-flavor lepton combination
 2689 (OSSF) with invariant mass within 15 GeV of the Z boson mass are discarded in order
 2690 to reject events from $WZ + \text{jets}$ production.

2691 The selection criteria in Table 6.6 represent a relatively loose selection that allows
 2692 to maintain a large signal efficiency while suppressing the main backgrounds; thus
 2693 that selection is called *pre-selection*. The events obtained from the pre-selection are
 2694 then used to extract the signal contribution in a second analysis step, using BDT dis-
 2695 criminators against the main backgrounds of $t\bar{t}W/t\bar{t}Z$ and non-prompt leptons from
 2696 $t\bar{t}$. The shape of the discriminator variables is then fit to the observed data distribu-
 2697 tion to estimate the signal and background yields, simultaneously for all channels.

2698 The expected and observed event yields of the pre-selection are shown in Table
 2699 6.7. For the tH and $t\bar{t}H$ processes, the largest contribution comes from Higgs decays
 2700 to WW (about 75%), followed by $\tau\tau$ (about 20%) and ZZ (about 5%). Other Higgs
 2701 production modes contribute negligible event yields (< 5% of the $tH+t\bar{t}H$ yield) as
 2702 shown in Table 6.8.

2703 6.7 Background modeling and predictions

2704 Irreducible backgrounds are reliably estimated from MC simulated events; therefore,
 2705 in this analysis all backgrounds involving prompt leptons are estimated in this way.
 2706 Reducible backgrounds, like non-prompt lepton backgrounds, are not well predicted

	3ℓ	$\mu^\pm \mu^\pm$		
tHq (Inclusive)	6.57	100.0%	17.38	100.0%
$tHq(H \rightarrow WW)$	4.84	73.9%	13.33	76.9%
$tHq(H \rightarrow \tau\tau)$	1.04	15.9%	3.62	20.6%
$tHq(H \rightarrow ZZ)$	0.48	7.2%	0.37	2.2%
$tHq(H \rightarrow \mu\mu)$	0.21	3.0%	0.04	0.2%
$tHq(H \rightarrow \gamma\gamma)$	< 0.01	0.1%	0.02	0.1%
$tHq(H \rightarrow bb)$	< 0.01	< 0.1%	0.01	< 0.1%
tHW (Inclusive)	7.32	100.0%	7.62	100.0%
$tHW(H \rightarrow WW)$	5.50	76.9%	5.60	74.1%
$tHW(H \rightarrow \tau\tau)$	1.40	20.6%	1.81	23.1%
$tHW(H \rightarrow ZZ)$	0.31	2.2%	0.21	2.7%
$tHW(H \rightarrow \mu\mu)$	0.12	0.2%	0.01	0.1%
$tHW(H \rightarrow \gamma\gamma)$	< 0.01	< 0.1%	< 0.01	< 0.1%
$tHW(H \rightarrow bb)$	< 0.01	< 0.1%	< 0.01	< 0.1%

Table 6.8: Signal yields split by decay channels of the Higgs boson. Forward jet p_T cut at 25 GeV.

2707 by simulation, hence, they are estimated using data-driven methods.

2708 6.7.1 $t\bar{t}W$ and diboson backgrounds

2709 Backgrounds from $t\bar{t}W$ and $t\bar{t}Z$ processes are estimated using simulated events, cor-
 2710 rected for data/MC differences and inefficiencies (trigger and lepton selection) in the
 2711 same way as signal events. Their production cross sections are calculated at NLO
 2712 order of QCD and EWK, considering theoretical uncertainties from unknown higher
 2713 orders of 12% for $t\bar{t}W$ and 10% for $t\bar{t}Z$. Additional uncertainties arise from the knowl-
 2714 edge of PDFs and α_s of about 4% each for $t\bar{t}W$ and $t\bar{t}Z$.

2715 The diboson contribution is also estimated from simulated events; however, the
 2716 overall normalization of this process is obtained from a dedicated control region.

2717 The motivation behind that strategy is that even though the measured inclusive
 2718 cross section for diboson processes (WZ, ZZ) is in good agreement with the NLO
 2719 calculations [149], that agreement is perturbed when leptonic Z decays and hadronic

2720 jets in the final state are required; those requirements are precisely the ones that
 2721 make the diboson production a background for the tHq signal. Thus, by using a
 2722 dedicated control region dominated by WZ production⁴, the overall normalization is
 2723 constrained.

2724 The control region is defined by the presence of at least three leptons, of which
 2725 one opposite-sign pair must be compatible with a Z boson decay, i.e., invert the Z-
 2726 veto which makes the control region orthogonal to signal region; the b-jet tagging
 2727 requirements is also inverted with respect to the signal region, i.e., require two not
 2728 b -jets. A scale factor is extracted from the predicted distribution of WZ events in the
 2729 control region, and the observed data, while keeping other processes fixed; this factor
 2730 is used to scale the diboson prediction in the signal selection region. More details
 2731 about the procedure used can be found in Reference [149] from where the scale factor
 2732 is taken.

2733 In order to test the usability of the diboson background scale factor in this analysis,
 2734 a Z-enriched control region⁵ was defined by inverting the Z-veto and requiring exactly
 2735 three tight leptons with $p_T > 25/15/15$ GeV, one or more jets passing the CSVv2 loose
 2736 working point and less than four central jets. Figure 6.10 shows the distribution of
 2737 three variables in the diboson control region; the good agreement between MC and
 2738 data motivates the adoption of the diboson background scale factor.

2739 Most of the diboson events passing the signal selection contain jets from light
 2740 quarks and gluons that are incorrectly tagged as b -jets; it makes the estimate mainly
 2741 sensitive to the experimental uncertainty in the mis-tag rate rather than the theore-
 2742 tical uncertainty in the jet flavor composition. The overall uncertainty assigned to
 2743 the diboson prediction is estimated from the statistical uncertainty due to the limited

⁴ ZZ background is strongly reduced by the cut on MET.

⁵ This control region is different to the one used to find the scale factor.

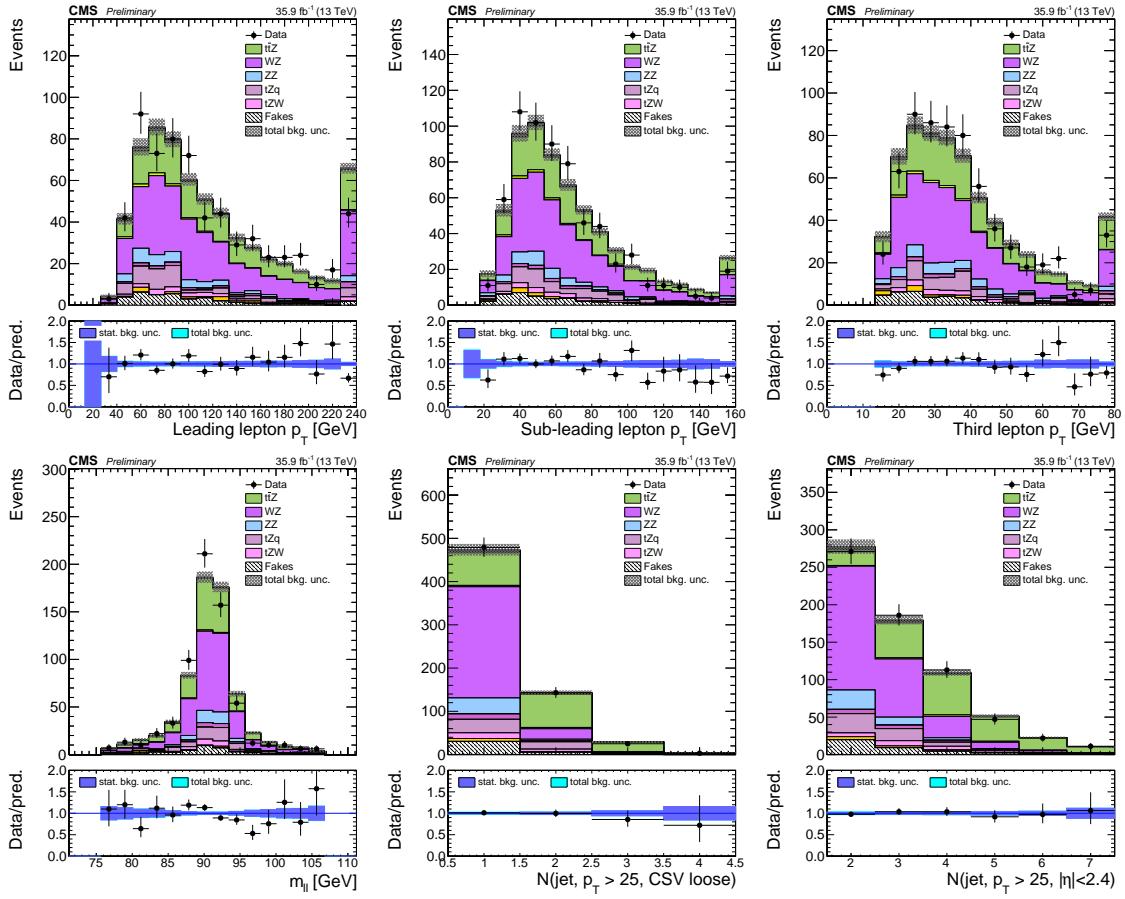


Figure 6.10: Kinematic distributions in the diboson control region.

sample size in the control region (30%), the residual background in the control region (20%), the uncertainties on the b -tagging rate (10–40%), and from the knowledge of PDFs and the theoretical uncertainties of the extrapolation (up to 10%).

6.7.2 Non-prompt and charge mis-ID backgrounds

The non-prompt lepton background contribution to the final selection is estimated using the fake factor method. The main idea of the method is to define a control region of events enriched in the background to estimate and determine a factor that relates (extrapolates) these events to those in the signal region. The method is data-driven in the sense that the control sample is selected from data, and the extrapolation

2753 factor is measured from data.

2754 In the signal region of this analysis, non-prompt leptons are predominantly pro-
 2755 duced in $t\bar{t}$ events, with a much smaller contribution, from Drell-Yan production;
 2756 therefore, the control region also known as *application region*, is defined by modifying
 2757 the event selection criteria in such a way that most of the events after selection are
 2758 $t\bar{t}$ events and thus the misidentification rate is increased. The application regions
 2759 for electrons and muons are defined by the fakeable object definitions in Tables 6.4
 2760 and 6.5. Since the fakeable definition is a loosened version of the tight definition, in
 2761 the context of fake rates the fakeable definition it becomes the loose selection.

2762 The number of events that pass both the loose and tight selections, divided by
 2763 the number of events that pass the loose selection but fail the tight one corresponds
 2764 to the *fake factor/fake rate*. The measurement of the fake factor is made using two
 2765 background dominated data samples collected with dedicated triggers (subtracting
 2766 the residual prompt lepton contribution using MC) as a function of p_T and $|\eta|$ and
 2767 separately for muons and electrons:

- 2768 • A sample dominated by QCD multijet events, collected using single lepton trig-
 2769 gers at relatively high p_T thresholds is used to extract ratios for lepton candi-
 2770 dates with p_T above 30 GeV.
- 2771 • A sample dominated by Z + jets events, where the two high p_T leptons resulting
 2772 from the Z decay are used to trigger the events without biasing the p_T spectrum
 2773 of a third lepton at low transverse momentum. It is used to determine the ratios
 2774 for low p_T leptons.

2775 Processes like W + jets, Z + jets, WZ and ZZ produce prompt leptons that
 2776 contaminate the samples; thus, they are suppressed by vetoing additional leptons in

2777 the selection, and the residual contamination is then subtracted using the transverse
 2778 mass as a discriminating variable.

2779 The extrapolation from the application region to the signal region is performed
 2780 by weighting the events in the application region using the fake factor according to
 2781 the following rules:

- 2782 • events with one lepton failing the tight criteria are weighted with the factor
 2783 $\frac{f}{(1-f)}$ for the estimate to the signal region.
- 2784 • events with two leptons (i,j) failing the tight criteria are weighted with the factor
 2785 $-\frac{f_i f_j}{(1-f_i)(1-f_j)}$ for the estimate to the signal region.
- 2786 • events with three leptons (i,j,k) failing the tight criteria are weighted with the
 2787 factor $\frac{f_i f_j f_k}{(1-f_i)(1-f_j)(1-f_k)}$ for the estimate to the signal region.

2788 The resulting prediction of the event yield in the signal selection carries an uncer-
 2789 tainty of 30-50% which is composed of the statistical uncertainty in the measurement
 2790 of the fake rates due to prescaling of lepton triggers, the uncertainty in the subtraction
 2791 of residual prompt leptons from the control region, and from testing the closure of the
 2792 method in simulated background events, hence, it is one of the dominant limitations
 2793 on the performance of multilepton analyses in general and this analysis in particular.

2794 In the $2lss$ channel case, additional background arises when the charge of a lepton
 2795 in events with an originally opposite-sign pair is misidentified; usually this happens
 2796 due to strongly asymmetric conversions of hard bremsstrahlung photons emitted from
 2797 the initial lepton, therefore, it is more likely to happen for electrons than for muons.

2798 Finally, backgrounds from electron charge mis-identification (muon charge mis-id.
 2799 is negligible) are estimated from the yield of opposite-sign event in the signal region
 2800 using a measured charge mis-identification probability. The mis-id. probability is

measured in same-sign and opposite-sign Drell–Yan events, in several bins of p_T and η . As for non-prompt leptons, the contribution from charge mis-identified electrons in our signal selection is predominantly from $t\bar{t}$ and Drell–Yan events. The systematic uncertainty of the normalization of the charge mis-id. estimate is evaluated at about 30%, stemming from a slight disagreement of the mis-id. probability between data and simulation. As it only affects the $e\mu$ channel, however, the impact of this background on the final sensitivity is very limited.

Figures 6.11 show the distributions of some relevant kinematic variables, normalized to the cross section of the respective processes and to the integrated luminosity.

A significant fraction of selected data events (about 50% in the dilepton channels, and about 80% in the trilepton channel) also passes the selection used in the dedicated search for ttH in multilepton channels [17]. The expected and observed event yields of this selection are shown in Tab. 2. For the tH and $t\bar{t}H$ processes, the largest contribution comes from Higgs decays to WW (about 75%), followed by $t\bar{t}H$ (about 20%) and ZZ (about 5%). Other Higgs production modes contribute negligible event yields (< 5% of the tH +ttH yield).

Multivariate techniques are used to discriminate the signal from the dominant backgrounds. The analysis yields a 95% confidence level (C.L.) upper limit on the combined tH + ttH production cross section times branching ratio of 0.64 pb, with an expected limit of 0.32 pb, for a scenario with $k_t = \sqrt{1.0}$ and $k_V = 1.0$. Values of k_t outside the range of $\sqrt{1.25}$ to $\sqrt{1.60}$ are excluded at 95% C.L., assuming $k_V = 1.0$.

Dont forget to mention previous constrains to ct check Reference ?? and References <https://link.springer.com/content/pdf/10.1007%2FJHEP01%282013%29088.pdf> (paragraph after eq 2)

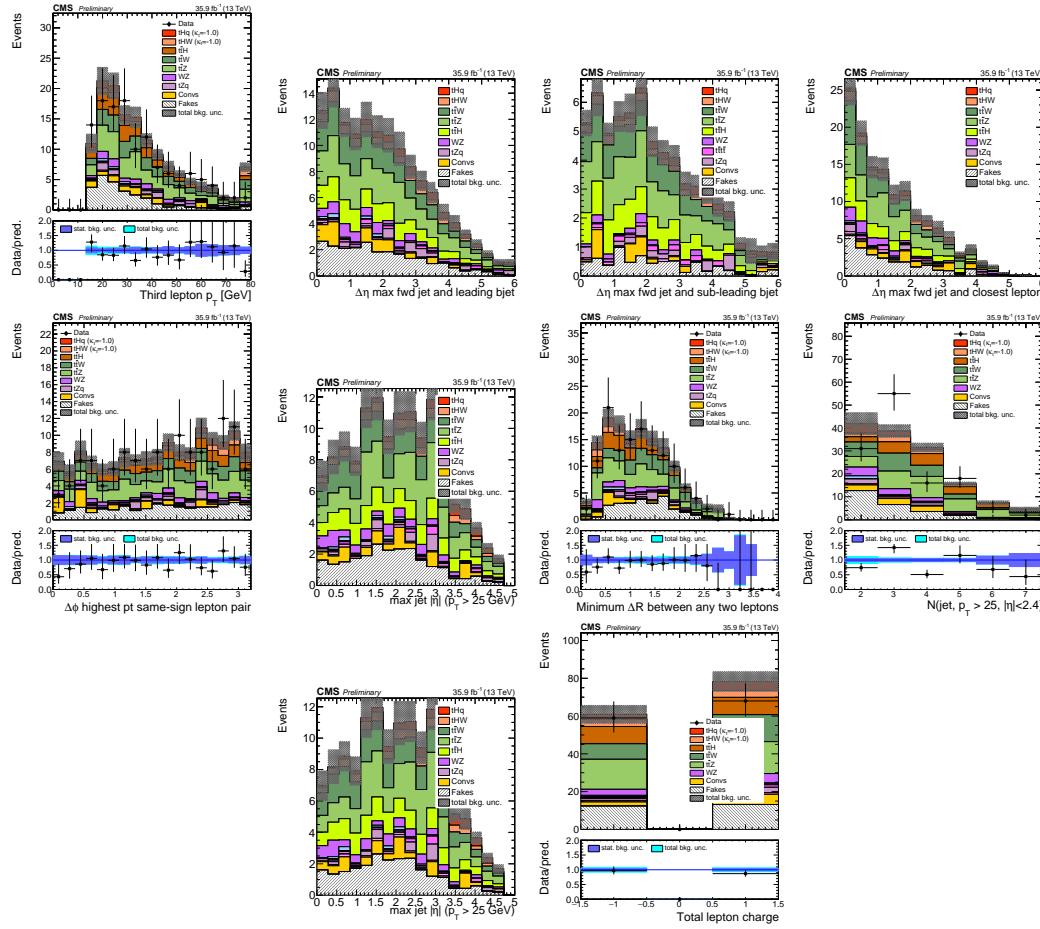


Figure 6.11: Distributions of input variables to the BDT for signal discrimination, three lepton channel, normalized to their cross section and to 35.9 fb^{-1} .

2827 6.8 Signal discrimination

2828 The production cross section for the signal processes tHq , tHW , and $t\bar{t}H$ is only
 2829 about 600 fb (the enhancement provided by inverted couplings, $\kappa_t = -1$ almost double
 2830 it), resulting in a small signal to background ratio even for a tight selection. A
 2831 multivariate method is hence employed to train discriminators to separate tH signal
 2832 events from the dominant background events.

2833 **6.8.1 MVA classifiers evaluation**

2834 Several MVA classifier algorithms were evaluated in order to determine the most
 2835 appropriate method for this analysis⁶. The comparison is based on the performance
 2836 of the classifiers, encoded in the plot of the background rejection as a function of the
 2837 signal efficiency (ROC curve). The top row of Figure 6.12 shows the ROC curves
 2838 for the several methods evaluated; two separated training were performed in the $3l$
 2839 channel: against $t\bar{t}$ (right) and $t\bar{t}V$ (left) processes.

2840 In both cases, the gradient boosted decision tree *BDTG* (*BDTA_GRAD* in the
 2841 plot) classifier offers the best results, followed by the adaptive BDT classifier (*BDTA*);
 2842 the several Fisher classifiers tested, which differ in their parameters and/or boosting
 2843 method, they offer similar performance among them, while the k-Nearest Neighbour
 2844 (kNN) classifier performance is below the rest of the classifiers. The corresponding
 2845 ROC curves and in the $2lss$ channel for trainings against $t\bar{t}V$ (left) and $t\bar{t}$ (right)
 2846 processes are shown in the bottom row of Figure 6.12; the BDTG performance is
 2847 similar to that in the $3l$ channel.

2848 **6.8.2 Discriminating variables**

2849 The classifier chosen to separate the tHq signal from the main backgrounds is the
 2850 *BDTG* classifier, trained on simulated signal and background events. The samples
 2851 used in the training are the tHq sample in Table 6.2, the samples in the third section
 2852 of table A.4 and the samples marked with an * in the same table.

2853 As explained in Section 5.1.1, a set of discriminating variables are given as input to
 2854 the *BDTG* which combines the individual discrimination power of each input variable

⁶ The choice of the tested algorithms was based on the recommendations provided by the official TMVA user guide, the experience from previous analyses and considering the expertise of the members of the tHq and $t\bar{t}H$ analyses groups. Only the BDT classifier is described in this thesis and a detailed description of all available methods can be found in Reference [127]

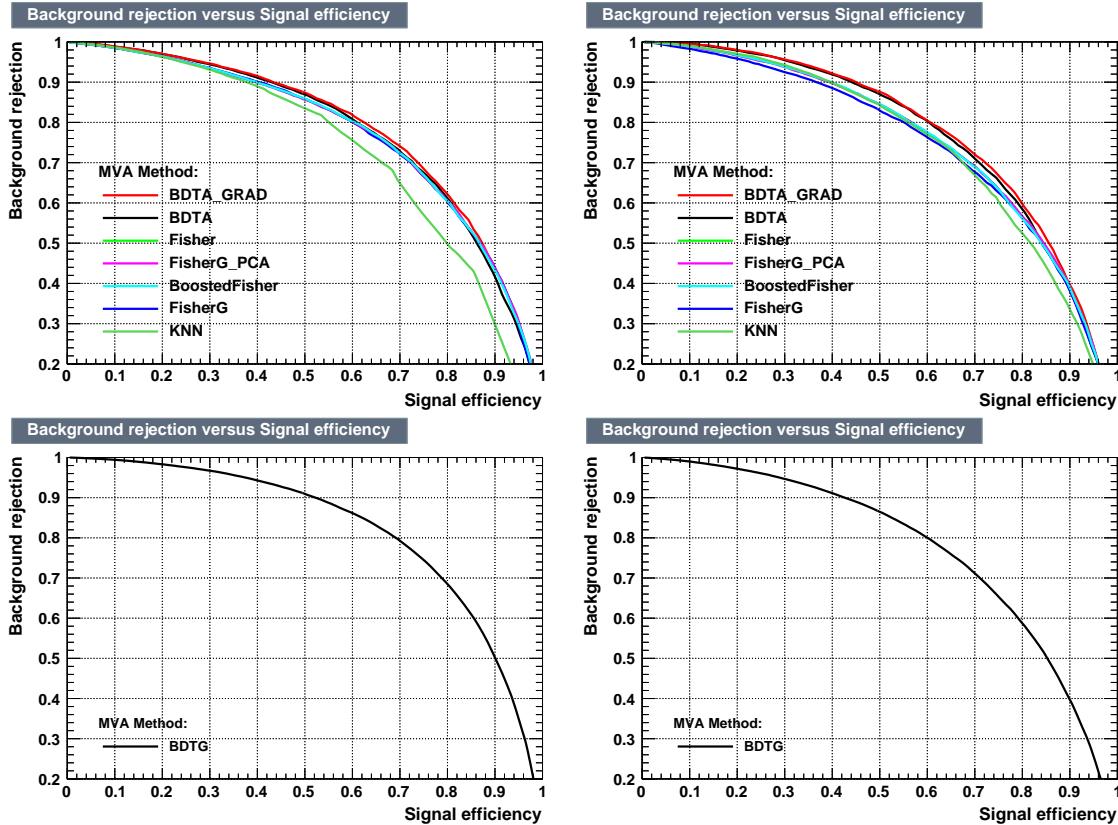


Figure 6.12: Top: Background rejection vs signal efficiency (ROC curves) for various MVA classifiers in the $3l$ channel for training against $t\bar{V}$ (left) and $t\bar{t}$ (right). Bottom: background rejection vs signal efficiency (ROC curve) in the $2lss$ channel for a single discriminator: BDTG, against $t\bar{V}$ (left) and $t\bar{t}$ (right).

2855 to produce a discriminator with the maximum discrimination power. Table 6.9 lists
 2856 the input variables used in the BDTG trainings for this analysis. The same set of
 2857 input variables was used to produce the plots for MVA classifiers evaluation.

2858 Plots in Figure 6.13 shows the BDTG input variables distributions for the signal
 2859 and background samples, in the $3l$ channels.

2860 All the input variables have some discrimination power, however, that power is
 2861 bigger for some of them; for instance, the third lepton p_T plot (top left in Figure 6.13)
 2862 shows some discrimination power against WZ and VVV backgrounds for which there
 2863 is a peak around 30 GeV while tHq peak around 18 GeV; although the discrimination

Variable name	Description
nJet25	Number of jets with $p_T > 25$ GeV, $ \eta < 2.4$
nJetEta1	Number of jets with $ \eta > 1.0$, non-CSV-loose
MaxEtaJet25	Max. $ \eta $ of any (non-CSV-loose) jet with $p_T > 25$ GeV
deltaFwdJetClosestLep	$\Delta\eta$ forward light jet and closest lepton
deltaFwdJetBJet	$\Delta\eta$ forward light jet and hardest CSV loose jet
deltaFwdJet2BJet	$\Delta\eta$ forward light jet and second hardest CSV loose jet
Lep3Pt/Lep2Pt	p_T of the 3 rd lepton (2 nd for ss2l)
totCharge	Sum of lepton charges
minDRll	Min ΔR any two leptons
dphiHighestPtSSPair	$\Delta\phi$ of highest p_T same-sign lepton pair

Table 6.9: BDTG input variables. First section lists variables related to jet multiplicities; second section lists variables related to forward jet activity, and third section lists variables related to lepton kinematics.

2864 power does not cover all the backgrounds, it counts for the final discriminator. A
 2865 similar situation can be seen in the plot for the number of jets (row three, column two);
 2866 $t\bar{t}W$, $t\bar{t}Z$ and $t\bar{t}H$ processes tend to have more jets compared to the tHq process. The
 2867 discrimination power is more evident in other plots like in the plot of the maximum
 2868 $|\eta|$ of the jets in the event (row two, column three). The same or equivalent input
 2869 variables are found to be performing well for both $3l$ and $2lss$ channels. Figure B.1
 2870 shows the corresponding input variables distribution plots for the $2lss$ channel.

2871 **Discrimination power from BDTG classifier**

2872 The Discrimination power of the input variables can also be evaluated from the BDTG
 2873 training, exclusively for the training samples, i.e., dominant backgrounds ($t\bar{t}$ and $t\bar{t}V$);
 2874 the training samples are submitted to the selection cuts on Table 6.6.

2875 Figure 6.14 shows the comparison between input variables for the two trainings
 2876 in the $3l$ channel; it reveals that some variables show opposite behavior for the two
 2877 background sources, which results in potentially screening the discrimination power
 2878 if they were to be used in a single discriminant, i.e., if the training would join $t\bar{t}$ and

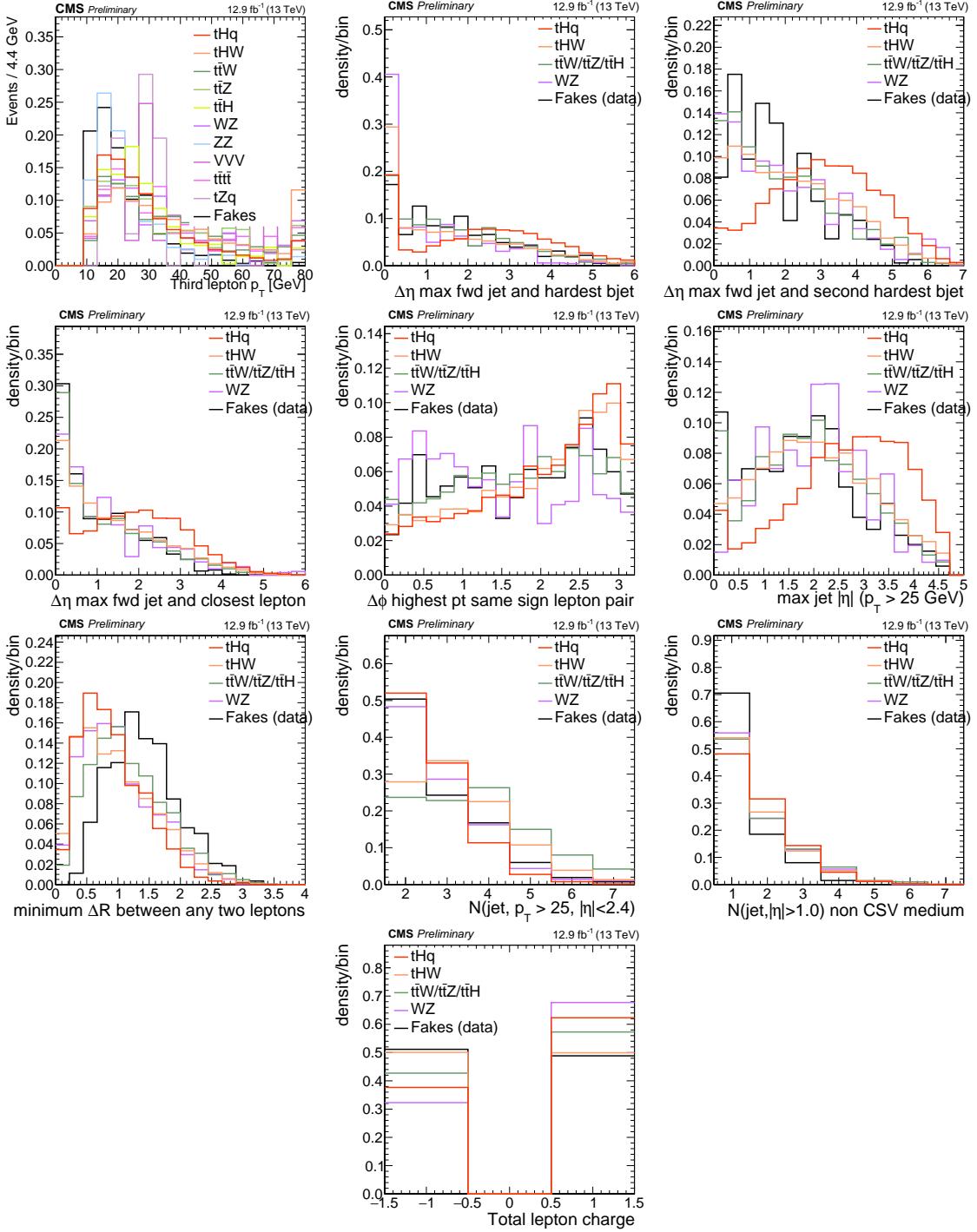


Figure 6.13: Distributions of the BDTG classifier input variables (not normalized) for signal discrimination in the $3l$ channel.

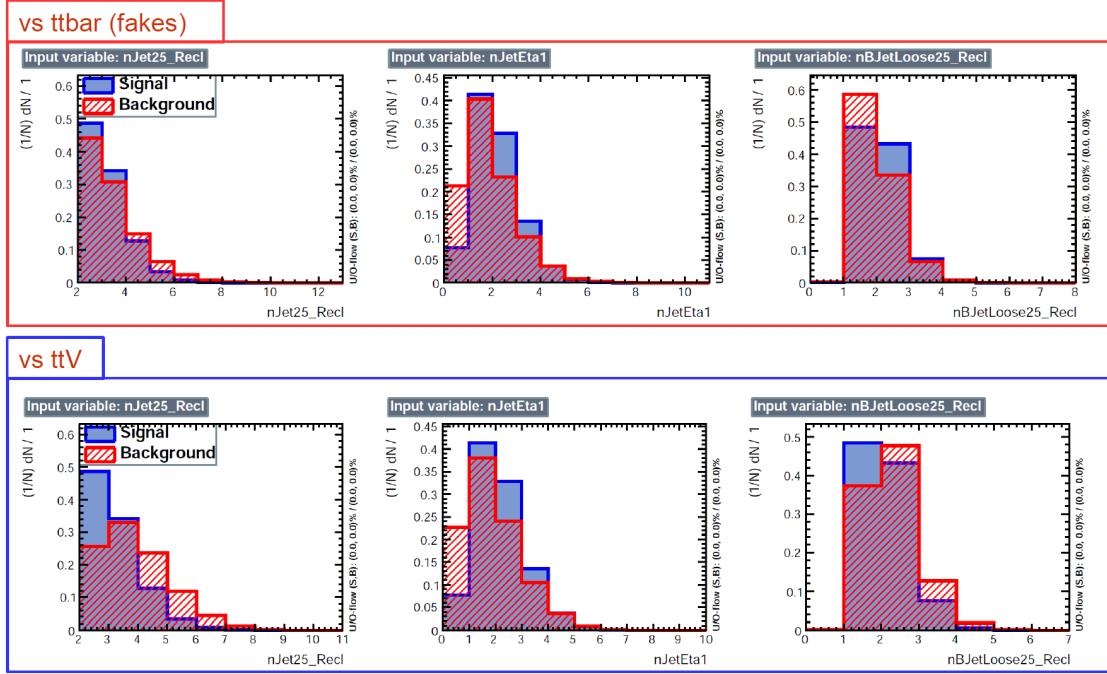


Figure 6.14: BDT input variables as seen by BDTG classifier for the $3l$ channel, tHq signal(blue) discriminated against ttV background (red).

2879 $tt\bar{V}$. For some other variables the distributions are similar in both background cases.

2880 In contrast to the distributions in Figure 6.13 only the dominant backgrounds are
2881 included; however, the discrimination power agrees among plots.

2882 Figures in the Appendix B.2, B.3, B.4, and B.5 show the input variables
2883 distributions for the $2lss$ and $3l$ channel as seen by the BDTG classifier.

2884 Input variables correlations

2885 From Table 6.9, it is clear that the input variables are correlated to some extend.
2886 These correlations play an important role for some MVA methods like the Fisher
2887 discriminant method in which the first step consist of performing a linear transfor-
2888 mation to an phase space where the correlations between variables are removed. In
2889 the case of BDT, correlations do not affect the performance. Figure 6.15 shows the

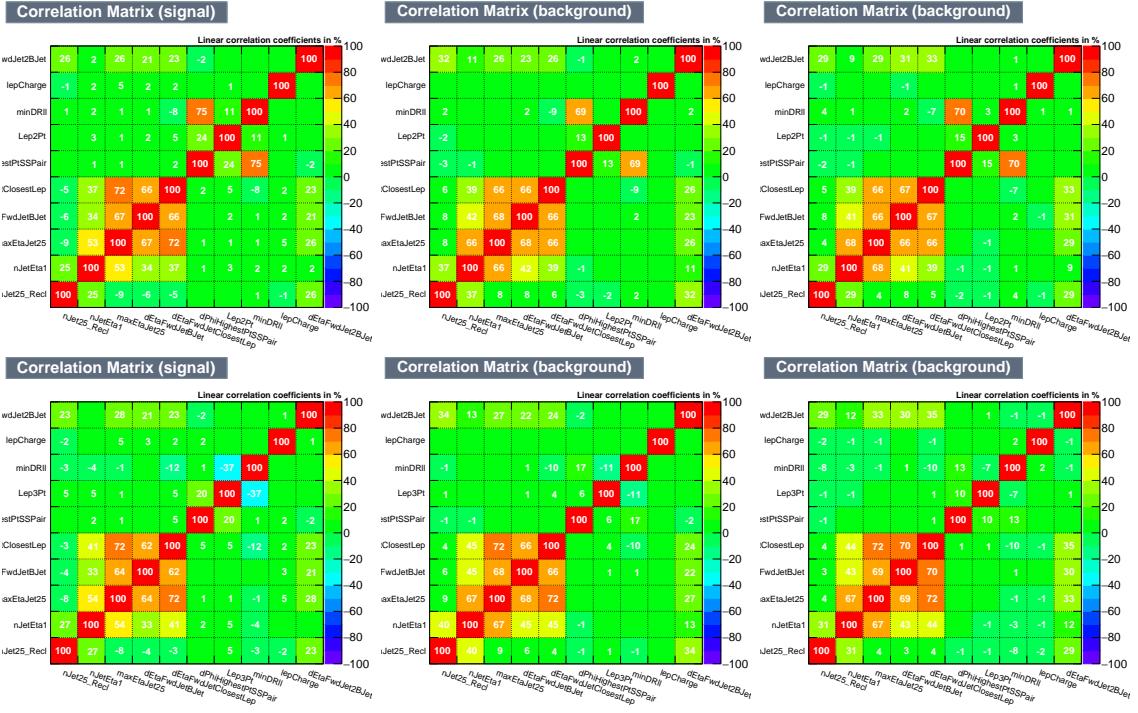


Figure 6.15: Signal (left), $t\bar{t}$ background (middle), and $t\bar{t}V$ background (right.) correlation matrices for the input variables in the BDTG classifier for the $2lss$ (top) and $3l$ (bottom) channels.

linear correlation coefficients for signal and background for the two training cases (the signal values are identical by construction). As expected, strong correlations appears for variables related to the forward jet activity.

6.8.3 BDTG classifiers response

After the training stage, the BDTG classifier is tested to ensure its ability to discriminate between simulated signal and background events. The BDTG classifier output distributions for signal and backgrounds in the $3l$ channel are shown in Figure 6.16. As expected, a good discrimination power is obtained using default discriminator parameter values; some overtraining is also visible.

In order to explore further optimization in the BDTG performance, several changes

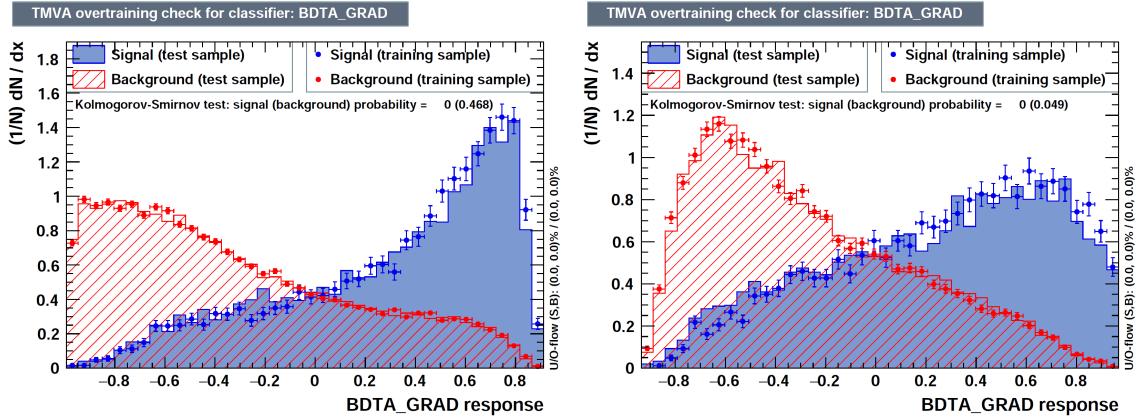


Figure 6.16: BDTG classifier output for trainings against $t\bar{t}V$ (left) and $t\bar{t}$ (right). Default BDTG parameters have been used.

from the default BDTG parameters were tested; Table 6.10 list the set of parameters found to be most discriminant with minimal overtraining as shown in Figure 6.17.

TMVA.Types.kBDT		
Option	Default	Used
NTrees	200	800
BoostType	AdaBoost	Grad
Shrinkage	1	0.1
nCuts	20	50
MaxDepth	3	

Table 6.10: Configuration used in the final BDTG training. Parameters not listed were not tested.

The ranking of the input variables by their importance in the classification process is shown in Table 6.11; for both trainings the rankings show almost the same 5 variables in the first places.

6.8.4 Additional discriminating variables

Given that the forward jet in background processes could be originated from pileup, two additional discriminating variables accounting for that were tested. These additional variables describe the forward jet momentum (`fwdJetPt25`) and the forward jet

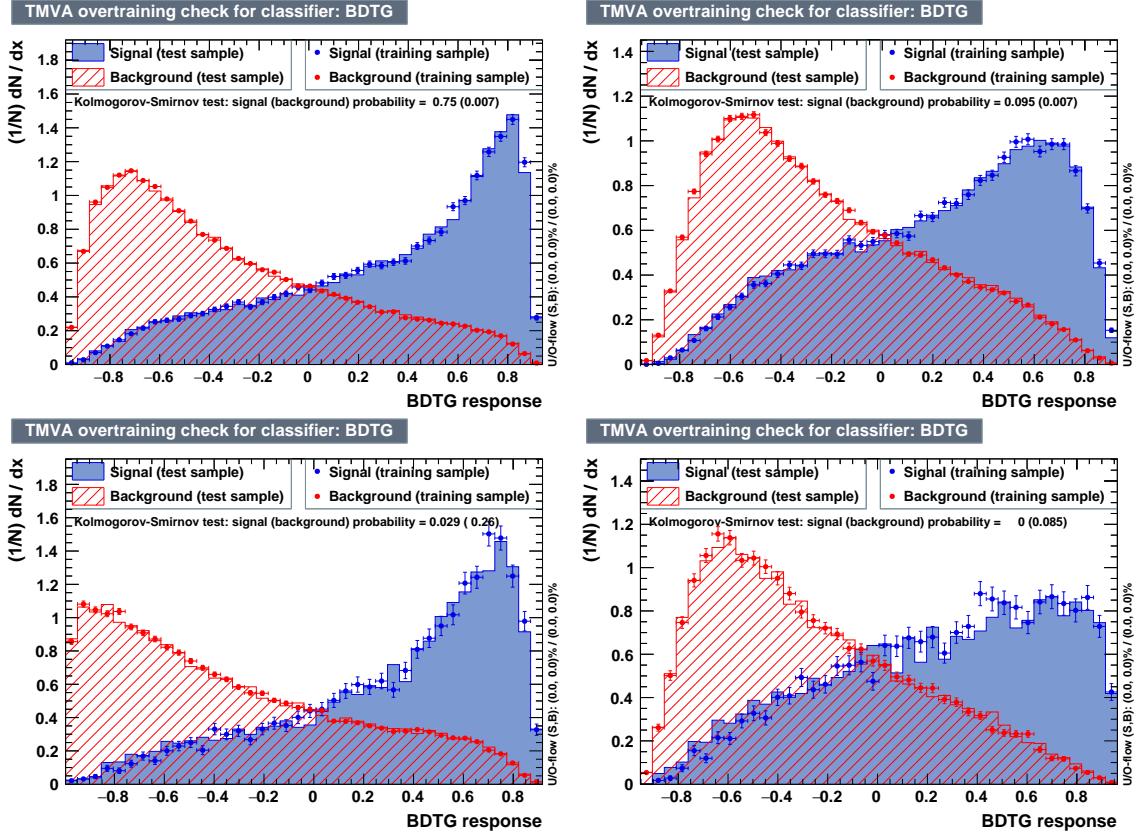


Figure 6.17: BDTG classifiers output for training against $t\bar{t}V$ (left) and $t\bar{t}$ (right) for $2lss$ channel(top) and $3l$ channel (bottom) .

2lss channel		3l channel	
Rank	$t\bar{t}$ training Variable	$t\bar{t}V$ training Variable	$t\bar{t}$ training Variable
1	minDRll	dEtaFwdJetBJet	dEtaFwdJetClosestLep
2	dEtaFwdJetClosestLep	Lep3Pt	minDRll
3	dEtaFwdJetBJet	maxEtaJet25	maxEtaJet25
4	dPhiHighestPtSSPair	dEtaFwdJet2BJet	dPhiHighestPtSSPair
5	Lep3Pt	dEtaFwdJetClosestLep	Lep2Pt
6	maxEtaJet25	minDRll	dEtaFwdJetClosestLep
7	dEtaFwdJet2BJet	dPhiHighestPtSSPair	dEtaFwdJet2BJet
8	nJetEta1	nJet25	nJetEta1
9	nJet25	nJetEta1	nJet25
10	lepCharge	lepCharge	lepCharge

Table 6.11: Input variables ranking for BDTG classifiers for the trainings in the $3l$ channel and $2lss$ channel. In both trainings the rankings show almost the same 5 variables in the first places.

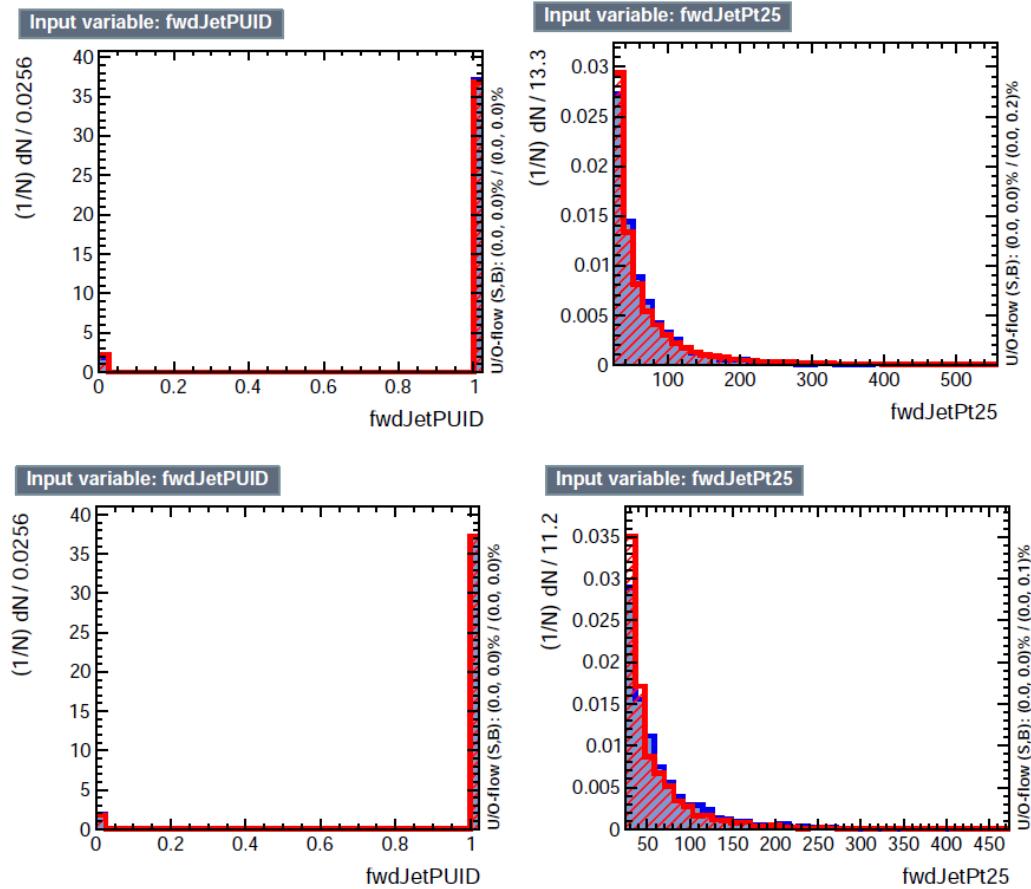


Figure 6.18: Additional discriminating variables distributions for $t\bar{t}V$ training (top row) and $t\bar{t}$ training (bottom row) in the $3l$ channel. The origin of the jets in the forward jet identification distribution is tagged as 0 for *pileup jets* while *primary vertex jets* are tagged as 1.

identification(fwdJetPUID); their distributions in the $3l$ channel are shown in Figure 6.18. The forward jet identification distribution show that for both, signal and background, jets are mostly originated in the primary vertex.

The testing was performed by including in the BDTG input one variable at a time, so the discrimination power of each variable can be evaluated individually, and then both simultaneously. fwdJetPUID was ranked in the last place in importance (11) in both training ($t\bar{t}V$ and $t\bar{t}$) while fwdJetPt25 was ranked 3 in the $t\bar{t}V$ training and 7 in the $t\bar{t}$ training. When training using 12 variables, fwdJetPt25 was ranked 5 and 7 in

2917 the $t\bar{t}V$ and $t\bar{t}$ trainings respectively, while fwdJetPUID was ranked 12 in both cases.

	ROC-integral	
	$t\bar{t}V$	$t\bar{t}$
base 10 var	0.848	0.777
+ fwdJetPUID	0.849	0.777
+ fwdJetPt25	0.856	0.787
12 var	0.856	0.787

Table 6.12: ROC-integral for all the testing cases performed in the evaluation of the additional variables discriminating power. The improvement in the discrimination performance provided by the additional variables is about 1% .

2918 The improvement in the discrimination performance provided by the additional
 2919 variables is about 1%, so it was decided not to include them in the procedure. Table
 2920 6.12 show the ROC-integral for all the testing cases performed.

2921 6.8.5 Signal extraction procedure

2922 Once the two BDTG classifiers, introduced in the previous section, are trained against
 2923 the dominant backgrounds in each channel, they are used to classify the events in the
 2924 samples; their outputs are then used to evaluate the signal cross section limits in a
 2925 fit to the classifier shape. Figure 6.19 shows the expected output distributions in a
 2926 2D plane of one training vs. the other, i.e., $t\bar{t}V$ vs. $t\bar{t}$. Top row shows the 2D planes
 2927 for tHq and tHW signals, while the bottom left plot shows the corresponding 2D
 2928 plane for the combined backgrounds, which are evaluated as in the final background
 2929 prediction, i.e., these are not the samples used in the BDTG training and this includes
 2930 data-driven backgrounds. The signal (combining of tHq and tHW) to background
 2931 ratio (S/B) is showed in the bottom right plot of Figure 6.19.

2932 Each event is now classified into one of ten 2D-bins according to its position in the
 2933 plane, as shown in Figure 6.20. The number of bins is chosen such that no bins are
 2934 entirely empty for any process. The bin boundary positions and number of bins have

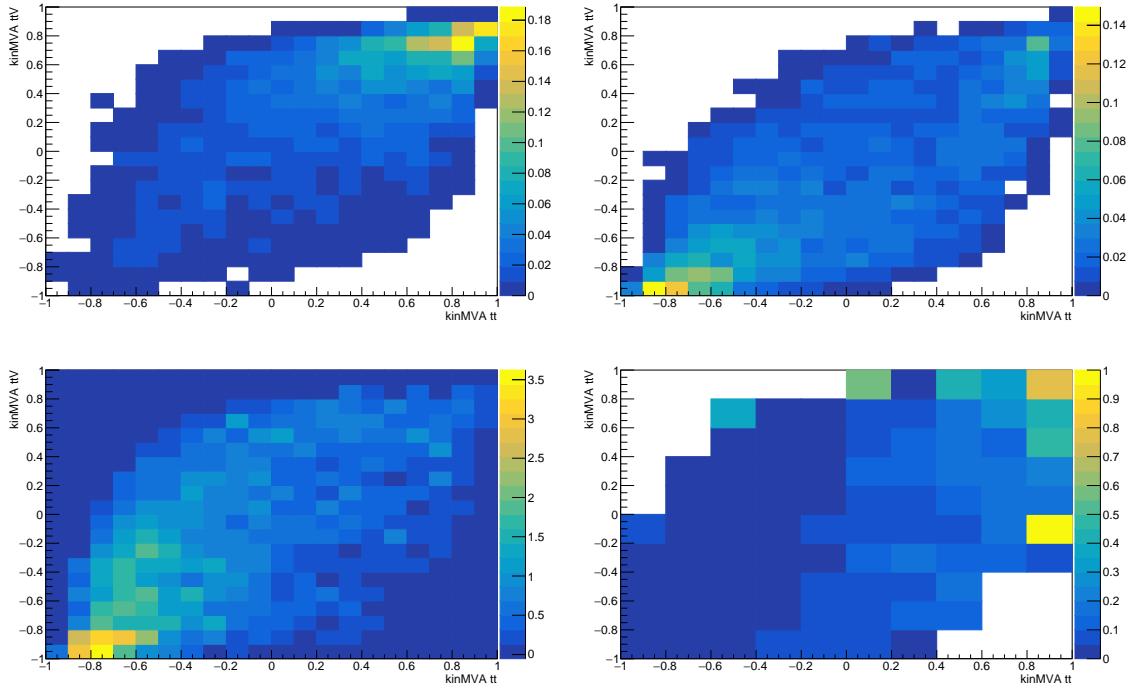


Figure 6.19: BDT classifier output planes (training vs $t\bar{t}$ on x-axis and vs $t\bar{t}V$ on y-axis) for the tHq and tHW signals (top row), and for the combined backgrounds (bottom left). Bottom right: S/B ratio (combining tHq and tHW) in the same plane. Plots are for $3l$ channel.

2935 been studied and optimized with respect to the expected limit on the signal strength
2936 (see Sec. 6.8.6).

2937 From this event categorization, a 1D histogram of expected distribution is pro-
2938 duced for each signal and background process, and fit to the observed data (or the
2939 Asimov dataset for expected limits).

2940 6.8.6 Binning and selection optimization

2941 The effect of the choice of pre-selection cuts and the number of bins of the 1D his-
2942 togram on the cross section limit is evaluated by varying the most important cuts and
2943 re-calculating the limit in each case. In this analysis, the optimization was performed
2944 in the $3l$ channel, by evaluating the upper limits on the $tHq + tHW$ expected signal

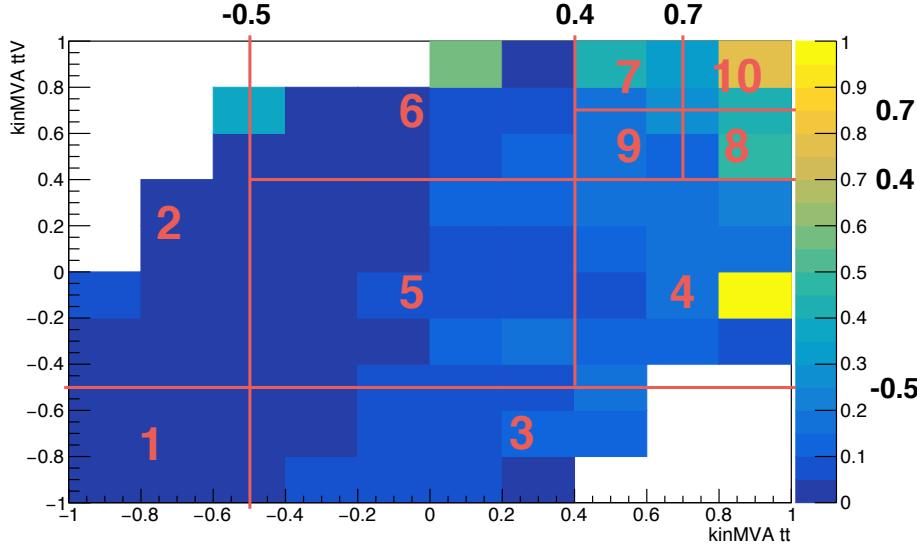


Figure 6.20: Binning overlaid on the S/B ratio map on the plane of classifier outputs.

2945 strength only (without $t\bar{t}H$ component), always evaluated at $\kappa_t = -1.0$, $\kappa_V = 1.0$.

2946 Table 6.13 shows the several variations explored, compared with a baseline; the
 2947 baseline is similar to the selection reported in Table 6.6 but only a loose CSV jet and
 2948 a Z veto of ± 10 GeV are required.

Selection	Variation	Expected limit
Baseline		< 2.93
Loose CSV tags	$\geq 1 \rightarrow \geq 2$	< 3.81
Medium CSV tags	$\geq 0 \rightarrow \geq 1$	< 2.76
Light forward jet η	$\geq 0 \rightarrow \geq 1$	< 2.94
Light forward jet η	$\geq 0 \rightarrow \geq 1.5$	< 3.00
MET > 30 GeV		< 2.91
Z veto ($ m_{\ell\ell} - m_Z $)	$> 10\text{GeV} \rightarrow > 15\text{GeV}$	< 2.79
One medium CSV + 15 GeV Z veto	combined	< 2.62

Table 6.13: Signal strength limit variation as a function of tighter cuts. The baseline selection corresponds to a looser selection compared to the one reported in Tab. 6.6 where only a CSV-loose b -jet is required, and the Z veto is loosened to ± 10 GeV. The optimal selection determined here corresponds to the baseline plus the two variations in the last row.

2949 The optimal limit is found when requiring a slightly tighter selection with respect
 2950 to the baseline. The optimal selection is reported in Table 6.6.

2951 The signal strength limit also depends on the chosen binning in the 2D plane as
 2952 the S/B ratio varies across the plane, hence, several sizes and binning combinations
 2953 were tested in order to improve the limit. Figure 6.21 shows some of the binning
 2954 combinations tested; in the default combination all the bins have the same size, while
 2955 the best limit was found for a set of 10 bins. The bin borders and the resulting limits
 2956 are shown in Table 6.14.

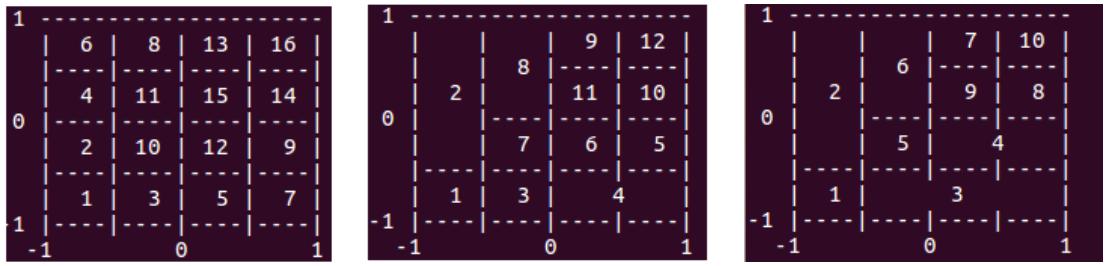


Figure 6.21: Binning combination scheme.

Number of bins	Bin borders						Expected limit
	x_1	x_2	x_3	y_1	y_2	y_3	
16 (default)	-0.5	0.0	0.5	-0.5	0.0	0.5	< 2.91
16	-0.5	0.3	0.7	-0.5	0.3	0.7	< 2.83
10	-0.5	0.0	0.5	-0.5	0.0	0.5	< 2.93
10	-0.5	0.0	0.7	-0.5	0.0	0.7	< 2.86
10	-0.5	0.0	0.7	-0.5	0.0	0.5	< 2.84
10	-0.5	0.0	0.5	-0.5	0.0	0.7	< 2.87
10	-0.5	0.4	0.7	-0.5	0.4	0.7	< 2.81

Table 6.14: Limit variation as a function of bin size. The final bin borders used in the $3l$ channel are indicated in bold.

2957 Combining the optimization of binning and using the tighter pre-selection cuts,
 2958 the expected limit in the $3l$ channel alone reaches **r<2.59**.

2959 A similar binning optimization was made for $2lss$ channel, including other binning
 2960 combinations. First, the $3l$ channel binning was used to estimate the expected limit,
 2961 then, bin borders were varied to obtain the best possible expected limit. The bin

2962 borders and the resulting signal strength limits for the same-sign dimuon channel are
 2963 shown in Table 6.15:

Number of bins	Bin borders						Expected limit
	x_1	x_2	x_3	y_1	y_2	y_3	
16	-0.5	0.4	0.7	-0.5	0.4	0.7	< 1.72
12	-0.5	0.4	0.7	-0.5	0.4	0.7	< 1.72
12	-0.3	0.4	0.7	-0.5	0.4	0.7	< 1.71
12	-0.3	0.3	0.7	-0.5	0.4	0.7	< 1.71
12	-0.3	0.3	0.7	-0.4	0.4	0.7	< 1.70
12	-0.3	0.3	0.7	-0.3	0.4	0.7	< 1.70
12	-0.3	0.3	0.7	-0.3	0.2	0.7	< 1.68
12	-0.3	0.3	0.7	-0.3	0.1	0.7	< 1.70
12	-0.3	0.3	0.7	-0.3	0.2	0.6	< 1.70
10	-0.5	0.4	0.7	-0.5	0.4	0.7	< 1.75
10	-0.3	0.3	0.7	-0.3	0.2	0.6	< 1.69

Table 6.15: Limit variation as a function of bin size in the same-sign dimuon channel. (In bold: the final bin borders used in the $2lss$ channel.)

2964 The expected limit was found to be **r<1.69** for optimized bin borders in 10 bins
 2965 and optimized pre-selection cuts.

2966 Two additional binning strategies were tested, however, the obtained limits are
 2967 degraded; they are documented in Appendix C.

2968 6.9 Signal model

2969 The goal of this analysis is to test the compatibility of points in the parameter space of
 2970 Higgs-to-vector boson and Higgs-to-top quark couplings. The simulated tHq , tHW ,
 2971 and $t\bar{t}H$ signal events are used with event-by-event weights to reflect the impact of the
 2972 couplings on kinematic distributions, and together with different predictions of the
 2973 respective production cross sections and branching ratios, we can produce limits for
 2974 different values of κ_V and κ_t . (See Tab. A.3 for the set of κ_t and κ_V values generated.)

2975 The slight shape-dependence of the BDT outputs as a function of the couplings is
 2976 documented in Appendix D.

2977 Apart from the κ_t/κ_V interference of the tHq and tHW production cross sections,
 2978 the cross section of $t\bar{t}H$ scales as κ_t^2 . Furthermore, the Higgs branching fractions to
 2979 vector bosons depend on κ_V , and the overall Higgs decay width depend both on κ_t
 2980 and κ_V when considering resolved top-quark loops in the $H \rightarrow \gamma\gamma$, $H \rightarrow Z\gamma$, and
 2981 $H \rightarrow gg$ decays. The relative contributions from $H \rightarrow WW$, $H \rightarrow ZZ$, and $H \rightarrow \tau\tau$
 2982 changes with changing κ_V .

2983 We hence set an upper limit on the combined cross section times branching ratio
 2984 of tHq , tHW , and $t\bar{t}H$.

2985 If we assume a modifier for the Higgs-to-tau coupling (κ_τ) to be equal to κ_t , the
 2986 relative fractions of WW , ZZ , and $\tau\tau$ in our selection will only depend on the ratio
 2987 of κ_t/κ_V . Any limit set at any given value of κ_t/κ_V is thus valid for all values of
 2988 κ_t and κ_V with that ratio, and could then be compared with theoretical predictions
 2989 of cross sections at different values of either modifier. Rather than as a function of
 2990 the κ_t/κ_V ratio, limits could (equivalently) be reported as a function of the relative
 2991 strength of Higgs-top and Higgs-vector-boson couplings, multiplied by the relative
 2992 sign. Such a parameter, further referred to as f_t , as defined in Equation 6.7, spans
 2993 the entire possible parameter space between -1.0 and 1.0 , with the SM expectation
 2994 at 0.5 . Absolute values of 1.0 or 0.0 would then correspond to purely Higgs-top and
 2995 purely Higgs-V couplings, respectively.

$$f_t = \text{sign}(\kappa_t/\kappa_V) \times \frac{\kappa_t^2}{\kappa_t^2 + \kappa_V^2}. \quad (6.7)$$

2996 Table 6.16 shows the points in the κ_t/κ_V and f_t parameter space that are mapped
 2997 by the 51 individual κ_t and κ_V points.

f_t	κ_t/κ_V	$\kappa_V = 0.5$	$\kappa_V = 1.0$	$\kappa_V = 1.5$
-0.973	-6.000	-3.00		
-0.941	-4.000	-2.00		
-0.900	-3.000	-1.50	-3.00	
-0.862	-2.500	-1.25		
-0.800	-2.000	-1.00	-2.00	-3.00
-0.692	-1.500	-0.75	-1.50	
-0.640	-1.333			-2.00
-0.610	-1.250		-1.25	
-0.500	-1.000	-0.50	-1.00	-1.50
-0.410	-0.833			-1.25
-0.360	-0.750		-0.75	
-0.308	-0.667			-1.00
-0.200	-0.500	-0.25	-0.50	-0.75
-0.100	-0.333			-0.50
-0.059	-0.250		-0.25	
-0.027	-0.167			-0.25
0.000	0.000	0.00	0.00	0.00
0.027	0.167			0.25
0.059	0.250		0.25	
0.100	0.333			0.50
0.200	0.500	0.25	0.50	0.75
0.308	0.667			1.00
0.360	0.750		0.75	
0.410	0.833			1.25
0.500	1.000	0.50	1.00	1.50
0.610	1.250		1.25	
0.640	1.333			2.00
0.692	1.500	0.75	1.50	
0.800	2.000	1.00	2.00	3.00
0.862	2.500	1.25		
0.900	3.000	1.50	3.00	
0.941	4.000	2.00		
0.973	6.000	3.00		

Table 6.16: The 33 distinct values of κ_t/κ_V and f_t as mapped by the 51 κ_t and κ_V points.

2998 The overall higgs decay width (modified by both κ_t and κ_V) becomes irrelevant
 2999 if limits are quoted as absolute cross sections rather than multiples of the expected
 3000 cross section (which depends on the overall Higgs decay width).

3001 The 1D histograms of events as categorized in regions of the 2D BDT plane is
 3002 then used in a maximum likelihood fit of signal and background shapes, where the
 3003 tHq , tHW , and $t\bar{t}H$ signals are floating with a common signal strength modifier r ,
 3004 producing a 95% C.L. upper limit the observed cross section of $tHq + tHW + t\bar{t}H$.

3005 This is done separately for each point of κ_t and κ_V , where the cross sections and
 3006 branching fractions are scaled accordingly in each point. Limits at fixed values of
 3007 κ_t/κ_V are by construction identical. Tables ??–?? and ??–?? in Appendix ?? show
 3008 the scalings of cross section times branching fraction, as well as branching fractions
 3009 alone for each of the Higgs decay modes and each of the signal components.

3010 Systematic uncertainties on the signal selection efficiency arise from correction
 3011 factors applied to the simulated events to better match the measured detector perfor-
 3012 mance and also from theoretical uncertainties in the modeling of the signal process.
 3013 Scale factors applied to correct for data/MC differences in the trigger efficiency, lepton
 3014 reconstruction and identification performance, and lepton selection efficiency carry a
 3015 combined uncertainty of about 5% from jet energy corrections is evaluated by varying
 3016 the correction factors within their uncertainty and propagating the effect to the final
 3017 result by recalculating all kinematic quantities. Effects on the overall normalization
 3018 of event yields and on the shape of kinematic properties are both taken into account.
 3019 Jet energy resolution effects have negligible impact on this

3020 analysis. Correction factors for data/MC differences in the b-tagging performance
 3021 are applied depending on the pT and \hat{t} , and on the flavor of the jet, and their effect
 3022 on the signal efficiency is evaluated by varying the factors within their measured
 3023 uncertainty and recalculating the overall event scale factors. The uncertainties from
 3024 unknown higher orders of tHq and tHW production are estimated from a change
 3025 in the Q2 scale of double and half the initial value, evaluated for each point of \hat{t}_Z
 3026 and \hat{t}_V . The tH signal component has an uncertainty of about $+5.8/-9.2$ scale

3027 variations and a further 3.6Uncertainties related to the choice of PDF set and its scale
3028 are estimated to be about 3.7tHq and about 4.0

³⁰²⁹ **Appendix A**

³⁰³⁰ **Datasets and triggers**

Dataset name
/JetHT/Run2016X-23Sep2016-vY/MINIAOD
/HTMHT/Run2016X-23Sep2016-vY/MINIAOD
/MET/Run2016X-23Sep2016-vY/MINIAOD
/SingleElectron/Run2016X-23Sep2016-vY/MINIAOD
/SingleMuon/Run2016X-23Sep2016-vY/MINIAOD
/SinglePhoton/Run2016X-23Sep2016-vY/MINIAOD
/DoubleEG/Run2016X-23Sep2016-vY/MINIAOD
/MuonEG/Run2016X-23Sep2016-vY/MINIAOD
/DoubleMuon/Run2016X-23Sep2016-vY/MINIAOD
/Tau/Run2016B-23Sep2016-v3/MINIAOD
/JetHT/Run2016H-PromptReco-v3/MINIAOD
/HTMHT/Run2016H-PromptReco-v3/MINIAOD
/MET/Run2016H-PromptReco-v3/MINIAOD
/SingleElectron/Run2016H-PromptReco-v3/MINIAOD
/SingleMuon/Run2016H-PromptReco-v3/MINIAOD
/SinglePhoton/Run2016H-PromptReco-v3/MINIAOD
/DoubleEG/Run2016H-PromptReco-v3/MINIAOD
/MuonEG/Run2016H-PromptReco-v3/MINIAOD
/DoubleMuon/Run2016H-PromptReco-v3/MINIAOD

Table A.1: Full 2016 dataset used in the analysis. In the first section of the table are listed the 23Sep2016 samples; in the Run2016X-23Sep-vY label, X:B-G tag the run while Y:1,3 tag the version of the data sample. Second section list the PromptReco version of the dataset.

Same-sign dilepton (==2 muons)
HLT_Mu17_TrkIsoVVL_Mu8_TrkIsoVVL_DZ_v*
HLT_Mu17_TrkIsoVVL_TkMu8_TrkIsoVVL_DZ_v*
HLT_IsoMu22_v*
HLT_IsoTkMu22_v*
HLT_IsoMu22_eta2p1_v*
HLT_IsoTkMu22_eta2p1_v*
HLT_IsoMu24_v*
HLT_IsoTkMu24_v*
Same-sign dilepton (==2 electrons)
HLT_Ele23_Ele12_CaloIdL_TrackIdL_IsoVL_DZ_v*
HLT_Ele27_eta2p1_WP Loose_Gsf_v*
HLT_Ele27_WPTight_Gsf_v*
HLT_Ele25_eta2p1_WPTight_Gsf_v*
Same-sign dilepton (==1 muon, ==1 electron)
HLT_Mu23_TrkIsoVVL_Ele8_CaloIdL_TrackIdL_IsoVL_v*
HLT_Mu8_TrkIsoVVL_Ele23_CaloIdL_TrackIdL_IsoVL_v*
HLT_Mu23_TrkIsoVVL_Ele8_CaloIdL_TrackIdL_IsoVL_DZ_v*
HLT_Mu8_TrkIsoVVL_Ele23_CaloIdL_TrackIdL_IsoVL_DZ_v*
HLT_IsoMu22_v*
HLT_IsoTkMu22_v*
HLT_IsoMu22_eta2p1_v*
HLT_IsoTkMu22_eta2p1_v*
HLT_IsoMu24_v*
HLT_IsoTkMu24_v*
HLT_Ele27_WPTight_Gsf_v*
HLT_Ele25_eta2p1_WPTight_Gsf_v*
HLT_Ele27_eta2p1_WP Loose_Gsf_v*
Three lepton
HLT_DiMu9_Ele9_CaloIdL_TrackIdL_v*
HLT_Mu8_DiEle12_CaloIdL_TrackIdL_v*
HLT_TripleMu_12_10_5_v*
HLT_Ele16_Ele12_Ele8_CaloIdL_TrackIdL_v*
HLT_Mu23_TrkIsoVVL_Ele8_CaloIdL_TrackIdL_IsoVL_v*
HLT_Mu23_TrkIsoVVL_Ele8_CaloIdL_TrackIdL_IsoVL_DZ_v*
HLT_Mu8_TrkIsoVVL_Ele23_CaloIdL_TrackIdL_IsoVL_v*
HLT_Mu8_TrkIsoVVL_Ele23_CaloIdL_TrackIdL_IsoVL_DZ_v*
HLT_Ele23_Ele12_CaloIdL_TrackIdL_IsoVL_DZ_v*
HLT_Mu17_TrkIsoVVL_Mu8_TrkIsoVVL_DZ_v*
HLT_Mu17_TrkIsoVVL_TkMu8_TrkIsoVVL_DZ_v*
HLT_IsoMu22_v*
HLT_IsoTkMu22_v*
HLT_IsoMu22_eta2p1_v*
HLT_IsoTkMu22_eta2p1_v*
HLT_IsoMu24_v*
HLT_IsoTkMu24_v*
HLT_Ele27_WPTight_Gsf_v*
HLT_Ele25_eta2p1_WPTight_Gsf_v*
HLT_Ele27_eta2p1_WP Loose_Gsf_v*

Table A.2: Table of high-level triggers considered in the analysis.

		<i>tHq</i>		<i>tHW</i>		
κ_V	κ_t	sum of weights	cross section [pb]	sum of weights	cross section [pb]	LHE weights
1.0	-3.0	35.700022	2.991	11.030445	0.6409	LHEweight_wgt[446]
1.0	-2.0	20.124298	1.706	5.967205	0.3458	LHEweight_wgt[447]
1.0	-1.5	14.043198	1.205	4.029093	0.2353	LHEweight_wgt[448]
1.0	-1.25	11.429338	0.9869	3.208415	0.1876	LHEweight_wgt[449]
1.0	-1.0		0.7927		0.1472	
1.0	-0.75	7.054998	0.6212	1.863811	0.1102	LHEweight_wgt[450]
1.0	-0.5	5.294518	0.4723	1.339886	0.07979	LHEweight_wgt[451]
1.0	-0.25	3.818499	0.3505	0.914880	0.05518	LHEweight_wgt[452]
1.0	0.0	2.627360	0.2482	0.588902	0.03881	LHEweight_wgt[453]
1.0	0.25	1.719841	0.1694	0.361621	0.02226	LHEweight_wgt[454]
1.0	0.5	1.097202	0.1133	0.233368	0.01444	LHEweight_wgt[455]
1.0	0.75	0.759024	0.08059	0.204034	0.01222	LHEweight_wgt[456]
1.0	1.0	0.705305	0.07096	0.273617	0.01561	LHEweight_wgt[457]
1.0	1.25	0.936047	0.0839	0.442119	0.02481	LHEweight_wgt[458]
1.0	1.5	1.451249	0.1199	0.709538	0.03935	LHEweight_wgt[459]
1.0	2.0	3.335034	0.2602	1.541132	0.08605	LHEweight_wgt[460]
1.0	3.0	10.516125	0.8210	4.391335	0.2465	LHEweight_wgt[461]
1.5	-3.0	45.281492	3.845	13.426212	0.7825	LHEweight_wgt[462]
1.5	-2.0	27.606715	2.371	7.809713	0.4574	LHEweight_wgt[463]
1.5	-1.5	20.476088	1.784	5.594971	0.3290	LHEweight_wgt[464]
1.5	-1.25	17.337465	1.518	4.635978	0.2749	LHEweight_wgt[465]
1.5	-1.0	14.483302	1.287	3.775902	0.2244	LHEweight_wgt[466]
1.5	-0.75	11.913599	1.067	3.014744	0.1799	LHEweight_wgt[467]
1.5	-0.5	9.628357	0.874	2.352505	0.1410	LHEweight_wgt[468]
1.5	-0.25	7.627574	0.702	1.789184	0.1081	LHEweight_wgt[469]
1.5	0.0	5.911882	0.5577	1.324946	0.08056	LHEweight_wgt[470]
1.5	0.25	4.479390	0.4365	0.959295	0.05893	LHEweight_wgt[471]
1.5	0.5	3.331988	0.3343	0.692727	0.04277	LHEweight_wgt[472]
1.5	0.75	2.469046	0.2558	0.525078	0.03263	LHEweight_wgt[473]
1.5	1.0	1.890565	0.2003	0.456347	0.02768	LHEweight_wgt[474]
1.5	1.25	1.596544	0.1689	0.486534	0.02864	LHEweight_wgt[475]
1.5	1.5	1.586983	0.1594	0.615638	0.03509	LHEweight_wgt[476]
1.5	2.0	2.421241	0.2105	1.170602	0.06515	LHEweight_wgt[477]
1.5	3.0	7.503280	0.5889	3.467546	0.1930	LHEweight_wgt[478]
0.5	-3.0	27.432685	2.260	8.929074	0.5136	LHEweight_wgt[479]
0.5	-2.0	13.956013	1.160	4.419093	0.2547	LHEweight_wgt[480]
0.5	-1.5	8.924438	0.7478	2.757611	0.1591	LHEweight_wgt[481]
0.5	-1.25	6.835341	0.5726	2.075247	0.1204	LHEweight_wgt[482]
0.5	-1.0	5.030704	0.4273	1.491801	0.08696	LHEweight_wgt[483]
0.5	-0.75	3.510528	0.2999	1.007273	0.05885	LHEweight_wgt[484]
0.5	-0.5	2.274811	0.1982	0.621663	0.03658	LHEweight_wgt[485]
0.5	-0.25	1.323555	0.1189	0.334972	0.01996	LHEweight_wgt[486]
0.5	0.0	0.656969	0.06223	0.147253	0.008986	LHEweight_wgt[487]
0.5	0.25	0.274423	0.02830	0.058342	0.003608	LHEweight_wgt[488]
0.5	0.5	0.176548	0.01778	0.068404	0.003902	LHEweight_wgt[489]
0.5	0.75	0.363132	0.03008	0.177385	0.009854	LHEweight_wgt[490]
0.5	1.0	0.834177	0.06550	0.385283	0.02145	LHEweight_wgt[491]
0.5	1.25	1.589682	0.1241	0.692099	0.03848	LHEweight_wgt[492]
0.5	1.5	2.629647	0.2047	1.097834	0.06136	LHEweight_wgt[493]
0.5	2.0	5.562958	0.4358	2.206057	0.1246	LHEweight_wgt[494]
0.5	3.0	14.843102	1.177	5.609519	0.3172	LHEweight_wgt[495]

Table A.3: κ_V and κ_t combinations generated for the two signal samples and their NLO cross sections. The *tHq* cross section is multiplied by the branching fraction of the enforced leptonic decay of the top quark (0.324) [150].

Sample	σ [pb]	*
TTWJetsToLNu_TuneCUETP8M1_13TeV-amcatnloFXFX-madspin-pythia8	0.2043	*
TTZToLLNuNu_M-10_TuneCUETP8M1_13TeV-amcatnlo-pythia8	0.2529	*
/store/cmst3/group/susy/gpetrucc/13TeV/u/TTLL_m1to10_L0_NoMS_for76X/	0.0283	
WGToLNuG_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	585.8	
ZGTo2LG_TuneCUETP8M1_13TeV-amcatnloFXFX-pythia8	131.3	
TGJets_TuneCUETP8M1_13TeV_amcatnlo_madspin_pythia8	2.967	
TGJets_TuneCUETP8M1_13TeV_amcatnlo_madspin_pythia8	2.967	
TTGJets_TuneCUETP8M1_13TeV-amcatnloFXFX-madspin-pythia8	3.697	
WpWpJJ_EWK-QCD_TuneCUETP8M1_13TeV-madgraph-pythia8	0.03711	
ZZZ_TuneCUETP8M1_13TeV-amcatnlo-pythia8	0.01398	
WWZ_TuneCUETP8M1_13TeV-amcatnlo-pythia8	0.1651	
WZZ_TuneCUETP8M1_13TeV-amcatnlo-pythia8	0.05565	
WW_DoubleScattering_13TeV-pythia8	1.64	
tZq_ll_4f_13TeV-amcatnlo-pythia8_TuneCUETP8M1	0.0758	
ST_tW1l_5f_L0_13TeV-MadGraph-pythia8	0.01123	
TTTT_TuneCUETP8M1_13TeV-amcatnlo-pythia8	0.009103	
WZTo3LNu_TuneCUETP8M1_13TeV-powheg-pythia8	4.4296	
ZZTo4L_13TeV_powheg_pythia8	1.256	
TTJets_SingleLeptFromTbar_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	182.1754	*
TTJets_SingleLeptFromT_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	182.1754	*
TTJets_DiLept_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	87.3	*
DYJetsToLL_M-10to50_TuneCUETP8M1_13TeV-amcatnloFXFX-pythia8	18610	
DYJetsToLL_M-50_TuneCUETP8M1_13TeV-madgraphMLM-pythia8	6024	
WJetsToLNu_TuneCUETP8M1_13TeV-amcatnloFXFX-pythia8	61526.7	
ST_tW_top_5f_inclusiveDecays_13TeV-powheg-pythia8_TuneCUETP8M1	35.6	
ST_tW_antitop_5f_inclusiveDecays_13TeV-powheg-pythia8_TuneCUETP8M1	35.6	
ST_t-channel_4f_leptonDecays_13TeV-amcatnlo-pythia8_TuneCUETP8M1	70.3144	
ST_t-channel_antitop_4f_leptonDecays_13TeV-powheg-pythia8_TuneCUETP8M1	26.2278	
ST_s-channel_4f_leptonDecays_13TeV-amcatnlo-pythia8_TuneCUETP8M1	3.68064	
WWTo2L2Nu_13TeV-powheg	10.481	
ttWJets_13TeV_madgraphMLM	0.6105	
ttZJets_13TeV_madgraphMLM	0.5297/0.692	

Table A.4: List of background samples used in this analysis (CMSSW 80X). The first section of the table lists the samples used in simulation to extract the final yields and shapes; the second section lists the samples of the processes for which the yields are estimated from data. The MC simulation is used to design the data driven methods and in the derivation of the associated systematic uncertainties. The third section list the leading order $t\bar{t}W$ and $t\bar{t}Z$ samples, which in addition to the ones marked with a *, where used in the BDT training.

³⁰³¹ **Appendix B**

³⁰³² **BDTG aditional plots**

³⁰³³ **B.1 BDTG input variables distributions for $2lss$**

³⁰³⁴ **channel**

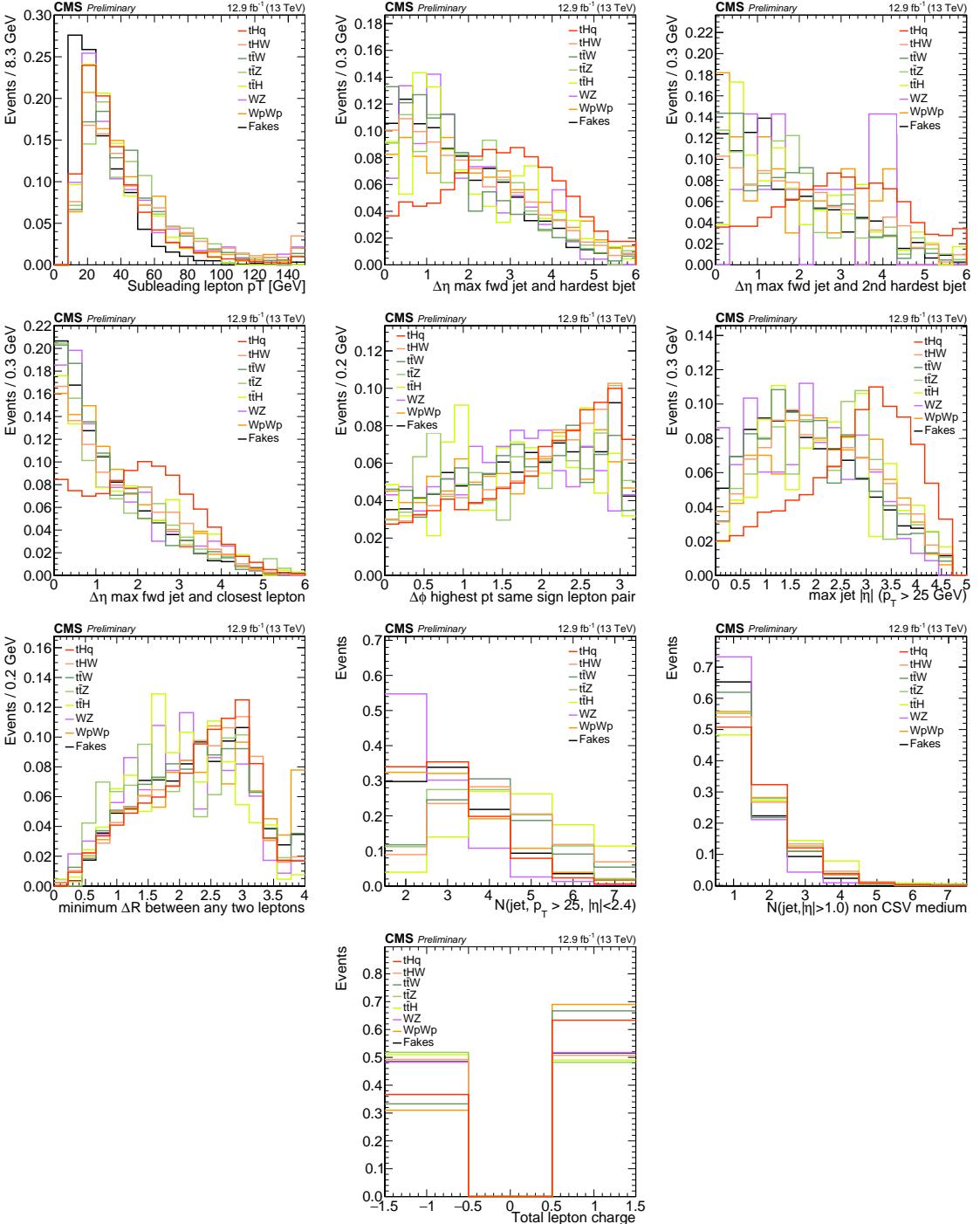


Figure B.1: Distributions of input variables to the BDT for signal discrimination, two lepton same sign channel.

3035 **B.2 Input variables distributions from BDTG**
 3036 classifiers

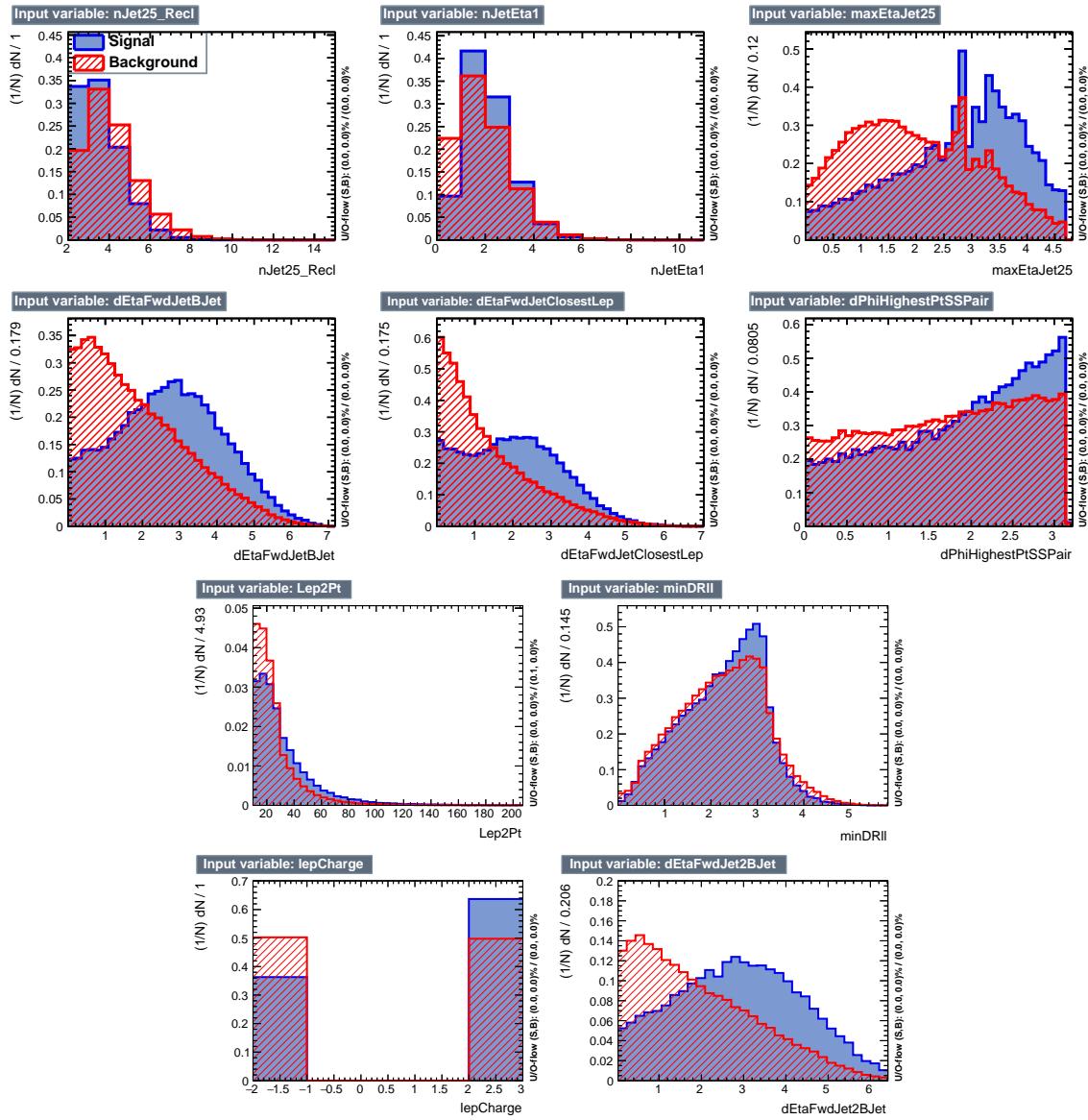


Figure B.2: BDT input variables as seen by BDTG classifier for the $2lss$ channel, tHq signal (blue) discriminated against $t\bar{t}$ background (red).

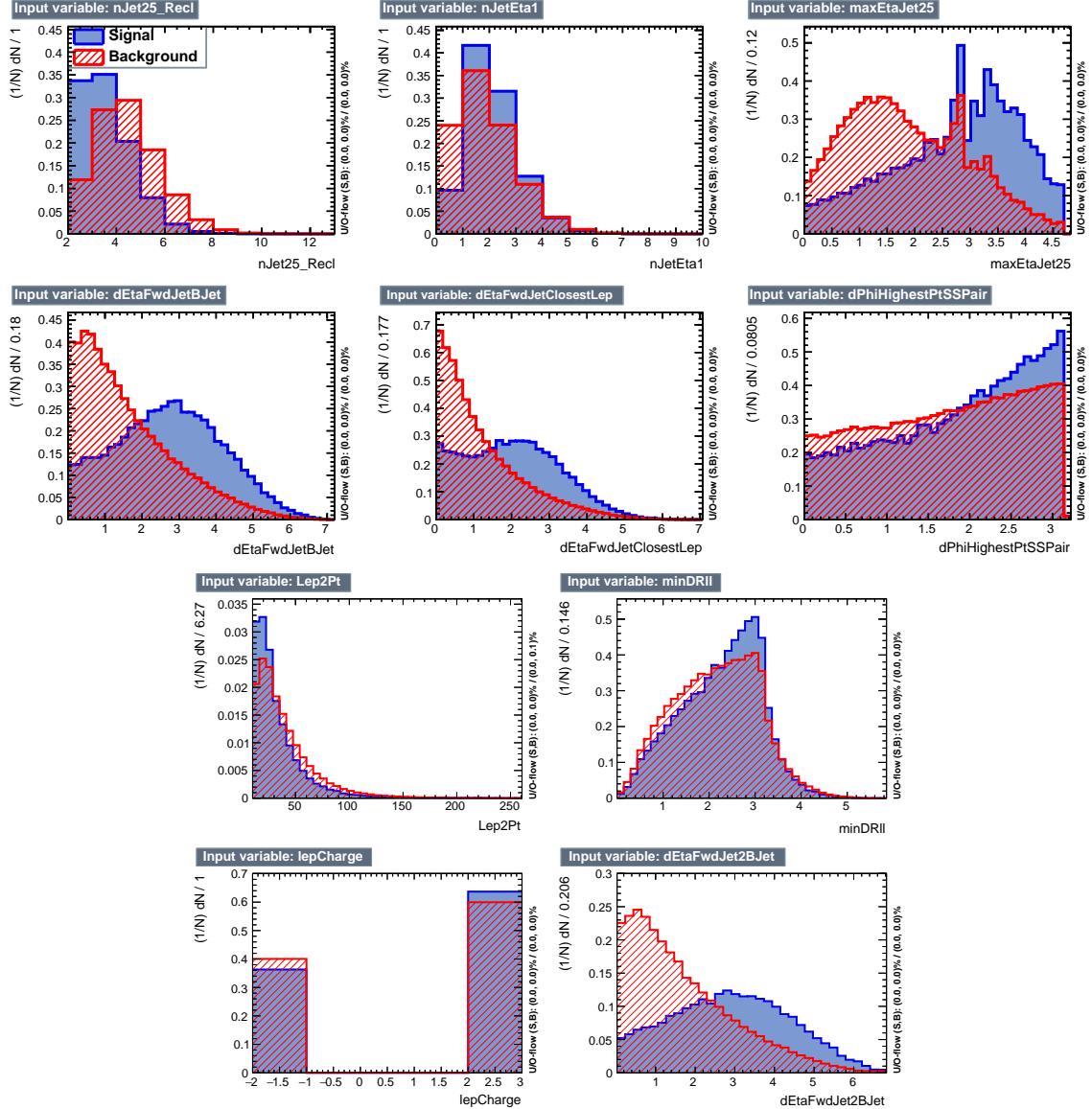


Figure B.3: BDT input variables as seen by BDTG classifier for the $2lss$ channel, tHq signal(blue) discriminated against $t\bar{t}V$ background (red).

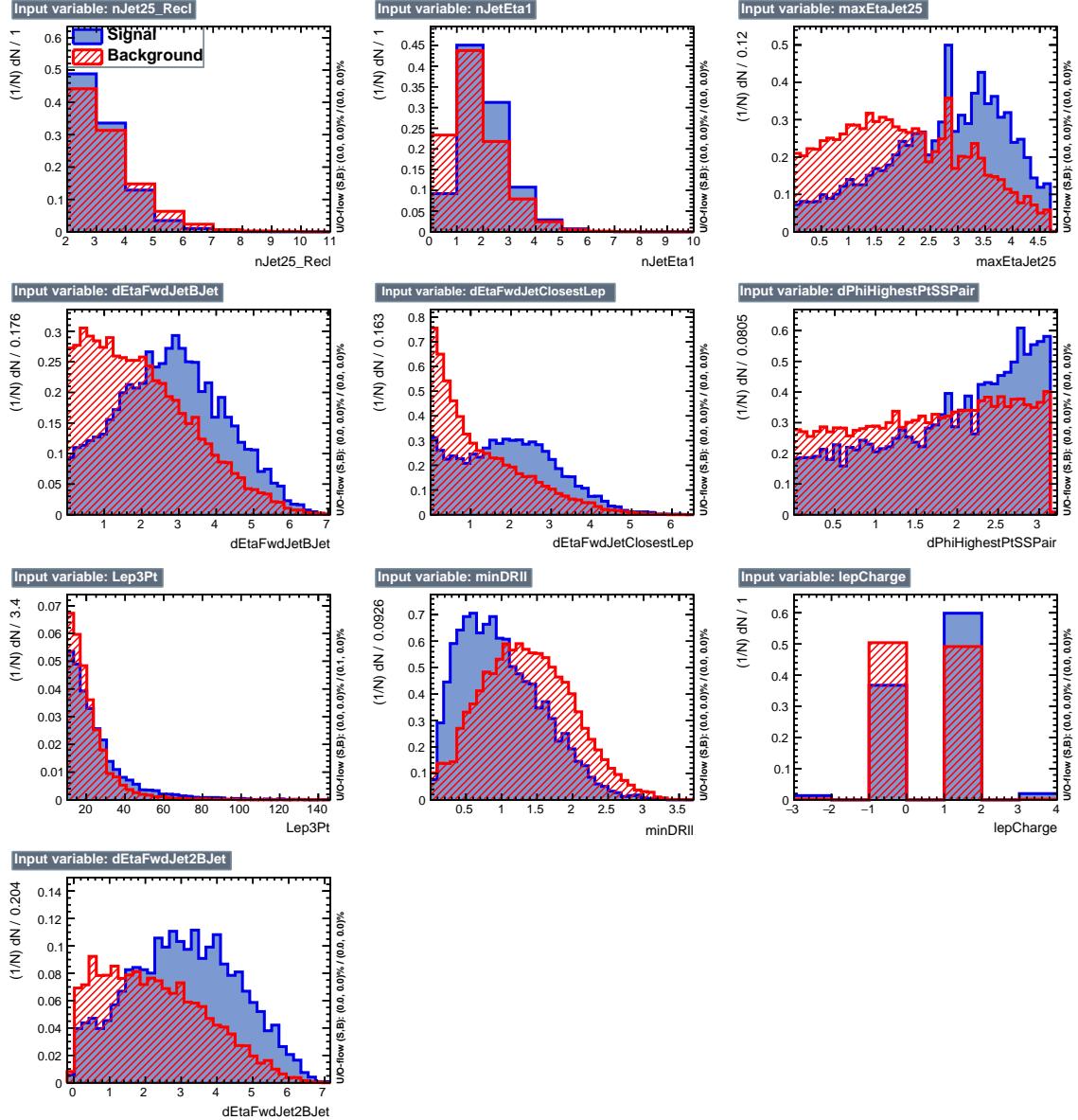


Figure B.4: BDT input variables as seen by BDTG classifier for the $3l$ channel, tHq signal (blue) discriminated against $t\bar{t}$ background (red).

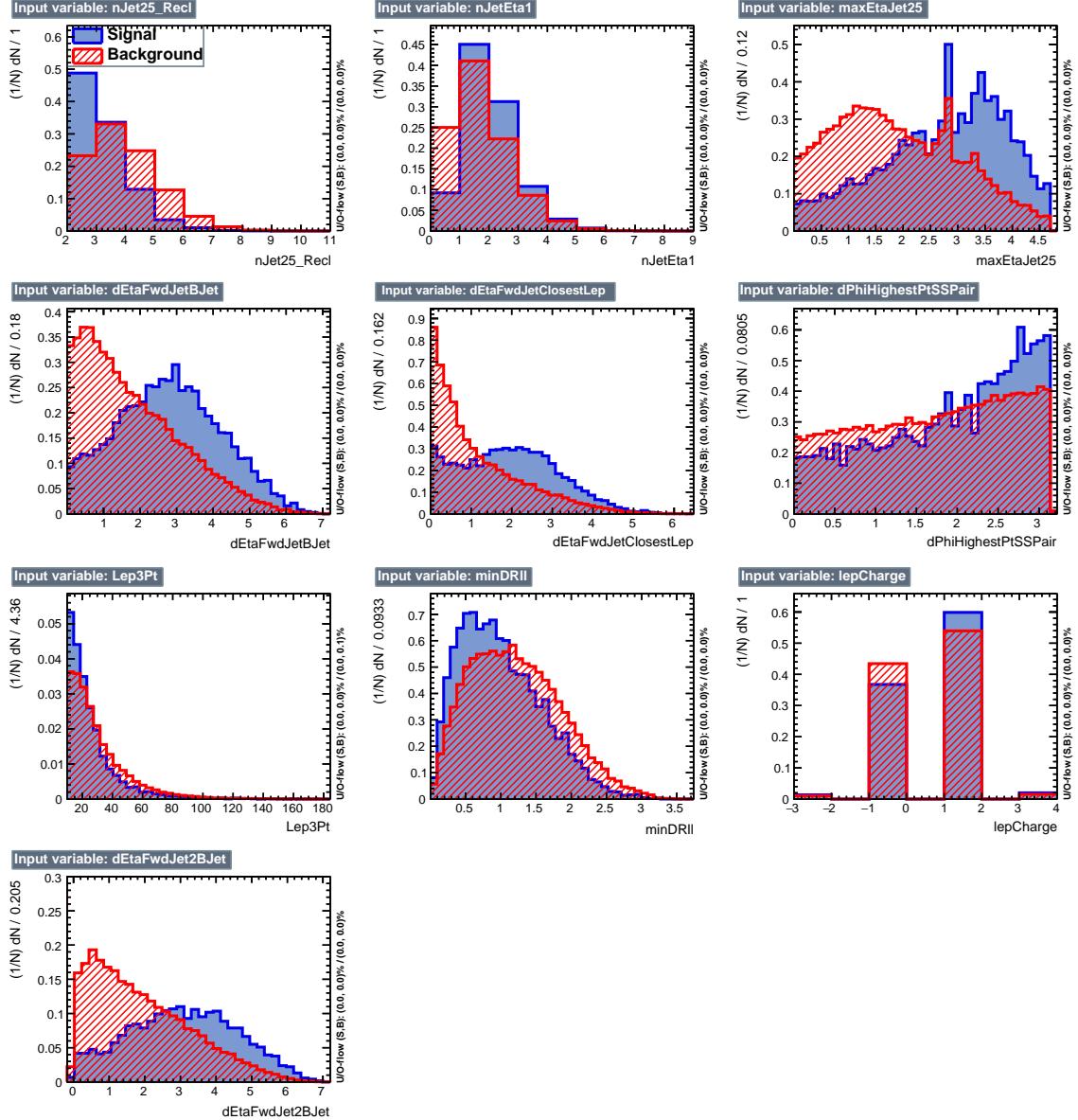


Figure B.5: BDT input variables as seen by BDTG classifier for the $3l$ channel, tHq signal (blue) discriminated against $t\bar{t}V$ background (red).

³⁰³⁷ **Appendix C**

³⁰³⁸ **Other binning strategies**

³⁰³⁹ Two additional strategies of clustering regions in the 2D plane of $BDTG_{tt}$ vs $BDTG_{ttV}$
³⁰⁴⁰ into bins were attempted, following studies done and documented in great detail in
³⁰⁴¹ Reference [149]. A brief description is provided in the following.

³⁰⁴² **Clustering by S/B ratio** In this method, the 2D plane is clustered into a given
³⁰⁴³ number of bins corresponding to regions where S/B is within a certain range. The
³⁰⁴⁴ bin borders are determined such that the number of background events in each bin is
³⁰⁴⁵ approximately equal. The resulting regions for $2lss$ and $3l$ events are shown in Figure
³⁰⁴⁶ C.1, while the expected distribution of signal and dominant backgrounds are shown
³⁰⁴⁷ in Figure C.2.

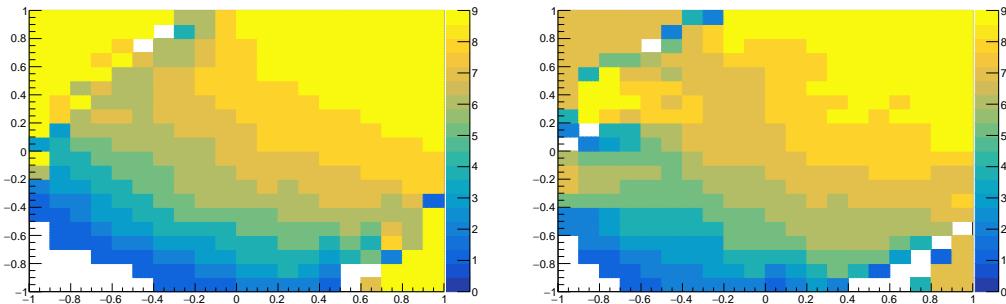


Figure C.1: Binning by S/B regions for $2lss$ (left) and $3l$ (right).

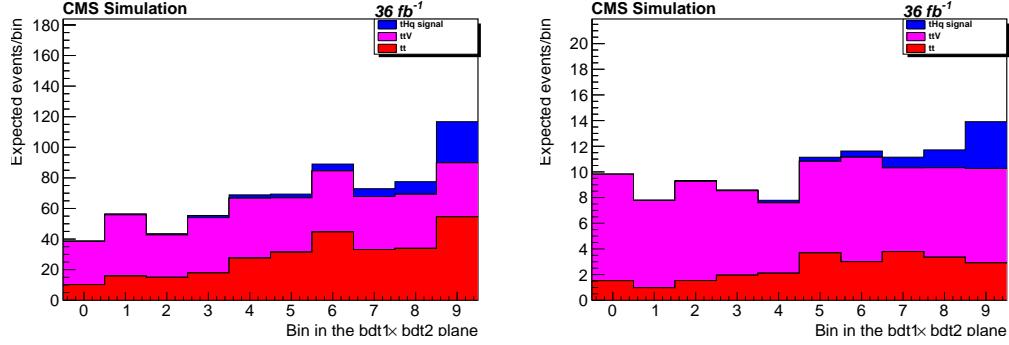


Figure C.2: Final bins (corresponding to S/B regions in the 2D plane) for $2lss$ and $3l$ (right).

Using this method, the resulting limits (for the $\kappa_t = -1, \kappa_V = 1$ scenario) are about 20% worse than with the binning in Section 6.8.6: $\mu^\pm\mu^\pm$ changed from 1.82 to 2.15, $3l$ changed from 1.52 to 1.75.

***k*-Means geometric clustering** This method employs a recursive application of the *k*-means algorithm (see Appendix D in Reference [149]) to separate the 2D plane into geometric regions. The resulting clustering (using the $t\bar{t}H$ multilepton code on tHq signal and $t\bar{t}$ and $t\bar{t}V$ background events) are shown in Figure C.3. The expected distribution of events for the signal and dominant backgrounds in these bins is shown in Fig. C.4.

Similarly to the S/B ratio binning, the limits using the *k*-means clustering are significantly worse than those of the bins described before. In the $\mu^\pm\mu^\pm$ channel, the limit deteriorates from 1.82 to 2.05, whereas in $3l$ it changes from 1.58 to 1.78.

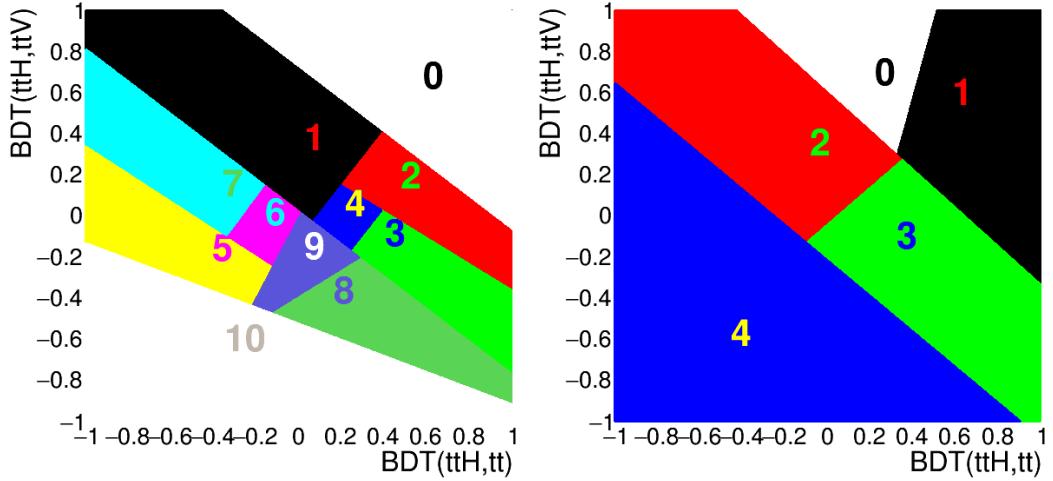


Figure C.3: Binning into geometric regions using a k -means algorithm for $2lss$ (left) and $3l$ (right).

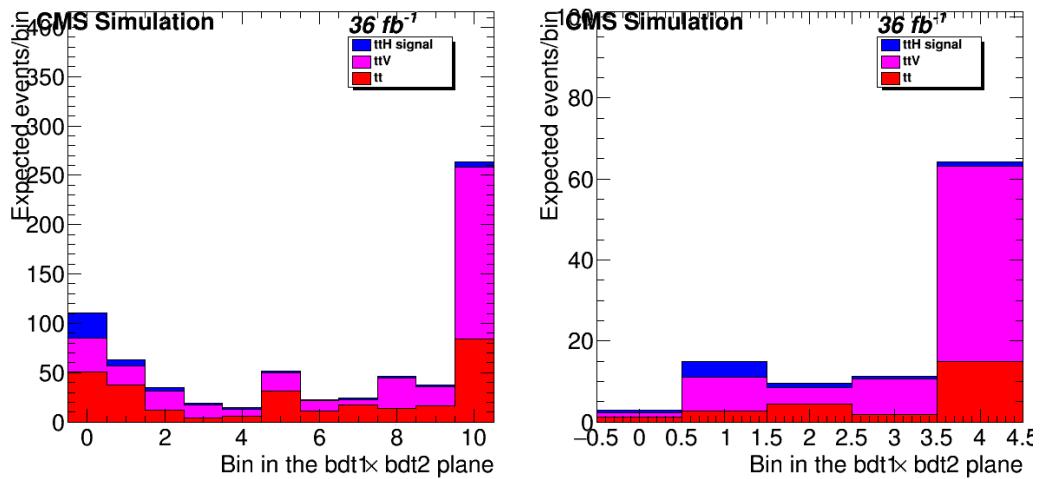


Figure C.4: Final bins using a k -means algorithm for $2lss$ (left) and $3l$ (right). Note that the bin numbering here is such that signal-like bins are lower.

3060 **Appendix D**

3061 **BDTG output variation with κ_V/κ_t**

3062 The BDTG classifier output was described in Section in the $\kappa_t = -1, \kappa_V = 1$ scenario;
 3063 the change of BDTG classifiers output shape when varying the κ_V/κ_t coupling sce-
 3064 nario is shown in Figure D.1 in the $3l$ channel for five different values of κ_t , with κ_V
 fixed at 1.0.

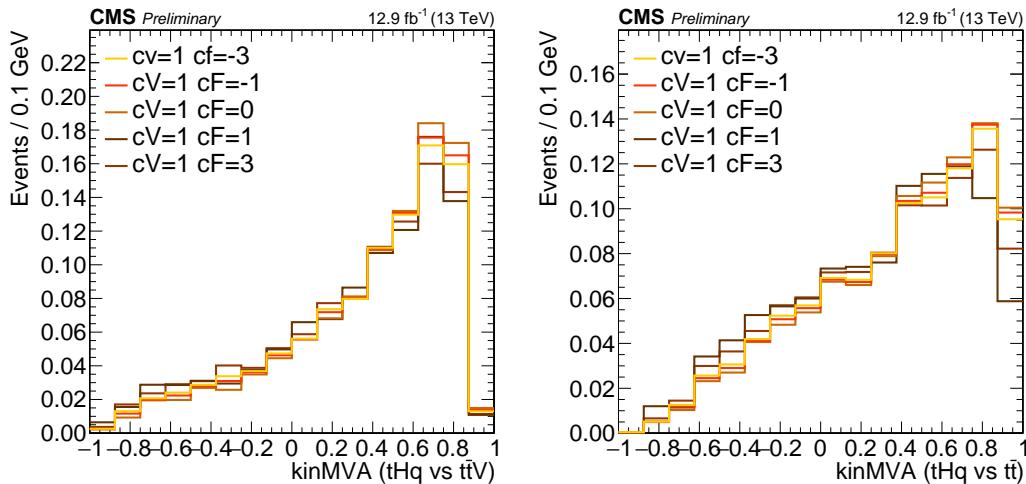


Figure D.1: Change of the BDTG classifiers output when varying κ_t coupling (κ_V is fixed at 1.0). Training vs. $t\bar{t}V$ (right) and vs. $t\bar{t}$ (left).

3065

3066 Complete this section !!!!!!! ask about this !

3067 **References**

- 3068 [1] J. Schwinger. "Quantum Electrodynamics. I. A Covariant Formulation". Physical
 3069 Review. 74 (10): 1439-61, (1948). <https://doi.org/10.1103/PhysRev.74.1439>
- 3070
- 3071 [2] R. P. Feynman. "Space-Time Approach to Quantum Electrodynamics". Physical
 3072 Review. 76 (6): 769-89, (1949). <https://doi.org/10.1103/PhysRev.76.769>
- 3073 [3] S. Tomonaga. "On a Relativistically Invariant Formulation of the Quantum
 3074 Theory of Wave Fields". Progress of Theoretical Physics. 1 (2): 27-42, (1946).
 3075 <https://doi.org/10.1143/PTP.1.27>
- 3076 [4] D.J. Griffiths, "Introduction to electrodynamics". 4th ed. Pearson, (2013).
- 3077 [5] F. Mandl, G. Shaw. "Quantum field theory." Chichester, Wiley (2009).
- 3078 [6] F. Halzen, and A.D. Martin, "Quarks and leptons: An introductory course in
 3079 modern particle physics". New York: Wiley, (1984) .
- 3080 [7] File: Standard_Model_of_Elementary_Particle_dark.svg. (2017, June 12)
 3081 Wikimedia Commons, the free media repository. Retrieved November 27, 2017
 3082 from <https://www.collegiate-advanced-electricity.com/single-post/2017/04/10/The-Standard-Model-of-Particle-Physics>.
- 3083

- 3084 [8] E. Noether, "Invariante Variationsprobleme", Nachrichten von der Gesellschaft
3085 der Wissenschaften zu Göttingen, mathematisch-physikalische Klasse, vol. 1918,
3086 pp. 235-257, (1918).
- 3087 [9] C. Patrignani et al. (Particle Data Group), Chin. Phys. C, 40, 100001 (2016)
3088 and 2017 update.
- 3089 [10] M. Goldhaber, L. Grodzins, A.W. Sunyar "Helicity of Neutrinos", Phys. Rev.
3090 109, 1015 (1958).
- 3091 [11] Palanque-Delabrouille N et al. "Neutrino masses and cosmology with Lyman-
3092 alpha forest power spectrum", JCAP 11 011 (2015).
- 3093 [12] M. Gell-Mann. "A Schematic Model of Baryons and Mesons". Physics Letters.
3094 8 (3): 214-215 (1964).
- 3095 [13] G. Zweig. "An SU(3) Model for Strong Interaction Symmetry and its Breaking"
3096 (PDF). CERN Report No.8182/TH.401 (1964).
- 3097 [14] G. Zweig. "An SU(3) Model for Strong Interaction Symmetry and its Breaking:
3098 II" (PDF). CERN Report No.8419/TH.412(1964).
- 3099 [15] M. Gell-Mann. "The Interpretation of the New Particles as Displaced Charged
3100 Multiplets". Il Nuovo Cimento 4: 848. (1956).
- 3101 [16] T. Nakano, K. Nishijima. "Charge Independence for V-particles". Progress of
3102 Theoretical Physics 10 (5): 581-582. (1953).
- 3103 [17] N. Cabibbo, "Unitary symmetry and leptonic decays" Physical Review Letters,
3104 vol. 10, no. 12, p. 531, (1963).

- 3105 [18] M.Kobayashi, T.Maskawa, “CP-violation in the renormalizable theory of weak
3106 interaction,” Progress of Theoretical Physics, vol. 49, no. 2, pp. 652-657, (1973).
- 3107 [19] File: Weak Decay (flipped).svg. (2017, June 12). Wikimedia Com-
3108 mons, the free media repository. Retrieved November 27, 2017
3109 from [https://commons.wikimedia.org/w/index.php?title=File:
3110 Weak_Decay_\(flipped\)\.svg&oldid=247498592](https://commons.wikimedia.org/w/index.php?title=File:Weak_Decay_(flipped)\.svg&oldid=247498592).
- 3111 [20] Georgia Tech University. Coupling Constants for the Fundamental Forces(2005).
3112 Retrieved January 10, 2018, from [http://hyperphysics.phy-astr.gsu.edu/
3113 hbase/Forces/couple.html#c2](http://hyperphysics.phy-astr.gsu.edu/hbase/Forces/couple.html#c2)
- 3114 [21] M. Strassler. (May 31, 2013).The Strengths of the Known Forces. Retrieved Jan-
3115 uary 10, 2018, from [https://profmattstrassler.com/articles-and-posts/
3116 particle-physics-basics/the-known-forces-of-nature/
3117 the-strength-of-the-known-forces/](https://profmattstrassler.com/articles-and-posts/particle-physics-basics/the-known-forces-of-nature/the-strength-of-the-known-forces/)
- 3118 [22] S.L. Glashow. “Partial symmetries of weak interactions”, Nucl. Phys. 22 579-
3119 588, (1961).
- 3120 [23] A. Salam, J.C. Ward. “Electromagnetic and weak interactions”, Physics Letters
3121 13 168-171, (1964).
- 3122 [24] S. Weinberg, “A model of leptons”, Physical Review Letters, vol. 19, no. 21, p.
3123 1264, (1967).
- 3124 [25] M. Peskin, D. Schroeder, “An introduction to quantum field theory”. Perseus
3125 Books Publishing L.L.C., (1995).
- 3126 [26] A. Pich. “The Standard Model of Electroweak Interactions” <https://arxiv.org/abs/1201.0537>
- 3127

- 3128 [27] G. Arnison et al. (UA1 Collaboration), Phys. Lett. B 122, 103 (1983).
- 3129 [28] M. Banner et al. (UA2 Collaboration), Phys. Lett. B 122, 476 (1983).
- 3130 [29] G. Arnison et al. (UA1 Collaboration), Phys. Lett. B 126, 398 (1983).
- 3131 [30] P. Bagnaia et al. (UA2 Collaboration), Phys. Lett. B 129, 130 (1983).
- 3132 [31] F.Bellaiche. (2012, 2 September). “What’s this Higgs boson anyway?”. Retrieved
3133 from: <https://www.quantum-bits.org/?p=233>
- 3134 [32] M. Endres et al. Nature 487, 454-458 (2012) doi:10.1038/nature11255
- 3135 [33] F. Englert, R. Brout. “Broken Symmetry and the Mass of Gauge
3136 Vector Mesons”. Physical Review Letters. 13 (9): 321-23.(1964)
3137 doi:10.1103/PhysRevLett.13.321
- 3138 [34] P.Higgs. “Broken Symmetries and the Masses of Gauge Bosons”. Physical Re-
3139 view Letters. 13 (16): 508-509,(1964). doi:10.1103/PhysRevLett.13.508.
- 3140 [35] G.Guralnik, C.R. Hagen and T.W.B. Kibble. “Global Conservation Laws
3141 and Massless Particles”. Physical Review Letters. 13 (20): 585-587, (1964).
3142 doi:10.1103/PhysRevLett.13.585.
- 3143 [36] CMS collaboration. “Observation of a new boson at a mass of 125 GeV with
3144 the CMS experiment at the LHC”. Physics Letters B. 716 (1): 30-61 (2012).
3145 arXiv:1207.7235. doi:10.1016/j.physletb.2012.08.021
- 3146 [37] ATLAS collaboration. “Observation of a New Particle in the Search for the Stan-
3147 dard Model Higgs Boson with the ATLAS Detector at the LHC”. Physics Letters
3148 B. 716 (1): 1-29 (2012). arXiv:1207.7214. doi:10.1016/j.physletb.2012.08.020.

- 3149 [38] ATLAS collaboration; CMS collaboration (26 March 2015). “Combined Mea-
 3150 surement of the Higgs Boson Mass in pp Collisions at $\sqrt{s}=7$ and 8 TeV with
 3151 the ATLAS and CMS Experiments”. Physical Review Letters. 114 (19): 191803.
 3152 arXiv:1503.07589. doi:10.1103/PhysRevLett.114.191803.
- 3153 [39] LHC InternationalMasterclasses“When protons collide”. Retrieved from http://atlas.physicsmasterclasses.org/en/zpath_protoncollisions.htm
- 3155 [40] CMS Collaboration, “SM Higgs Branching Ratios and Total Decay Widths (up-
 3156 date in CERN Report4 2016)”. <https://twiki.cern.ch/twiki/bin/view/LHCPhysics/CERNYellowReportPageBR> , last accessed on 17.12.2017.
- 3158 [41] R.Grant V. “Determination of Higgs branching ratios in $H \rightarrow W^+W^- \rightarrow l\nu jj$
 3159 and $H \rightarrow ZZ \rightarrow l^+l^-jj$ channels”. Physics Department, University of Ten-
 3160 nessee (Dated: October 31, 2012). Retrieved from <http://aesop.phys.utk.edu/ph611/2012/projects/Riley.pdf>
- 3162 [42] LHC Higgs Cross Section Working Group, Denner, A., Heinemeyer, S. et al.
 3163 “Standard model Higgs-boson branching ratios with uncertainties”. Eur. Phys.
 3164 J. C (2011) 71: 1753. <https://doi.org/10.1140/epjc/s10052-011-1753-8>
- 3165 [43] D. de Florian et al., LHC Higgs Cross Section Working Group,
 3166 CERNâš2017âš002-M, arXiv:1610.07922[hep-ph] (2016).
- 3167 [44] ATLAS and CMS Collaborations, “Measurements of the Higgs boson produc-
 3168 tion and decay rates and constraints on its couplings from a combined ATLAS
 3169 and CMS analysis of the LHC pp collision data at $\sqrt{s} = 7$ and 8 TeV,” (2016).
 3170 CERN-EP-2016-100, ATLAS-HIGG-2015-07, CMS-HIG-15-002.

- 3171 [45] J. A. Aguilar-Saavedra, R. Benbrik, S. Heinemeyer, and M. Perez-Victoria,
 3172 “Handbook of vector-like quarks: Mixing and single production”, Phys. Rev. D
 3173 88 (2013) 094010, doi:10.1103/PhysRevD.88.094010, arXiv:1306.0572.
- 3174 [46] A. Greljo, J. F. Kamenik, and J. Kopp, “Disentangling flavor vio-
 3175 lation in the top-Higgs sector at the LHC”, JHEP 07 (2014) 046,
 3176 doi:10.1007/JHEP07(2014)046, arXiv:1404.1278.
- 3177 [47] F. Demartin, F. Maltoni, K. Mawatari, and M. Zaro, “Higgs production in
 3178 association with a single top quark at the LHC,” European Physical Journal C,
 3179 vol. 75, p. 267, (2015). doi:10.1140/epjc/s10052-015-3475-9, arXiv:1504.00611.
- 3180 [48] F. Demartin, B. Maier, F. Maltoni, K. Mawatari, and M. Zaro, “tWH associated
 3181 production at the LHC”, European Physical Journal C, vol. 77, p. 34, (2017).
 3182 arXiv:1607.05862
- 3183 [49] F. Maltoni, K. Paul, T. Stelzer, and S. Willenbrock, “Associated production
 3184 of Higgs and single top at hadron colliders”, Phys.Rev. D64 (2001) 094023,
 3185 [hep-ph/0106293].
- 3186 [50] S. Biswas, E. Gabrielli, F. Margaroli, and B. Mele, “Direct constraints on the
 3187 top-Higgs coupling from the 8 TeV LHC data,” Journal of High Energy Physics,
 3188 vol. 07, p. 073, (2013).
- 3189 [51] M. Farina, C. Grojean, F. Maltoni, E. Salvioni, and A. Thamm, “Lifting de-
 3190 generacies in Higgs couplings using single top production in association with a
 3191 Higgs boson,” Journal of High Energy Physics, vol. 05, p. 022, (2013).
- 3192 [52] T.M. Tait and C.-P. Yuan, “Single top quark production as a window to physics
 3193 beyond the standard model”, Phys. Rev. D 63 (2000) 014018 [hep-ph/0007298].

- 3194 [53] CMS Collaboration, “Modelling of the single top-quark production in associa-
3195 tion with the Higgs boson at 13 TeV.” [https://twiki.cern.ch/twiki/bin/](https://twiki.cern.ch/twiki/bin/viewauth/CMS/SingleTopHiggsGeneration13TeV)
3196 [viewauth/CMS/SingleTopHiggsGeneration13TeV](#), last accessed on 16.01.2018.
- 3197 [54] CMS Collaboration, “SM Higgs production cross sections at $\sqrt{s} =$
3198 13 TeV.” [https://twiki.cern.ch/twiki/bin/](https://twiki.cern.ch/twiki/bin/view/LHCPhysics/CERNYellowReportPageAt13TeV)
3199 [view/LHCPhysics/CERNYellowReportPageAt13TeV](#), last accessed on 16.01.2018.
- 3200 [55] S. Dawson, The effective W approximation, Nucl. Phys. B 249 (1985) 42.
- 3201 [56] S. Biswas, E. Gabrielli and B. Mele, JHEP 1301 (2013) 088 [[arXiv:1211.0499](https://arxiv.org/abs/1211.0499)
3202 [hep-ph]].
- 3203 [57] LHC Higgs Cross Section Working Group, “Handbook of LHC Higgs Cross
3204 Sections: 4.Deciphering the Nature of the Higgs Sector”, [arXiv:1610.07922](https://arxiv.org/abs/1610.07922).
- 3205 [58] J. Ellis, D. S. Hwang, K. Sakurai, and M. Takeuchi.“Disentangling Higgs-Top
3206 Couplings in Associated Production”, JHEP 1404 (2014) 004, [[arXiv:1312.5736](https://arxiv.org/abs/1312.5736)].
- 3207 [59] CMS Collaboration, V. Khachatryan et al., “Precise determination of the mass
3208 of the Higgs boson and tests of compatibility of its couplings with the standard
3209 model predictions using proton collisions at 7 and 8 TeV,” [arXiv:1412.8662](https://arxiv.org/abs/1412.8662).
- 3210 [60] ATLAS Collaboration, G. Aad et al., “Updated coupling measurements of the
3211 Higgs boson with the ATLAS detector using up to 25 fb^{-1} of proton-proton
3212 collision data”, ATLAS-CONF-2014-009.
- 3213 [61] File:Cern-accelerator-complex.svg. Wikimedia Commons, the free media repos-
3214 itory. Retrieved January, 2018 from <https://commons.wikimedia.org/wiki/>
3215 [File:Cern-accelerator-complex.svg](#)

- 3216 [62] J.L. Caron , “Layout of the LEP tunnel including future LHC infrastructures.”,
3217 (Nov, 1993). A C Collection. Legacy of AC. Pictures from 1992 to 2002. Re-
3218 trieved from <https://cds.cern.ch/record/841542>
- 3219 [63] M. Vretenar, “The radio-frequency quadrupole”. CERN Yellow Report CERN-
3220 2013-001, pp.207-223 DOI:10.5170/CERN-2013-001.207. arXiv:1303.6762
- 3221 [64] L.Evans. P. Bryant (editors). “LHC Machine”. JINST 3 S08001 (2008).
- 3222 [65] CERN Photographic Service.“Radio-frequency quadrupole, RFQ-1”, March
3223 1983, CERN-AC-8303511. Retrieved from <https://cds.cern.ch/record/615852>.
- 3225 [66] CERN Photographic Service “Animation of CERN’s accelerator network”, 14
3226 October 2013. DOI: 10.17181/cds.1610170 Retrieved from <https://videos.cern.ch/record/1610170>
- 3228 [67] C.Sutton. “Particle accelerator”.Encyclopedia Britannica. July 17,
3229 2013. Retrieved from <https://www.britannica.com/technology/particle-accelerator>.
- 3231 [68] L.Guiraud. “Installation of LHC cavity in vacuum tank.”. July 27 2000. CERN-
3232 AC-0007016. Retrieved from <https://cds.cern.ch/record/41567>.
- 3233 [69] J.L. Caron, “Magnetic field induced by the LHC dipole’s superconducting coils”.
3234 March 1998. AC Collection. Legacy of AC. Pictures from 1992 to 2002. LHC-
3235 PHO-1998-325. Retrieved from <https://cds.cern.ch/record/841511>.
- 3236 [70] AC Team. “Diagram of an LHC dipole magnet”. June 1999. CERN-DI-9906025
3237 retrieved from <https://cds.cern.ch/record/40524>.

- 3238 [71] CMS Collaboration “Public CMS Luminosity Information”. https://twiki.cern.ch/twiki/bin/view/CMSPublic/LumiPublicResults#2016__proton_proton_13_TeV_collis, last accessed 24.01.2018
- 3239
- 3240
- 3241 [72] J.L Caron. “LHC Layout” AC Collection. Legacy of AC. Pictures from 1992
3242 to 2002. September 1997, LHC-PHO-1997-060. Retrieved from <https://cds.cern.ch/record/841573>.
- 3243
- 3244 [73] J.A. Coarasa. “The CMS Online Cluster:Setup, Operation and Maintenance
3245 of an Evolving Cluster”. ISGC 2012, 26 February - 2 March 2012, Academia
3246 Sinica, Taipei, Taiwan.
- 3247 [74] CMS Collaboration. “The CMS experiment at the CERN LHC” JINST 3 S08004
3248 (2008).
- 3249 [75] CMS Collaboration. “CMS detector drawings 2012” CMS-PHO-GEN-2012-002.
3250 Retrieved from <http://cds.cern.ch/record/1433717>.
- 3251 [76] Davis, Siona Ruth. “Interactive Slice of the CMS detector”, Aug. 2016,
3252 CMS-OUTREACH-2016-027, retrieved from <https://cds.cern.ch/record/2205172>
- 3253
- 3254 [77] R. Breedon. “View through the CMS detector during the cooldown of the
3255 solenoid on February 2006. CMS Collection”, February 2006, CMS-PHO-
3256 OREACH-2005-004, Retrieved from <https://cds.cern.ch/record/930094>.
- 3257 [78] Halyo, V. and LeGresley, P. and Lujan, P. “Massively Parallel Computing and
3258 the Search for Jets and Black Holes at the LHC”, Nucl.Instrum.Meth. A744
3259 (2014) 54-60, DOI: 10.1016/j.nima.2014.01.038”

- 3260 [79] A. Dominguez et. al. “CMS Technical Design Report for the Pixel Detector
3261 Upgrade”, CERN-LHCC-2012-016. CMS-TDR-11.
- 3262 [80] CMS Collaboration. “Description and performance of track and primary-vertex
3263 reconstruction with the CMS tracker,” Journal of Instrumentation, vol. 9, no.
3264 10, p. P10009,(2014).
- 3265 [81] CMS Collaboration and M. Brice. “Images of the CMS Tracker Inner Bar-
3266 rel”, November 2008, CMS-PHO-TRACKER-2008-002. Retrieved from <https://cds.cern.ch/record/1431467>.
- 3268 [82] M. Weber. “The CMS tracker”. 6th international conference on hyperons, charm
3269 and beauty hadrons Chicago, June 28-July 3 2004.
- 3270 [83] CMS Collaboration. “Projected Performance of an Upgraded CMS Detector at
3271 the LHC and HL-LHC: Contribution to the Snowmass Process”. Jul 26, 2013.
3272 arXiv:1307.7135
- 3273 [84] L. Veillet. “End assembly of HB with EB rails and rotation inside SX ”,Jan-
3274 uary 2002. CMS-PHO-HCAL-2002-002. Retrieved from <https://cds.cern.ch/record/42594>.
- 3276 [85] J. Puerta-Pelayo.“First DT+RPC chambers installation round in the UX5 cav-
3277 ern.”. January 2007, CMS-PHO-OREACH-2007-001. Retrieved from <https://cds.cern.ch/record/1019185>
- 3279 [86] X. Cid Vidal and R. Cid Manzano. “CMS Global Muon Trigger” web site:
3280 Taking a closer look at LHC. Retrieved from https://www.lhc-closer.es/taking_a_closer_look_at_lhc/0.lhc_trigger

- 3282 [87] WLCG Project Office, “Documents & Reference - Tiers - Structure,”
3283 (2014). <http://wlcg.web.cern.ch/documents-reference> , last accessed on
3284 30.01.2018.
- 3285 [88] CMS Collaboration. “CMSSW Application Framework”, <https://twiki.cern.ch/twiki/bin/view/CMSPublic/WorkBookCMSSWFramework>,
3286 last accesses 06.02.2018
- 3288 [89] A. Buckleya, J. Butterworthb, S. Giesekec, et. al. “General-purpose event gen-
3289 erators for LHC physics”. arXiv:1101.2599v1 [hep-ph] 13 Jan 2011
- 3290 [90] A. Quadt. “Top Quark Physics at Hadron Colliders”. Advances in the Physics
3291 of Particles and Nuclei. Springer-Verlag Berlin Heidelberg. DOI: 10.1007/978-
3292 3-540-71060-8 (2007)
- 3293 [91] DurhamHep Data Project, “The Durham HepData Project - PDF Plotter.”
3294 <http://hepdata.cedar.ac.uk/pdf/pdf3.html> , last accessed on 02.02.2018.
- 3295 [92] G. Altarelli and G. Parisi. “ASYMPTOTIC FREEDOM IN PARTON LAN-
3296 GUAGE”, Nucl.Phys. B126:298 (1977).
- 3297 [93] Yu.L. Dokshitzer. Sov.Phys. JETP 46:641 (1977)
- 3298 [94] V.N. Gribov, L.N. Lipatov. “Deep inelastic e p scattering in perturbation the-
3299 ory”, Sov.J.Nucl.Phys. 15:438 (1972)
- 3300 [95] F. Maltoni, G. Ridolfi, and M. Ubiali, “b-initiated processes at the LHC: a
3301 reappraisal,” Journal of High Energy Physics, vol. 07, p. 022, (2012).
- 3302 [96] B. Andersson, G. Gustafson, G.Ingelman and T. Sjostrand, “Parton fragmen-
3303 tation and string dynamics”, Physics Reports, Vol. 97, No. 2-3, pp. 31-145,
3304 1983.

- 3305 [97] CMS Collaboration, “Event generator tunes obtained from underlying event
3306 and multiparton scattering measurements;” European Physical Journal C, vol.
3307 76, no. 3, p. 155, (2016).
- 3308 [98] J. Alwall et. al., “The automated computation of tree-level and next-to-leading
3309 order differential cross sections, and their matching to parton shower simula-
3310 tions,” Journal of High Energy Physics, vol. 07, p. 079, (2014).
- 3311 [99] T. Sjöstrand and P. Z. Skands, “Transverse-momentum-ordered showers and
3312 interleaved multiple interactions,” European Physical Journal C, vol. 39, pp.
3313 129–154, (2005).
- 3314 [100] S. Frixione, P. Nason, and C. Oleari, “Matching NLO QCD computations with
3315 Parton Shower simulations: the POWHEG method,” Journal of High Energy
3316 Physics, vol. 11, p. 070, (2007).
- 3317 [101] S. Agostinelli et al., “GEANT4: A Simulation toolkit,” Nuclear Instruments
3318 and Methods in Physics, vol. A506, pp. 250–303, (2003).
- 3319 [102] J.Allison et.al.,“Recent developments in Geant4”, Nuclear Instruments and
3320 Methods in Physics Research A 835 (2016) 186-225.
- 3321 [103] CMS Collaboration “Full Simulation Offline Guide”, <https://twiki.cern.ch/twiki/bin/view/CMSPublic/SWGuideSimulation>, last accessed 04.02.2018
- 3323 [104] A. Giammanco. “The Fast Simulation of the CMS Experiment” J. Phys.: Conf.
3324 Ser. 513 022012 (2014)
- 3325 [105] A.M. Sirunyan et. al. “Particle-flow reconstruction and global event description
3326 with the CMS detector”, JINST 12 P10003 (2017) <https://doi.org/10.1088/1748-0221/12/10/P10003>.

- 3328 [106] The CMS Collaboration. “ Description and performance of track and pri-
 3329 mary vertex reconstruction with the CMS tracker”. JINST 9 P10009 (2014).
 3330 doi:10.1088/1748-0221/9/10/P10009
- 3331 [107] J. Incandela. “Status of the CMS SM Higgs Search” July 4, 2012. Pdf slides.
 3332 Retrieved from https://indico.cern.ch/event/197461/contributions/1478917/attachments/290954/406673/CMS_4July2012_Final.pdf
- 3334 [108] P. Billoir and S. Qian, “Simultaneous pattern recognition and track fitting by
 3335 the Kalman filtering method”, Nucl. Instrum. Meth. A 294 219. (1990).
- 3336 [109] W. Adam, R. Fruhwirth, A. Strandlie and T. Todorov, “Reconstruction of
 3337 electrons with the Gaussian sum filter in the CMS tracker at LHC”, eConf
 3338 C 0303241 (2003) TULT009 [physics/0306087].
- 3339 [110] K. Rose, “Deterministic Annealing for Clustering, Compression, Classification,
 3340 Regression and related Optimisation Problems”, Proc. IEEE 86 (1998) 2210.
- 3341 [111] R. Fruhwirth, W. Waltenberger and P. Vanlaer, “ Adaptive Vertex Fitting”,
 3342 CMS Note 2007-008 (2007).
- 3343 [112] CMS collaboration, “Performance of CMS muon reconstruction in pp collision
 3344 events at $\sqrt{s} = 7 \text{ TeV}$ ”, JINST 7 P10002 2012, [arXiv:1206.4071].
- 3345 [113] Coco, Victor and Delsart, Pierre-Antoine and Rojo-Chacon, Juan and Soyez,
 3346 Gregory and Sander, Christian, “Jets and jet algorithms”, Proceedings,
 3347 HERA and the LHC Workshop Series on the implications of HERA for LHC
 3348 physics: 2006-2008, pag. 182-204. <http://inspirehep.net/record/866539/files/access.pdf>, (2009), doi:10.3204/DESY-PROC-2009-02/54

- 3350 [114] M. Cacciari, G. P. Salam, and G. Soyez, “The anti- k_t jet clustering algorithm,”
3351 Journal of High Energy Physics, vol. 04, p. 063, (2008).
- 3352 [115] S. Catani, Y. L. Dokshitzer, M. H. Seymour, and B. R. Webber, “Longitudi-
3353 nally invariant K_t clustering algorithms for hadron hadron collisions”, Nuclear
3354 Physics B, vol. 406, pp. 187–224, (1993).
- 3355 [116] Y.L. Dokshitzer, G.D. Leder, S.Moretti, and B.R. Webber, “Better jet clustering
3356 algorithms,” Journal of High Energy Physics, vol. 08, p. 001, (1997).
- 3357 [117] B. Dorney. “Anatomy of a Jet in CMS”. Quantum Diaries. June
3358 1st, 2011. Retrieved from [https://www.quantumdiaries.org/2011/06/01/
3359 anatomy-of-a-jet-in-cms/](https://www.quantumdiaries.org/2011/06/01/anatomy-of-a-jet-in-cms/)
- 3360 [118] The CMS Collaboration.“Event Displays from the high-energy collisions at 7
3361 TeV”, May 2010, CMS-PHO-EVENTS-2010-007, Retrieved from [https://cds.
3362 cern.ch/record/1429614](https://cds.cern.ch/record/1429614).
- 3363 [119] The CMS collaboration. “Determination of jet energy calibration and transverse
3364 momentum resolution in CMS”. JINST 6 P11002 (2011). [http://dx.doi.org/
3365 10.1088/1748-0221/6/11/P11002](http://dx.doi.org/10.1088/1748-0221/6/11/P11002)
- 3366 [120] The CMS Collaboration, “Introduction to Jet Energy Corrections at
3367 CMS.”. <https://twiki.cern.ch/twiki/bin/view/CMS/IntroToJEC>, last ac-
3368 cessed 10.02.2018.
- 3369 [121] CMS Collaboration Collaboration. “Identification of b quark jets at the CMS
3370 Experiment in the LHC Run 2”. Tech. rep. CMS-PAS-BTV-15-001. Geneva:
3371 CERN, (2016). <https://cds.cern.ch/record/2138504>.

- 3372 [122] CMS Collaboration Collaboration. “Performance of missing energy reconstruction
3373 in 13 TeV pp collision data using the CMS detector”. Tech. rep. CMS-PAS-
3374 JME16-004. Geneva: CERN, 2016. <https://cds.cern.ch/record/2205284>.
- 3375 [123] CMS Collaboration, “New CMS results at Moriond (Electroweak) 2013”,
3376 Retrieved from http://cms.web.cern.ch/sites/cms.web.cern.ch/files/styles/large/public/field/image/HIG13004_Event01_0.png?itok=LAWZzPHR
- 3379 [124] CMS Collaboration, “New CMS results at Moriond (Electroweak) 2013”,
3380 Retrieved from http://cms.web.cern.ch/sites/cms.web.cern.ch/files/styles/large/public/field/image/TOP12035_Event01.png?itok=uMdnSqzC
- 3383 [125] K. Skovpen. “Event displays highlighting the main properties of heavy flavour
3384 jets in the CMS Experiment”, Aug 2017, CMS-PHO-EVENTS-2017-006. Re-
3385 trieval from <https://cds.cern.ch/record/2280025>.
- 3386 [126] G. Cowan. “Topics in statistical data analysis for high-energy physics”.
3387 arXiv:1012.3589v1
- 3388 [127] A. Hoecker et al., “TMVA-Toolkit for multivariate data analysis”
3389 arXiv:physics/0703039v5 (2009)
- 3390 [128] L. Lista. “Statistical Methods for Data Analysis in Particle Physics”, 2nd
3391 ed. Springer International Publishing. (2017) <https://dx.doi.org/10.1007/978-3-319-62840-0>

- 3393 [129] I. Antcheva et al., “ROOT-A C++ framework for petabyte data storage, sta-
3394 tistical analysis and visualization ,” Computer Physics Communications, vol.
3395 182, no. 6, pp. 1384â€¢1385, (2011).
- 3396 [130] Y. Coadou. “Boosted decision trees”, ESIPAP, Archamps, 9 Febru-
3397 ary 2016. Lecture. Retrieved from https://indico.cern.ch/event/472305/contributions/1982360/attachments/1224979/1792797/ESIPAP_MVA160208-BDT.pdf
- 3400 [131] J.H. Friedman. “Greedy function approximation: A gradient boosting ma-
3401 chine”. Ann. Statist. Volume 29, Number 5 (2001), 1189-1232. https://projecteuclid.org/download/pdf_1/euclid-aos/1013203451.
- 3403 [132] W. Verkerke and D. Kirkby, “The RooFit toolkit for data modeling,” arXiv
3404 preprint physics, (2003).
- 3405 [133] CMS Collaboration, “Documentation of the RooStats-based statistics
3406 tools for Higgs PAG”. <https://twiki.cern.ch/twiki/bin/view/CMS/SWGuideHiggsAnalysisCombinedLimit>, last accessed on 08.04.2018.
- 3408 [134] F. James, M. Roos, “MINUIT: Function minimization and error analysis”. Cern
3409 Computer Centre Program Library, Geneve Long Write-up No. D506, 1989
- 3410 [135] J. Neyman and E. S. Pearson, “On the problem of the most efficient tests of
3411 statistical hypotheses”. Springer-Verlag, (1992).
- 3412 [136] A.L. Read. “Modified frequentist analysis of search results (the CL_s method),”
3413 (2000). CERN-OPEN-2000-205.
- 3414 [137] C. Palmer. “Searches for a Light Higgs with CMS”, CMS-CR-2012-215. <https://cds.cern.ch/record/1560435>.

- 3416 [138] A. Wald, “Tests of statistical hypotheses concerning several parameters when
 3417 the number of observations is large”, Transactions of the American Mathematical
 3418 society, vol. 54, no. 3, pp. 426–482, (1943).
- 3419 [139] G. Cowan, K. Cranmer, E. Gross, and O. Vitells, “Asymptotic formulae for
 3420 likelihood-based tests of new physics”, European Physical Journal C, vol. 71,
 3421 p. 1554, (2011).
- 3422 [140] S. S. Wilks, “The Large-Sample Distribution of the Likelihood Ratio for Testing
 3423 Composite Hypotheses”, Annals of Mathematical Statistics, vol. 9, pp. 60–62,
 3424 (03, 1938).
- 3425 [141] B. Hespel, F. Maltoni, and E. Vryonidou, “Higgs and Z boson associated pro-
 3426 duction via gluon fusion in the SM and the 2HDM”, JHEP 06 (2015) 065,
 3427 [https://dx.doi.org/10.1007/JHEP06\(2015\)065](https://dx.doi.org/10.1007/JHEP06(2015)065), arXiv:1503.01656.
- 3428 [142] ATLAS Collaboration, “Measurements of Higgs boson pro-
 3429 duction and couplings in diboson final states with the AT-
 3430 LAS detector at the LHC”, Phys. Lett. B726 (2013) 88–119,
 3431 doi:10.1016/j.physletb.2014.05.011, 10.1016/j.physletb.2013.08.010,
 3432 arXiv:1307.1427. [Erratum: Phys. Lett.B734,406(2014)].
- 3433 [143] CMS Collaboration, “Search for the associated production of a Higgs boson
 3434 with a single top quark in proton-proton collisions at $\sqrt{s} = 8$ TeV”, JHEP 06
 3435 (2016) 177, doi:10.1007/JHEP06(2016)177, arXiv:1509.08159.
- 3436 [144] B. Stieger, C. Jorda Lope et al., “Search for Associated Production of a Single
 3437 Top Quark and a Higgs Boson in Leptonic Channels”, CMS Analysis Note CMS
 3438 AN-14-140, 2014.

- 3439 [145] M. Peruzzi, C. Mueller, B. Stieger et al., “Search for ttH in multilepton final
 3440 states at $\sqrt{s} = 13$ TeV”, CMS Analysis Note CMS AN-16-211, 2016.
- 3441 [146] CMS Collaboration, “Search for H to bbar in association with a single top quark
 3442 as a test of Higgs boson couplings at $\sqrt{s} = 13$ TeV”, CMS Physics Analysis
 3443 Summary CMS-PAS-HIG-16-019, 2016.
- 3444 [147] CMS Collaboration, “Search for production of a Higgs boson and a single top
 3445 quark in multilepton final states in proton collisions at $\sqrt{s} = 13$ TeV”, CMS
 3446 Physics Analysis Summary CMS-PAS-HIG-17-005, 2016.
- 3447 [148] CMS Collaboration, “PdmV2016Analysis,” (2016). <https://twiki.cern.ch/twiki/bin/viewauth/CMS/PdmV2016Analysis#DATA>, last accessed 11.04.2016.
- 3449 [149] M. Peruzzi, F. Romeo, B. Stieger et al., “Search for ttH in multilepton final1
 3450 states at $\sqrt{s} = 13$ TeV”, CMS Analysis Note CMS AN-17-029, 2017.
- 3451 [150] B. Maier, “SingleTopHiggProduction13TeV”, February, 2016. <https://twiki.cern.ch/twiki/bin/viewauth/CMS/SingleTopHiggsGeneration13TeV>.
- 3453 [151] B. WG, “BtagRecommendation80XReReco”, February, 2017. <https://twiki.cern.ch/twiki/bin/view/CMS/BtagRecommendation80XReReco>.
- 3455 [152] CMS Collaboration, “Identification of b quark jets at the CMS Experiment
 3456 in the LHC Run 2”, CMS Physics Analysis Summary CMS-PAS-BTV-15-001,
 3457 2016.
- 3458 [153] CMS Collaboration, “Baseline muon selections for Run-II.” <https://twiki.cern.ch/twiki/bin/view/CMSPublic/SWGuideMuonIdRun2>, last accessed on
 3459 24.02.2018.

- 3461 [154] G. Petrucciani and C. Botta, “Two step prompt muon identification”, January,
3462 2015. [https://indico.cern.ch/event/368007/contribution/2/material/
3463 slides/0.pdf](https://indico.cern.ch/event/368007/contribution/2/material/slides/0.pdf).
- 3464 [155] H. Brun and C. Ochando, “Updated Results on MVA eID with 13 TeV samples”,
3465 October, 2014. [https://indico.cern.ch/event/298249/contribution/3/
3466 material/slides/0.pdf](https://indico.cern.ch/event/298249/contribution/3/material/slides/0.pdf).
- 3467 [156] K. Rehermann and B. Tweedie, “Efficient Identification of Boosted Semileptonic
3468 Top Quarks at the LHC”, JHEP 03 (2011) 059, [https://dx.doi:10.1007/
3469 JHEP03\(2011\)059](https://dx.doi.org/10.1007/JHEP03(2011)059), arXiv:1007.2221.
- 3470 [157] CMS Collaboration. “Tag and Probe”, [https://twiki.cern.ch/twiki/bin/
3471 view/CMS/TagAndProbe](https://twiki.cern.ch/twiki/bin/view/CMS/TagAndProbe), last accessed on 02.03.2018.